# Flight DataSet Prediction

## Flight DataSet- An Overview

**Task 1:** Login to MS Azure and identify the Flight delay experiment which has to be used in this Project.

**Login to Microsoft Azure:**

Open your web browser and navigate to the Azure portal and enter your credentials to log in.

**Navigate to Azure Machine Learning Studio:**

In the Azure portal, find and click on "Machine Learning" in the left-hand menu.

**Locate the Flight Delay Experiment:**

Explore your workspace to find the experiment related to flight delay. Look for experiments in the Azure Machine Learning Studio that might be named or related to flight delay prediction.
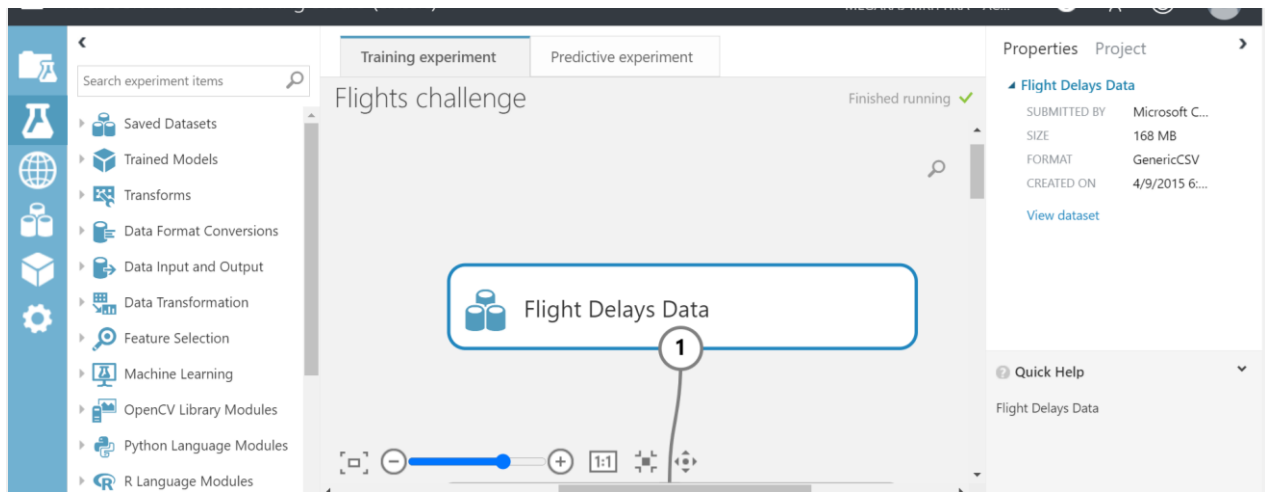
**Task 2:** Run the experiment

**Open the Experiment:** Click on the experiment name to open it in Azure Machine Learning Studio.

**Review Experiment Details:** Look for any parameters or configurations that need to be set before running the experiment.

**Run the Experiment:** Find and click on a "Run" or "Start" button to execute the experiment.

**Monitor Progress:** Monitor the progress of the experiment. This might involve checking logs, visualizations, or any output generated by the experiment.

**Review Results:**Once the experiment is complete, review the results. This could include performance metrics, charts, or any other relevant information.

**Task 3:** Explain the experiment and Regression/ Clustering model used in a single Page and include it in the Project Report

**Regression Model:**

In the Flight Delay Prediction experiment, regression is a kind of supervised learning used to forecast continuous outcomes. Its application here is to the estimation of the length of flight delays, a continuous range of values variable. With the help of predictive modeling, we can comprehend and predict how different factors may affect flight delays.

**Algorithm:**

The Boosted Decision Tree regression model that was selected uses an ensemble learning methodology. An essential component of this technique is decision trees, which produce a structure that partitions the input space into regions. Boosting is an ensemble strategy that systematically combines several weak learners (shallow decision trees). By correcting mistakes created by the previous trees, this procedure creates a reliable and accurate predictive model. The group dynamic aids in capturing intricate relationships within the data.

**Model Configuration:**

- **Learning Rate:**

    Each tree's contribution to the ensemble is determined by the learning rate (0.1). Striking a precise balance between training efficiency and accuracy was necessary.

- **Maximum Tree Depth:**

    To avoid overfitting, set each tree's maximum depth to five. This guarantees that the model can identify intricate patterns in the data without compromising its capacity to generalize to previously undiscovered data.

- **Number of Boosting Iterations:**

    After experimentation, 100 boosting iterations were found optimal. This number strikes a balance between the model's performance and computational efficiency.

**Training Process:**

- **Initial Training:**

    The complete dataset is used to train a basic decision tree. The overall patterns and trends in the flight delay data are represented by this tree.

- **Sequential Boosting:**

    Next, successive trees undergo sequential training. The goal of every tree is to reduce the faults that the ones before it were introduced. Over time, the predictive power of the model is strengthened by this iterative process.

**Evaluation Metrics:**

- **Mean Absolute Error(MAE):** The average absolute difference between the anticipated and actual flight delay periods is measured by the Mean Absolute Error or MAE. Better precision is indicated by a lower MAE.

- **Root Mean Squared Error (RMSE):** For big prediction mistakes, RMSE offers a more substantial penalty. By taking into account the square of the variations between the expected and actual values, it provides a thorough understanding of prediction accuracy.

**Boosted Decision Tree:**

**Overview:**

    Boosted Decision Trees are an ensemble learning technique that combines multiple weak learners, usually shallow decision trees, to form a robust and accurate predictive model. It's particularly effective for tasks like regression, where the goal is to predict continuous outcomes.

**Algorithm:**

- **Decision Trees:**

    Decision trees make sequential decisions based on input features, forming a tree-like structure.

- **Boosting:**

    Boosting involves training a series of decision trees sequentially, with each tree correcting errors made by its predecessors. The final prediction is a weighted sum of individual tree predictions.

**Results:**

The Boosted Decision Tree Regression model demonstrated remarkable performance in predicting flight delay durations. The low MAE and RMSE values indicate its suitability for accurately estimating the duration of flight delays. These results validate the effectiveness of the chosen algorithm and its configuration for this specific prediction task.
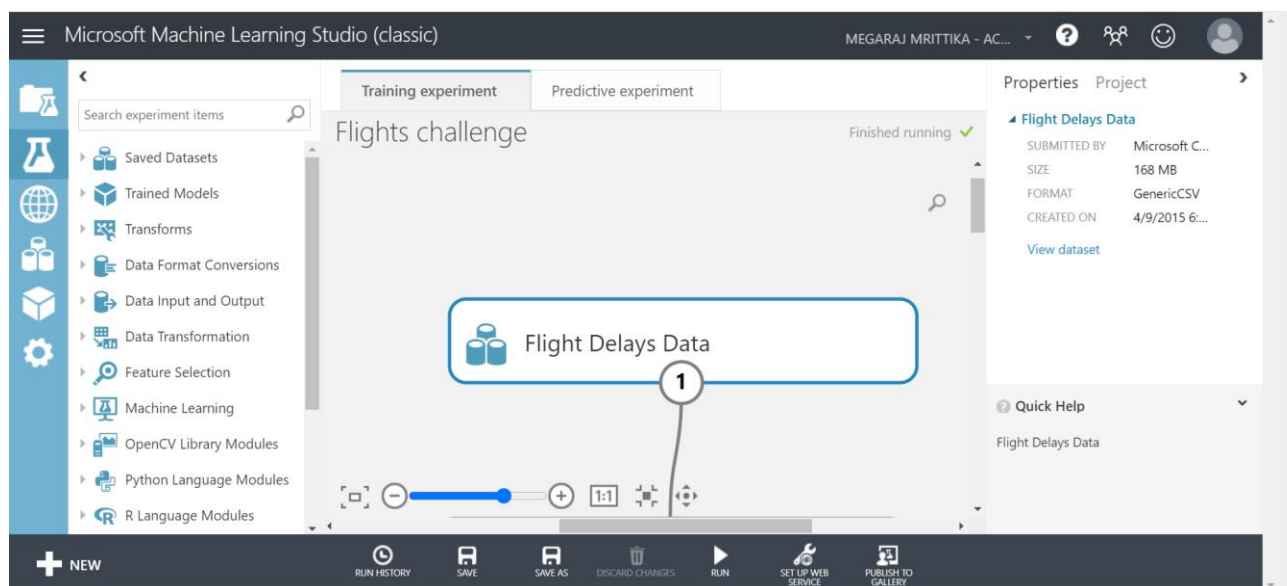
This comprehensive explanation highlights the key aspects of the regression model used in the Flight Delay Prediction experiment, covering its algorithm, configuration, training process, evaluation metrics, and results.

## Create an Experiment with existing data storage mechanisms.

**Task 1:** Sign in to Azure Machine Learning Workspace consider the need to standardize data definitions and sources within the Azure Machine Learning Workspace.

**Task 2:** Create a new experiment, with an appropriate name like "Flights Challenge".

**Task 3:** Add the Flights Delay Data sample dataset to the experiment, and then visualize its contents.

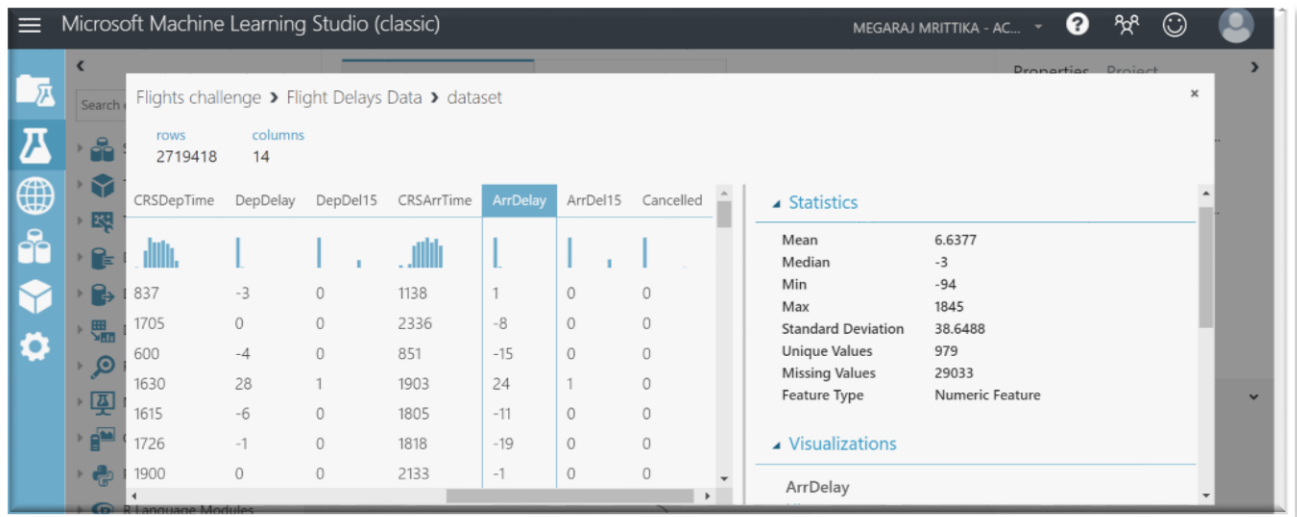**Task 4:** Answer the following questions: -

**How many rows are in the dataset?**

The number of rows in the dataset is **2719418** and the number of columns in the dataset is **14**
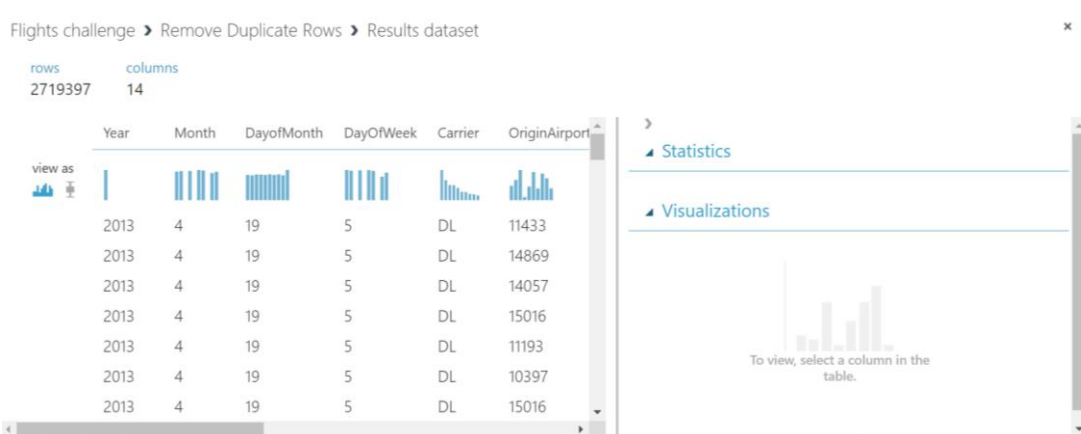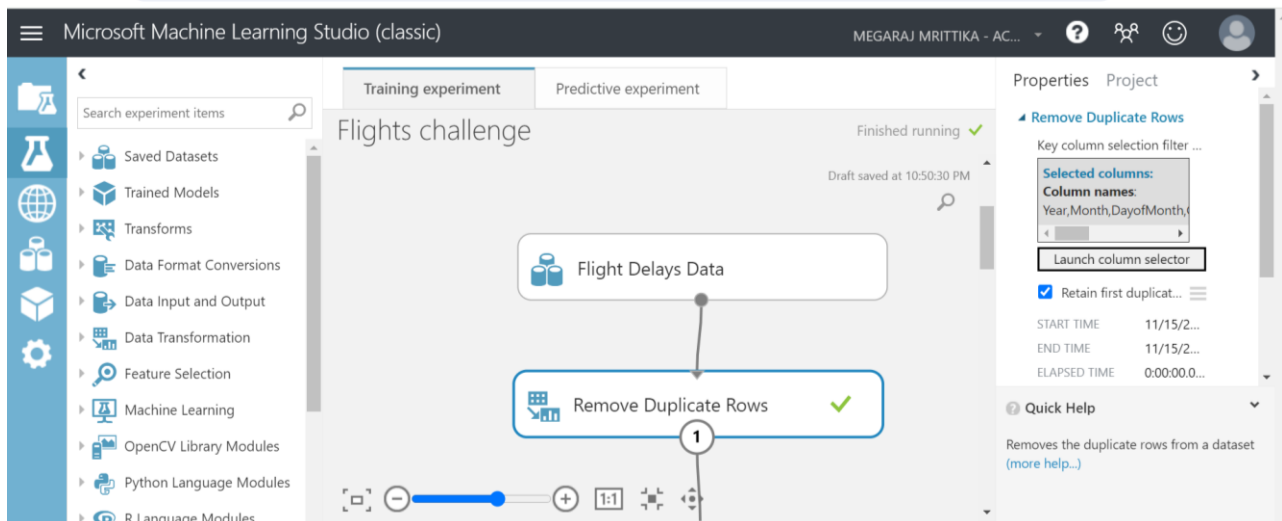


