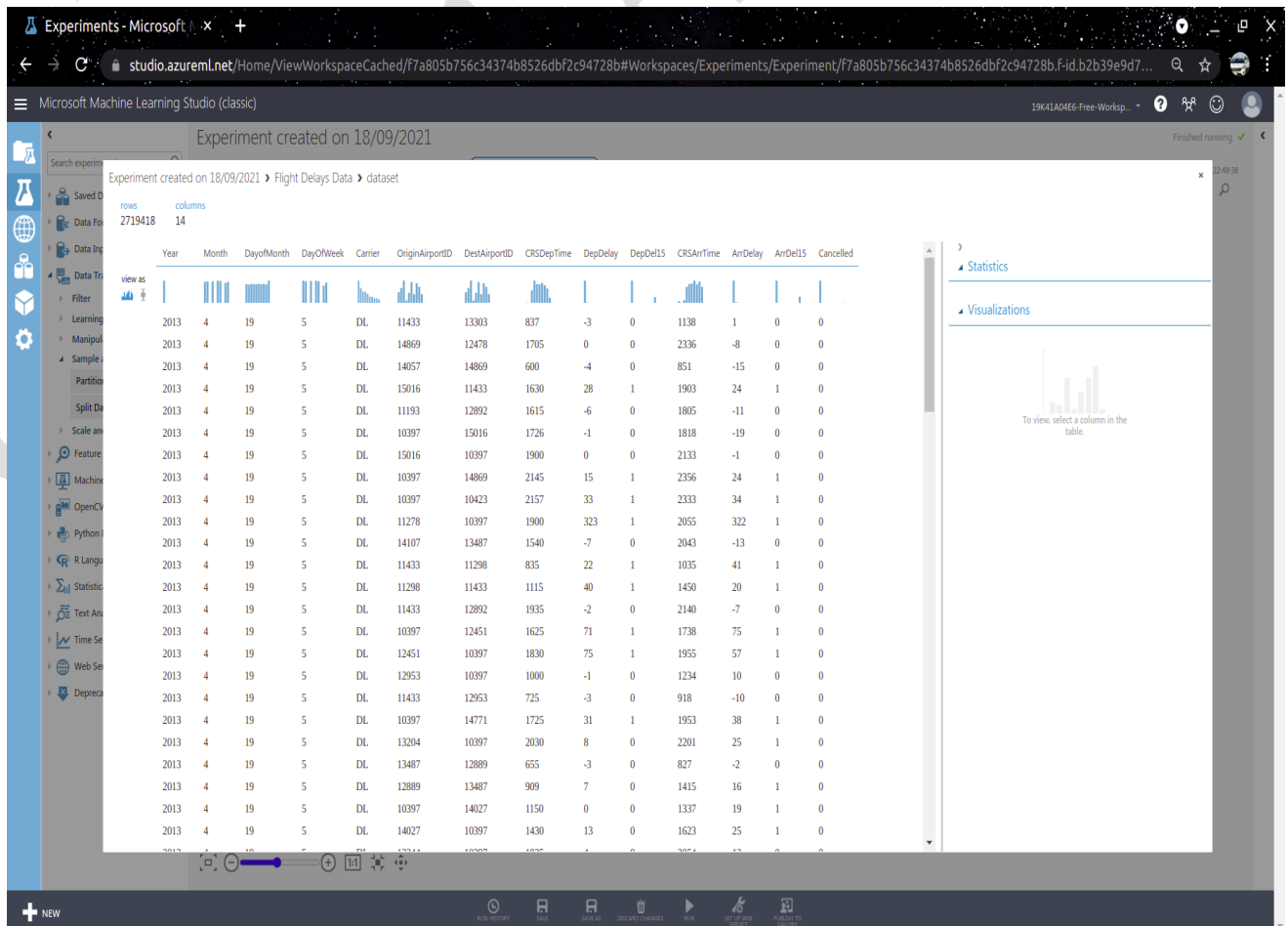# Artificial Intelligence

## Assignment 4

**Software Used:-** Microsoft Azure Machine Learning Studio

**About the assignment:-** The model used in this assignment is Two Class Logistic Regression Model. The trained model finds the probability of a flight to be delayed by more than 15 minutes based on certain parameters. The steps followed in the workflow are:

1. Load the data
2. Study the data and find out the necessary variables
3. Pre Process the data
4. Choose the model to be used (Two Class Logistic Regression Model in this case)
5. Split the data for training and testing
6. Train the model with the training data
7. Score the model using the testing data
8. Evaluate the model based on the results

