

MIE1628HY Cloud-based Data Analytics

Assignment 5

Faculty of Engineering – University of Toronto – Graduate Studies

Spring/Summer 2023

Submitted by: Zorina Song

UTORid: songzir1

Student#: 1005637198

Email: zo.song@mail.utoronto.ca

Submission date: 08/13/2023

Team member: Jingcheng Zhao,

UTORid: zhaoj116

Student#: 1003797965

Email: Jensen.zhao@mail.utoronto.ca

Part A:

1. [Marks: 5] Explain below the 5 components shown in orange boxes. Explain which Azure components you will use where in this big data architecture and why.

Explanation of the 5 components:

- Azure Data Lake: cloud platform, which allow us to store and analyze large volumes of data in various formats, such as structured, semi-structured, and unstructured data. It provides capabilities for big data processing, analytics, and machine learning.
- Azure Databricks: fully managed first-party service which could enables an open data lakehouse in Azure. Offers an integrated environment to enable multiple departments work together on big data and machine learning projects. It provides tools for data preprocessing, analytics, and machine learning model development, leveraging the power of distributed computing and parallel processing.
- Azure Data Factory: a cloud ETL service which can be used to scale-out serverless data integration. It can also support data movement across on-premises and cloud environments and facilitates data transformations.
- Azure Synapse Analytics: a cloud-based analytics service that combines big data and data warehousing capabilities into a single unified platform. It allows us to analyze large volumes of data using both on-demand and provisioned resources. Synapse Analytics integrates data storage, data processing, and data integration, enabling us to perform complex analytics tasks on massive datasets.
- Azure Cosmos DB: a globally distributed, multi-model database service, which can handle massive amounts of data with low-latency and high

availability for managing data at large scales. Various data models, including document, key-value, graph, and column-family, can be supported, allowing us to choose the right model for our application.

Which Azure components you will use in this big data architecture:

- Ingest Data: Azure Data Factory

Ideal for ingesting data from various sources into the data processing pipeline, as it provides connectors to a wide range of data sources and destinations and allows us to define data pipelines, schedule data movement, and manage data workflows effectively.

- Data Store: Azure Data Lake Storage

It is optimized for storing and analyzing large amounts of data in its raw, unprocessed form. We can accommodate both structured and unstructured data, making it suitable for big data scenarios.

- Prepare and transform data: Azure Databricks

Due to its distributed processing capabilities and interactive notebooks, and it can handle complex data processing tasks at scale, making it well-suited for cleaning, transforming, and enriching large datasets before analysis.

- Model and serve data: Azure Synapse Analytics

As a fully managed analytics service that brings together big data and data warehousing into a single unified platform. It's designed for querying and analyzing large datasets at scale. Allows us to build data warehouses and analytical pipelines that integrate data from various sources.

2. [Marks: 5] Explain how Stream Analytics works in Azure.

Ans:

As a real-time data processing service in Azure, it takes in streaming data from various sources, applies SQL-like queries for transformations, aggregations, and filtering, and then outputs the results to different destinations. We can use it for analyzing live data, providing real-time insights and integration with other Azure services.

Data Ingestion ==> processing ==> query execution

==> windowing and time Management (allowing us to perform operations on data over specific time intervals)

==> output sinks (destinations where the processed data is stored or sent for further analysis)

==> scaling and monitoring

3. [Marks: 10] Deploy all the resources in Azure Portal. Implement a Stream Analytics job by using the Azure portal.

ANS:

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

Home > All resources > a5spstorage | Containers > spcontainer >

spcontainer

Container

Search

Upload | Change access level

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Azure AD User Account)

Location: spcontainer

Search blobs by prefix (case-...)

Show deleted blobs

Add filter

Name

0_c77cce6f023b4e4ba39ce30e...

0_c77cce6f023b4e4ba39ce30e44b19f4d_1.json

Blob

Save | Discard | Download | Refresh | Delete

Overview | Versions | Snapshots | Edit | Generate SAS

```
1 {"messageId":32,"deviceId":"Raspberry Pi Web Client","temperature":27.3964672290422}
2 {"messageId":33,"deviceId":"Raspberry Pi Web Client","temperature":28.84907488052965,
3 {"messageId":34,"deviceId":"Raspberry Pi Web Client","temperature":25.865944918925354
4 {"messageId":35,"deviceId":"Raspberry Pi Web Client","temperature":31.198984760666903
5 {"messageId":37,"deviceId":"Raspberry Pi Web Client","temperature":26.99186964008029,
6 {"messageId":38,"deviceId":"Raspberry Pi Web Client","temperature":30.911526978040158
7 {"messageId":39,"deviceId":"Raspberry Pi Web Client","temperature":31.042709619878053
8 {"messageId":41,"deviceId":"Raspberry Pi Web Client","temperature":28.16153376967659,
9 {"messageId":43,"deviceId":"Raspberry Pi Web Client","temperature":30.481543604547063
10 {"messageId":45,"deviceId":"Raspberry Pi Web Client","temperature":26.748834078244864
11 {"messageId":46,"deviceId":"Raspberry Pi Web Client","temperature":28.64476494613961,
12 {"messageId":47,"deviceId":"Raspberry Pi Web Client","temperature":26.972108141197804
13 {"messageId":48,"deviceId":"Raspberry Pi Web Client","temperature":27.00405726511784
14 {"messageId":49,"deviceId":"Raspberry Pi Web Client","temperature":30.93753297689195
15 {"messageId":50,"deviceId":"Raspberry Pi Web Client","temperature":29.73429106314599,
16 {"messageId":51,"deviceId":"Raspberry Pi Web Client","temperature":30.186477628782978
```

Json | Preview

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

Home > All resources > a5spstorage | Containers > spcontainer >

spcontainer

Container

Search

Upload | Change access level

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Azure AD User Account)

Location: spcontainer

Search blobs by prefix (case-...)

Show deleted blobs

Add filter

Name

0_c77cce6f023b4e4ba39ce30e...

0_c77cce6f023b4e4ba39ce30e44b19f4d_1.json

Blob

Save | Discard | Download | Refresh | Delete

Overview | Versions | Snapshots | Edit | Generate SAS

```
16 {"messageId":51,"deviceId":"Raspberry Pi Web Client","temperature":30.186477628782978,
17 {"messageId":52,"deviceId":"Raspberry Pi Web Client","temperature":30.265776672581218,
18 {"messageId":53,"deviceId":"Raspberry Pi Web Client","temperature":27.57819747915143,
19 {"messageId":54,"deviceId":"Raspberry Pi Web Client","temperature":27.38036701500195,
20 {"messageId":55,"deviceId":"Raspberry Pi Web Client","temperature":31.629805187276086,
21 {"messageId":56,"deviceId":"Raspberry Pi Web Client","temperature":25.00996525075396,
22 {"messageId":57,"deviceId":"Raspberry Pi Web Client","temperature":29.394093675905644,
23 {"messageId":58,"deviceId":"Raspberry Pi Web Client","temperature":31.64712323746891,
24 {"messageId":59,"deviceId":"Raspberry Pi Web Client","temperature":30.985907370669672,
25 {"messageId":61,"deviceId":"Raspberry Pi Web Client","temperature":28.715057822215353,
26 {"messageId":62,"deviceId":"Raspberry Pi Web Client","temperature":31.1940869545853523,
27 {"messageId":63,"deviceId":"Raspberry Pi Web Client","temperature":26.749062494901697,
28 {"messageId":64,"deviceId":"Raspberry Pi Web Client","temperature":25.855635069857534,
29 {"messageId":65,"deviceId":"Raspberry Pi Web Client","temperature":27.42693310064791,
30 {"messageId":70,"deviceId":"Raspberry Pi Web Client","temperature":29.01849034623413,
31 {"messageId":71,"deviceId":"Raspberry Pi Web Client","temperature":29.67444990456427,
```

Json | Preview

Deployed Resources:

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

Home >

All resources

University of Toronto

Create | Manage view | Refresh | Export to CSV | Open query | Assign tags | Delete

Filter for any field...