**What is the mean value of the ArrDelay column?**

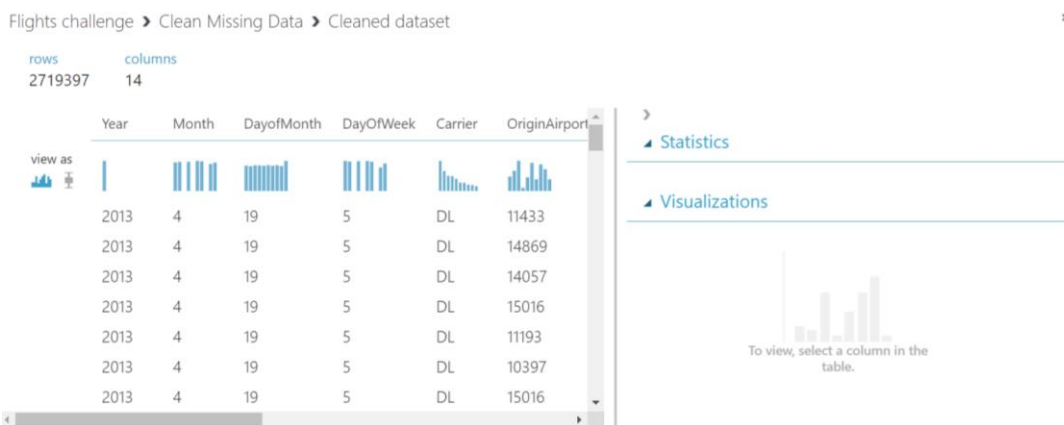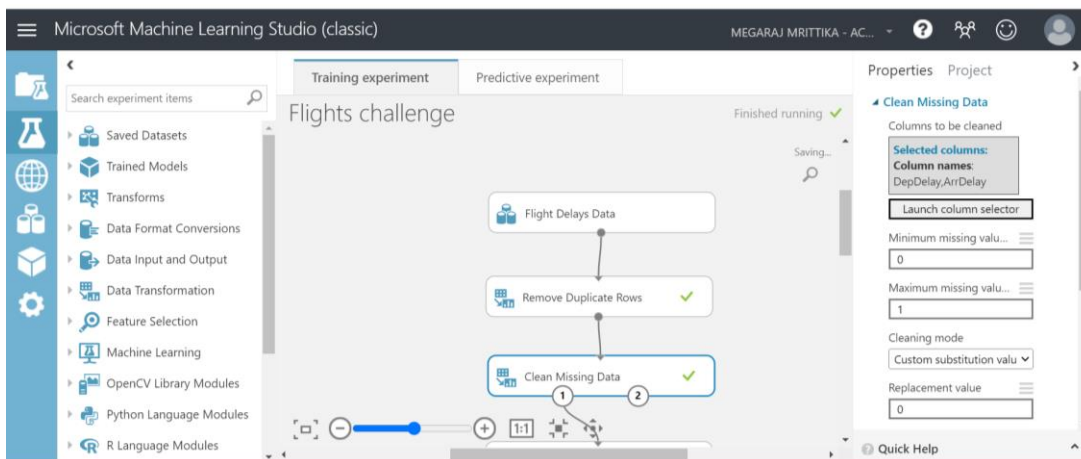The Mean value of ArrDelay column is **6.6377**

# Remove Duplicates and Replace Missing Values [ ETL process]

**Task 1:** Remove duplicate rows (retaining the first instance of each row). Rows are considered duplicates in this dataset if they have matching values for all the following fields: Year, Month , DayofMonth, Carrier, Origin Airport ID,  Dest Airport ID, CRS Dep Time,CRS Arr Time Use the built-in Azure Machine Learning module
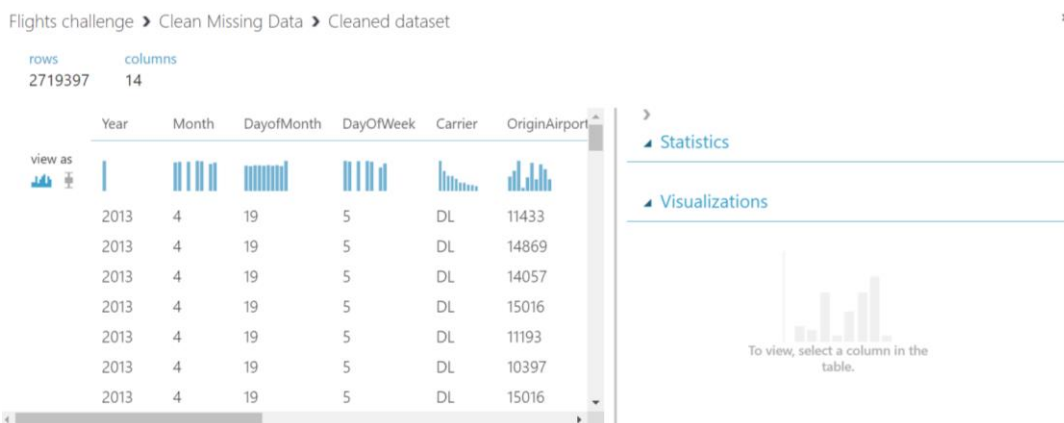
**Task 2:** After removing the duplicate rows, replace missing values in the DepDelay and ArrDelay columns with the value 0 (zero). Use the built-in Azure Machine Learning module



**Task 3:** Answer the following questions, after you have removed duplicate rows and replaced missing values
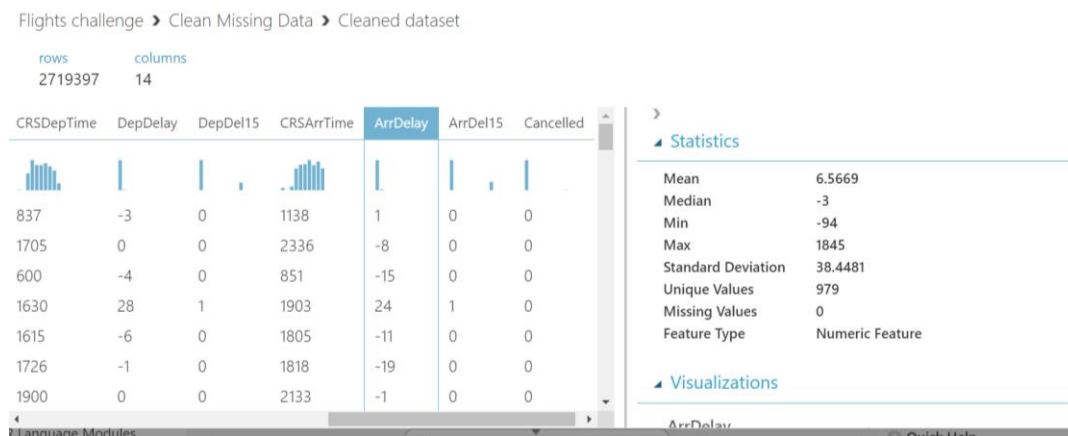
**How many rows remain in the dataset?**

The number of rows remaining in the dataset is **2719397**



**What is the mean value of the ArrDelay column?**

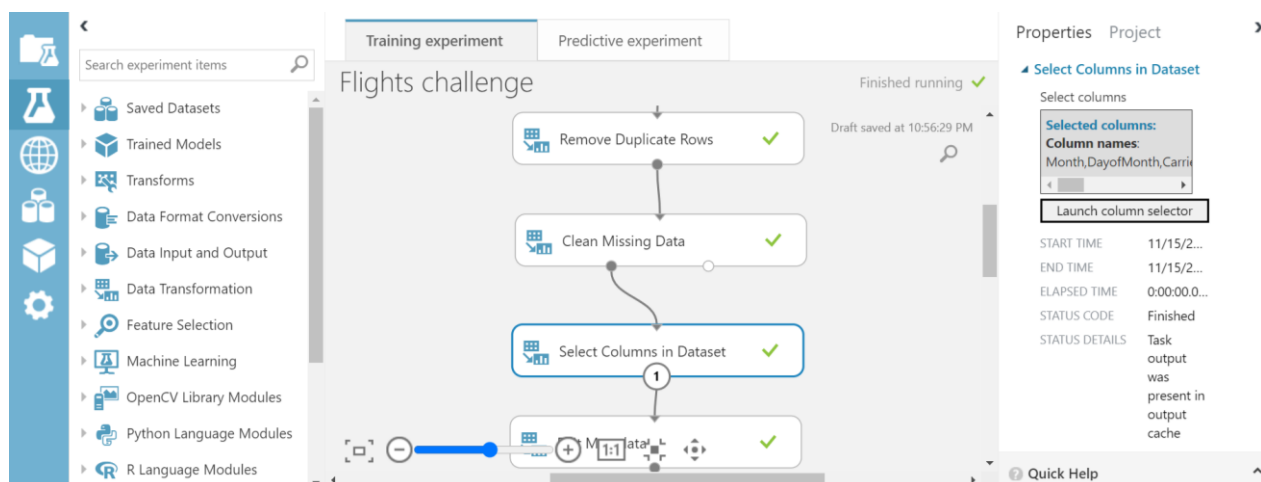The mean value of the ArrDelay column is **6.5669**



## Train a Regression Model

To predict a numeric value, such as the number of minutes delayed or early a flight arrives, train a regression model. Perform the following tasks to train a regression model

**Task 1**: Return to the Azure Machine Learning experiment you created in Part 1.

**Task 2:** Add a Select Columns in the Dataset module, and use it to select only the Month, DayofMonth, DayOfWeek, Carrier, OriginAirportID, DestAirportID, CRSDepTime, DepDelay, CRSArrTime, andArrDelay columns.



**Task 3:** Add an Edit Metadata module and use it to make the OriginAirportID, DestAirportID, and Carrier columns Categorical.

**Task 4:** Add a Normalize Data module and use it to standardize the CRSDepTime, CRSArrTime, and DepDelay columns using the ZScoretransformation method.



**Task 5:** Add a Split Data module and use it to split the rows into 70% / 30% subsets. Use a random seed value of 0.

**Train dataset:**



Flights challenge > Split Data > Results dataset1

| rows | columns |
|---|---|
| 1903578 | 10 |

| Month | DayofMonth | DayOfWeek | Carrier | OriginAirportID | DestA |
|---|---|---|---|---|---|
| 7 | 24 | 3 | US | 14107 | 14570 |
| 5 | 16 | 4 | DL | 11433 | 10721 |
| 5 | 28 | 2 | FL | 13204 | 11433 |
| 6 | 3 | 1 | WN | 15016 | 10821 |
| 7 | 23 | 2 | DL | 11278 | 11433 |
| 4 | 25 | 4 | UA | 14771 | 12478 |
| 6 | 25 | 2 | WN | 14831 | 12889 |

◢ Statistics

◢ Visualizations

To view, select a column in the table.

**Test Dataset:**



Flights challenge > Split Data > Results dataset2

| rows | columns |
|---|---|
| 815819 | 10 |

| Month | DayofMonth | DayOfWeek | Carrier | OriginAirportID | DestA |
|---|---|---|---|---|---|
| 9 | 8 | 7 | AS | 10721 | 14747 |
| 8 | 30 | 5 | DL | 14771 | 10397 |
| 6 | 10 | 1 | EV | 13198 | 11042 |
| 8 | 18 | 7 | DL | 11697 | 12953 |
| 10 | 23 | 3 | UA | 12892 | 14771 |
| 4 | 14 | 7 | DL | 10397 | 14107 |
| 9 | 3 | 2 | UA | 14771 | 12264 |

◢ Statistics

◢ Visualizations

To view, select a column in the table.

