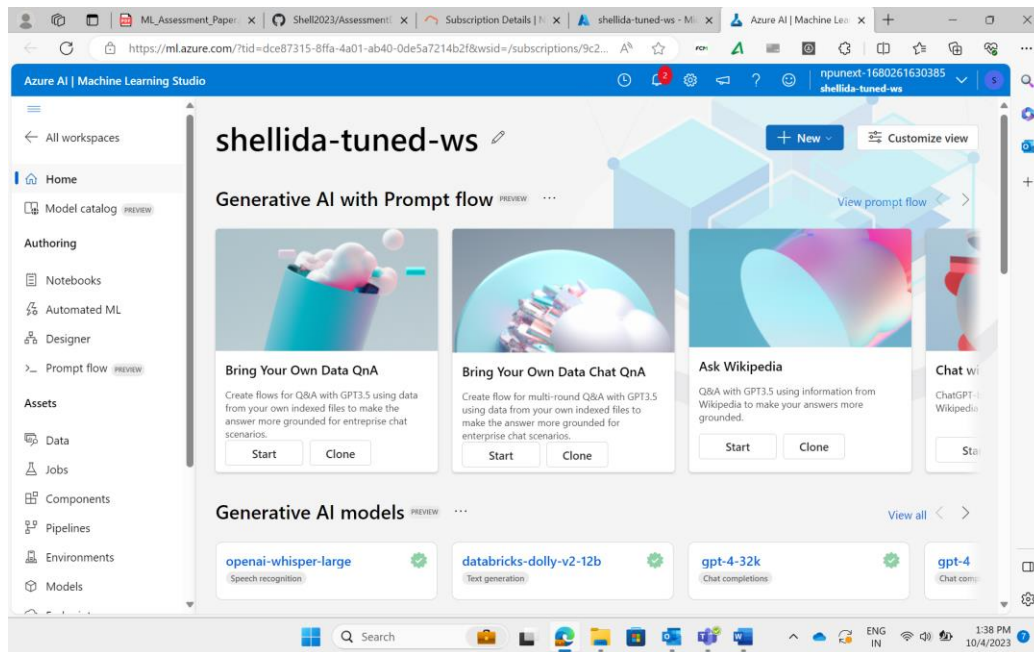


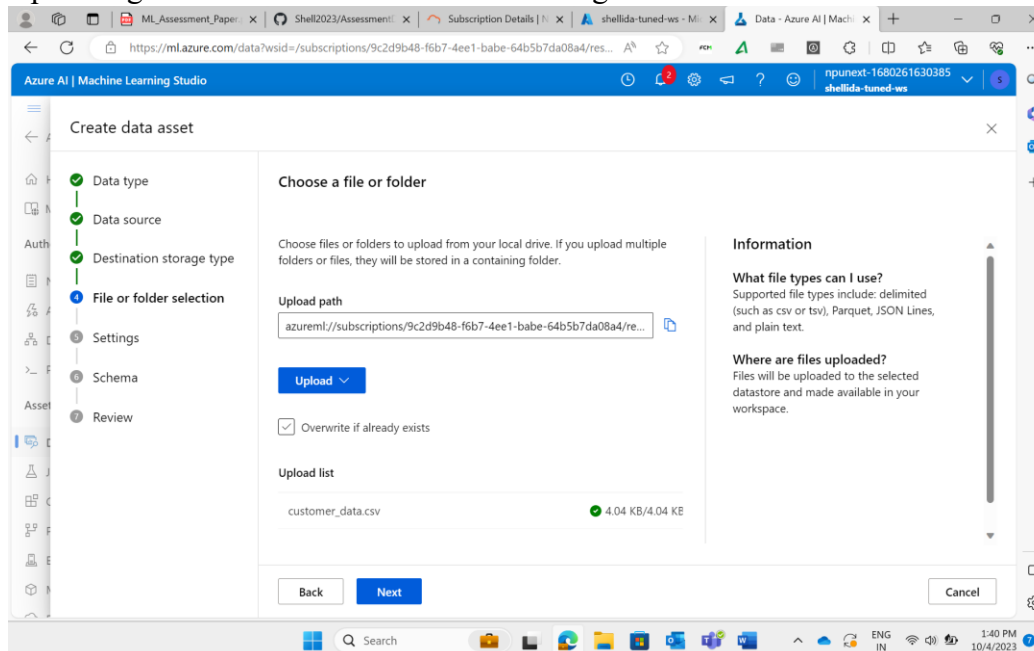
AZURE MACHINE LEARNING ASSESSMENT

- ANISHA GUPTA
EMP ID: 655517

Opened Azure ML Studio



Uploading customer data to Azure Blob Storage



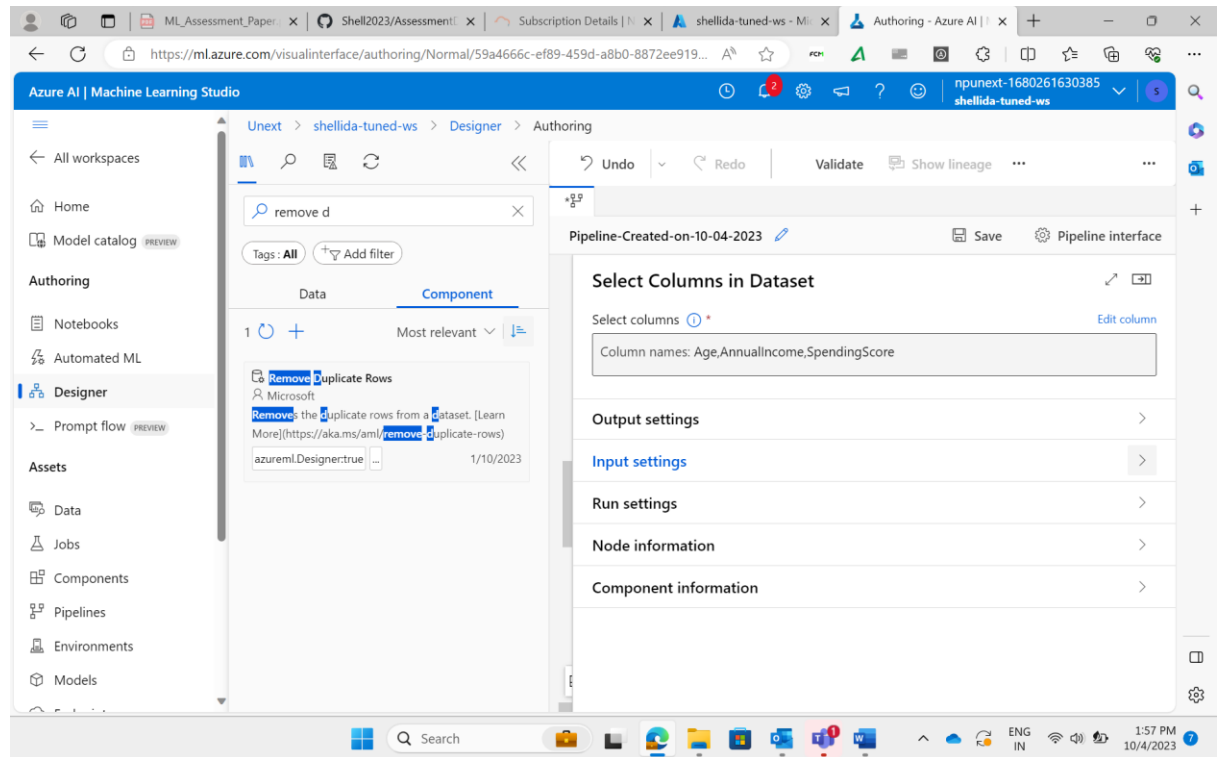
Uploaded data from Azure Blob Storage:

The screenshot shows the Azure AI Machine Learning Studio interface. The left sidebar contains navigation options: All workspaces, Home, Model catalog, Authoring, Notebooks, Automated ML, Designer (selected), Prompt flow, Assets, Data, Jobs, Components, Pipelines, Environments, and Models. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. It features a search bar, a 'Tags: All' filter, and a 'Data' tab. A table lists three data assets: 'customer' (Version 1, Shellunext unextIDA29, 10/4/2023), 'Orders' (Version 1, Shellunext unextIDA29, 10/4/2023), and 'customer-analysis' (Version 1, Shellunext unextIDA29, 10/4/2023). The pipeline canvas shows a 'customer' data source connected to a 'Data output' node. The bottom status bar indicates the time is 1:45 PM on 10/4/2023.