Subscription equals all | Resource group equals all | Type equals all | Location equals all | Add filter

0 Recommendations | 0 Unsecure resources

No grouping | List view

Name	Type	Resource group	Location	Subscription
1628iot	IoT Hub	A5	East Asia	Azure for Students
a5spstorage	Storage account	A5	East Asia	Azure for Students
spasajob	Stream Analytics job	A5	East Asia	Azure for Students

Part B:

Question 1:

The dataset contains the following features for each concrete mixture:

1. Cement (component 1): The amount of cement(kg) in the mixture(m3).
2. Blast Furnace Slag (component 2): The amount of blast furnace slag(kg) in the mixture(m3).
3. Fly Ash (component 3): The amount of fly ash(kg) in the mixture(m3).
4. Water (component 4): The amount of water(kg) in the mixture(m3).
5. Superplasticizer (component 5): The amount of superplasticizer(kg) in the mixture(m3).
6. Coarse Aggregate (component 6): The amount of coarse aggregate(kg) in the mixture(m3).
7. Fine Aggregate (component 7): The amount of fine aggregate(kg) in the mixture(m3).
8. Age: The age of the concrete in days.

Problem Statement:

Given the dataset of concrete mixtures with their corresponding components and their actual compressive strength (measured in MPa), our task is to build a regression model capable of accurately predicting the concrete's compressive strength for any given mixture of components and its age. The successful development of such a model can be valuable in construction and civil engineering, as it would help engineers, contractors, and researchers estimate the strength of concrete structures without conducting time-consuming and costly physical tests. The model can aid in optimizing concrete mixtures for specific applications, ensuring safety, durability, and cost-effectiveness in construction projects.

Question 2:

Explore Dataset:

Azure AI | Machine Learning Studio

University of Toronto > azureml > Notebooks

ME1620_Part8.ipynb

Kernel: jupyterlab - Running

Python 3.10 - SDK v2

```

1 %pip install openpyxl
2 import pandas as pd
3 df = pd.read_excel('Concrete_Data.xlsx', engine='openpyxl')

```

Requirement already satisfied: openpyxl in /anaconda/envs/azureml_py310_sdkv2/11b/python3.10/site-packages (3.1.2)
Requirement already satisfied: et-xmlfile in /anaconda/envs/azureml_py310_sdkv2/11b/python3.10/site-packages (from openpyxl) (1.1.0)
Note: you may need to restart the kernel to use updated packages.

```

1 df.head(10)

```

	Cement (component 1)(kg in a m ³ mixture)	Blast Furnace Slag (component 2)(kg in a m ³ mixture)	Fly Ash (component 3)(kg in a m ³ mixture)	Water (component 4)(kg in a m ³ mixture)	Superplasticizer (component 5)(kg in a m ³ mixture)	Coarse Aggregate (component 6)(kg in a m ³ mixture)	Fine Aggregate (component 7)(kg in a m ³ mixture)	Age (day)	Concrete compressive strength(MPa, megapascals)
0	540.0	0.0	0.0	162.0	2.5	1040.0	676.0	28	79.86111
1	540.0	0.0	0.0	162.0	2.5	1055.0	676.0	28	61.887366
2	332.5	142.5	0.0	228.0	0.0	932.0	594.0	270	40.299535
3	332.5	142.5	0.0	228.0	0.0	932.0	594.0	365	41.052780
4	198.6	132.4	0.0	192.0	0.0	978.4	825.5	360	44.296075
5	266.0	114.0	0.0	228.0	0.0	932.0	670.0	90	47.029847
6	380.0	95.0	0.0	228.0	0.0	932.0	594.0	365	43.696299
7	380.0	95.0	0.0	228.0	0.0	932.0	594.0	28	36.447770
8	266.0	114.0	0.0	228.0	0.0	932.0	670.0	28	45.854291
9	475.0	0.0	0.0	228.0	0.0	932.0	594.0	28	39.289790

```

1 #check the size of the dataframe
2 num_rows, num_columns = df.shape
3
4 print(f"Number of rows: {num_rows}")
5 print(f"Number of columns: {num_columns}")

```

Number of rows: 1030
Number of columns: 9

```

1 #get the basic information of the dataframe
2 print(df.describe())

```

```

Cement (component 1)(kg in a m3 mixture) \
count      1030.000000
mean       281.165631
std        104.587142
min        102.000000
25%        192.375000
50%        272.900000
75%        350.000000
max        540.000000

Blast Furnace Slag (component 2)(kg in a m3 mixture) \
count      1030.000000
mean        73.895485
std         86.279104
min          0.000000
25%          0.000000
50%          0.000000
75%          0.000000
max          0.000000

```

```

1 new_column_names = ['Cement', 'Blast_Furnace_Slag', 'Fly_Ash', 'Water', 'Superplasticizer', 'Coarse_Aggregate', 'Fine_Aggregate', 'Age', 'Concrete_Strength']
2 df_new = df
3 df_new.columns = new_column_names

```

There is a description of the features in the above section.

Through our exploration, we have identified 8 input features and 1 target feature in our dataset. These features are all in numerical and continuous formats. The only input feature in integer is the 'age' column, which represents the number of. There are a total of 1030 entries in our sample dataset.

Upon analyzing the initial 10 rows of data, we observed instances where certain columns contained values of 0. These columns include measurements like the amount of blast furnace slag (in kg/m³) and fly ash (in kg/m³).

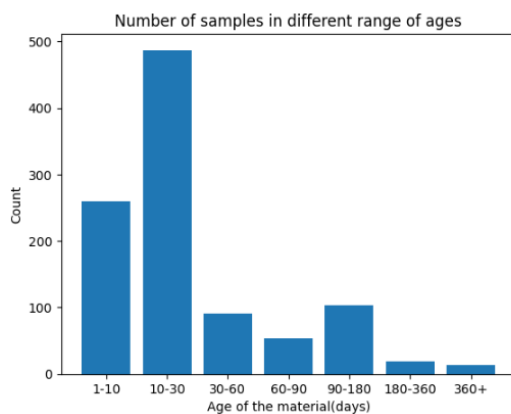
Furthermore, we have computed fundamental statistics such as the mean, standard deviation, minimum, and maximum values for each feature. These statistics provide valuable insights into the data distribution, aiding in our understanding of the dataset and guiding us in our exploration for subsequent steps.

Here are the 6 plots with explanations:

Plot1: Number of samples in different range of ages

```
1 import matplotlib.pyplot as plt
2 bins = [1, 10, 30, 60, 90, 180, 360, float('inf')]
3 labels = ['1-10', '10-30', '30-60', '60-90', '90-180', '180-360', '360+']
4 df_new['Age_Count'] = pd.cut(df_new['Age'], bins=bins, labels=labels)
5 # Count the occurrences of each range
6 range_counts = df_new['Age_Count'].value_counts().sort_index()
7
8 # Plotting
9 plt.bar(range_counts.index, range_counts.values)
10 plt.xlabel('Age of the material(days)')
11 plt.ylabel('Count')
12 plt.title('Number of samples in different range of ages')
13 plt.show()
```

[6] ✓



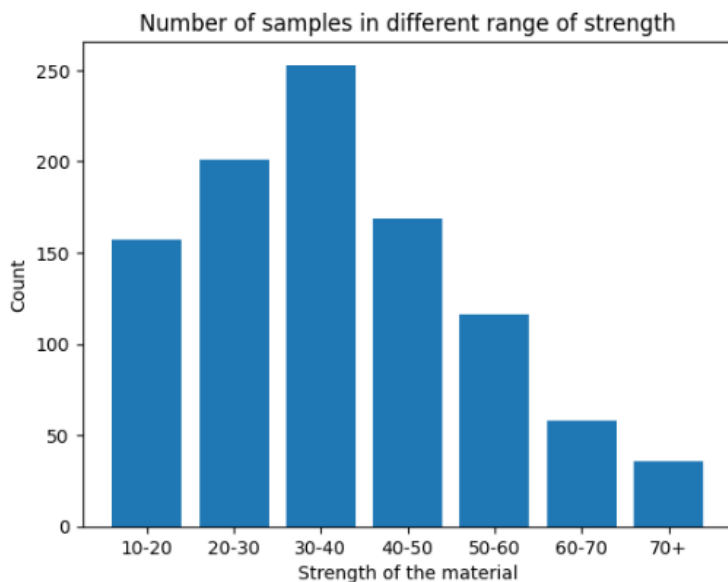
Explain:

In summary, the graph shows that most samples are concentrated in the 1 to 30-day age range, with the peak count in the 10-30 days range which is almost 500 counts and there are over 250 counts in the lower 1 to 10-day age range. The count generally diminishes with greater age ranges beyond 30 days, but the "90-180 days" range presents itself as the third most frequent range, with around 100 counts. This prominence might be influenced by the wider span of this interval compared to the previous ones, potentially encompassing a broader range of samples.

Plot2: Number of samples in different range of strength

```
1 bins = [10, 20, 30, 40, 50, 60, 70, float('inf')]
2 labels = ['10-20', '20-30', '30-40', '40-50', '50-60', '60-70', '70+']
3 df_new['Strength_Count'] = pd.cut(df_new['Concrete_Strength'], bins=bins, labels=labels, right=False)
4 # Count the occurrences of each range
5 range_counts = df_new['Strength_Count'].value_counts().sort_index()
6
7 # Plotting
8 plt.bar(range_counts.index, range_counts.values)
9 plt.xlabel('Strength of the material')
10 plt.ylabel('Count')
11 plt.title('Number of samples in different range of strength')
12 plt.show()
```

[7]



Explain:

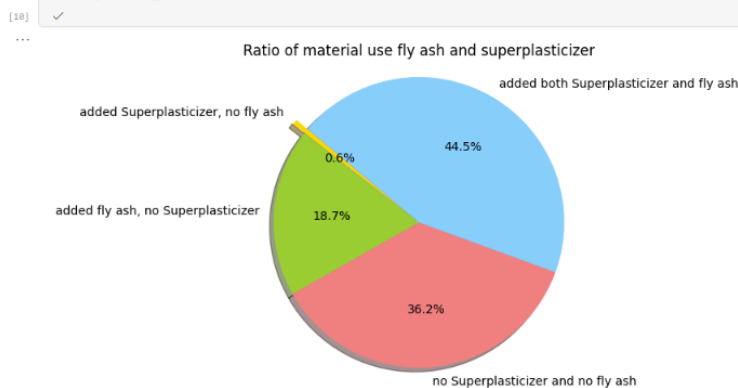
The plot shows a slight right-skewed distribution. There are three main phases with increased strength:

1. Increasing strength: As the strength value increases from 10 MPa, the number of samples also rises. This indicates that, within the strength 10-40 MPa, higher strength values are more common among the samples.
2. Peak frequency: The most frequent strength values are found in the range of 30-40 MPa. Around 250 samples have strength within this range, which makes it the peak of the distribution.
3. Decreasing trend: Following the 30-40 MPa range, the counts of samples start to decrease. This means that fewer samples have strength values beyond this point in our samples.

Plot3: Ratio of material use fly ash and superplasticizer

```
[9] ✓
1 ash_nonzero_sup_zero = ((df_new['Fly_Ash'] != 0) & (df_new['Superplasticizer'] == 0)).sum()
2 sup_nonzero_ash_zero = ((df_new['Superplasticizer'] != 0) & (df_new['Fly_Ash'] == 0)).sum()
3 both_zero = ((df_new['Fly_Ash'] == 0) & (df_new['Superplasticizer'] == 0)).sum()
4 both_nonzero = ((df_new['Fly_Ash'] != 0) & (df_new['Superplasticizer'] != 0)).sum()
```

```
1 # Create a pie chart
2 labels = ['added Superplasticizer, no fly ash', 'added fly ash, no Superplasticizer', 'no Superplasticizer and no fly ash', 'added both Superplasticizer and fly ash']
3 sizes = [ash_nonzero_sup_zero, sup_nonzero_ash_zero, both_zero, both_nonzero]
4 colors = ['gold', 'yellowgreen', 'lightcoral', 'lightskyblue']
5 explode = (0.1, 0, 0, 0) # explode 1st slice
6
7 plt.pie(sizes, explode=explode, labels=labels, colors=colors,
8         autopct='%1.1f%%', shadow=True, startangle=140)
9
10 plt.axis('equal')
11 plt.title('Ratio of material use fly ash and superplasticizer')
12 plt.show()
```



Explain:

The most common material usage scenario involves adding both fly ash and superplasticizer components to the mixture, constituting approximately 44.5% of the samples.

Interestingly, the second-largest group of samples consists of those that do not include either fly ash or superplasticizer, which accounts for a substantial portion of 36.2%.

Around 18.7% of the samples added fly ash only, without the addition of superplasticizer, while only a light fraction, about 0.6%, of the samples use superplasticizer as the exclusive component in the mixture.

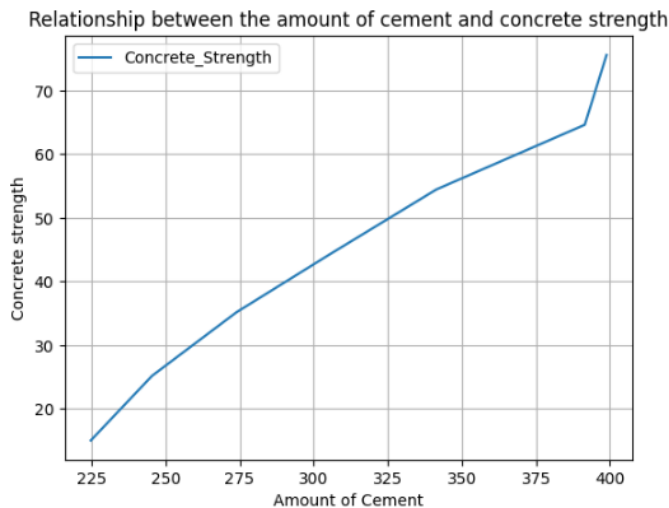
Plot4. Relationship between the amount of cement and concrete strength

11]

```
1 df_s = df_new.drop('Age_Count',axis = 1).groupby("Strength_Count").mean()
2 df_s = df_s.reset_index(drop=False)
```

12]

```
1 df_s.plot(x='Cement', y='Concrete_Strength', kind='line')
2
3 plt.ylabel('Concrete strength')
4 plt.xlabel('Amount of Cement')
5 plt.title('Relationship between the amount of cement and concrete strength')
6 plt.grid(True)
7 plt.show()
```



Explain:

Most data points form a pattern that ascends as the amount of cement increases, which indicates a positive linear relationship between the amount of cement(kg/m^3) and the concrete strength (Mpa). In other words, as the amount of cement used in the mixture increases, the concrete strength tends to increase as well. The pattern observed aligns with our existing knowledge that an increased amount of cement is likely to result in higher concrete strength.

However, the plot shows a turning point around 380 kg/m^3 of cement. Beyond this point, the rate of strength increases experiences a sudden and dramatic rise. This suggests that the rate of strength improvement becomes notably higher for cement quantities exceeding 380 kg/m^3 .

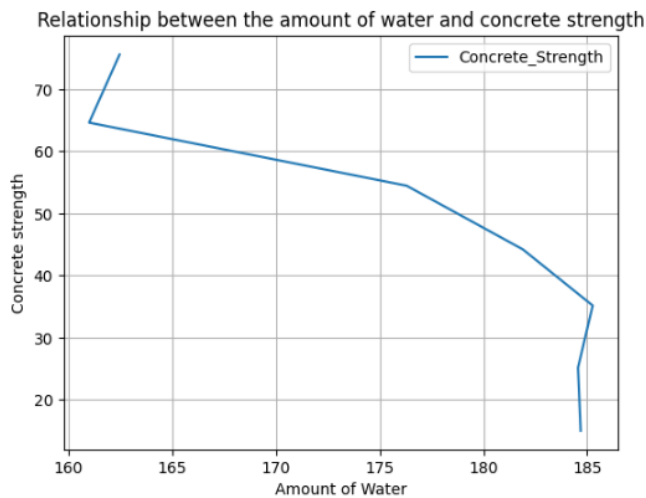
Plot5. Relationship between the amount of water and concrete strength

+ Code + Markdown

▶ | ▼

```
1 df_s.plot(x='Water', y='Concrete_Strength', kind='line')
2
3 plt.ylabel('Concrete strength')
4 plt.xlabel('Amount of Water')
5 plt.title('Relationship between the amount of water and concrete strength')
6 plt.grid(True)
7 plt.show()
```

[13]



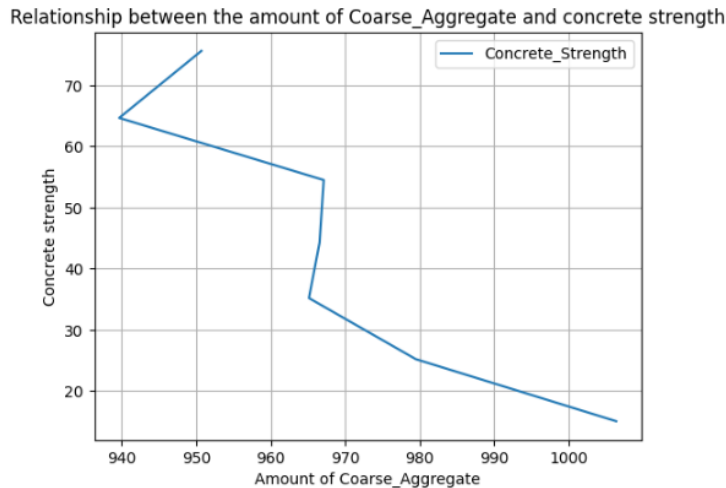
Explain:

In general, there is a negative linear relationship between the amount of water and concrete strength. As the amount of water increases, the concrete strength tends to decrease. Some data points at the extreme ends of the water amount spectrum do not conform to the general trend. These outliers or noises might represent unusual or non-standard scenarios that deviate from the typical relationship. Around 176 kg/m³ of water, a turning point was observed that beyond this point, higher amounts of water seem to lead to a more rapid decrease in concrete strength.

Plot6. Relationship between the amount of coarse aggregate and concrete strength

```
1 df_s.plot(x='Coarse_Aggregate', y='Concrete_Strength', kind='line')
2
3 plt.ylabel('Concrete strength')
4 plt.xlabel('Amount of Coarse_Aggregate')
5 plt.title('Relationship between the amount of Coarse_Aggregate and concrete strength')
6 plt.grid(True)
7 plt.show()
```

[44]



Explain:

Overall, there seems to be a negative relationship between the amount of coarse aggregate and concrete strength, such that the concrete strength tends to decrease as the amount of coarse aggregate increases. Within the interval of 940 to 950 kg/m³ of coarse aggregate, there are fluctuations in concrete strength, suggesting there are variations in strength within this specific range of coarse aggregate amounts. Besides, a dramatically decreases in the slope of the line occurs between 960 and 970 kg/m³ of coarse aggregate, indicating that changes in coarse aggregate beyond this range might have a more pronounced impact on concrete strength. Beyond the drop point, the trend indicates a more steady and consistent decrease in concrete strength as coarse aggregate amounts increase.

Question 3:

3.1 Replacing NaN Values with Column Averages:

Question3

Change the columns name to more readable form

```
[23] 1 df_predict = pd.read_excel('Concrete_Data.xlsx',engine='openpyxl')
      2
      3 new_column_names = ['Cement', 'Blast_Furnace_Slag', 'Fly_Ash', 'Water', 'Superplasticizer', 'Coarse_Aggregate', 'Fine_Aggregate', 'Age', 'Concrete_Strength']
      4 df_predict.columns = new_column_names
```

Fill all NAN values with average values

```
[24] 1 # Replace NaN values with the average of their respective columns
      2 df_predict.fillna(df_predict.mean(), inplace=True)
```

NaN values (missing values) can disrupt the training of the model in the further step. Replacing them with the average of their respective columns is a common strategy to ensure that the data maintains its general distribution while filling in the gaps.

3.2 Removing Duplicates:

Removing Duplicates

```
[25] 1 # Remove duplicates
      2 df_predict = df_predict.drop_duplicates()
```

Removing duplicates ensures that each data point is unique and representative, and result in an overly optimistic view of the model's performance.

3.3 Standardizing the Dataset:

```
1 from sklearn.preprocessing import StandardScaler
2
3 # Standardize the dataset
4 scaler = StandardScaler()
5 scaled_data = scaler.fit_transform(df_predict)
6
7 # Convert the scaled data array back to a DataFrame
8 scaled_df = pd.DataFrame(scaled_data, columns=df_predict.columns)
9
10 print("Original DataFrame:")
11 print(df_predict.head(5))
12 print("\nStandardized DataFrame:")
13 print(scaled_df.head(5))
```

[26] ✓

Original DataFrame:

	Cement	Blast_Furnace_Slag	Fly_Ash	Water	Superplasticizer \
0	540.0	0.0	0.0	162.0	2.5
1	540.0	0.0	0.0	162.0	2.5
2	332.5	142.5	0.0	228.0	0.0
3	332.5	142.5	0.0	228.0	0.0
4	198.6	132.4	0.0	192.0	0.0

	Coarse_Aggregate	Fine_Aggregate	Age	Concrete_Strength
0	1040.0	676.0	28	79.986111
1	1055.0	676.0	28	61.887366
2	932.0	594.0	270	40.269535
3	932.0	594.0	365	41.052780
4	978.4	825.5	360	44.296075

Standardization scales the features to have a mean of 0 and a standard deviation of 1. As the differences in the amounts of components (measured in kg/m^3) across the samples are noticeable, this step ensures that each feature contributes to the model's learning process equally and prevents features with larger numerical values from dominating features with smaller values.

After standardizing the data, working with it in a DataFrame format is a common data structure used in Pandas, as it's often more convenient to work with DataFrames in further analysis or use in models.

3.4 Train-Test Split:

Train-Test Split

```
1 from sklearn.model_selection import train_test_split
2
3 y = scaled_df['Concrete_Strength']
4 X = scaled_df.drop(columns = 'Concrete_Strength')
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
7
8 print("X_train:")
9 print(X_train)
10 print("\nX_test:")
11 print(X_test)
12 print("\ny_train:")
13 print(y_train)
14 print("\ny_test:")
15 print(y_test)
16
```

[31] ✓

... X_train:

	Cement	Blast_Furnace_Slag	Fly_Ash	Water	Superplasticizer	\
78	2.422701	-0.836469	-0.865363	-1.888146	3.746798	
29	1.882876	-0.836469	-0.865363	2.153088	-1.019442	
280	-0.015141	-0.836469	0.652195	-1.006763	0.980012	
507	1.508930	-0.836469	-0.865363	1.684267	-1.019442	
652	-1.072258	2.068521	-0.865363	1.004477	-1.019442	
..	
106	0.568311	1.357948	-0.865363	-0.336350	0.586206	
270	-1.052314	-0.346267	1.070735	-1.113185	0.810997	
860	-1.204577	0.847078	0.895434	-0.191015	0.332682	
435	-1.089517	0.655502	1.193055	-0.328849	0.346204	
102	0.073552	1.496115	-0.865363	-1.752188	0.873532	

Coarse Aggregate Fine Aggregate Age

A train-test split allows us to train the model on one subset of the data (training set) and test it on another subset (testing set) that it has not seen during training. It is essential to evaluate the model's performance on unseen data to understand its generalization ability. Here we choose 80% as training set and the remaining 20% as testing set, which is a common choice.

Question 4:

Machine learning model1: Decision Tree Regressor

```
1  from sklearn.model_selection import train_test_split
2  from sklearn.linear_model import LinearRegression
3  from sklearn.tree import DecisionTreeRegressor
4  from sklearn.ensemble import RandomForestRegressor
5  from sklearn.metrics import mean_squared_error, mean_absolute_error
6  # Train Decision Tree Regressor model
7  dt_model = DecisionTreeRegressor()
8  dt_model.fit(X_train, y_train)
9
10
11 # Make predictions
12 dt_pred = dt_model.predict(X_test)
13
14 # Calculate RMSE
15 dt_rmse = mean_squared_error(y_test, dt_pred, squared=False)
16
17 # Calculate MSE
18 dt_mse = mean_squared_error(y_test, dt_pred)
19
20 # Calculate MAE
21 dt_mae = mean_absolute_error(y_test, dt_pred)
22
23 # Display MSE and MAE values
24 print(f"Decision Tree MSE: {dt_mse:.2f}")
25 print(f"\nDecision Tree RMSE: {dt_rmse:.2f}")
26 print(f"\nDecision Tree MAE: {dt_mae:.2f}")
```

Decision Tree MSE: 0.15

Decision Tree RMSE: 0.39

Decision Tree MAE: 0.24

Machine learning model2: Random Forest Regressor

```
1  # Train Random Forest Regressor model
2  rf_model = RandomForestRegressor()
3  rf_model.fit(X_train, y_train)
4
5  # Make predictions
6  rf_pred = rf_model.predict(X_test)
7
8  # Calculate RMSE
9  rf_rmse = mean_squared_error(y_test, rf_pred, squared=False)
10
11 # Calculate MSE
12 rf_mse = mean_squared_error(y_test, rf_pred)
13
14 # Calculate MAE
15 rf_mae = mean_absolute_error(y_test, rf_pred)
16
17
18 # Display MSE, RMSE and MAE values
19 print(f"Random Forest MSE: {rf_mse:.2f}")
20 print(f"\nRandom Forest RMSE: {rf_rmse:.2f}")
21 print(f"\nRandom Forest MAE: {rf_mae:.2f}")
```

[43] Random Forest MSE: 0.10

Random Forest RMSE: 0.32

Random Forest MAE: 0.21

The reasons why these 2 algorithms are chosen:

We are facing a regression problem; therefore, the selection of the algorithms is limited to those can be used on continuous data type. Among the available machine learning algorithms. I selected decision tree regressor and random forest regressor.

The reasons why the decision tree regression is selected:

1. **Interpretability:** Decision trees provide a clear visualization of how decisions are made based on component values. This helps engineers and researchers understand the relationships between individual components and concrete strength.
2. **Component Interaction:** Decision trees can capture complex interactions between components. They allow for the identification of specific component thresholds that lead to changes in material strength.
3. **Feature Importance:** Engineers can easily identify the most critical components that influence material strength, aiding in material design and optimization.

The reasons why the random forest regression is selected:

1. **Improved Accuracy:** Random forests combine predictions from multiple decision trees, leading to better prediction accuracy by capturing various interactions in the data.
2. **Robustness:** Material strength prediction might involve noisy data or outliers. Random forests handle such variations well due to the ensemble averaging of multiple trees.
3. **Generalization:** Random forests generalize better to new data by reducing overfitting that can occur with a single decision tree.

The comparison between decision tree regression and random forest regression:

We evaluated Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE) for both Decision Tree Regression and Random Forest Regression models. These metrics are commonly used to evaluate the performance of regression models.

1. Mean Squared Error (MSE):

- MSE measures the average squared difference between predicted and actual values. A lower MSE indicates a better fit to the data.
- The Random Forest model has a lower MSE (0.1) compared to the Decision Tree model (0.15). This suggests that the Random Forest's ensemble approach has produced more accurate predictions on average.

2. Root Mean Squared Error (RMSE):

- RMSE is the square root of MSE and represents the average prediction error in the same units as the target variable.
- The Random Forest model has a lower RMSE (0.32) than the Decision Tree model (0.39). This again indicates that the Random Forest model's predictions are closer to the actual values.

3. Mean Absolute Error (MAE):

- MAE measures the average absolute difference between predicted and actual values.
- The Random Forest model has a slightly lower MAE (0.21) compared to the Decision Tree model (0.24), which means that the Random Forest predictions have a smaller absolute error on average.

In summary, the results suggest that the Random Forest Regression model is better than the Decision Tree Regression model in terms of prediction accuracy. The lower values of MSE, RMSE, and MAE for the Random Forest model indicate that it provides more accurate predictions on average. This improvement can be attributed to the ensemble nature of the Random Forest, which combines multiple decision trees to achieve better generalization and mitigate overfitting.

Question 5.

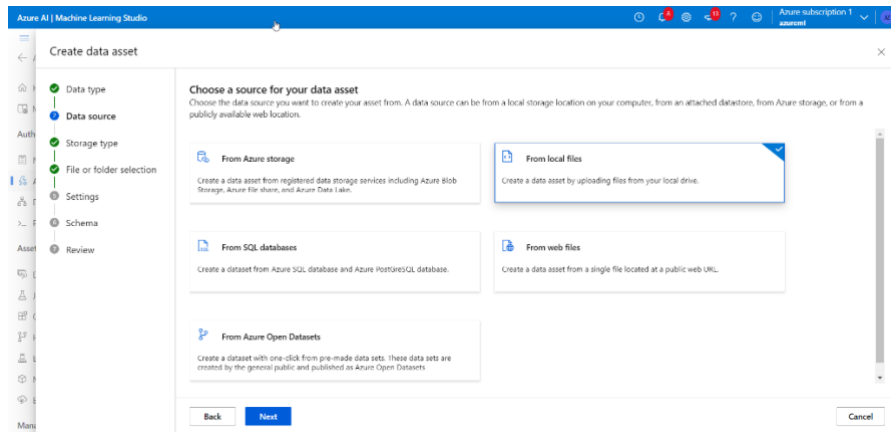
We choose to apply Automated ML on Microsoft Azure for our data set, here are the steps and results:

Setup procedure:

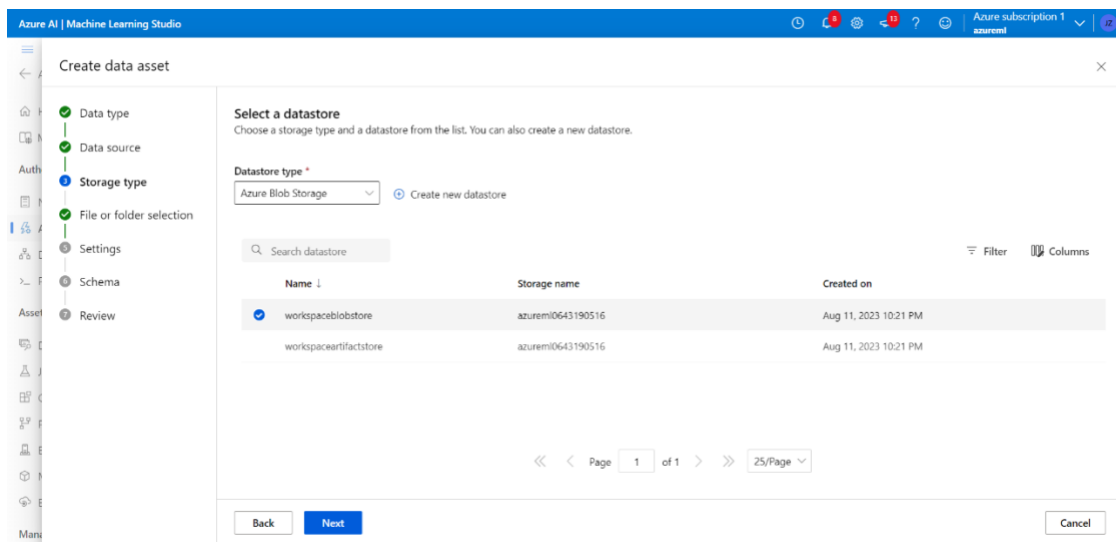
Step1: Create Data Asset (Data Type)

The screenshot shows the 'Create data asset' dialog in Azure AI Machine Learning Studio. The 'Data type' tab is active. The 'Name' field contains 'concrete_strength.csv'. The 'Description' field contains 'Csv file for automated Machine learning to predict the concrete strength'. The 'Type' dropdown is set to 'Tabular'. On the right, the 'Use cases for data types' section provides guidance: 'When should I use File type?' (recommended for single data files), 'When should I use Folder type?' (same as File type but for folders), and 'When should I use Table type?' (useful for advanced scenarios requiring schema definitions). Navigation buttons 'Back', 'Next', and 'Cancel' are at the bottom.

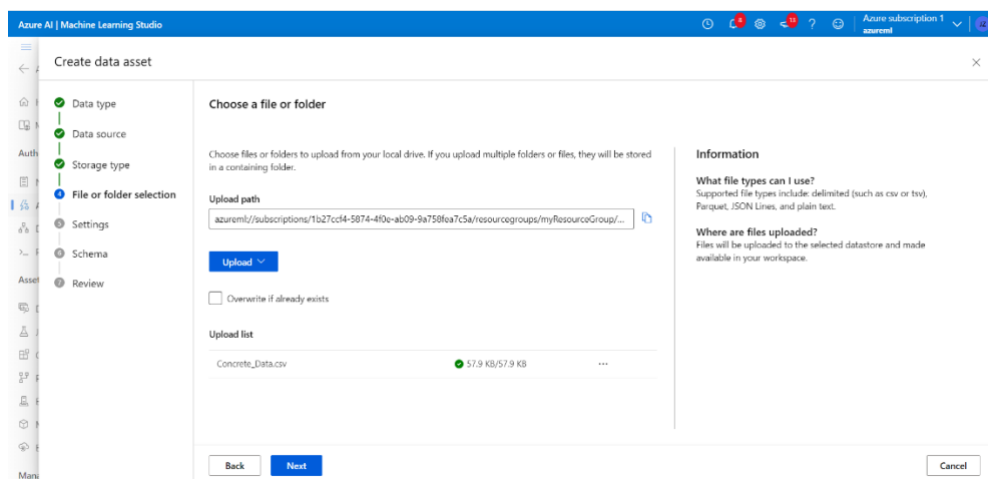
Step2: Create Data Asset (Data Source)



Step3: Create Data Asset (Storage type)



Step4: Create Data Asset ()