**Workflow:-**

- Study the data

- Remove all the cancelled flights from the dataset



- Drop the columns that are not needed for the training and testing

- Remove rows with any NULL or missing values



- The dataset before splitting for training and testing is shown below

- The data is split 50-50 for training and testing



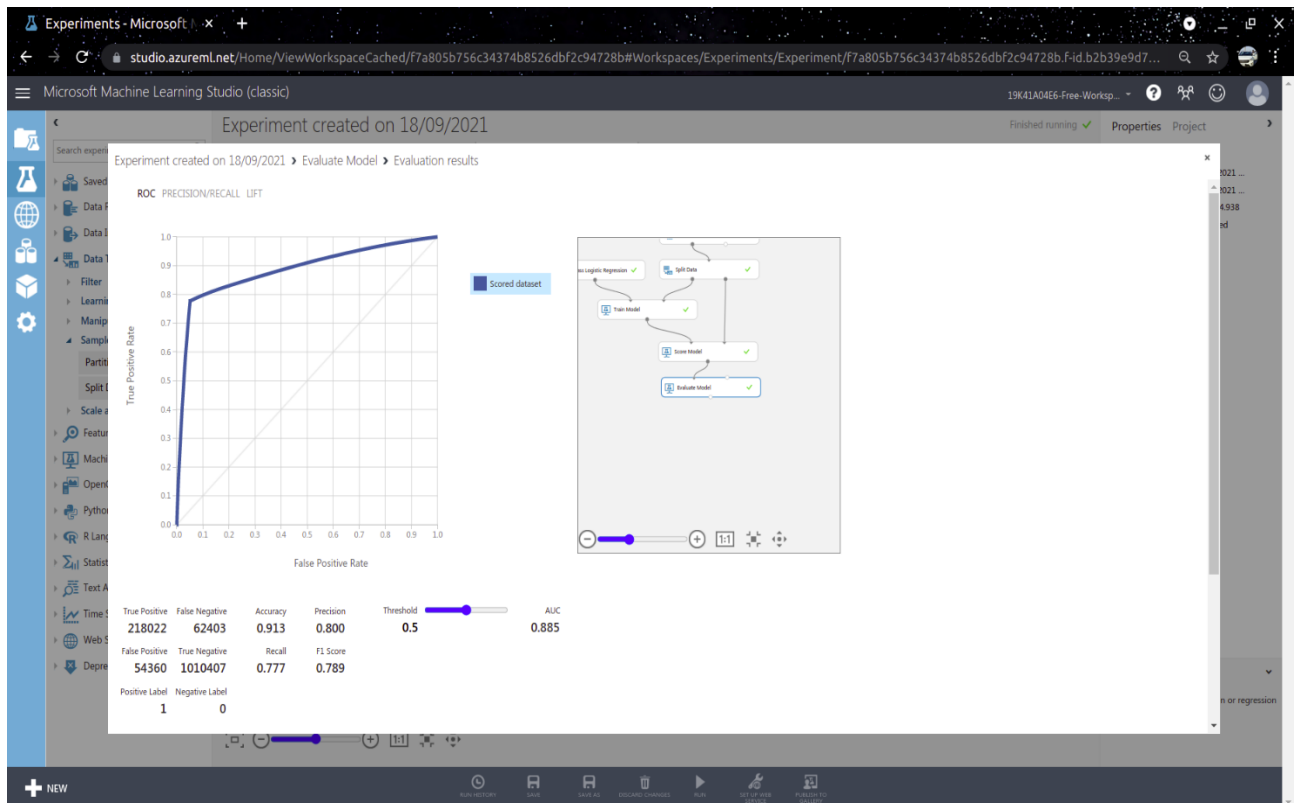- The model is initialized

- Train the model using the training data



- The scored dataset is shown below

- The evaluation of model is done and the results are as follows



**Conclusion:** From the evaluation, we can see that the accuracy of the model is 91.3% and the model is able to predict the outcome properly in most cases. The parameters used to predict the delays are:

1. Year of flight
2. Month of flight
3. Date of flight
4. Day of the week of the flight
5. The Carrier/Airline
6. Origin Airport
7. Destination Airport
8. Scheduled Departure Time
9. Departure Delay (atleast 15 minutes)
10. Scheduled Arrival Time
11. Arrival Delay (atleast 15 minutes)

This means that on giving the first 10 parameters as input, the model can predicts the probability of the flight getting delayed by more than 15 minutes for arriving at the destination airport.

The final workflow is as follows: