

Azure Machine Learning – Hands On

Mohammad Aadil Sehrawat

Employee Code: 654985

The screenshot shows the 'Create new workspace' wizard in the Azure portal, specifically the 'Basics' tab. The page is titled 'Azure Machine Learning' with a subtitle 'Create a machine learning workspace'. The navigation bar includes 'Home > Azure Machine Learning >'. The 'Basics' tab is selected, with other tabs being 'Networking', 'Encryption', 'Identity', 'Tags', and 'Review + create'. The 'Resource details' section explains that every workspace must be assigned to an Azure subscription and resource group. It shows a dropdown for 'Subscription' with the value 'npunext-1680262051301' and an empty dropdown for 'Resource group'. A 'Create new' link is present below the resource group dropdown. The 'Workspace details' section prompts the user to configure basic settings like storage connection, authentication, container, and more, with a 'Learn more' link. It includes input fields for 'Name', 'Region' (set to 'East US'), 'Storage account' (with a 'Create new' link), and 'Key vault' (with a 'Create new' link'). At the bottom, there are three buttons: 'Review + create' (highlighted in blue), '< Previous', and 'Next : Networking'.

portal.azure.com/#view/Microsoft_Azure_MLTeamAccounts/CreateMachineLearningServicesBladeV2/_provisioningContext~/~/7B*initialVal...

Microsoft Azure Search resources, services, and docs (G+)

Home > Azure Machine Learning >

Azure Machine Learning

Create a machine learning workspace

Basics Networking Encryption Identity Tags Review + create

Resource details

Every workspace must be assigned to an Azure subscription, which is where billing happens. You use resource groups like folders to organize and manage resources, including the workspace you're about to create. [Learn more about Azure resource groups](#)

Subscription * ⓘ npunext-1680262051301

Resource group * ⓘ [Create new](#)

Workspace details

Configure your basic workspace settings like its storage connection, authentication, container, and more. [Learn more](#)

Name * ⓘ

Region * ⓘ East US

Storage account * ⓘ [Create new](#)

Key vault * ⓘ [Create new](#)

[Review + create](#) < Previous Next : Networking

Creating a new Azure Machine Learning Instance

Lumen x ML_Assessment_Paper.pdf x Azure Machine Learning - Micro x +

portal.azure.com/#view/Microsoft_Azure_MLTeamAccounts/CreateMachineLearningServicesBladeV2/_provisioningContext~/~/7B*initialVal...

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

Home > Azure Machine Learning >

Azure Machine Learning

Create a machine learning workspace

Resource details

Every workspace must be assigned to an Azure subscription, which is where billing happens. You use resource groups like folders to organize and manage resources, including the workspace you're about to create. [Learn more about Azure resource groups](#)

Subscription * ⓘ npunext-1680262051301

Resource group * ⓘ aadil-rg [Create new](#)

Workspace details

Configure your basic workspace settings like its storage connection, authentication, container, and more. [Learn more](#)

Name * ⓘ aadil-aml ✓

Region * ⓘ East US

Storage account * ⓘ (new) aadilaml9069656161 [Create new](#)

Key vault * ⓘ (new) aadilaml0534132963 [Create new](#)

Application insights * ⓘ (new) aadilaml4959900418 [Create new](#)

[Review + create](#) < Previous Next : Networking

Configuration for new Azure Machine Learning Instance

Lumen x ML_Assessment_Paper.pdf x Azure Machine Learning - Micro x +

portal.azure.com/#view/Microsoft_Azure_MLTeamAccounts/CreateMachineLearningServicesBladeV2/_provisioningContext~/~/7B*initialVal...

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

Home > Azure Machine Learning >

Azure Machine Learning

Create a machine learning workspace

✓ Validation passed

Basics Networking Encryption Identity Tags Review + create

Basics

Subscription	npunext-1680262051301
Resource group	aadil-rg
Region	East US
Name	aadil-aml
Storage account	(new) aadilaml9069656161
Key vault	(new) aadilaml0534132963
Application insights	(new) aadilaml4959900418
Container registry	None

Networking

Connectivity method	Enable public access from all networks
Network isolation	Public

Encryption

[Create](#) < Previous Next > [Download a template for automation](#)

Creation under progress

Lumen ML_Assessment_Paper.pdf Microsoft.MachineLearningServices

portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions%2Fa017e82e-69e8-4283-b4a5-6ac9e21...

Microsoft Azure Search resources, services, and docs (G+)

Shellunext_1693422709... UNEXT (NPUNEXT.ONMICROSO...

Microsoft.MachineLearningServices | Overview

Deployment

Search

Delete Cancel Redeploy Download Refresh

Overview Inputs Outputs Template

Deployment is in progress

Deployment name : Microsoft.MachineLearningSer... Start time : 10/4/2023, 1:33:06 PM
Subscription : npunext-1680262051301 Correlation ID : fe351e15-d873-46c6-a553-592...
Resource group : aadil-rg

Deployment details

Resource	Type	Status	Operation
There are no resources to display.			

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 >](#)

Deployment in progress

Lumen ML_Assessment_Paper.pdf Shell2023/AssessmentData/cust Microsoft.MachineLearningServices

portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions%2Fa017e82e-69e8-4283-b4a5-6ac9e21...

Microsoft Azure Search resources, services, and docs (G+)

Shellunext_1693422709... UNEXT (NPUNEXT.ONMICROSO...

Microsoft.MachineLearningServices | Overview

Deployment

Search

Delete Cancel Redeploy Download Refresh

Overview Inputs Outputs Template

Your deployment is complete

Deployment name : Microsoft.MachineLearningSer... Start time : 10/4/2023, 1:33:06 PM
Subscription : npunext-1680262051301 Correlation ID : fe351e15-d873-46c6-a553-592...
Resource group : aadil-rg

Deployment details

Next steps

[Go to resource](#)

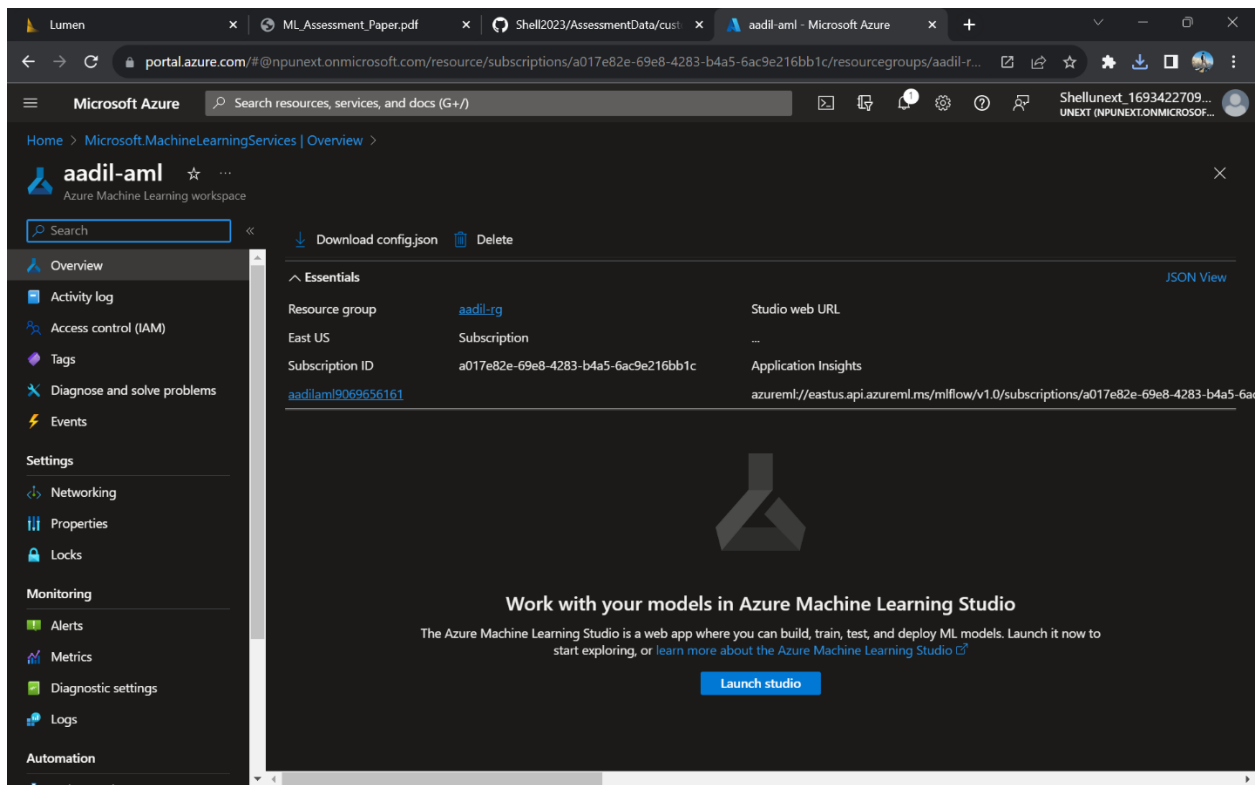
Cost management
Get notified to stay within your budget and prevent unexpected charges on your bill.
[Set up cost alerts >](#)