Step5: Create Data Asset (Setting)

Azure AI | Machine Learning Studio

Create data asset

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

Cement (comp...	Blast Furnace S...	Fly Ash (comp...	Water (compo...	Superplasticize...	Coarse Aggreg...	Fine Aggregat...	Age (day)	Concrete comp...
540	0	0	162	2.5	1,040	676	28	79.99
540	0	0	162	2.5	1,055	676	28	61.89
332.5	142.5	0	228	0	932	594	270	40.27
332.5	142.5	0	228	0	932	594	365	41.05
198.6	132.4	0	192	0	978.4	825.5	360	44.3
266	114	0	228	0	932	670	90	47.03

Back Next Review Cancel

Step6: Create Data Asset (Schema)

Azure AI | Machine Learning Studio

Create data asset

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
<input type="checkbox"/>	Path	String		Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Cement (component 1)(kg in a m...	Decimal (dot :)	540, 540, 332.5	Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Blast Furnace Slag (component ...	Decimal (dot :)	0, 0, 142.5	Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Fly Ash (component 3)(kg in a m...	Decimal (dot :)	0, 0, 0	Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Water (component 4)(kg in a m...	Decimal (dot :)	162, 162, 228	Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Superplasticizer (component 5)...	Decimal (dot :)	2.5, 2.5, 0	Not applicable to selected type	Not applicable to selected type
<input checked="" type="checkbox"/>	Coarse Aggregate (component ...	Decimal (dot :)	1040, 1055, 932	Not applicable to selected type	Not applicable to selected type

Back Next Cancel

Step7: Select Data Asset

Azure AI | Machine Learning Studio

University of Toronto > azureml > Automated ML > Start job

Create a new Automated ML job

1 Select data asset
2 Configure job
3 Select task and settings
4 Hyperparameter configuration (Computer Vision only)
5 Validate and test

Select data asset
Select an input data asset from the list below, or create a new data asset. AutomatedML currently only supports tabular data for authoring jobs.

+ Create Refresh Show supported data assets only

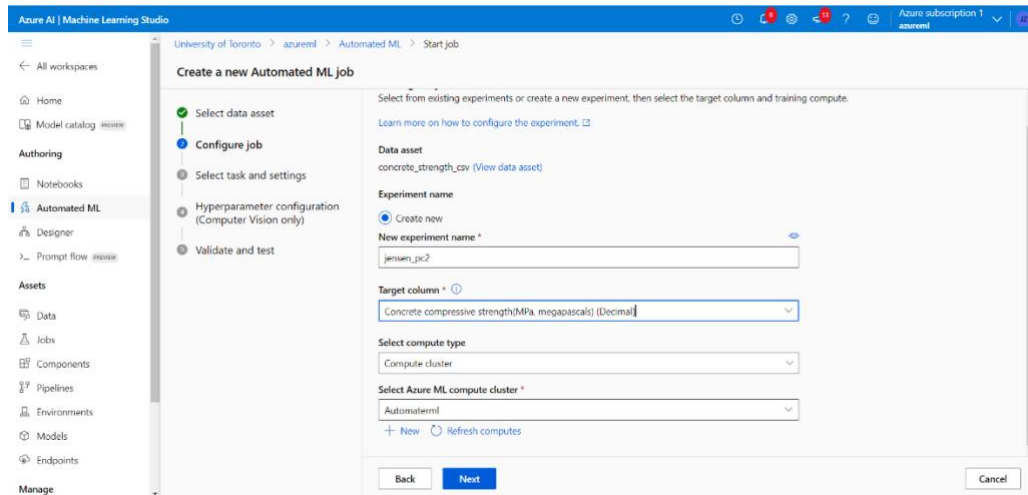
Search Filter

Name	Dataset type	Created on	Modified on
concrete_strength.csv	Tabular	Aug 13, 2023 3:59 PM	Aug 13, 2023 3:59 PM
concrete_strength_pred	Tabular	Aug 13, 2023 3:36 PM	Aug 13, 2023 3:36 PM
Concrete_Strength_2	Tabular	Aug 13, 2023 3:32 PM	Aug 13, 2023 3:32 PM
Concrete_Strength	Tabular	Aug 13, 2023 1:50 PM	Aug 13, 2023 1:50 PM

<< < Page 1 of 1 > >> 25/Page

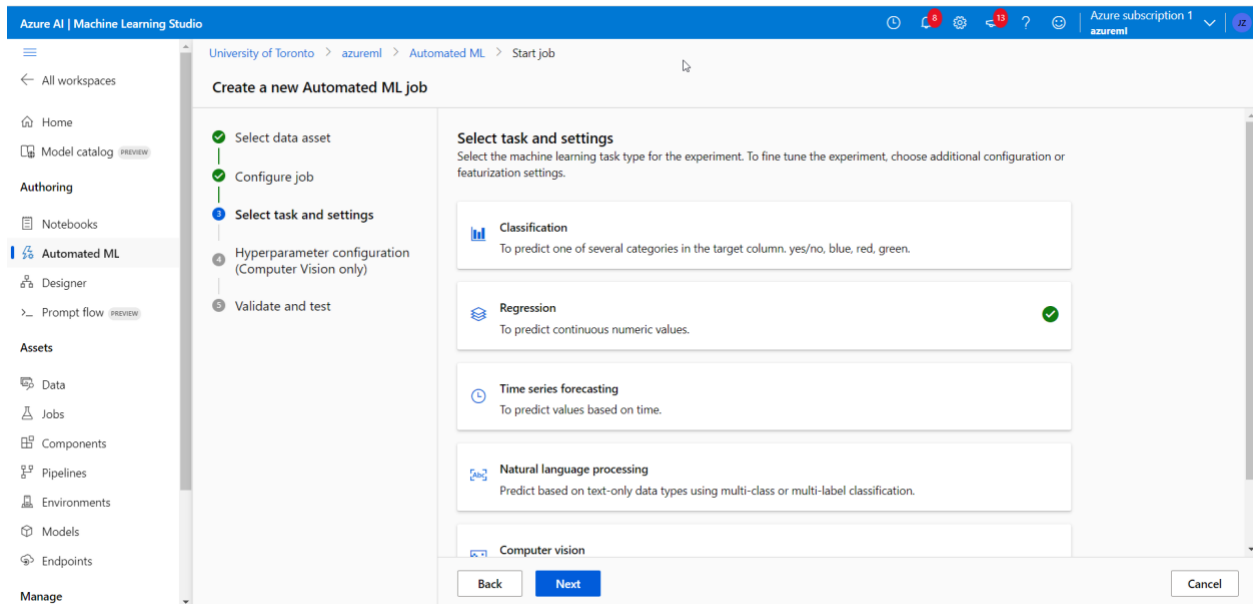
Back Next Cancel

Step8: Configure job



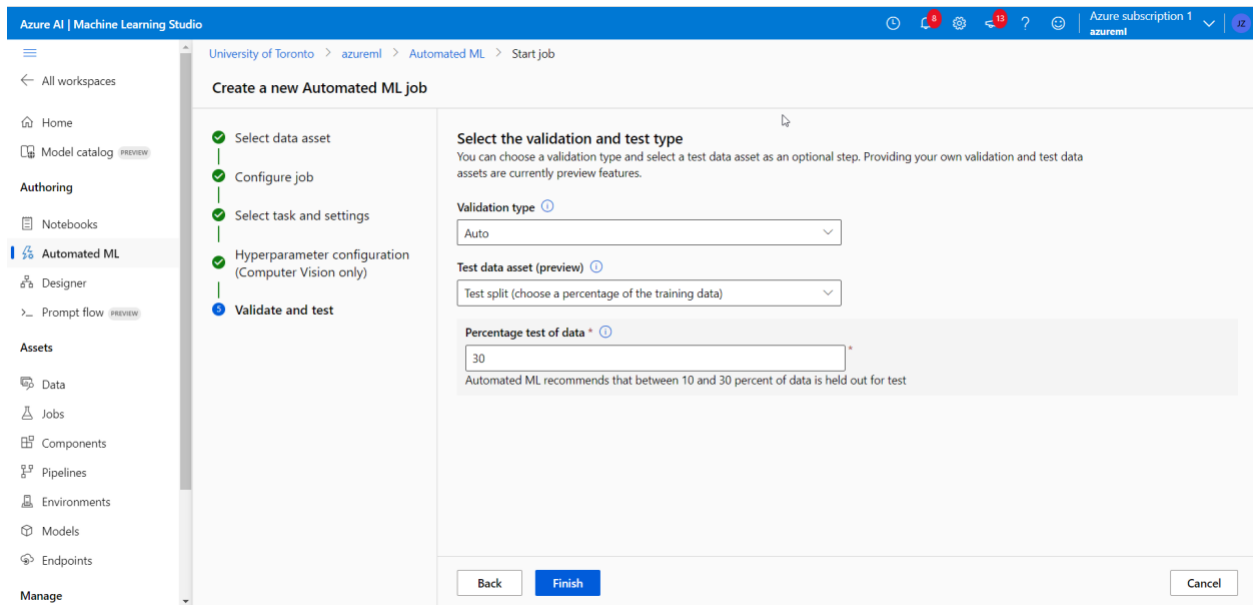
Step9: Select task and settings.

