Megha Sharma

Megha.sharma2@shell.com

Hands On assignment

1 a.

Data type



Data source

Blob storage destination



File selection

Review data asset

**Screenshot 1 — Settings page:**

Azure AI | Machine Learning Studio

npunext-1673505489600
MSAML

Create data asset

- Data type
- Data source
- Destination storage type
- File or folder selection
- Settings
- Schema
- 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 |
| 6 | 41 | 120,433 | 52 |
| 7 | 18 | 52,885 | null |
| 8 | 63 | 108,250 | 95 |

Back    Next    Review    Cancel

1:46 PM 10/4/2023

---

**Screenshot 2 — Schema page:**

Azure AI | Machine Learning Studio

npunext-1673505489600
MSAML

Create data asset

- Data type
- Data source
- Destination storage type
- File or folder selection
- Settings
- Schema
- Review

## Schema

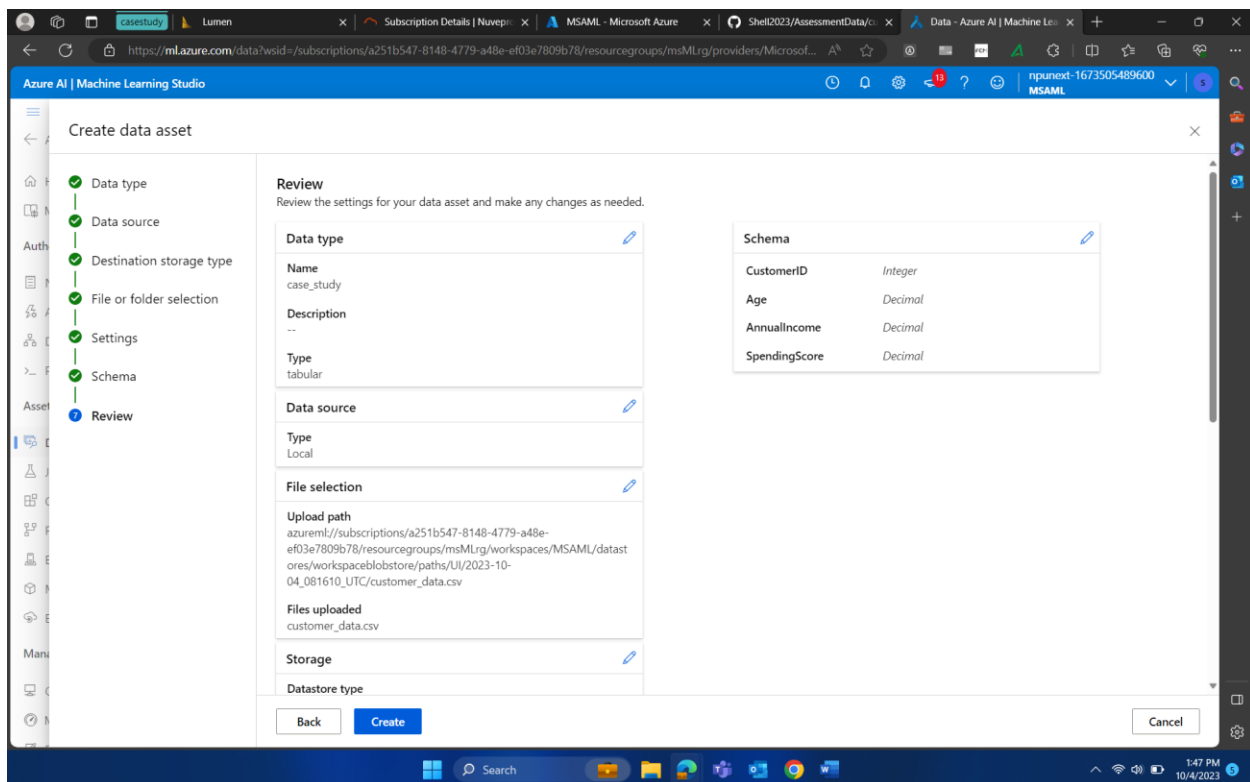Column types are auto-detected based on the initial subset of the data and can be updated here. Values not aligning with the specified column type will fail conversion and would be either null-filled or replaced with error value. Any conversions preview errors are non-blocking and you can proceed.

🔍 Search column name

| Include | Column name | Type | Example values | Date format ⓘ | Properties ⓘ |
|---|---|---|---|---|---|
| ⊘ | Path | String | | Not applicable to selected ... | Not applicable t... |
| ⊙ | CustomerID | Integer | 1, 2, 3 | Not applicable to selected ... | Not applicable t... |
| ⊙ | Age | Decimal (dot '.') | 46, 43, 48 | Not applicable to selected ... | Not applicable t... |
| ⊙ | AnnualIncome | Decimal (dot '.') | 371045, 45194, 111465 | Not applicable to selected ... | Not applicable t... |
| ⊙ | SpendingScore | Decimal (dot '.') | 99, 24, 59 | Not applicable to selected ... | Not applicable t... |

Back    Next    Cancel

1:46 PM 10/4/2023

1 b.



Select columns

Clean missing data

2 a.



Linear regression model

2 b.



Split data

3 a.



Hyper parameter tuning

Assessment Questions.