**Task 6:** Add a Boosted Decision Tree Regression module and a Train Model module. Then use the default settings to train the model with the 70% data split to predict the ArrDelay label column.
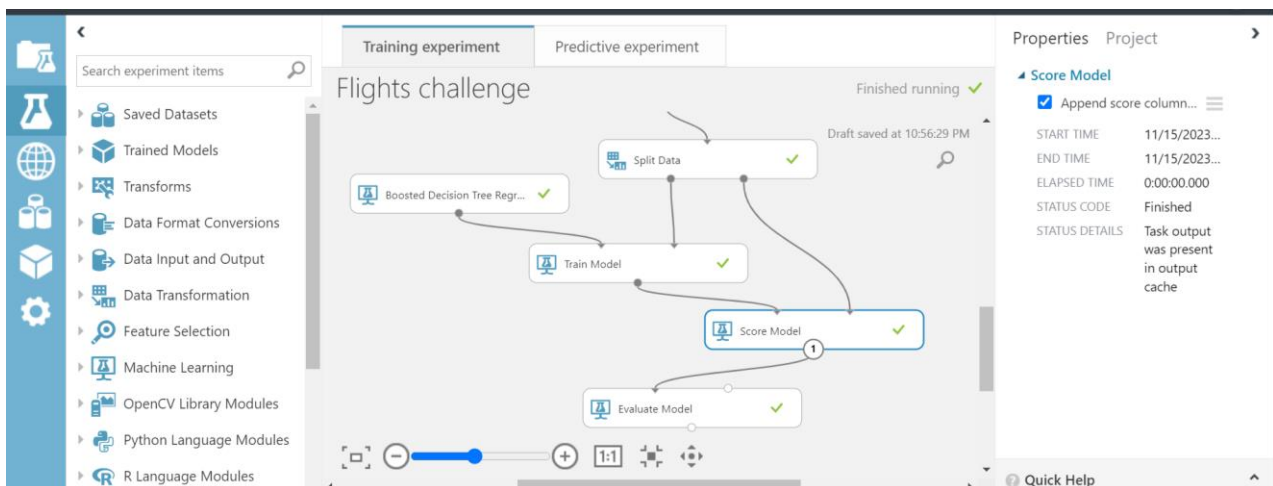
**Boosted Decision Tree Regression module:**

**Train Model Module:**



**Task 7:** Add a Score Model module, and use it to score the trained model using the 30% split of data.

rows 815819  columns 11

| DestAirportID | CRSDepTime | DepDelay | CRSArrTime | ArrDelay | Scored Labels |
|---|---|---|---|---|---|
| 14747 | -1.117256 | -0.318114 | -0.786029 | -15 | -8.918402 |
| 10397 | -0.247461 | 3.049315 | 0.906395 | 123 | 117.610634 |
| 11042 | -1.265758 | 0.154996 | -0.962154 | 24 | 12.127501 |
| 12953 | -1.318794 | 2.103095 | -1.016814 | 67 | 84.300438 |
| 14771 | -1.424867 | -0.540754 | -1.397407 | 0 | -13.408275 |
| 14107 | -1.106649 | -0.401604 | -1.176744 | -27 | -13.241999 |

**Statistics**

| | |
|---|---|
| Mean | 6.6024 |
| Median | -4.4085 |
| Min | -32.9918 |
| Max | 1180.4 |
| Standard Deviation | 36.357 |
| Unique Values | 769006 |
| Missing Values | 0 |
| Feature Type | Numeric Score |

**Visualizations**

**Task 8:** Add an Evaluate Model module and use it to evaluate the results from the Score Model module.

Training experiment | Predictive experiment

Flights challenge   Finished running ✓

Draft saved at 10:56:29 PM

Search experiment items

- Saved Datasets
- Trained Models
- Transforms
- Data Format Conversions
- Data Input and Output
- Data Transformation
- Feature Selection
- Machine Learning
- OpenCV Library Modules
- Python Language Modules
- R Language Modules

Boosted Decision Tree Regr... ✓

Split Data ✓

Train Model ✓

Score Model ✓

Evaluate Model ✓

**Properties**  Project

**Evaluate Model**

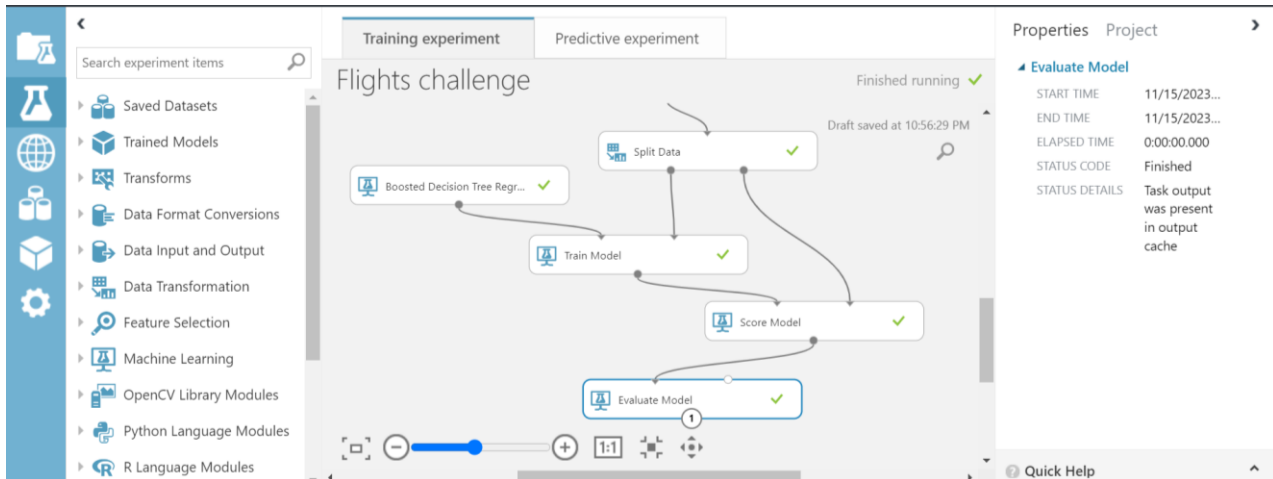| | |
|---|---|
| START TIME | 11/15/2023... |
| END TIME | 11/15/2023... |
| ELAPSED TIME | 0:00:00.000 |
| STATUS CODE | Finished |
| STATUS DETAILS | Task output was present in output cache |

Quick Help

# Test and Evaluate the Model

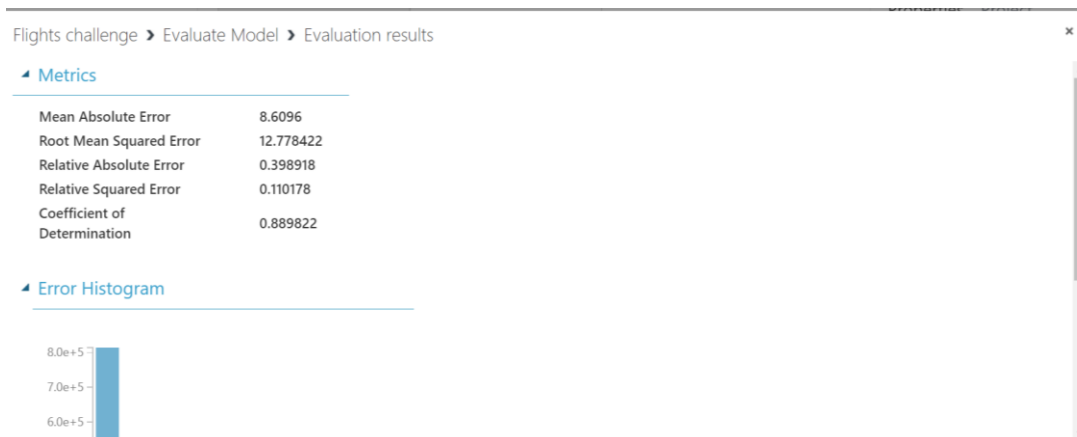Perform the following tasks

**Task 1:** Run the experiment

**Task 2**: When it has finished, visualize the output of the Evaluate Model module.
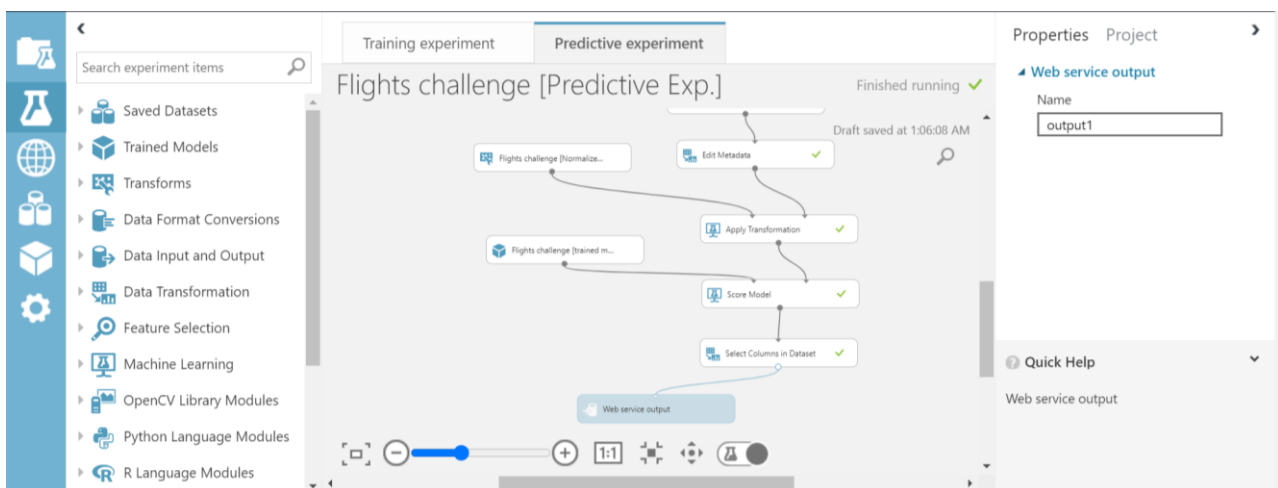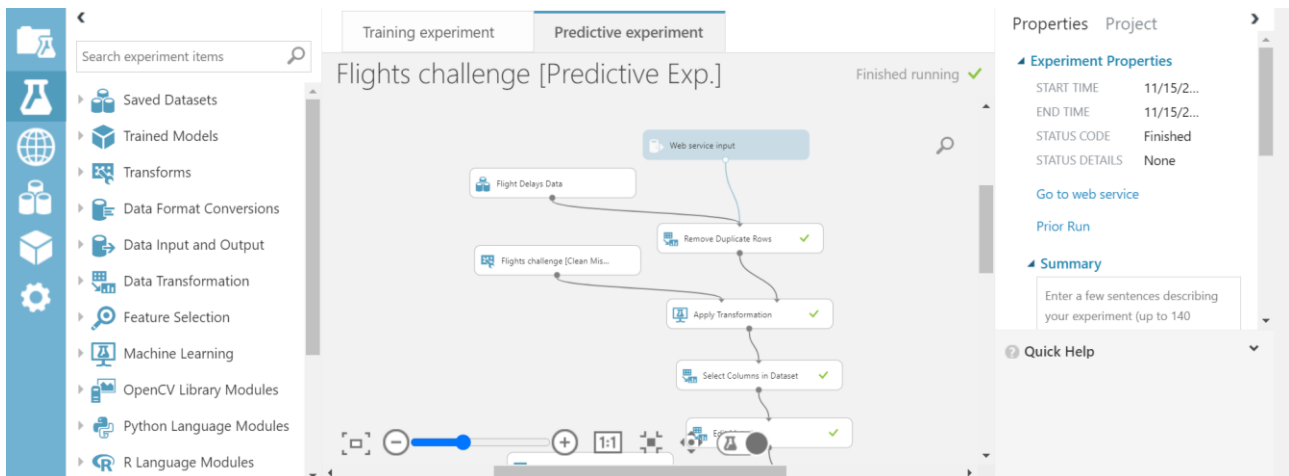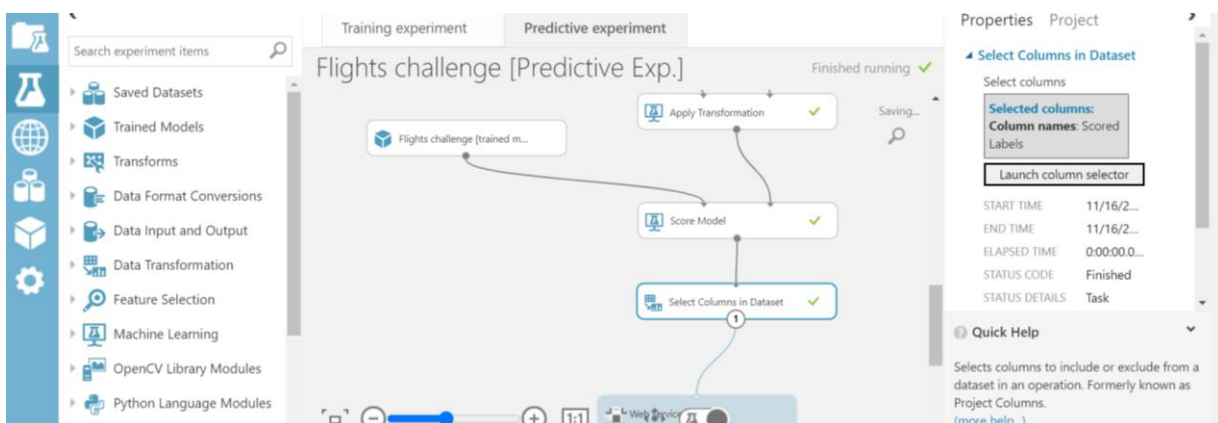


**Evaluate result:**



| Metrics | |
| --- | --- |
| Mean Absolute Error | 8.6096 |
| Root Mean Squared Error | 12.778422 |
| Relative Absolute Error | 0.398918 |
| Relative Squared Error | 0.110178 |
| Coefficient of Determination | 0.889822 |

**Error Histogram:**

## Publish and Use the Model

**Task 1:** Set up the experiment as a web service, creating a predictive experiment (if the option to do this is not available, save and re-run the experiment).





**Task 2:** In the predictive experiment, add a Select Columns in the Dataset module and place it between the Score Model and Web service output modules. Use this to select only the Scored Labels column.

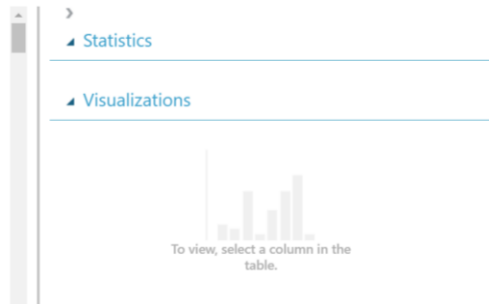Flights challenge [Predictive Exp.] ➤ Select Columns in Dataset ➤ Results dataset                    ✕

rows        columns
2719397     1

Scored Labels                                    ➤
                                                 ◢ Statistics
view as
◩ ⟟ |
                                                 ◢ Visualizations
        -7.463635
        -4.074299
        -7.719961
        24.673752
        -9.544199                                To view, select a column in the
        -8.366618                                          table.
        -0.930374

**Task 3**: Save and run the modified predictive experiment, and then deploy the web service.



**Task 4:** Enter the values predicted by your model for each input row