Microsoft Defender for Cloud
Secure your apps and infrastructure
[Go to Microsoft Defender for Cloud >](#)

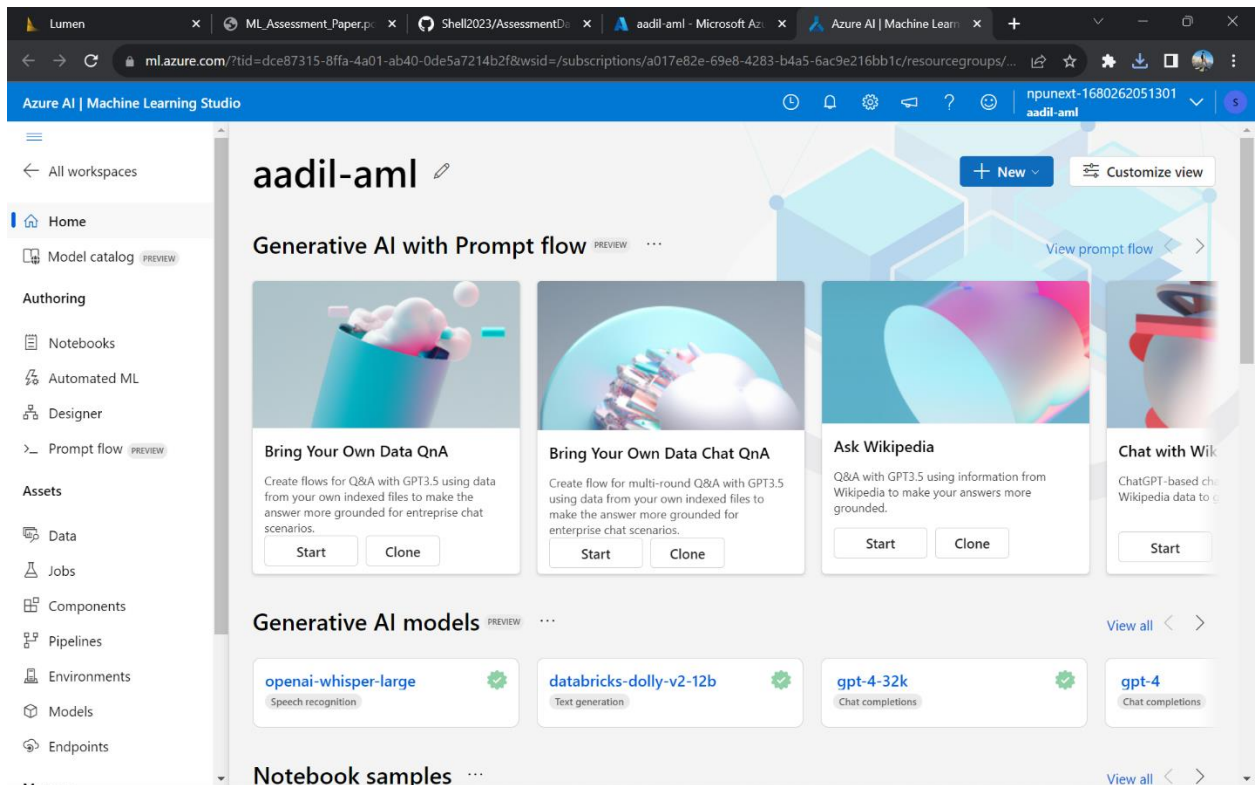
Free Microsoft tutorials
[Start learning today >](#)

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 >](#)

Deployment Complete

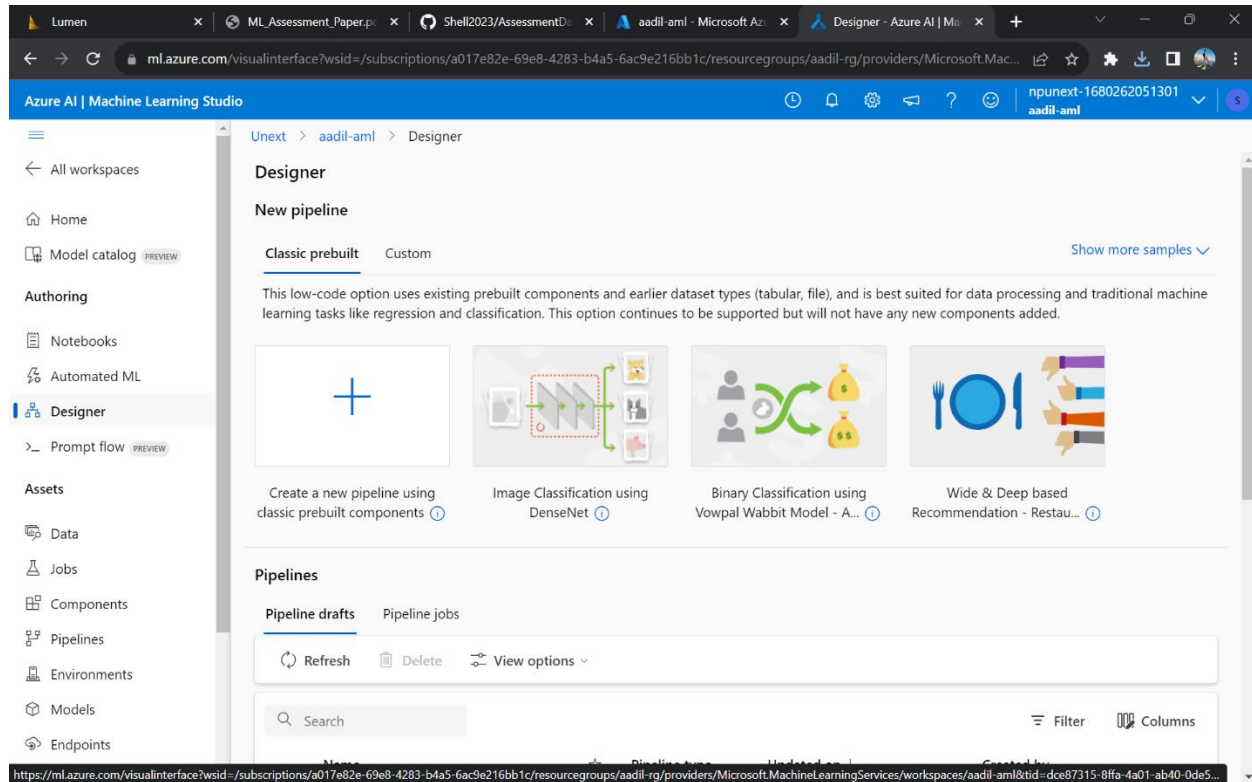


Resource

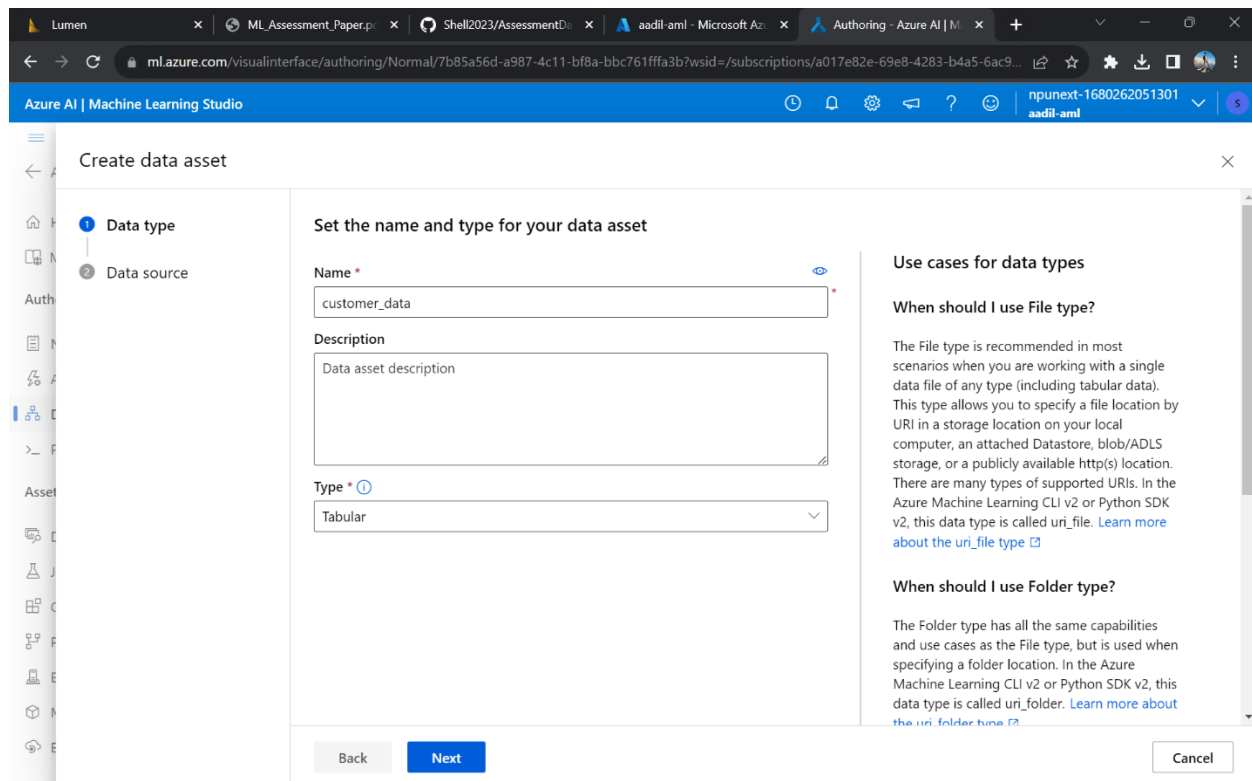


Azure ML Studio

1. Data Preparation



Azure ML Studio Designer



Creating Data-Asset

ml.azure.com/visualinterface/authoring/Normal/7b85a56d-a987-4c11-bf8a-bbc761ffa3b7?wsid=/subscriptions/a017e82e-69e8-4283-b4a5-6ac9...

Azure AI | Machine Learning Studio

Create data asset

- ✓ Data type
- ✓ Data source
- ✓ Destination storage type
- 4 File or folder selection
- 5 Settings
- 6 Schema
- 7 Review

Choose a file or folder

Choose files or folders to upload from your local drive. If you upload multiple folders or files, they will be stored in a containing folder.

Upload path
azureml://subscriptions/a017e82e-69e8-4283-b4a5-6ac9e216bb1c/resource...

Upload

- Upload files
- Upload folder

Upload list

File Types supported are delimited (i.e. csv, tsv), Parquet, JSON Lines, and plain text.

Information

What file types can I use?
Supported file types include: delimited (such as csv or tsv), Parquet, JSON Lines, and plain text.

Where are files uploaded?
Files will be uploaded to the selected datastore and made available in your workspace.

Back Next Cancel

Uploading dataset

ml.azure.com/visualinterface/authoring/Normal/7b85a56d-a987-4c11-bf8a-bbc761ffa3b7?wsid=/subscriptions/a017e82e-69e8-4283-b4a5-6ac9...

Azure AI | Machine Learning Studio

Create data asset

- ✓ Data type
- ✓ Data source
- ✓ Destination storage type
- ✓ File or folder selection
- 5 Settings
- 6 Schema
- 7 Review

Settings

These settings determine how the data is parsed. The initial settings are automatically detected; you can change them as needed to reparse the data.

File format: Delimited Delimiter: Comma Example: Field1,Field2,Field3 Encoding: UTF-8

Column headers: All files have same headers Skip rows: None

☐ Dataset contains multi-line data

Note: Processing tabular files with multi-line data is slower because multiple CPU cores cannot be used to ingest the data in parallel. Checking this option may result in slower processing times.

Data preview

CustomerID	Age	AnnualIncome	SpendingScore
1	46	371,045	99
2	43	45,194	24
3	48	111,465	59
4	61	null	21
5	39	191,670	43

Back Next Review Cancel

Validation

ml.azure.com/visualinterface/authoring/Normal/7b85a56d-a987-4c11-bf8a-bbc761ffa3b7wsid=/subscriptions/a017e82e-69e8-4283-b4a5-6ac9...
Azure AI | Machine Learning Studio

Create data asset

✓ Data type

✓ Data source

✓ Destination storage type

✓ File or folder selection

✓ Settings

✓ Schema

1 Review

Review

Review the settings for your data asset and make any changes as needed.

Data type

Name
customer_data

Description
--

Type
tabular

Data source

Type
Local

File selection

Upload path
azureml://subscriptions/a017e82e-69e8-4283-b4a5-6ac9e216bb1c/resourcegroups/aadil-rg/workspaces/aadil-aml/datastores/workspaceblobstore/paths/UI/2023-10-04_080536.UTC/customer_data.csv

Schema

CustomerID	Integer
Age	Decimal
AnnualIncome	Decimal
SpendingScore	Decimal

Back

Create

Cancel

Create Data Asset

ml.azure.com/visualinterface/authoring/Normal/7b85a56d-a987-4c11-bf8a-bbc761ffa3b7wsid=/subscriptions/a017e82e-69e8-4283-b4a5-6ac9...
Azure AI | Machine Learning Studio

All workspaces

Home

Model catalog

Authoring

Notebooks

Automated ML

Designer

Prompt flow

Assets

Data

Jobs

Components

Pipelines

Environments

Models

Endpoints

Unext > aadil-aml > Designer > Authoring

Search by name, tags and description

Tags: All Add filter

Data

Component

✓ Success: customer_data data asset create...

