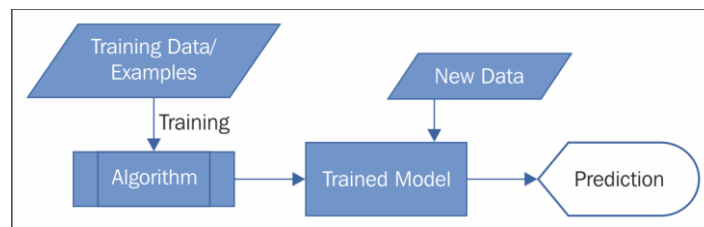
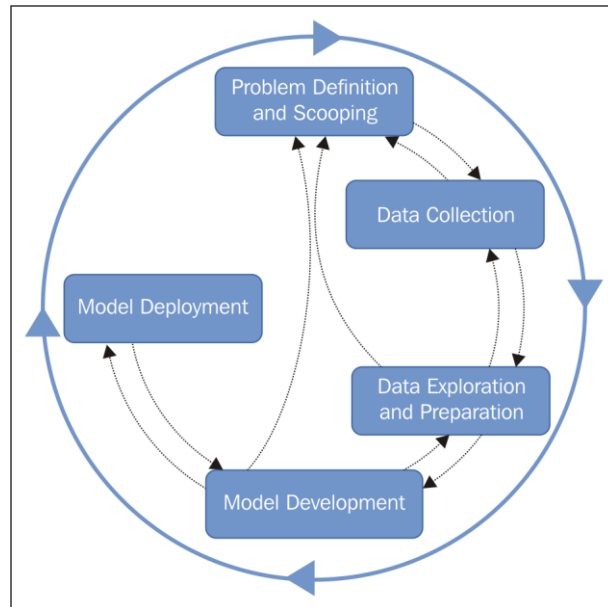
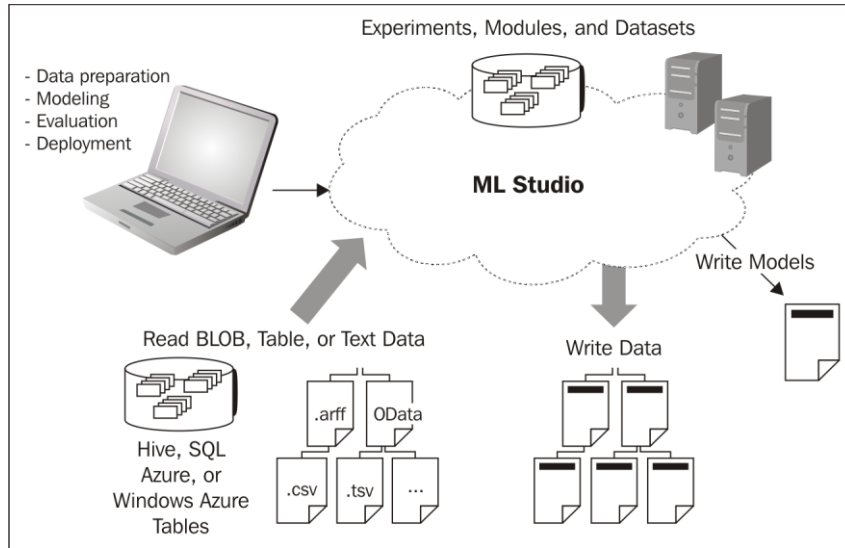


# Chapter 1



# Chapter 2



Microsoft Azure

- CDN 0
- AUTOMATION 0
- SCHEDULER 0
- API MANAGEMENT 0
- MACHINE LEARNING** 0
- NETWORKS 0

machine learning **PREVIEW**

You have no ML workspaces.

CREATE AN ML WORKSPACE →

Microsoft Azure

- CDN 0
- AUTOMATION 0
- SCHEDULER 0
- API MANAGEMENT 0
- MACHINE LEARNING** 1

machine learning **PREVIEW**

NAME	STATUS	OWNER	SUBSCRIPTION	LOCATION
dvanatest	→ ✓	@live.com		

mlWorkspace

# mlworkspace

[DASHBOARD](#) [CONFIGURE](#) [WEB SERVICES](#)

NAME	STORAGE	STATUS	OWNER	SUBSCRIPTION
mlWorkspace	mlstorageXXXXXX	✓ Online	XXXXXXXX@XXXXXX.com	XXXXXX

+ NEW
MANAGE KEYS
**OPEN IN STUDIO**
DELETE

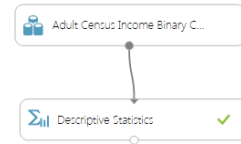
- EXPERIMENTS
- WEB SERVICES
- DATASETS
- TRAINED MODELS
- SETTINGS

## experiments

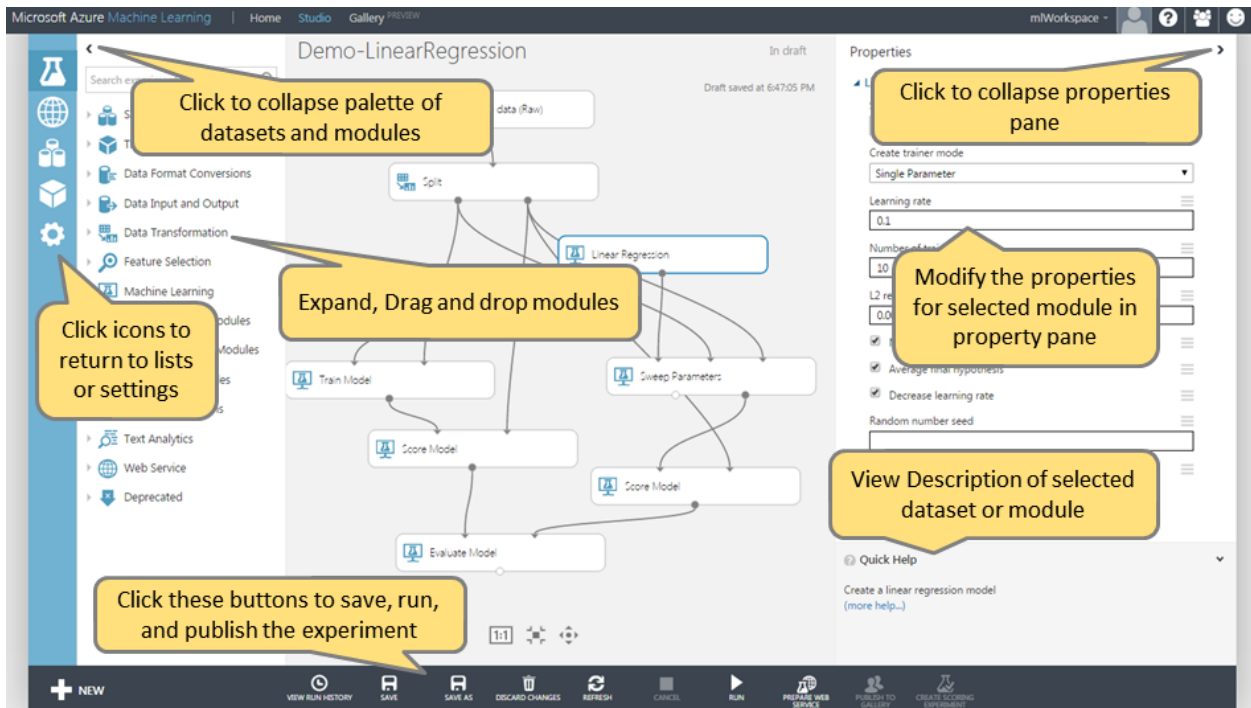
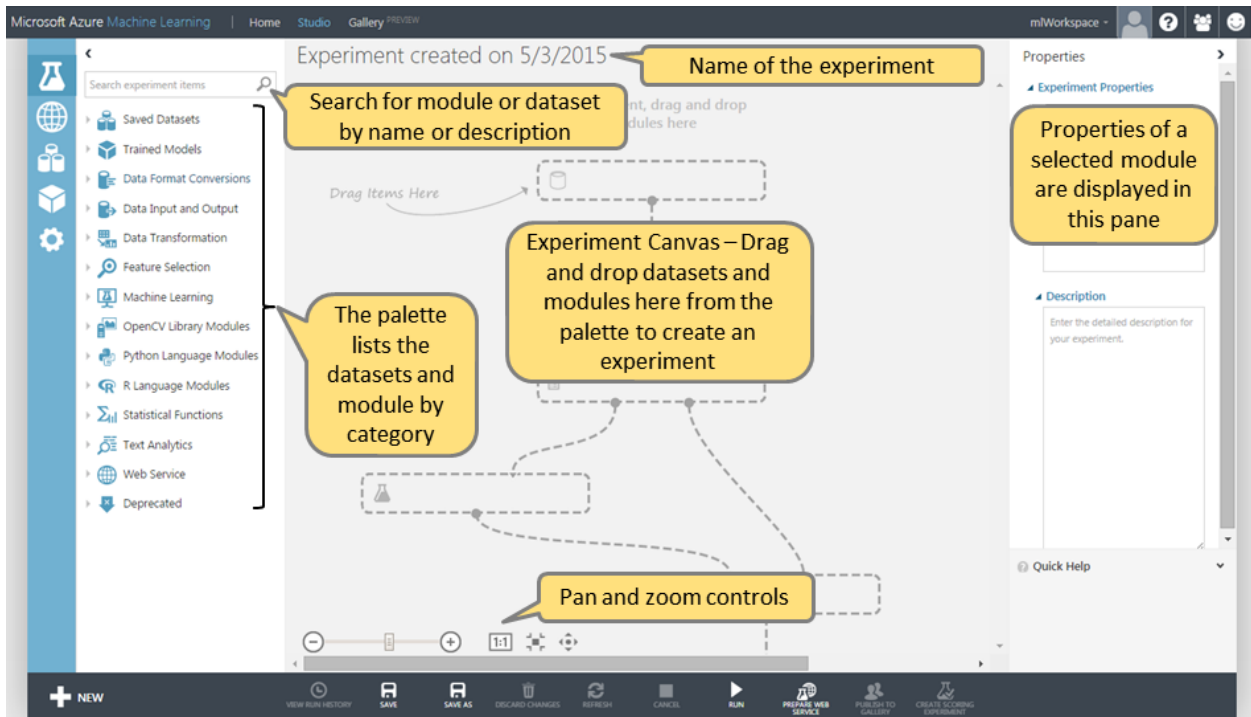
MY EXPERIMENTS SAMPLES

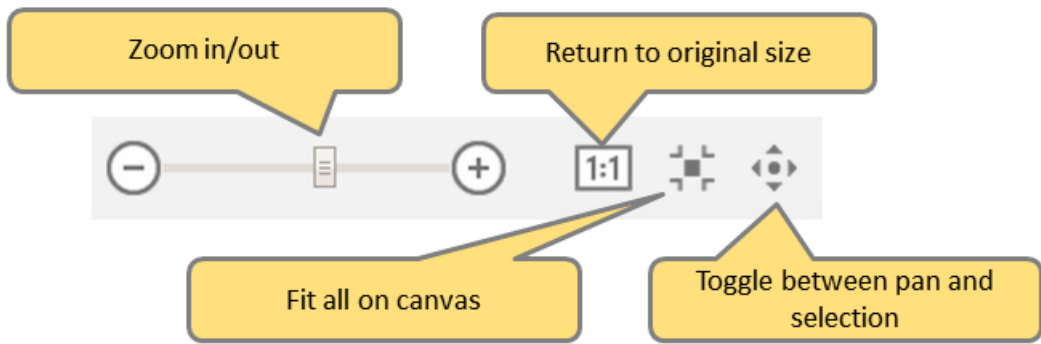
	NAME	AUTHOR	STATUS	LAST EDITED
<input type="checkbox"/>	Demo-LinearRegression	sumitmund	Draft	11/26/2014 2:11:53 PM
<input type="checkbox"/>	Demo-RegressionDF_Sweep	sumitmund	Draft	11/14/2014 11:46:43 AM
<input type="checkbox"/>	Regression-ForestFire	sumitmund	Draft	11/7/2014 6:36:13 PM
<input type="checkbox"/>	Ch6-DF	sumitmund	Draft	10/27/2014 4:53:12 AM
<input type="checkbox"/>	Ch-3	sumitmund	Draft	10/26/2014 7:52:56 PM
<input type="checkbox"/>	CH - Data Input Output	sumitmund	Draft	10/22/2014 6:57:18 PM
<input type="checkbox"/>	mnist	sumitmund	Finished	10/6/2014 9:27:08 PM
<input checked="" type="checkbox"/>	Ch-2	sumitmund	Finished	10/5/2014 6:44:36 AM
<input type="checkbox"/>	Demo-EvaluateModel	sumitmund	Draft	9/21/2014 2:44:35 PM
<input type="checkbox"/>	Demo-R	sumitmund	Finished	8/29/2014 3:44:06 PM

1 2 3 ← →



+ NEW
DELETE
COPY TO WORKSPACE





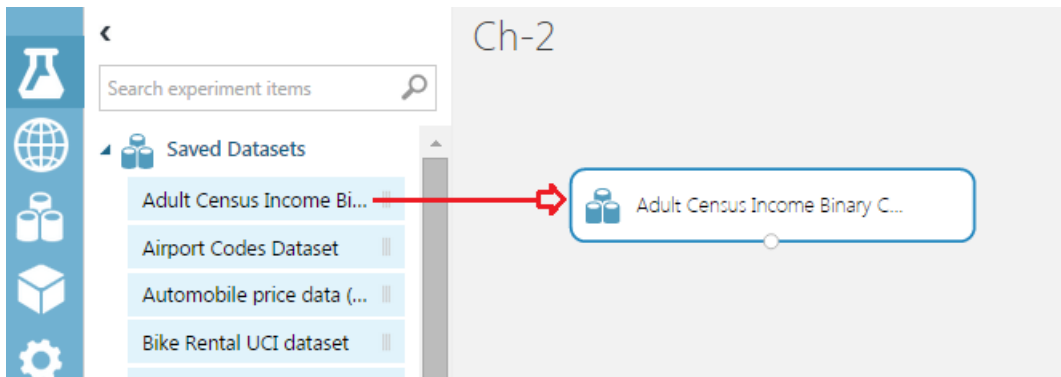
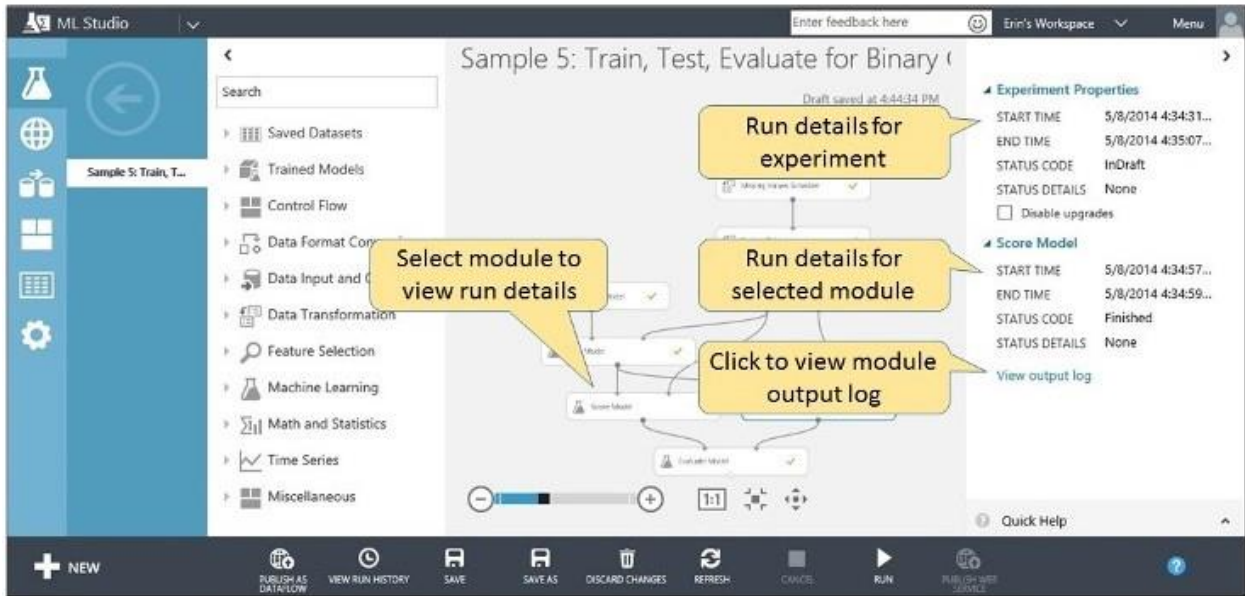
Pima Indian (NN-ParamSweep) Running (0:00:23)

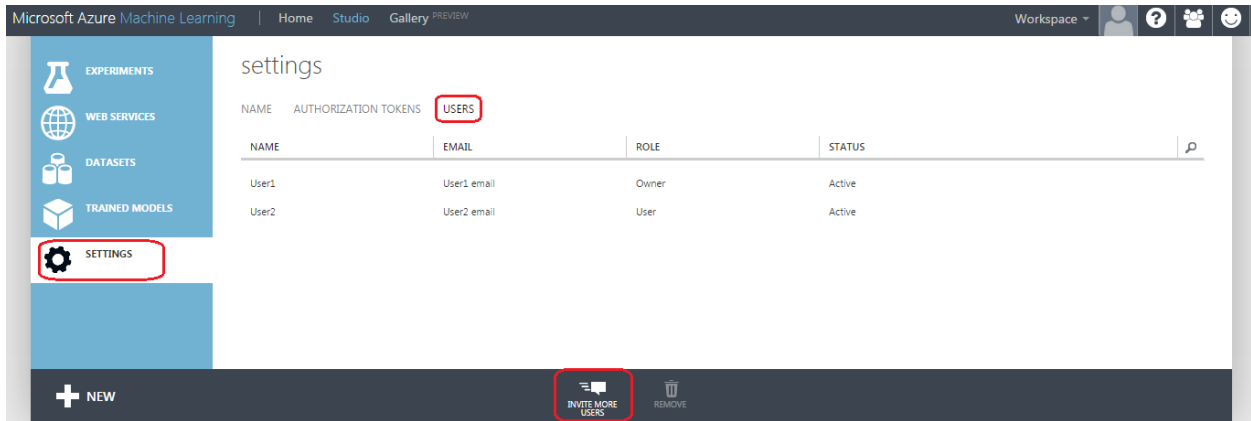
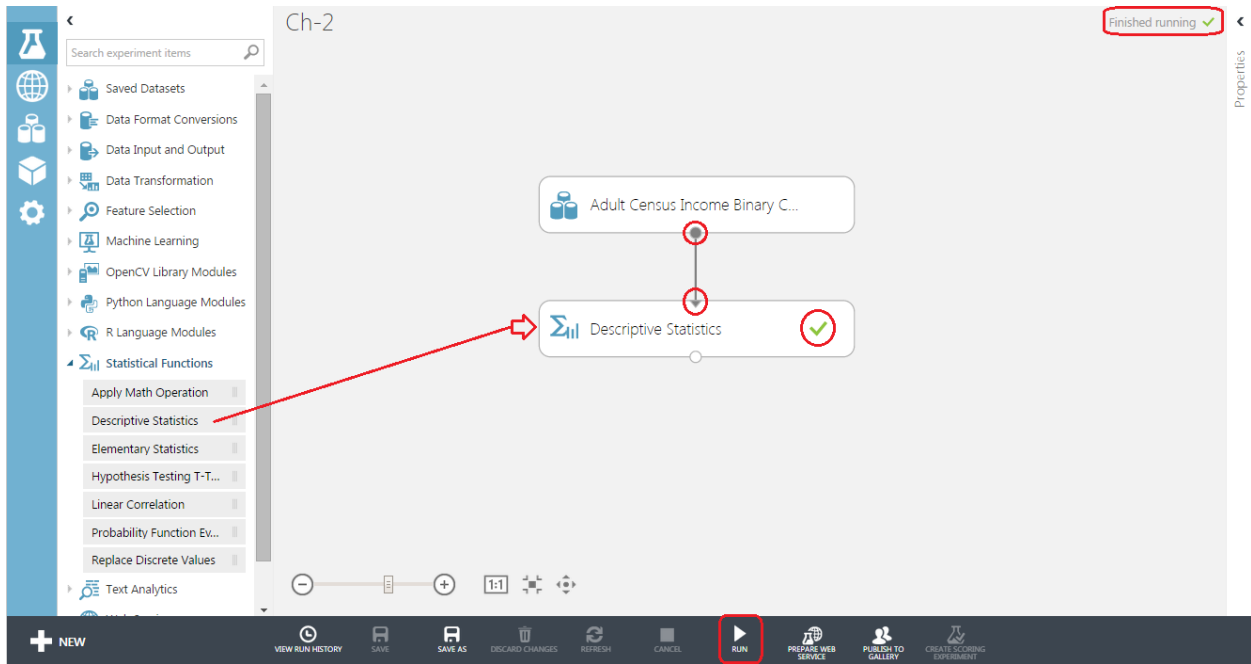
Module icon indicates status

- Finished
- Running
- Scheduled

Click "Refresh" to update status of the experiment

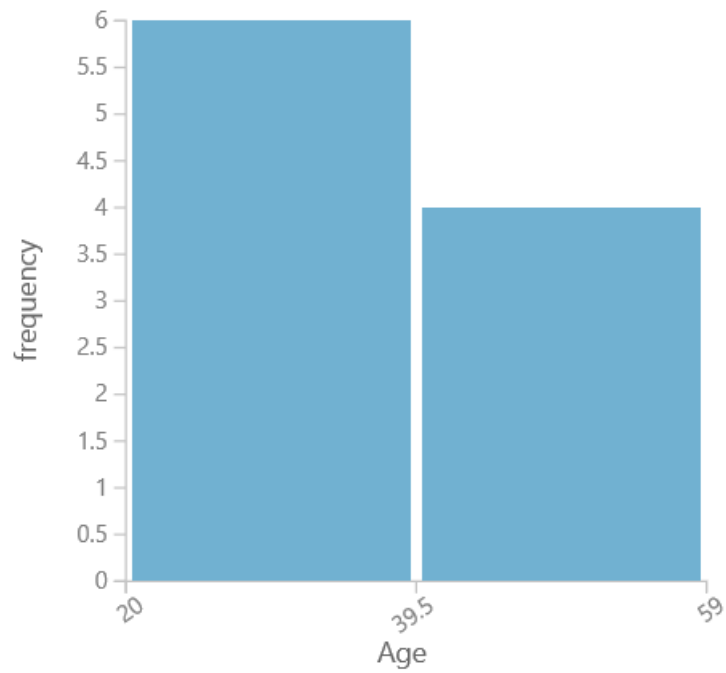
NEW | VIEW RUN HISTORY | SAVE | SAVE AS | DISCARD CHANGES | REFRESH | CANCEL | RUN | PREPARE WEB SERVICE | PUBLISH TO GALLERY | CREATE SCORING EXPERIMENT