Removing duplicate rows:

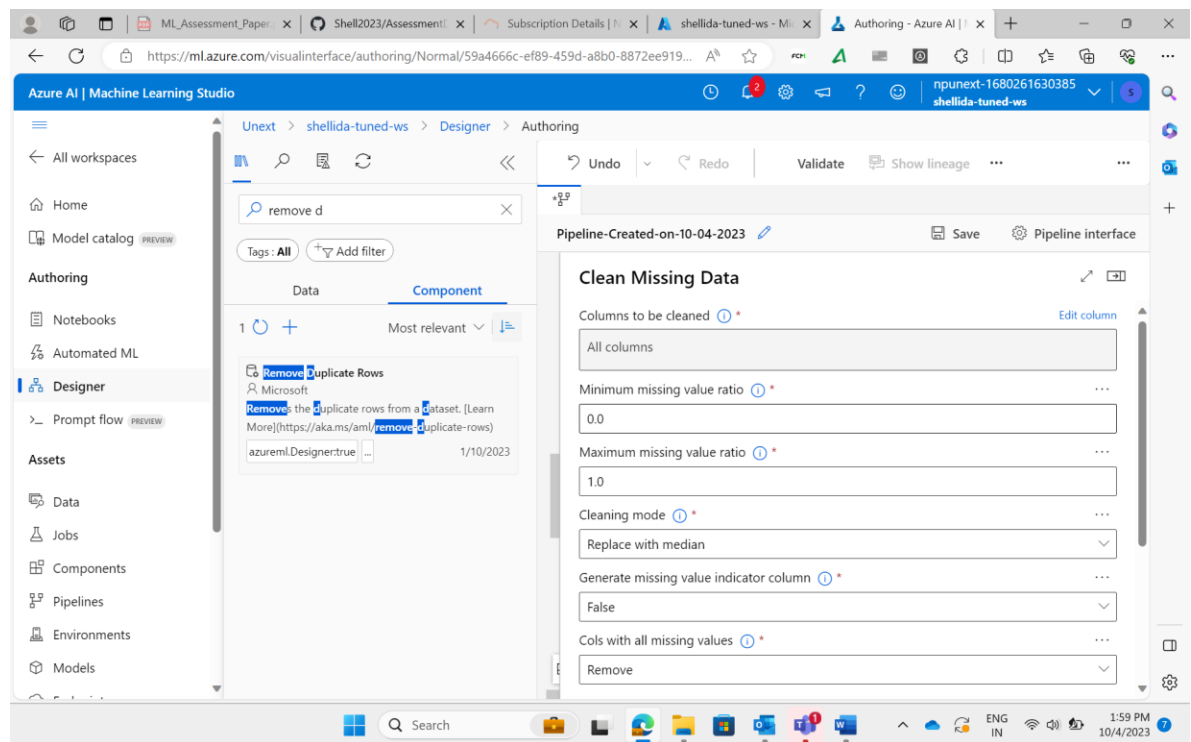
The screenshot shows the Azure AI Machine Learning Studio interface with the 'Remove Duplicate Rows' component selected. The left sidebar is the same as the previous screenshot. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. The search bar contains 'remove d'. The 'Component' tab is active, showing a table with one component: 'Remove Duplicate Rows' (Microsoft, 1/10/2023). The component description states: 'Removes the duplicate rows from a dataset. [Learn More](https://aka.ms/aml/remove-duplicate-rows)'. The right sidebar shows the configuration for the 'Remove Duplicate Rows' component. The 'Key column selection filter expression' is set to 'All columns'. The 'Retain first duplicate row' checkbox is checked. The bottom status bar indicates the time is 1:54 PM on 10/4/2023.

Selected all columns from the dataset:



Removed Id, and kept the rest of the columns

Cleaning the missing data:



Splitting Data into Training (70%) and Testing (30%)

The screenshot shows the Azure ML Studio interface. The left sidebar contains navigation options: All workspaces, Home, Model catalog, Authoring, Notebooks, Automated ML, Designer (selected), Prompt flow, Assets, Data, Jobs, Components, Pipelines, Environments, and Models. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. A search bar at the top of the component list shows 'split'. The 'Split Data' component is selected, and its configuration panel is open on the right. The configuration includes: Splitting mode (Split Rows), Fraction of rows in the first output dataset (0.7), Randomized split (True), Random seed (0), and Stratified split (False). The output settings and input settings sections are also visible. The bottom status bar shows the time as 2:02 PM on 10/4/2023.

Train Model: Linear Regression

The screenshot shows the Azure ML Studio interface. The left sidebar is the same as the previous screenshot. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. A search bar at the top of the component list shows 'hyperpara'. The 'Tune Model Hyperparameters' component is selected, and its configuration panel is open on the right. The configuration includes: Perform a parameter sweep on the model to determine the optimum parameter settings. The bottom status bar shows the time as 2:07 PM on 10/4/2023.