1 Last update...

You can find the prebuilt sample data under Component tab. [Click here](#)

customer_data

Version 1

Shellunext unextIDA27

10/4/2023

Pipeline-Created-on-10-04-2023

Save

Pipeline interface

customer_data

Shellunext unextIDA27

Description

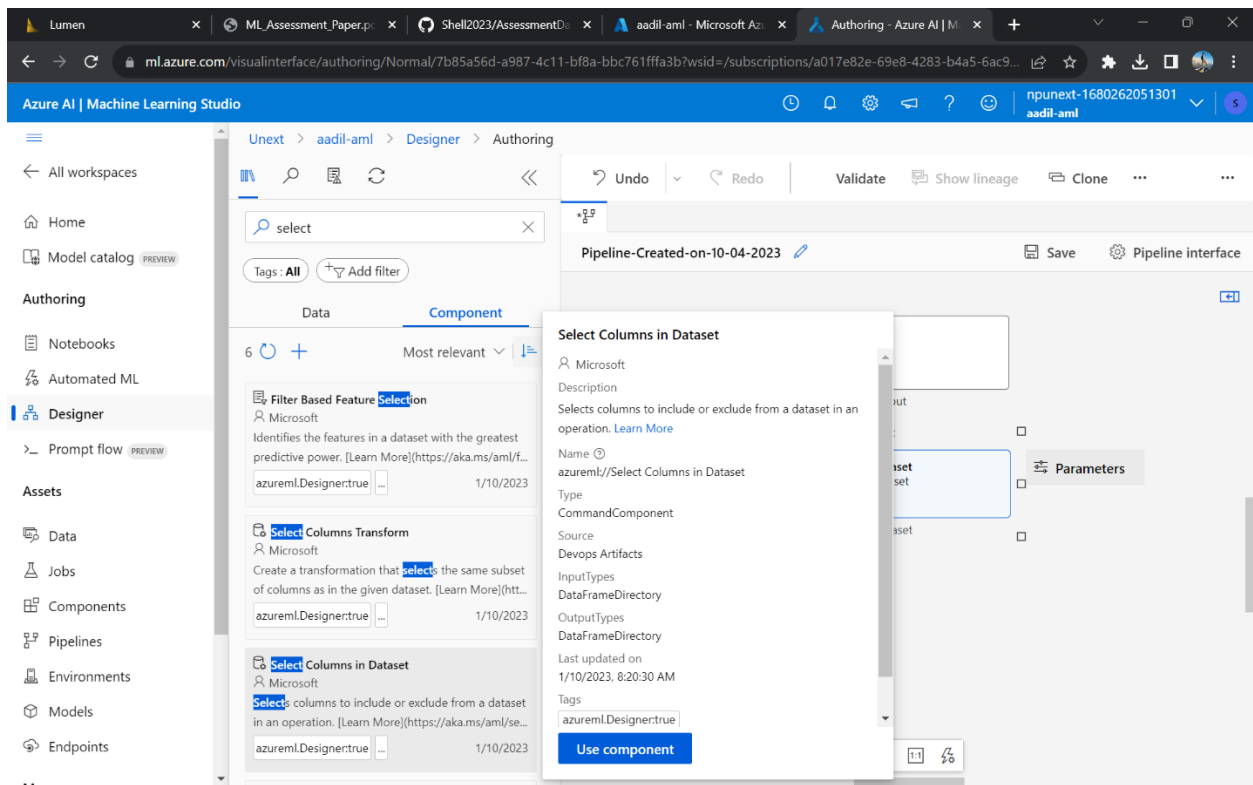
...

