



# Two Class Logistic Regression with Azure ML



# Contents

## Table of Contents

Problem Statement .....	4
Dataset for Loan Approval .....	5
Cleaning the Data .....	6
Cleaning Numeric Data .....	9
Select Columns in Dataset .....	11
Split Data .....	12
Usage of Two-Class Logistic Regression.....	15
Score and Evaluate Model .....	17
Understanding the Output .....	23
Impact Analysis & Stratification.....	29

# Goals and Requirements

**Estimated time to complete lab is 40-45 minutes**

## Goals

1. Implement and design a model for automating loan approval.
2. Approach of using Two Class Logistic Regression

## Requirements:

1. Access to an Azure Machine Learning Studio

# Logistic Regression

## Problem Statement

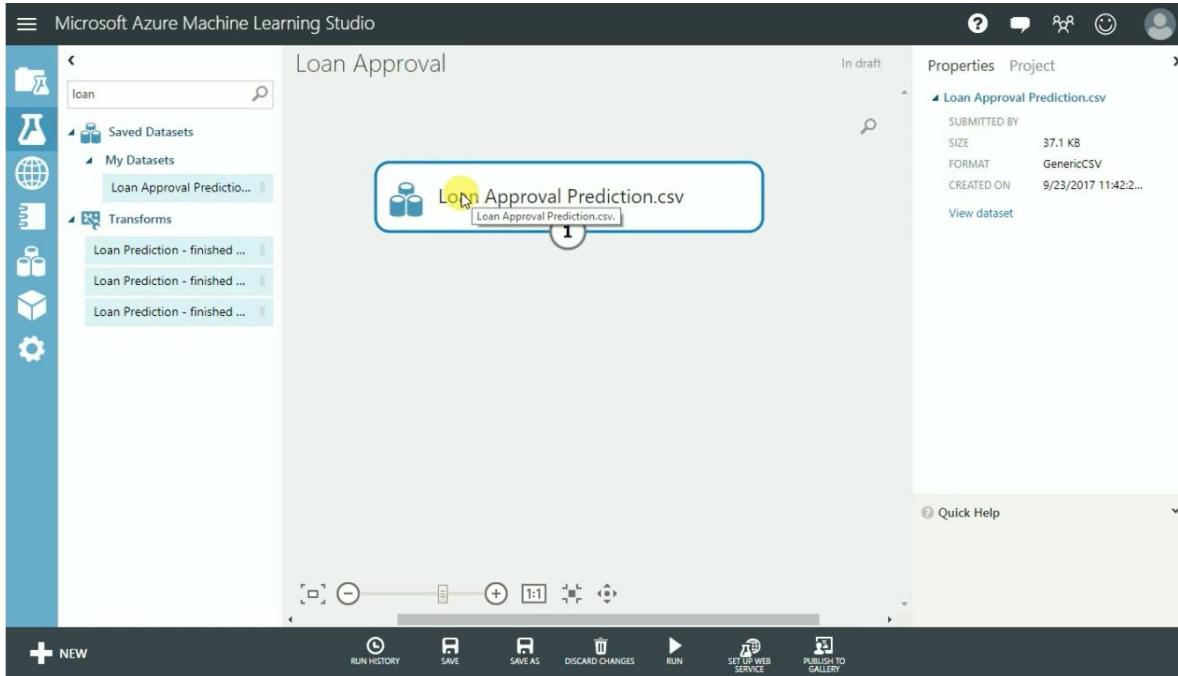
### Problem Statement

- Automate loan eligibility process
- Identify customers whose loan will be approved

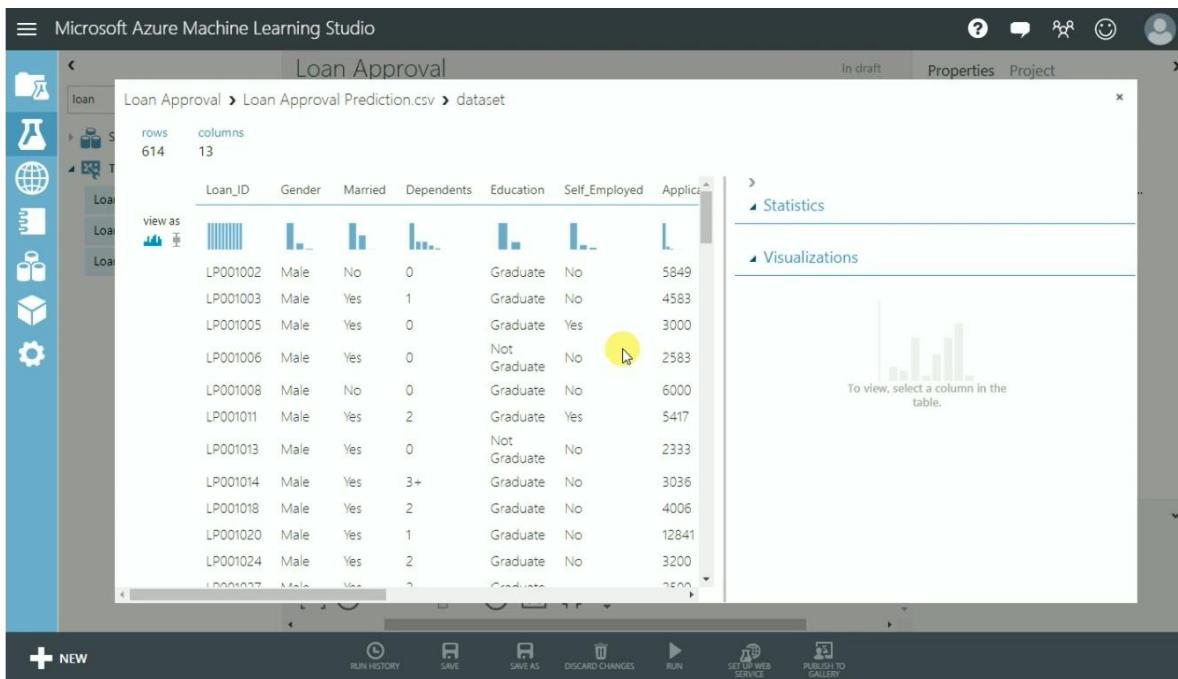
Loan_ID	Gender	Married	Dependents	Self_Employed	Income	LoanAmt	Term	CreditHistory	Property_Area	Status
LP001002	Male	No	0	No	\$5,849.00		60	1	Urban	Y
LP001003	Male	Yes	1	No	\$4,583.00	\$128.00	120	1	Rural	N
LP001005	Male	Yes	0	Yes	\$3,000.00	\$66.00	60	1	Urban	Y
LP001006	Male	Yes	2	No	\$2,583.00	\$120.00	60	1	Urban	Y

## Dataset for Loan Approval

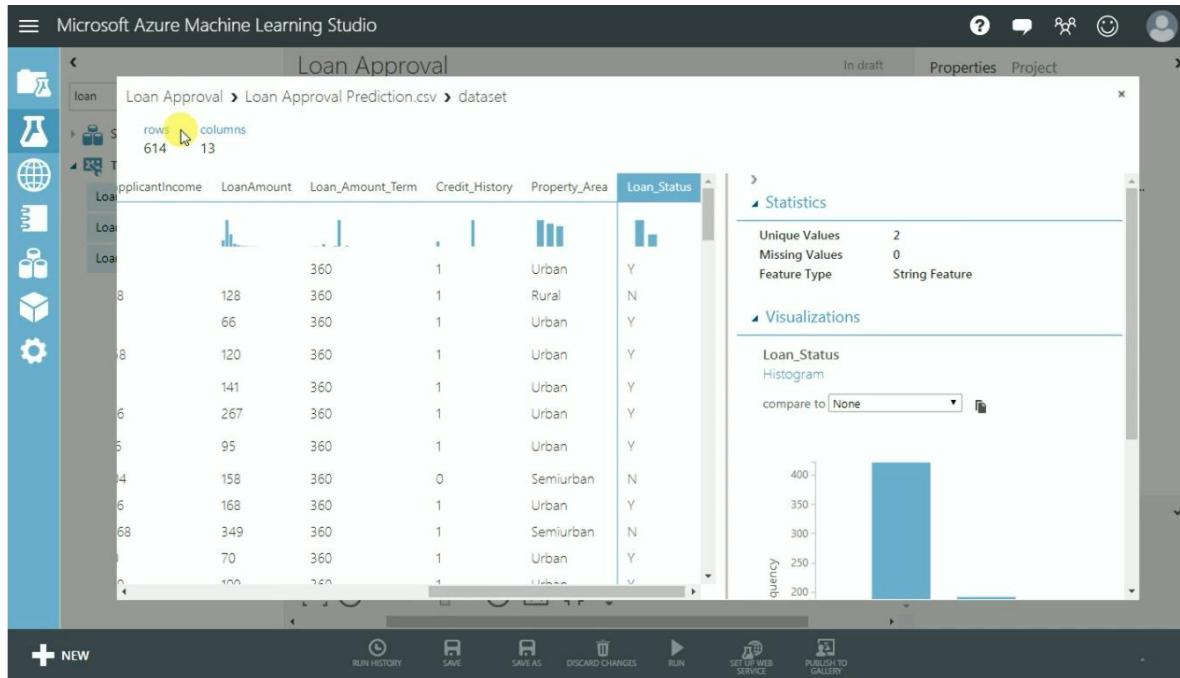
Place the dataset in canvas



Visualize the data that we need to deal which involved loan application

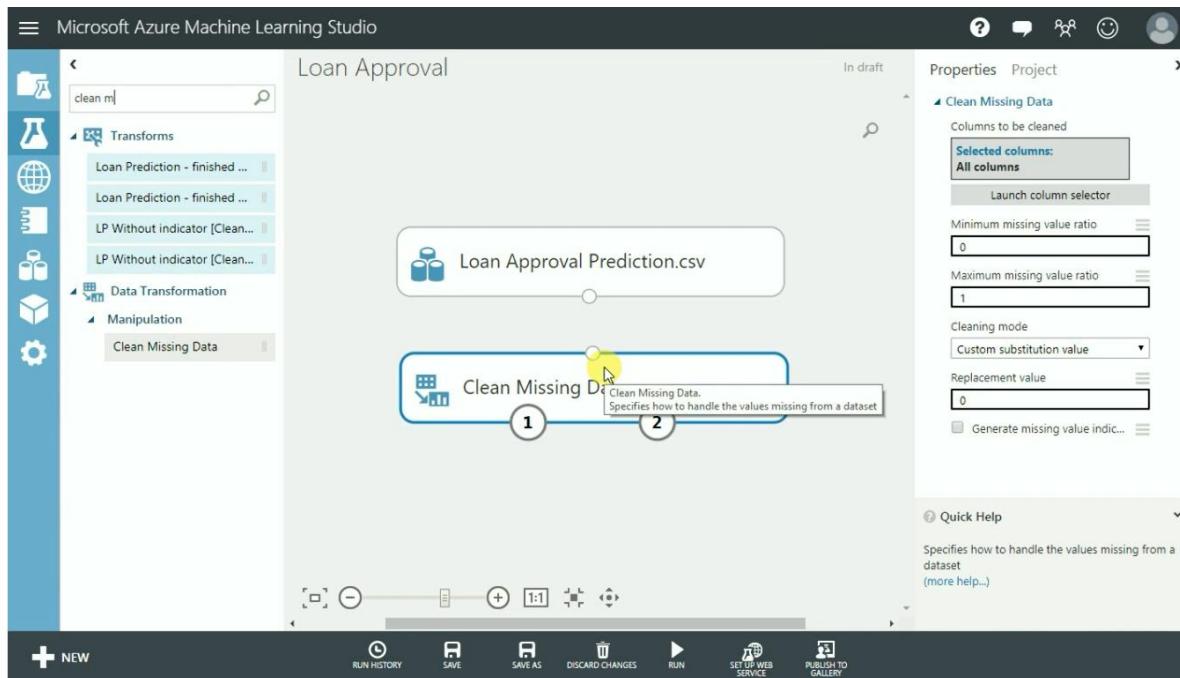


This dataset got 614 rows and 13 columns

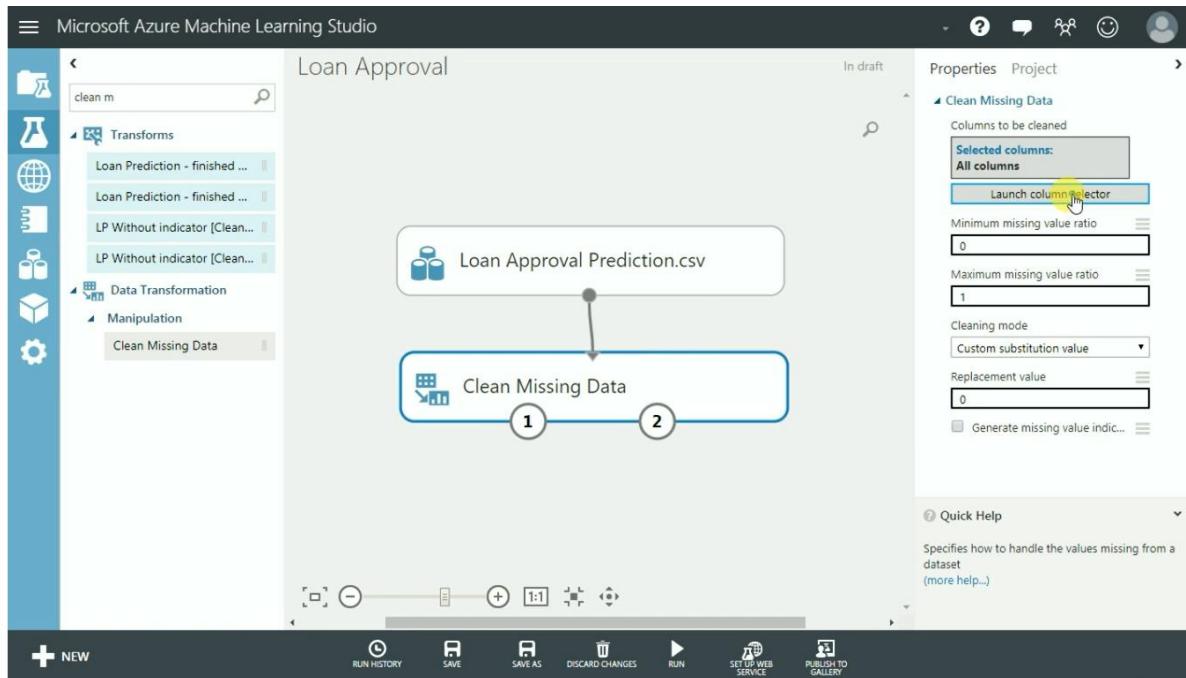


## Cleaning the Data

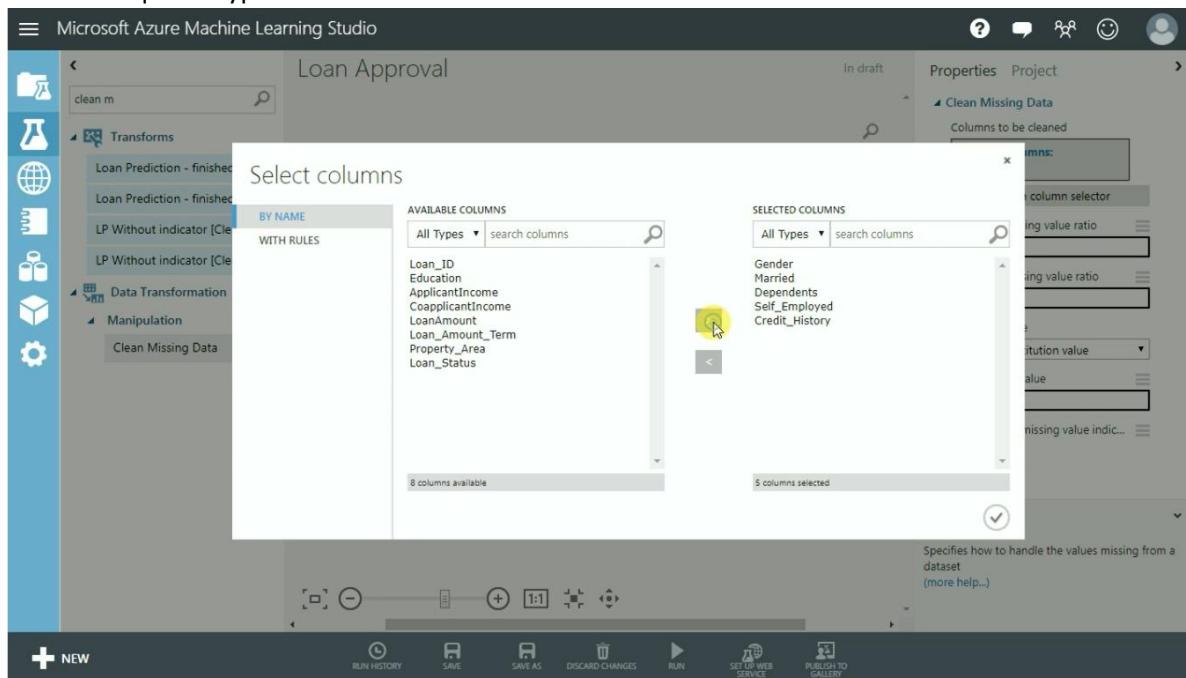
Visualize each column and check for missing values and perform clean missing data action