Without Hyperparameter Tuning: We directly train data

This screenshot shows the Azure Machine Learning Studio interface. The left sidebar contains navigation options: All workspaces, Home, Model catalog, Authoring (Notebooks, Automated ML, Designer), Assets (Data, Jobs, Components, Pipelines, Environments, Models), and Prompt flow. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. A search bar at the top of the component list shows 'trainin'. Below it, a list of components is displayed, including 'Train Clustering Model', 'Train Model', and 'Train PyTorch Model'. The 'Train Model' component is selected. The right pane shows a pipeline diagram titled 'Pipeline-Created-on-10-04-2023'. The pipeline consists of several steps: 'Data Input Dataset', 'Remove Duplicate Rows', 'Select Columns in Dataset', 'Clean Missing Data', 'Split Data', 'Train Model', 'Evaluate Model', 'Score Model', and 'Scored dataset'. The 'Train Model' step is highlighted with a red triangle. The bottom status bar shows the time as 2:23 PM on 10/4/2023.

This screenshot shows the configuration for the 'Train Model' component in the Azure Machine Learning Studio. The left sidebar is the same as the previous screenshot. The main workspace is titled 'Unext > shellida-tuned-ws > Designer > Authoring'. The search bar at the top of the component list shows 'trainin'. The 'Train Model' component is selected. The right pane shows the configuration for the 'Train Model' component. The 'Label column' is set to 'SpendingScore'. The 'Model explanations' are set to 'False'. The 'Output settings', 'Input settings', 'Run settings', 'Node information', and 'Component information' sections are visible. The bottom status bar shows the time as 2:26 PM on 10/4/2023.

Hyperparameter Tuning of Model: Using Mean Absolute Error for regression

The screenshot displays the Azure ML Studio interface. The left sidebar shows the 'Designer' tab with a search bar containing 'hyperpara'. The main workspace shows the 'Tune Model Hyperparameters' component. The configuration panel on the right is titled 'Tune Model Hyperparameters' and includes the following settings:

- Specify parameter sweeping mode: Random sweep
- Maximum number of runs on random sweep: 5
- Random seed: 123
- Metric for measuring performance for classification: Accuracy
- Metric for measuring performance for regression: Mean absolute error
- Label column: Column names: SpendingScore

Score Model:

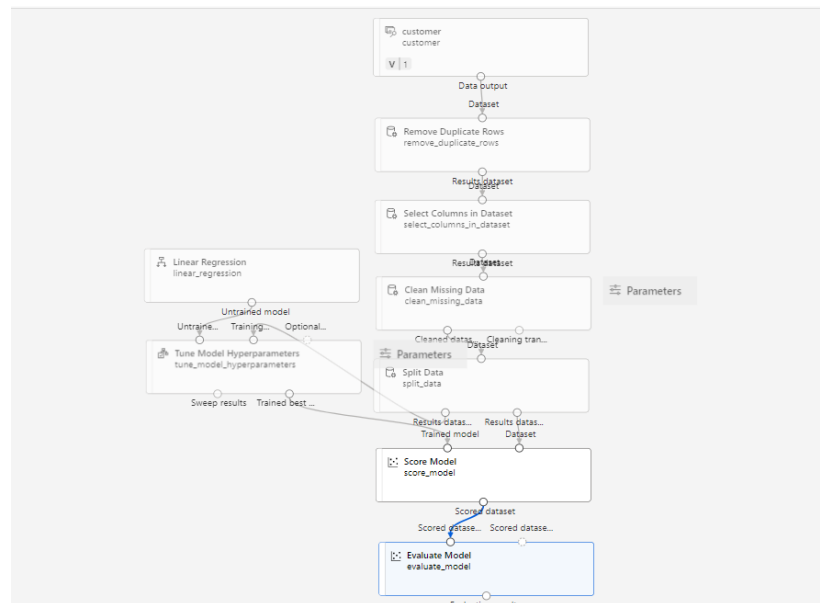
The screenshot displays the Azure ML Studio interface. The left sidebar shows the 'Designer' tab with a search bar containing 'score'. The main workspace shows the 'Score Model' component. The configuration panel on the right is titled 'Score Model' and includes the following settings:

- Append score columns to output: True
- Output settings: >
- Input settings: >
- Run settings: >
- Node information: >
- Component information: >