Use data

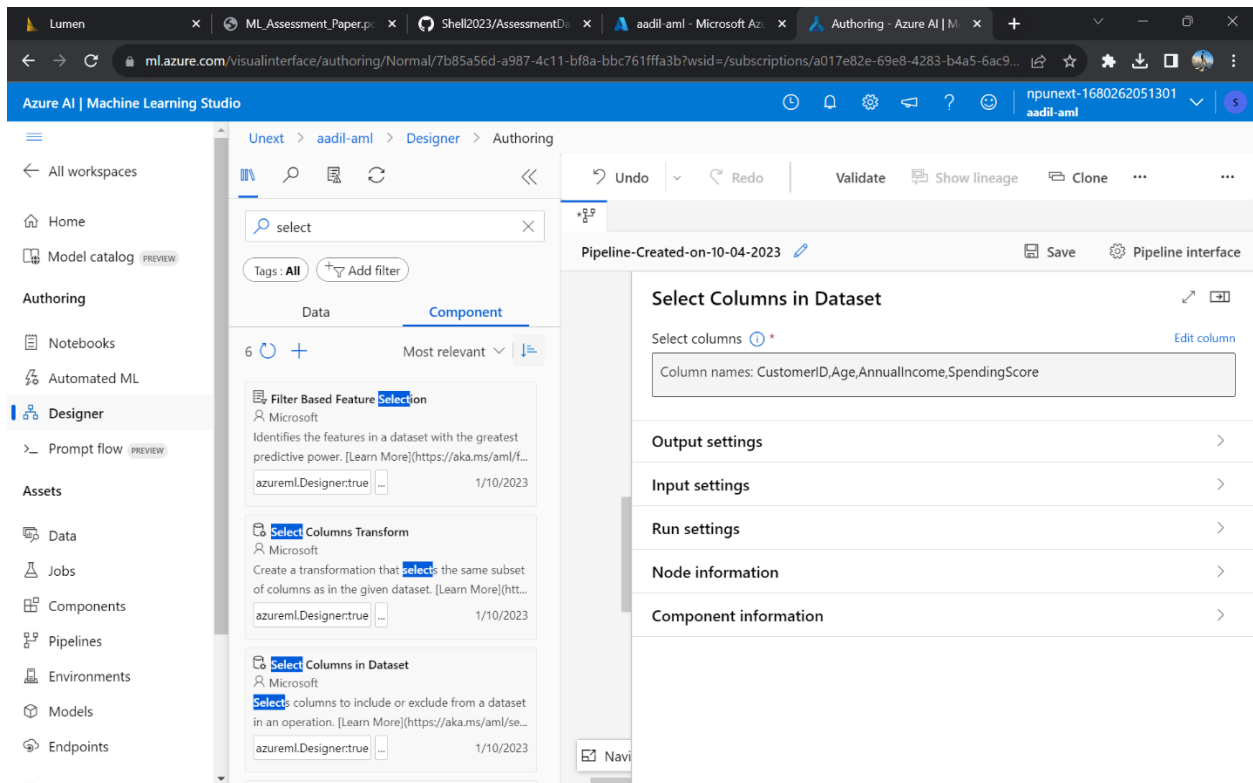
Navigator

100%

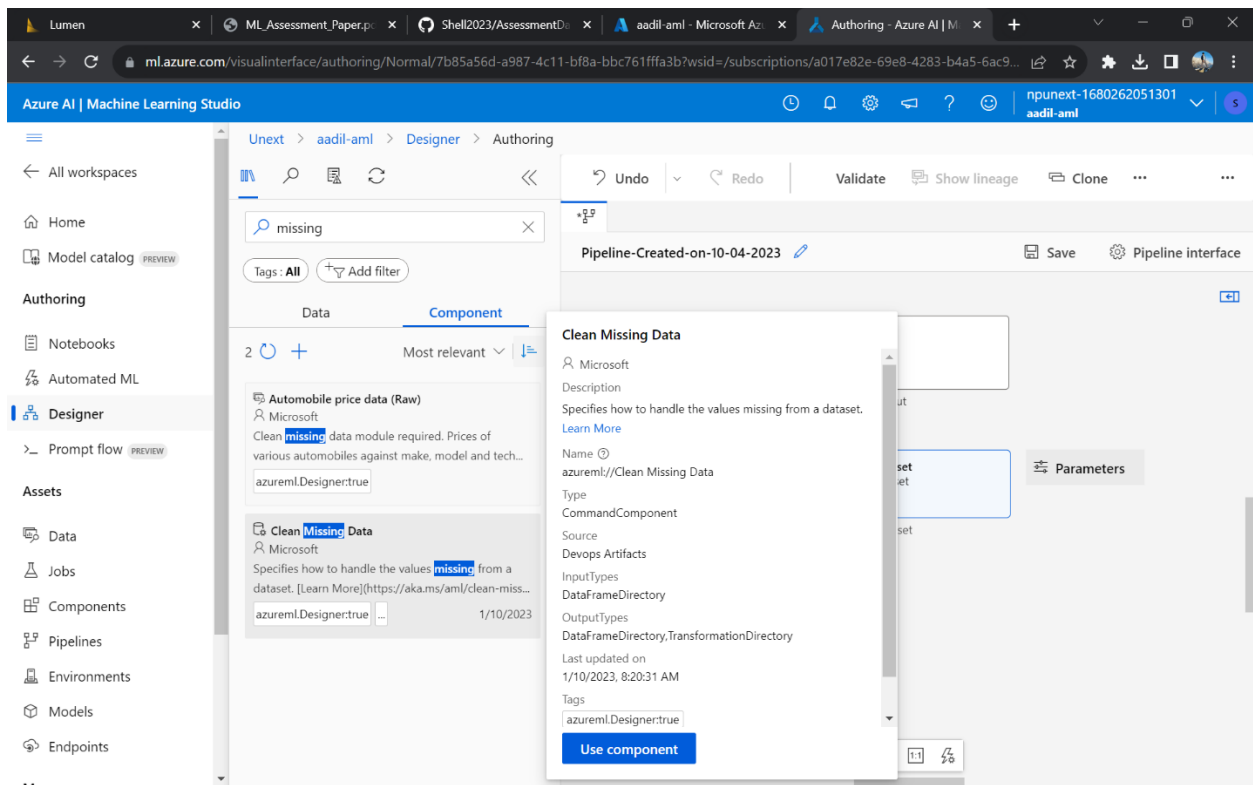
Select Dataset



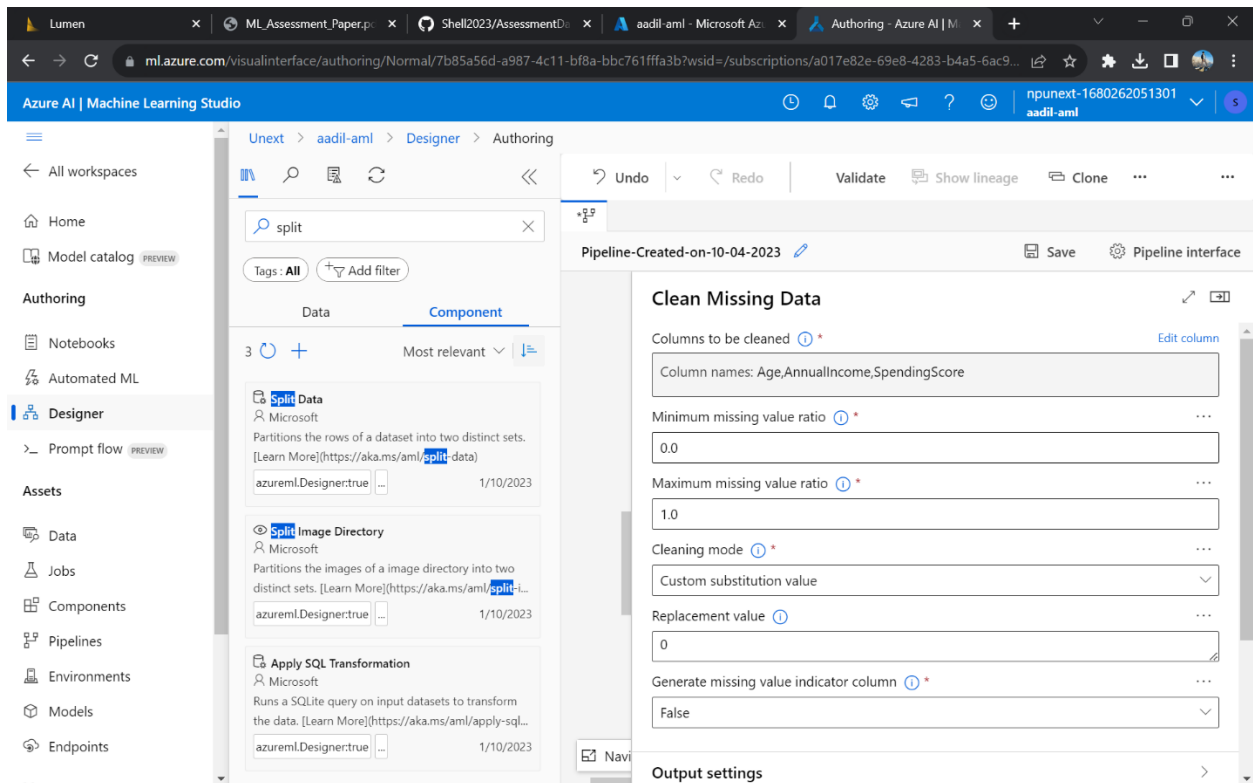
Select Column in Dataset



Select Columns

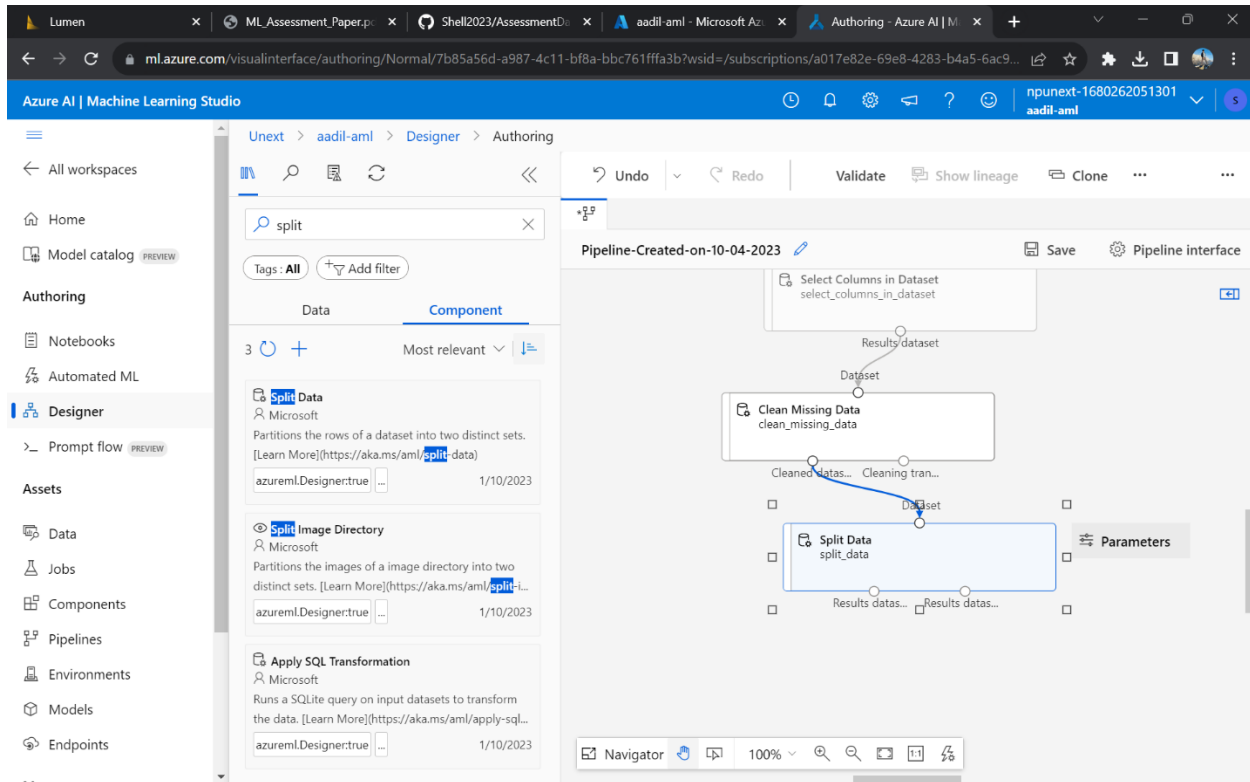


Clean Missing Data Component

