



## CIS5560 Term Project Tutorial



**Authors:** [Anjali Baldawa; Dhruvi Patel; Verusca Capuno; Digeshkumar Kansara](#)

**Instructor:** [Jongwook Woo](#)

**Date:** [05/17/2020](#)

### Lab Tutorial

Anjali Baldawa ([abaldaw@calstatela.edu](mailto:abaldaw@calstatela.edu))

Dhruvi Patel([dpatel86@calstatela.edu](mailto:dpatel86@calstatela.edu))

Verusca Capuno([vcapuno@calstatela.edu](mailto:vcapuno@calstatela.edu))

Digeshkumar Kansara([dkansar@calstatela.edu](mailto:dkansar@calstatela.edu))

05/17/2020

### Outbrain Ad Prediction using Spark

---

### Objectives

In this hands-on lab, you will learn how to:

- Get data from Kaggle
- Sampling data using Jupyter Notebook and Pandas
- Upload data, create Experiment, join data, feature engineering, train, test data and measure accuracy using different algorithm in Azure ML Studio
- Upload data, create cluster, use notebook and Pyspark file, create schema, join data, cross validation, tuning parameter, train, test data, measure accuracy of algorithm using different measures and using different algorithms in data bricks.

- Upload data, use Pyspark file, create schema, join data, cross validation, tuning parameter, train, test data, measure accuracy of algorithm using different measures and using different algorithms In Oracle BDCE.

## **Platform Specification**

### **Azure ML**

- Free workspace
- 10 GB storage
- Single node

### **Data Bricks community Edition**

- Databricks subscription
- Cluster 6.5 ML (includes Scala 2.11, Spark 2.4.5)
- 15.3 GB Memory, 2 Cores, 1 DBU
- Python Version 3

### **Oracle BDCE**

- Cluster 5.2 (includes Apache Spark 2.4)
- Python Version 3
- Cluster size 1000.3 GB
- Number of Processors 32
- Memory Size 247.625312GB
- CPU Speed 2.20 GHz

## **Step 1: Get data from Kaggle Website**

---

- 1) The following steps explain how to download data from Kaggle Website.
  - To download Outbrain ad click prediction data we need to navigate to Kaggle website using following link.  
<https://www.kaggle.com/c/outbrain-click-prediction/data>

**Note:** Following are additional steps to download data.

1. Navigate to Rule Tab from Top menu.
2. Press I Understand Button.
3. Enter your Phone number.
4. It will send you a 6 digit password.

5. Enter password and Now you are ready to download data.

- Now move to data tab from Top menu and download .csv files.

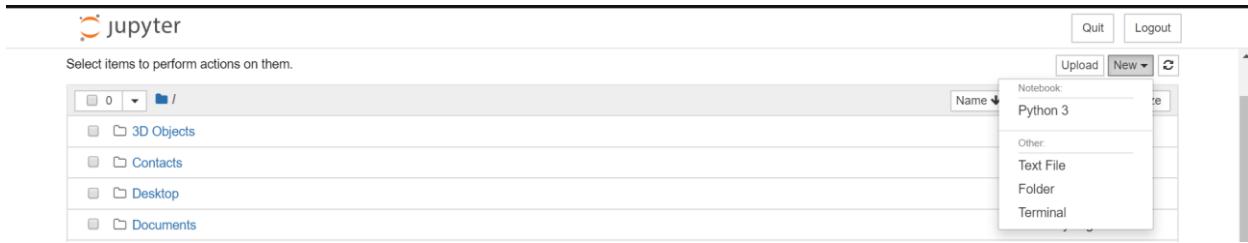
**Note:** The size of the files is huge so you need to have a good internet connection to download the dataset. Approximately, it would take 1 hour to download 1 file. The download time depends on the internet speed. Once you download files into your local pc, you need to unzip the file in your local system.

## **Step 2: Sample Data using Jupyter Notebook**

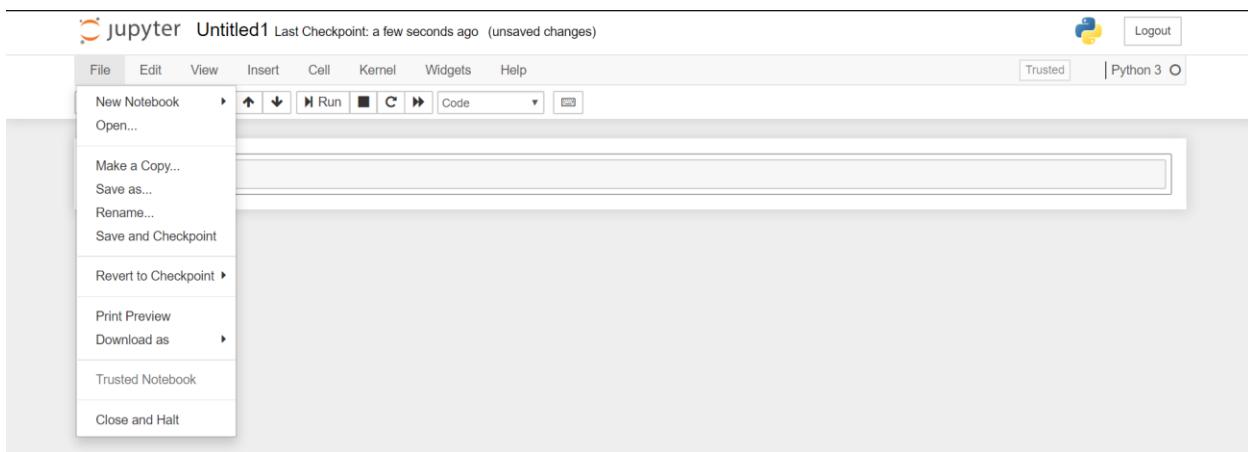
---

To sample data in smaller CSV file We need have Python3 and Anaconda. We are using pandas library to sample dataset.

1. Download Python.
  1. Download Python from following link:  
<https://www.python.org/downloads/>
  2. You need to choose proper version according to your system requirement MAC OS, Windows or Linux.
  3. Press Downloads button (Current latest version is 3.8.2)
  4. After Downloading it is simple install on your system.
2. Download Anaconda. (Minimum 4GB, Recommended 8GB)
  1. Download Anaconda from following link:  
<https://www.anaconda.com/products/individual>
  2. Scroll Down to Bottom.
  3. Choose your Operating system and download appropriate version.
  4. It will take I little while as size is 465 MB
  5. Unzip downloaded anaconda file.
  6. Install on your system as it is simple steps just press next and in last press finish button.
3. Launch Jupyter notebook and sample data.
  1. Start Anaconda in your system.
  2. You will see jupyter notebook option and press Launch button it will open in your default browser. It will open following screen for you. From right side top pane click on new and select Python which will create new notebook with python3.



It will open following window. From file choose **open** option.

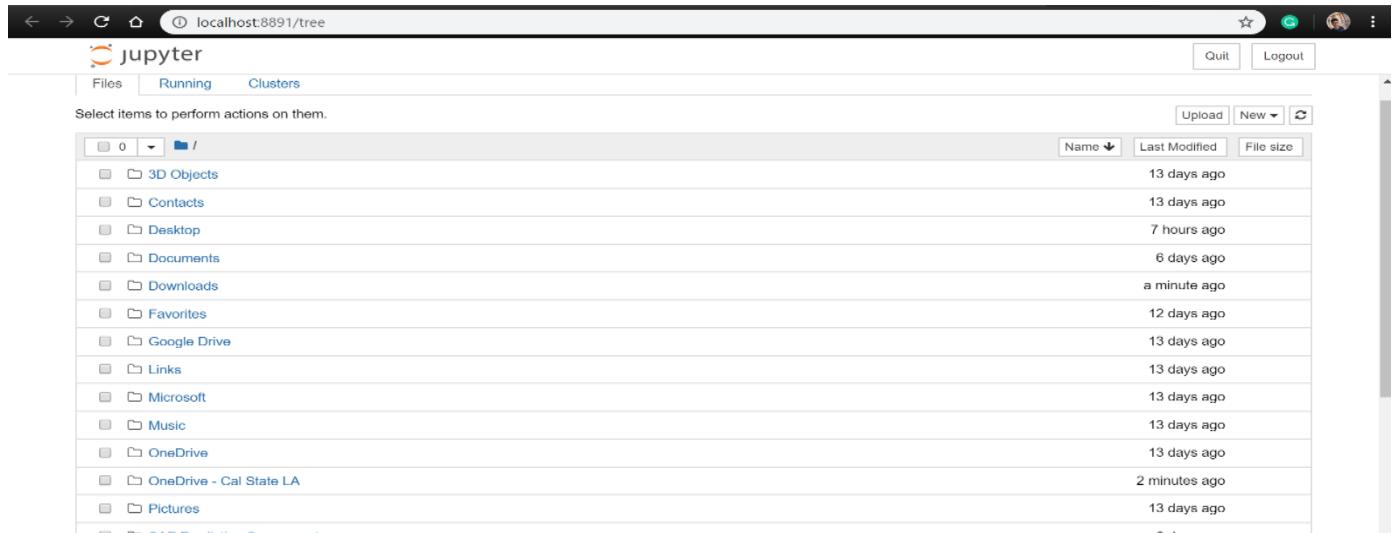


It will show you all directories in your system.

Now you need to download **Project\_sample.ipynb** file provided to you.

Now navigate to directory where your file is located and press open.

It will open Notebook file for you.



Then follow steps in file using which you can easily sample data.

## Step 3: Azure ML Studio

In step 2, we sampled data using jupyter notebook. Now we have our sampled files (Around 40 MB) to import in Azure ML studio.

Create an Azure ML Account

Azure ML offers a free-tier account, which you can use to complete the experiment.

Sign Up for a Microsoft Account

1. If you do not already have a Microsoft account, sign up for one at <https://signup.live.com/>. You don't need to use your school email account to sign up but you can use any email account.

**Sign Up for a Free Azure ML Account**

1. Browse to <https://azure.microsoft.com/en-us/services/machine-learning/> login and click Get started now
2. When prompted, choose the option to sign in, and sign in with your Microsoft account credentials.
3. On the **Welcome** page, watch the overview video if you want to see an introduction to Azure ML Studio. Then close the **Welcome** page by clicking the checkmark icon.

## Create Experiment and Upload Data Set (CSV files) in ML studio

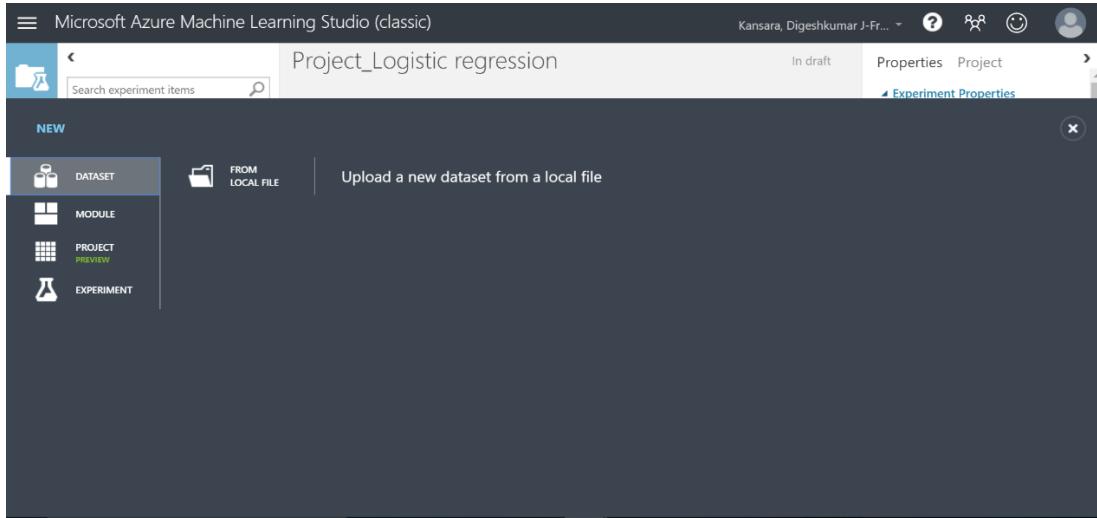
1. You should now be in Azure ML Studio with the Experiments page selected, which looks like the following image (if not, click the Studio tab at the top of the page).

The screenshot shows the Microsoft Azure Machine Learning Studio (classic) interface. On the left, there's a sidebar with icons for Projects, Experiments (selected), Web Services, Datasets, Trained Models, and Settings. Below the sidebar, there's a 'NEW' button. The main area has tabs for 'MY EXPERIMENTS' and 'SAMPLES'. A table lists various experiments with columns for Name, Author, Status, Last Edited, Project, and a delete icon. To the right of the table is a visual experiment workspace. It shows two CSV files, 'eventdata.csv' and 'traind.csv', being joined. This leads to a series of modules: Two-Class Logistic Regression, Two-Class Boosted Decision Tree, Two-Class Support Vector Machine, and Two-Class Decision Forest. These are connected via 'Select Columns in Dataset' and 'Split Data' modules. The workspace has a dark background with light-colored boxes for components.

2. In the Studio, at the bottom left, click NEW. Then in the collection of Microsoft samples, select Blank Experiment. This creates a blank experiment, which looks similar to the following image.

The screenshot shows a new blank experiment in Microsoft Azure Machine Learning Studio (classic). The interface is similar to the previous one, with a sidebar on the left listing various sample modules: Saved Datasets, Data Format Conversions, Data Input and Output, Data Transformation, Feature Selection, Machine Learning, OpenCV Library Modules, Python Language Modules, R Language Modules, Statistical Functions, Text Analytics, and Time Series. The main workspace is currently empty, with a placeholder message: 'Experiment created on 5/5/2020' and 'To create your experiment, drag and drop datasets and modules here'. There are dashed boxes indicating where to drop items. On the right side, there are sections for 'Properties' (Status Code: In Draft), 'Experiment Properties' (Summary and Description fields), and 'Quick Help'. At the bottom, there are buttons for RUN HISTORY, SAVE, SAVE AS, DISCARD CHANGES, RUN, SET UP WEB SERVICE, and PUBLISH TO GALLERY.

3. Change the title of your experiment from “Experiment created on “today’s date” to “Project\_Logistic regression”.
4. Now we have to upload the dataset for that we have to click NEW at the bottom left. Then in the NEW dialog box, click the DATASET tab as shown in the following image.



5. Click FROM LOCAL FILE then in the upload a new dataset dialog box, browse to select the eventd.csv file from the folder where you extracted the lab files on your local computer and enter the following details as shown in the image below, and then click the OK icon.
  - This is a new version of an existing dataset: Unselected
  - Enter a name for the new dataset: eventd
  - Select a type for the new dataset: Generic CSV file with a header (.csv)
  - Provide an optional description: Diabetes patient appointments.

Note: May be file names are different in your case. Depending on what you gave during sampling of data in jupyter notebook.

**Upload a new dataset**

SELECT THE DATA TO UPLOAD:

eventd.csv

This is the new version of an existing dataset

ENTER A NAME FOR THE NEW DATASET:

eventd

SELECT A TYPE FOR THE NEW DATASET:

Generic CSV File with a header (.csv)

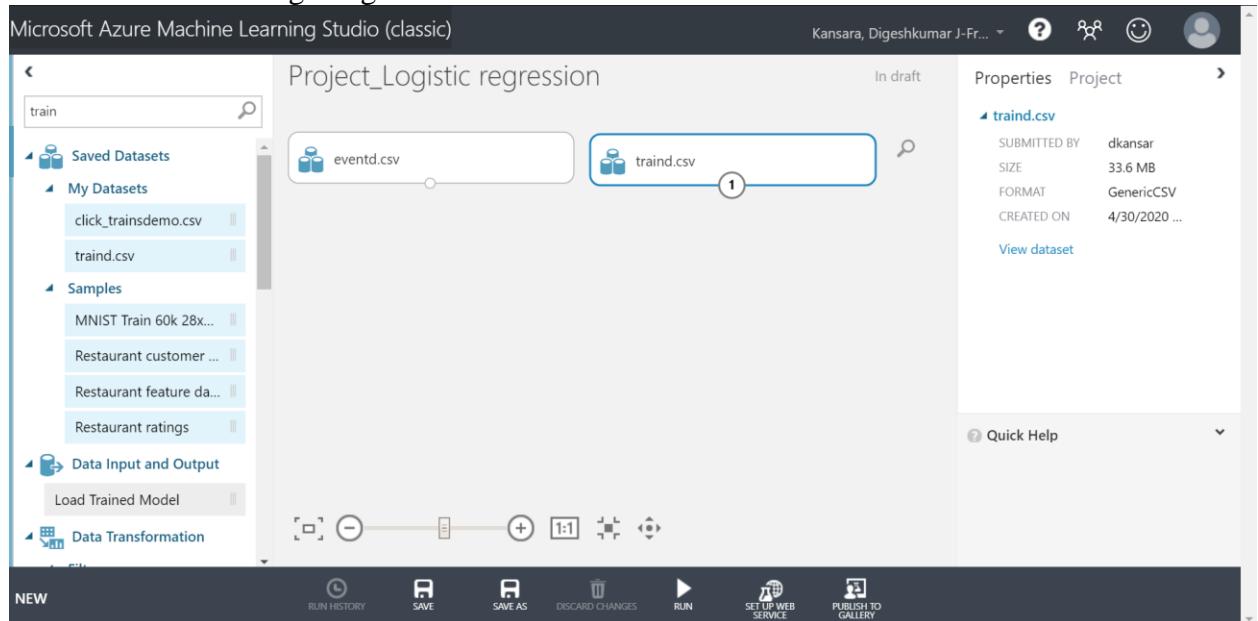
PROVIDE AN OPTIONAL DESCRIPTION:

Diabetes patient appointments

6. Wait for the upload of the dataset to be completed, and then on the experiment items pane, expand Saved Datasets and My Datasets to verify that the eventd dataset is listed.
7. Follow step 4-6 to upload traind and Promoted\_content.csv file.

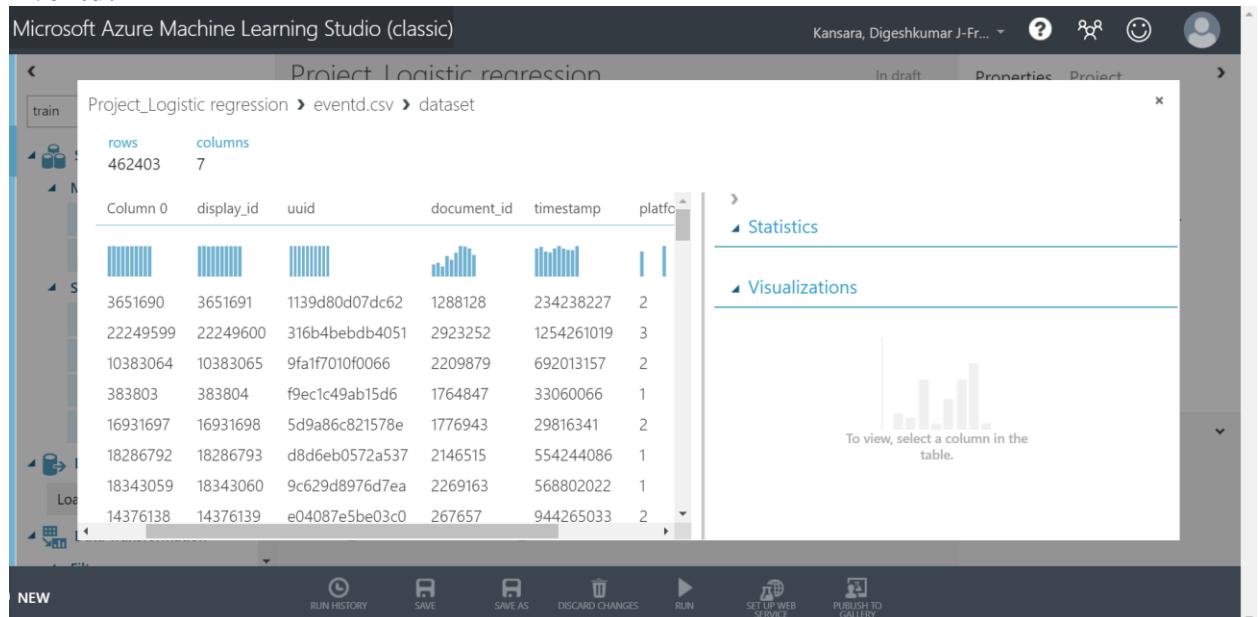
### Create Experiment with and joining three CSV files into one.

1. In the experiment items pane on the left, expand Saved Datasets, expand Samples, and drag eventd and traind (One by One) to the experiment canvas in the middle of the page, as shown in the following image.



2. Right click on eventd, traind data set and click visualize to view the data in dataset.

### Eventd:



## Traind:

Microsoft Azure Machine Learning Studio (classic)

Kansara, Digeshkumar J-Fr... Properties Project

Project\_Logistic regression > traind.csv > dataset

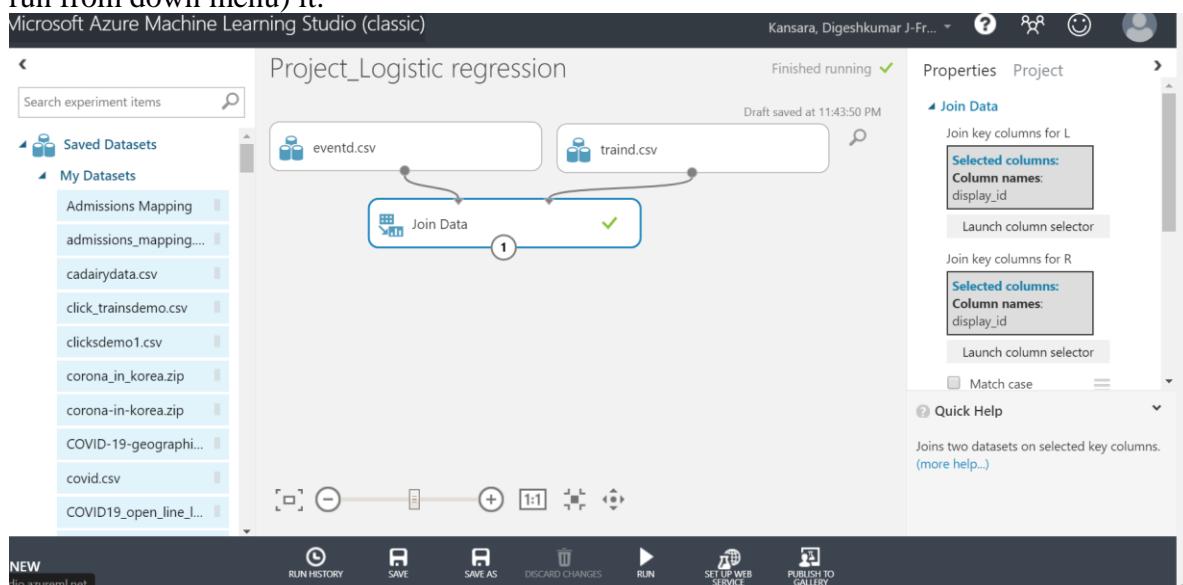
rows: 1307126 columns: 4

	Column 0	display_id	ad_id	clicked
58055587	11254709	252432	0	
48047511	9333749	405548	0	
66417980	12862694	332011	0	
59585503	11546676	284799	0	
64836639	12561207	45105	0	
68616222	13291651	270384	0	
27520623	5326485	141024	0	
22768702	4408986	235808	0	

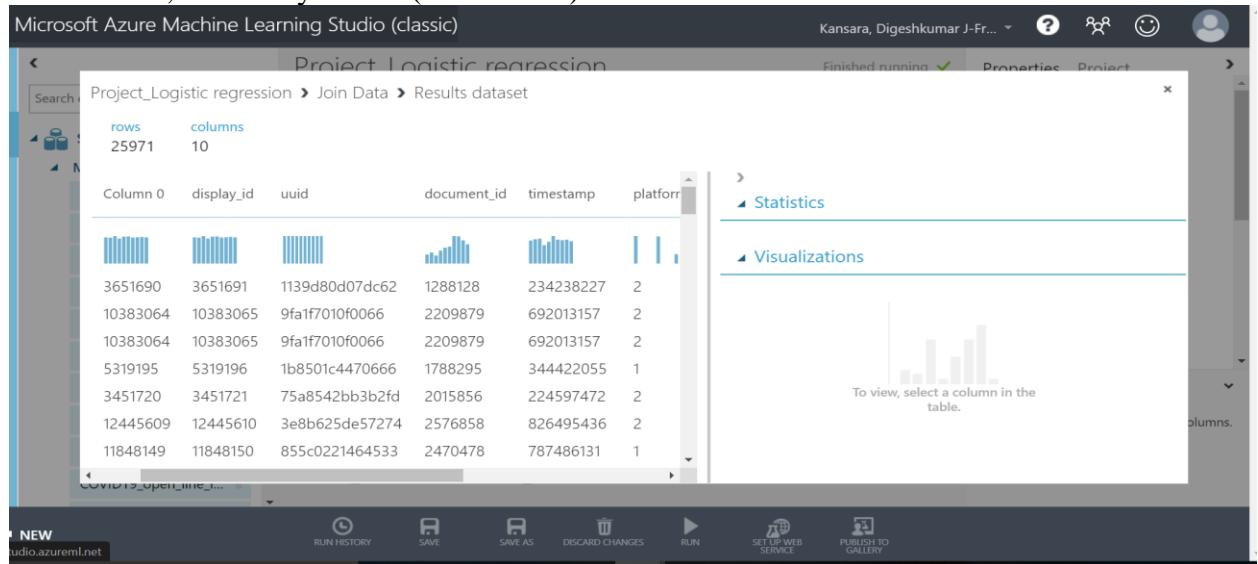
To view, select a column in the table.

Run History Save Discard Changes Run Set up Web Service Publish to Gallery

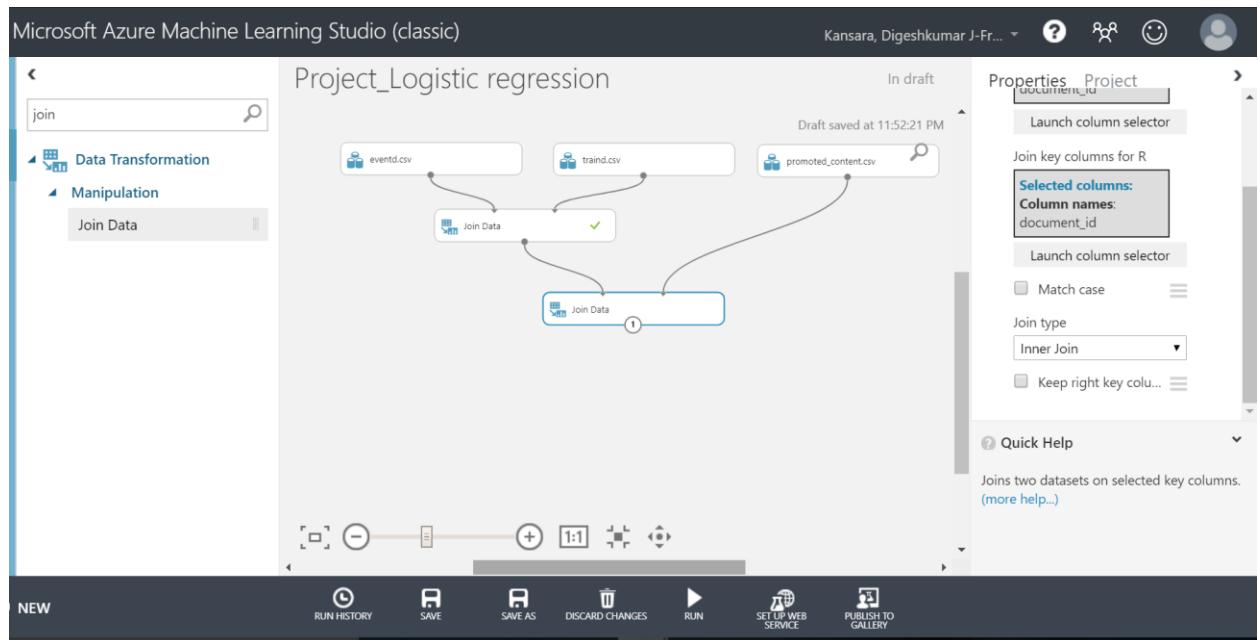
3. Close the dataset.
4. Search for the **Join** module and drag it to the canvas below both the **eventd** dataset and the **traind** dataset. Then connect the output ports from the **eventd** dataset and the **traind** dataset to the Dataset1 and Dataset2 input ports of the Join module respectively.
5. Select the **Join** module, and in the Properties pane, set the following properties:
  - a. **Join key columns for L:** Launch the column selector and select **display\_id**.
  - b. **Join key columns for R:** Launch the column selector and select **display\_id**.
  - c. **Join type:** Inner Join
  - d. **Keep right key column:** Unselected
6. Verify that your experiment resembles the following image, and then save and run (press run from down menu) it.



7. When the experiment has finished running, visualize the Results dataset output port of the Join module, and verify that it. (10 columns)



8. In the experiment items pane on the left, expand **Saved Datasets**, expand **My dataset**, and drag **Promoted\_content** to the experiment canvas in the right side of the page and bring join module in after first join module. Connect output of first join module to left input of second join module and output of promoted\_content module to right input of second join module. Verify your experiment connection as shown in following image.



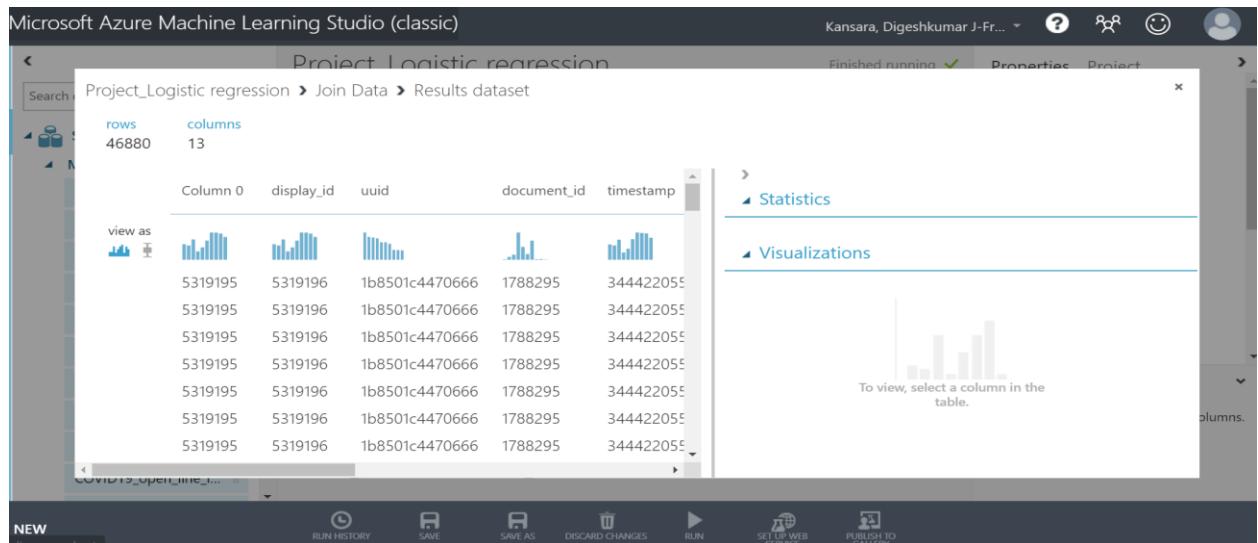
9. Select the Join module, and in the Properties pane, set the following properties:
- Join key columns for L:** Launch the column selector and select **document\_id**.
  - Join key columns for R:** Launch the column selector and select **document\_id**.

### c. Join type: Inner Join

d. **Keep right key column:** Unselected

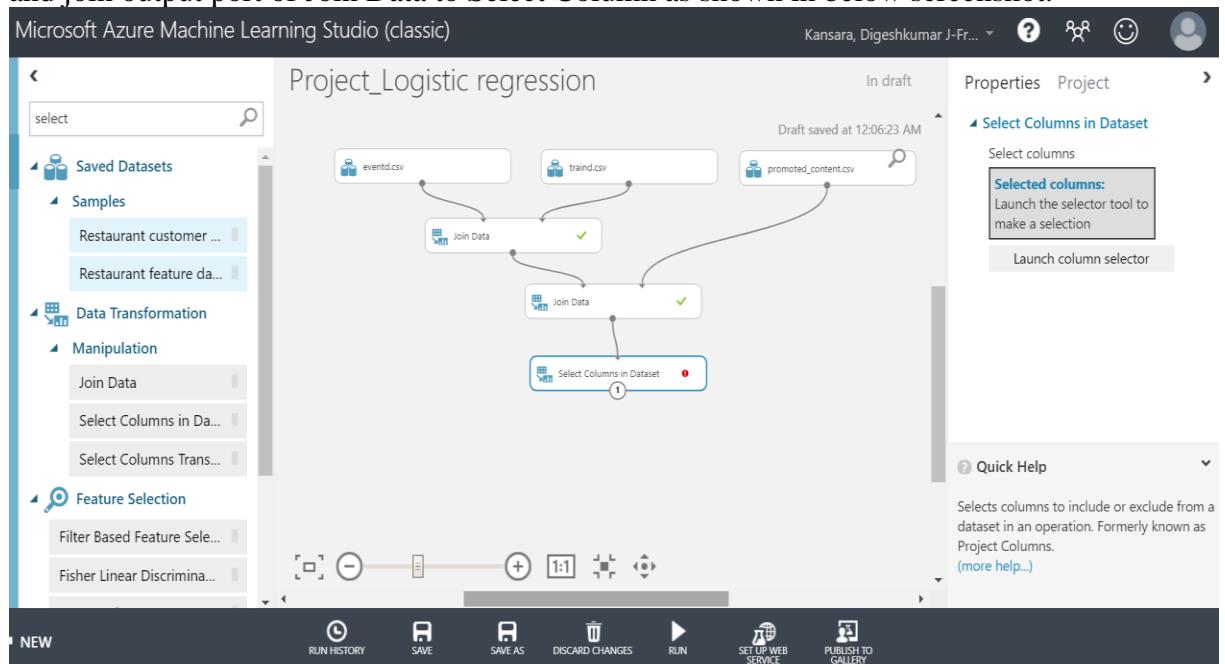
#### **10. Save and Run Experiment.**

11. When the experiment has finished running, visualize the result dataset output port of the Join module, and verify that it. (13 columns)

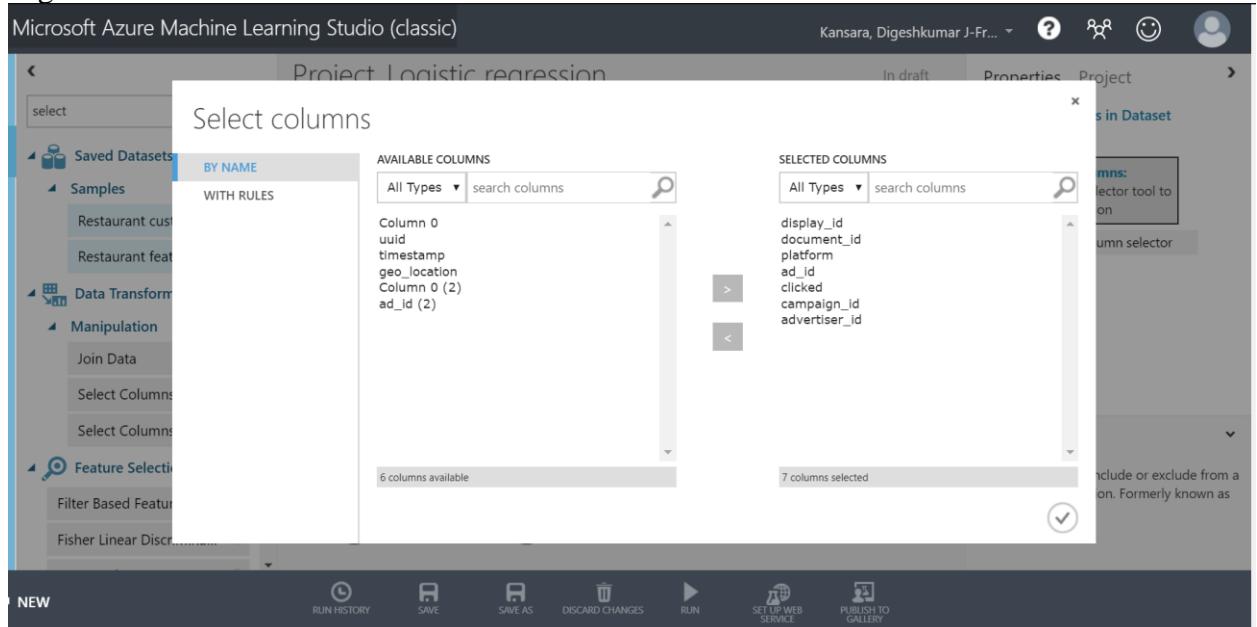


## Prepare data for training and testing

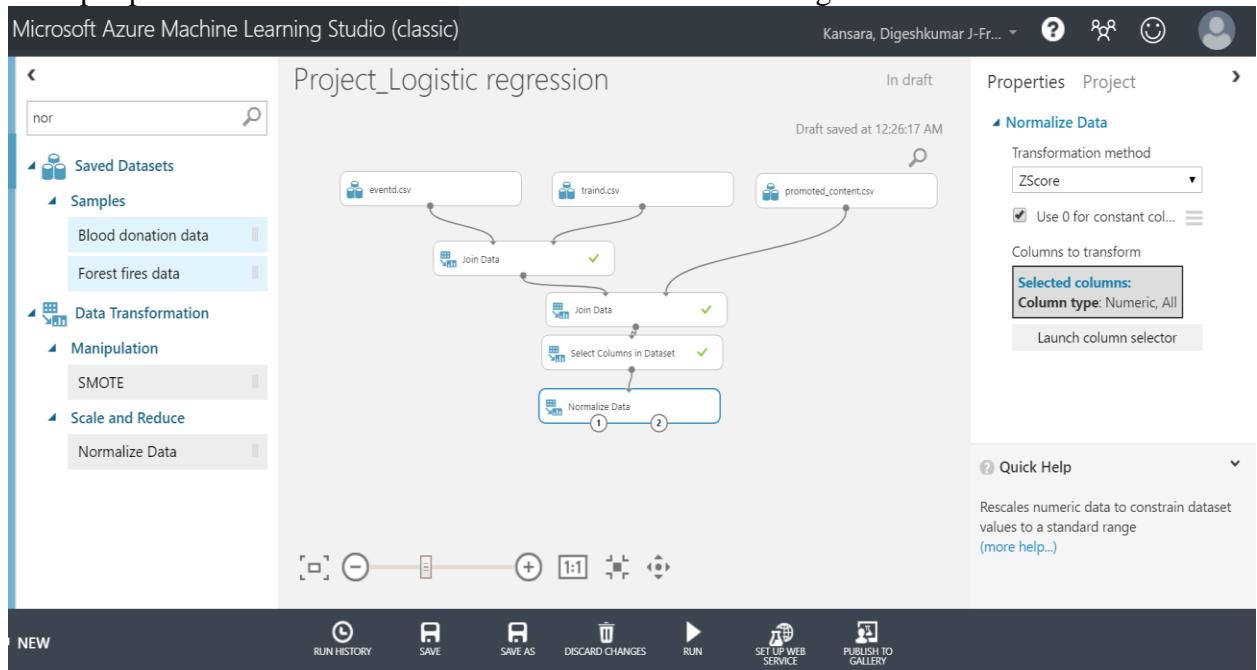
1. Search for **Project Columns (Select Columns in Data Set)** and drag it to the experiment and join output port of Join Data to Select Column as shown in below screenshot.



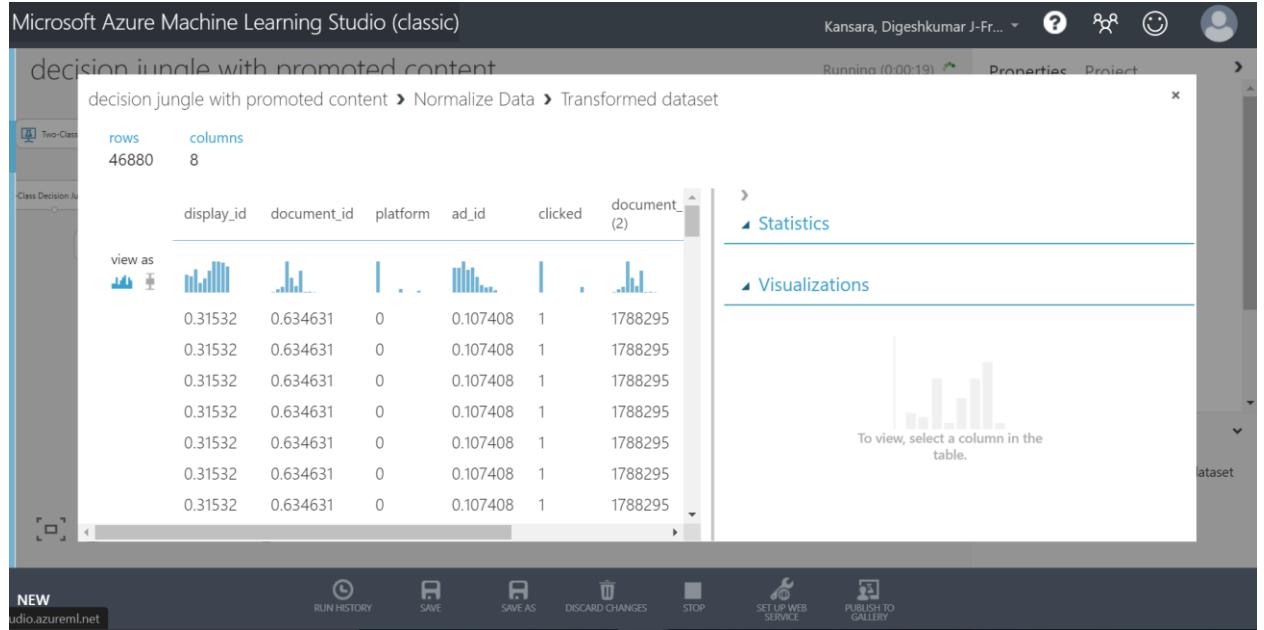
2. In the **Select columns** dialog box, select option “**BY Name**” and select **display\_id**, **document\_id**, **platform**, **ad\_id**, **clicked**, **campgin\_id**, **advertiser\_id** column names as shown in the image below. Then click the **OK** icon to close the column selector. Click on Right click button.



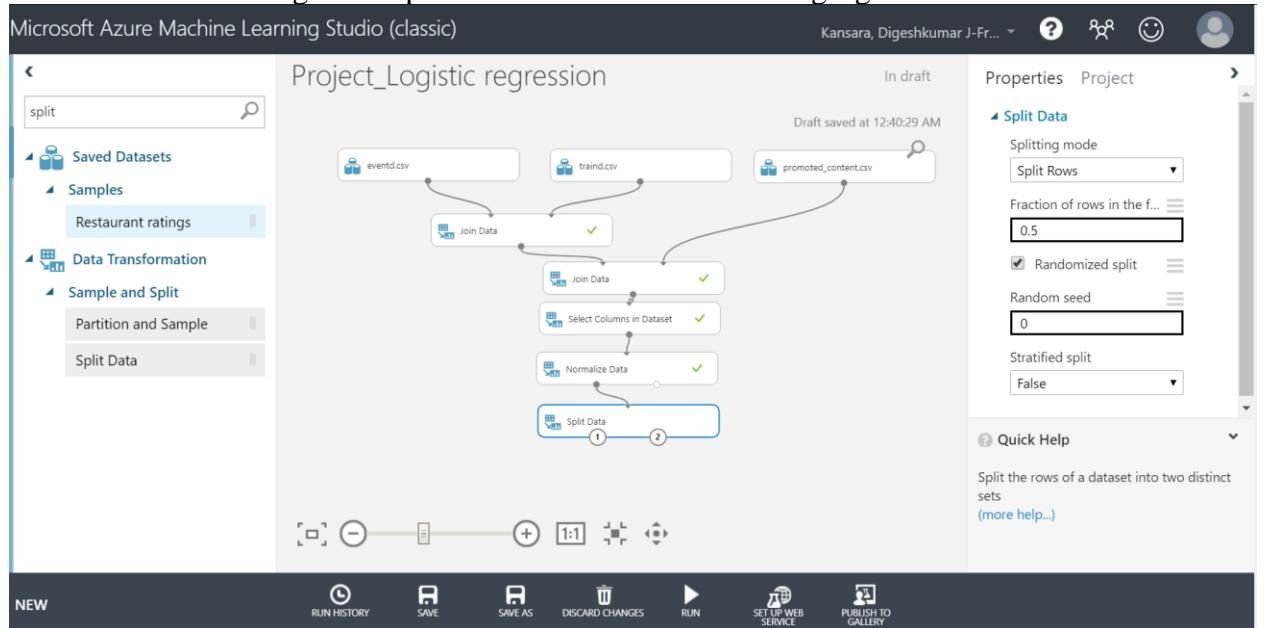
3. Save and Run Experiment.  
 4. Now we need to Normalize our dataset. From the left experiment pane search for **Normalize data** and drag to experiment. Connect output port of **select column dataset** to the input port of Normalize data module. As shown in following module.



- Select the **Normalize data** module, and in the Properties pane, set the following properties:
  - Transformation method:** Minmax.
  - Launch column selector:** select **By Name** and select **all columns**.
- Save and Run the experiment. Right click on Normalize and select visualize to analyze the data.



- Now we need to split data for training and testing. In left side experiment pane select **Split Data** module and bring it to experiment. As shown in following figure.



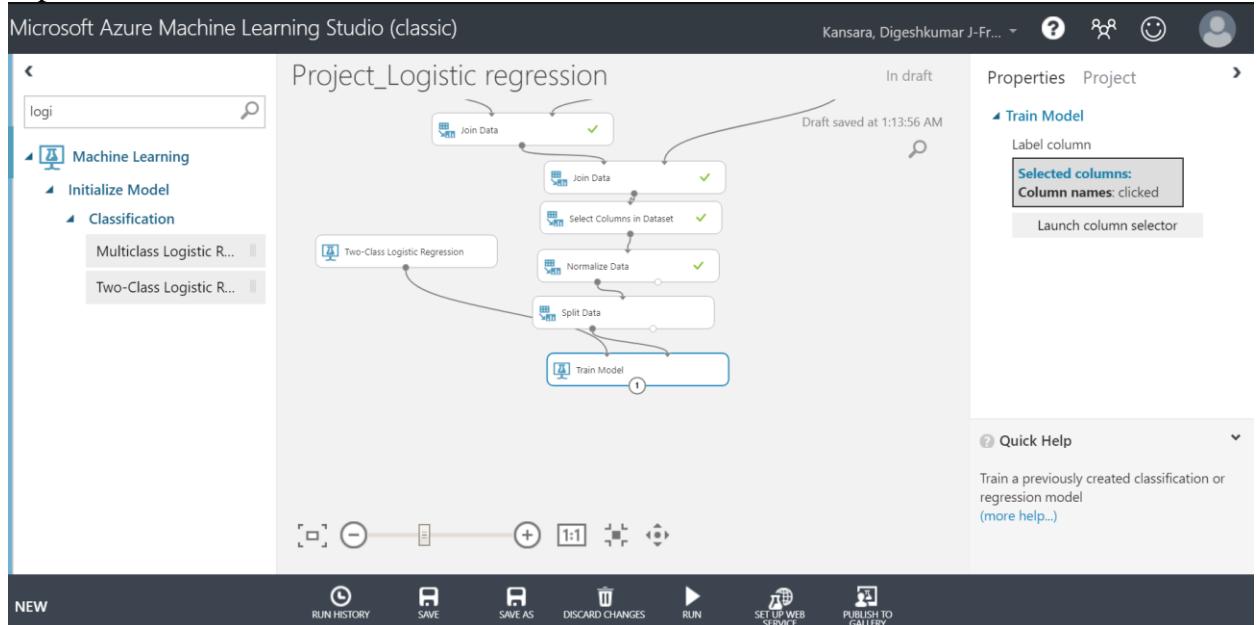
8. Select the **split module**, and in the Properties pane, set the following properties:
  - a. **Splitting mode:** splitting Rows.
  - b. **Fraction:** 0.7
  - c. **Randomized Split:** Selected
  - d. **Random Split:** 1234
9. Save and Run Experiment.

## Train data and measure accuracy and AUC of different Algorithm

### Algorithm 1: Two class Logistic Regression

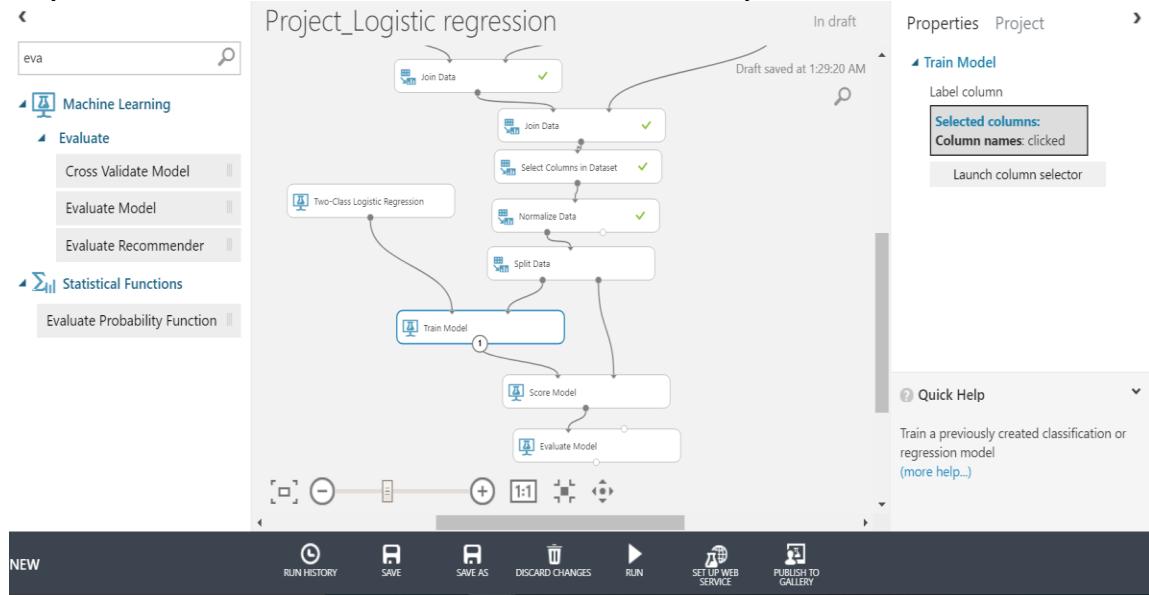
In previous step, we have created our training and test dataset. Next step is to measure accuracy of different Algorithm.

1. From left side experiment pane search for Two-class Logistic Regression module in the experiment and put it left side.
2. Bring **Train data** module in the experiment. Select launch column selector and select **clicked** column which we are going to predict.
3. Now bring Output of Two class Logistic Regression to left side of Train Data module and left output of split data module to right side of Train data module input. Now verify experiment as shown below.



4. Save and Run the Experiment.
5. Now we need to Score and Evaluate our model.
6. From left side experiment pane drag Score module to the experiment. Connect output of train module to left input of score module and right output of split module to right input of score module.

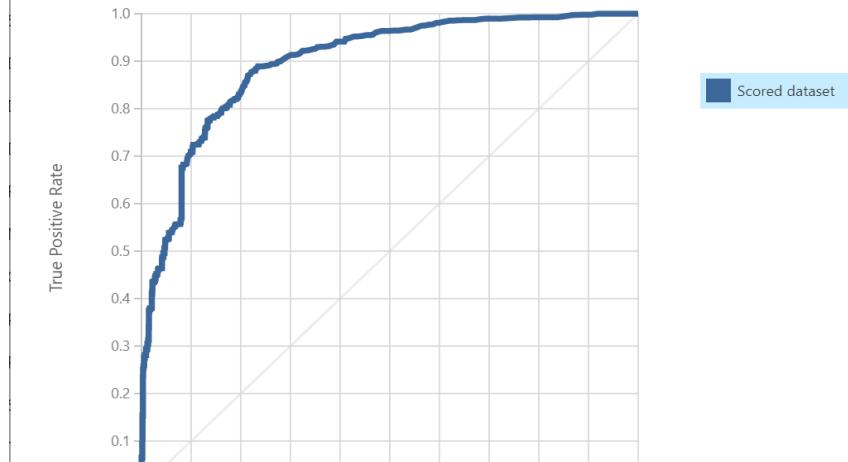
7. From Left side experiment pane drag Evaluate Model module in experiment and connect output module of score module to Evaluate model. The experiment should look like below.



8. Save and Run the experiment.  
 9. Now right click on Evaluate model and select visualize and check output as follows.

### AUC of Algorithm

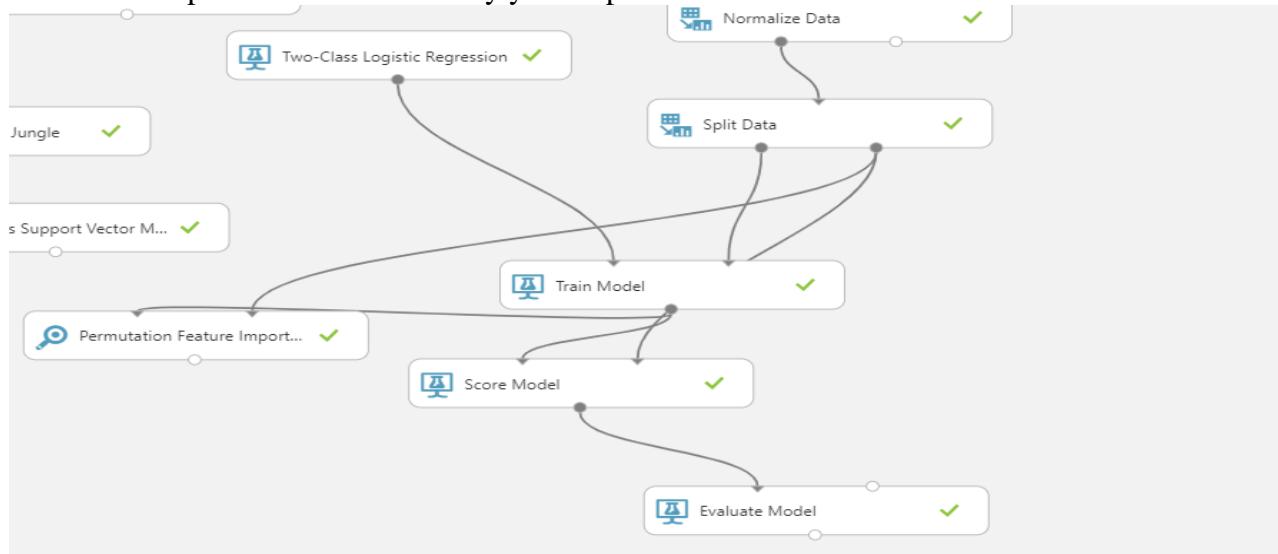
Logistic Regression > Evaluate Model > Evaluation results



### Accuracy of algorithm.

True Positive	False Negative	Accuracy	Precision	Threshold	AUC		
866	1251	0.893	0.775	0.5	0.898		
False Positive	True Negative	Recall	F1 Score				
252	11695	0.409	0.535				
Positive Label		Negative Label					
1	0						

**10. Now Bring Permutation Feature Importance** module to left side of experiment after training data. From Train Data Model output to right input of permutation feature importance module and connect right side of split data output to left input of permutation feature importance module. Verify your experiment as follows.



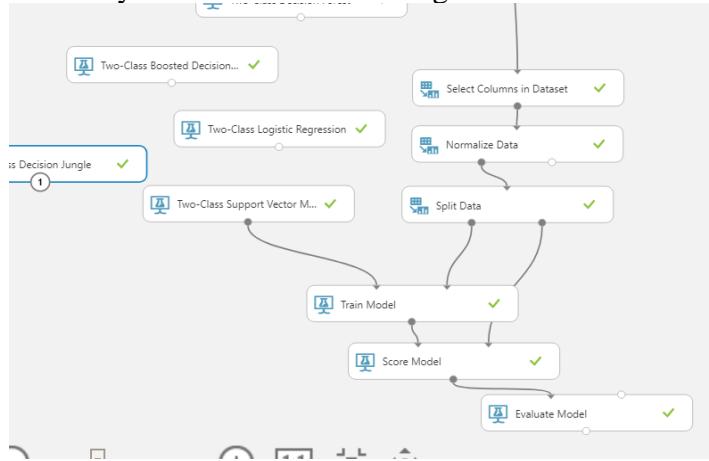
**11. Save and Run Experiment.** Now right click on permutation feature important module and click visualize. You can see result as following.

rows	columns
8	2
<hr/>	
	Feature
	Score
<hr/>	
view as	
	geo_location 0.096701
	timestamp 0.005759
	platform 0.003057
	display_id 0.00128
	ad_id 0.000427
	document_id -0.000071
	campaign_id -0.000711
	advertiser_id -0.000924

#### **Algorithm 2: Support Vector Machine (SVM)**

1. We have completed Ad prediction Using Logistic Regression. Now we are going to use support vector machine algorithm.
2. Open Logistic\_Algorithm that we have created in previous step from left side experiment side pane find two class support vector machine and bring it to experiment module.

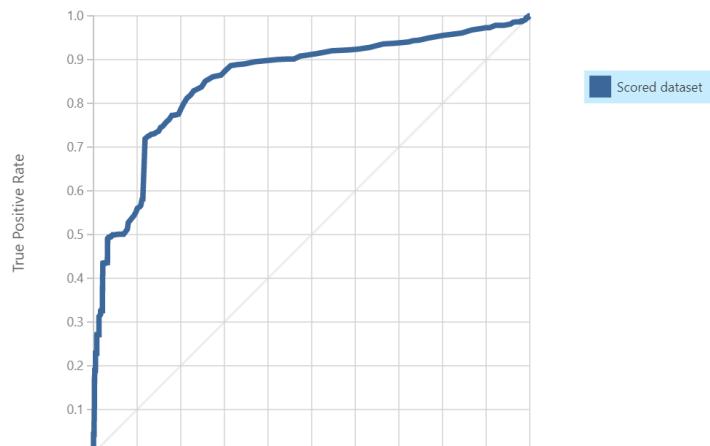
3. Now disconnect between Two class logistic regression with train data.
4. Connect output of Two class Support vector machine to left side of train data. You can verify connection as following.



5. Save and run Experiment. Visualize Accuracy and AUC as following.

#### AUC:

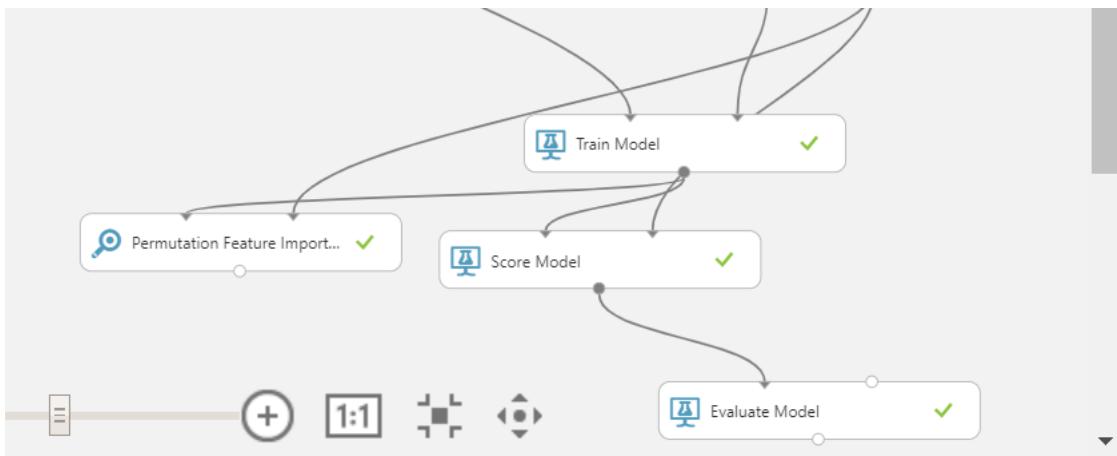
Support Vector Machine > Evaluate Model > Evaluation results



#### Accuracy of Algorithm:

True Positive	False Negative	Accuracy	Precision	Threshold	AUC
704	1413	0.882	0.742	0.5	0.855
Positive Label	Negative Label				
1	0				
False Positive	True Negative	Recall	F1 Score		
245	11702	0.333	0.459		

6. Now Bring Permutation Feature Importance module to left side of experiment after training data. From Train Data Model output to right input of permutation feature importance module and connect right side of split data output to left input of permutation feature importance module. Verify your experiment as follows.

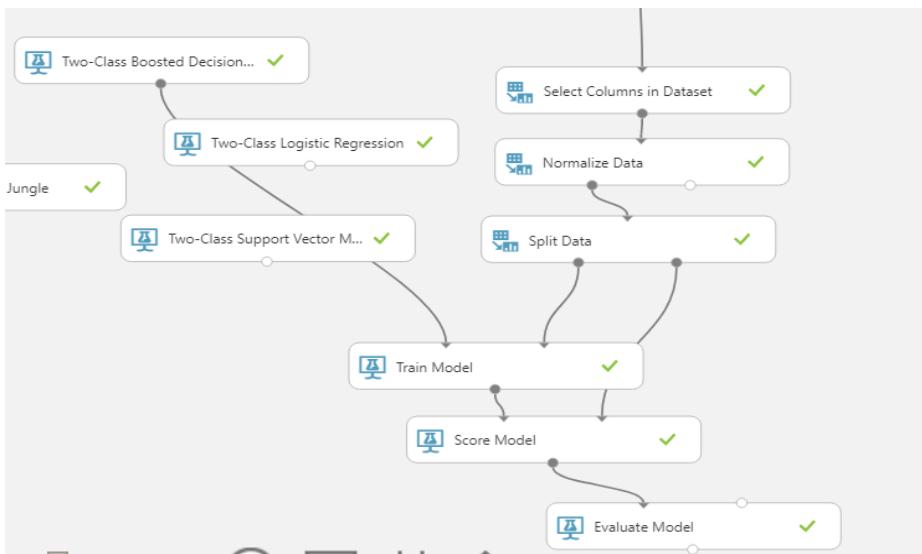


7. Save and Run Experiment. Now right click on permutation feature important module and click visualize. You can see result as following.

rows	columns
8	2
<hr/>	
	Feature      Score
view as	
	geo_location      0.078214
	display_id      0
	document_id      0
	timestamp      0
	platform      0
	ad_id      0
	campaign_id      0
	advertiser_id      0

### Algorithm 3: Two class Boosted Decision Tree

1. We have completed Ad prediction using Support Vector Machine. Now We are going to use Two class Boosted Decision Tree algorithm.
2. Open Logistic\_Algorithm that we have created in previous step from left side experiment side pane find two class Boosted Decision Tree and bring it to experiment module.
3. Now disconnect between Two class Boosted Decision Tree with train data.
4. Connect output of Two class Boosted Decision Tree to left side of train data. You can verify connection as following.



5. Check Property of Two class Boosted Decision Tree from right properties pane as shown in follow.

**Two-Class Boosted Decision ...**

Create trainer mode

Single Parameter ▾

Maximum number of le...

20

Minimum number of sa...

10

Learning rate

0.2

Number of trees constr...

100

Learning rate

0.2

Number of trees constr...