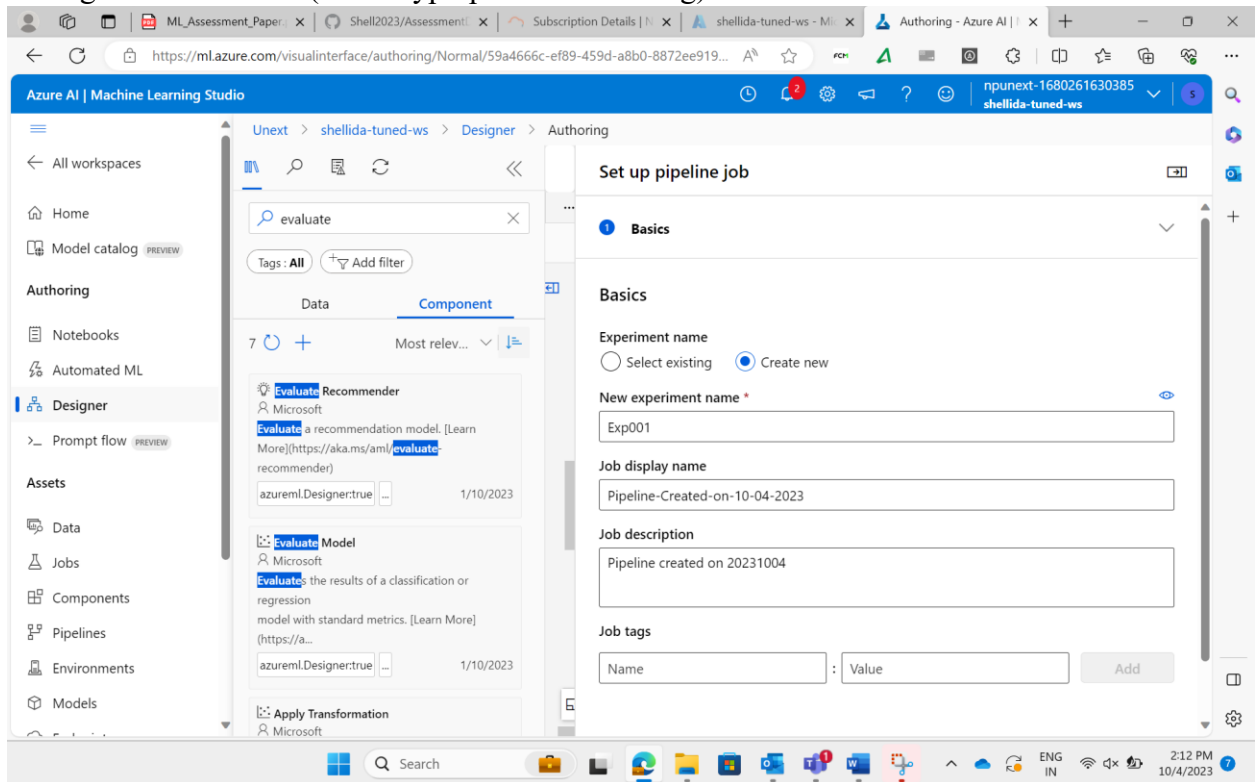
Evaluate Model:

The screenshot displays the Azure AI Machine Learning Studio interface. The top navigation bar shows the workspace name 'shellida-tuned-ws' and the 'Designer' tab. The left sidebar contains a navigation menu with options like 'All workspaces', 'Home', 'Model catalog', 'Authoring', 'Notebooks', 'Automated ML', 'Designer', 'Prompt flow', 'Assets', 'Data', 'Jobs', 'Components', 'Pipelines', 'Environments', and 'Models'. The 'Designer' tab is active, showing a search bar with 'evaluate' and a list of components. The 'Evaluate Model' component is selected, showing its description: 'Evaluates the results of a classification or regression model with standard metrics. [Learn More](https://aka.ms/aml/evaluate-model)'. The main canvas displays a pipeline titled 'Pipeline-Created-on-10-04-2023'. The pipeline steps include: 'Linear Regression linear_regression', 'Untrained model', 'Tune Model Hyperparameters tune_model_hyperparameters', 'Sweep results', 'Trained best', 'Clean Missing Data clean_missing_data', 'Split Data split_data', 'Score Model score_model', and 'Evaluate Model evaluate_model'. The 'Evaluate Model' component is highlighted in blue. The bottom status bar shows the search bar, a zoom level of 86%, and the system clock indicating 2:09 PM on 10/4/2023.

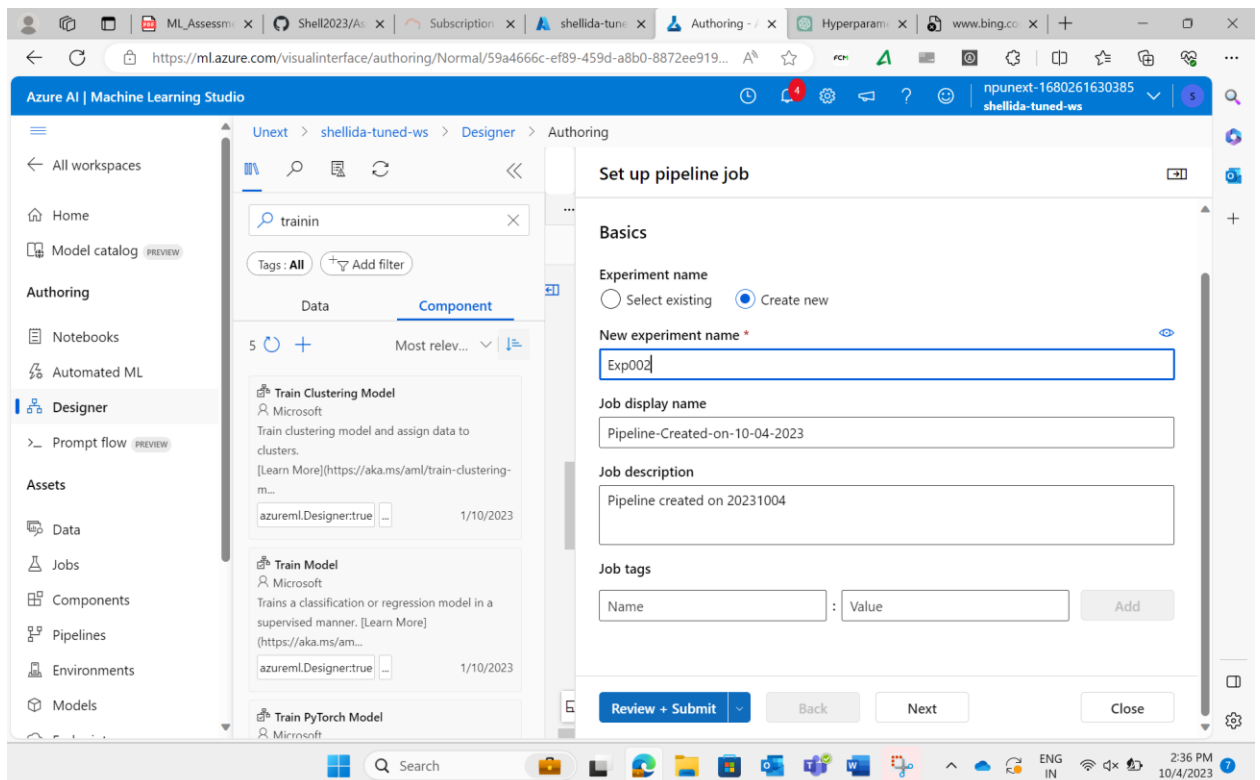
The entire pipeline:



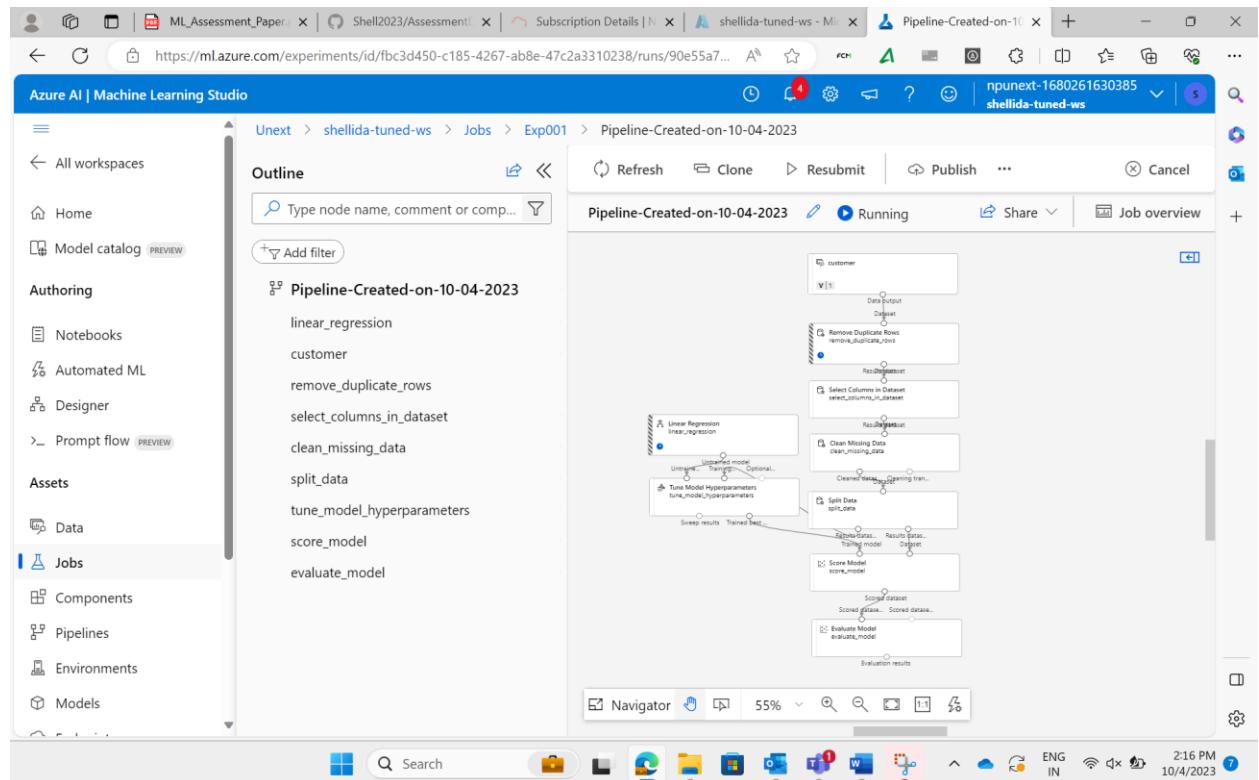
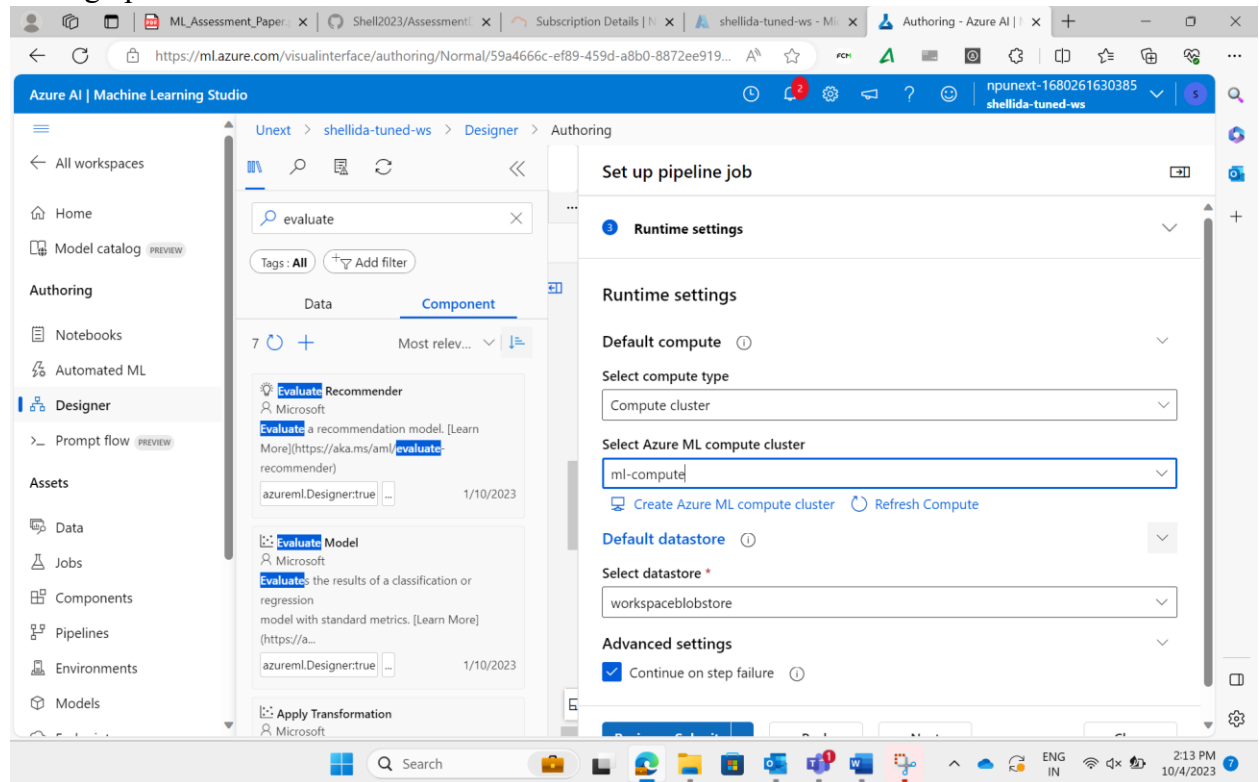
Configure and Submit: (With Hyperparameter Tuning)