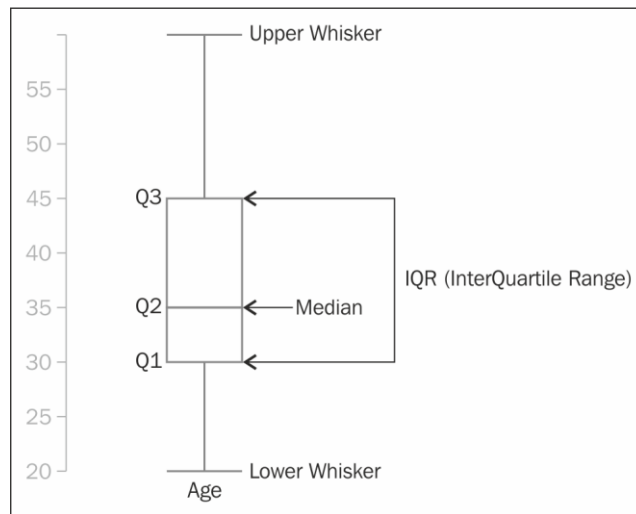
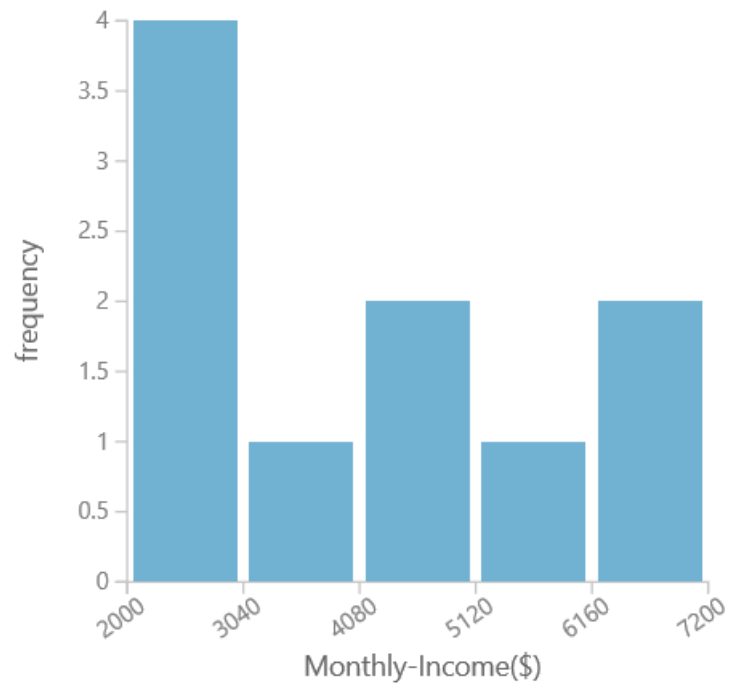


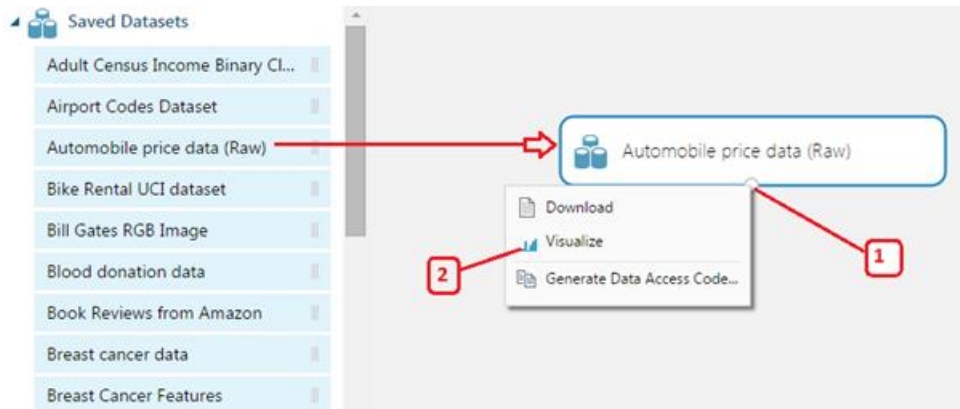
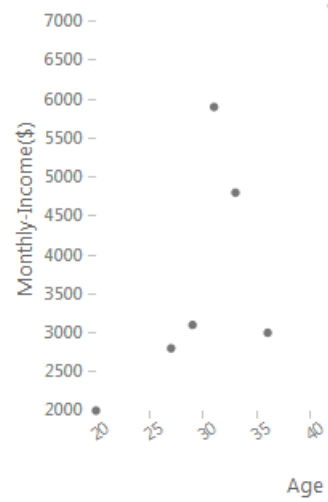
## Chapter 3

	Name	Age	Gender	Monthly-Income(\$)
Mean		37.7		4270
Median		34.5		3950
Min		20		2000
Max		59		7200
Standard Deviation		12.4637		1850.5555
Unique Values	10	10	2	10
Missing Values	0	0	0	0
Feature Type	String	Numeric	String	Numeric









rows 205  
columns 26

symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	width	height
3		alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.0
3		alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8	64.1	48.0
1		alfa-romero	gas	std	two	hatchback	rwd	front	94.5	171.2	65.5	52.4
2	164	audi	gas	std	four	sedan	fwd	front	99.8	176.6	66.2	54.0
2	164	audi	gas	std	four	sedan	4wd	front	99.4	176.6	66.4	54.0
2		audi	gas	std	two	sedan	fwd	front	99.8	177.3	66.3	53.2
1	158	audi	gas	std	four	sedan	fwd	front	105.8	192.7	71.4	55.0
1		audi	gas	std	four	wagon	fwd	front	105.8	192.7	71.4	55.0
1	158	audi	gas	turbo	four	sedan	fwd	front	105.8	192.7	71.4	55.0
0		audi	gas	turbo	two	hatchback	4wd	front	99.5	178.2	67.9	52.0
2	192	bmw	gas	std	two	sedan	rwd	front	101.2	176.8	64.8	54.0
0	192	bmw	gas	std	four	sedan	rwd	front	101.2	176.8	64.8	54.0
0	188	bmw	gas	std	two	sedan	rwd	front	101.2	176.8	64.8	54.0
0	188	bmw	gas	std	four	sedan	rwd	front	101.2	176.8	64.8	54.0
1		bmw	gas	std	four	sedan	rwd	front	103.5	189	66.9	55.0
0		bmw	gas	std	four	sedan	rwd	front	103.5	189	66.9	55.0
0		bmw	gas	std	two	sedan	rwd	front	103.5	193.8	67.9	53.0
0		bmw	gas	std	four	sedan	rwd	front	110	197	70.9	56.0
2	121	chevrolet	gas	std	two	hatchback	fwd	front	88.4	141.1	60.3	53.0
1	98	chevrolet	gas	std	two	hatchback	fwd	front	94.5	150.9	63.6	52.0

#### Statistics

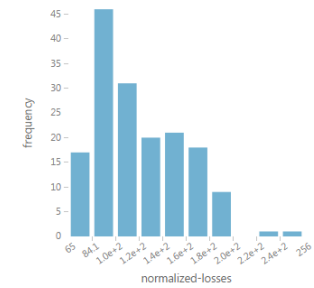
Mean 122  
Median 115  
Min 65  
Max 256  
Standard Deviation 35.4422  
Unique Values 51  
Missing Values 41  
Feature Type Numeric Feature

#### Visualizations

normalized-losses

Histogram

compare to None

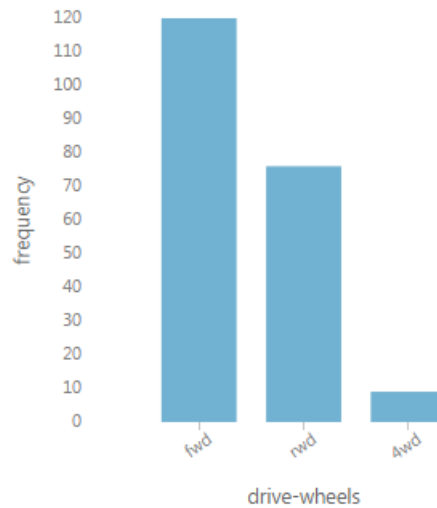


#### Visualizations

drive-wheels

Histogram

compare to None



frequency log scale

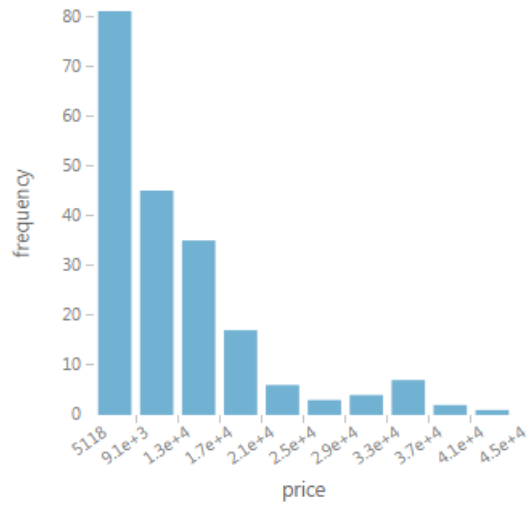
bins

10

price

Histogram

compare to  



- price log scale
- frequency log scale

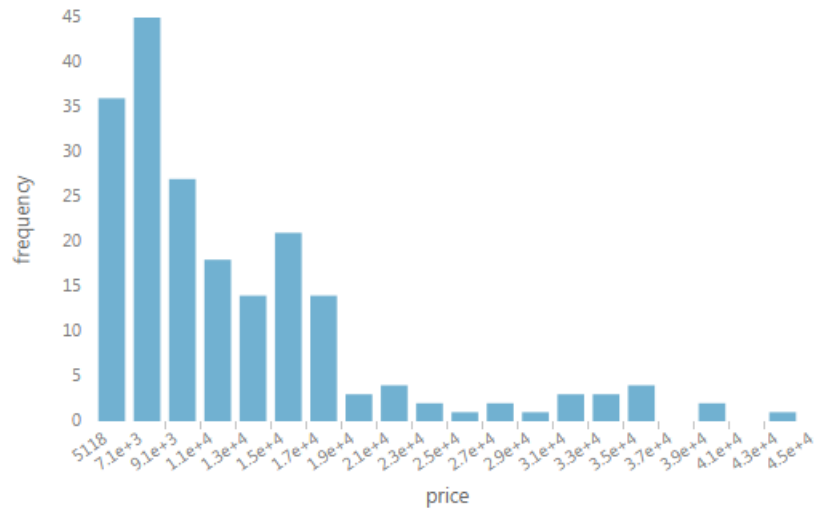
bins

- cumulative distribution
- probability density

price

Histogram

compare to



- price log scale
- frequency log scale

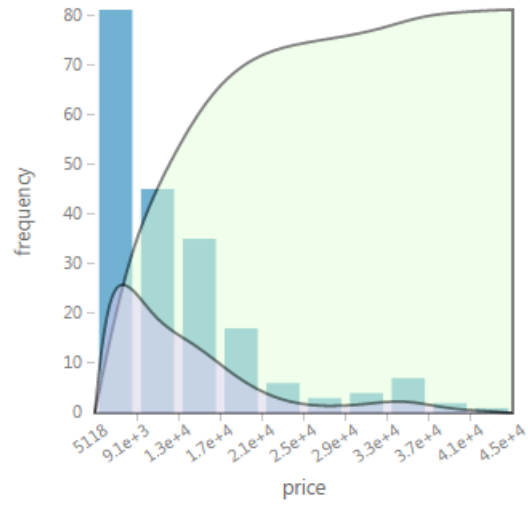
bins

- cumulative distribution
- probability density

price

Histogram

compare to




- price log scale
- frequency log scale

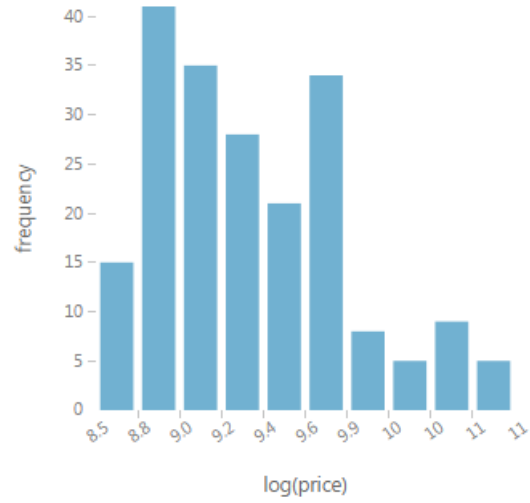
bins

- cumulative distribution
- probability density

price

Histogram

compare to  








price log scale

frequency log scale

bins

cumulative distribution

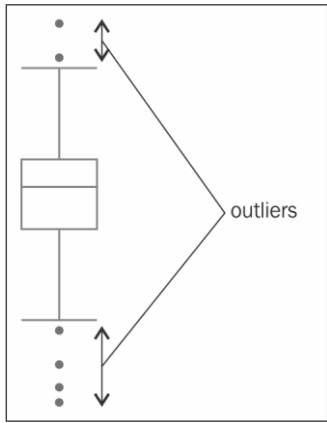
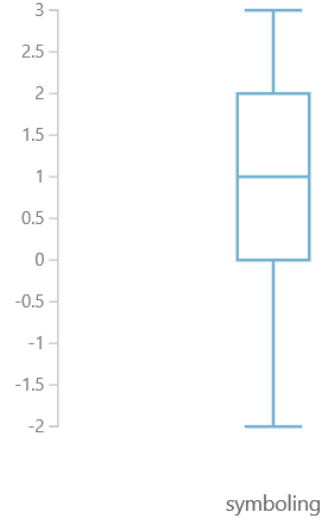
probability density

	symboling	normalized-losses	make	fuel-type	as
view as					
	3		alfa-romero	gas	sto
	3		alfa-romero	gas	sto
	1		alfa-romero	gas	sto
	2	164	audi	gas	sto
	2	164	audi	gas	sto
	2		audi	gas	sto
	1	158	audi	gas	sto
	1		audi	gas	sto
	1	158	audi	gas	tu
	0		audi	gas	tu
	2	192	bmw	gas	sto
	0	192	bmw	gas	sto

Visualizations

symboling  
BoxPlot

compare to



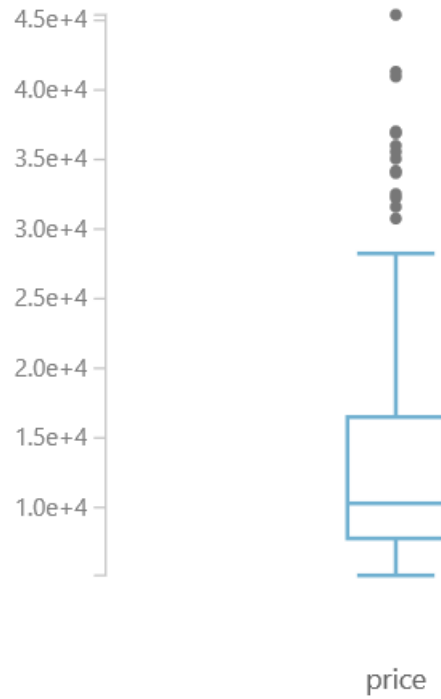


## Visualizations

price

BoxPlot

compare to   



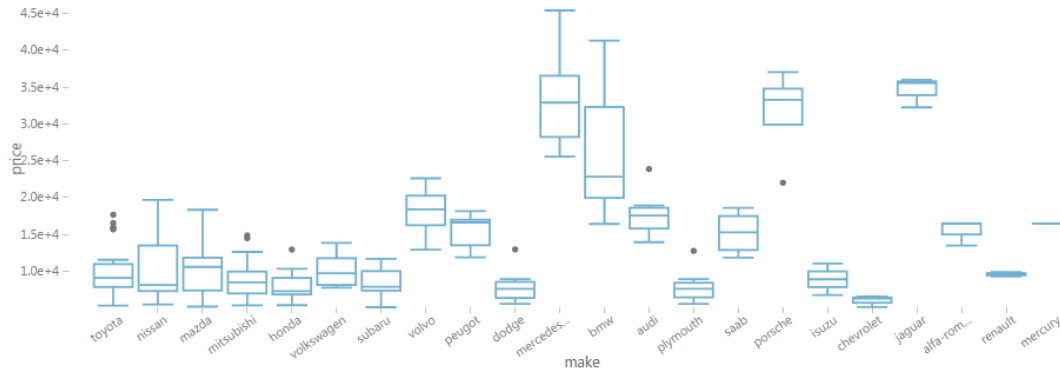
price log scale

Visualizations

price

MultiboxPlot

compare to



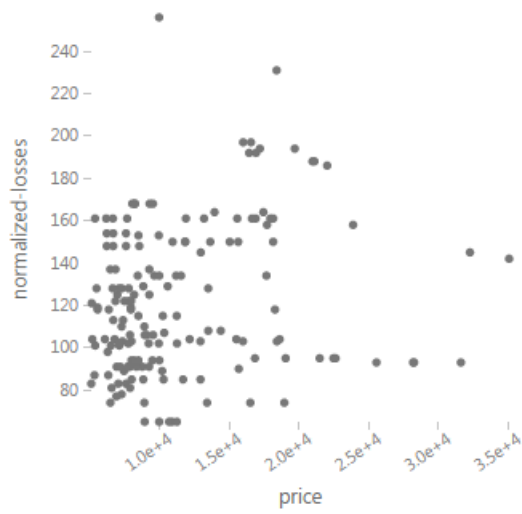
price log scale

categories

price

ScatterPlot

compare to



price log scale

normalized-losses log scale

## Visualizations

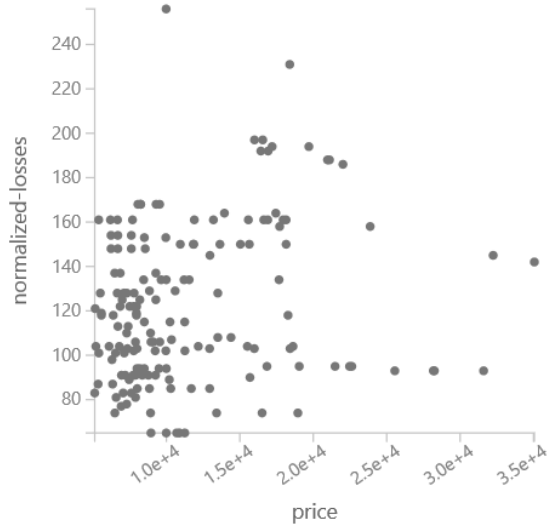
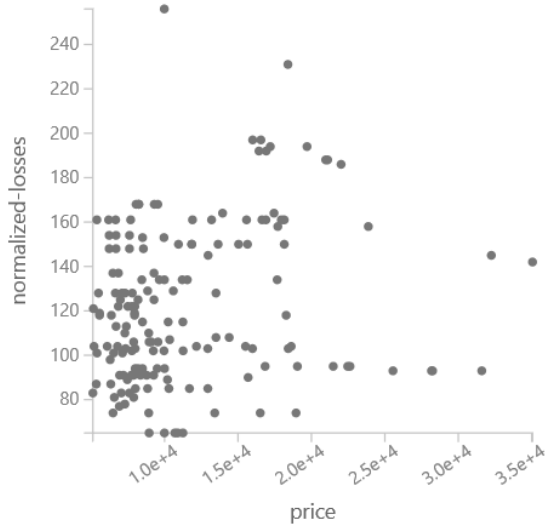
price

ScatterPlot

compare to



price



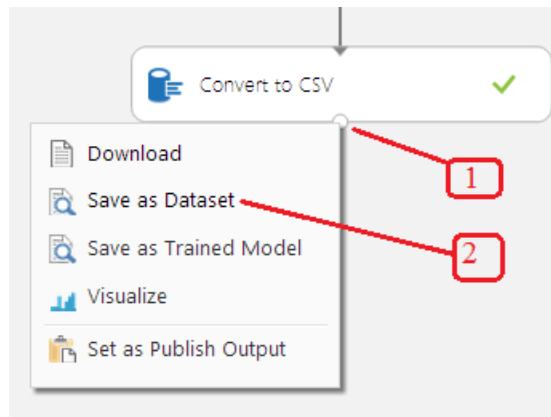
# Chapter 4

The screenshot shows the Orange3 software interface. On the left is a sidebar with icons for various data processing steps. The main workspace contains a workflow with a 'Reader' widget. A red box labeled '2' points to the 'Reader' widget in the sidebar, and another red box labeled '8' points to the 'Reader' widget in the workflow. The 'Properties' panel on the right is open for the 'Reader' widget, showing settings for 'Data source' (Web URL via HTTP), 'URL' (http://archive.ics.uci.edu/ml/machine-learning-databases/...), 'Data format' (CSV), and a checked option for 'CSV or TSV has header row'. A red box labeled '3' points to the 'Data source' dropdown, '4' to the 'URL' text field, '5' to the 'Data format' dropdown, and '6' to the 'CSV or TSV has header row' checkbox. A red box labeled '7' points to the 'RUN' button in the bottom toolbar.

This image shows a close-up of a workflow with two widgets: 'Reader' and 'Convert to CSV'. A context menu is open over the 'Convert to CSV' widget. A red box labeled '1' points to the context menu, and another red box labeled '2' points to the 'Download' option in the menu.

### Data Format Conversions

- Convert to ARFF
- Convert to CSV
- Convert to Dataset
- Convert to SVMlight
- Convert to TSV



#### Writer

Please specify data destination

Azure Blob Storage ▼

Please specify authentication type

Account ▼

Azure account name

YourAccountName

Azure account key

.....

Path to blob beginning with container

path/fileName.csv

Azure blob storage write mode

Overwrite ▼

File format for blob file

CSV ▼

Write blob header row

Search experiment items

Draft saved at 07:34:34 AM

**Data Input and Output**

- Enter Data
- Reader
- Writer

Data Transformation

Feature Selection

**Enter Data**

DataFormat: CSV

HasHeader

Data

```

1 Name, Age, Annual Income
2 Robin, 26, 30000
3 Harry, 34, 28000
4 David, 45, 56000
  
```

**Saved Datasets**

**My Datasets**

ch4\_example1.csv

ch4\_example1.csv

WEB SERVICES

**DATASETS**

TRAINED MODELS

MY DATASETS    SAMPLES

NAME	SUBMITTED BY	DESCRIPTION	DATA TYPE
ch4_example1.csv	sumitmund		GenericCSV
perfume_data.txt	sumitmund		GenericCSVNoHeader
mnist_test.csv	sumitmund		GenericCSV

NEW    DOWNLOAD    DELETE    GENERATE DATA ACCESS CODE...



#### Writer

Please specify data destination

Azure Blob Storage

Please specify authentication type

Account

Azure account name

YourAccountName

Azure account key

.....

Path to blob beginning with container

path/fileName.csv

Azure blob storage write mode

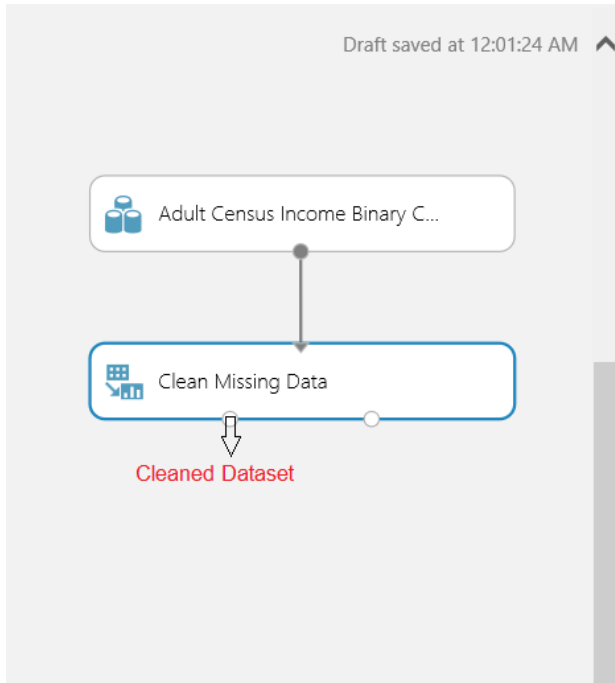
Overwrite

File format for blob file

CSV

Write blob header row

## Chapter 5



### Clean Missing Data

Columns to be cleaned

**Selected columns:**

All columns

Launch column selector

Minimum missing value ratio

0

Maximum missing value ratio

1

Cleaning mode

Custom substitution value

Replacement value

0

Generate missing value indicator column

Project Columns

Select columns

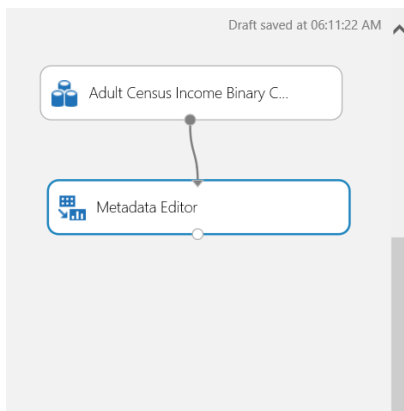
Allow duplicates and preserve column order in selection

Begin With: All columns

Exclude: column names

workclass X

Launch column selector



### Metadata Editor

Column

**Selected columns:**

Column names: education-num

Launch column selector

Data type

Unchanged

Categorical

Categorical

Fields

Features

New column names

Comma separated list of column names or zero-based indices selected using *column selector*

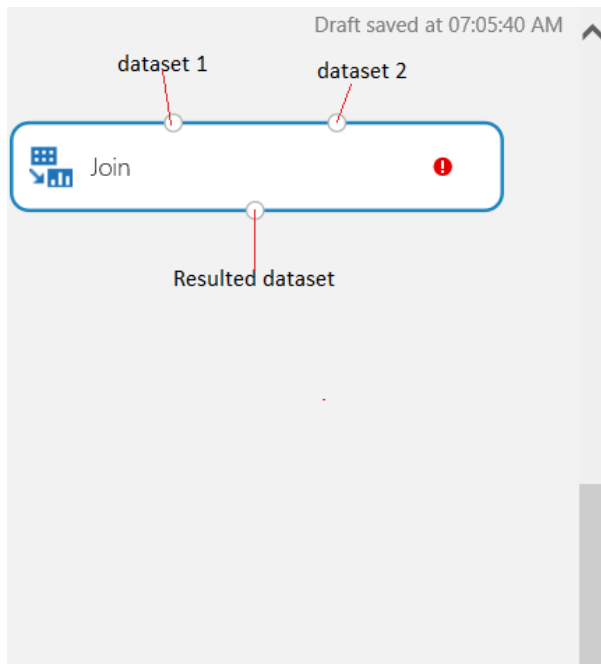
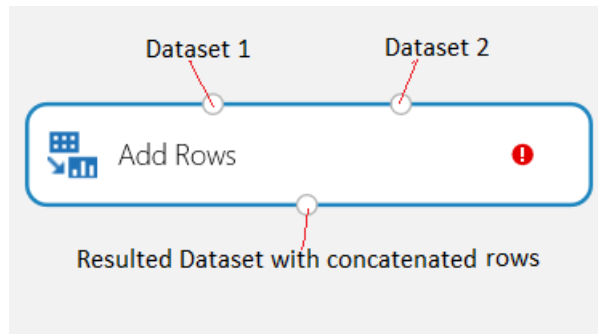
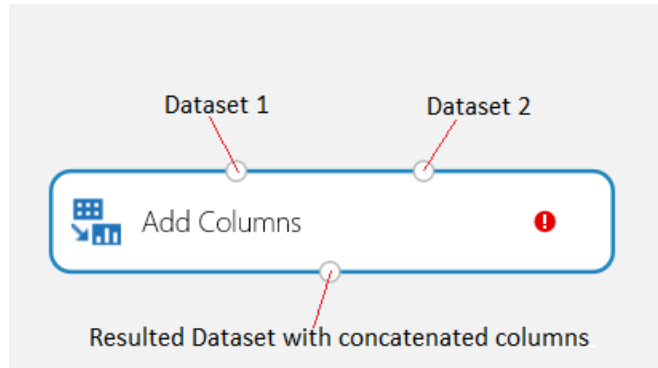
The new data type of the column(s)

Whether the column(s) should be considered categorical

Whether the column(s) should be considered features or labels by learning algorithms

Comma separated list of new column name(s)





#### Join

Join key columns for L

Key columns for Left dataset (dataset 1)

**Selected columns:**

Launch the selector tool to make a selection

Launch column selector

Join key columns for R

Key columns for Right dataset (dataset 2)

**Selected columns:**

Launch the selector tool to make a selection

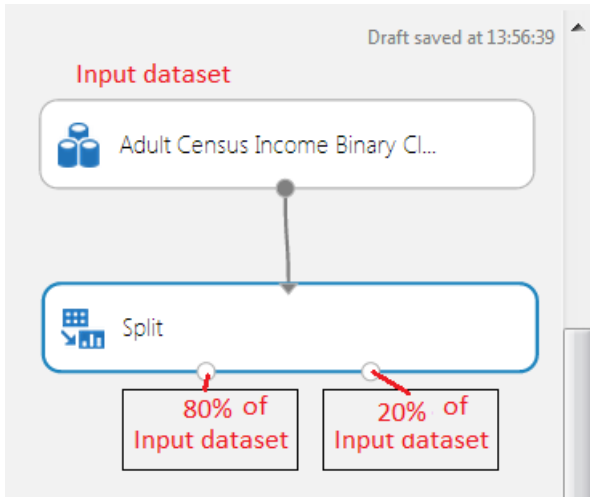
Launch column selector

Match case

Join type

Left Outer Join

Keep right key columns in joined table



**Split**

Splitting mode

Split Rows

Fraction of rows in the first output dataset

0.8

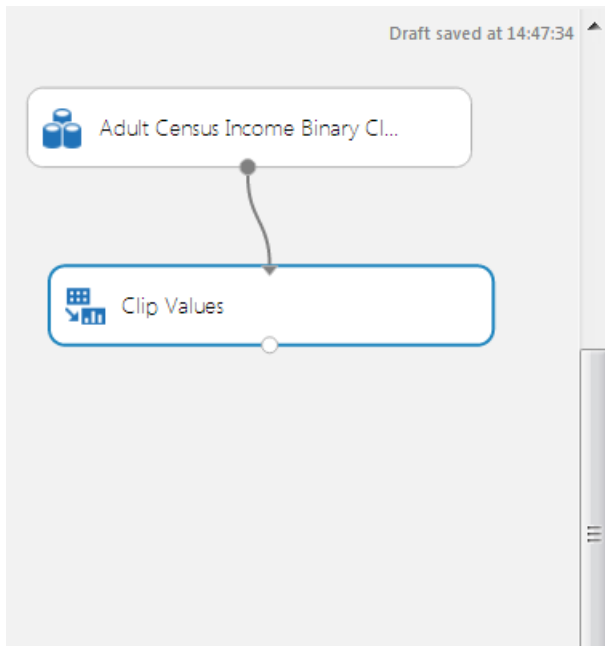
Randomized split

Random seed

0

Stratified split

False



**Clip Values**

Set of thresholds

ClipPeaks

Upper threshold

Percentile

Percentile number for upper threshold

99

Upper substitute value

Threshold

List of columns

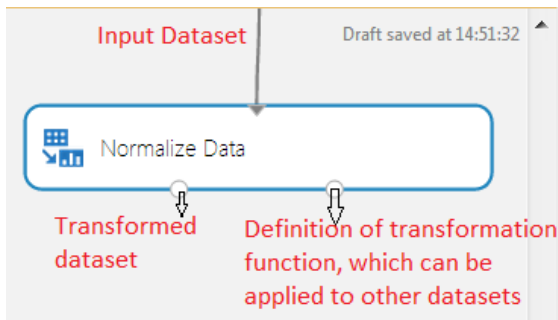
**Selected columns:**

**Column type:** Numeric, All

Launch column selector

Overwrite flag

Add indicator columns



**Normalize Data**

Transformation method

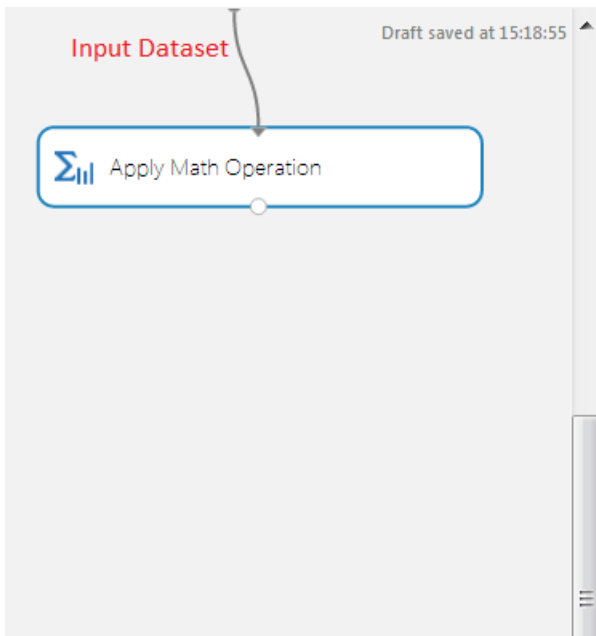
ZScore

Columns to transform

**Selected columns:**

**Column type:** Numeric, All

Launch column selector



#### Apply Math Operation

Category: Basic

Basic math function: Pow

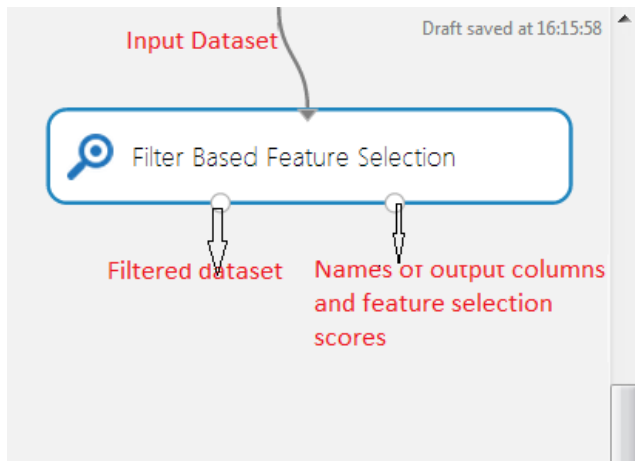
Second argument type: Constant

Constant second argument: 2

Column set: Selected columns: Column type: Numeric, All

Launch column selector

Output mode: ResultOnly



#### Filter Based Feature Selection

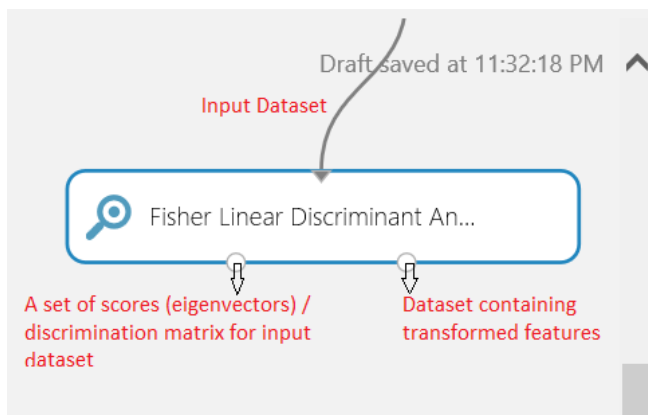
Feature scoring method: Pearson Correlation

Operate on feature columns only

Target column: Selected columns: Column names: income

Launch column selector

Number of desired features: 3

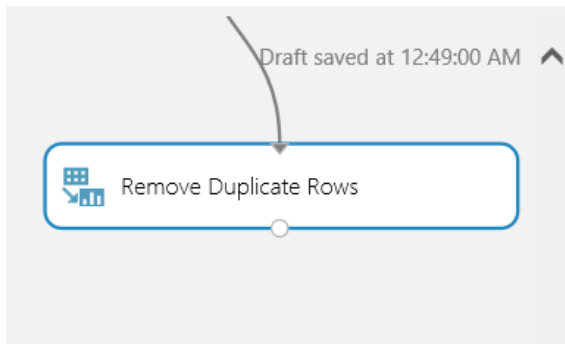
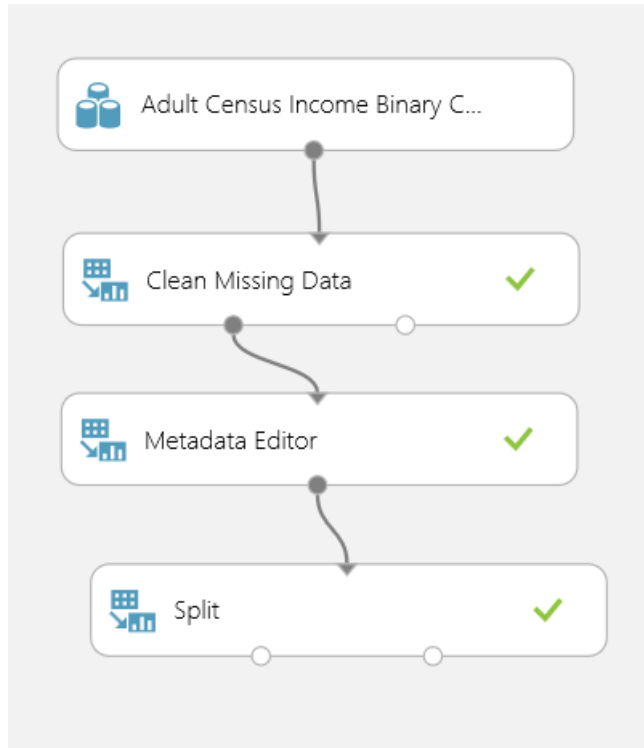


#### Fisher Linear Discriminant Analysis

Class labels column: Selected columns: Column names: income

Launch column selector

Number of feature extractors: 3



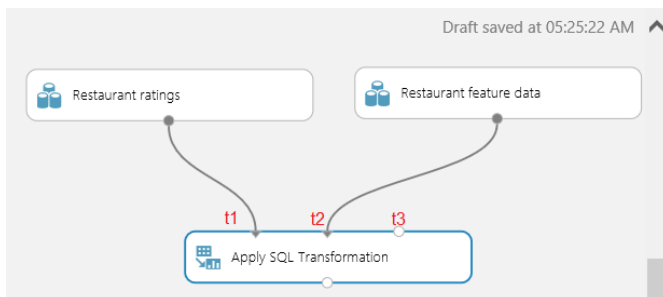
#### Remove Duplicate Rows

Key column selection filter expression

**Selected columns:**  
All columns

Launch column selector

Retain first duplicate row



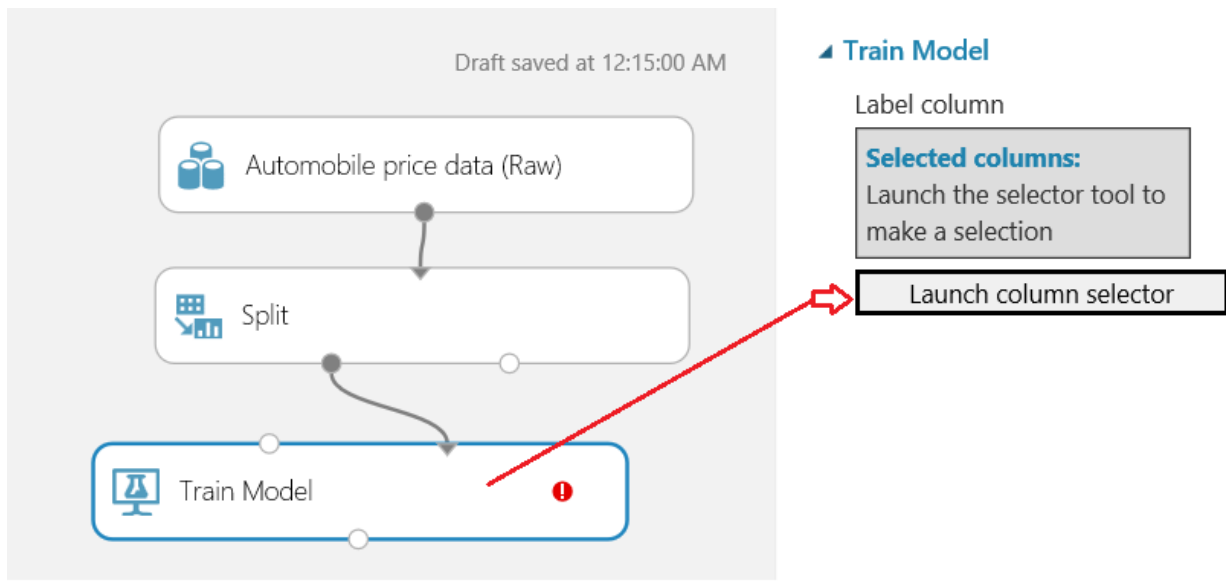
#### Apply SQL Transformation

SQL Query Script

```

1 SELECT DISTINCT t1.placeid, t2.Name,
2     t2.City, AVG(rating) as 'AvgRating'
3 from t1 JOIN t2 ON t1.placeid = t2.placeid
4 group by t1.placeid;
  
```

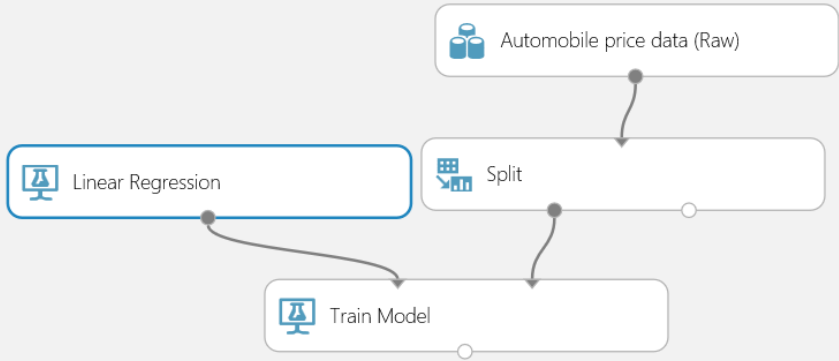
# Chapter 6



Select a single column ✕

Include  column names  price

Draft saved at 06:22:57 AM



#### Linear Regression

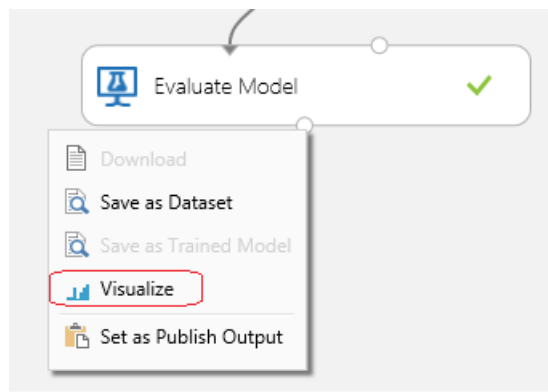
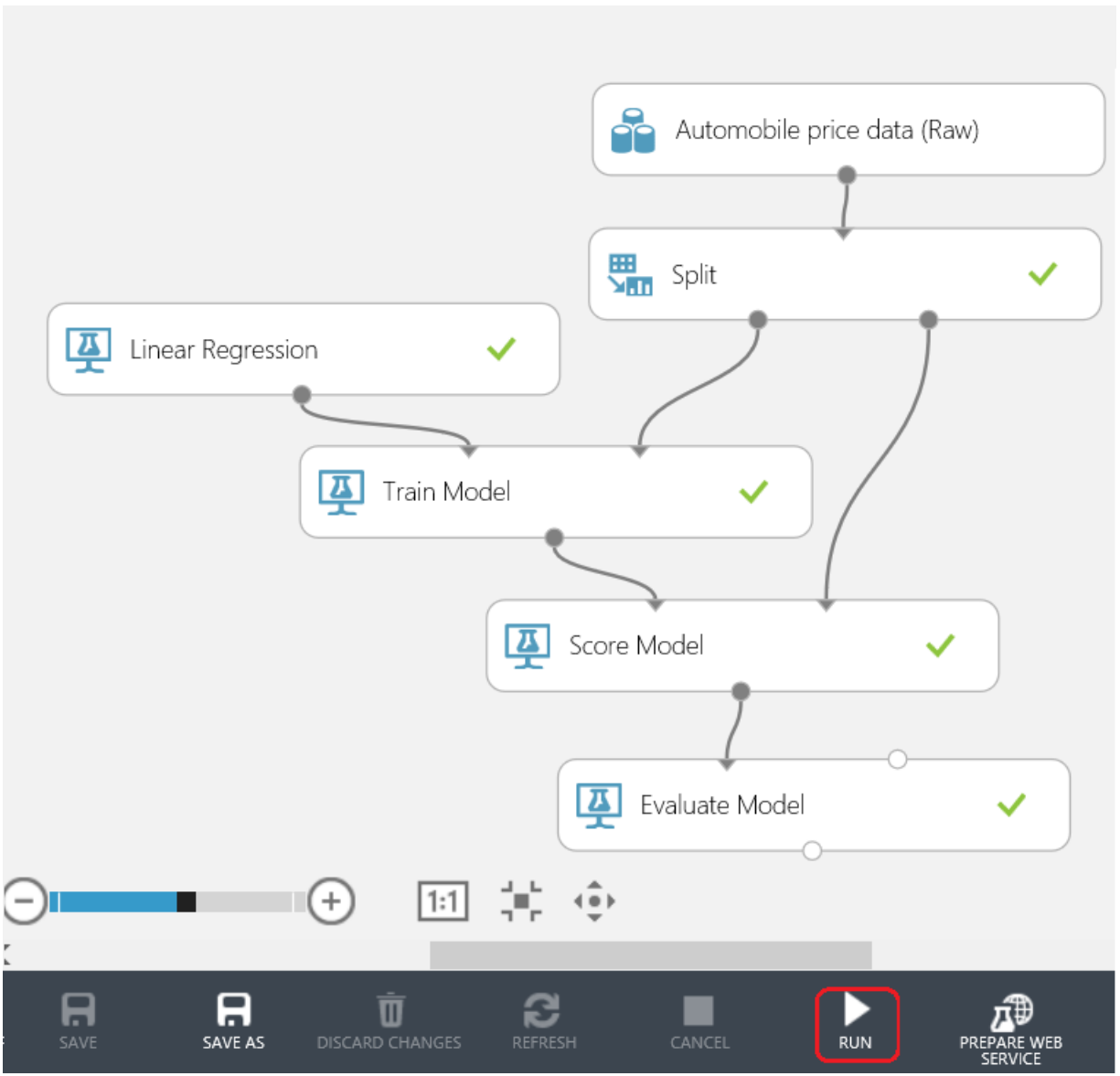
Solution method  
Ordinary Least Squares

L2 regularization weight  
0.001

Include intercept term

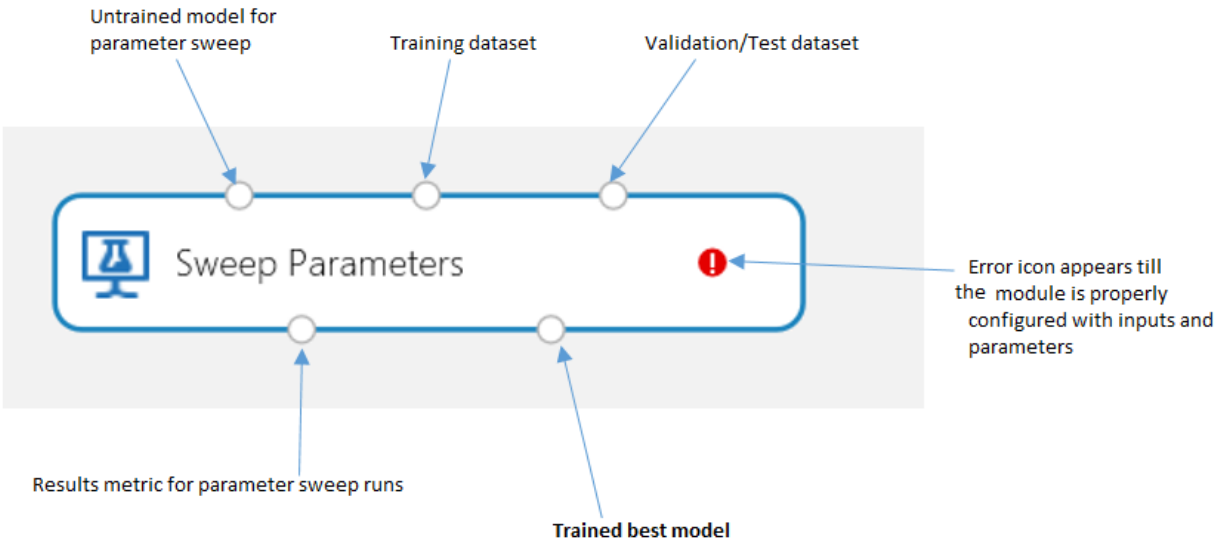
Random number seed  
0

Allow unknown categori...



Demo-LinearRegression > Evaluate Model > Evaluation results

rows	columns					
1	5					
		Mean Absolute Error	Root Mean Squared Error	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
view as						
		1700.399749	2484.018779	0.189728	0.058297	0.941703



In draft  
Draft saved at 07:07:28 AM

### Properties

**Sweep Parameters**

Specify parameter sweeping mode

Label column

**Selected columns:**  
 Column names:

Metric for measuring performance for classification

Metric for measuring performance for regression



Draft saved at 07:15:02 AM

Decision Forest Regression

### Decision Forest Regression

Resampling method: **Bagging**

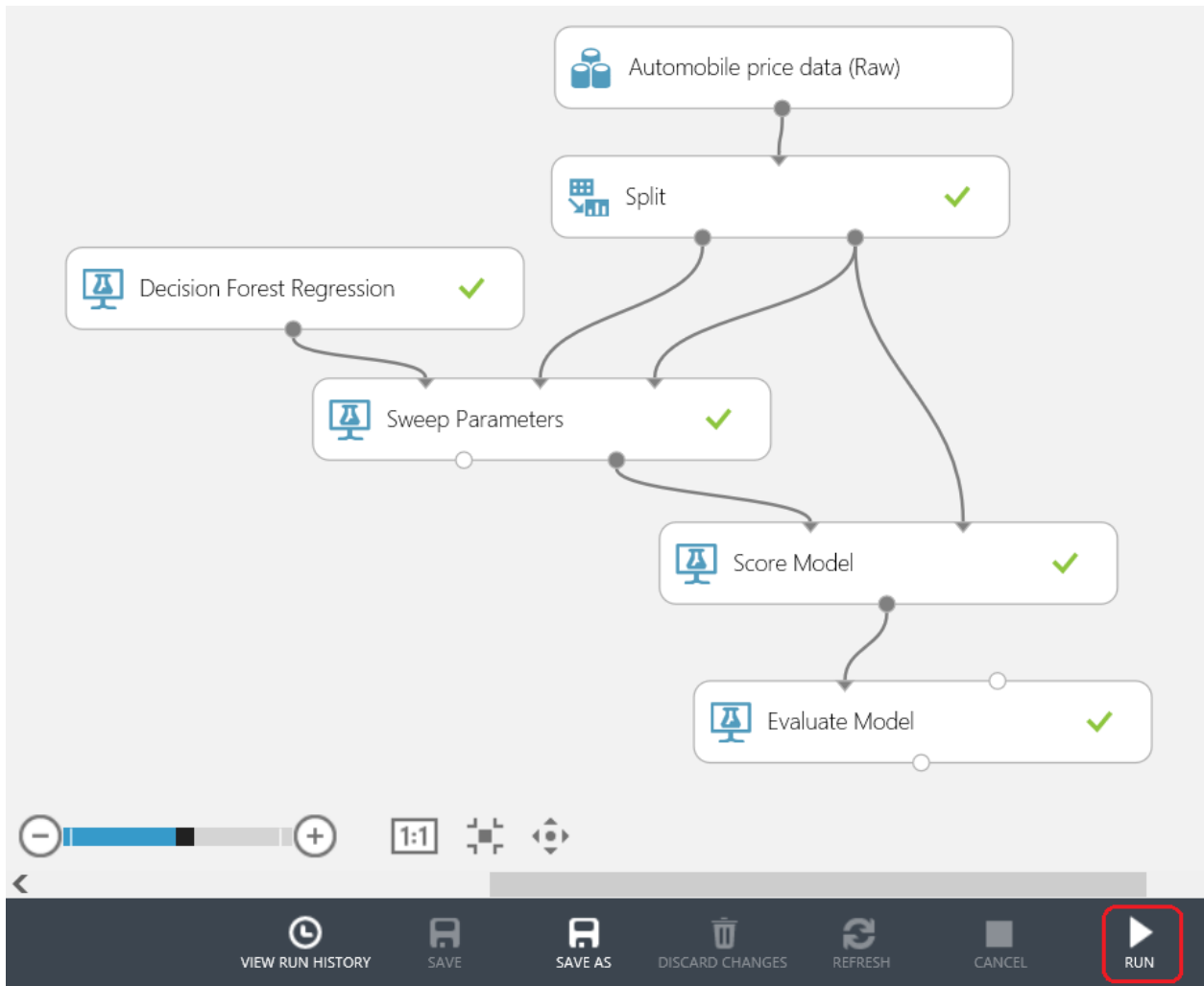
Number of decision trees: 8

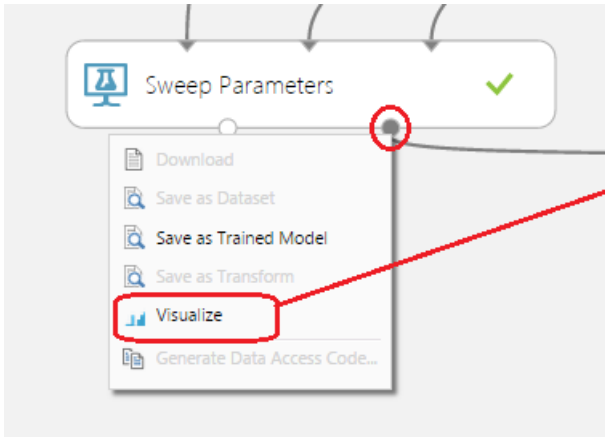
Maximum depth of the decision trees: 32

Number of random splits per node: 128

Minimum number of samples per leaf node: 1

Allow unknown values for categorical feat...

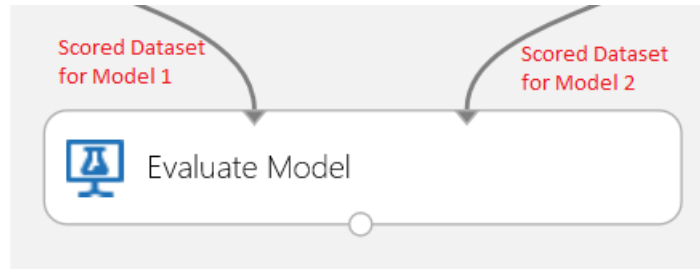


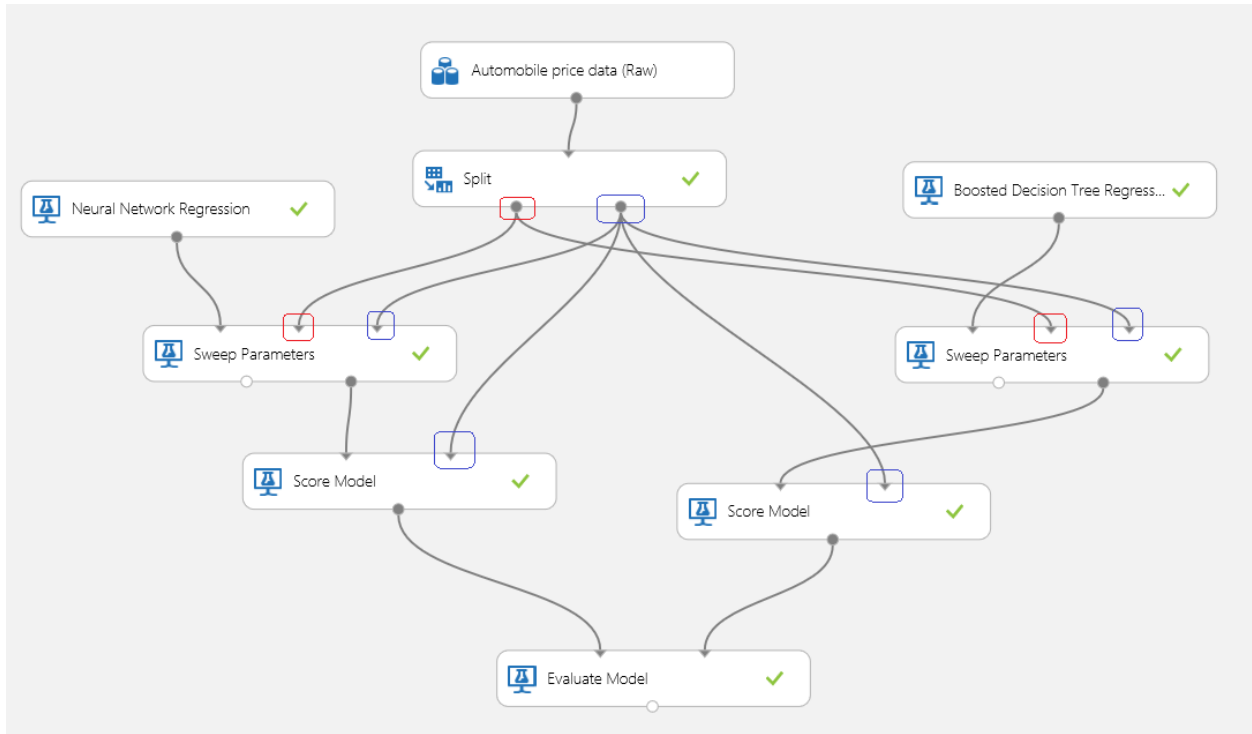


### Gemini Decision Forest Regressor

#### Settings

Setting	Value
Min Leaf Sample Count	4
Random Split Count	128
Max Depth	16
Ensemble Element Count	32
Class Count	1
Resampling Method	Bagging
Random Number Seed	5
Allow Unknown Levels	True





Mean Absolute Error	Root Mean Squared Error	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
5071.800876	7664.589306	0.565904	0.555025	0.444975
1340.164661	1958.367584	0.149534	0.036235	0.963765

# Chapter 7

Search experiment items

► Data Format Conversions

◄ Data Input and Output

- Enter Data
- Reader
- Writer

► Data Transformation

Draft saved at 01:17:16 AM

Reader

Reader

Web URL via HTTP

URL  
http://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data

Data format  
CSV

CSV or TSV has header row

Draft saved at 01:27:00 AM

Reader

Clean Missing Data

### Clean Missing Data

Columns to be cleaned

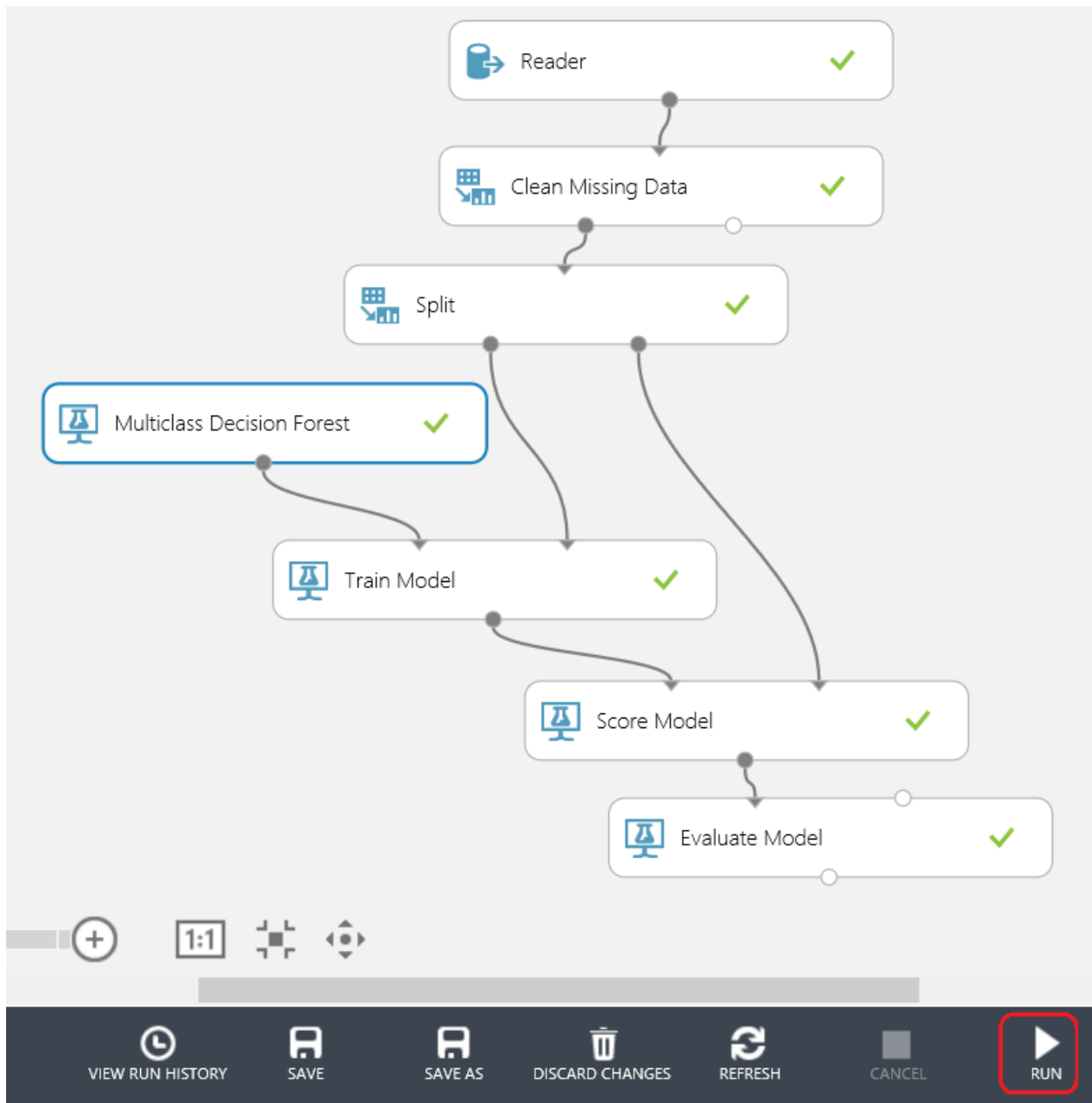
**Selected columns:**  
All columns

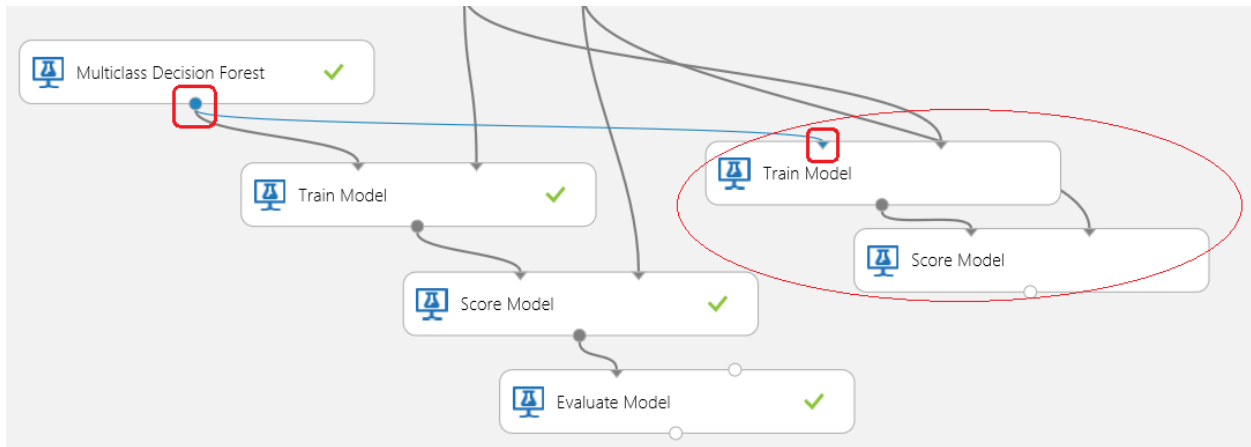
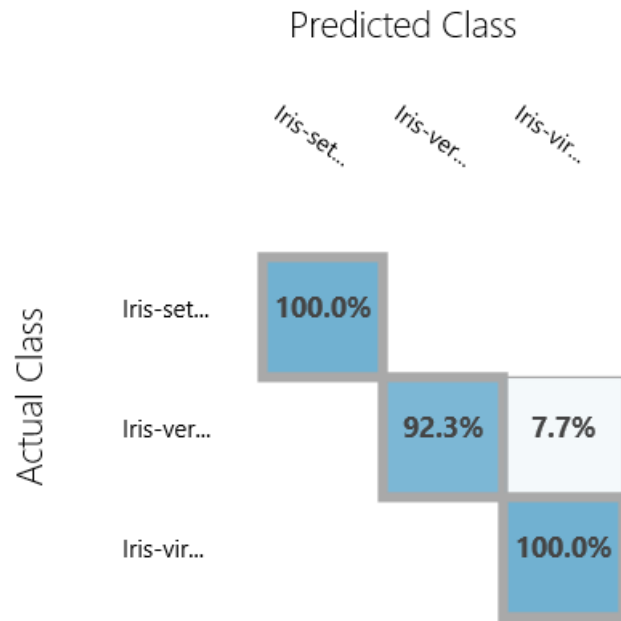
Launch column selector

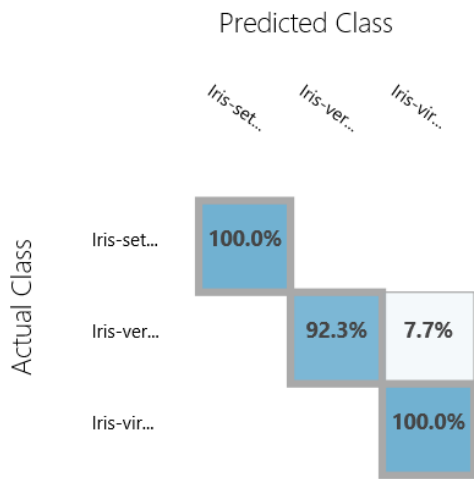
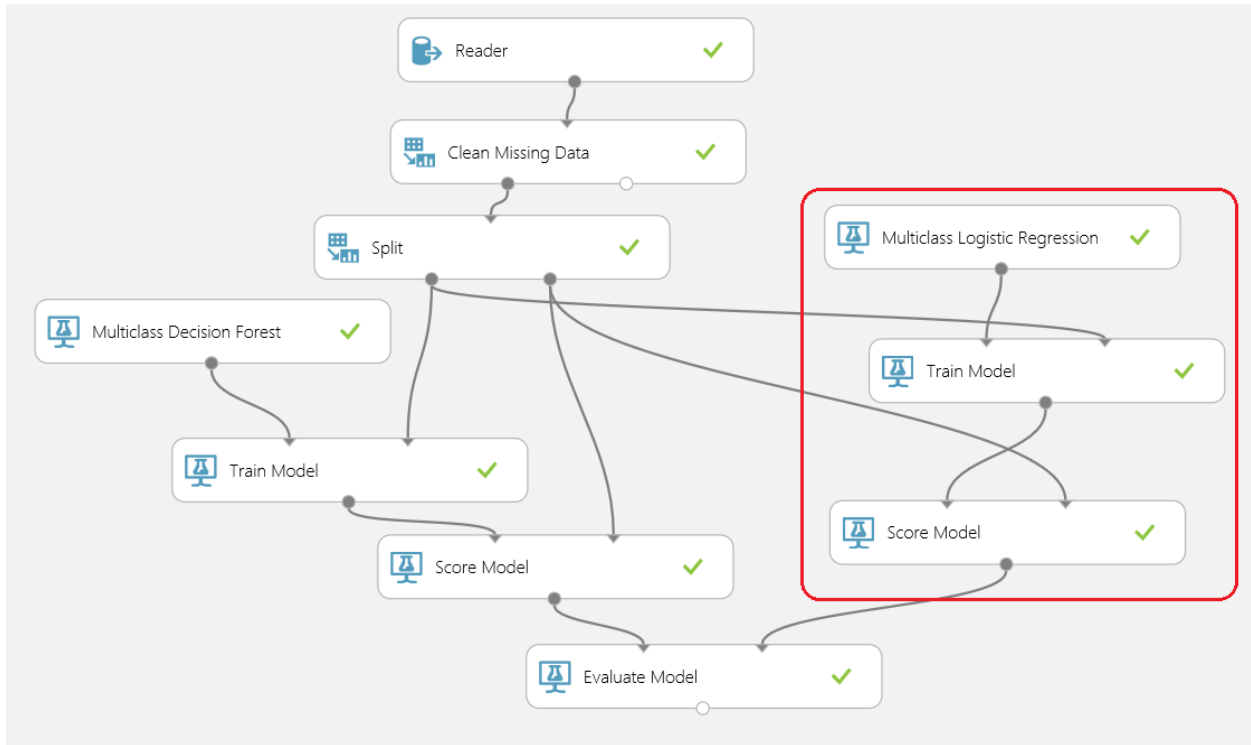
Minimum missing value ratio  
0

Maximum missing value ratio  
1

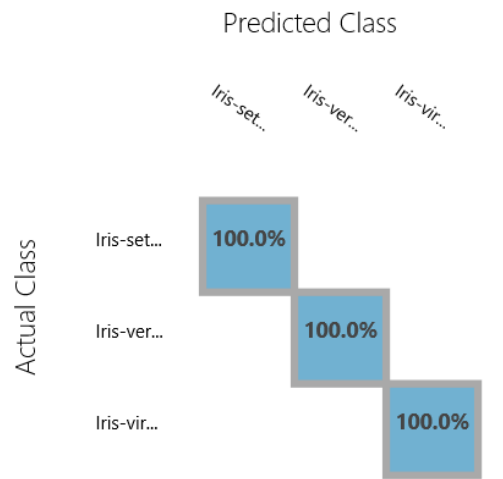
Cleaning mode  
Remove entire row



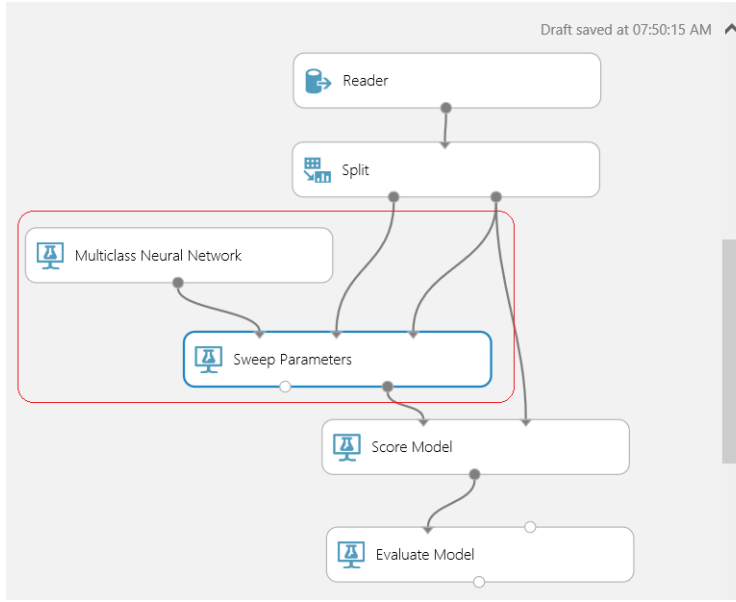




↑  
Multiclass Decision Forest



↑  
Multiclass Logistic Regression



#### Sweep Parameters

Specify parameter sweeping mode

Entire grid

Label column

Selected columns:

Column names: Col1

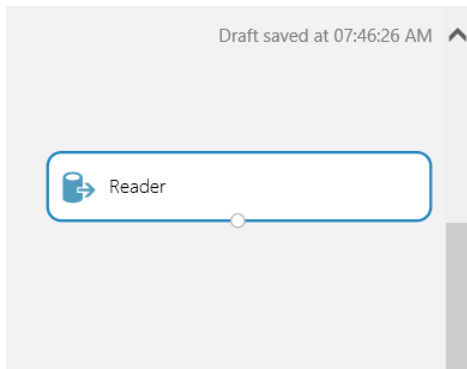
Launch column selector

Metric for measuring performance for classification

Accuracy

Metric for measuring performance for regression

Mean absolute error



#### Reader

Data source

Web URL via HTTP

URL

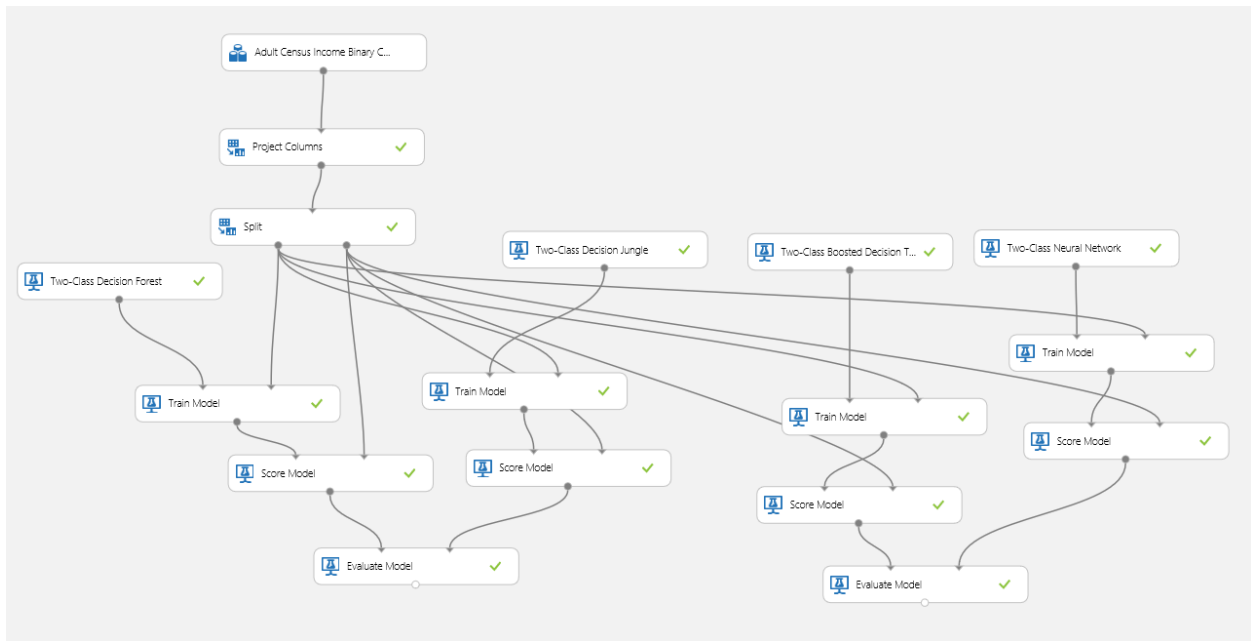
http://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data

Data format

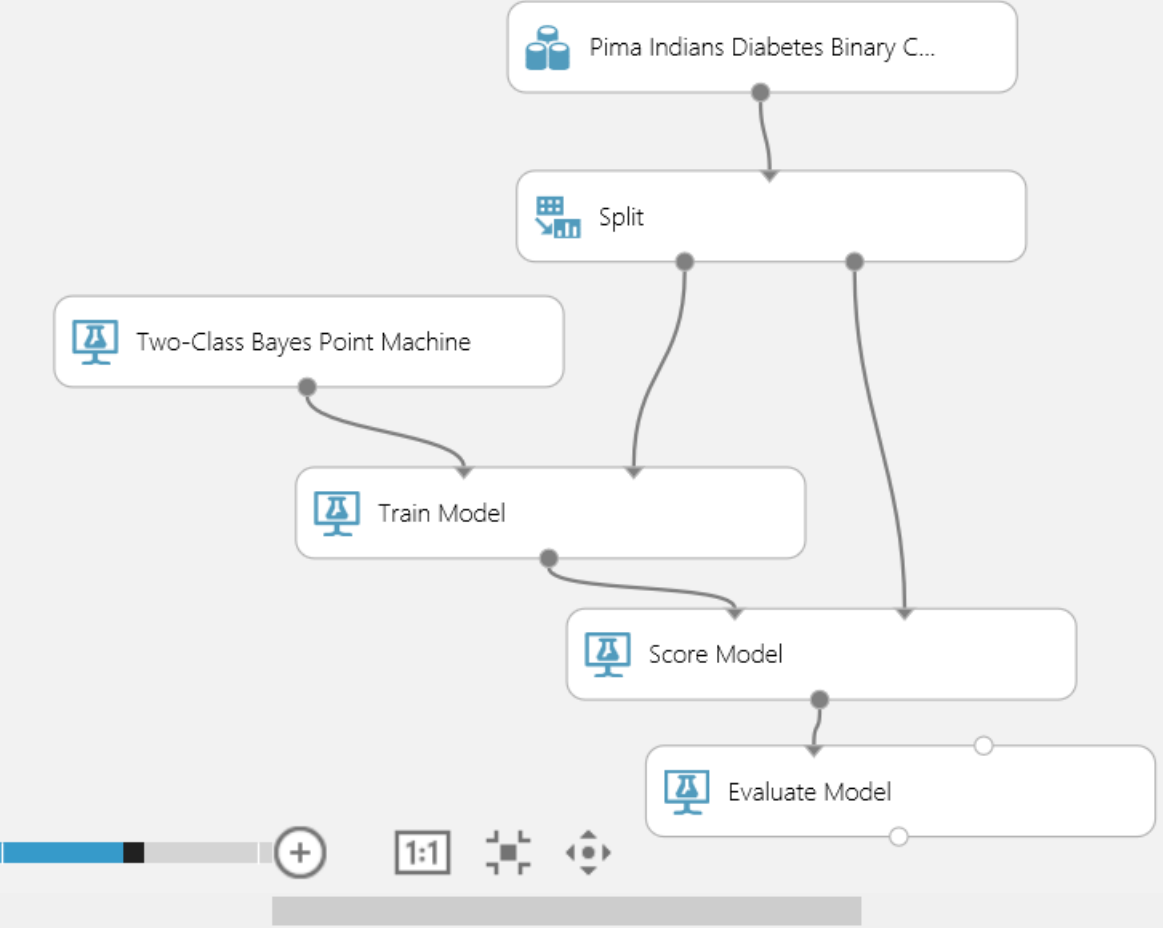
CSV

CSV or TSV has header row





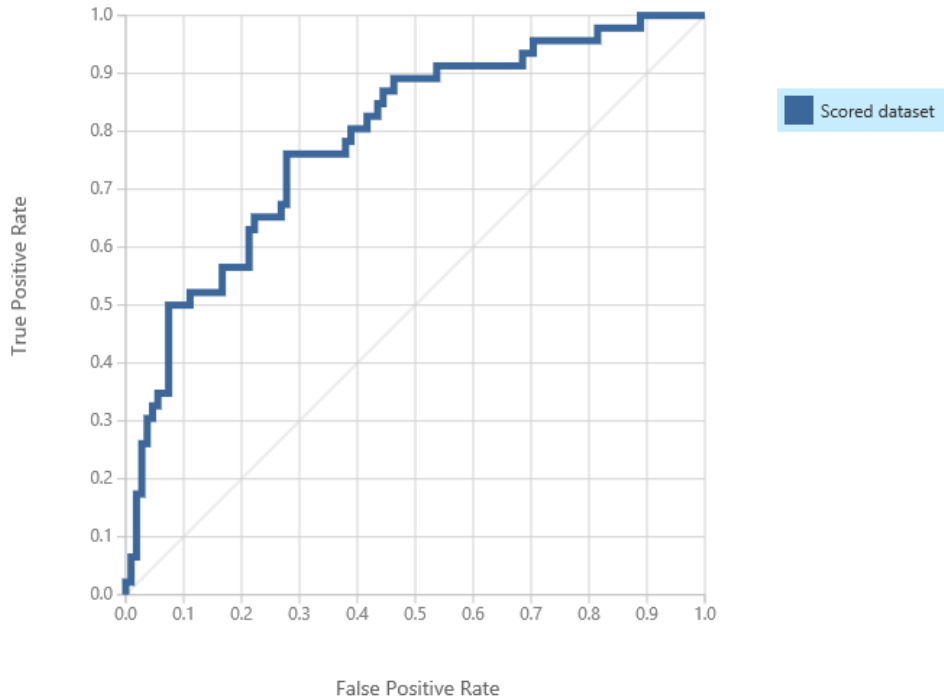
Draft saved at 12:29



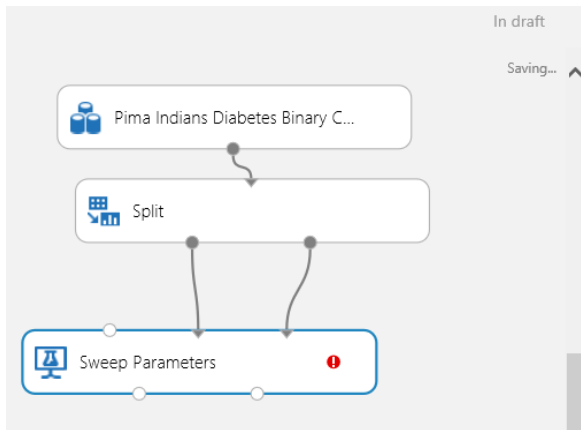
SAVE    SAVE AS    DISCARD CHANGES    REFRESH    CANCEL    **RUN**    PREPARE WEB SERVICE

Classification - Pima Indians Diabetes Dataset > Evaluate Model > Evaluation results

ROC PRECISION/RECALL LIFT



True Positive	False Negative	Accuracy	Precision	Threshold	AUC
23	23	0.779	0.676	0.51	0.788
False Positive	True Negative	Recall	F1 Score		
11	97	0.500	0.575		



**Properties**

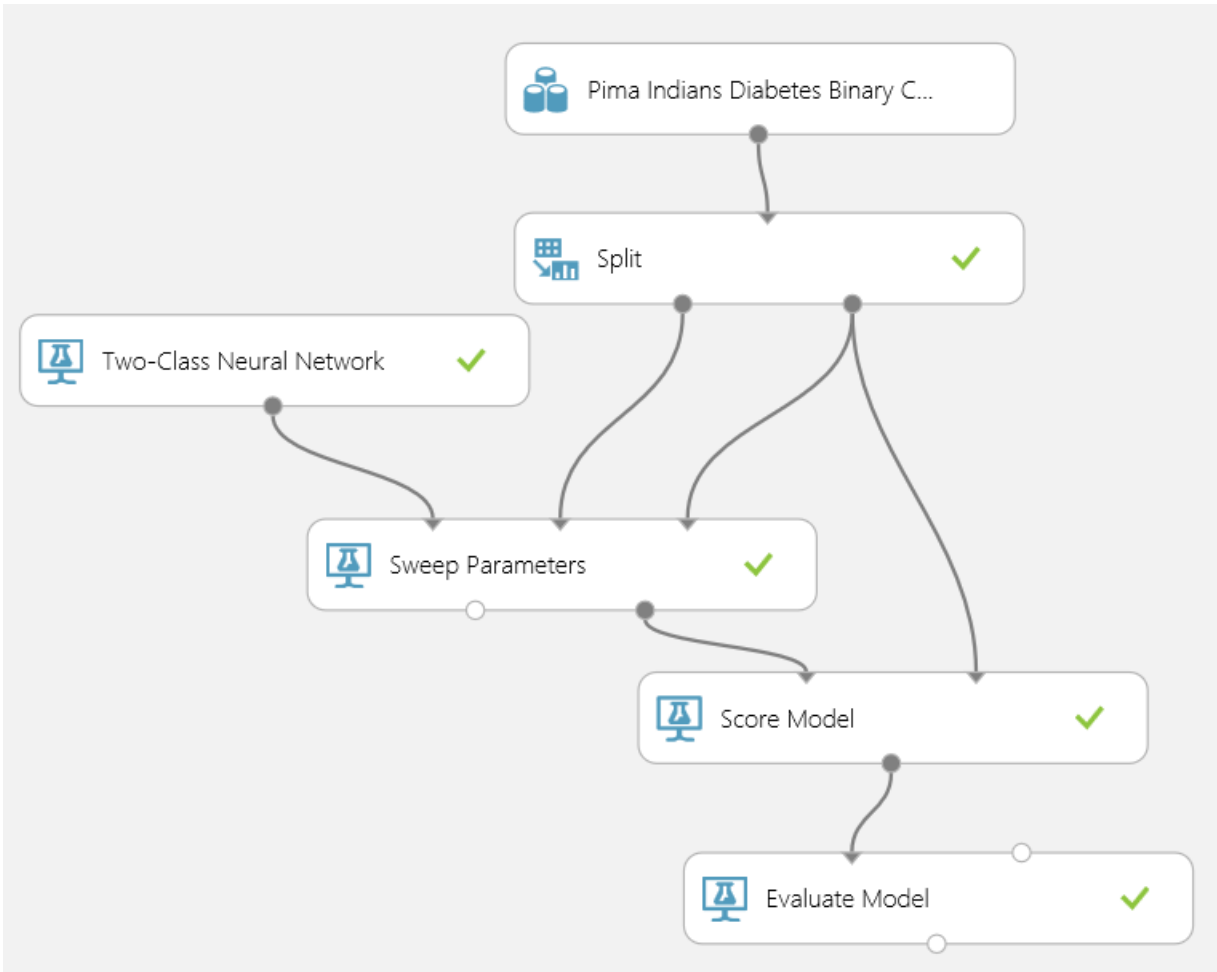
**Sweep Parameters**

Specify parameter sweeping mode: Entire grid

Label column: Selected columns: Column names: Class variable (0 or 1)

Metric for measuring performance for classification: Accuracy

Metric for measuring performance for regression: Mean absolute error



Classification - Adult Census Income > Adult Census Income Binary Classification dataset > dataset

rows 32561  
columns 15

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	ca	ge
view as												
	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	21	
	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male	0	

Statistics

Unique Values	9
Missing Values	1836
Feature Type	String

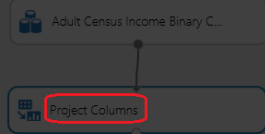
Visualizations

# Classification - Adult Census Income (DF)

In draft

Properties

Draft saved at 11:56:17 AM



Select columns

Allow duplicates and preserve column order in selection

Begin With: All columns

Exclude: column names

workclass X occupation X native-country X

Project Columns

Select columns

Selected columns: All columns

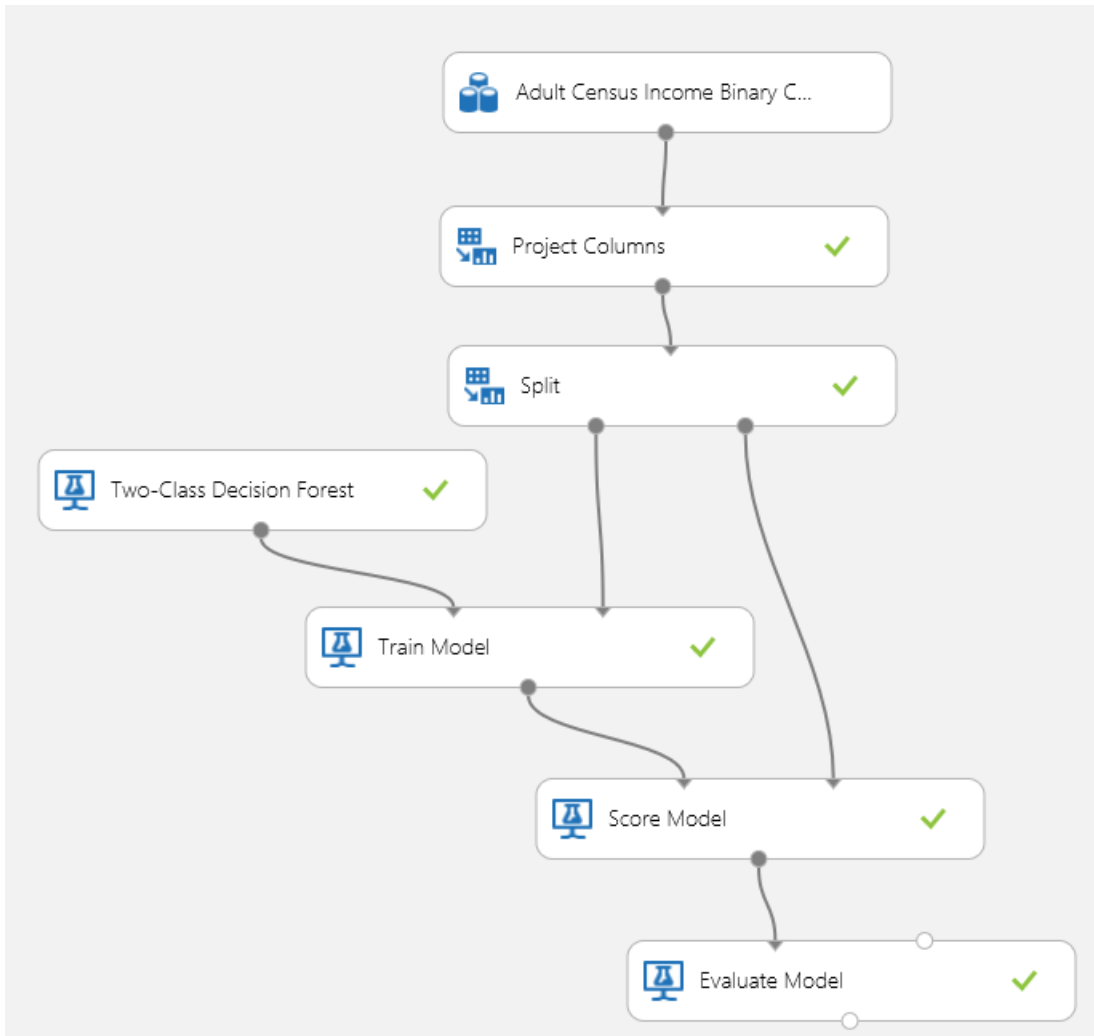
Exclude column names: workclass, occupation, native-country

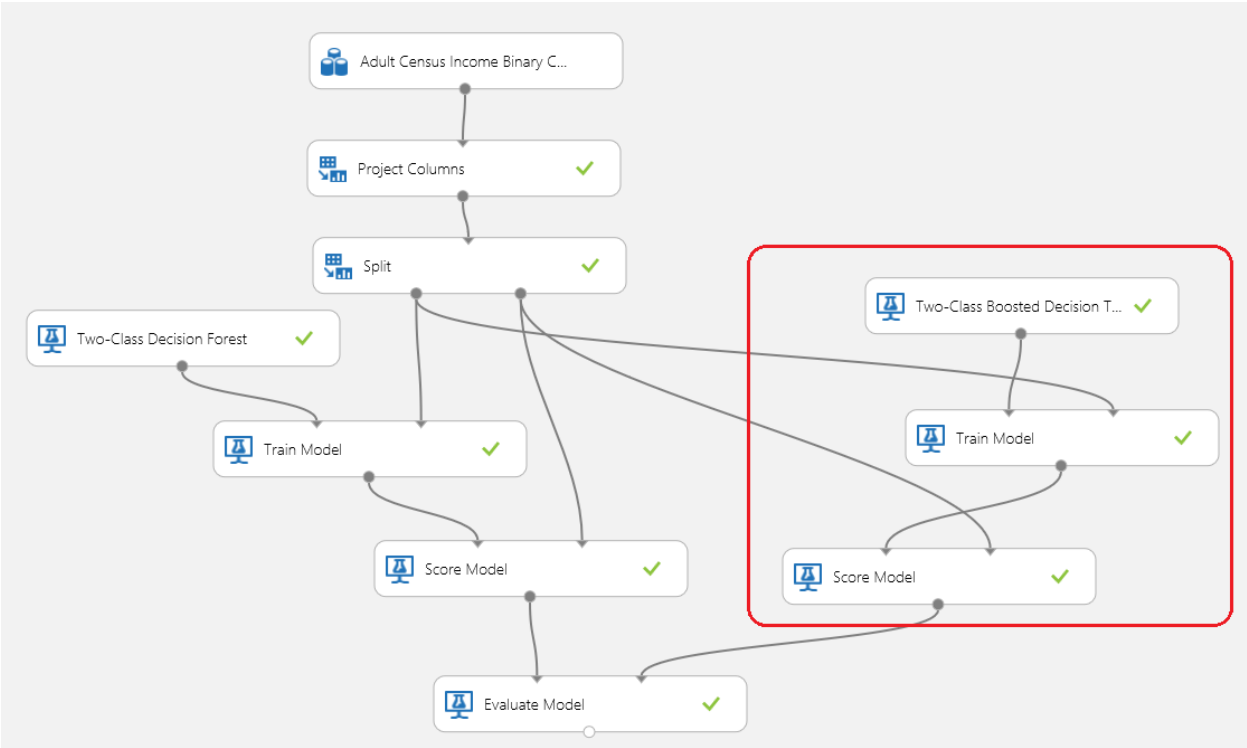
Launch column selector

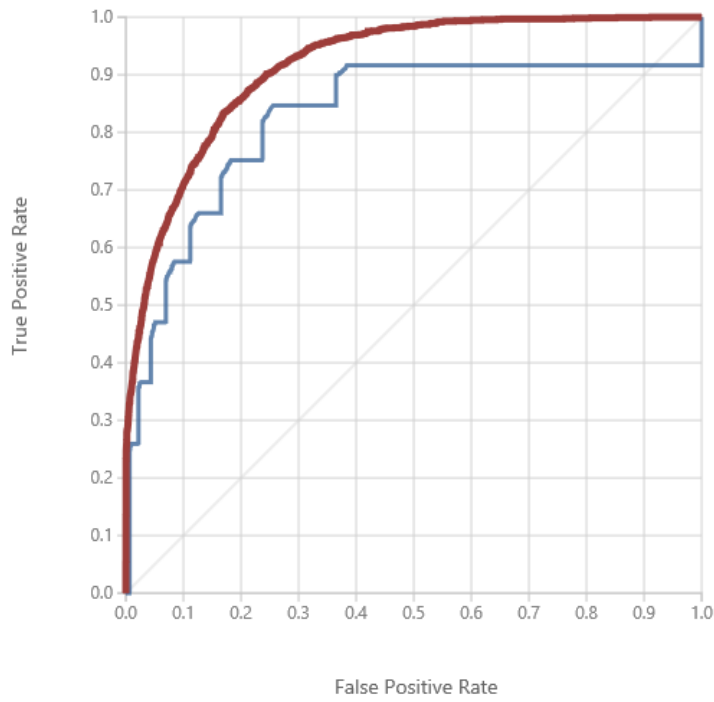
Experiment Properties

STATUS CODE: InDraft

Disable upgrades






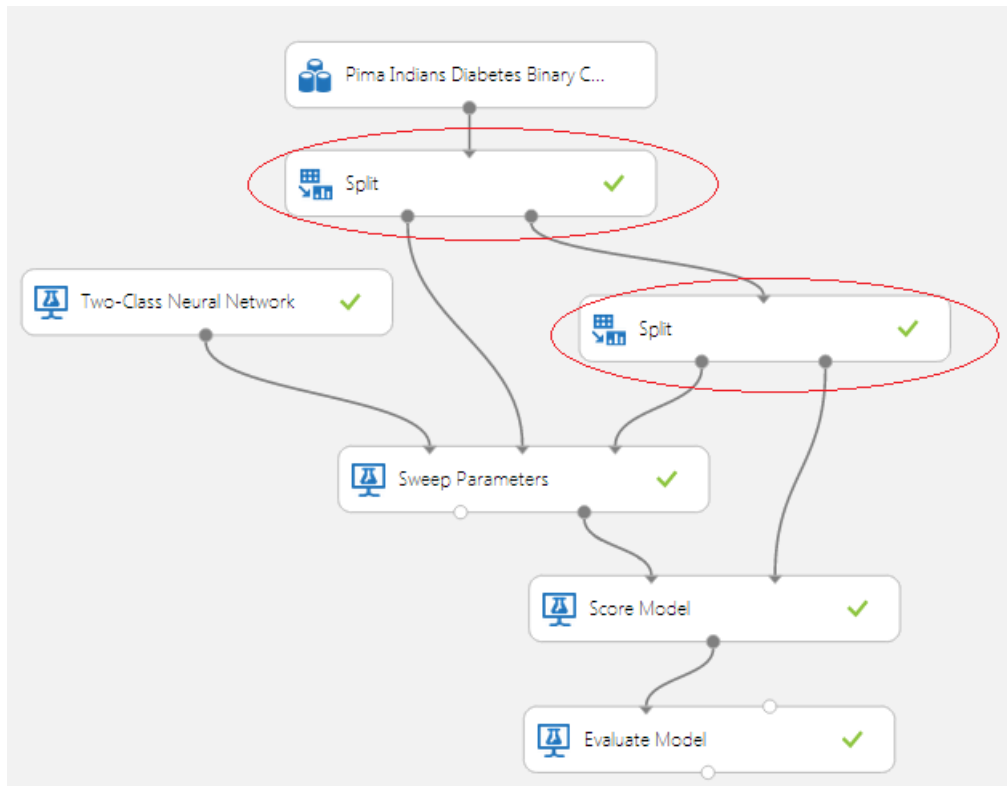
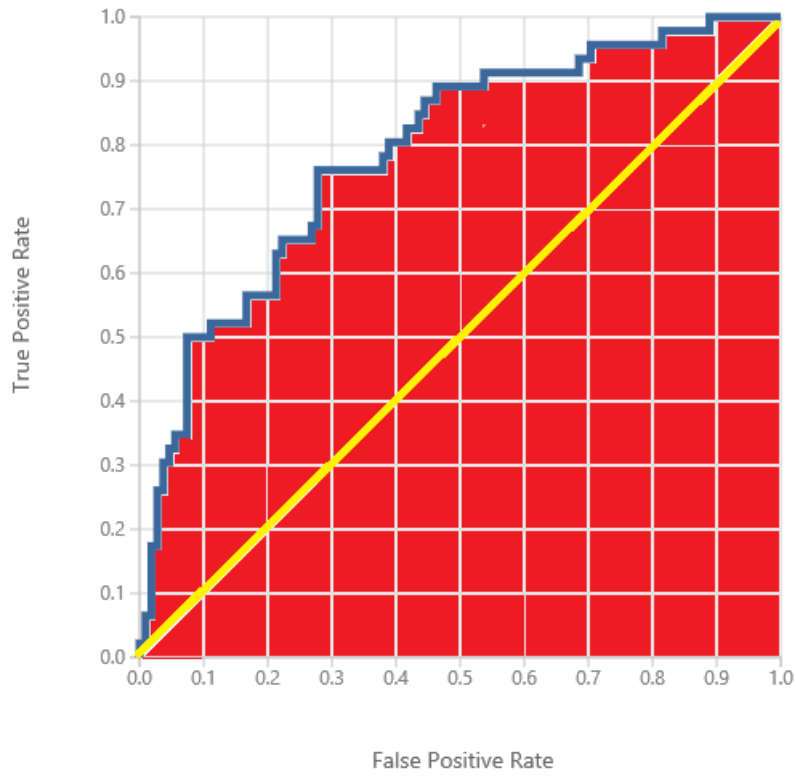


Model with *Two-Class Decision Forest*

Scored dataset  
Scored dataset to compare

Model with *Two-Class Boosted Decision Tree*

True Positive	False Negative	Accuracy	Precision	Threshold		AUC
<b>1017</b>	<b>561</b>	<b>0.859</b>	<b>0.741</b>	<b>0.5</b>		<b>0.916</b>
False Positive	True Negative	Recall	F1 Score			
<b>356</b>	<b>4578</b>	<b>0.644</b>	<b>0.689</b>			





Draft saved at 16:16:53

Two-Class Bayes Point Machine

Number of training iterations: 30

Include bias

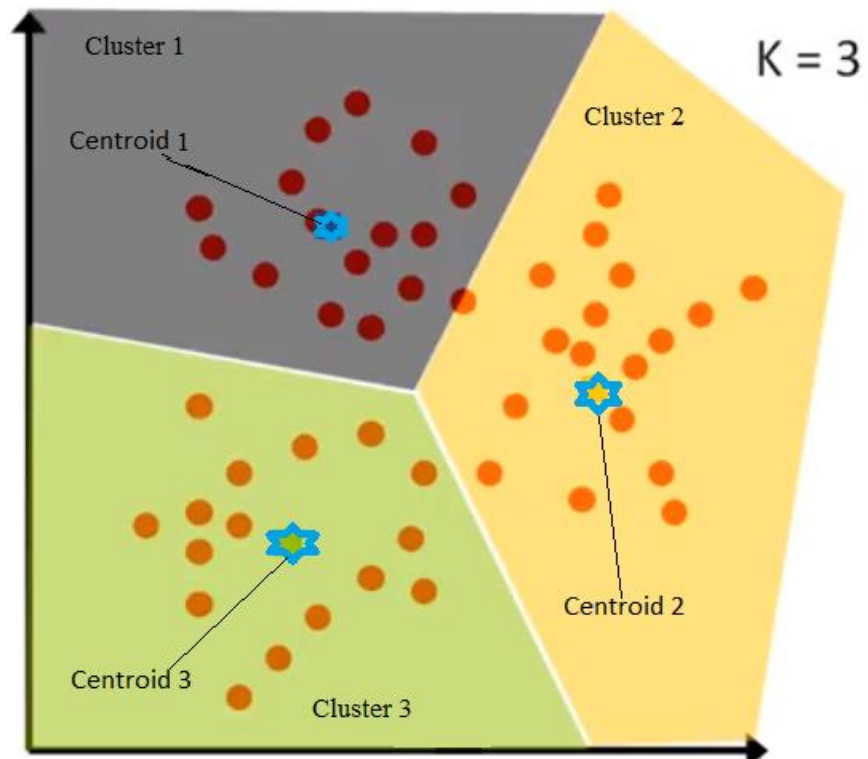
Allow unknown values in cate...

Experiment Properties

		Predicted →	
		P	N
Actual ↓	T	1	1
	F	2	1

		Predicted →		
		Low	Medium	High
Actual ↓	Low	<b>0 (0%)</b>	1 (100%)	
	Medium		<b>2 (66.6%)</b>	1(33.3%)
	High			<b>1 (100%)</b>

# Chapter 8



Country	RedMeat	WhiteMeat	Eggs	Milk	Fish	Cereals	Starch	Nuts	Fr&Veg
Albania	10.1	1.4	0.5	8.9	0.2	42.3	0.6	5.5	1.7
Austria	8.9	14	4.3	19.9	2.1	28	3.6	1.3	4.3
Belgium	13.5	9.3	4.1	17.5	4.5	26.6	5.7	2.1	4
Bulgaria	7.8	6	1.6	8.3	1.2	56.7	1.1	3.7	4.2
Czechoslovakia	9.7	11.4	2.8	12.5	2	34.3	5	1.1	4
Denmark	10.6	10.8	3.7	25	9.9	21.9	4.8	0.7	2.4
E Germany	8.4	11.6	3.7	11.1	5.4	24.6	6.5	0.8	3.6
Finland	9.5	4.9	2.7	33.7	5.8	26.3	5.1	1	1.4
France	18	9.9	3.3	19.5	5.7	28.1	4.8	2.4	6.5

Search experiment items

Draft saved at 8:06:48 PM

Reader

Data source: Web URL via HTTP

URL: <http://www.biz.uiowa.edu/faculty/jledolter/DataMining/protein.csv>

Data format: CSV

CSV or TSV has header row

Search experiment items

Draft saved at 10:46:26 PM

K-Means Clustering

Number of Centroids: 3

Metric: Euclidean

Initialization: K-Means++

Iterations: 100

Search experiment items

Draft saved at 10:51:03 PM

Saved Datasets  
 Data Format Conversions  
 Data Input and Output  
 Data Transformation  
 Feature Selection  
 Machine Learning  
 Evaluate  
 Initialize Model  
 Score  
 Train  
 Sweep Parameters  
 Train Anomaly Detecto...  
 Train Clustering Model  
 Train Matchbox Recom...  
 Train Model

Reader

K-Means Clustering

Train Clustering Model

Train Clustering Model

Column Set

Selected columns:  
 All columns  
 Exclude column names: Country

Launch column selector

Check for Append or Uncheck for Result O...

Saved Datasets  
 Data Format Conversions  
 Data Input and Output  
 Data Transformation  
 Feature Selection  
 Machine Learning  
 Evaluate  
 Initialize Model  
 Score  
 Apply Transformation  
 Assign to Clusters  
 Score Matchbox Recom...  
 Score Model  
 Train  
 OpenCV Library Modules

Reader

K-Means Clustering

Train Clustering Model

Assign to Clusters

Column Set

Selected columns:  
 All columns  
 Exclude column names: Country

Launch column selector

Check for Append or Uncheck for Re...

	Country	RedMeat	WhiteMeat	Eggs	Milk	Fish	Cereals	Starch	Nuts	Fr&Veg	Assignments
view as											
	Albania	10.1	1.4	0.5	8.9	0.2	42.3	0.6	5.5	1.7	0
	Austria	8.9	14	4.3	19.9	2.1	28	3.6	1.3	4.3	1
	Belgium	13.5	9.3	4.1	17.5	4.5	26.6	5.7	2.1	4	1
	Bulgaria	7.8	6	1.6	8.3	1.2	56.7	1.1	3.7	4.2	0
	Czechoslovakia	9.7	11.4	2.8	12.5	2	34.3	5	1.1	4	2
	Denmark	10.6	10.8	3.7	25	9.9	21.9	4.8	0.7	2.4	1
	E Germany	8.4	11.6	3.7	11.1	5.4	24.6	6.5	0.8	3.6	2
	Finland	9.5	4.9	2.7	33.7	5.8	26.3	5.1	1	1.4	1
	France	18	9.9	3.3	19.5	5.7	28.1	4.8	2.4	6.5	1
	Greece	10.2	3	2.8	17.6	5.9	41.7	2.2	7.8	6.5	0
	Hungary	5.3	12.4	2.9	9.7	0.3	40.1	4	5.4	4.2	2

**Project Columns**

Select columns

**Selected columns:**  
All columns

Launch column selector

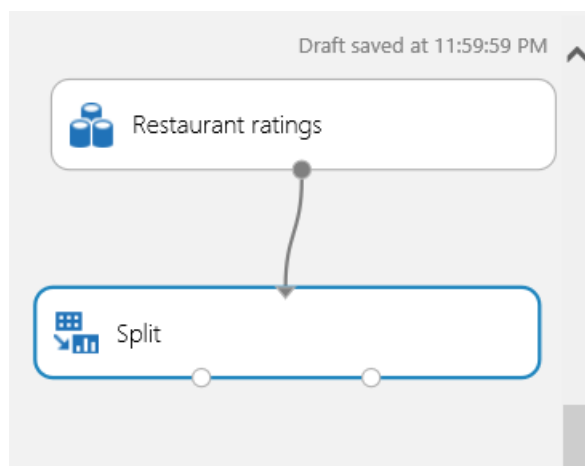
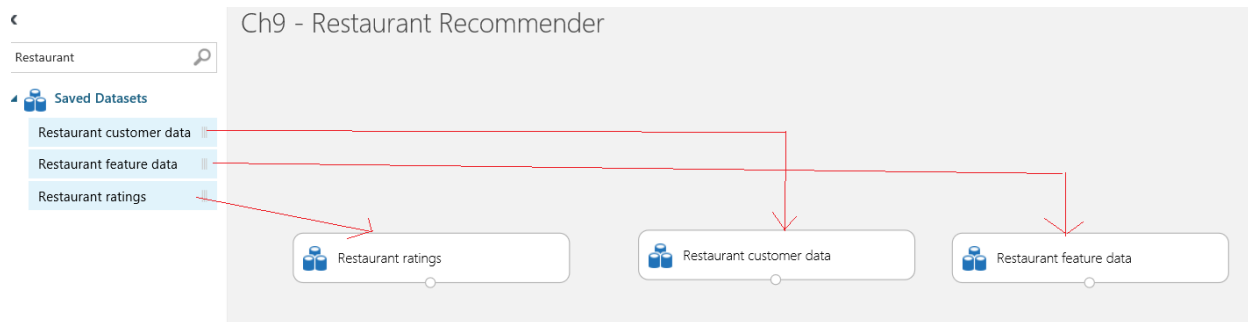
START TIME	5/2/2015 1...
END TIME	5/2/2015 1...
ELAPSED TIME	0:00:02.913
STATUS CODE	Finished
STATUS DETAILS	None

[View output log](#)

Quick Help

DISCARD CHANGES   REFRESH   CANCEL   **RUN**   PREPARE WEB SERVICE   PUBLISH TO GALLERY   CREATE SCORING EXPERIMENT

## Chapter 9



#### Split

Splitting mode

Recommender Split

Fraction of training-only users

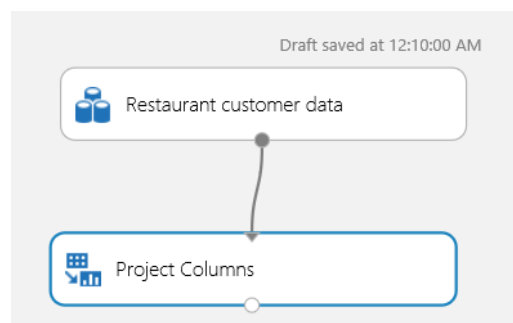
0.5

Fraction of test user ratings for train...

0.25

Fraction of cold users

0



#### Project Columns

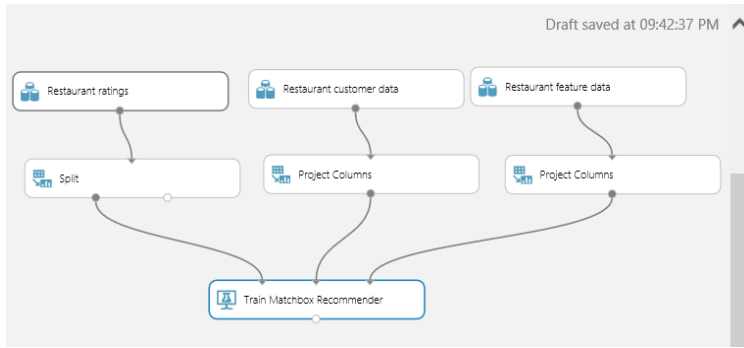
Select columns

Selected columns:

Column names: userID,latitude,longitude,interest,personality

Launch column selector

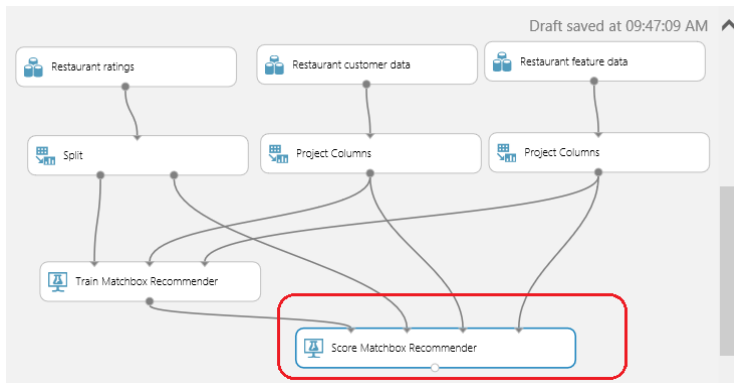
#### Experiment Properties



#### Train Matchbox Recommender

Number of traits

Number of recommendation algorithm iterations





#### Score Matchbox Recommender

Recommender prediction kind

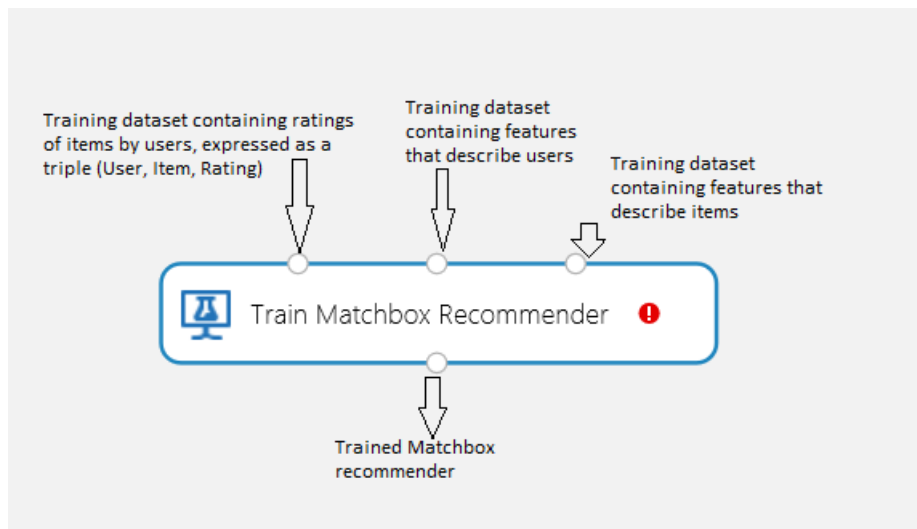
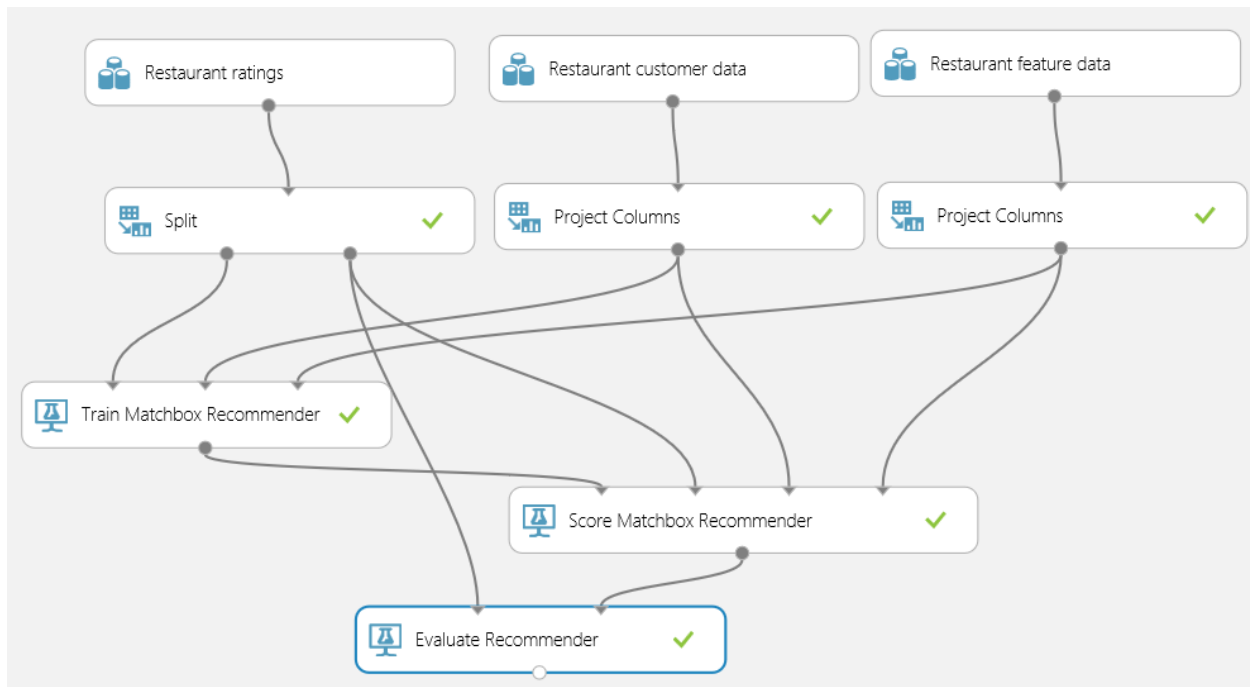
Recommended item selection

Maximum number of items to recommend to a user

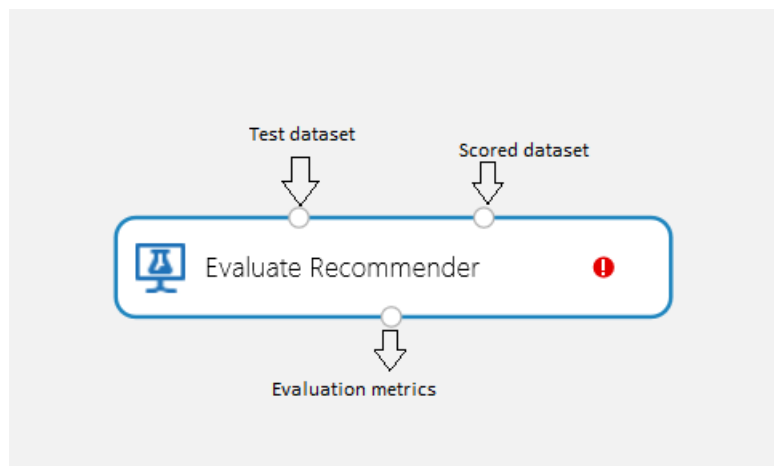
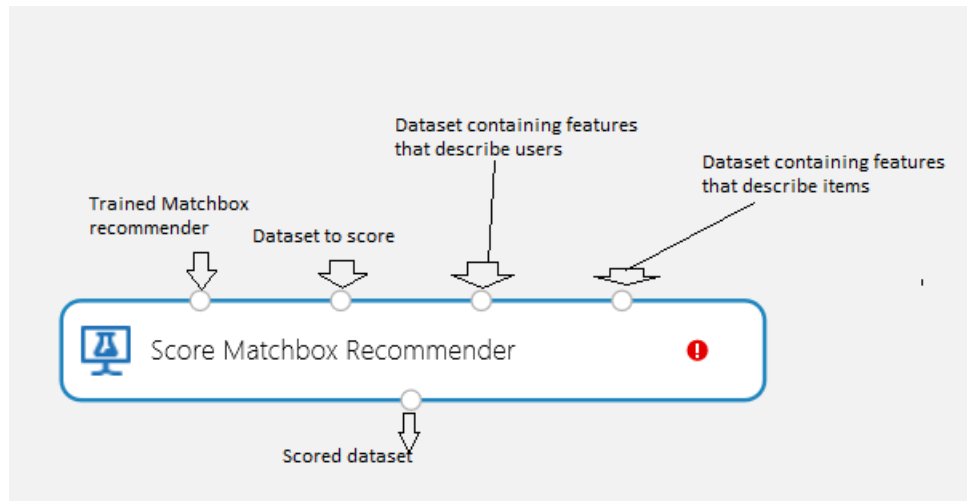
Minimum size of the recommendation pool for a single...

view as  

User	Item 1	Item 2	Item 3
U1048	135034	135026	135065
U1117	135018	132766	135088
U1049	135052	132862	135051
U1088	135057	135071	135032
U1062	135052	135045	135062
U1035	134986	135018	132773
U1125	135062	135076	135038
U1013	135075	135079	132921

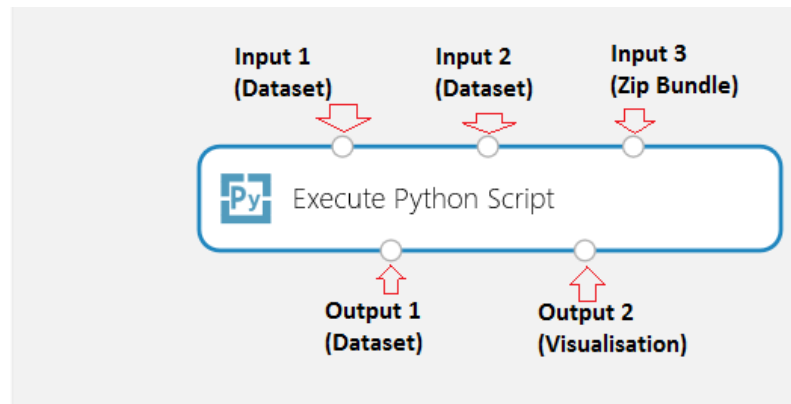






## Chapter 10

- Feature Selection
- Machine Learning
- OpenCV Library Modules
- **Python Language Modules**
  - Execute Python Script
- R Language Modules



Python script

```

1 # The script MUST contain a function named azureml_main
2 # which is the entry point for this module.
3 #
4 # The entry point function can contain up to two input arguments:
5 # Param<dataframe1>: a pandas.DataFrame
6 # Param<dataframe2>: a pandas.DataFrame
7
8 def azureml_main(dataframe1 = None, dataframe2 = None):
9
10 ##### Execution logic goes here #####
11 # Code to get the result dataFrame
12     resultDataFrame = ...
13
14 # Code to generate the visualisation (if any!)
15     ...
16 # Return value must be of a sequence of pandas.DataFrame
17     return resultDataFrame,

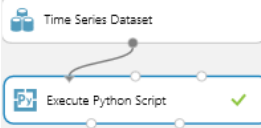
```

Annotations in the code block: A red arrow points to 'dataframe1' in line 8, labeled 'Input 1'. Another red arrow points to 'dataframe2' in line 8, labeled 'Input 2'. A third red arrow points to 'resultDataFrame' in line 17, labeled 'Output 1'.

In draft

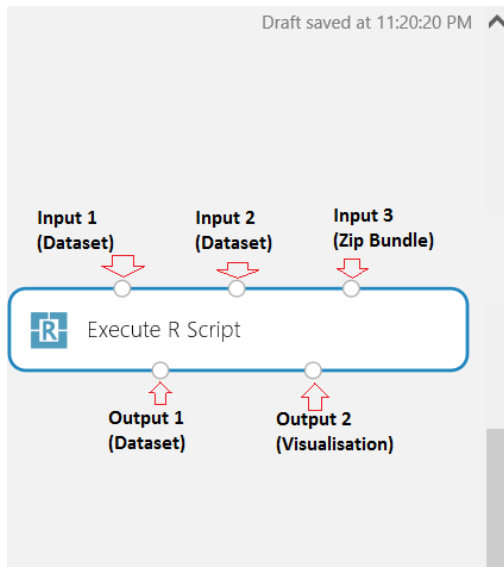
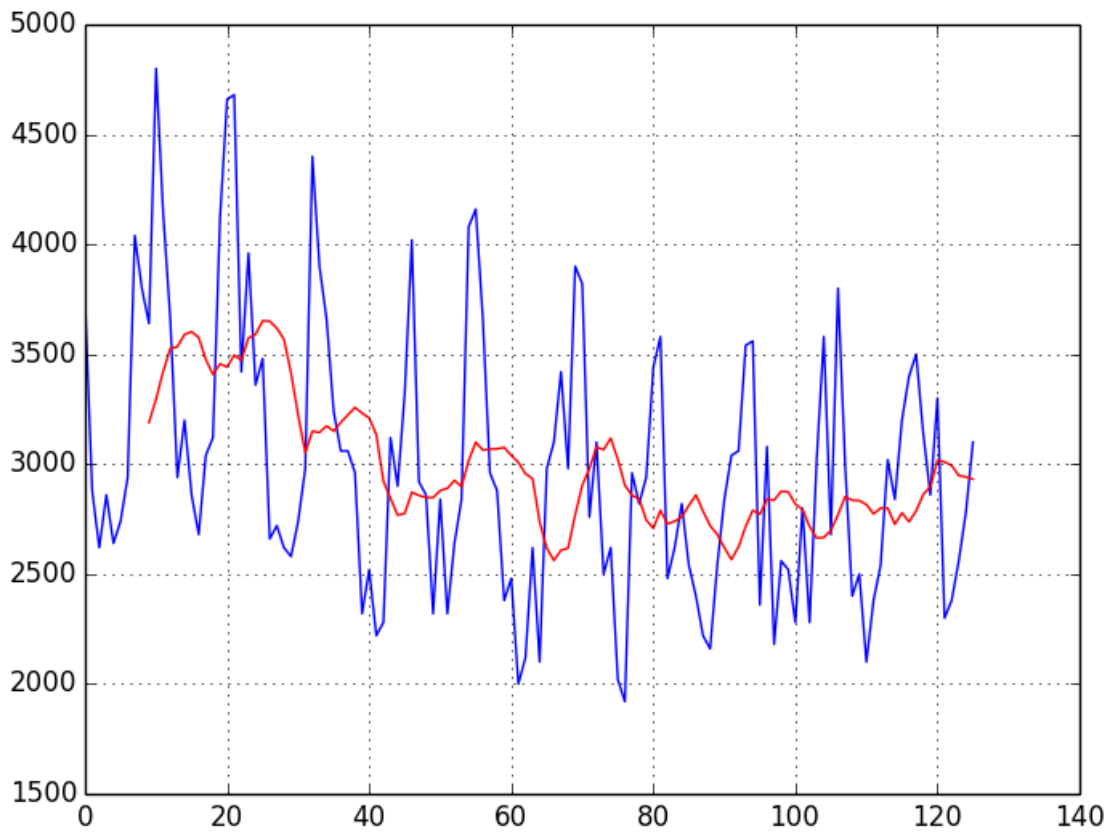
## Properties

Draft saved at 07:31:53 PM



### Python script

```
1 def azureml_main(dataframe1 = None, dataframe2 = None):
2     ##### Import package(s) #####
3     import pandas as pd
4     import numpy as np
5     import matplotlib
6     # Change to use a non-interactive backend, Agg (for PNGs)
7     matplotlib.use("agg")
8     import matplotlib.pyplot as plt
9     # Assign the input data frame to a new one
10    sampleTimeSeries = dataframe1
11    # Calculate moving average. Results in a series
12    movingAverage = pd.rolling_mean(sampleTimeSeries['N1725'], 10)
13    # Add a new column for moving average data and assigne the value
14    sampleTimeSeries['MovingAverage'] = movingAverage
15
16    ##### Visualisation #####
17    # Create New figure
18    fig = plt.figure()
19    ax = fig.gca()
20    ##### Plot into specified axis
21    # Plot the time series
22    sampleTimeSeries['N1725'].plot(ax=ax)
23    # Add line for moving average
24    sampleTimeSeries['MovingAverage'].plot(ax=ax, color='red')
25    # Save the figure to image
26    fig.savefig('TimeSeriesWithMovingAverage.png')
27    # Return value(pandas.DataFrame) - time series data with moving average
28    return sampleTimeSeries,
```



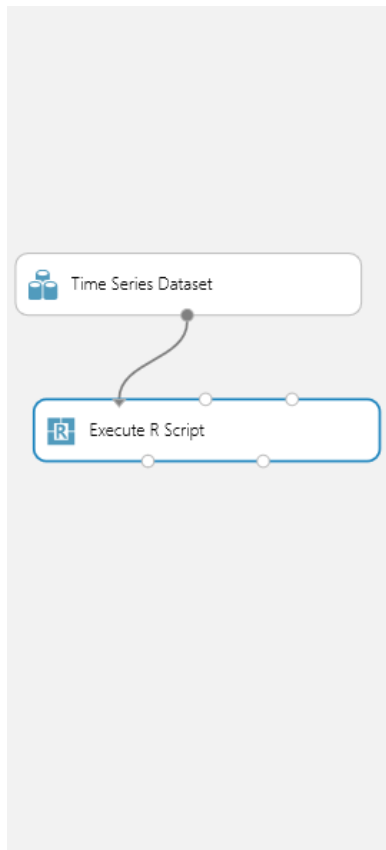
#### Execute R Script

R Script

```

1 # Map 1-based optional input ports to variables
2 dataset1 <- mam1.mapInputPort(1) # class: data.frame
3 dataset2 <- mam1.mapInputPort(2) # class: data.frame
4
5 # Contents of optional Zip port are in ./src/
6 # source("src/yourfile.R");
7 # load("src/yourData.rdata");
8
9 # Sample operation
10 data.set = rbind(dataset1, dataset2);
11
12 # You'll see this output in the R Device port.
13 # It'll have your stdout, stderr and PNG graphics device(s).
14 plot(data.set);
15
16 # Select data.frame to be sent to the output Dataset port
17 mam1.mapOutputPort("data.set");

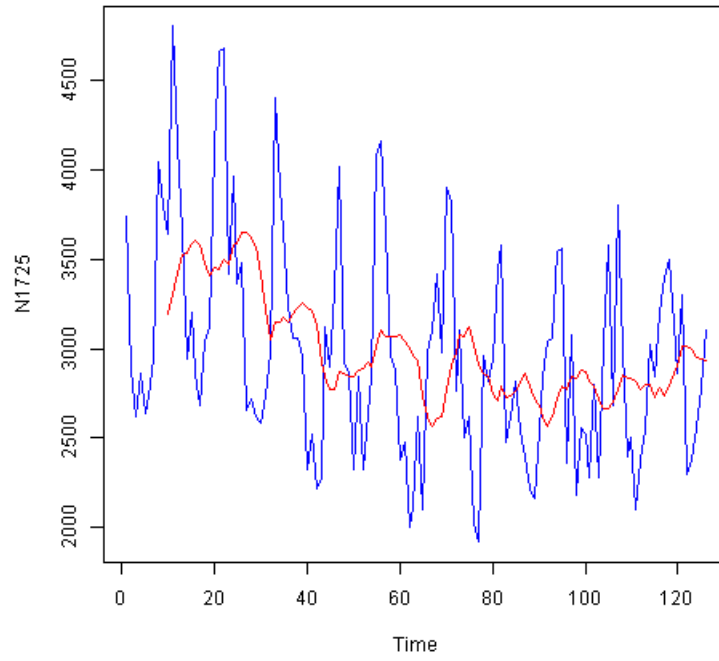
```



## R Script

```
1 # Map 1-based optional input ports to variables
2 dataset1 <- maml.mapInputPort(1) # class: data.frame
3
4 # Assign input dataset to new one
5 timeSeries <- dataset1
6
7 # Load the zoo package to the session
8 library(zoo)
9
10 # Calculate moving average with moving window 10
11 ma <- rollmean(timeSeries[2],10)
12 # Add a new column and assign moving average values to it
13 timeSeries$MovingAverage[10:126] <- ma
14
15 ##### Visualisation #####
16 # Plot time series
17 plot.ts(timeSeries['N1725'], col='blue')
18 # Add a line for moving average
19 lines(timeSeries['MovingAverage'], col='red')
20
21 # Select data.frame to be sent to the output Dataset port -
22 # time series data with moving average
23 maml.mapOutputPort("timeSeries");
```

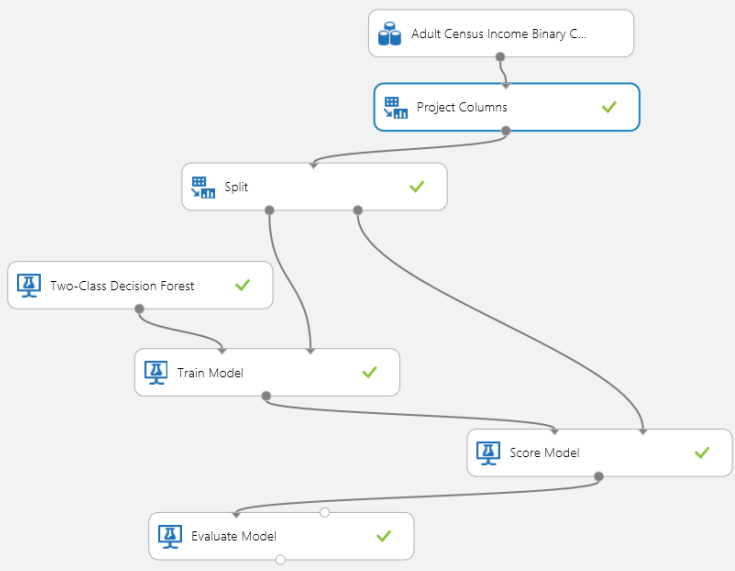
- ▶ Standard Output
- ▶ Standard Error
- ◀ Graphics Device



## Chapter 11

# Income Predictor

Finished running ✓



## Properties

### Project Columns

Select columns

#### Selected columns:

#### Column names:

age,education,sex,race,income

Launch column selector

START TIME 1/31/2015 2:13:10 PM

END TIME 1/31/2015 2:13:10 PM

ELAPSED TIME 0:00:00.000

STATUS CODE Finished

STATUS DETAILS Task output was present in output cache

### Experiment Properties

START TIME 1/31/2015 2:13:10 PM

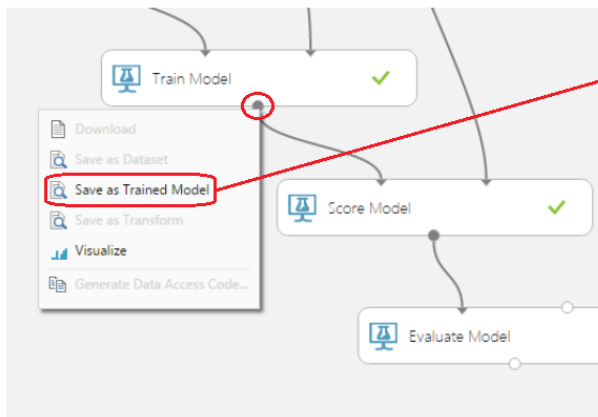
END TIME 1/31/2015 2:13:19 PM

STATUS CODE Finished

STATUS DETAILS None

### Project Columns

Create a projection of a dataset  
[\(more...\)](#)