100

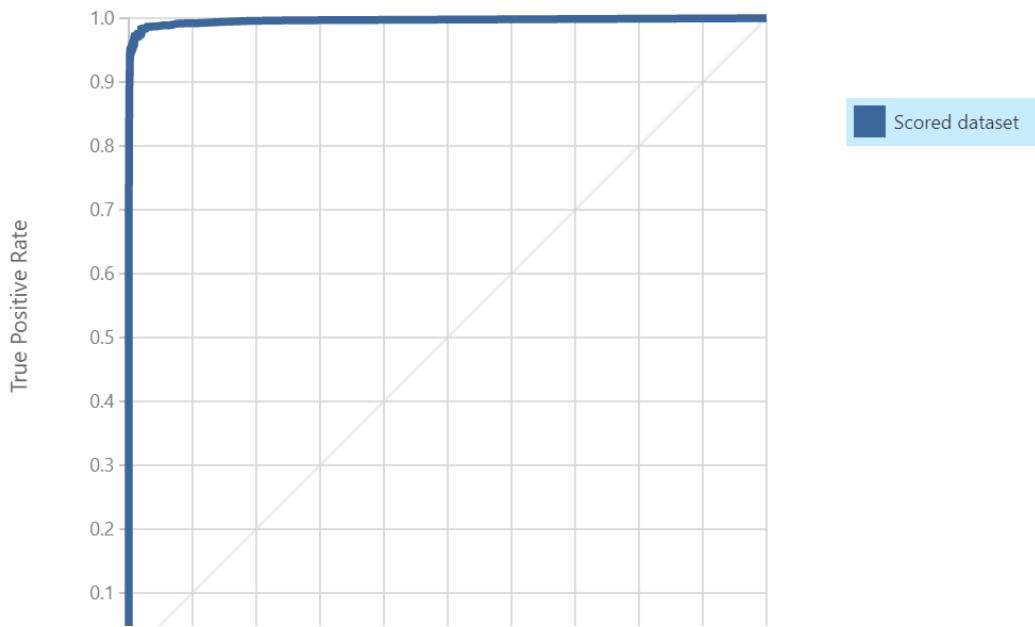
Random number seed

Allow unknown cat...

6. Save and run the experiment. Visualize Accuracy and AUC as following.

**AUC:**

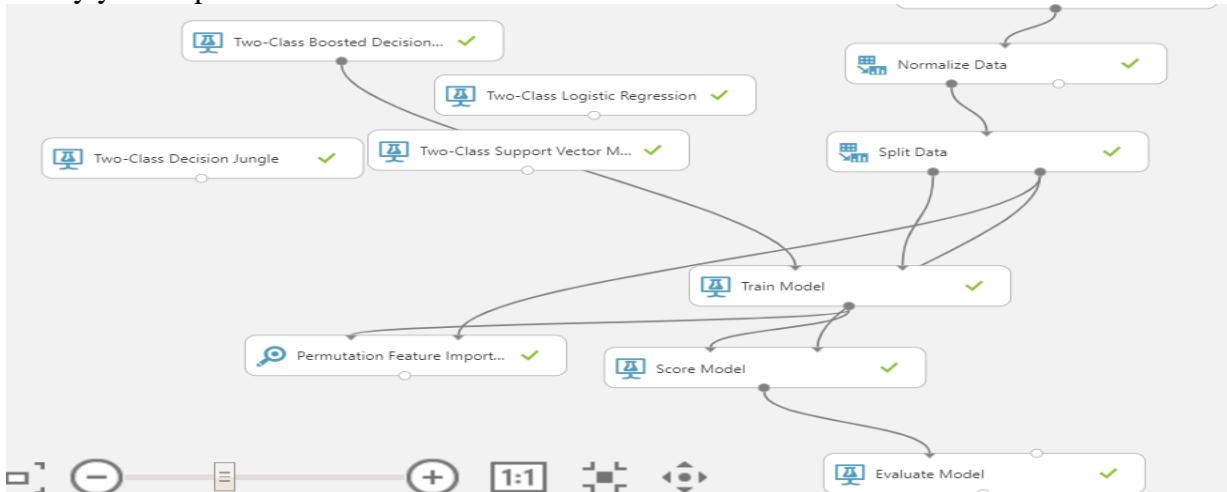
Two class Boosted Decision → Evaluate Model → Evaluation results



### Accuracy:

True Positive	False Negative	Accuracy	Precision	Threshold	AUC
2020	97	<b>0.989</b>	<b>0.969</b>	0.5	0.996
False Positive	True Negative	Recall	F1 Score		
64	11883	<b>0.954</b>	<b>0.962</b>		
Positive Label	Negative Label				
1	0				

**7. Now Bring Permutation Feature Importance** module to left side of experiment after training data. From Train Data Model output to right input of permutation feature importance module and connect right side of split data output to left input of permutation feature importance module. Verify your experiment as follows.

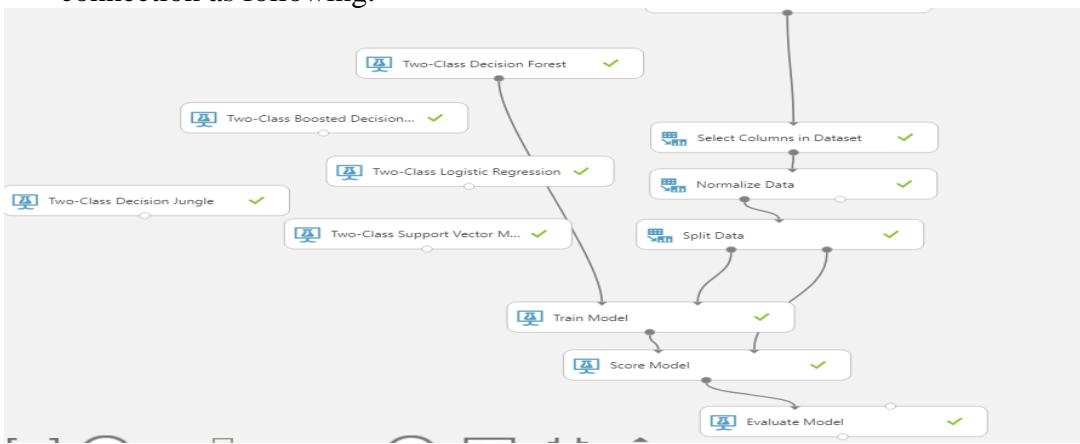


8. Save and Run Experiment. Now right click on permutation feature important module and click visualize. You can see result as following.

rows	columns
7	2
Feature	Score
ad_id	0.157992
display_id	0.096061
document_id	0.046786
advertiser_id	0.014078
campaign_id	0.00583
platform	0.002702
timestamp	0

#### Algorithm 4: Two Class Decision Forest

1. Now we are going to use Two class Boosted Decision Tree algorithm.
2. Open Logistic\_Algorithm that we have created in previous step from left side experiment side pane find two class Decision Forest and bring it to experiment module.
3. Now disconnect between two class Decision Forest with train data.
4. Connect output of two class Decision Forest to left side of train data. You can verify connection as following.



5. Check Property of Two class Decision Forest from right Properties pane as shown in follow.

Properties Project

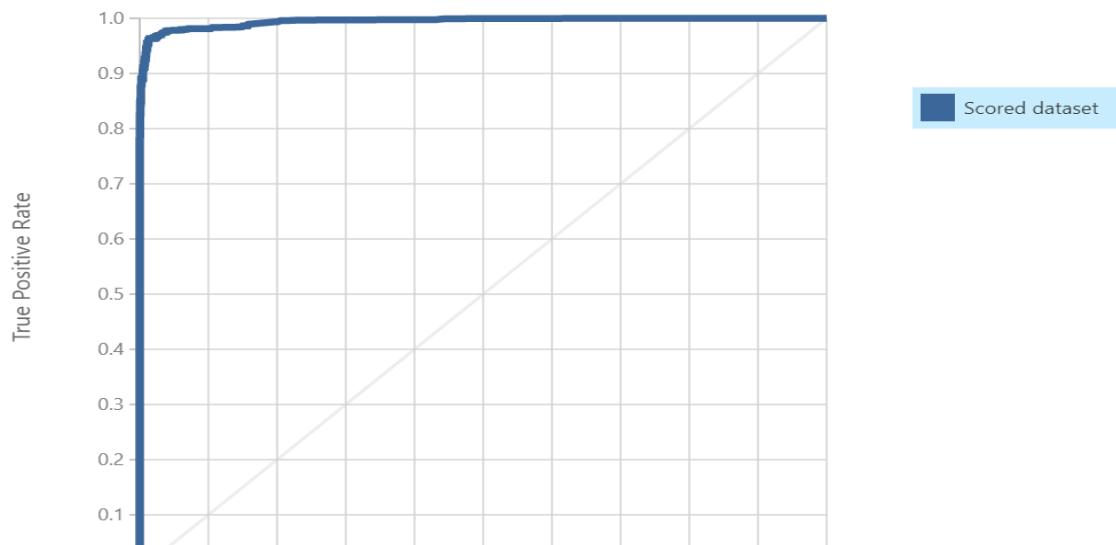
Two-Class Decision Forest

Resampling method	Bagging
Create trainer mode	Single Parameter
Number of decision trees	8
Maximum depth of the ...	32
Number of random spli...	128
Number of random spli...	128
Minimum number of sa...	1
<input checked="" type="checkbox"/> Allow unknown val...	

6. Save and run Experiment. Visualize Accuracy and AUC as following.

#### AUC:

Two Class Decision Forest ➤ Evaluate Model ➤ Evaluation results

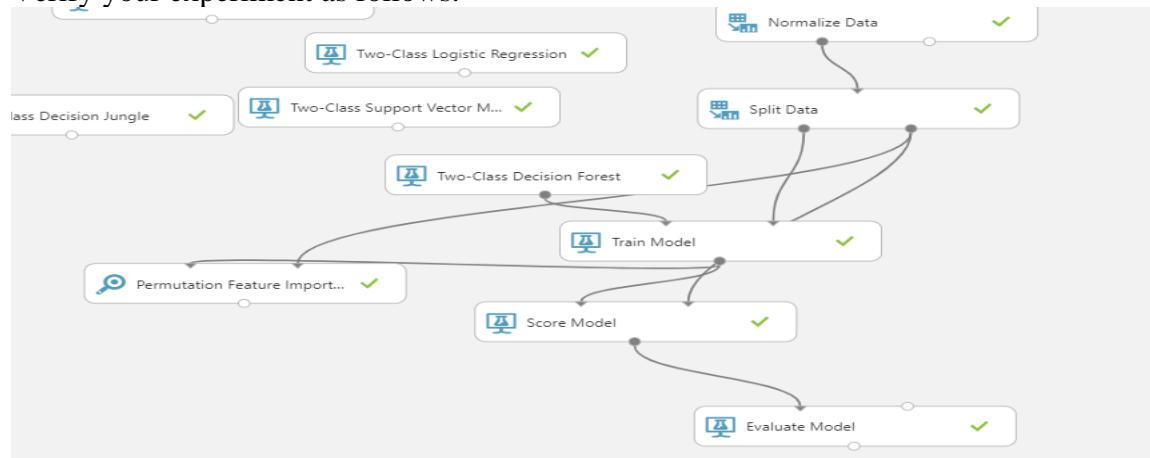


## Accuracy:

7.

True Positive	False Negative	Accuracy	Precision	Threshold	AUC
1753	364	0.974	0.996	0.5	0.995
False Positive	True Negative	Recall	F1 Score		
7	11940	0.828	0.904		
Positive Label	Negative Label				
1	0				

**7. Now Bring Permutation Feature Importance** module to left side of experiment after training data. From Train Data Model output to right input of permutation feature importance module and connect right side of split data output to left input of permutation feature importance module. Verify your experiment as follows.



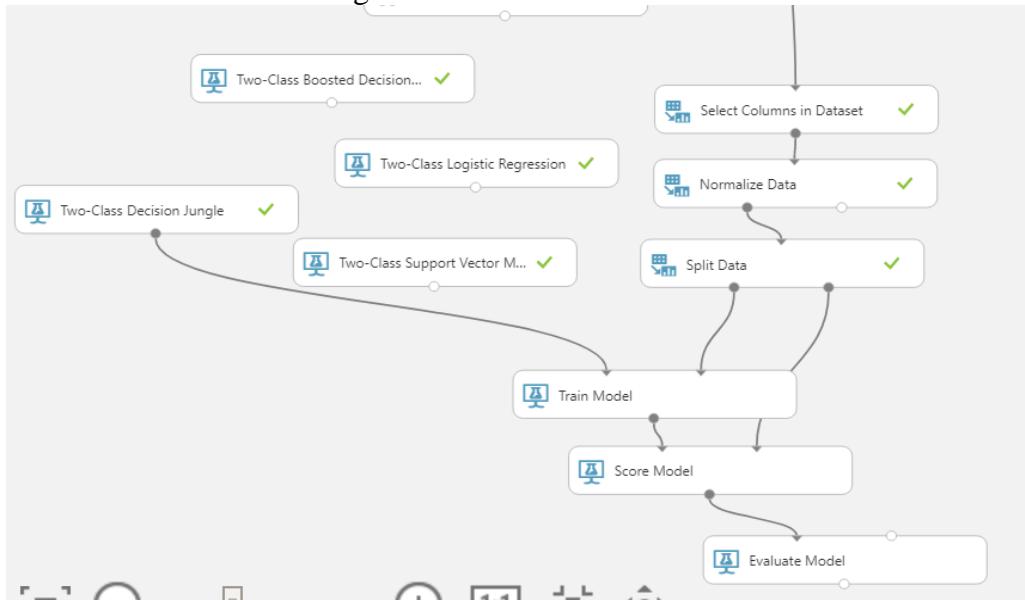
**8. Save and Run Experiment.** Now right click on permutation feature important module and click visualize. You can see result as following.

rows	columns
7	2
<hr/>	
Feature	Score
ad_id	0.140785
timestamp	0.084044
display_id	0.054892
document_id	0.047426
advertiser_id	0.045577
platform	0.01351
campaign_id	0.002062

---

### Algorithm 5: Two Class Decision Jungle

1. We have completed Ad prediction using Two Class Decision Forest. Now we are going to use Two Class Decision Jungle algorithm.
2. Open Logistic\_Algorithm that we have created in previous step from left side experiment side pane find two class Decision Jungle and bring it to experiment module.
3. Now disconnect between two class Decision Jungle with train data.
4. Connect output of two class Decision Jungle to left side of train data. You can verify connection as following.



5. Check Property of Two class Decision Jungle from right Properties pane as shown in follow.

Properties Project

Two-Class Decision Jungle

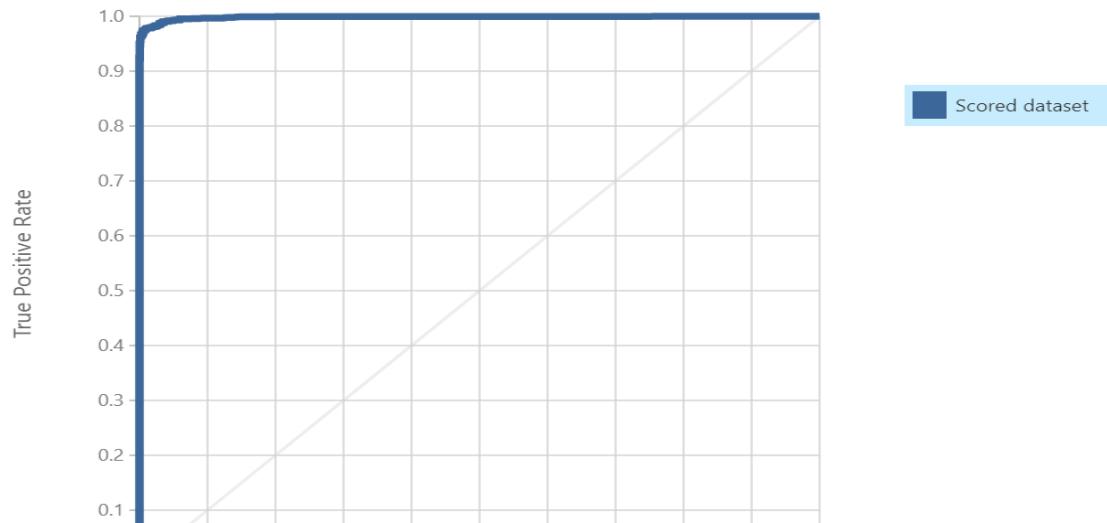
Resampling method	Bagging
Create trainer mode	Single Parameter
Number of decision DA...	8
Maximum depth of the ...	32
Maximum width of the ...	128

Maximum depth of the ...	<input type="text" value="32"/>
Maximum width of the ...	<input type="text" value="128"/>
Number of optimizatio...	<input type="text" value="2048"/>
<input checked="" type="checkbox"/> Allow unknown val...	<input type="text"/>

6. Save and run the experiment. Visualize Accuracy and AUC as following.

### AUC:

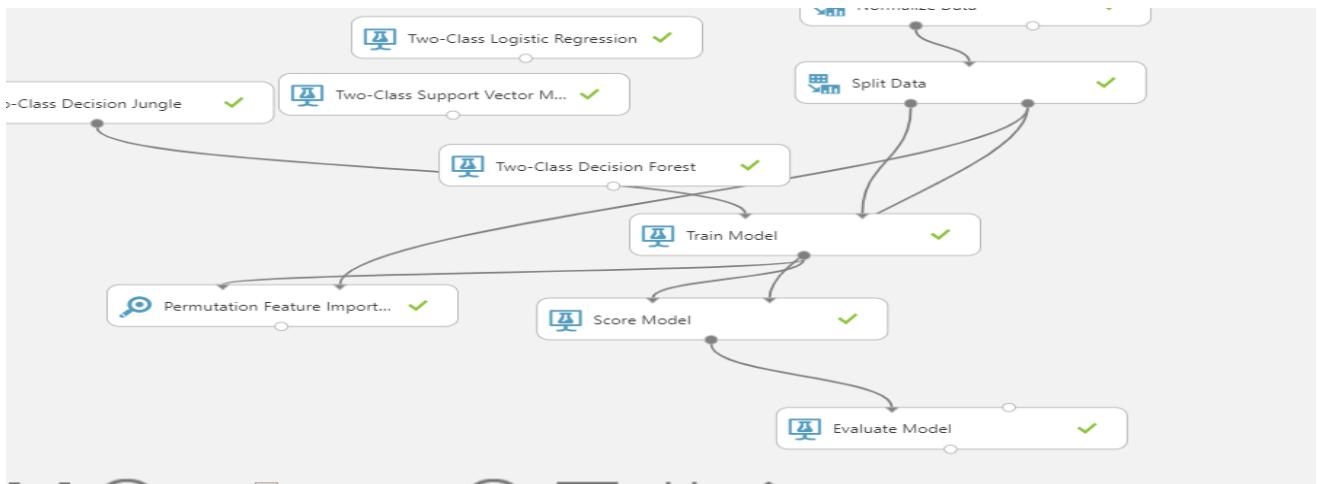
Two Class Decision Jungle ➔ Evaluate Model ➔ Evaluation results



### Accuracy:

True Positive	False Negative	Accuracy	Precision	Threshold	AUC
1974	143	0.989	0.997	0.5	0.998
False Positive	True Negative	Recall	F1 Score		
6	11941	0.932	0.964		
Positive Label		Negative Label			
1	0				

7. Now Bring Permutation Feature Importance module to left side of experiment after training data. From Train Data Model output to right input of permutation feature importance module and connect right side of split data output to left input of permutation feature importance module. Verify your experiment as follows.