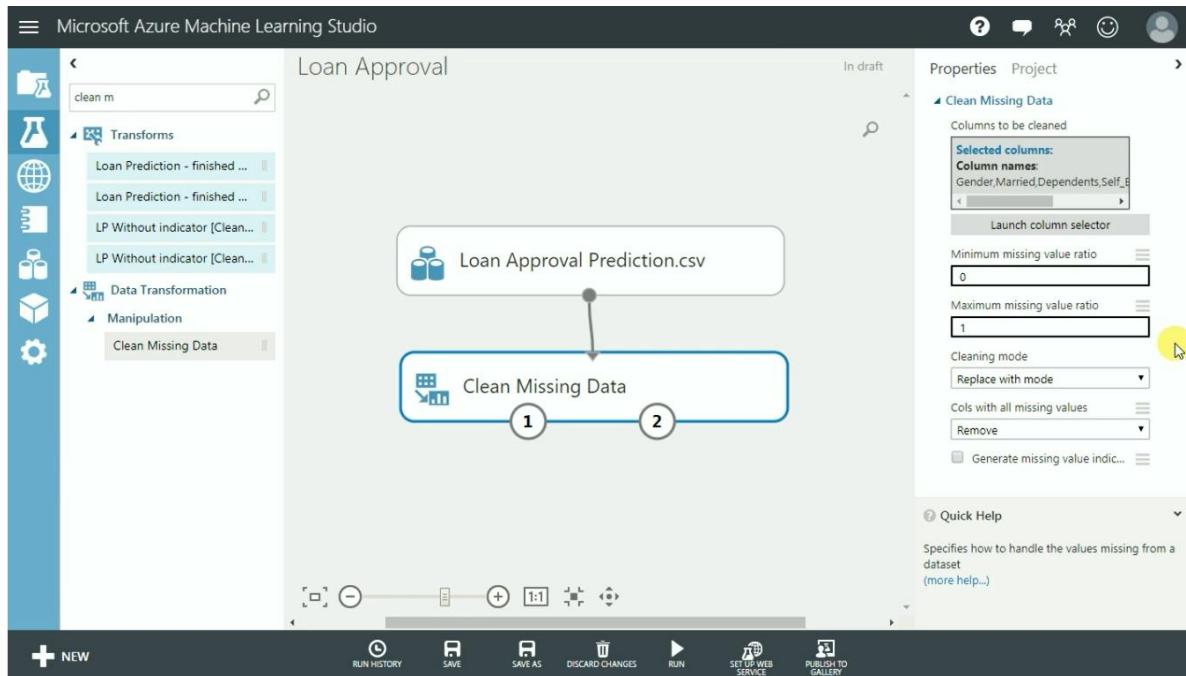
Attach both the datasets and launch column selector



Select the required types and click ok

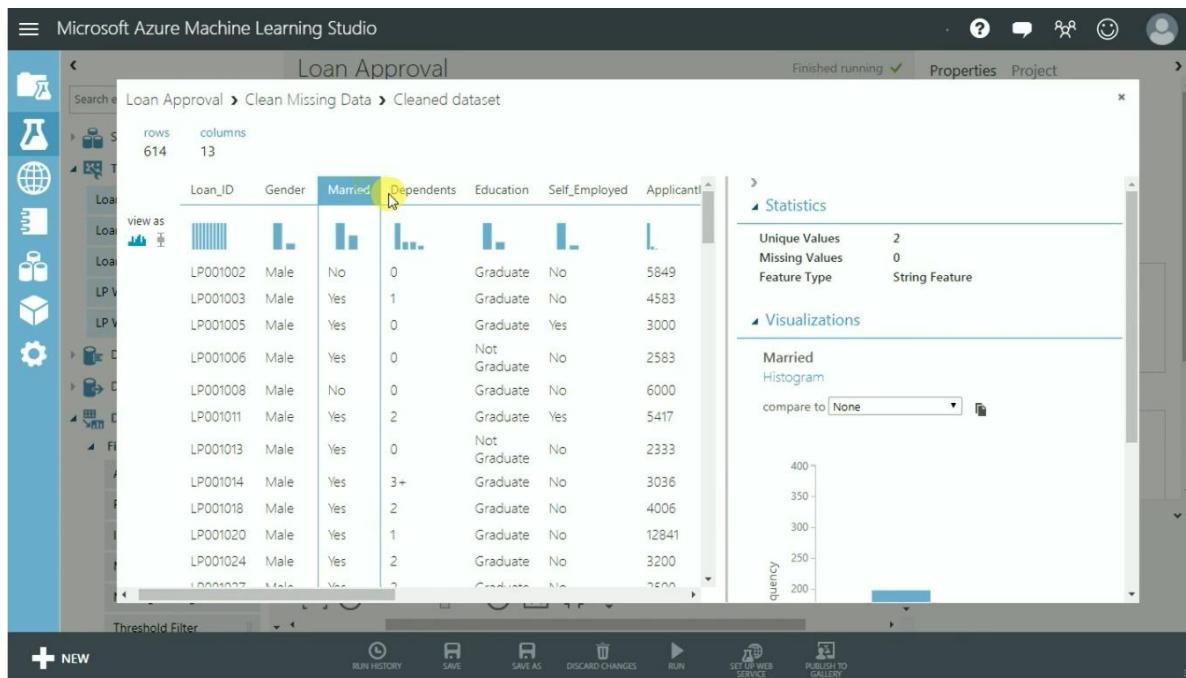


Select the required parameters to clean missing data for string features with replacement of mode



Run the module and visualize to check the result.

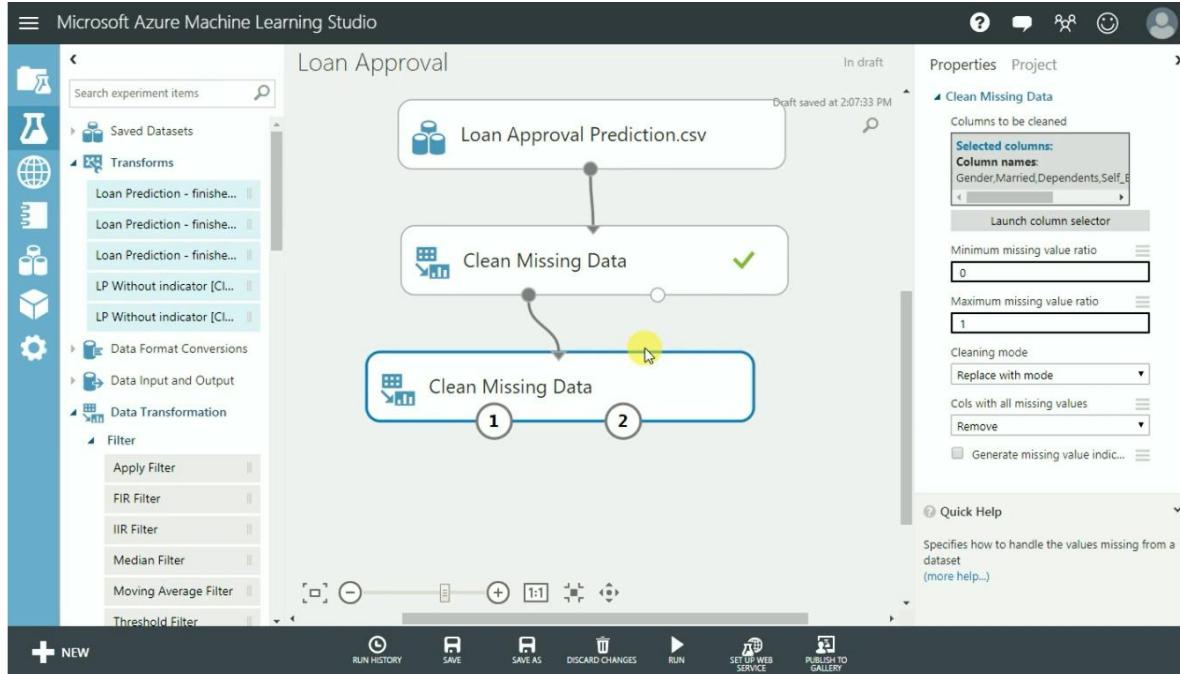
The result is successful



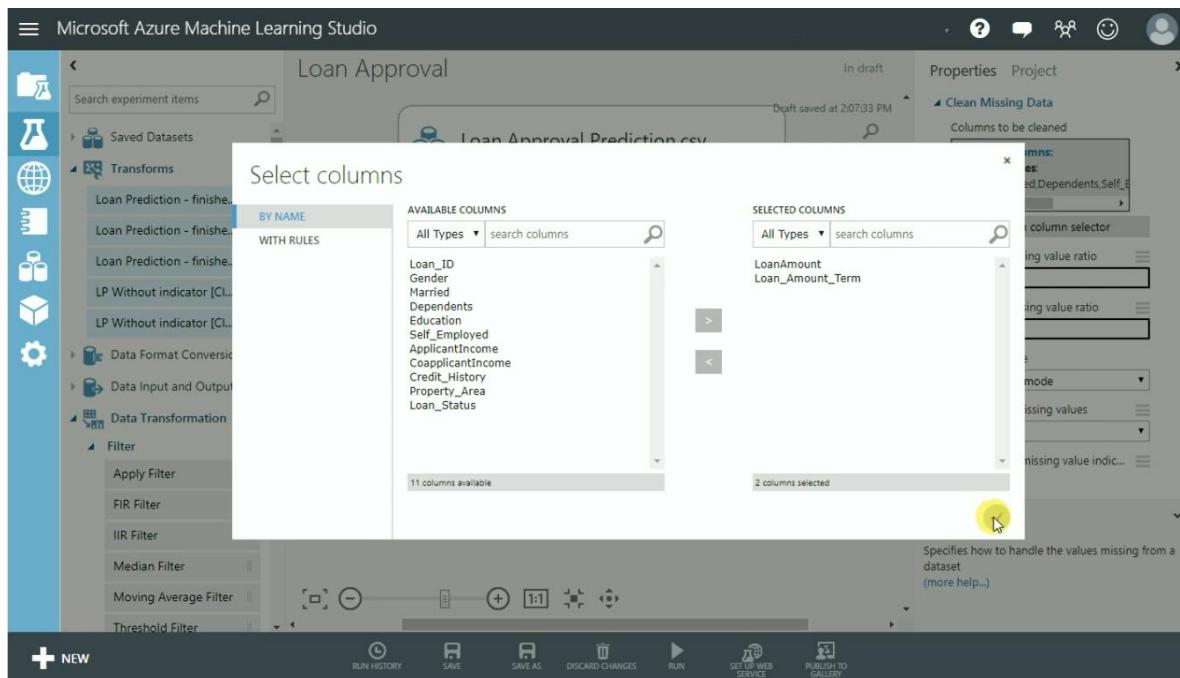
# Cleaning Numeric Data

Now experiment for numeric features for clean missing data

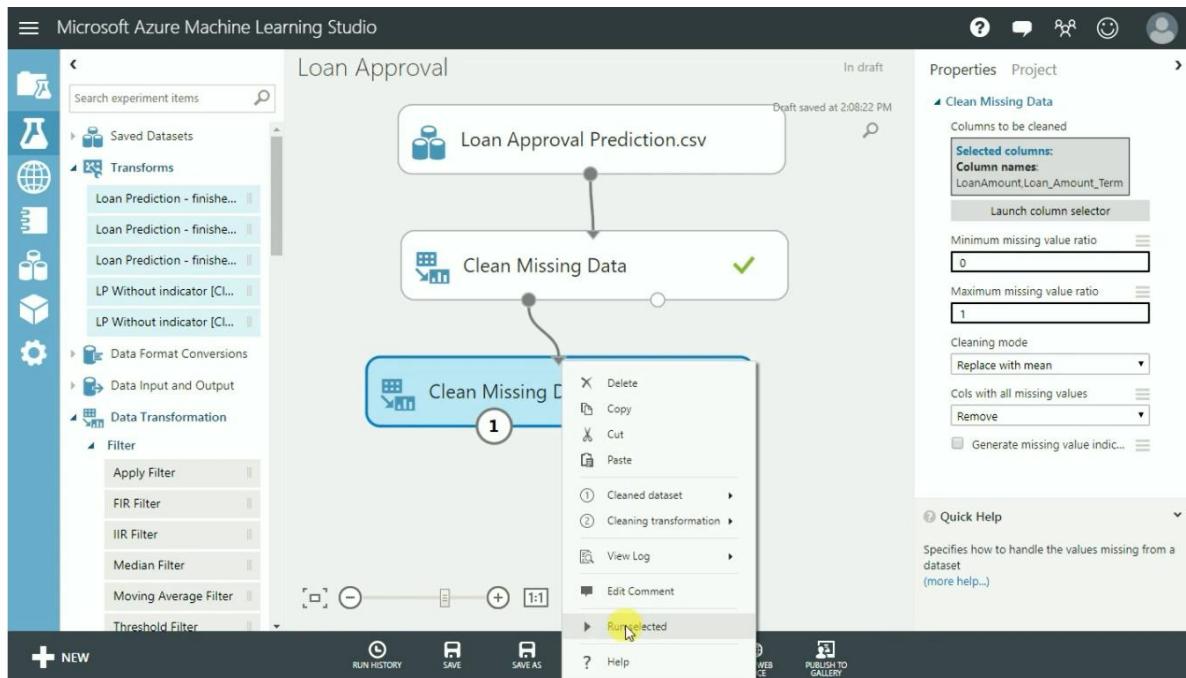
Simply Copy and paste clean missing dataset and connect the output node to input



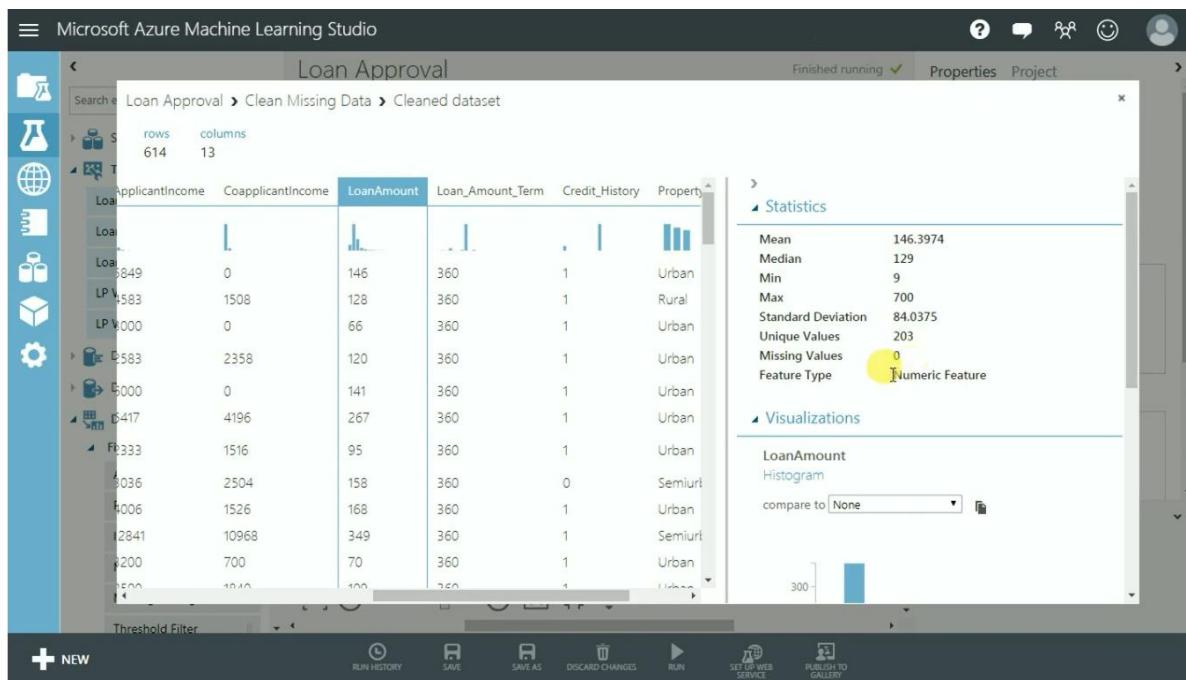
Launch column selector and select numerical features and press ok



Change the parameters cleaning mode by replace with mean and run the module

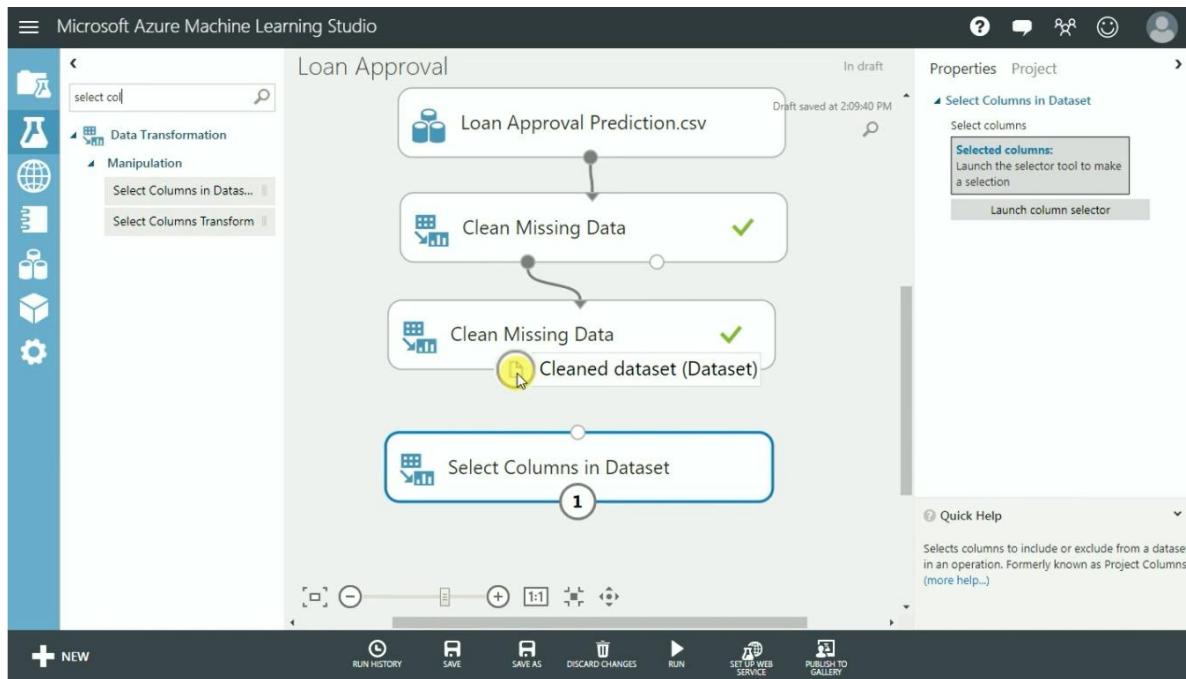


Visualize the module now and can see the numeric features showing missing values 0



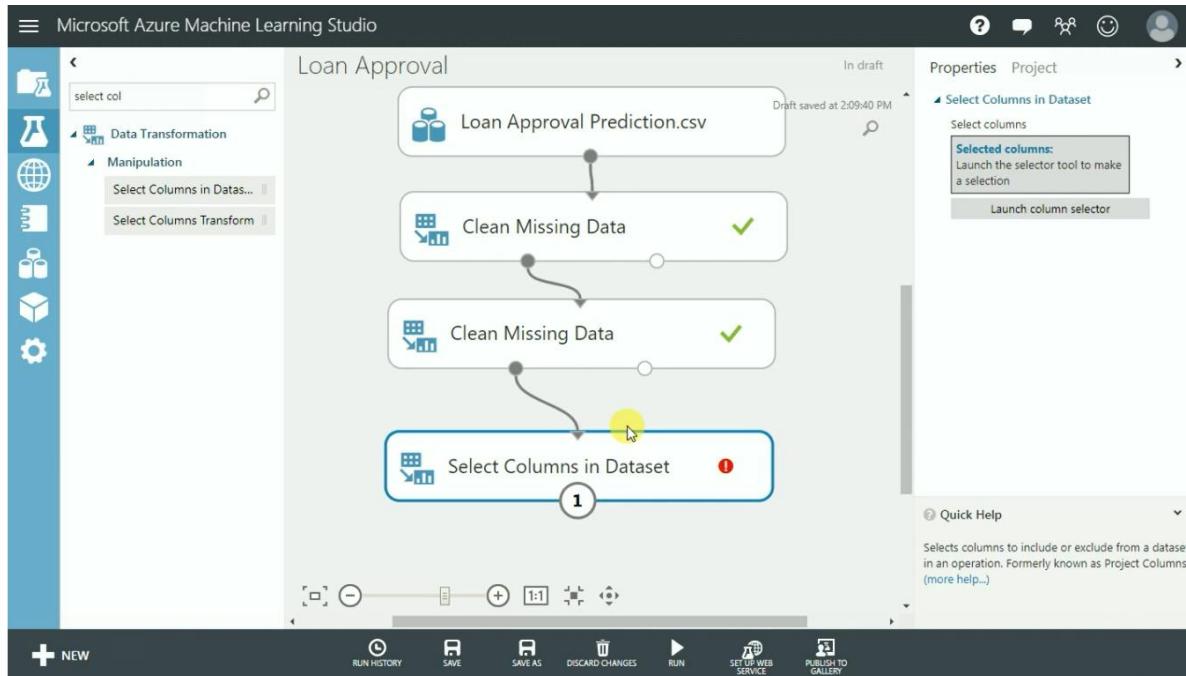
## Select Columns in Dataset

Search columns dataset and drop in canvas



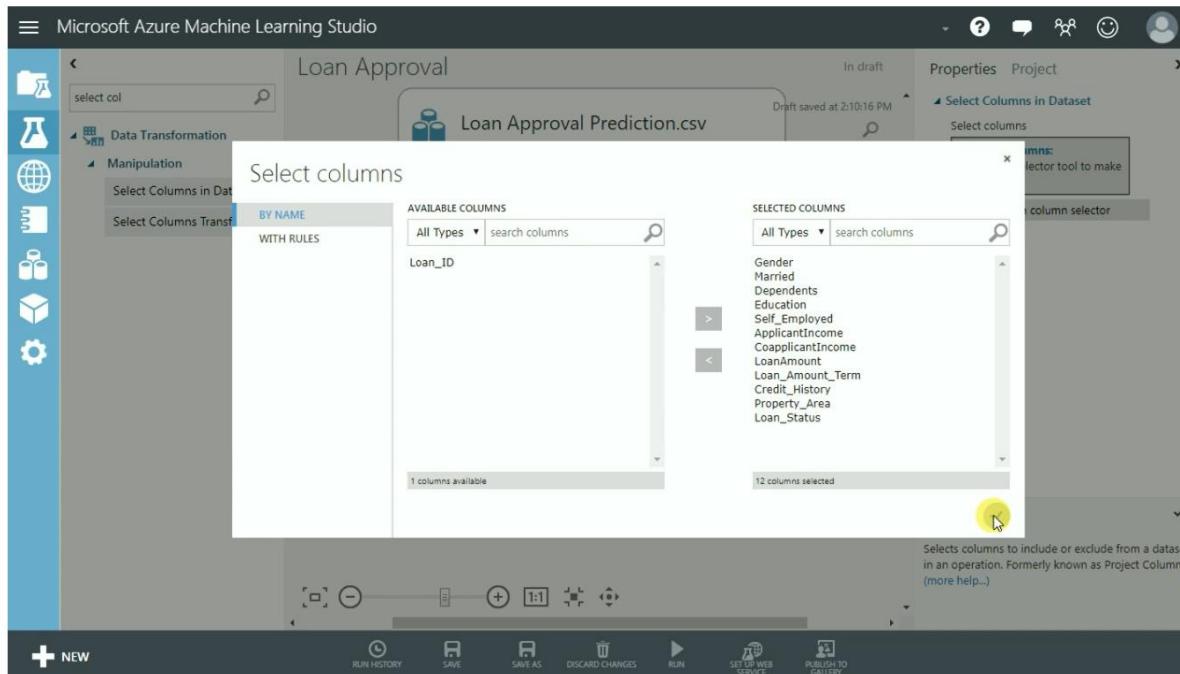
Connect the clean missing dataset output node to select columns in dataset input node

And launch column selector



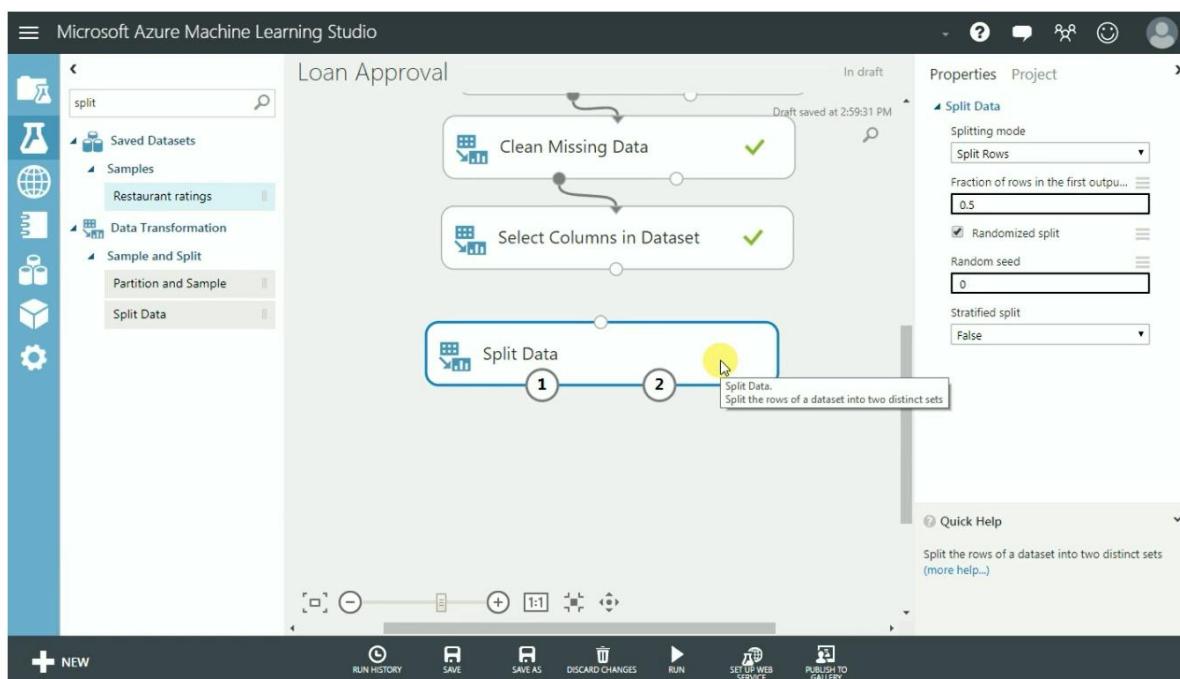
Select all the types except loan ID and press ok

Run and visualize

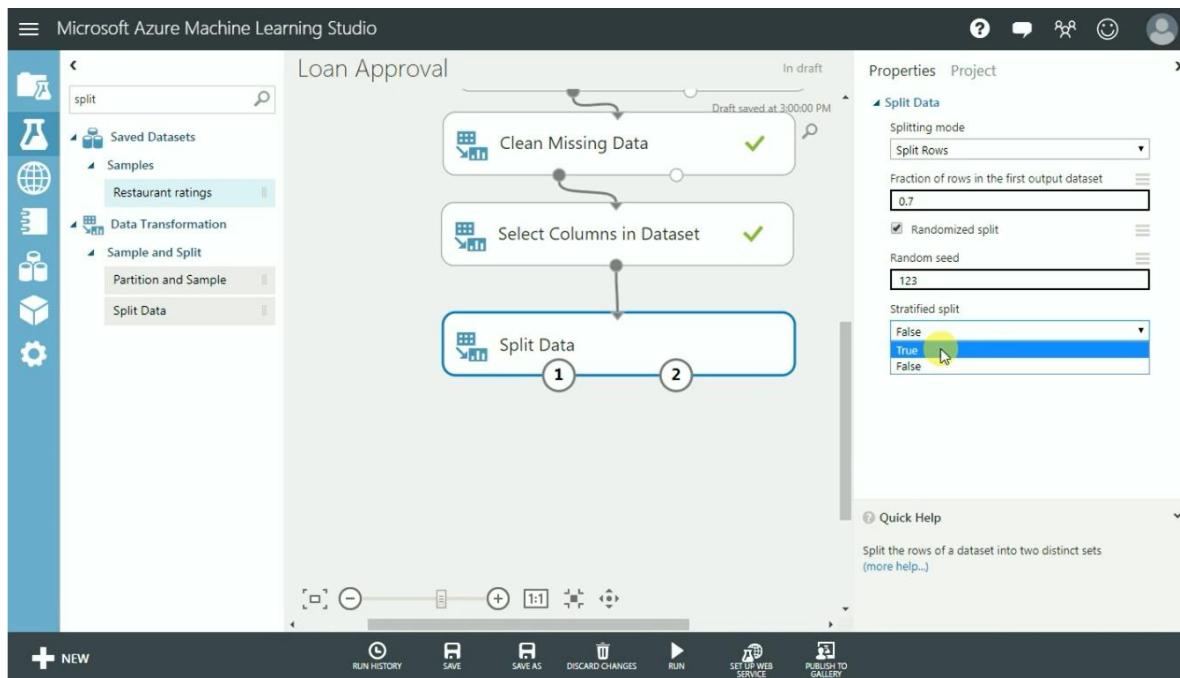


## Split Data

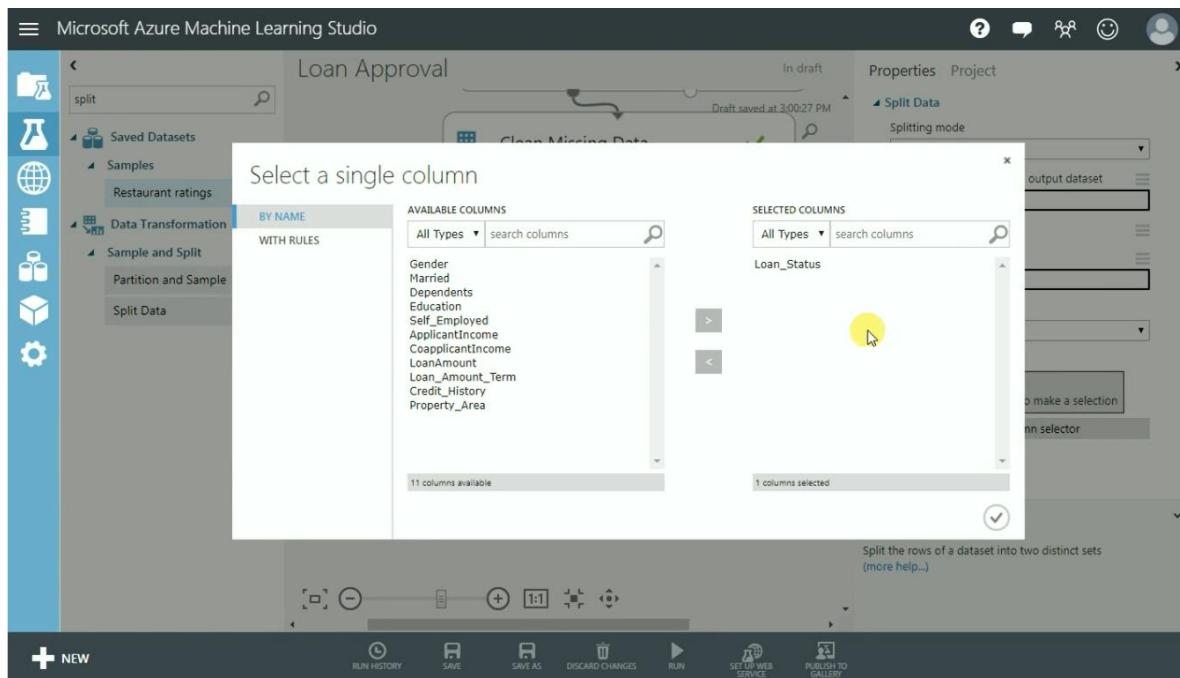
Search for split data and drop in canvas



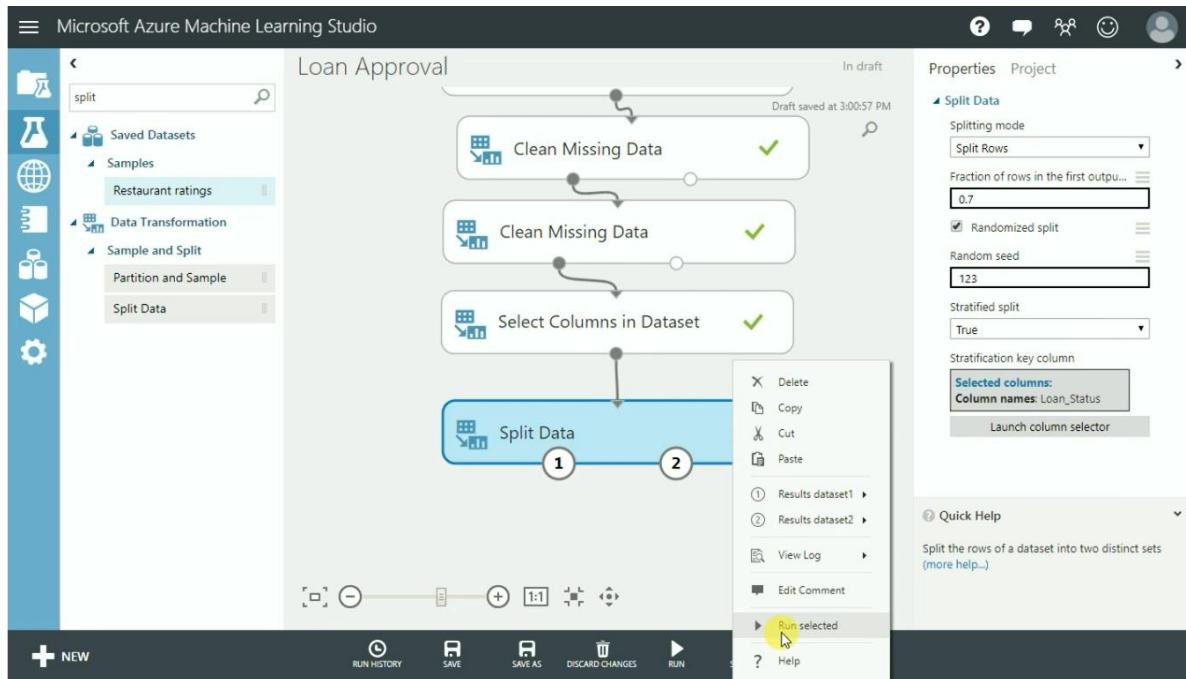
Attach both the dataset and input parameters as required



Now launch the column selector and select loan status and press ok



## Run the module



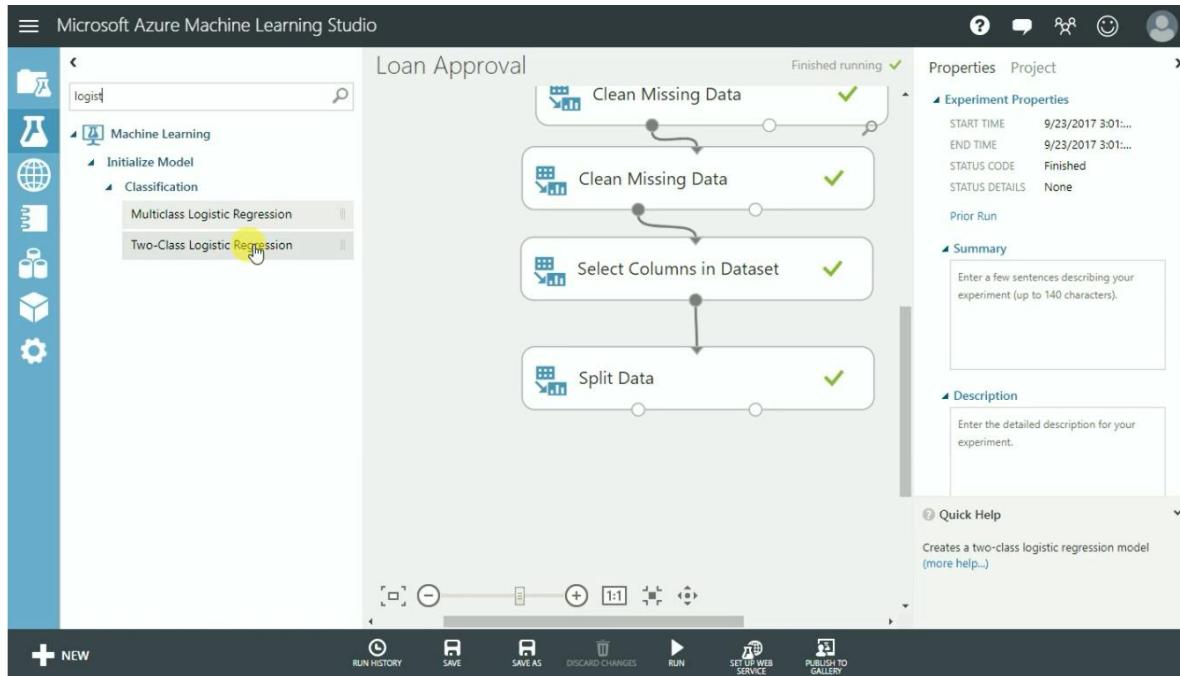
Visualize node 1 post execution. Now you can find 70 percent data here

And similarly, you can visualize node 2.

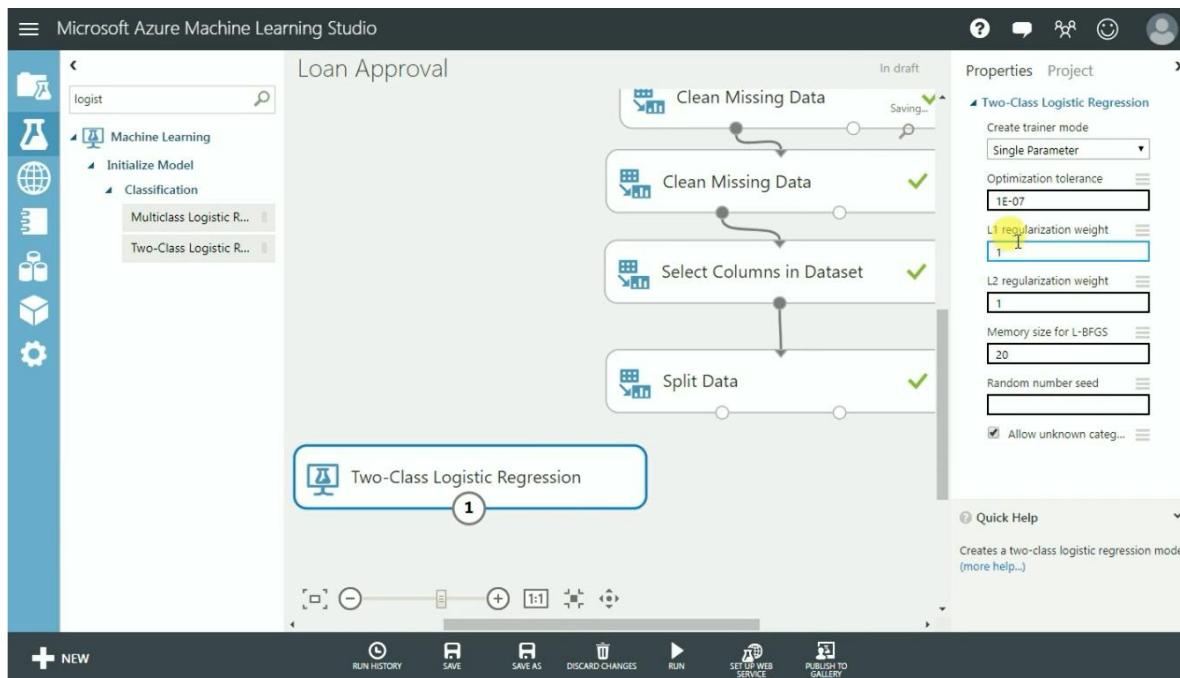
The screenshot shows the results of the "Split Data" module. The results dataset is named "Results dataset1" and contains 429 rows and 12 columns. The columns listed are Gender, Married, Dependents, Education, Self\_Employed, and ApplicantIncome. Below the table, there are visualizations for each column, including histograms and summary statistics. The status bar at the bottom indicates "Finished running".

## Usage of Two-Class Logistic Regression

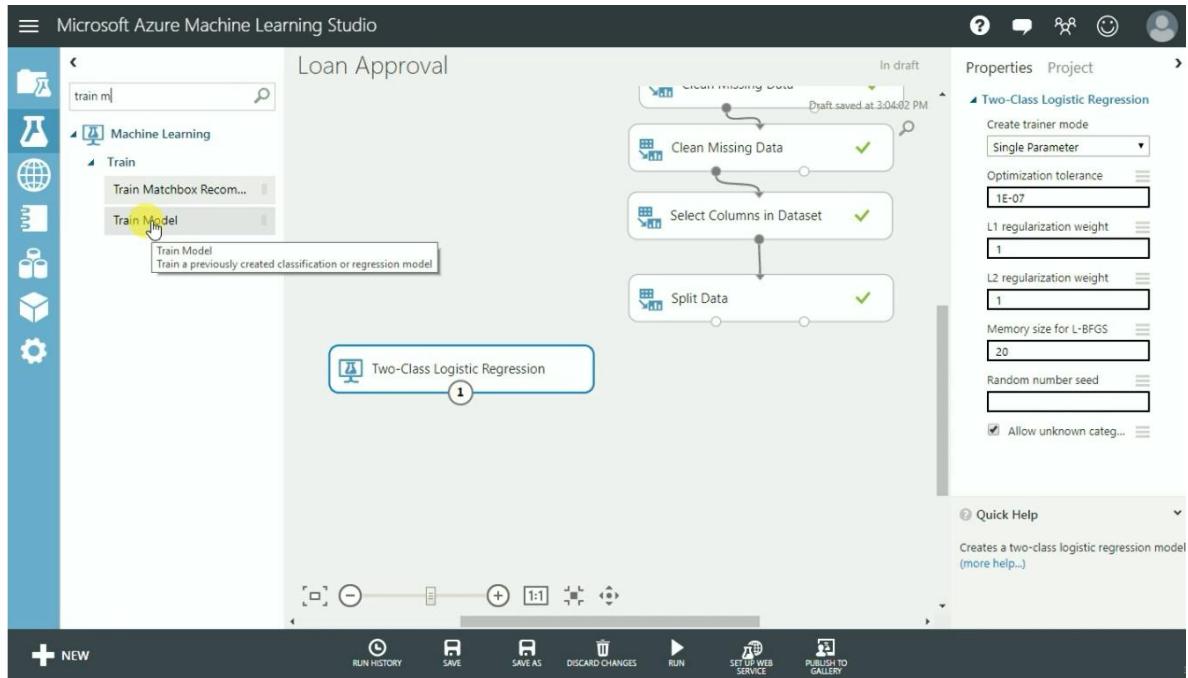
Search for two class logistic regression



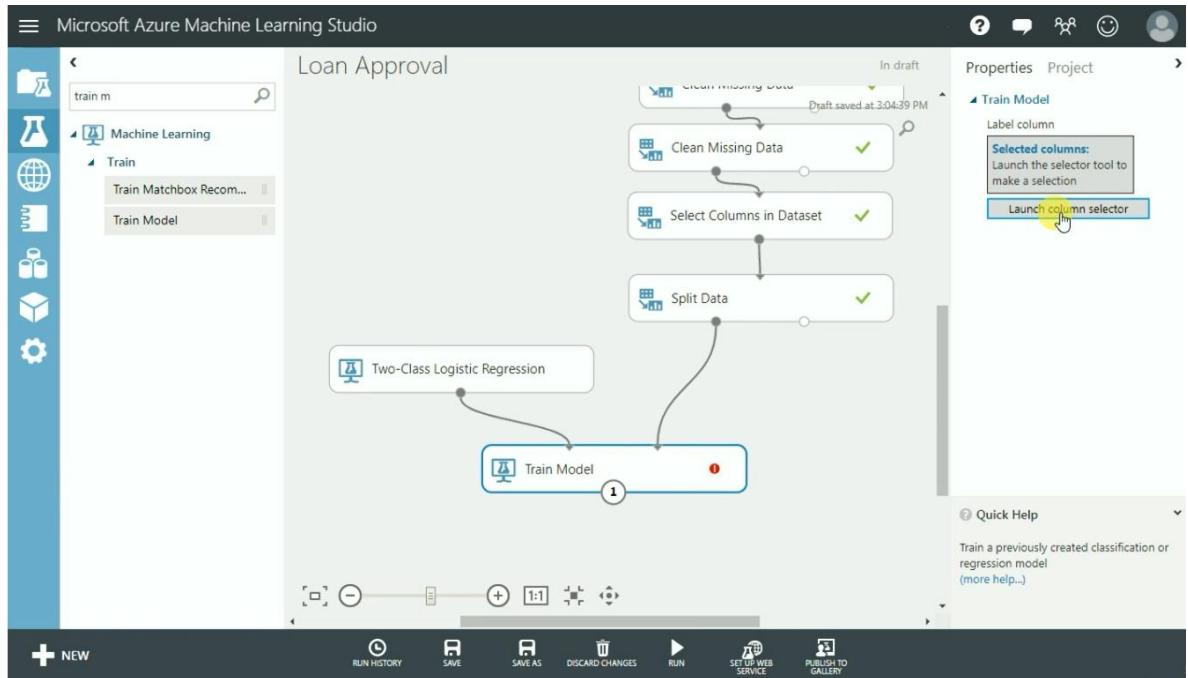
Drop the same in canvas and input required parameters



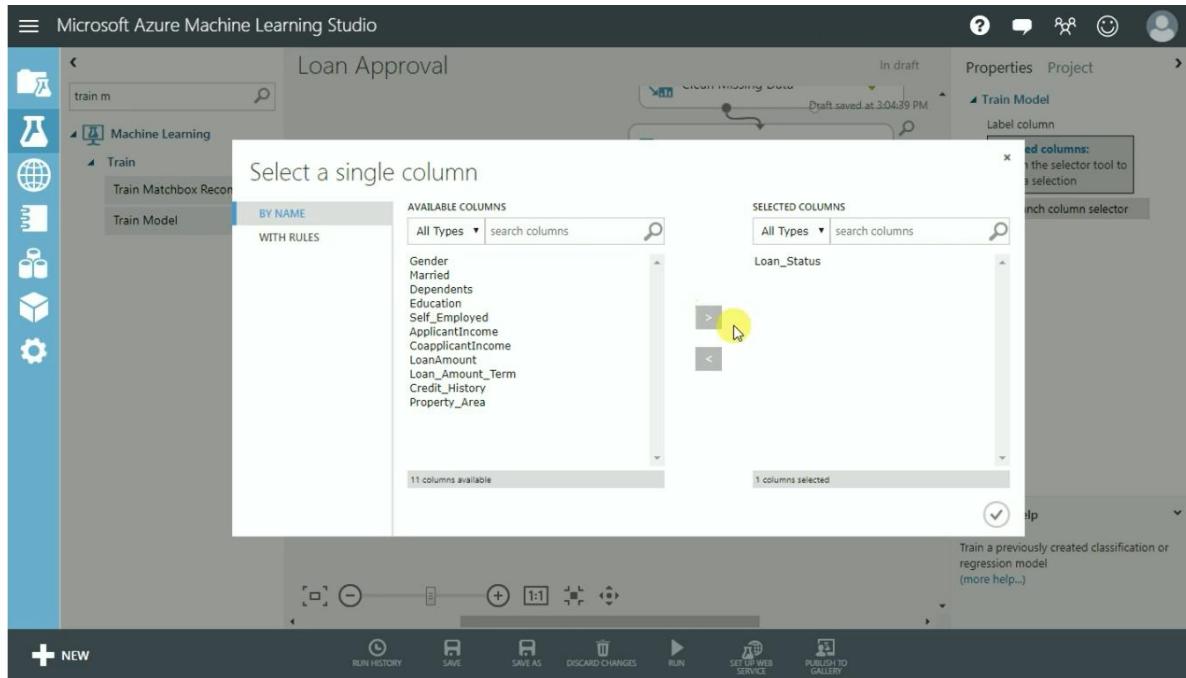
Search for model and drop in the canvas



Connect output node from split data and two class logistic regression to the train model input node

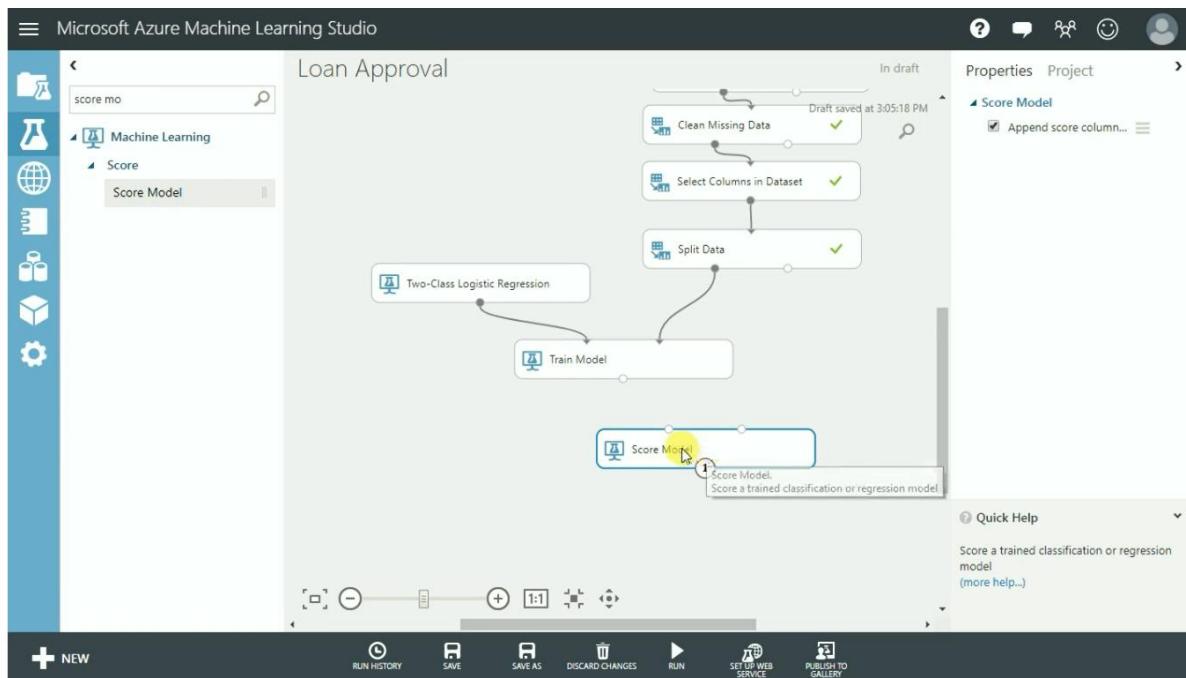


Now select launch column selector and select loan status and click ok

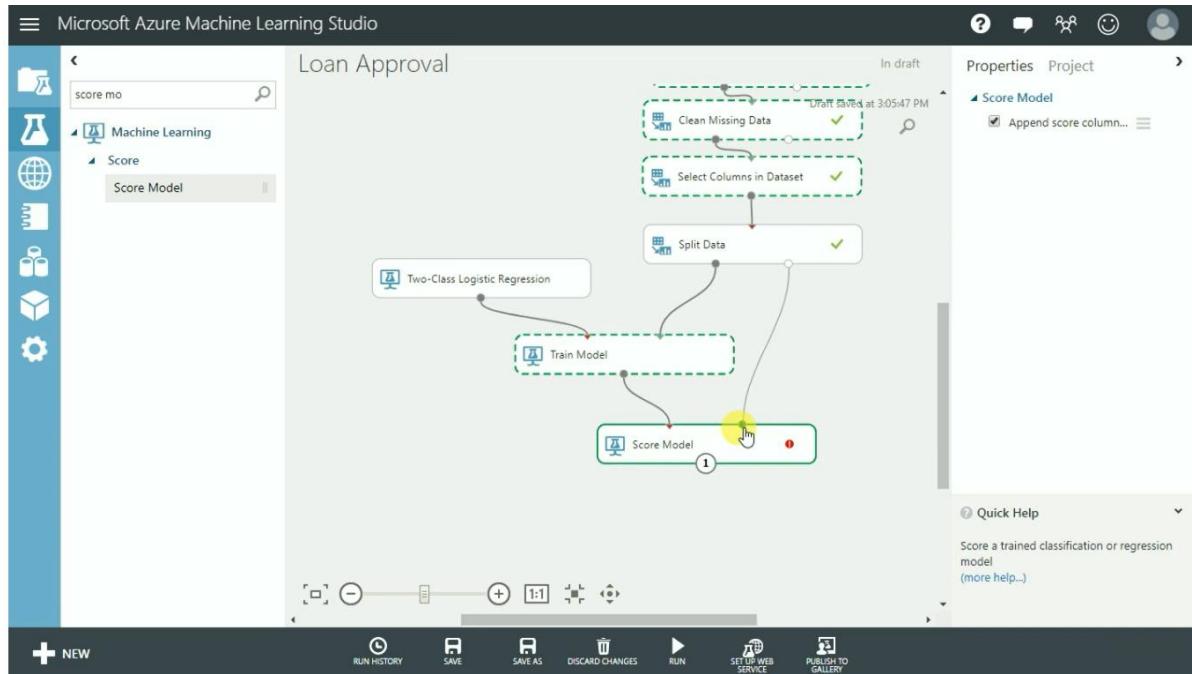


## Score and Evaluate Model

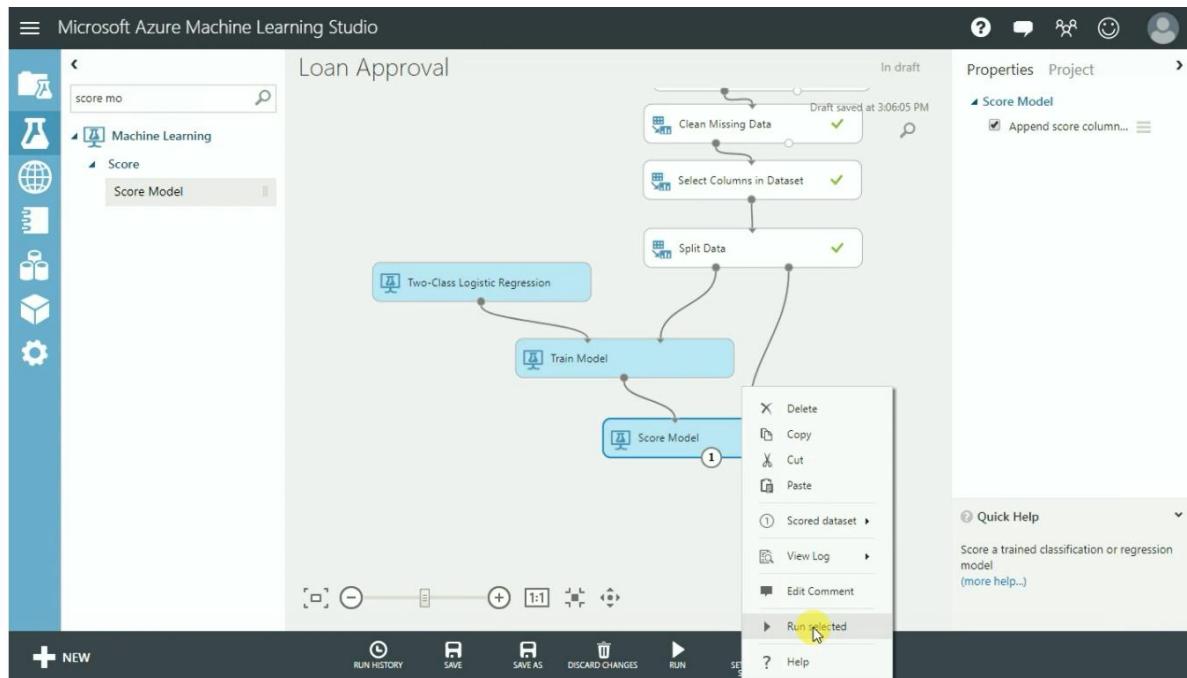
Search for score model and drop in canvas



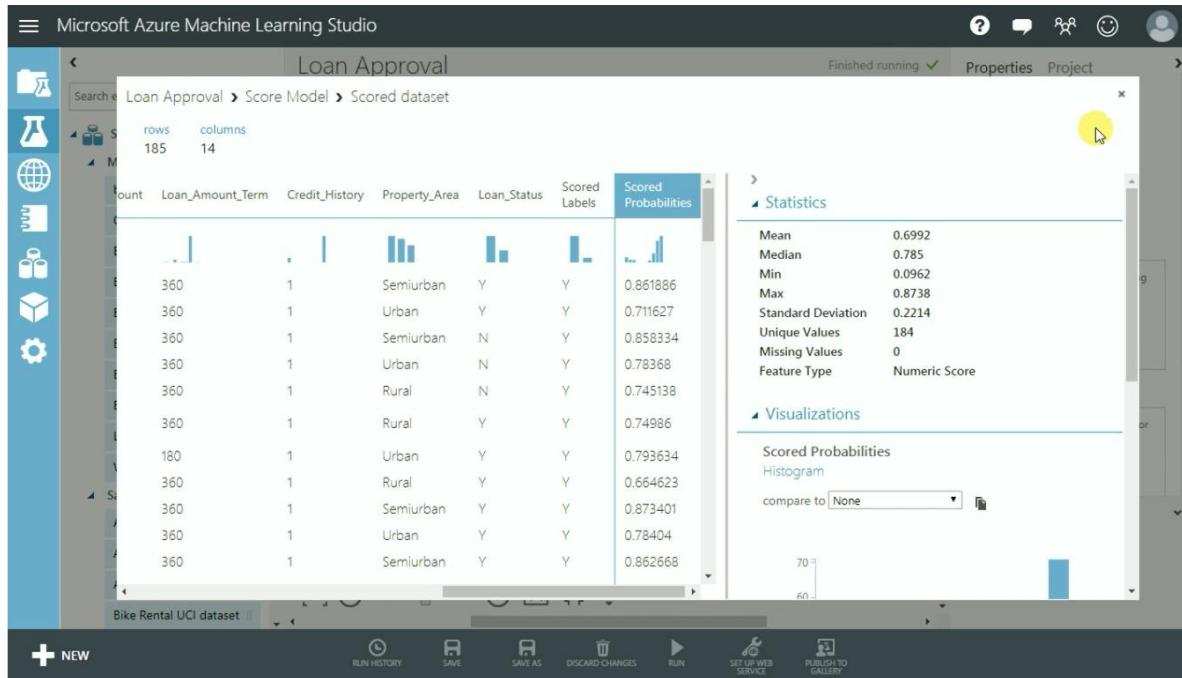
Connect output node from split data and connect to score model



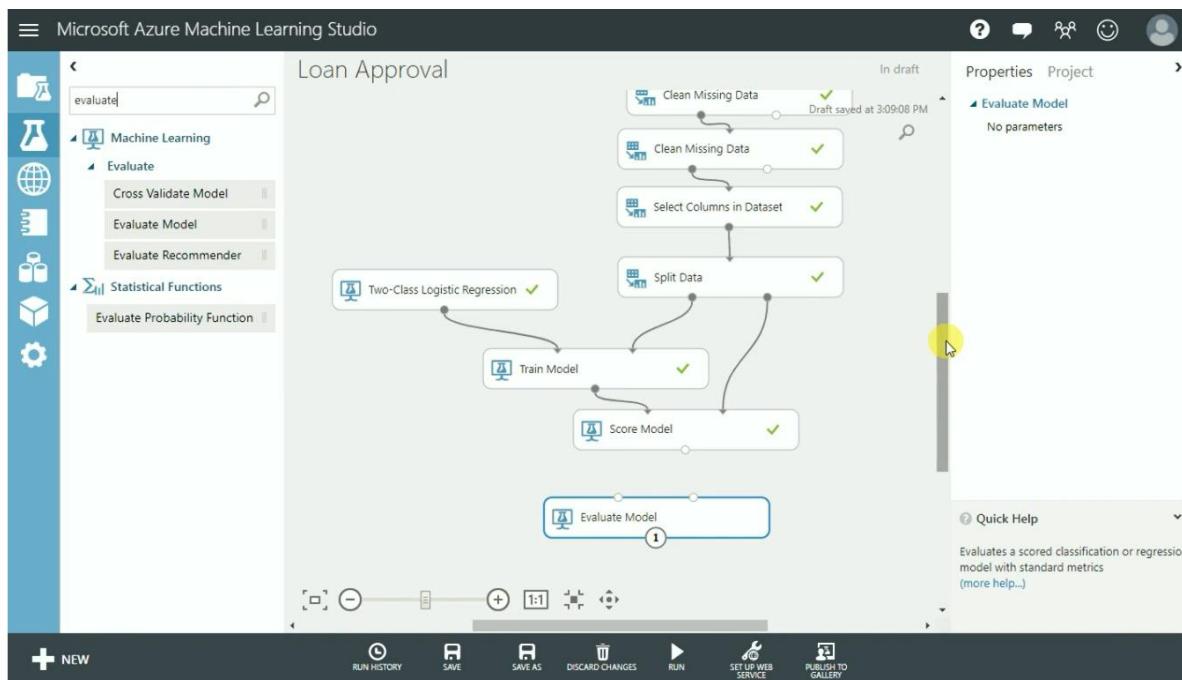
Right click and run selected



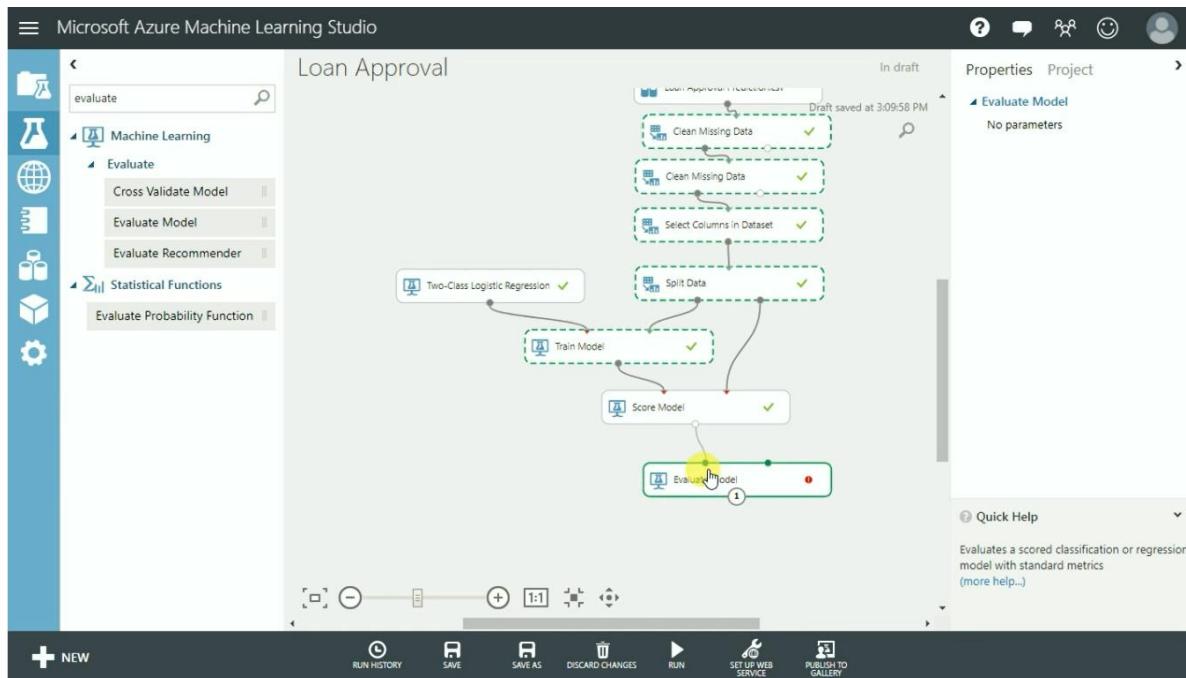
Visualize the result and check the accuracy



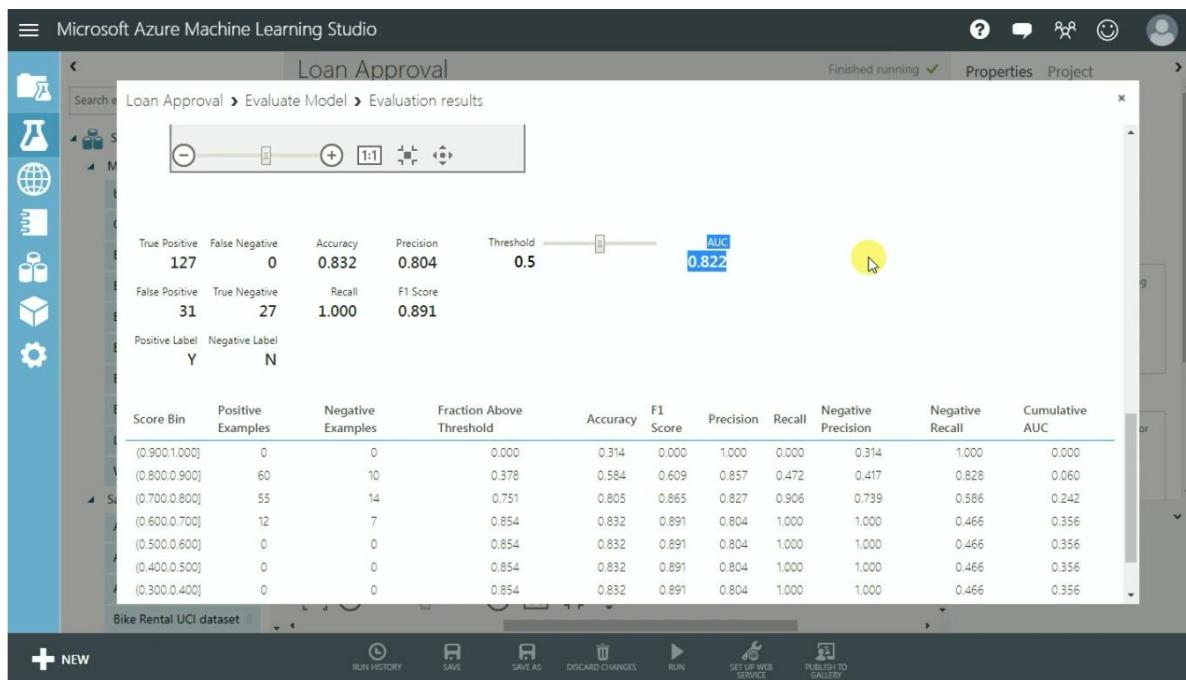
Since there is a compromise in accuracy search for evaluate model and drop in canvas



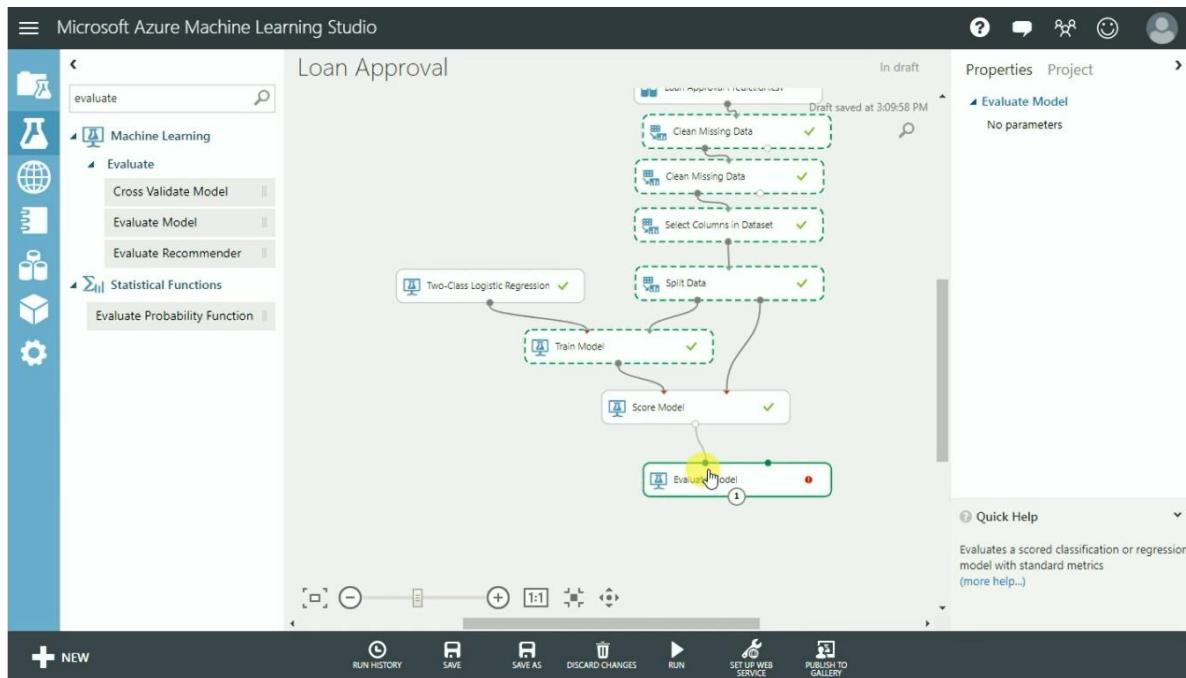
Connect the score model to evaluate model node1 and run selected



Visualize the output showing good accuracy



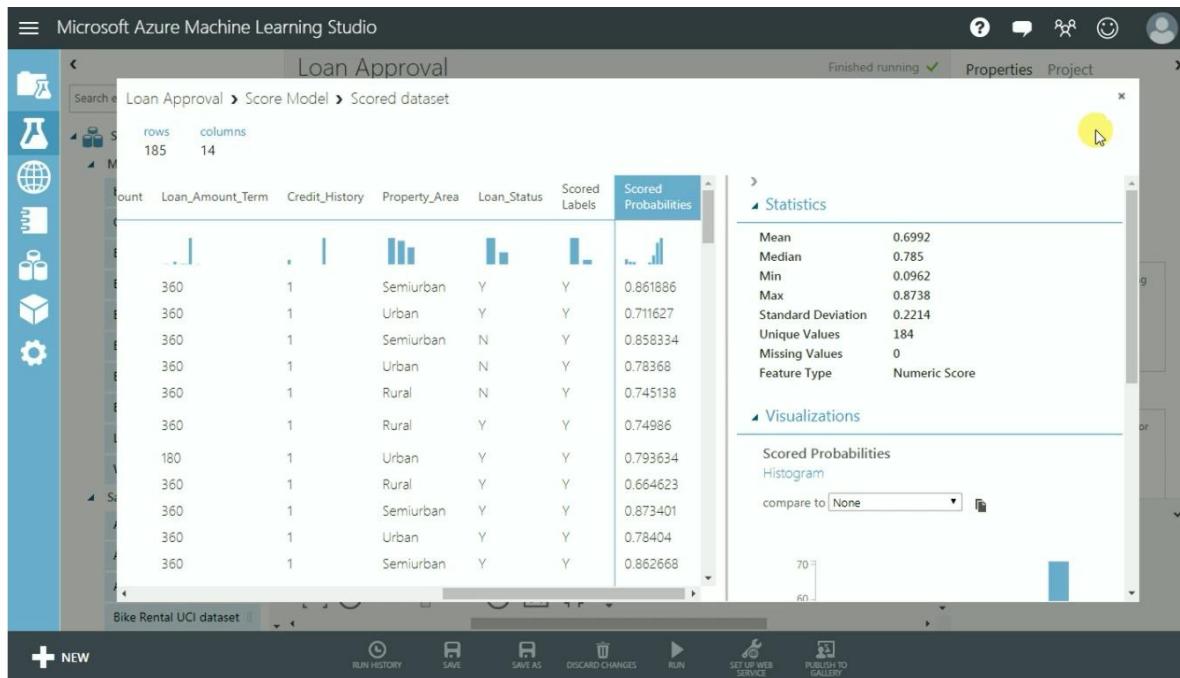
## UNDERSTANDING THE RESULTS



Right click and visualize Evaluate model



Right click and visualize Score model

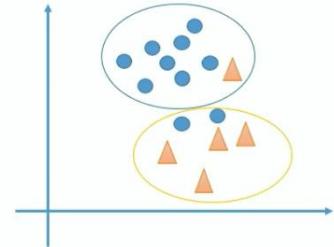


## Understanding the Output

Prediction Outcome			
	Predicted Positives	Predicted Negatives	
Actual Positives	True Positives	False Negatives	
Actual Negative	False Positives	True Negatives	

	Predicted Positives	Predicted Negatives	
Actual Positives	8	2	10
Actual Negative	1	4	5
		9	6



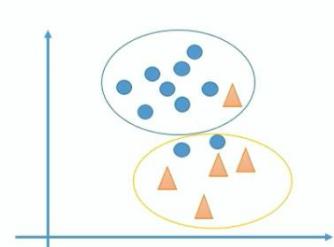
Accuracy – Proportions of total number of correct results

$$\text{Accuracy} = (8 + 4) / 15 = 0.8 \text{ or } 80\%$$

Prediction Outcome			
	Predicted Positives	Predicted Negatives	
Actual Positives	True Positives	False Negatives	
Actual Negative	False Positives	True Negatives	

	Predicted Positives	Predicted Negatives	
Actual Positives	8	2	10
Actual Negative	1	4	5
		9	6



Precision – Proportion of correct positive results out of all predicted positive results

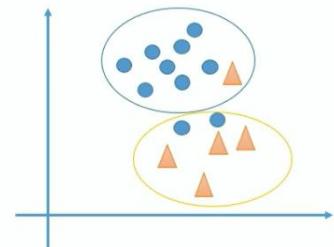
$$\text{Precision} = 8 / 9 = 0.889 \text{ or } 88.9\%$$

## Prediction Outcome

	Predicted Positives	Predicted Negatives
Actual Positives	True Positives	False Negatives
Actual Negative	False Positives	True Negatives

	Predicted Positives	Predicted Negatives	
Actual Positives	8	2	10
Actual Negative	1	4	5
	9	6	



Recall – Proportion of actual positive cases

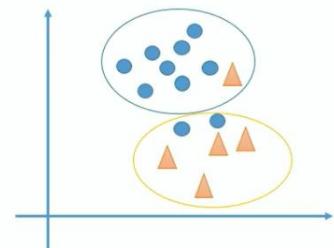
$$\text{Recall} = 8 / (8 + 2) = 0.8 \text{ or } 80\%$$

## Prediction Outcome

	Predicted Positives	Predicted Negatives
Actual Positives	True Positives	False Negatives
Actual Negative	False Positives	True Negatives

	Predicted Positives	Predicted Negatives	
Actual Positives	8	2	10
Actual Negative	1	4	5
	9	6	



F1-Score – Weighted Average (Harmonic Mean) of Precision and Recall

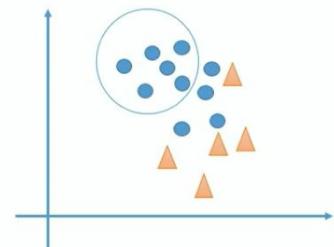
$$\text{F1Score} = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall}) = 0.84$$

## Prediction Outcome

	Predicted Positives	Predicted Negatives	
Actual Positives	True Positives	False Negatives	
Actual Negative	False Positives	True Negatives	

	Predicted Positives	Predicted Negatives	
Actual Positives	6	4	10
Actual Negative	0	5	5



In the Previous example

$$\text{Precision} = 6 / 6 = 1 \text{ or } 100\%$$

$$\text{Recall} = 6 / (6 + 4) = 0.6 \text{ or } 60\%$$

$$\text{Average} = 0.8$$

May lead to false interpretation

$$\text{Precision} = 0.889$$

$$\text{Recall} = 0.8$$

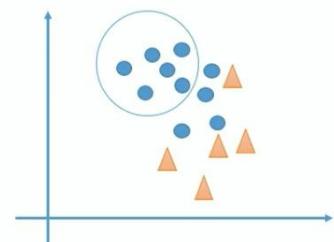
$$\text{Average} = 0.84$$

## Prediction Outcome

	Predicted Positives	Predicted Negatives	
Actual Positives	True Positives	False Negatives	
Actual Negative	False Positives	True Negatives	

	Predicted Positives	Predicted Negatives	
Actual Positives	6	4	10
Actual Negative	0	5	5



In the first example

$$\text{Precision} = 6 / 6 = 1 \text{ or } 100\%$$

$$\text{Recall} = 6 / (8 + 2) = 0.6 \text{ or } 60\%$$

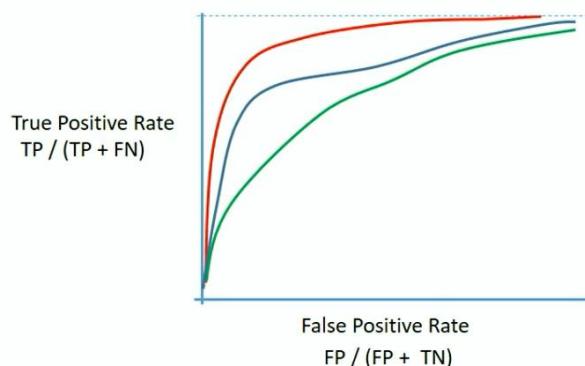
$$\text{F1Score} = 0.75$$

$$\text{Precision} = 0.889$$

$$\text{Recall} = 0.8$$

$$\text{F1Score} = 0.84$$

## AUC ROC



AUC – Area Under the Curve

ROC – Receiver Operating Characteristics

First used during World War II for the analysis of radar signals

Following the attack on Pearl Harbor in 1941, the United States army began new research to increase the prediction of correctly detected Japanese aircraft from their radar signals.

For this purposes they measured the ability of radar receiver operators to make these important distinctions, which was called the Receiver Operating Characteristics

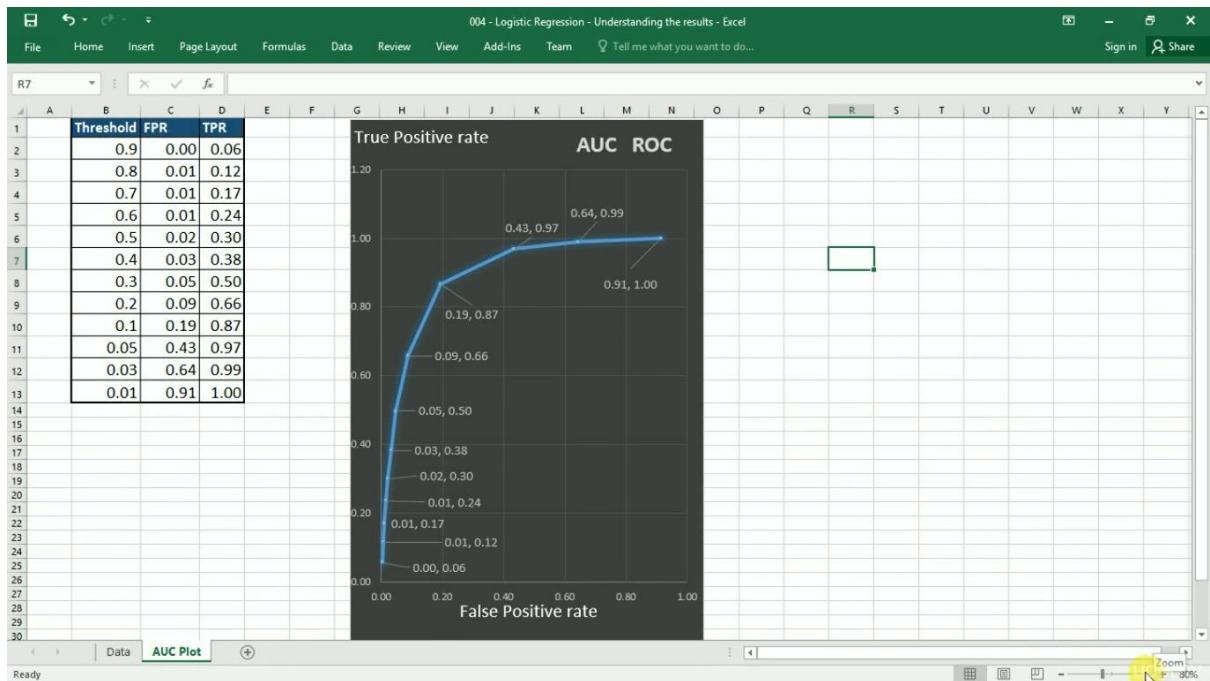
Provides a single number that lets you compare models of different types.