I selected regression.

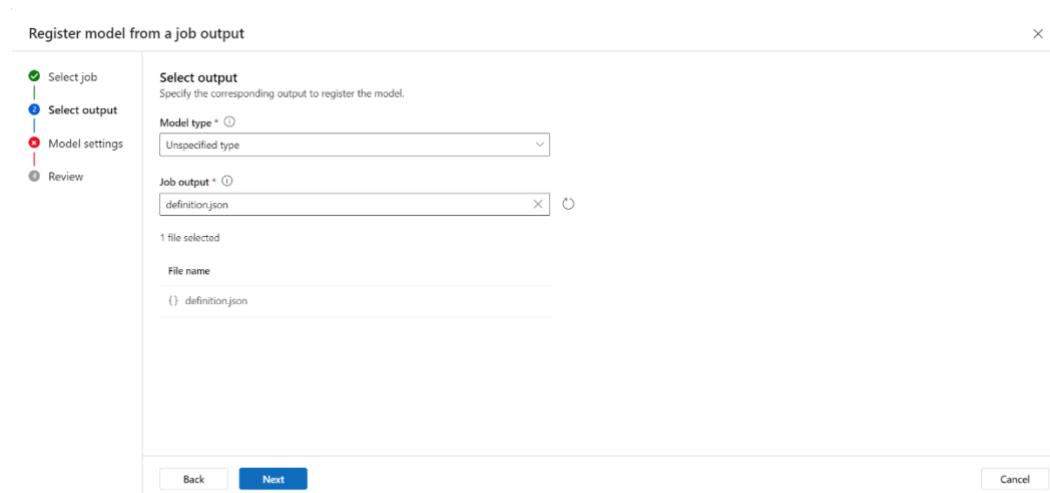


Step10: Hyperparameter configuration

The train test split ration is set to 70:30.



Step11: Register Model (Select output)



Step12: Register Model (Model Settings)

Register model from a job output

Model settings
Configure the settings for your model.

Name *
concrete_strength_predict

Description
predict concrete strength

Version
1

Tags
No tags

Back Next Cancel

Step13: Register Model (Review)

Register model from a job output

Review
Review or make changes to your selections.

Select job
Job: elated_rat_bic9ghy42

Select output
Model type: Unspecified type
Job output: definition.json
Files: 1 file selected
{ } definition.json

Model settings
Name: concrete_strength_predict
Description: predict concrete strength
Version: 1
Tags: No tags

Back Register Cancel

Automated ML Results:

Results Overview:

Azure AI | Machine Learning Studio

University of Toronto > azureml > Automated ML > jensen_pc2 > elated_rat_bk9ghj42

elated_rat_bk9ghj42 Completed

Overview Data guardrails Models Outputs + logs Child jobs

Refresh Edit and submit (preview) Register model Cancel Delete Compare (preview) >

Properties

Status Completed

Script name --

Created by Jingcheng ZHAO

Job type Automated ML

Experiment jensen_pc2

Arguments None

See all properties Raw JSON

See VAML job definition Job VAML

Created on Aug 13, 2023 4:01 PM

Start time Aug 13, 2023 4:02 PM

Duration 40m 51.78s

Compute duration 40m 51.78s

Compute target Automaterml

Name AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b

Tags

dynamic_allowlisting_iterations: <25><30><35>

fit_time_000:
0.0743110000000001;0.1093136000000001;0.0987280;0.0707090000000001;0.068815;0.0688065;0.0429945;0.0052428;0.0747303;0.019
754900000000002;0.01671620000000004;0.03042879999999995;0.0355677;0.007385199999999998;0.004311099999999999;0.004426200
000000001;0.0049356;0.0052074;0.005876299999999999;4.2263250999999999;0.024093299999999998;0.0118655;0.1182744;0.720476;0.42
5077100000000004;0.6301659000000001;0.21621839999999998;0.0398775;0.16840700000000003;0.027261199999999996;0.0901240000000
0001;0.0293037;0.005529499999999999;1.0241297;0.0859463;0.147534700000000002;0.4950977;0.6755931;0.4732145;0.0441416;21.23

Inputs

Input name: training_data
Dataset: concrete_strength_csv1

Outputs

Output name: best_model
Model: azureml_AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b_40_output_mlflow_log_model_2064710023:1

Best model summary

Algorithm name VotingEnsemble

Ensemble details View ensemble details

Normalized root mean squared error 0.05868 View all other metrics

Sampling 100.00 %

Registered models No registration yet

Deploy status No deployment yet

Run summary

Task type Regression View configuration settings

Featureization Auto

Best Model:

Azure AI | Machine Learning Studio

University of Toronto > azureml > Models > azureml_AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b_40_output_mlflow_log_model_2064710023:1

azureml_AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b_40_output_mlflow_log_model_2064710023:1

Details Artifacts Jobs Data Feature sets Responsible AI Explanations (preview) Fairness (preview)

Refresh Download all Register Share model

Attributes

Created on Aug 13, 2023 4:41 PM

Created by Jingcheng ZHAO

Type MLFLOW

Created by job AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b_40

Asset ID azureml:/locations/eastus2/workspaces/b08fa8a-0bba-4322-bf3c-b6c03f236e57/models/azureml_AutoML_0d044ff2-69aa-465b-ba35-127791c9c65b_40_output_mlflow_log_model_2064710023/versions/1

Build your models responsibly.
Try out the Responsible AI dashboard Try it out

Ensemble Details:

Ensemble details



Select an ensemble algorithm to see the ensemble weights and hyperparameters.

- ☒ StandardScalerWrapper, XGBoo...
- ☐ StandardScalerWrapper, XGBoo...
- ☐ SparseNormalizer, XGBoostReg...
- ☐ StandardScalerWrapper, LightG...
- ☐ MaxAbsScaler, ExtremeRandom...
- ☐ StandardScalerWrapper, Elastic...
- ☐ MaxAbsScaler, XGBoostRegress...
- ☐ MaxAbsScaler, LightGBM

Close

Ensemble weight: 0.13333333333333333

Data transformation:

```
1 {  
2   "class_name":  
   "StandardScaler",  
3   "module": "sklearn.  
preprocessing",  
4   "param_args": [],  
5   "param_kwargs": {  
6     "with_mean": false,  
7     "with_std": false  
8   },  
9   "prepared_kwargs": {},  
10  "spec_class": "preproc"  
11 }
```

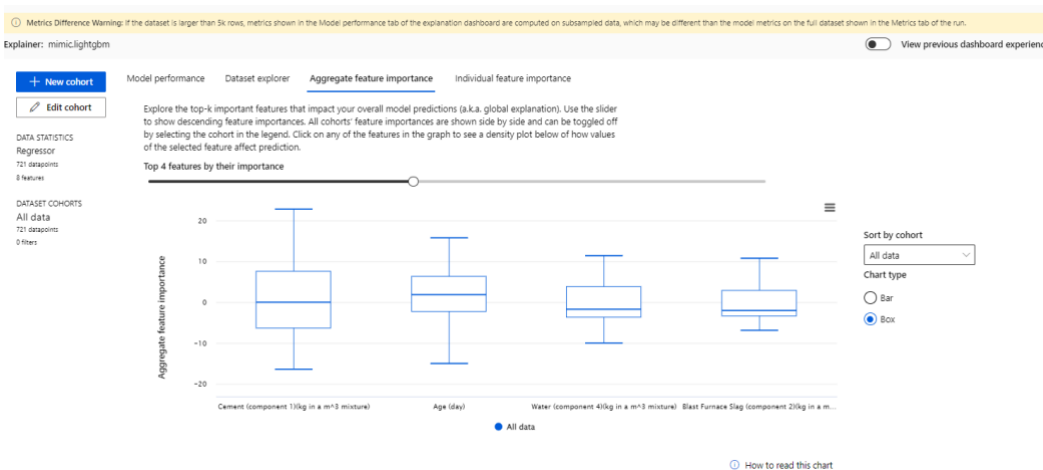
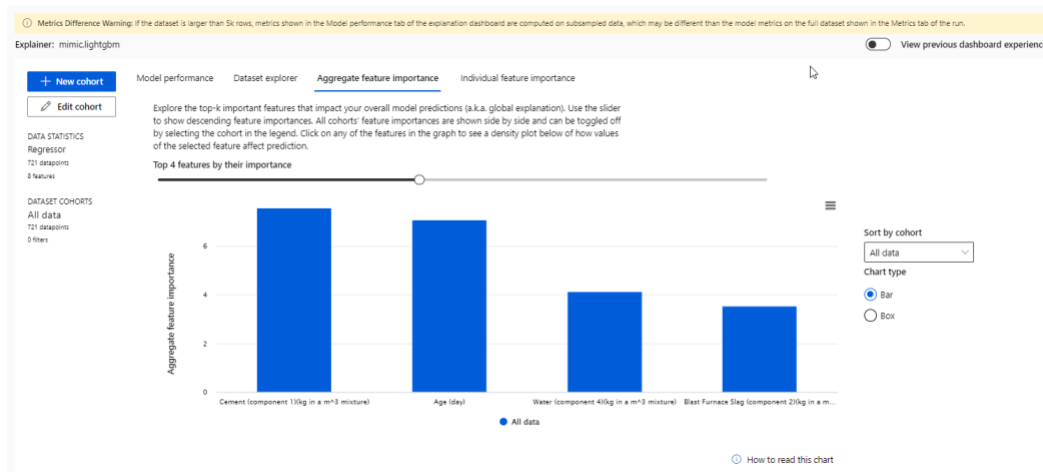
Training algorithm:

```
1 {  
2   "class_name":  
   "XGBoostRegressor",  
3   "module": "automl.  
client.core.common.  
model_wrappers",  
4   "param_args": [],  
5   "param_kwargs": {  
6     "booster": "gbtree",  
7     "colsample_bytree":  
1, 1,  
8     "eta": 0.3,  
9     "gamma": 0.1,  
10    "max_depth": 8,  
11    "max_leaves": 3,  
12    "n_estimators": 100,  
13    "objective":  
   "reg:linear",  
14    "reg_alpha": 2.5,  
15    "reg_lambda": 0.  
10416666666666667.
```

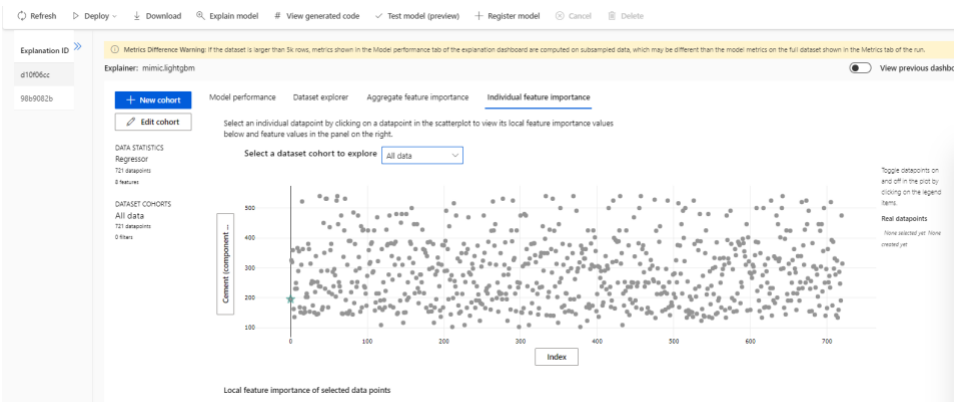
Best Model Performance:



Feature Importance:



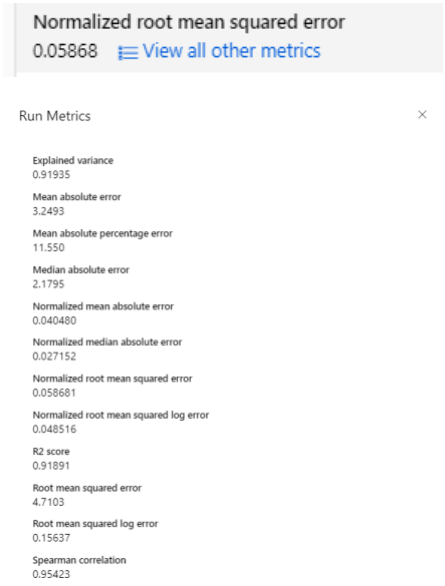
Individual feature importance:



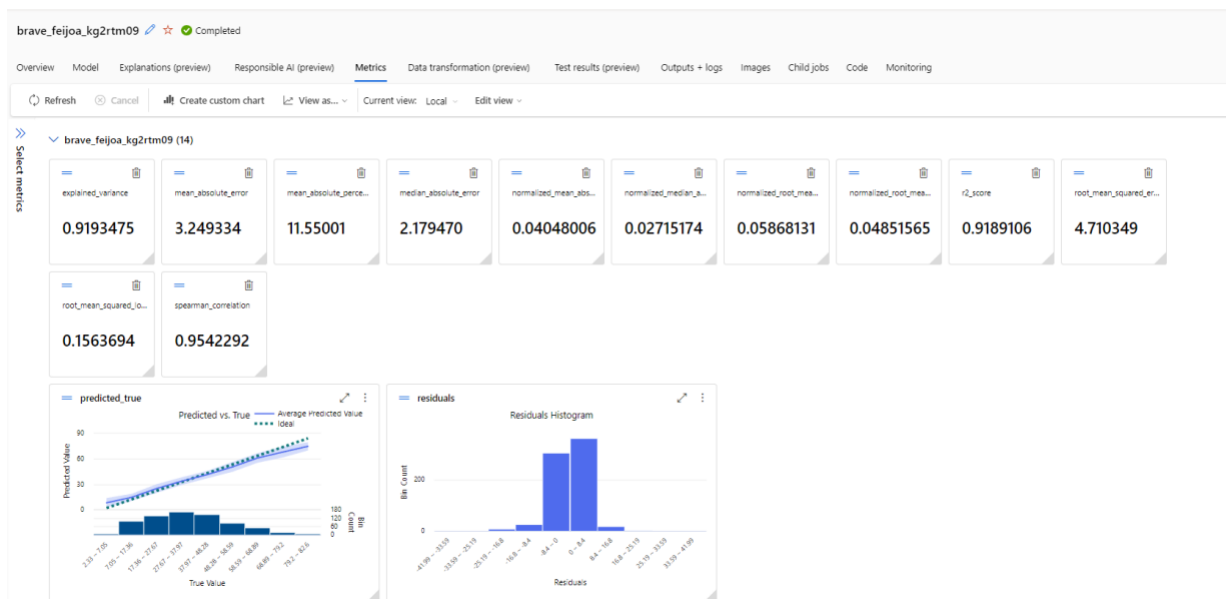
Data Transformation:



Metrics Results:



Final Results Visualization:



Best model results explanation:

The best model algorithm is **Voting Ensemble**. Explained Variance is 0.9193. A higher explained variance indicates that the model's predictions closely match the actual target values. In this case, 0.9193 is quite high, suggesting that the model's predictions capture a substantial amount of variance in the target variable. The mean absolute error is 3.249, a lower MAE indicates that, on average, the model's predictions are close to the actual values. In this case, 3.249 is not very ideal compared to the 2 models we trained in question4. The model's RMSE is 4.71, not excellent but it indicates that the model's predictions are, on average, closer to the actual values. The R-squared value of this model is 0.91891, which is very good since it's very close to 1. A higher R-squared value (close to 1) indicates that the model's predictions closely match the actual variability in the data. A lower MAPE indicates that, on average, the model's predictions have a smaller percentage error. In this case, the average percentage error is 11.55%.

Based on the metrics, the **Voting Ensemble** model seems to perform well. It captures a significant portion of the variance in the target variable, has low errors, and indicates a strong fit to the data. It's essential to consider the context of your problem and domain-specific requirements when interpreting these results. Meanwhile, as we can see in the predicted vs true line graph, the prediction line is very close to the true value line. However, the performance of the 2 models we trained in question4 is even better in terms of RMSE and MAE, which means there is still room for improvement in Microsoft Azure Automated ML.