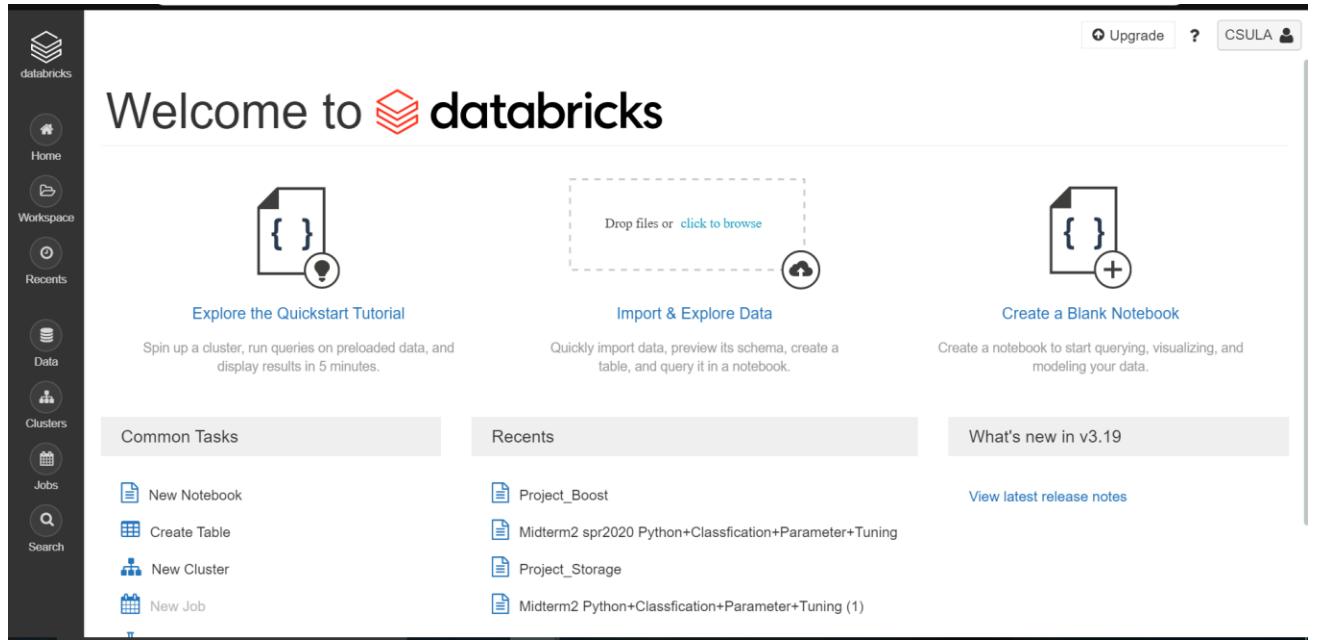
8. Save and Run Experiment. Now right click on permutation feature important module and click visualize. You can see result as following.

Two Class Decision Jungle ➤ Permutation Feat

rows	columns
7	2
<hr/>	
Feature	Score
ad_id	0.130902
display_id	0.066695
document_id	0.058092
timestamp	0.055958
advertiser_id	0.044937
platform	0.023251
campaign_id	0.008604

## Step 4: Ad prediction using spark in Databricks

- To Work in Databricks we need to navigate to Databricks Community Edition to login. If you **don't have** databricks community edition account, you need to first Sign Up for Community Edition.  
<https://community.cloud.databricks.com/login.html?o=1035240001367446#notebook/2698979388235870/command/2698979388235882>
- When we login to databricks. We need to create a cluster. To create a cluster, we need to click on clusters button from left side pane as shown in following image.



3. Now it will navigate to **create cluster page** as shown in below. click on create cluster button.

The screenshot shows the 'Create Cluster' page. At the top, it says 'Create Cluster' and 'New Cluster'. There is a 'Cancel' button and a 'Create Cluster' button. Below that, there is a 'Cluster Name' input field with a placeholder 'Please enter a cluster name'. To the right of the input field are 'UI' and 'JSON' buttons. Under 'Databricks Runtime Version', it says 'Runtime: 6.5 (Scala 2.11, Spark 2.4.5)'. A note says 'New This Runtime version supports only Python 3.' Below that is an 'Instance' section with a note: 'Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours. For more configuration options, please upgrade your Databricks subscription.' At the bottom, there are tabs for 'Instances' (selected) and 'Spark'. An 'Availability Zone' dropdown is set to 'us-west-2c'.

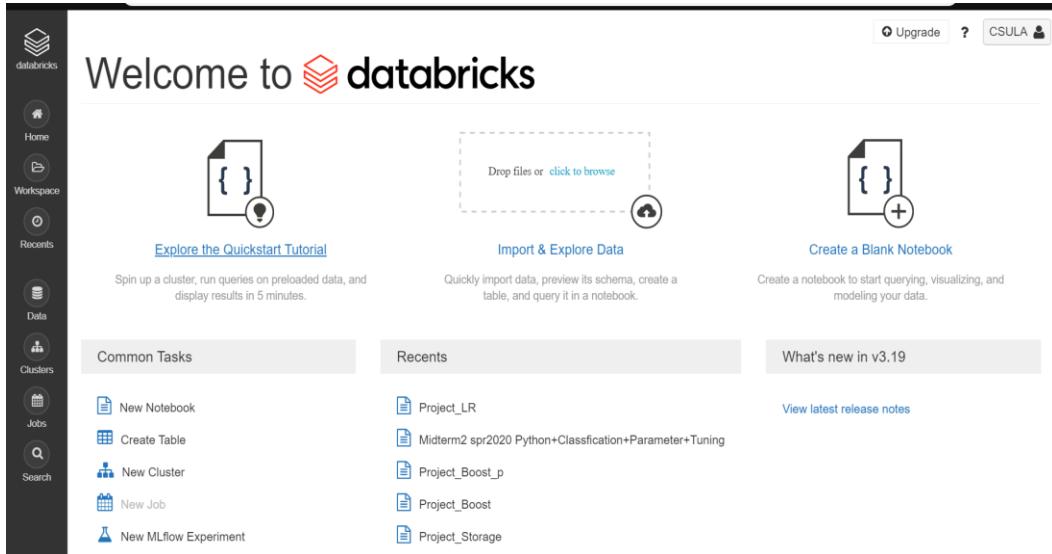
**Cluster Name:** Choose combination of any letter.

**Databricks Runtime Version:** Choose compatible to your project and system and click on create cluster button from top as it become visible after we enter cluster name and cluster version.

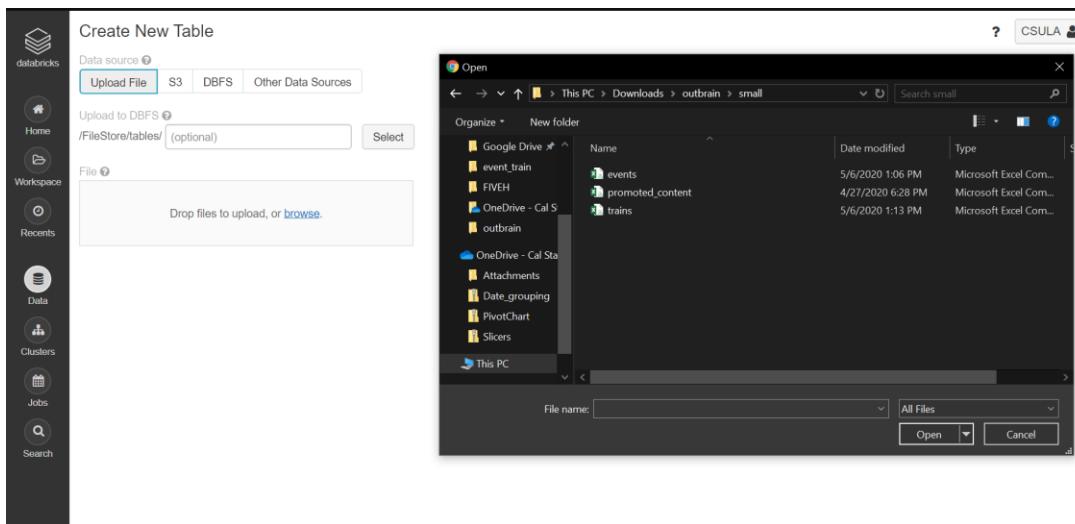
It will take 15-20 minute to create a cluster.

- After launching a cluster, select “Create Table” at the Databricks UI/GUI in order to upload the file **trains.csv**, **events.csv** and **Promoted\_content.csv**.

**Note:** I changed my file name trains.csv for some reasons. You can upload csv files with name which you gave during sampling in jupyter notebook in stage 2.



- It will go to the following page. You need to select “browse” and locate the file to upload at the file explorer:



**Note:** File Size is 35MB it will take about 15 minutes.

- Create a table from the uploaded csv file.
  - Select “Create Table with UI” then select the cluster from the drop-down box, which you created above.
  - Then, click Preview Table button

Create New Table

Upload to DBFS [?](#)

/FileStore/tables/ (optional) [Select](#)

File [?](#)

events.csv ✓  
35.7 MB [Remove file](#)

✓ File uploaded to /FileStore/tables/events.csv

[Create Table with UI](#) [Create Table in Notebook](#) [?](#)

Select a Cluster to Preview the Table

Choose a cluster with which you will read and preview the data.

Cluster [?](#) Project [▼](#)

[Preview Table](#)

7. It will display the portion of table columns, data types, and the values. You have to set up the data type and mark the head option as you have seen at ipynb at for example:
  - Table Name: table2015\_csv
  - Mark “First row is header”
  - Mark “Infer Schema”
  - Then, select **Create Table** button

8. It will show the following table “events\_csv” created.

col_name	data_type	comment
_c0	int	null
display_id	int	null
uuid	string	null
document_id	int	null
timestamp	int	null
platform	int	null
geo_location	string	null

_c0	display_id	uuid	document_id	timestamp	platform	geo_location
3651690	3651691	1139d80d07dc62	1288128	234238227	2	US>MO
22249599	22249600	316b4bebdb4051	2923252	1254261019	3	US>NY>501
10383064	10383065	9fa1f7010f0066	2209879	692013157	2	US>CA>807
383803	383804	f9ec1c49ab15d6	1764847	33060066	1	US>TX>635

**Now follow step from 4 to 8 and create table for trains, Promoted\_content.**

9. You have to create a folder named cis5560 of Data Bricks account, which should be at the left frame of the page at: Workspace > Users > Your Account Email > cis5560.

The screenshot shows the Databricks workspace interface. On the left is a dark sidebar with icons for Home, Workspace, Recents, Data, and Clusters. The main area is titled "Workspace" and contains two dropdown menus: "Workspace" and "Users". The "Workspace" menu has options for "Shared" and "Users", with "Users" currently selected. The "Users" menu shows one entry: "dkansar@calstatela.edu". Below these menus is a list of notebook files:

- final
- Fine Grained Demand ...
- joinplastic
- Loginistcn
- midterm
- Midterm2 spr2020 Python...
- oracle
- oracle join

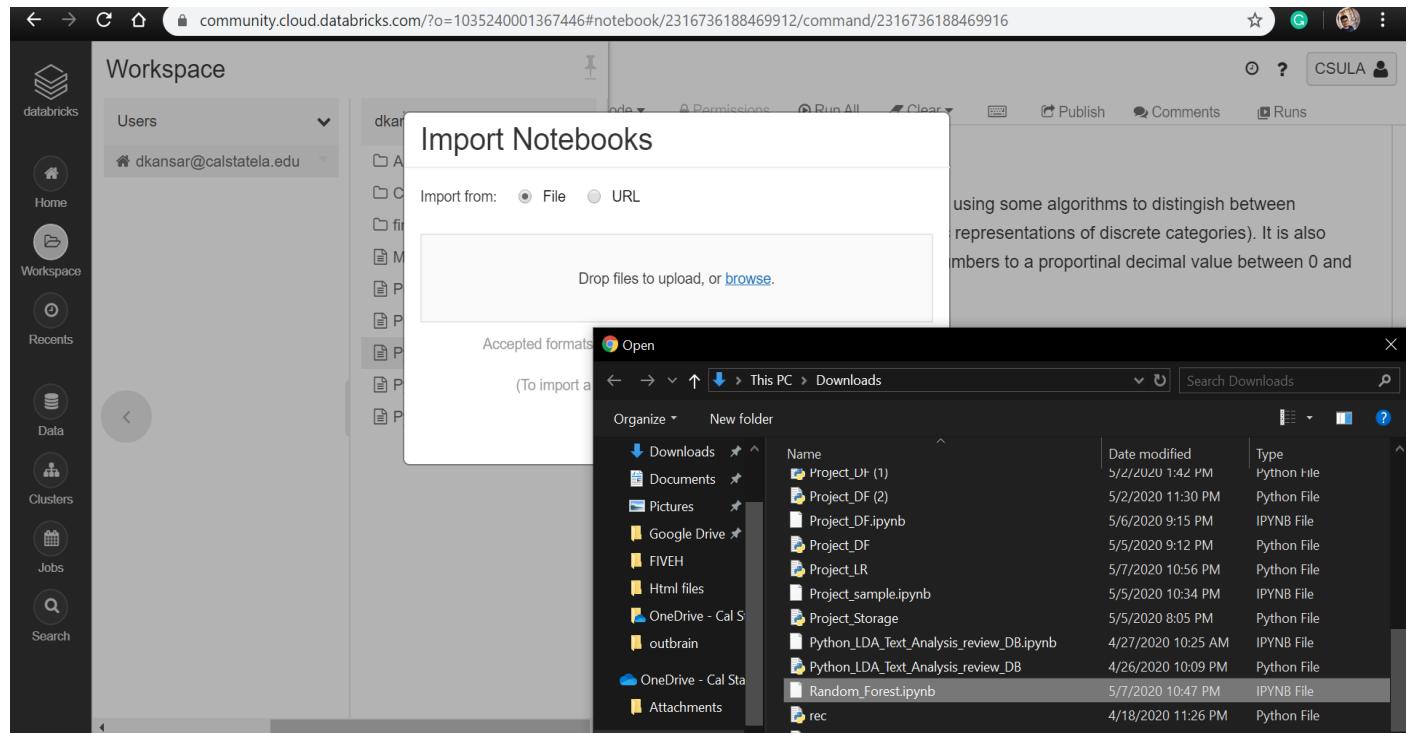
## Random Forest Algorithm:

10. Import “Random\_forest.ipynb” to the folder **cis5560: cis5560 > import**

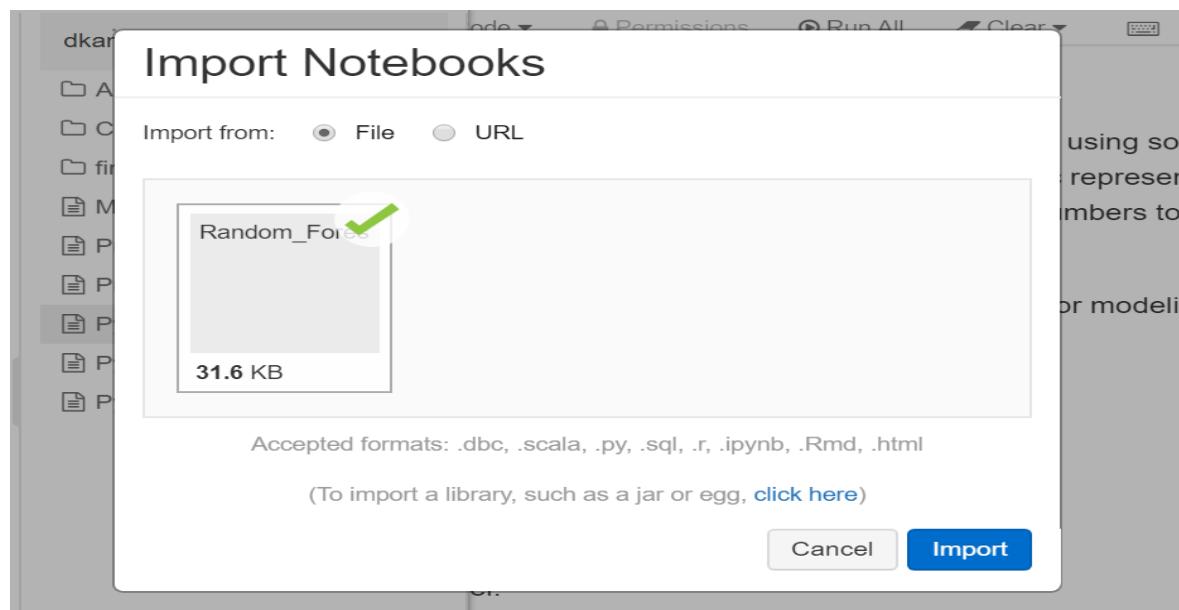
The screenshot shows the Databricks workspace interface. On the left is a dark sidebar with icons for Home, Workspace, Recents, Data, and Clusters. The main area is titled "Workspace" and contains two dropdown menus: "Users" and "Permissions". The "Users" menu has one entry: "dkansar@calstatela.edu". The "Permissions" menu has several options: Create, Clone, Rename, Move, Delete, Import, Export, and Permissions. A tooltip for "Create" says: "For example, it is features (wh or example, t". A tooltip for "Clone" says: "ally prepare orical feature". A tooltip for "Permissions" says: "gical features or numeric features". Below these menus is a list of notebook files:

- ADS
- Cis5560
- final
- Midterm2 Python+Classif
- Project\_DF\_AUC
- Project2
- Python\_Python\_Pipeline
- Python\_Recommendatio
- Python\_Text\_Analysis

11. Now click on import and browse “Random\_forest.ipynb” file.



12. Select Import.



13. You are directed to the page below. **Play all codes** in the cells by selecting the **play** button:

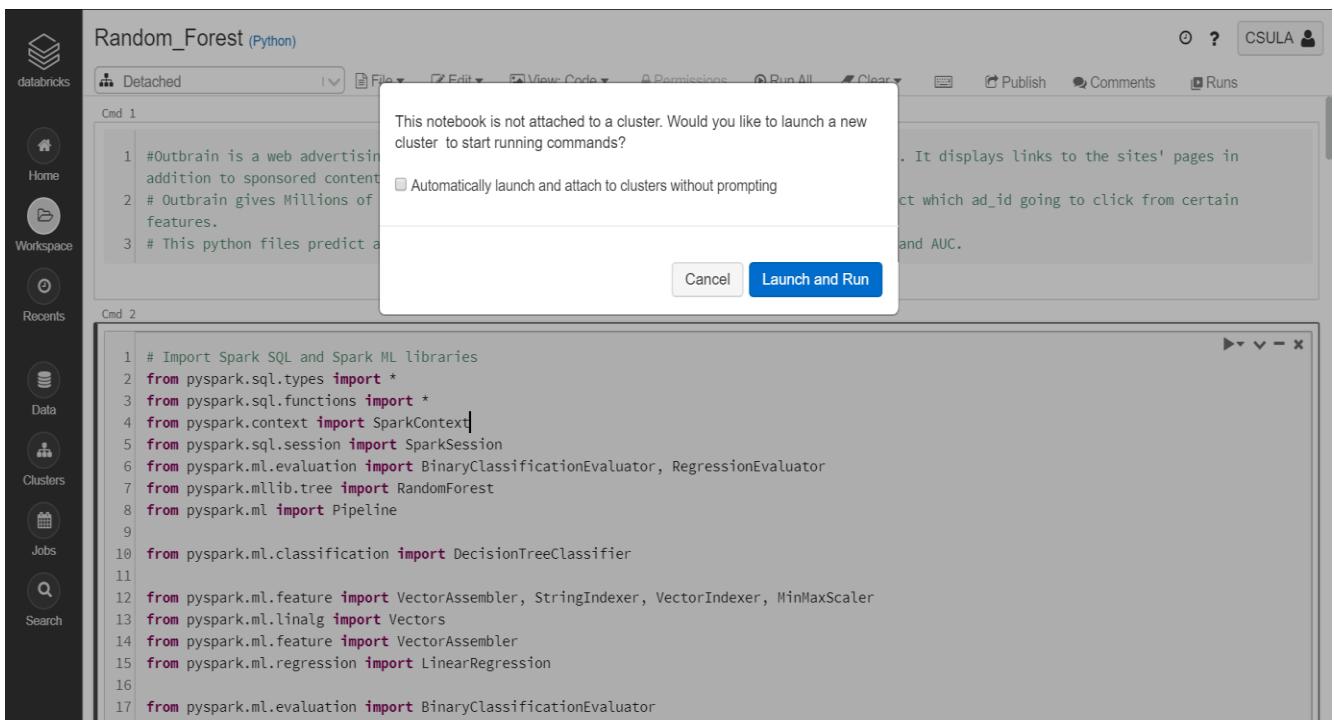
```

1 #Outbrain is a web advertising platform that displays boxes of links to pages within websites. It displays links to the sites' pages in
addition to sponsored content, generating revenue from the latter.
2 # Outbrain gives Millions of advertiser id click by user in 15 days and our click is to predict which ad_id going to click from certain
features.
3 # This python files predict ad using random forest classifier algorithm and measure Accuracy and AUC.

1 # Import Spark SQL and Spark ML libraries
2 from pyspark.sql.types import *
3 from pyspark.sql.functions import *
4 from pyspark.context import SparkContext
5 from pyspark.sql.session import SparkSession
6 from pyspark.ml.evaluation import BinaryClassificationEvaluator, RegressionEvaluator
7 from pyspark.mllib.tree import RandomForest
8 from pyspark.ml import Pipeline
9
10 from pyspark.ml.classification import DecisionTreeClassifier
11
12 from pyspark.ml.feature import VectorAssembler, StringIndexer, VectorIndexer, MinMaxScaler
13 from pyspark.ml.linalg import Vectors
14 from pyspark.ml.feature import VectorAssembler
15 from pyspark.ml.regression import LinearRegression
16
17 from pyspark.ml.evaluation import BinaryClassificationEvaluator

```

14. If spark cluster is not attached, you will see the following window popped up. Mark the “Automatically Launch ...” and select “Attach and Run”.



15. Now, it runs all cells of the notebook imported. You will see the following at the end of the code:

```

Cmd 30
1 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
2 evaluation = MulticlassClassificationEvaluator(
3     labelCol="trueLabel", predictionCol="prediction", metricName="accuracy")
4 accuracy = evaluation.evaluate(prediction)
5 print(accuracy)

▶ (2) Spark Jobs
0.9520460358056266

Command took 2.89 minutes -- by dkansar@calstatela.edu at 5/12/2020, 11:53:35 AM on Project

Cmd 31
Following Code will give total Error in project

Cmd 32
1 print("Test Error = %g" % (1.0 - accuracy))

Test Error = 0.047954

Command took 0.03 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:53:35 AM on Project

Cmd 33
Random Forest Evaluator

Calculate AUC using Random Forest Evaluator.

Cmd 34
1 rf_evaluator = MulticlassClassificationEvaluator(labelCol="trueLabel", predictionCol="prediction")
2 rf_auc = rf_evaluator.evaluate(prediction)
3 print("AUC for Random Forest is= ", rf_auc)

▶ (3) Spark Jobs
AUC for Random Forest is= 0.9499206293049177

Command took 4.14 minutes -- by dkansar@calstatela.edu at 5/12/2020, 11:53:35 AM on Project

Cmd 35
1 # Only for Classification Logistic Regression not for Linear Regression
2
3 tp = float(predicted.filter("prediction == 1.0 AND truelabel == 1").count())
4 fp = float(predicted.filter("prediction == 1.0 AND truelabel == 0").count())
5 tn = float(predicted.filter("prediction == 0.0 AND truelabel == 0").count())
6 fn = float(predicted.filter("prediction == 0.0 AND truelabel == 1").count())
7 metrics = spark.createDataFrame([
8     ("TP", tp),
9     ("FP", fp),
10    ("TN", tn),
11    ("FN", fn),
12    ("Precision", tp / (tp + fp)),
13    ("Recall", tp / (tp + fn))],["metric", "value"])
14 metrics.show()
15
16
▶ (7) Spark Jobs
▶ └── metrics: pyspark.sql.dataframe.DataFrame = [metric: string, value: double]
+-----+-----+
| metric |      value |
+-----+-----+
|   TP |      1562.0 |
|   FP |      138.0 |
|   TN |      11839.0 |
|   FN |      537.0 |
| Precision | 0.9188235294117647 |
| Recall | 0.7441638875655073 |
+-----+-----+

Command took 4.42 minutes -- by dkansar@calstatela.edu at 5/12/2020, 12:49:16 PM on Project

```

## Logistic Regression:

16. Now Follow **Step 10-14** and import Logistic\_Regression.ipynb file.

17. Now run the notebook imported. You will see the following at the end of the code:

```
Cmd 26
1 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
2 evaluation = MulticlassClassificationEvaluator(
3     labelCol="trueLabel", predictionCol="prediction", metricName="accuracy")
4 accuracy = evaluation.evaluate(prediction)
5 print("Accuracy of Logistic Regression is: ",accuracy)

▶ (2) Spark Jobs
Accuracy of Logistic Regression is:  0.8530731429385371
Command took 17.39 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:31:11 PM on Project

Cmd 27
Following Code will give total Error in project

Cmd 28
1 print("Test Error = %g" % (1.0 - accuracy))

Test Error = 0.146927
Command took 0.03 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:31:11 PM on Project

Cmd 29
}
Logistic Regression Evaluator

Calculate AUC using Logistic Regression Evaluator.

Cmd 30
1 rf_evaluator = MulticlassClassificationEvaluator(labelCol="trueLabel", predictionCol="prediction")
2 rf_auc = rf_evaluator.evaluate(prediction)
3 print("AUC for Logistic Regression is= ", rf_auc)

Cmd 31
1 tp = float(predicted.filter("prediction == 1.0 AND truelabel == 1").count())
2 fp = float(predicted.filter("prediction == 1.0 AND truelabel == 0").count())
3 tn = float(predicted.filter("prediction == 0.0 AND truelabel == 0").count())
4 fn = float(predicted.filter("prediction == 0.0 AND truelabel == 1").count())
5 metrics = spark.createDataFrame([
6     ("TP", tp),
7     ("FP", fp),
8     ("TN", tn),
9     ("FN", fn),
10    ("Precision", tp),
11    ("Recall", tp / (tp + fn))], ["metric", "value"])
12 metrics.show()
13

▶ (7) Spark Jobs
▶   metrics: pyspark.sql.dataframe.DataFrame = [metric: string, value: double]
+-----+
| metric | value |
+-----+-----+
|      TP |    0.0 |
|      FP |    0.0 |
|      TN | 11978.0 |
|      FN |  2063.0 |
| Precision |    0.0 |
| Recall |    0.0 |
+-----+-----+

Command took 41.49 seconds -- by dkansar@calstatela.edu at 5/13/2020, 12:24:09 AM on Project
```

## Gradient\_Boosted Algorithm:

18. Now Follow **Step 10-14** and import Gradient\_boost.ipynb file.

19. Now, it runs all cells of the notebook imported. You will see the following at the end of the code:

```

1 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
2 evaluation = MulticlassClassificationEvaluator(
3     labelCol="trueLabel", predictionCol="prediction", metricName="accuracy")
4 accuracy = evaluation.evaluate(prediction)
5 print(accuracy)

▶ (2) Spark Jobs
0.8984530229917683

Command took 25.96 seconds -- by dkansar@calstatela.edu at 5/12/2020, 12:55:10 PM on Project

```

Cmd 27

Following Code will give total Error in project

Cmd 28

```
1 print("Test Error = %g" % (1.0 - accuracy))
```

Test Error = 0.101547

Command took 0.01 seconds -- by dkansar@calstatela.edu at 5/12/2020, 12:55:11 PM on Project

Cmd 29

## Gradient Boosting Evaluator

Calculate AUC using Gradient Boosting Evaluator.

Cmd 30

```
1 rf_evaluator = MulticlassClassificationEvaluator(labelCol="trueLabel", predictionCol="prediction")
2 rf_auc = rf_evaluator.evaluate(prediction)
3 print("AUC for Gradient Boost is= ", rf_auc)
```

▶ (3) Spark Jobs

AUC for Gradient Boost is= 0.8750931883418767

Cmd 31

```
1 # Only for Classification Logistic Regression not for Linear Regression
2
3 tp = float(predicted.filter("prediction == 1.0 AND truelabel == 1").count())
4 fp = float(predicted.filter("prediction == 1.0 AND truelabel == 0").count())
5 tn = float(predicted.filter("prediction == 0.0 AND truelabel == 0").count())
6 fn = float(predicted.filter("prediction == 0.0 AND truelabel == 1").count())
7 metrics = spark.createDataFrame([
8     ("TP", tp),
9     ("FP", fp),
10    ("TN", tn),
11    ("FN", fn),
12    ("Precision", tp / (tp + fp)),
13    ("Recall", tp / (tp + fn))],["metric", "value"])
14 metrics.show()
```

▶ (7) Spark Jobs

▶ [metrics: pyspark.sql.dataframe.DataFrame = [metric: string, value: double]]

metric	value
TP	653.0
FP	10.0
TN	12008.0
FN	1421.0
Precision	0.9849170437405732
Recall	0.31485053037608485

Command took 53.06 seconds -- by dkansar@calstatela.edu at 5/12/2020, 12:55:11 PM on Project

## Decision Tree Classifier:

20. Now Follow Step 10-14 and import Decision Tree Classifier.ipynb file.

21. Now, it runs all cells of the notebook imported. You will see the following at the end of the code:

```
ed to cluster. Project, 15.25 GB | 2 Cores | DBR | Spark 2.4.5 | Scala 2.11
import MulticlassClassificationEvaluator
| MulticlassClassificationEvaluator(
|   predictionCol="prediction", metricName="accuracy")
4 accuracy = evaluation.evaluate(prediction)
5 print("Accuracy for Decision Tree classifier is=",accuracy)

▶ (2) Spark Jobs
Accuracy for Decision Tree classifier is= 0.943342776203966
Command took 34.98 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:11:53 PM on Project
Cmd 31
Following Code will give total Error in project

Cmd 32
1 print("Test Error = %g" % (1.0 - accuracy))
Test Error = 0.0566572
Command took 0.03 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:11:53 PM on Project
Cmd 33
1 %md ### Decision Tree Evaluator
2 Calculate AUC using Decision Tree Evaluator.

Cmd 34
1 rf_evaluator = MulticlassClassificationEvaluator(labelCol="trueLabel", predictionCol="prediction")
2 rf_auc = rf_evaluator.evaluate(prediction)
3 print("AUC for Decision Tree is= ", rf_auc)

▶ (3) Spark Jobs
AUC for Decision Tree is=  0.9410911194259339
Command took 54.48 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:11:53 PM on Project
|
▶ (3) Spark Jobs
AUC for Decision Tree is=  0.9410911194259339
Command took 54.48 seconds -- by dkansar@calstatela.edu at 5/12/2020, 11:11:53 PM on Project

Cmd 35
1 # Only for Classification Logistic Regression not for Linear Regression
2
3 tp = float(predicted.filter("prediction == 1.0 AND truelabel == 1").count())
4 fp = float(predicted.filter("prediction == 1.0 AND truelabel == 0").count())
5 tn = float(predicted.filter("prediction == 0.0 AND truelabel == 0").count())
6 fn = float(predicted.filter("prediction == 0.0 AND truelabel == 1").count())
7 metrics = spark.createDataFrame([
8   ("TP", tp),
9   ("FP", fp),
10  ("TN", tn),
11  ("FN", fn),
12  ("Precision", tp / (tp + fp)),
13  ("Recall", tp / (tp + fn))],["metric", "value"])
14 metrics.show()
15

▶ (7) Spark Jobs
▶ └── metrics: pyspark.sql.dataframe.DataFrame = [metric: string, value: double]
+-----+-----+
| metric|      value|
+-----+-----+
|    TP|      1571.0|
|    FP|      213.0|
|    TN|     11749.0|
|    FN|      587.0|
|Precision|0.8806053811659192|
|  Recall|0.7279888785912882|
+-----+-----+
```

## Step 5: Ad prediction using spark in Oracle BDE

---

1. We need to import our dataset (event, train and Promoted\_content) to oracle BDE cluster.  
In order to do that we need to navigate to directory where we stored our csv files.
2. We can login to Oracle BDE by any terminal.

Ssh [dkansar@129.150.76.160](ssh dkansar@129.150.76.160)

```
C:\Users\pdmuser>ssh dkansar@129.150.76.160
-- WARNING -- This system is for the use of authorized users only. Individuals
using this computer system without authority or in excess of their authority
are subject to having all their activities on this system monitored and
recorded by system personnel. Anyone using this system expressly consents to
such monitoring and is advised that if such monitoring reveals possible
evidence of criminal activity system personnel may provide the evidence of such
monitoring to law enforcement officials.

dkansar@129.150.76.160's password:
Last login: Fri May  8 04:45:41 2020 from 47.155.196.173
```

**Note:** Username (dkansar) and IP address may be different in your case.

3. After navigating to directory open terminal your system and write following command.

scp \*.csv [dkansar@129.150.76.160:~/](scp *.csv dkansar@129.150.76.160:~/)

**Note:** 1: Instead of \* you need to put different file names event, train and Promoted content. File name may also be different in your case.

2: Ip Address and Username (dkansar) may be different in your case.

4. Now we need to file our csv files to hdfs directory we can do it by following.

- a. Hdfs dfs -put event.csv (File name may be different.)
- b. Hdfs dfs -put train.csv (File name may be different.)
- c. Hdfs dfs -put Promoted\_content.csv (File name may be different.)

5. Now we can check our file uploaded in hdfs.

Hdfs dfs -ls

```
-bash-4.1$ hdfs dfs -ls
Found 9 items
-rw-r--r--  2 dkansar hdfs  285719991 2020-05-04 07:50 event.csv
-rw-r--r--  2 dkansar hdfs 1208549589 2020-05-02 18:23 eventd.csv
-rw-r--r--  2 dkansar hdfs  35715640 2020-05-06 20:22 events.csv
-rw-r--r--  2 dkansar hdfs  72088113 2020-05-05 21:19 flights.csv
drwxr-xr-x - dkansar hdfs          0 2020-04-27 05:19 lda
-rw-r--r--  2 dkansar hdfs 13886609 2020-05-02 18:23 promoted_content.csv
-rw-r--r--  2 dkansar hdfs 234703736 2020-05-04 07:25 train.csv
-rw-r--r--  2 dkansar hdfs 1159782400 2020-04-30 08:11 traind.csv
-rw-r--r--  2 dkansar hdfs 36395707 2020-05-06 20:18 trains.csv
-bash-4.1$
```

6. Now we have our csv files in hdfs.

7. In step 3, we have predicted data in data bricks now you need to login to databricks and open Logistic\_Regression file. Now browse File>Export>source file as shown in following image.

```

Random_Forest (Python)
Cmd 1
Outbrain Ad Predicti
1 # Outbrain is a web advert addition to sponsored content
2 # Outbrain gives Millions features.
3 # This python files predi
Cmd 2
1 # Import Spark SQL and Spark ML libraries
2 from pyspark.sql.types import *
3 from pyspark.sql.functions import *
4 from pyspark.context import SparkContext
5 from pyspark.sql.session import SparkSession

```

Using the same step, you need to download Gradient\_Boosted.py and Random\_Forest.py file.

- Now open Python notebook file in any editor and make change as follows.

**IS\_SPARK\_SUBMIT\_CLI = True (change from false to True)**

- After navigating to directory open terminal your system and write following command.

scp \*.py dkansar@129.150.76.160:~/

**Note:** 1: Instead of \* you need to put different file names Logistic Regression, Gradient boost and Random\_Forest, Decision Tree Classifier. File name may also be different in your case.  
2: Ip Address and Username (dkansar) may be different in your case.

- Now we need run Logistic Regression in Oracle BDE as follows.

**Spark-submit Logistic\_Regression.py**

Now we can measure Accuracy and AUC as follows.

**Note:** In Oracle BDCE you can only able to see Result of last command Confusion Matrix.

- If you want to see result of AUC Comment out Confusion Matrix Code (you can Comment line by just putting # in front of line). Then upload and save file using Step 5 and 6 you can able to see AUC of algorithm.
- If you want to see result of Accuracy Comment out Confusion Matrix Code (you and AUC Commands Comment line by just putting # in front of line). Then upload and run file using step 5 and 6 you able to see Accuracy of Algorithm.

**Accuracy:**

```
20/05/04 22:55:21 INFO Executor: Finished task 48.0 in stage 143.0 (TID 6257). 1827 bytes resu
20/05/04 22:55:21 INFO TaskSetManager: Finished task 48.0 in stage 143.0 (TID 6257) in 16 ms c
20/05/04 22:55:21 INFO TaskSchedulerImpl: Removed TaskSet 143.0, whose tasks have all complete
20/05/04 22:55:21 INFO DAGScheduler: ResultStage 143 (countByValue at MulticlassMetrics.scala:42)
20/05/04 22:55:21 INFO DAGScheduler: Job 31 finished: countByValue at MulticlassMetrics.scala:42
0.830304235215
Test Error = 0.169696
20/05/04 22:55:21 INFO SparkContext: Invoking stop() from shutdown hook
20/05/04 22:55:21 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/04 22:55:21 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/04 22:55:21 INFO MemoryStore: MemoryStore cleared
```

**AUC:**

```
20/05/07 05:20:27 INFO TaskSchedulerImpl: Removed TaskSet 161.0, whose tasks have all completed, from pool
20/05/07 05:20:27 INFO DAGScheduler: ResultStage 161 (collectAsMap at MulticlassMetrics.scala:53) finished in 0.202
20/05/07 05:20:27 INFO DAGScheduler: Job 34 finished: collectAsMap at MulticlassMetrics.scala:53, took 4.343779 s
('AUC for Logistic Regression is= ', 0.7539270092721175)
20/05/07 05:20:27 INFO SparkContext: Invoking stop() from shutdown hook
```

**Confusion Matrix:**

```
+-----+
| metric | value |
+-----+
| TP    | 0.0 |
| FP    | 0.0 |
| TN    | 13086.0 |
| FN    | 2508.0 |
| Precision | 0.0 |
| Recall   | 0.0 |
+-----+
```

```
20/05/13 05:15:48 INFO SparkContext: Invoking stop() from shutdown hook
20/05/13 05:15:48 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/13 05:15:48 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

7. Now we need run Gradient\_Boosted in Oracle BDE as follows.

**Spark-submit Gradient\_Boosted.py**

Now we can measure Accuracy and AUC as follows.

**Accuracy:**

```
20/05/08 06:46:39 INFO TaskSchedulerImpl: Removed TaskSet 17078.0, whose tasks have all completed, from pool
20/05/08 06:46:39 INFO DAGScheduler: ResultStage 17078 (countByValue at MulticlassMetrics.scala:42) finished in 0.108 s
20/05/08 06:46:39 INFO DAGScheduler: Job 2980 finished: countByValue at MulticlassMetrics.scala:42, took 0.528726 s
('Accuracy of for Gradient Boosted is= ', 0.9247810743693635)
Test Error = 0.0752189
20/05/08 06:46:39 INFO SparkContext: Invoking stop() from shutdown hook
```

## AUC:

```
20/05/08 18:28:01 INFO TaskSetManager: FINISHED task 199.0 in stage 17120.0 (TID 920145) in 12 ms on localhost (executor d)
20/05/08 18:28:01 INFO TaskSetManager: Finished task 48.0 in stage 17120.0 (TID 920145) in 12 ms on localhost (executor d)
20/05/08 18:28:01 INFO TaskSchedulerImpl: Removed TaskSet 17120.0, whose tasks have all completed, from pool
20/05/08 18:28:01 INFO DAGScheduler: ResultStage 17120 (collectAsMap at MulticlassMetrics.scala:53) finished in 0.094 s
20/05/08 18:28:01 INFO DAGScheduler: Job 2987 finished: collectAsMap at MulticlassMetrics.scala:53, took 0.356720 s
('AUC for Gradient Boosted is= ', 0.9208955514357123)
20/05/08 18:28:01 INFO SparkContext: Invoking stop() from shutdown hook
20/05/08 18:28:01 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/08 18:28:01 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/08 18:28:10 INFO MemoryStore: MemoryStore cleared
20/05/08 18:28:10 INFO BlockManager: BlockManager stopped
```

## Confusion Matrix:

```
20/05/12 22:17:13 INFO DAGScheduler: ResultStage 16938 (showString at NativeMethodAccessorImpl.java:17) finished: showString at NativeMethodAccessorImpl.java:17
20/05/12 22:17:13 INFO DAGScheduler: Job 2960 finished: showString at NativeMethodAccessorImpl.java:17
20/05/12 22:17:13 INFO CodeGenerator: Code generated in 10.815248 ms
+-----+-----+
| metric | value |
+-----+-----+
| TP | 1474.0 |
| FP | 354.0 |
| TN | 12713.0 |
| FN | 946.0 |
| Precision | 0.8063457330415755 |
| Recall | 0.6090909090909091 |
+-----+-----+
20/05/12 22:17:13 INFO SparkContext: Invoking stop() from shutdown hook
20/05/12 22:17:13 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/12 22:17:13 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/12 22:17:22 INFO MemoryStore: MemoryStore cleared
```

- Now we need run Random\_Forest in Oracle BDE as follows.

### Spark-submit Random\_Forest.py

Now we can measure Accuracy and AUC as follows.

#### Accuracy:

```
20/05/04 04:37:15 INFO TaskSchedulerImpl: Removed TaskSet 82.0, whose tasks have all completed, from pool
20/05/04 04:37:15 INFO DAGScheduler: ResultStage 82 (countByValue at MulticlassMetrics.scala:42) finished in 0.932862190813
20/05/04 04:37:15 INFO DAGScheduler: Job 34 finished: countByValue at MulticlassMetrics.scala:42, took 0.444
Test Error = 0.0671378
20/05/04 04:37:15 INFO SparkContext: Invoking stop() from shutdown hook
20/05/04 04:37:15 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
```

#### AUC:

```
20/05/07 05:31:08 INFO TaskSchedulerImpl: Removed TaskSet 276.0, whose tasks have all completed, from pool
20/05/07 05:31:08 INFO DAGScheduler: ResultStage 276 (collectAsMap at MulticlassMetrics.scala:53) finished in 0.116 s
20/05/07 05:31:08 INFO DAGScheduler: Job 54 finished: collectAsMap at MulticlassMetrics.scala:53, took 1.664514 s
('AUC for Random forest is= ', 0.9285742872348733)
20/05/07 05:31:08 INFO SparkContext: Invoking stop() from shutdown hook
20/05/07 05:31:08 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4042
20/05/07 05:31:08 ERROR LivelistenerBus: SparkListenerBus has already stopped! Dropping event SparkListenerExecutorMet
```

### Confusion Matrix:

```
20/05/12 20:14:51 INFO DAGScheduler: Job 62 finished: showString at NativeMethodAccessorImpl.java:0, took 0.062228 s
20/05/12 20:14:51 INFO CodeGenerator: Code generated in 5.842159 ms
+-----+-----+
| metric | value |
+-----+-----+
| TP | 1736.0 |
| FP | 149.0 |
| TN | 12859.0 |
| FN | 717.0 |
| Precision | 0.9209549071618037 |
| Recall | 0.7077048512026091 |
+-----+-----+
20/05/12 20:14:51 INFO SparkContext: Invoking stop() from shutdown hook
20/05/12 20:14:51 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/12 20:14:51 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/12 20:14:51 INFO MemoryStore: MemoryStore cleared
20/05/12 20:14:51 INFO BlockManager: BlockManager stopped
```

9. Now we need run Random\_Forest in Oracle BDE as follows.

### Spark-submit Decision Tree Classifier.py

Now we can measure Accuracy and AUC as follows.

#### Accuracy:

```
20/05/13 06:04:07 INFO TaskSchedulerImpl: Removed TaskSet 253.0, whose tasks have all completed, from
20/05/13 06:04:07 INFO DAGScheduler: ResultStage 253 (countByValue at MulticlassMetrics.scala:42) finished in 0.003 s
20/05/13 06:04:07 INFO DAGScheduler: Job 50 finished: countByValue at MulticlassMetrics.scala:42, took 0.003 s
('Accuracy for Decision Tree classifier is=', 0.9387385053749514)
Test Error = 0.0612615
20/05/13 06:04:07 INFO SparkContext: Invoking stop() from shutdown hook
20/05/13 06:04:07 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/13 06:04:07 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/13 06:04:07 INFO MemoryStore: MemoryStore cleared
20/05/13 06:04:07 INFO BlockManager: BlockManager stopped
```

#### AUC:

```
20/05/13 06:01:31 INFO TaskSetManager: Finished task 0.0 in stage 271.0 (TID 16623) in 29 ms on localhost
20/05/13 06:01:31 INFO TaskSchedulerImpl: Removed TaskSet 271.0, whose tasks have all completed, from pool 0
20/05/13 06:01:31 INFO DAGScheduler: ResultStage 271 (collectAsMap at MulticlassMetrics.scala:53) finished in 0.003 s
20/05/13 06:01:31 INFO DAGScheduler: Job 53 finished: collectAsMap at MulticlassMetrics.scala:53, took 0.003 s
('AUC for Decision Tree is= ', 0.9358037291398622)
20/05/13 06:01:31 INFO SparkContext: Invoking stop() from shutdown hook
20/05/13 06:01:31 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/13 06:01:31 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
20/05/13 06:01:32 INFO MemoryStore: MemoryStore cleared
20/05/13 06:01:32 INFO BlockManager: BlockManager stopped
20/05/13 06:01:32 INFO BlockManagerMaster: BlockManagerMaster stopped
20/05/13 06:01:32 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
```

### Confusion Matrix:

```
20/05/13 05:58:05 INFO DAGScheduler: Job 61 finished: showString at NativeMethodAccessorImpl.java:0, took 0.062228 s
20/05/13 05:58:05 INFO CodeGenerator: Code generated in 5.842159 ms
+-----+-----+
| metric | value |
+-----+-----+
| TP | 1815.0 |
| FP | 252.0 |
| TN | 12735.0 |
| FN | 645.0 |
| Precision | 0.8780841799709724 |
| Recall | 0.7378048780487805 |
+-----+-----+
20/05/13 05:58:05 INFO SparkContext: Invoking stop() from shutdown hook
20/05/13 05:58:05 INFO SparkUI: Stopped Spark web UI at http://10.129.201.54:4041
20/05/13 05:58:05 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