004 - Logistic Regression - Understanding the results - Excel															
File Home Insert Page Layout Formulas Data Review View Add-Ins Team Tell me what you want to do... Sign in Share															
H17															
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
4	Threshold	0.05				Threshold	0.03				Threshold	0.01			
5	TP	FN				TP	FN				TP	FN			
6	1019	33	TPR	0.96863		1042	10	TPR	0.99049		1051	1	TPR	0.99905	
7	FP	TN	FPR	0.43342		FP	TN	FPR	0.64431		FP	TN	FPR	0.91277	
8	3463	4527				5148	2842				7293	697			
9	Threshold	0.1				Threshold	0.2				Threshold	0.31			
11	TP	FN				TP	FN				TP	FN			
12	913	139	TPR	0.86787		692	360	TPR	0.65779		524	528	TPR	0.4981	
13	FP	TN	FPR	0.19312		FP	TN	FPR	0.08811		FP	TN	FPR	0.04781	
14	1543	6447				704	7286				382	7608			
15	Threshold	0.4				Threshold	0.5				Threshold	0.6			
17	TP	FN				TP	FN				TP	FN			
18	404	648	TPR	0.38403		317	735	TPR	0.30133		249	803	TPR	0.23669	
19	FP	TN	FPR	0.03179		FP	TN	FPR	0.0219		FP	TN	FPR	0.01452	
20	254	7736				175	7815				116	7874			
21	Threshold	0.7				Threshold	0.8				Threshold	0.9			
23	TP	FN				TP	FN				TP	FN			
24	179	873	TPR	0.17015		122	930	TPR	0.11597		60	992	TPR	0.05703	
25	FP	TN	FPR	0.00914		FP	TN	FPR	0.00526		FP	TN	FPR	0.00275	
26	73	7917				42	7948				22	7968			
27	Data	AUC Plot	+												

004 - Logistic Regression - Understanding the results - Excel

File Home Insert Page Layout Formulas Data Review View Add-ins Team Tell me what you want to do... Sign in Share Row: 3

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
3		FPR = FP/(FP + TN)															
4																	
5	Threshold	0.05				Threshold	0.03				Threshold	0.01					
6	TP	FN				TP	FN				TP	FN					
7	1019	33	TPR	0.96863		1042	10	TPR	0.99049		1051	1	TPR	0.99905			
8	FP	TN	FPR	0.43342		FP	TN	FPR	0.64431		FP	TN	FPR	0.91277			
9	3463	4527				5148	2842				7293	697					
11	Threshold	0.1				Threshold	0.2				Threshold	0.31					
12	TP	FN				TP	FN				TP	FN					
13	913	139	TPR	0.86787		692	360	TPR	0.65779		524	528	TPR	0.4981			
14	FP	TN	FPR	0.19312		FP	TN	FPR	0.08811		FP	TN	FPR	0.04781			
15	1543	6447				704	7286				382	7608					
17	Threshold	0.4				Threshold	0.5				Threshold	0.6					
18	TP	FN				TP	FN				TP	FN					
19	404	648	TPR	0.38403		317	735	TPR	0.30133		249	803	TPR	0.23669			
20	FP	TN	FPR	0.03179		FP	TN	FPR	0.0219		FP	TN	FPR	0.01452			
21	254	7736				175	7815				116	7874					
23	Threshold	0.7				Threshold	0.8				Threshold	0.9					
24	TP	FN				TP	FN				TP	FN					
25	179	873	TPR	0.17015		122	930	TPR	0.11597		60	992	TPR	0.05703			
26	FP	TN	FPR	0.00914		FP	TN	FPR	0.00526		FP	TN	FPR	0.00275			
27	73	7917				42	7948				22	7968					
28	Data	AUC Plot															



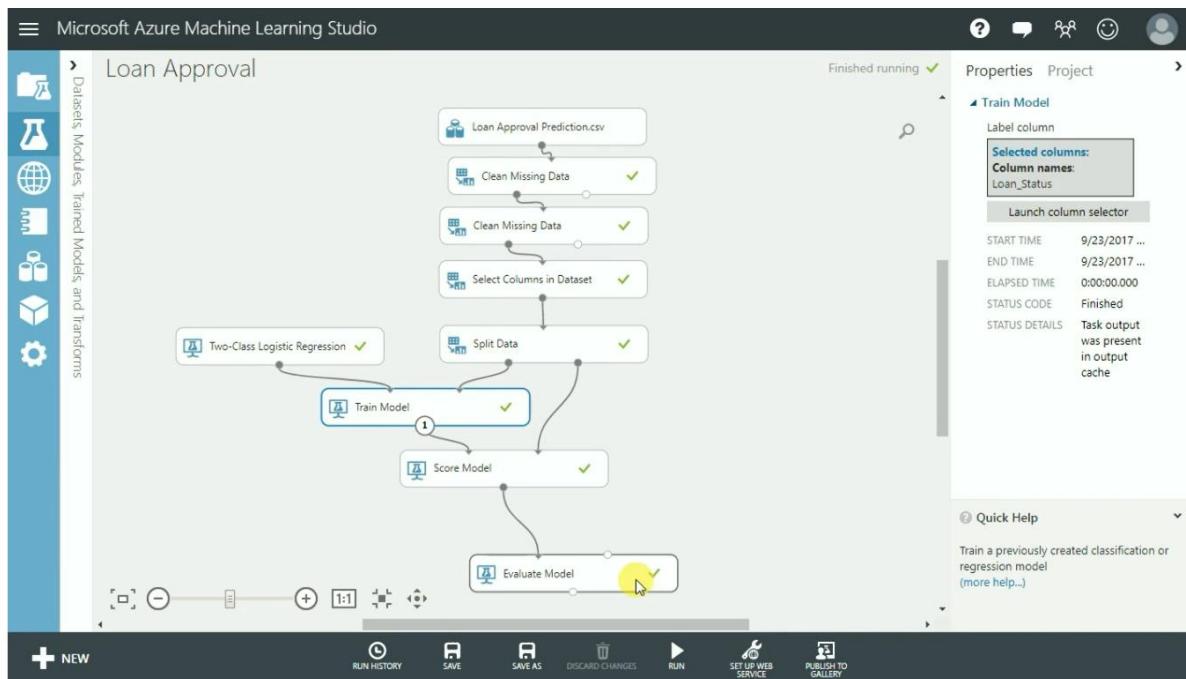
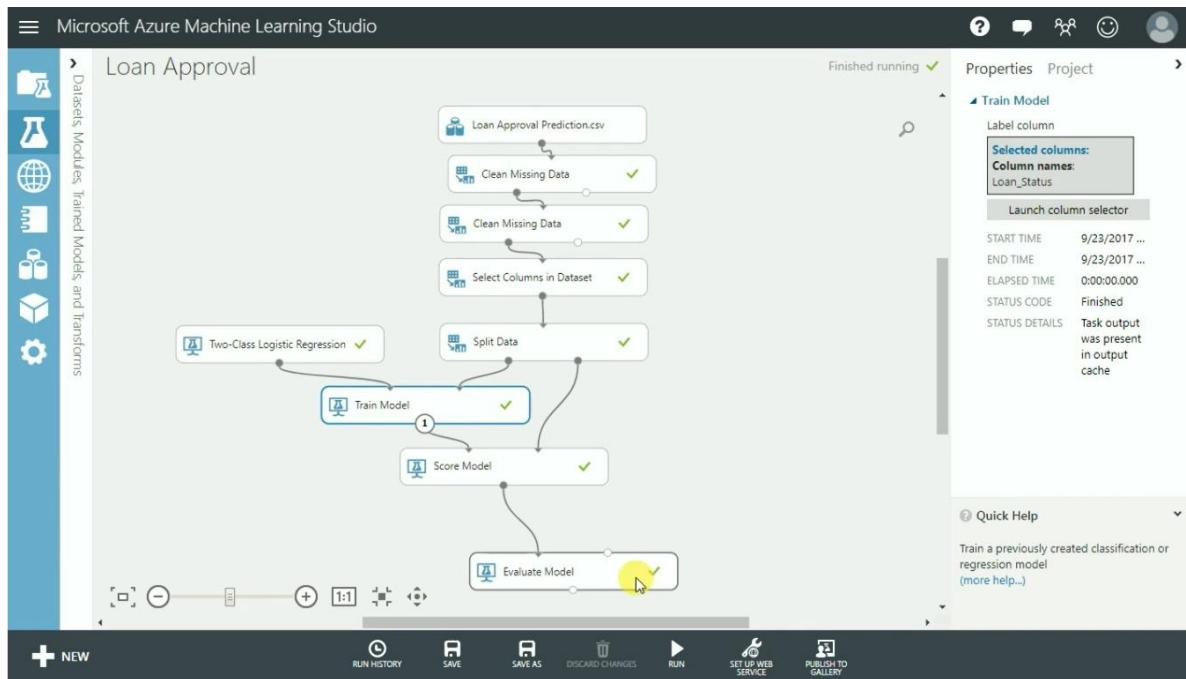
File Home Insert Page Layout Formulas Data Review View Add-Ins Team Tell me what you want to do... Sign in Share

K8 A B C D E F G H I J K L M N O P Q

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
4																	
5	Threshold	0.05				Threshold	0.03				Threshold	0.01					
6	TP	FN				TP	FN				TP	FN					
7	1019	33	TPR	0.96863		1042	10	TPR	0.99049		1051	1	TPR	0.99905			
8	FP	TN	FPR	0.43342		FP	TN	FPR	0.64431		FP	TN	FPR	0.91277			
9	3463	4527				5148	2842				7293	697					
11	Threshold	0.1				Threshold	0.2				Threshold	0.31					
12	TP	FN				TP	FN				TP	FN					
13	913	139	TPR	0.86787		692	360	TPR	0.65779		524	528	TPR	0.4981			
14	FP	TN	FPR	0.19312		FP	TN	FPR	0.08811		FP	TN	FPR	0.04781			
15	1543	6447				704	7286				382	7608					
17	Threshold	0.4				Threshold	0.5				Threshold	0.6					
18	TP	FN				TP	FN				TP	FN					
19	404	648	TPR	0.38403		317	735	TPR	0.30133		249	803	TPR	0.23669			
20	FP	TN	FPR	0.03179		FP	TN	FPR	0.0219		FP	TN	FPR	0.01452			
21	254	7736				175	7815				116	7874					
23	Threshold	0.7				Threshold	0.8				Threshold	0.9					
24	TP	FN				TP	FN				TP	FN					
25	179	873	TPR	0.17015		122	930	TPR	0.11597		60	992	TPR	0.05703			
26	FP	TN	FPR	0.00914		FP	TN	FPR	0.00526		FP	TN	FPR	0.00275			
27	73	7917				42	7948				22	7968					
28																	
29	Data	AUC Plot															

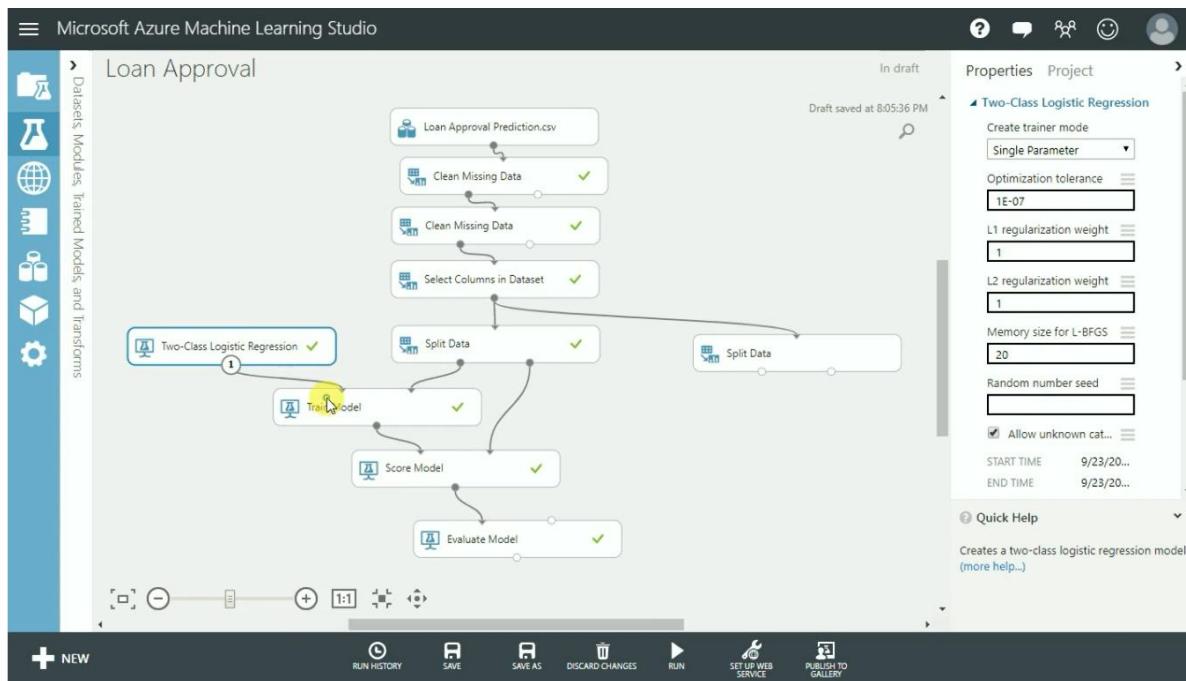
## Impact Analysis & Stratification

Select the Loan prediction model and visualize the evaluate model

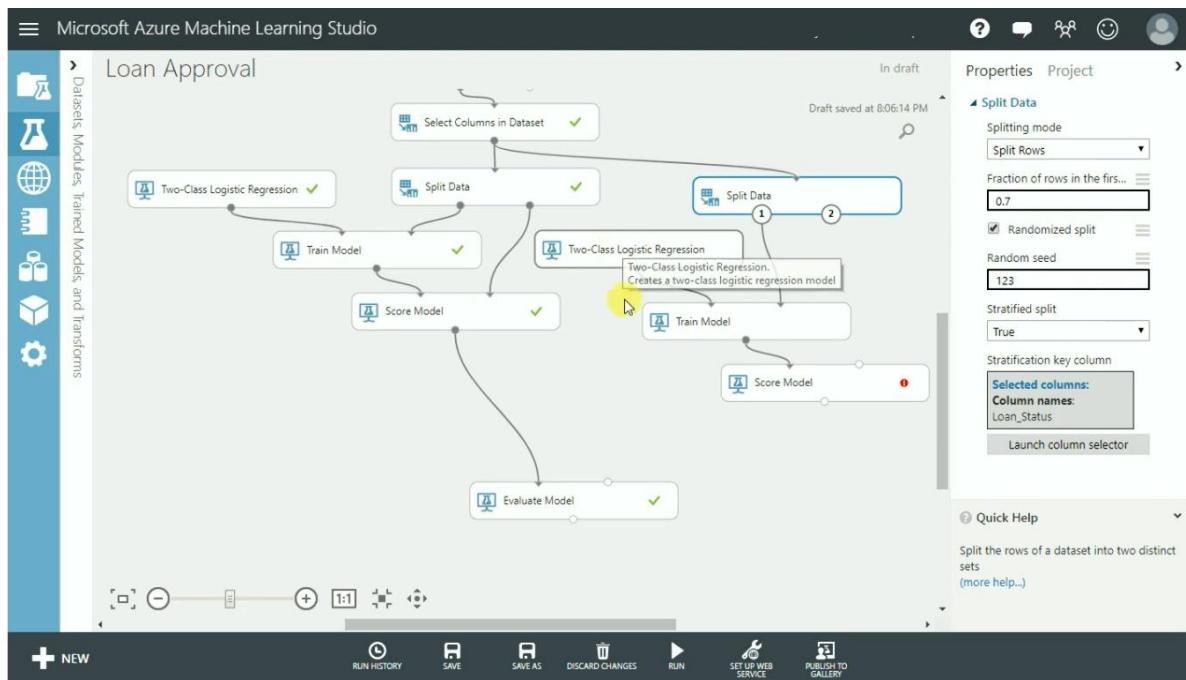


To analyse the impact of stratification which done at split data

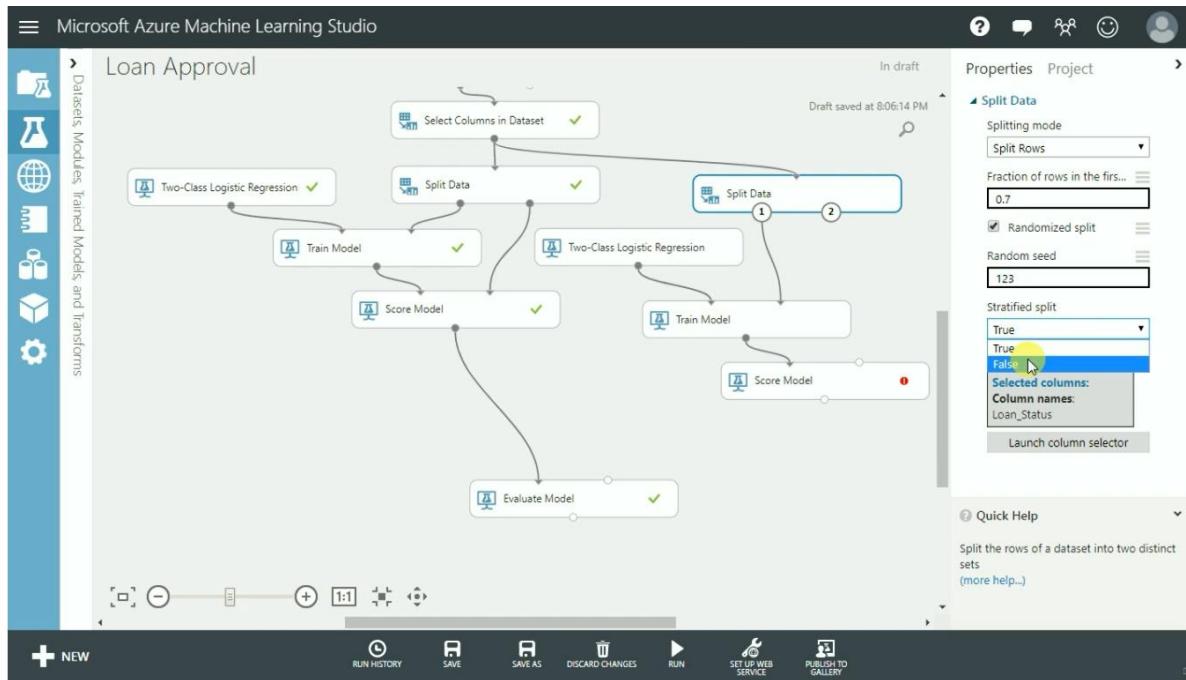
Repeat the datasets from split data again as below



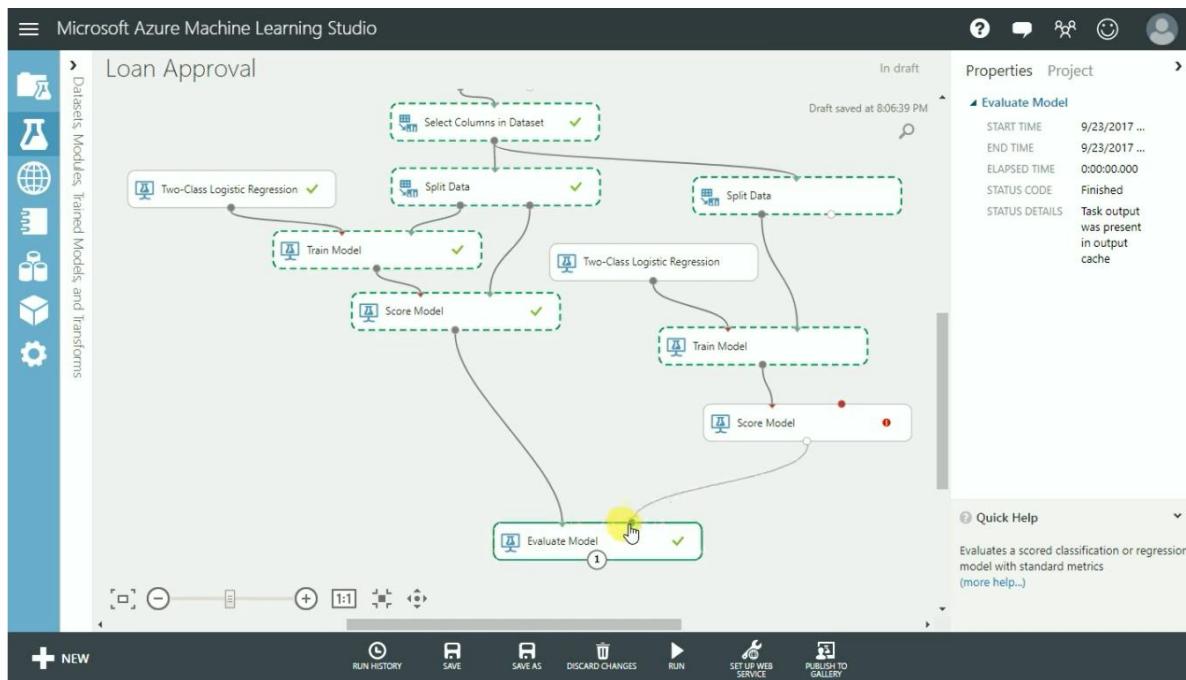
Copy and paste the data sets in split data as shown



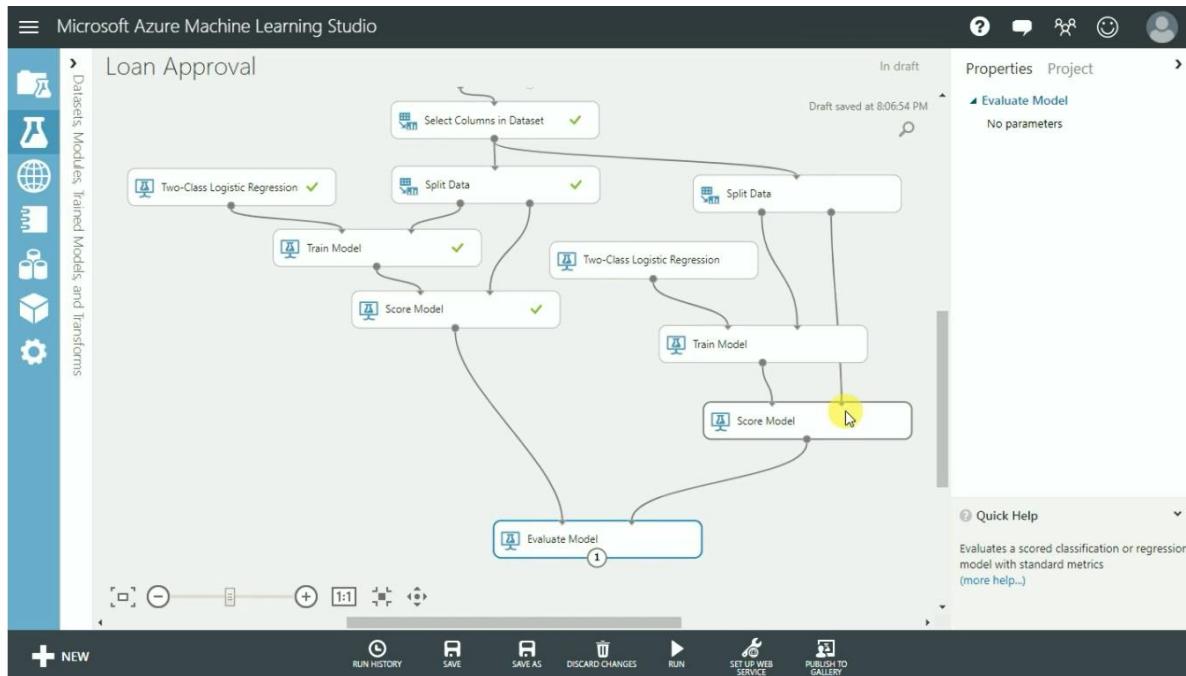
## Change parameter as false in stratified split



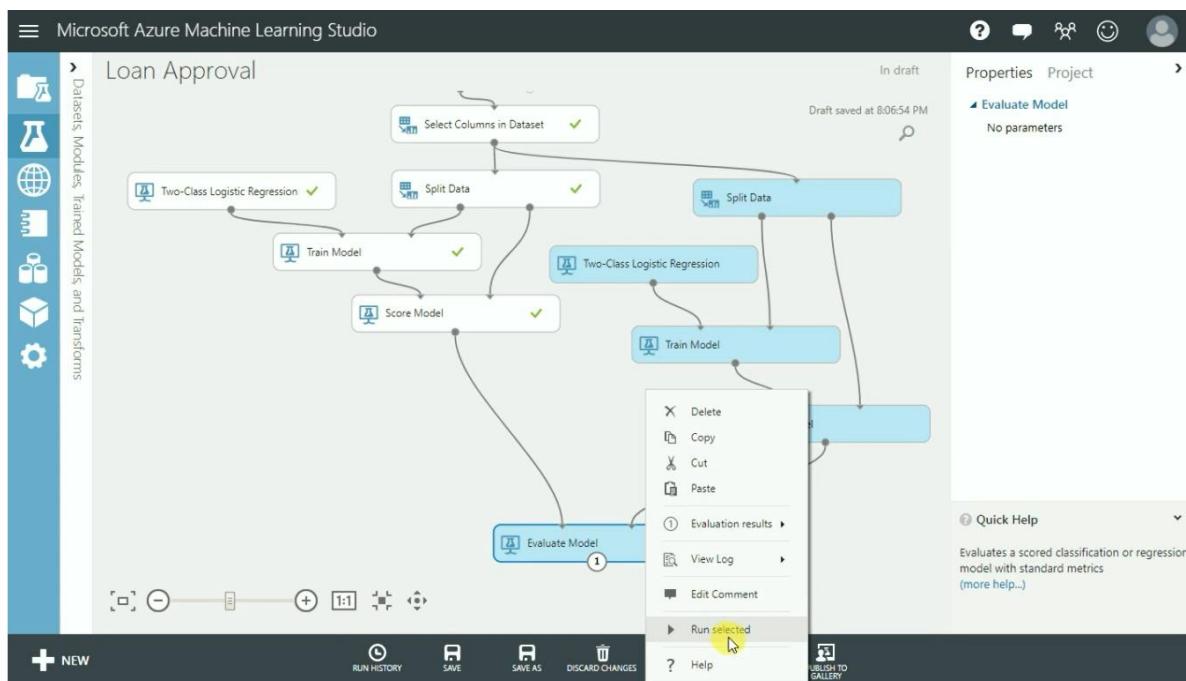
Connect output node from score model to input node of evaluate model



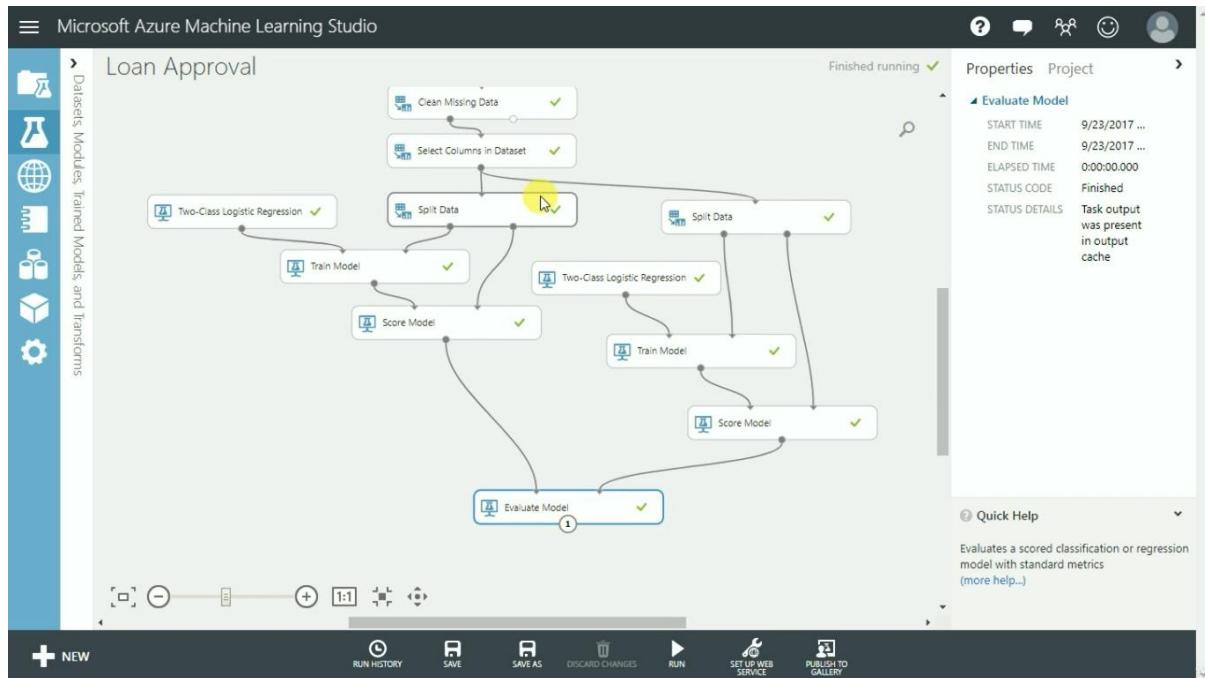
Connect split data output node 2 to score model input node 2



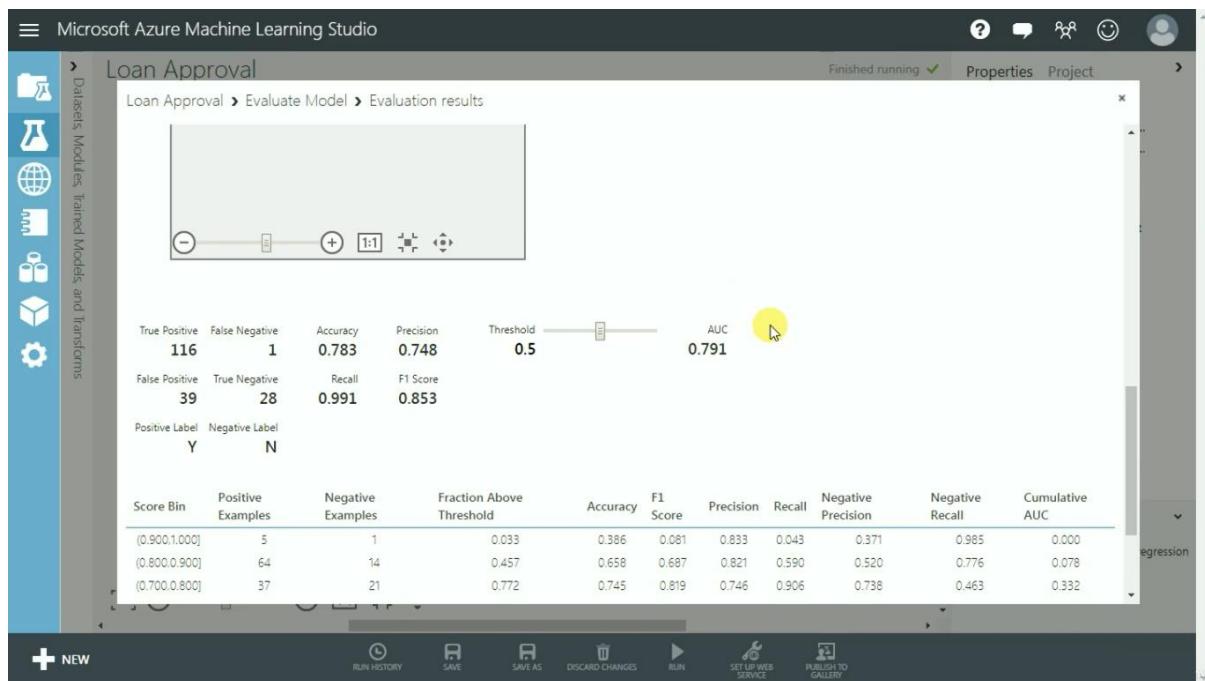
Now time to run the evaluate model



After successful execution visualize the evaluate model



After this action accuracy and auc has come down and thus found impact of stratification



## Stratification Impact

True Positive: 127, False Negative: 0, Accuracy: 0.832, Precision: 0.804, Threshold: 0.5, AUC: 0.822

False Positive: 31, True Negative: 27, Recall: 1.000, F1 Score: 0.891

Positive Label: Y, Negative Label: N

True Positive: 116, False Negative: 1, Accuracy: 0.783, Precision: 0.748, Threshold: 0.5, AUC: 0.791

False Positive: 39, True Negative: 28, Recall: 0.991, F1 Score: 0.853

Positive Label: Y, Negative Label: N



1/16