1.
   a. Data Collection: Gather and acquire the data needed for training from various sources, including databases, files, web services, or external datasets.
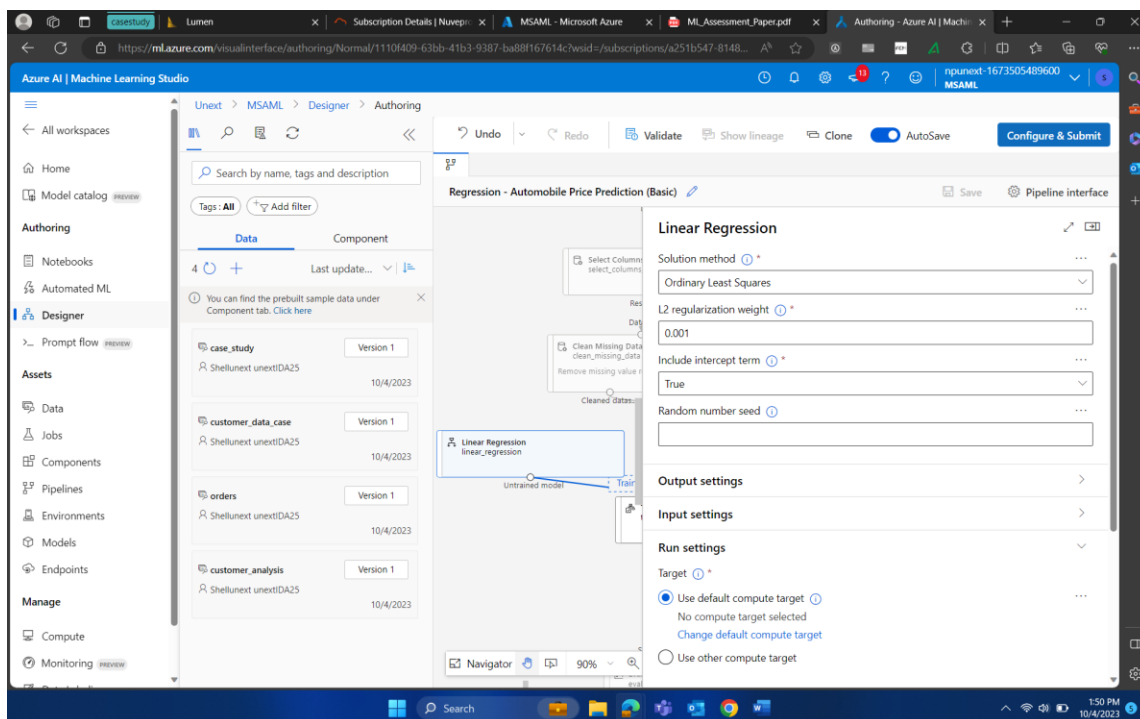   b. Data Cleaning: Identify and address issues like missing values, outliers, and inconsistencies within the dataset. This often involves actions such as filling in missing data, eliminating outlier values, or rectifying data errors.
   c. Data Transformation: Modify the data by actions like converting categorical variables into numerical forms, scaling numerical attributes, and encoding text or images.
   d. Feature Engineering: Enhance the dataset by creating new features or altering existing ones to improve model performance. This includes selecting relevant features and transforming them meaningfully.
   e. Data Splitting: Divide the dataset into distinct portions—typically training, validation, and test sets. The training set is used to train the model, the validation set aids in hyperparameter tuning, and the test set assesses model performance.
   f. Data Uploading: Upload the cleaned and processed dataset to Azure Machine Learning Workspace or storage, making it accessible for model training and experimentation.
   g. It is important to split dataset into training and validation because The training set is used to train the model, the validation set aids in hyperparameter tuning, and the test set assesses model performance.


2. Splitting a dataset into training and testing sets is essential in machine learning for various reasons:
   a. Model Assessment: It allows evaluating how well a model performs on new, unseen data.
   b. Overfitting Detection: Helps identify if a model is too closely tailored to the training data, indicating overfitting.
   c. Hyperparameter Tuning: Enables systematic tuning of model settings to optimize generalization.
   d. Performance Metrics: Provides data for calculating performance metrics like accuracy, precision, and recall.
   e.  Bias and Variance Analysis: Aids in assessing bias and variance to strike a balance.
   f. Model Selection: Facilitates choosing the best model among candidates.
   g. Validation of Assumptions: Ensures model behavior aligns with expectations and real-world scenarios.
   h. Real-world Simulation: Mimics the model's performance in production environments.


3. We can see that there is a linear trend or direct relation between the income and spending score and age, we can try using linear regression model.
   a. Interpretability: Linear Regression provides clear interpretability of coefficients. This means we can easily understand the direction and magnitude of the impact of each predictor variable on the purchasing behavior. This is important when we want to explain why certain factors influence buying decisions.

b. Simplicity: Linear Regression is a simple and straightforward algorithm. It's easy to implement and understand, making it a good choice when we want to quickly analyze the relationship between a few key variables and purchasing behavior.

Linear Regression is a suitable algorithm for predicting customer purchasing behavior when we have a reasonable expectation of a linear relationship between predictor variables and the target variable. It offers interpretability and simplicity, making it a valuable tool in certain scenarios. However, it's essential to assess whether its assumptions hold and whether it provides adequate predictive performance for your specific dataset and goals.

4. Hyperparameter tuning involves optimizing a machine learning model's settings, called hyperparameters, for better performance. It's vital because:
   a. Improved Performance: Proper hyperparameter choices enhance model accuracy and prevent overfitting.
   b. Resource Efficiency: Efficient hyperparameters save training time and resources.
   c. Robustness: Tuned hyperparameters make models adaptable to diverse data

Grid Search is a systematic technique:

   a. It explores hyperparameter combinations, such as learning rates and tree numbers, exhaustively.
   b. Cross-validation assesses each combination's performance, finding the best setup.
   c. Benefits include thorough exploration, reproducibility, and ease of implementation. However, it can be computationally intensive.