# Azure ML Studio - Hands-On Assessment

## -Musaddiq Shariff

## Emp ID: 655023

1. Create Data set



2. Clean Missing Data

3. Split the dataset



4. Selecting appropriate model

5. Tune model Hyperparameters



6. Model Running

Assessment Questions:

1. 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: Steps involved in preparing dataset are:

1) Data collection: In this we use several types of sources like Gather relevant data from various sources, which may include databases, files, or external APIs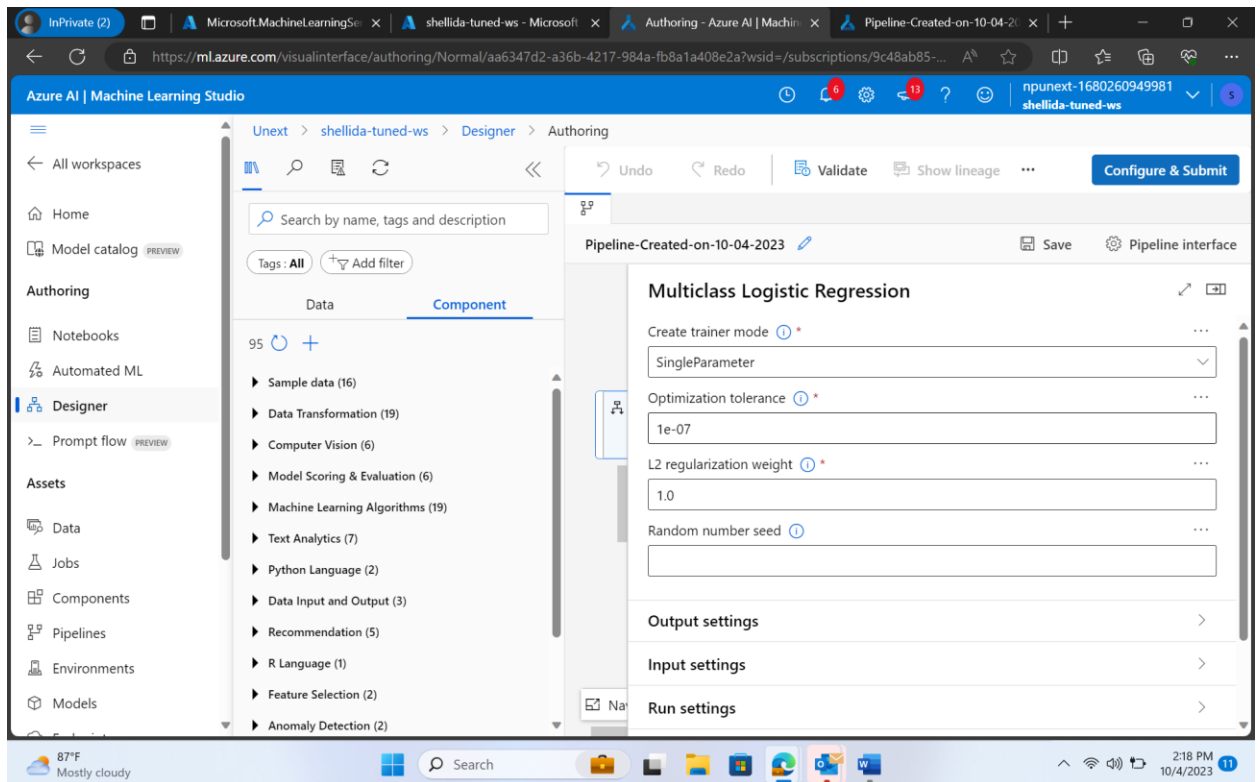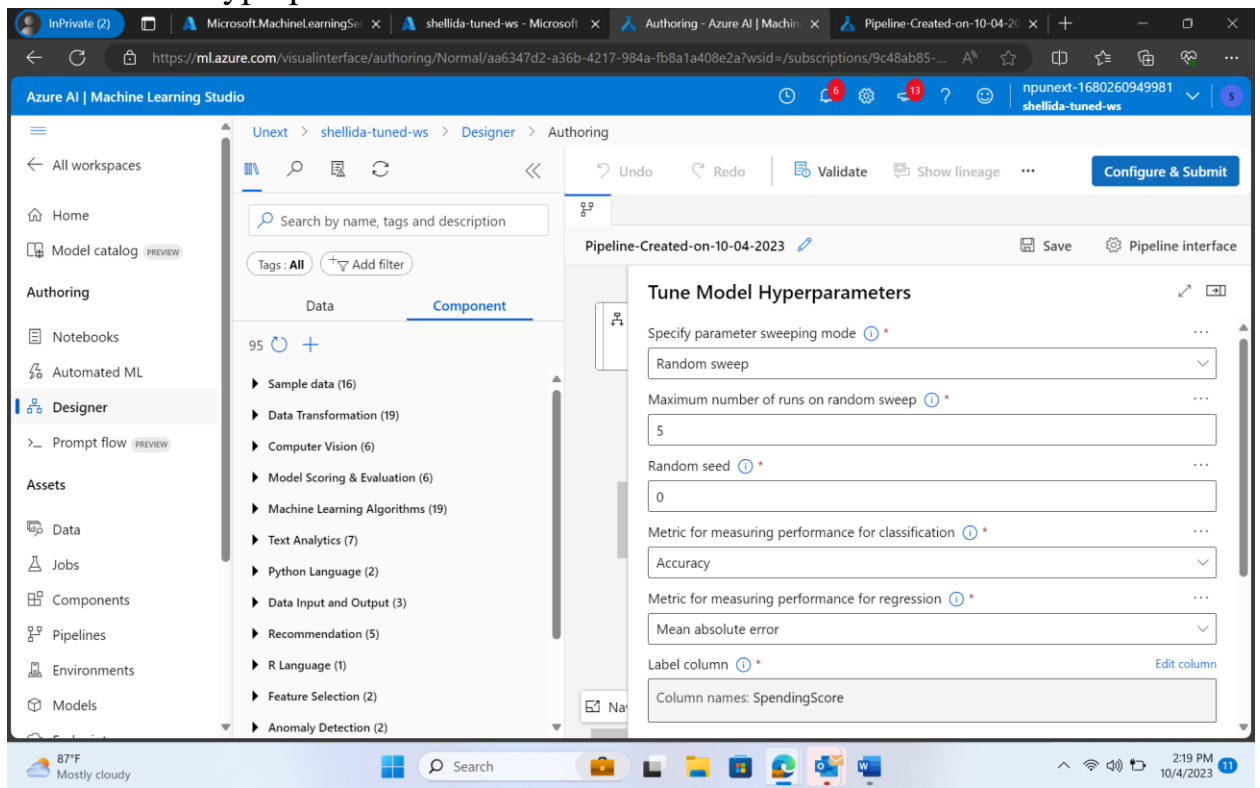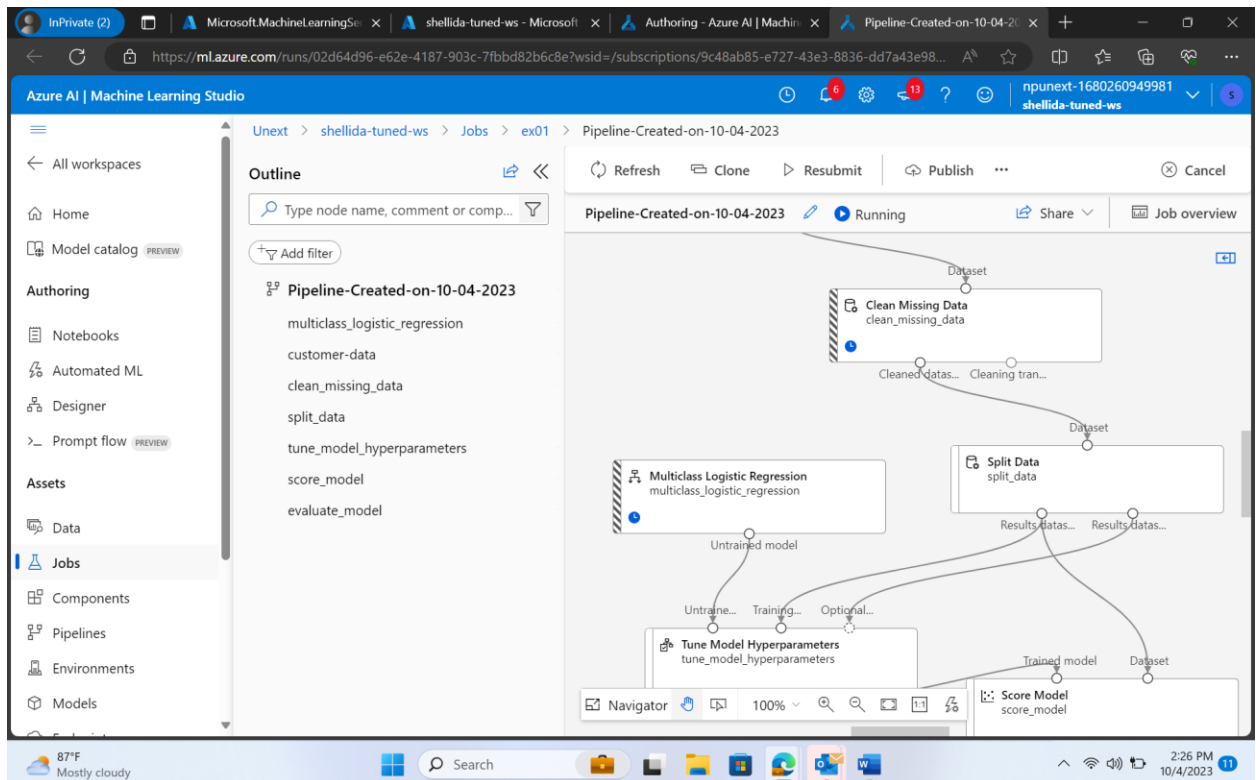. Azure Component: Utilize Azure Data Factory, Azure Databricks, or other services to ingest and collect data. We may also use Azure Data Factory, Azure Databricks, or other services to ingest and collect data.

2) Data Exploration: Analyze the dataset to understand its structure, identify missing values, outliers, and explore the distribution of features. For this we can use Use Azure Notebooks, Azure Databricks, or Jupyter notebooks in Azure Machine Learning to perform exploratory data analysis (EDA).

3) Data Cleaning: Handle missing values, outliers, and correct any inconsistencies in the dataset.t

4) Data Splitting: Divide the dataset into training and testing sets to assess the model's performance on unseen data.

2. 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. It is necessary to split the dataset into training and testing sets because the model has to be trained with a new dataset and we refer here it as training dataset and again for the validation purposed we use the testing dataset to check the accuracy of the model because the data available in the testing dataset will be unique for the model means its new for the model. The best results are obtained by using this technique.

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

Ans. I chose Random Forest Algorithm

The reasons are:

1) Random Forest is an ensemble learning method, meaning it combines the predictions of multiple individual models (decision trees) to improve overall performance
2) Non-Linearity and Complex Relationships
3) Handling Categorical Variables:
4) Robust to Overfitting:
5) Handling Imbalanced Datasets

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

Ans. In machine learning, a model's hyperparameters are configuration settings that are not learned from the data but are set prior to training. Hyperparameter tuning, also known as hyperparameter optimization, is the process of selecting the best set of hyperparameters for a machine learning model to improve its performance. The goal is to find the hyperparameters that result in the best generalization performance on unseen data.

Hyperparameters are configuration settings that are predefined before training and can significantly impact a model's ability to generalize to new data. Grid Search is a commonly employed technique for hyperparameter tuning, where a predefined grid of hyperparameter values is systematically searched. For each combination in the grid, the model is trained and evaluated, and the set of hyperparameters yielding the best performance is chosen. Grid Search offers a transparent and intuitive approach, ensuring a thorough exploration of the hyperparameter space. It is beneficial for improving model accuracy, preventing overfitting, and enhancing overall robustness, contributing to the reliability and effectiveness of machine learning models.