Without Hyperparameter Tuning



Setting up the runtime environment:



Similarly Exp002 is also in running state

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.

Data Ingestion:

Ingest or import the data into Azure ML

Data Exploration and Analysis:

Perform exploratory data analysis (EDA) to gain a better understanding of your dataset. Explore summary statistics, data distributions, missing values, and potential outliers. Visualization tools in Azure ML can help with this step.

Data Cleaning:

Handle missing values, duplicate records, and outliers. Impute missing values, remove duplicates, and apply appropriate techniques to address outliers. Data cleaning is essential to ensure the quality of the dataset.

Feature Engineering:

Create, transform, or select features (variables) that are relevant to your machine learning task. Feature engineering can involve techniques such as one-hot encoding, scaling, or creating new features from existing ones.

Data Splitting:

Split the dataset into multiple subsets for training, validation, and testing purposes. Common splits include a training set for model training, a validation set for hyperparameter tuning and model selection, and a test set for evaluating final model performance.

Data Preprocessing:

Apply preprocessing techniques such as normalization, standardization, or feature scaling to make the data suitable for machine learning algorithms. Preprocessing ensures that all features have the same scale and distribution.

Data Transformation:

Perform any necessary data transformations, including encoding categorical variables, handling time series data, or applying dimensionality reduction techniques like Principal Component Analysis (PCA).

Data Validation:

Validate the prepared dataset to ensure that it meets the requirements of your machine learning problem. Check for any data leakage or issues that might affect model performance.

Data Storage:

Store the cleaned and preprocessed dataset in Azure ML's data storage solutions for easy access during model training and deployment.

Data Feeding to Model:

Once your dataset is prepared, you can feed it into the machine learning model training process in Azure ML, using the appropriate model training algorithm and techniques. Evaluate the result from the model using the scores.

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?

Splitting the dataset into training and testing sets is a fundamental step in developing a machine learning model, and it serves several important purposes in the model development process. Here's why it's crucial and how it helps in model evaluation:

Assessing Generalization Performance: The primary goal of a machine learning model is to make accurate predictions on new, unseen data. By splitting the dataset into two distinct sets, you can evaluate how well your model generalizes to unseen data

Avoiding Overfitting: Overfitting occurs when a model learns to memorize the training data rather than capturing its underlying patterns. The testing set acts as a check against overfitting by evaluating the model's performance on data it has not seen during training.

Model Selection: Splitting the data allows you to compare different machine learning models or hyperparameter settings. This helps you choose the best-performing model or configuration.

Performance Evaluation: By having a separate testing set, you can calculate various performance metrics, such as accuracy, precision, recall, F1 score, or mean squared error, to quantitatively assess how well your model performs on real-world data. These metrics provide valuable insights into the model's strengths and weaknesses.

Detecting Data Leakage: Data leakage is a common issue where information from the testing data unintentionally influences the model's training. By keeping the testing set

separate from the training set, you can prevent data leakage and ensure that the model's performance evaluation is unbiased.

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

I have used Linear Regression.

Linear regressions can be used in business to evaluate trends and make estimates or forecasts. For example, if a company's sales have increased steadily every month for the past few years, by conducting a linear analysis on the sales data with monthly sales, the company could forecast sales in future months.

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

Hyperparameter tuning is the process of finding the best set of hyperparameters (configurations not learned from data) for a machine learning algorithm to optimize its performance. It's important because the right hyperparameters can significantly impact a model's accuracy and generalization.

One technique for hyperparameter tuning is "Grid Search," where you specify a range of values for each hyperparameter, and the algorithm systematically evaluates all combinations to find the best one. Benefits include simplicity and thorough exploration of hyperparameter space.