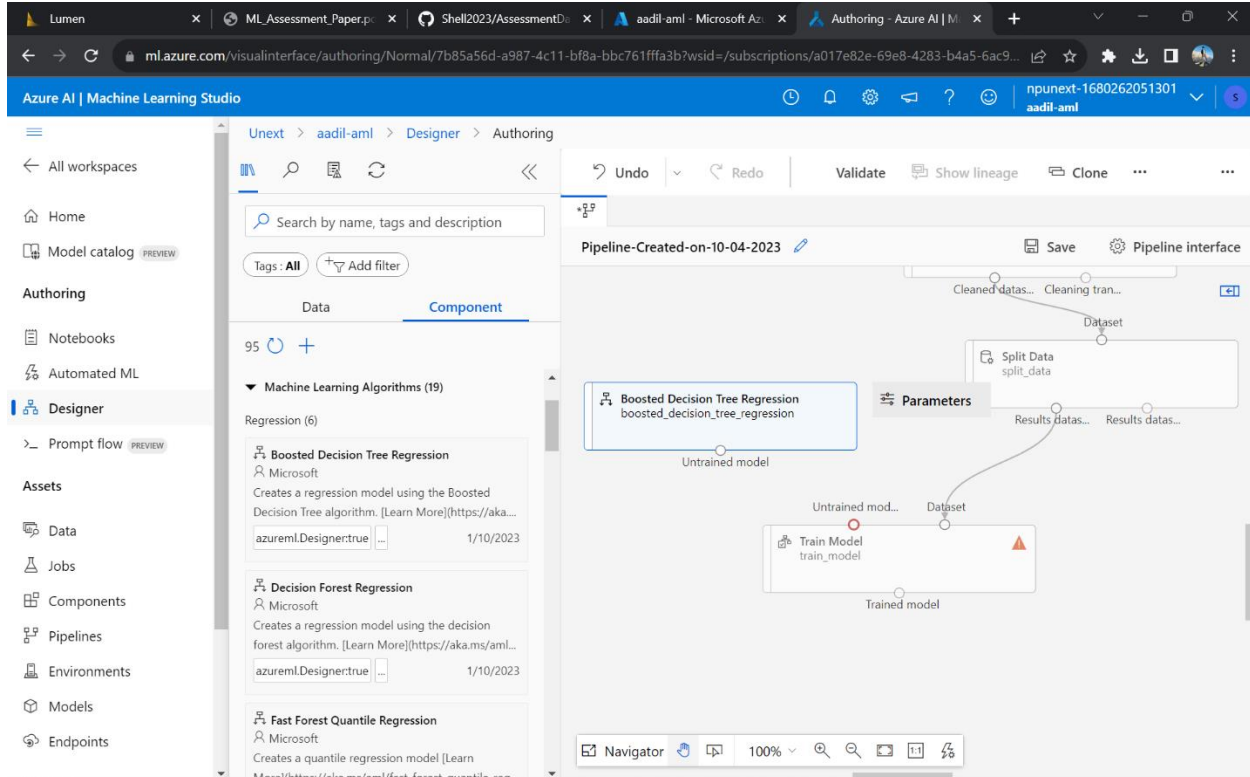


Configure Clean Data

2. Model Development

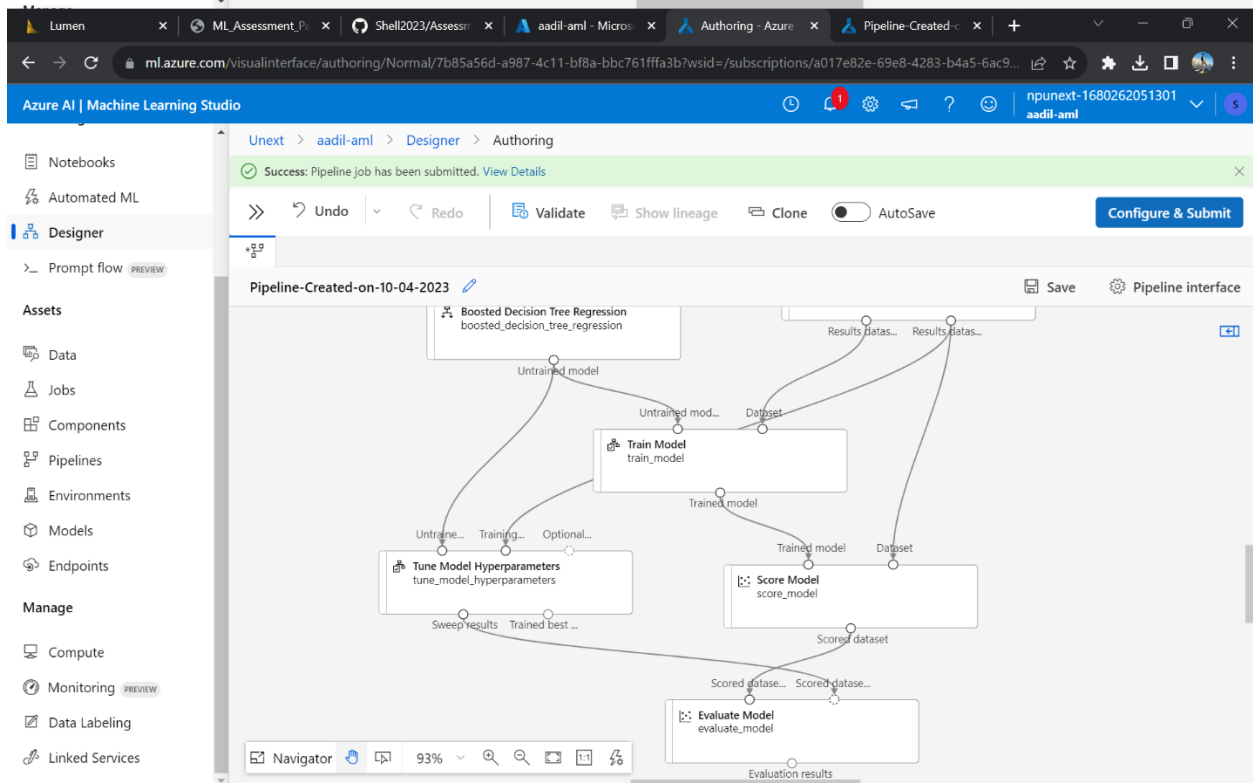
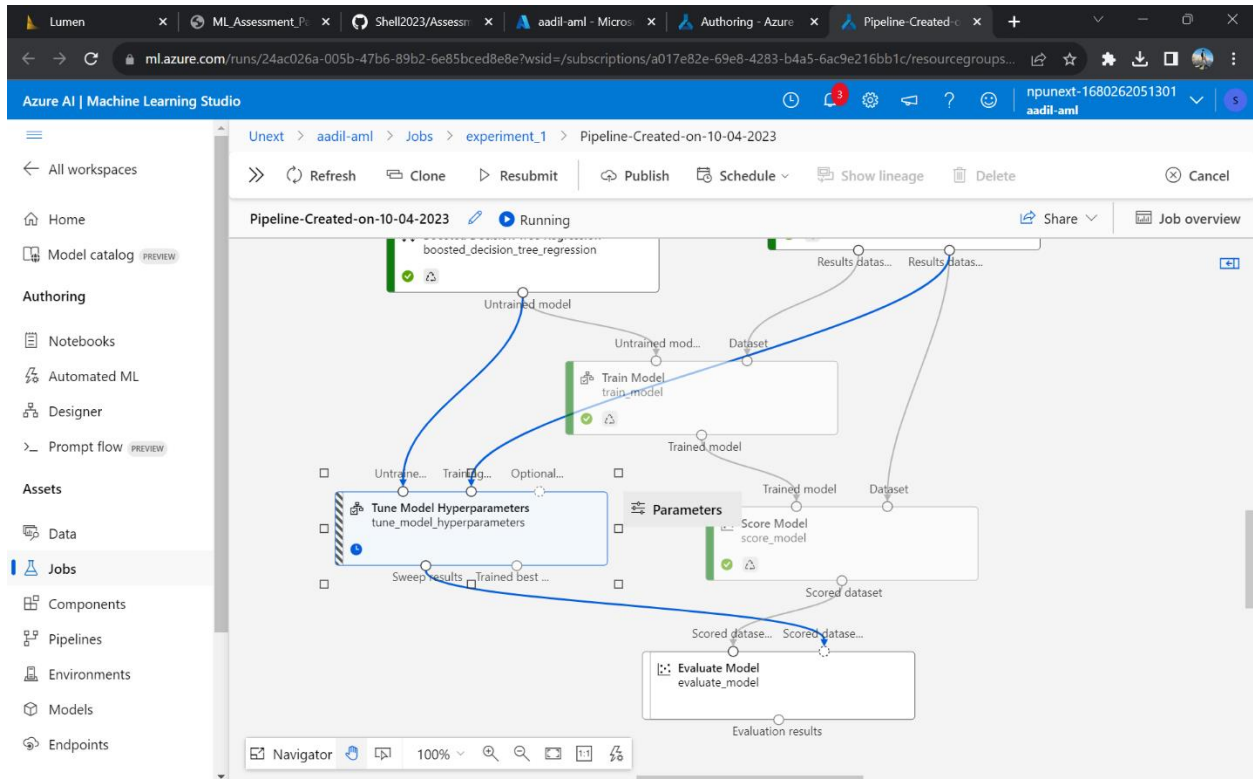


Split Data

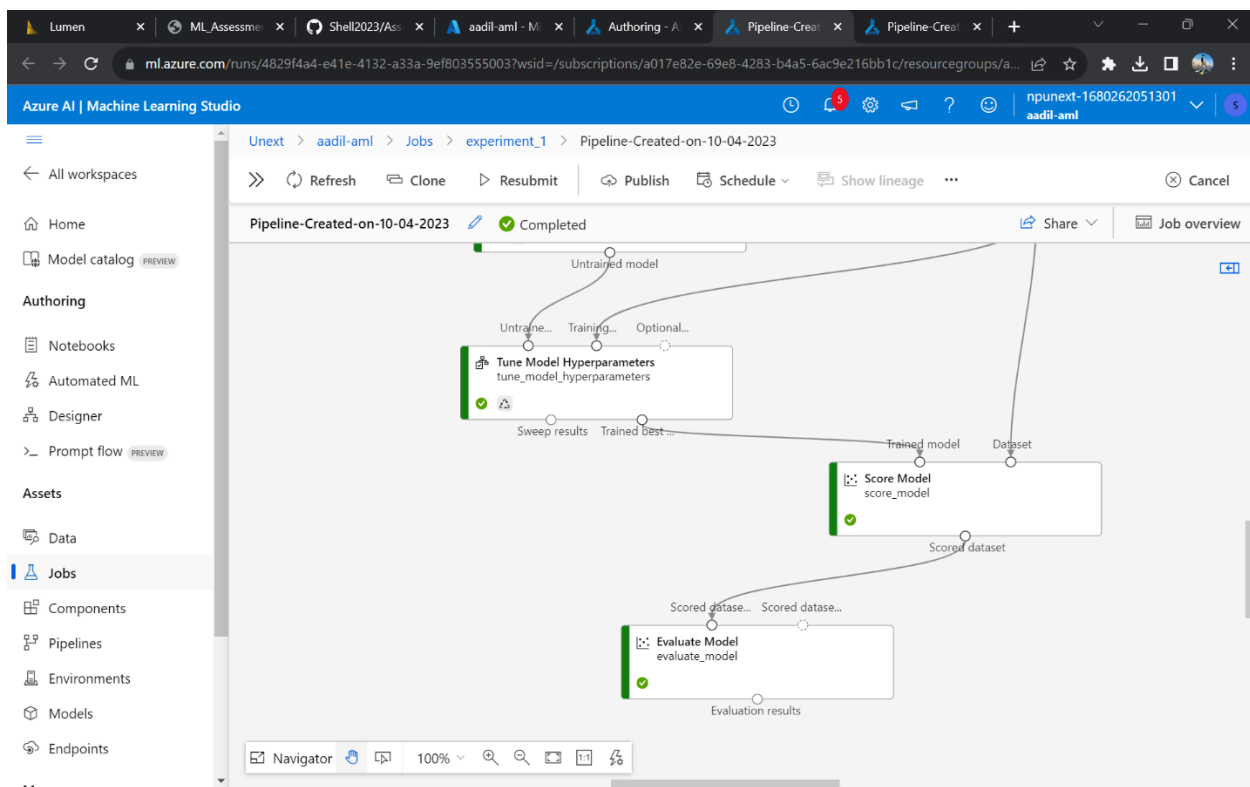


Machine Learning Algorithm

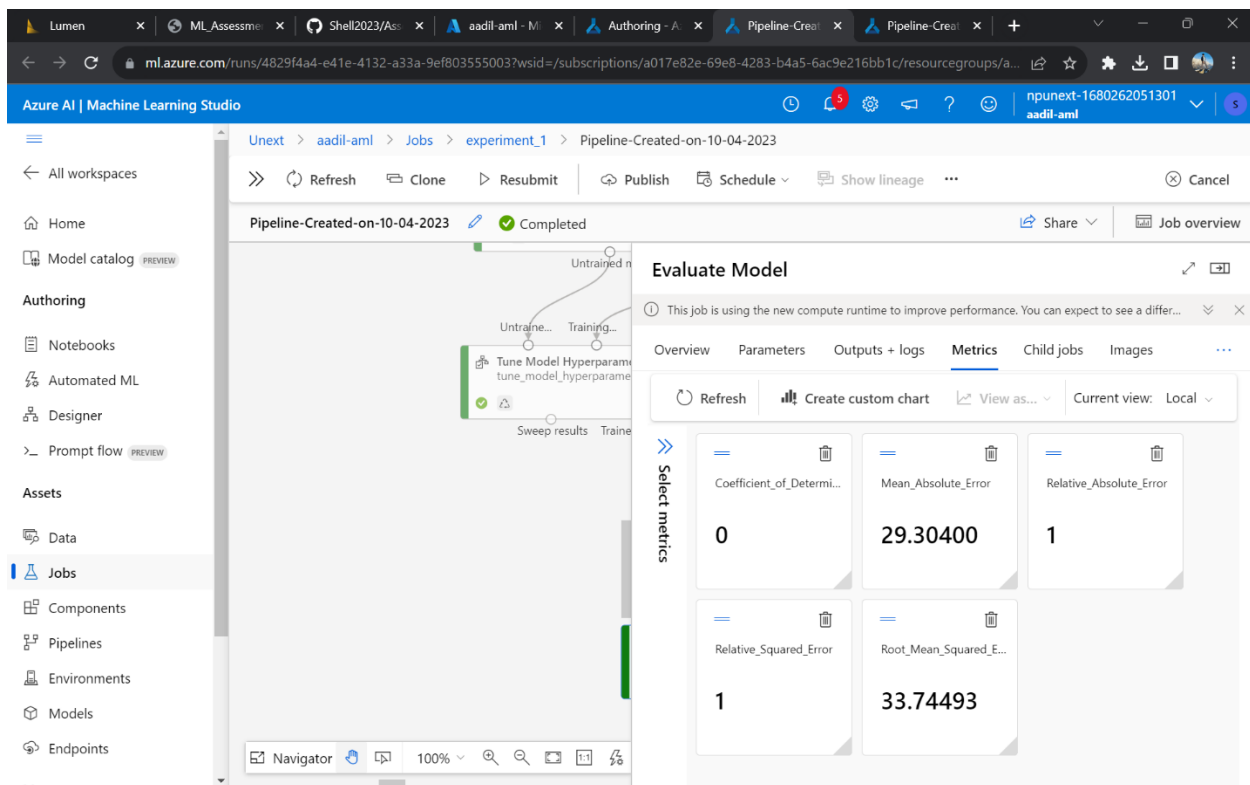
3. Hyperparameter Tuning



Train Model Hyperparameter, Score Model, Evaluate Model



Pipeline Ran Successfully



Evaluate Model Metrics

HANDS-ON COMPLETED

ASSESSMENT → QUESTION-ANSWERS

What are the key steps involved in preparing the dataset for training a machine learning model using Azure Machine Learning? Briefly explain each step.

Ans:

Preparing a dataset for training a machine learning model using Azure Machine Learning involves several key steps:

1. Data Collection and Ingestion:

- This initial step involves gathering the data from various sources, such as databases, files, or external APIs.

- In Azure Machine Learning, you can use Azure Data Factory or Azure Data Lake Storage to ingest and store your data.

2. Data Exploration:

- Before proceeding, it's essential to understand your data. This includes examining the data's structure, distribution, and quality.

- Use tools like Azure Databricks, Jupyter notebooks, or Azure Machine Learning Studio for data exploration.

3. Data Cleaning and Preprocessing:

- Clean the data by handling missing values, outliers, and duplicates.

- Perform preprocessing tasks like feature scaling, one-hot encoding, or text tokenization as needed.

4. Feature Engineering:

- Create new features or transform existing ones to improve the model's performance. This might involve domain-specific knowledge.

- Feature selection can also be part of this step to reduce dimensionality.

5. Data Splitting:

- Split the dataset into training, validation, and test sets. This helps evaluate the model's performance accurately.

- Azure Machine Learning provides built-in tools for data splitting and cross-validation.

6. Data Storage:

- Store the prepared dataset in a versioned and organized manner, typically in Azure Blob Storage or Azure Data Lake.
- This ensures traceability and reproducibility throughout the machine learning workflow.

7. Data Pipeline Creation:

- Build a data pipeline using Azure Machine Learning Dataflows or other orchestration tools to automate data preparation steps.

8. Data Validation and Monitoring:

- Continuously monitor the quality of your dataset, set up alerts for data drift, and retrain models when necessary.

1. Why is it important to split the dataset into training and testing sets when developing a machine learning model? How does this help in model evaluation?

Ans

Splitting the dataset into training and testing sets is a fundamental practice in machine learning for several important reasons:

1. Model Evaluation and Validation:

- By splitting the dataset, you can evaluate your machine learning model's performance on data it has never seen during training.
- This helps assess how well your model generalizes to unseen or new data, which is a crucial measure of its effectiveness.

2. Avoiding Overfitting:

- When a machine learning model is trained on a single dataset without a separate testing set, it can "overfit" to the training data.
- Overfitting means the model has learned to fit the training data's noise and anomalies rather than the underlying patterns. As a result, it performs poorly on new data.

- The testing set serves as an independent validation to detect overfitting. If the model performs well on the training set but poorly on the testing set, it's a sign of overfitting.

3. Hyperparameter Tuning:

- Splitting the dataset allows you to tune hyperparameters (e.g., learning rate, regularization strength) without using the testing data for this purpose.

- You can adjust hyperparameters on the training set and use the testing set to validate the chosen hyperparameter values.

4. Model Selection:

- When you have multiple candidate models, splitting the data helps you compare their performance on the same testing set.

- This allows you to select the best-performing model for your specific problem.

5. Bias and Variance Analysis:

- By evaluating a model on both training and testing sets, you can gain insights into its bias and variance.

- High training accuracy but low testing accuracy indicates high variance (overfitting), while low training and testing accuracy suggest high bias (underfitting).

6. Confidence in Results:

- Splitting the data and evaluating on a separate set provides a more accurate assessment of your model's true performance.

- It gives you confidence that the reported performance metrics are representative of how the model will perform in practice.

3. Describe a machine learning algorithm suitable for predicting customer purchasing behaviour in the given scenario. Explain why you chose this algorithm.

Ans

Predicting customer purchasing behavior is a common use case in marketing and e-commerce. One suitable machine learning algorithm for this scenario is the Random Forest – Boosted algorithm.

1. Ensemble Method: Random Forest is an ensemble learning method that combines multiple decision trees to make predictions. Each decision tree is trained on a subset of the data, and their predictions are averaged or voted upon to produce the final result. This ensemble approach often results in more robust and accurate predictions.

2. Handling Complex Relationships: Customer purchasing behavior can be influenced by a multitude of factors, some of which may have complex and non-linear relationships. Random Forests can capture these complex relationships as they can model both linear and non-linear patterns in the data.

3. Feature Importance: Random Forests provide a measure of feature importance, which is valuable in understanding which factors are most influential in predicting customer behavior. This information can be used for business insights and decision-making.

4. Reduced Overfitting: Random Forests are less prone to overfitting compared to single decision trees. By aggregating the results of multiple trees, they reduce the variance and improve generalization to new, unseen data.

5. Robustness to Outliers: Random Forests are robust to outliers and noisy data, which can be common in real-world customer behavior datasets.

6. Scalability: Random Forests can handle large datasets efficiently, making them suitable for scenarios with a significant amount of customer data.

4. What is hyperparameter tuning, and why is it important in machine learning? Explain a technique used for hyperparameter tuning and its benefits.

Ans

Hyperparameter tuning is the process of optimizing the hyperparameters of a machine learning algorithm to achieve the best possible model performance. Hyperparameters are parameters that are not learned from the data during training but are set before the training process begins. They control aspects of the learning process, such as the model's complexity, convergence speed, and generalization ability.

Hyperparameter tuning is essential in machine learning for several reasons:

1. Performance Improvement: Properly tuned hyperparameters can significantly improve a model's performance. They can make the difference between a model that underfits (too simple) or overfits (too complex) the data and one that generalizes well to new, unseen data.

2. Generalization: Optimized hyperparameters help the model generalize better to new data, reducing the risk of overfitting or underfitting. This leads to more reliable and robust models.

3. Efficiency: Hyperparameter tuning can lead to faster convergence during training, making the model learn more efficiently. This is particularly important when dealing with large datasets or complex models.

4. Model Selection: In cases where multiple machine learning algorithms are available, tuning hyperparameters can help select the best algorithm for a specific problem by comparing their performance under optimized settings.