## Save trained model

This is the new version of an existing trained model

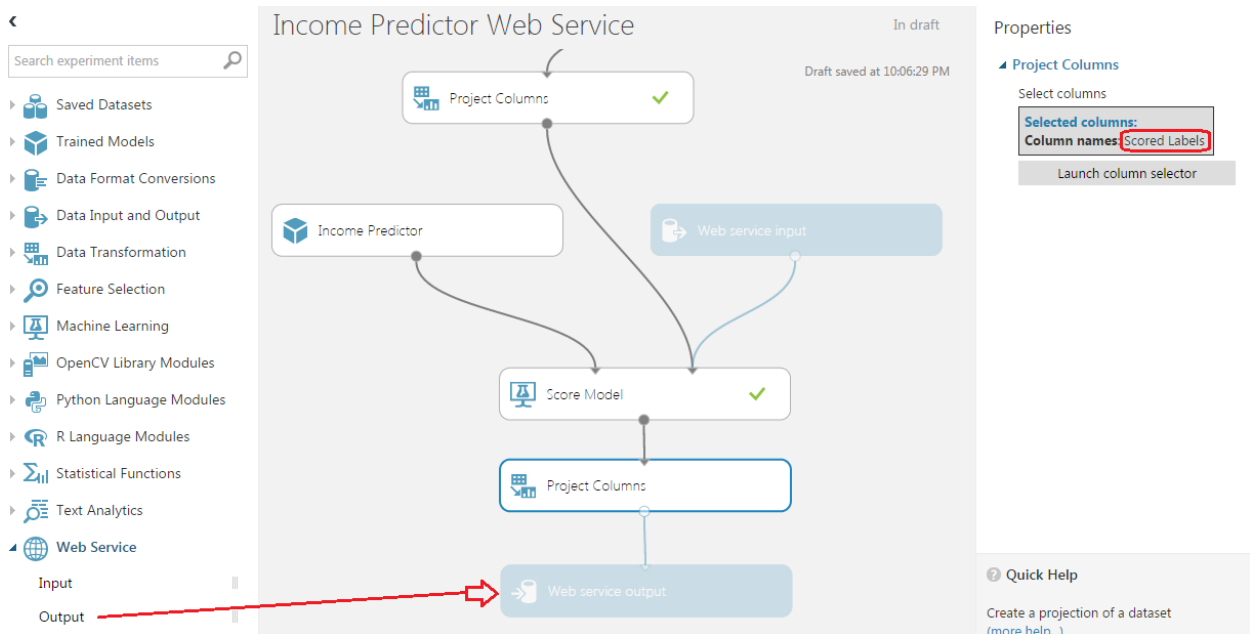
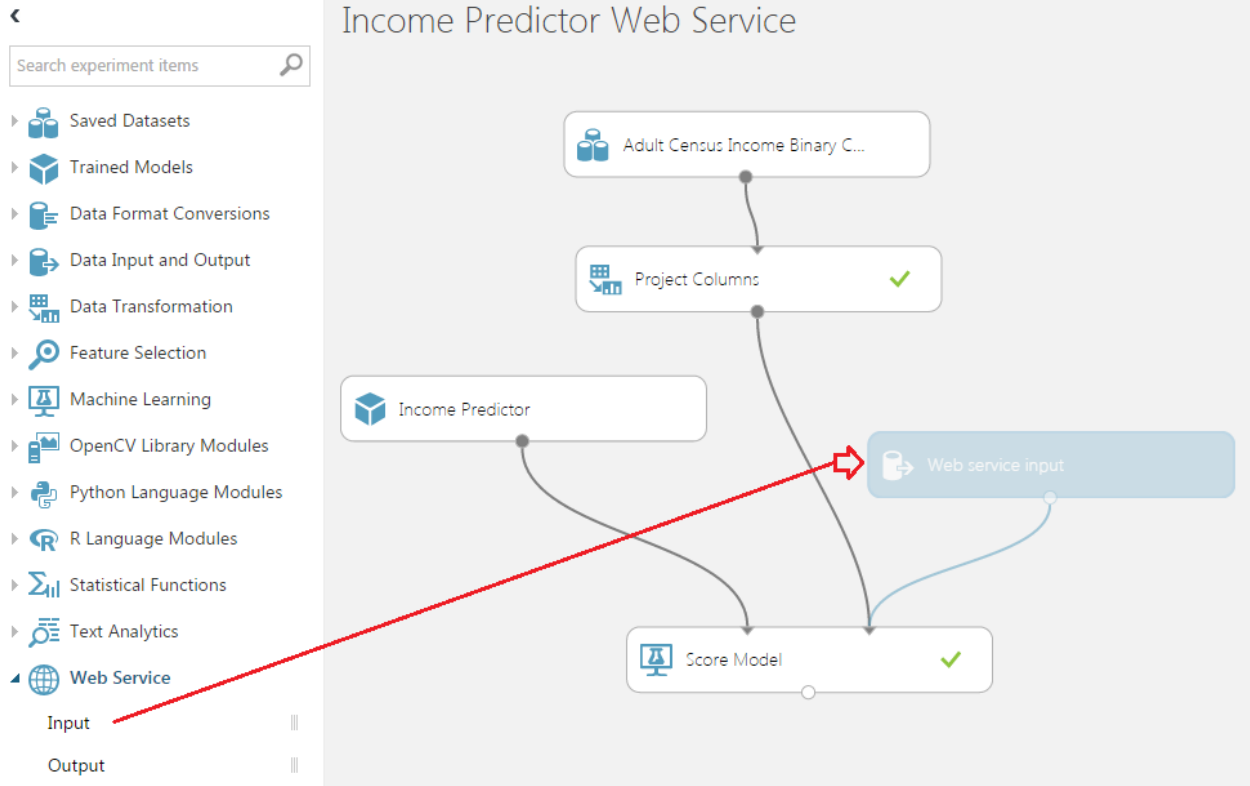
Enter a name for the new trained model:

Income Predictor

Provide an optional description:

Input: age, education, sex, race

✓



Would you like to publish the web service 'Income Predictor Web Service'?

YES  NO

NEW VIEW RUN HISTORY SAVE SAVE AS DISCARD CHANGES UNDO REDO CANCEL RUN PUBLISH WEB SERVICE PUBLISH TO GALLERY CREATE CLONING PROJECT



# income predictor web service

DASHBOARD CONFIGURATION

General



Published experiment

[View snapshot](#)

[View latest](#)



Description

No description provided for this web service.

API key



oyeUahaQll5OBv/KXngd1SneS85OQ5EEVQMcoYDC6LoTzpu5eOIMIOY6N8QFR+Q8HPqJX/M+z4k+zPHYPk

Default Endpoint

API HELP PAGE	TEST	APPS
<a href="#">REQUEST/RESPONSE</a>	<a href="#">Test</a>	<a href="#">Download Excel Workbook</a>
<a href="#">BATCH EXECUTION</a>		



Additional endpoints

Number of additional endpoints created for this web service: 0

[Manage endpoints in Azure management portal](#)



Test Income Predictor Web Service Service

## Enter data to predict

AGE

EDUCATION


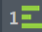
RACE

SEX



✓ 'Income Predictor Web Service' test returned [ ">50K" ]...

DETAILS ⓘ CLOSE ✕

+ NEW  DELETE 

← 'Income Predictor Web Service' test returned [ ">50K" ]... CLOSE ✕

✓ Result: {"Results":{"output1":{"type":"table","value":{"ColumnNames":["Scored Labels"],"ColumnTypes":["String"],"Values":{"[">50K"]}}}}}

- EXPERIMENTS
- WEB SERVICES**
- DATASETS
- TRAINED MODELS

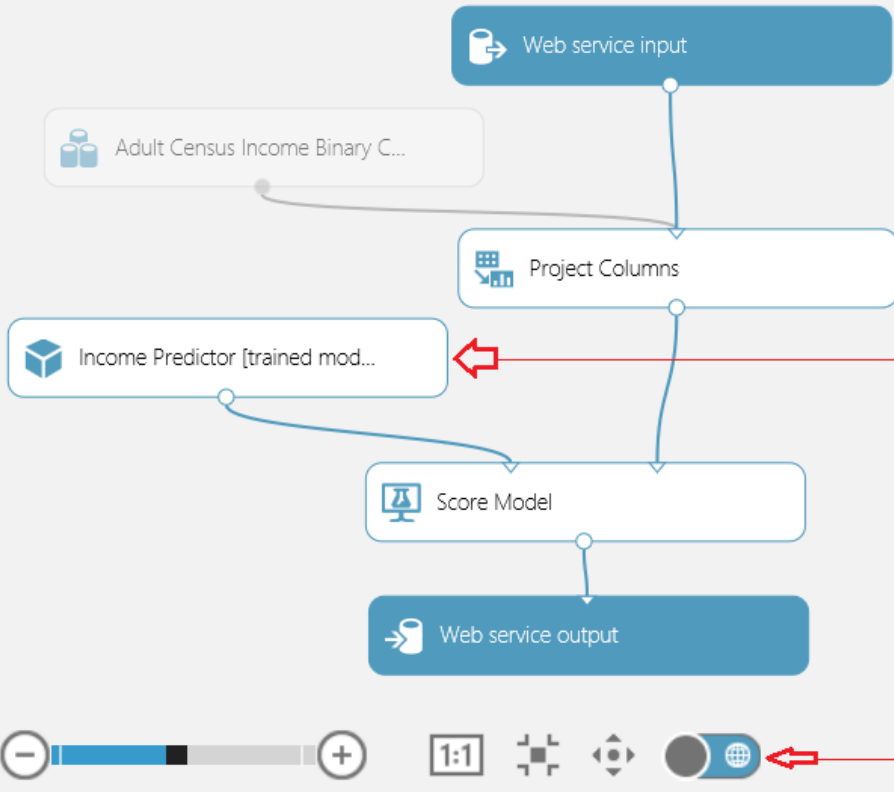
### web services

NAME	CREATED ON
Income Predictor Web Service	5/8/2015 10:17:24 PM



Training experiment    Scoring experiment

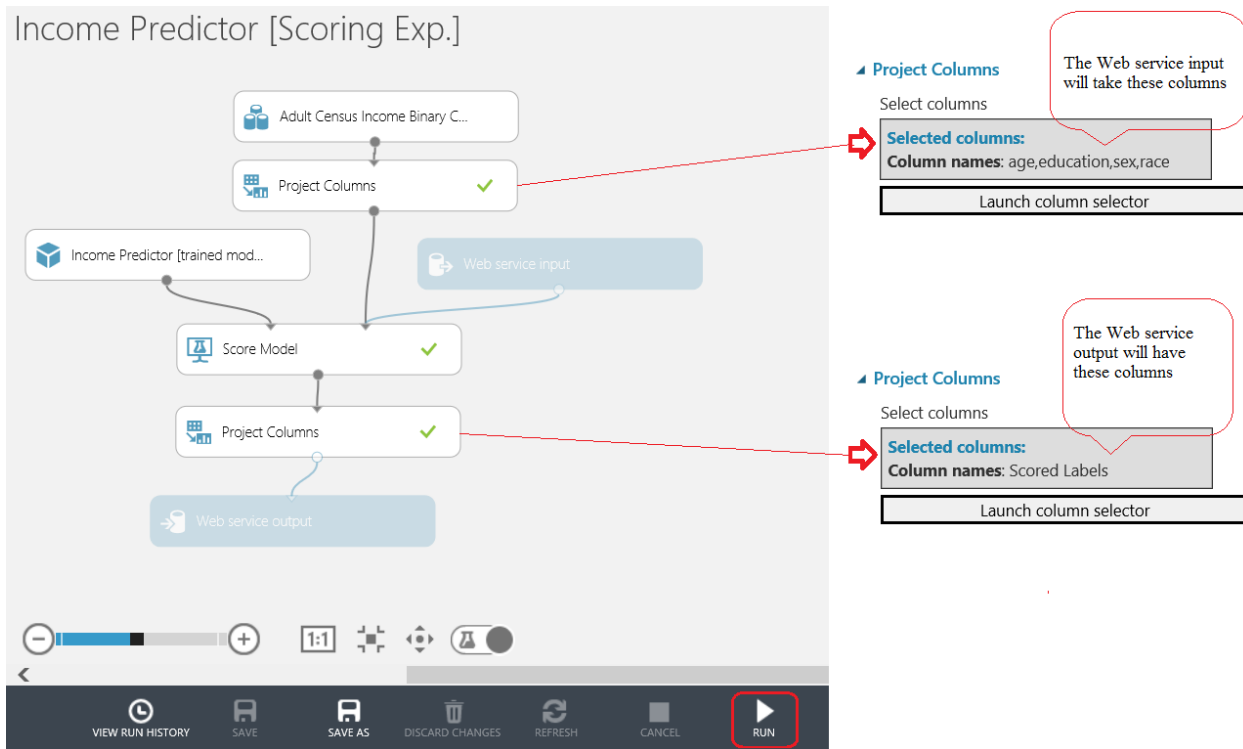
# Income Predictor [Scoring Exp.]



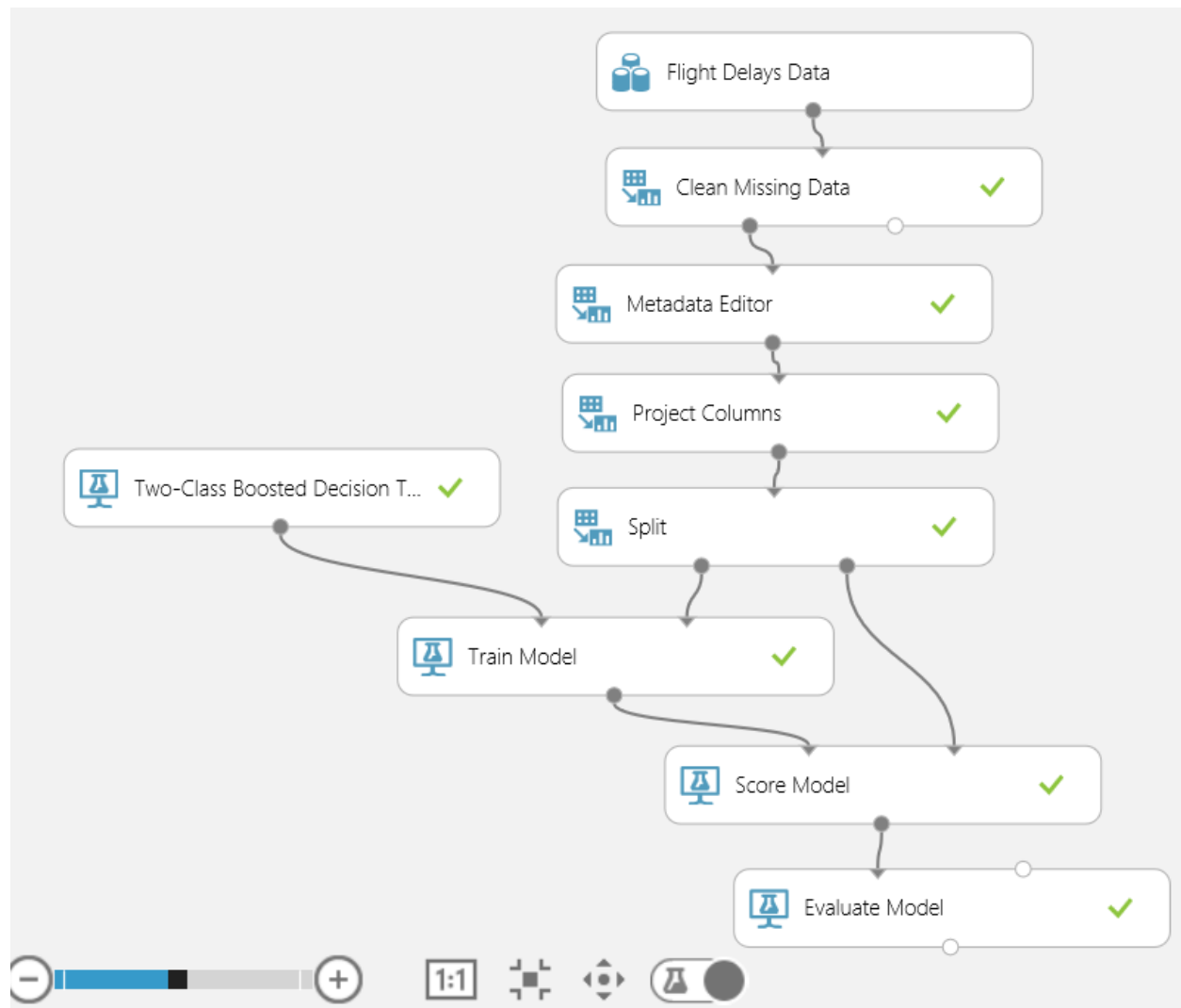
Link to corresponding training experiment

Trained Model Module

Switch between experiment and web service view



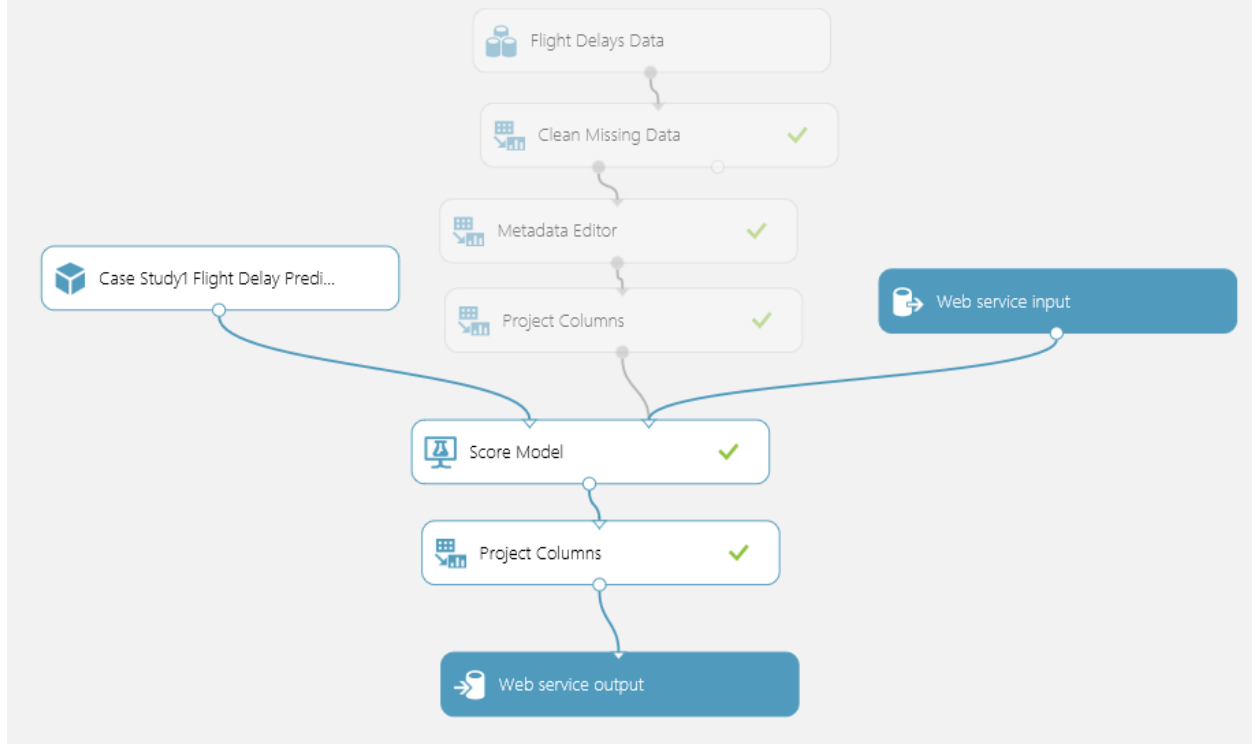
## Chapter 12



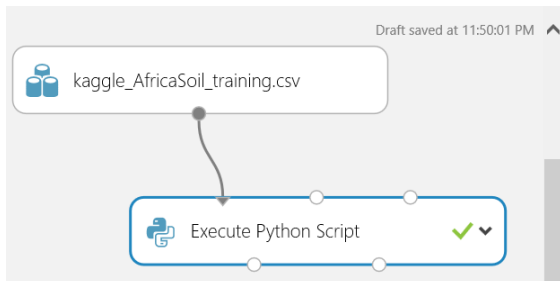
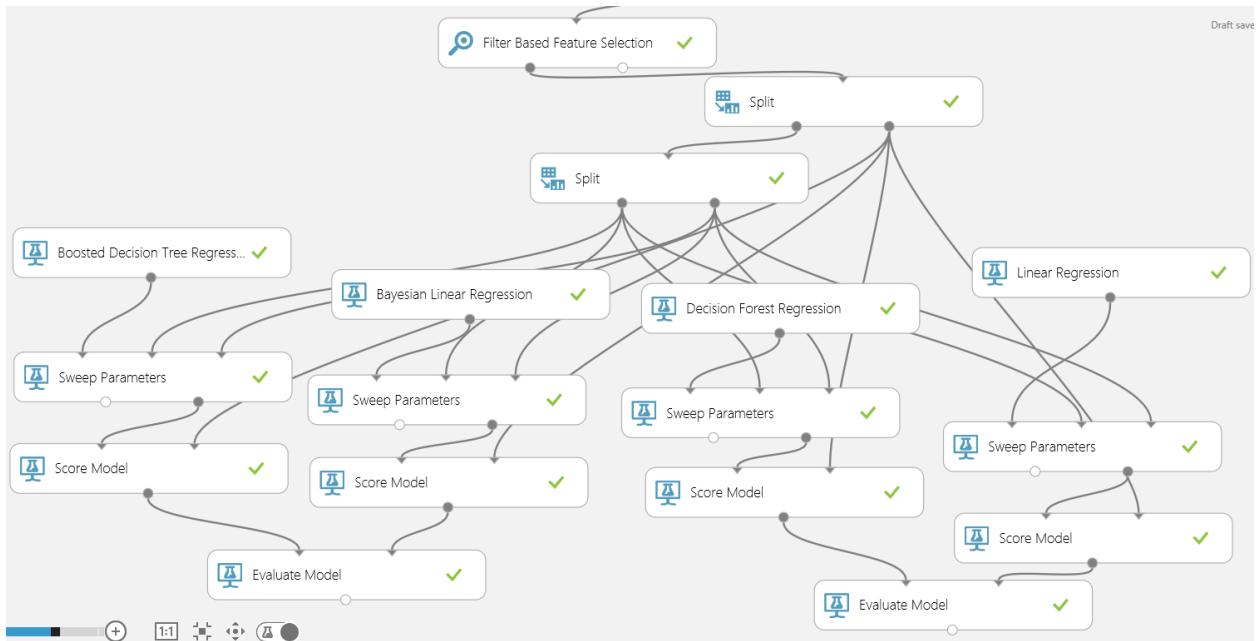
Training experiment

Scoring experiment

## Case Study1: Flight Delay Prediction [Scoring Exp.]



## Chapter 13



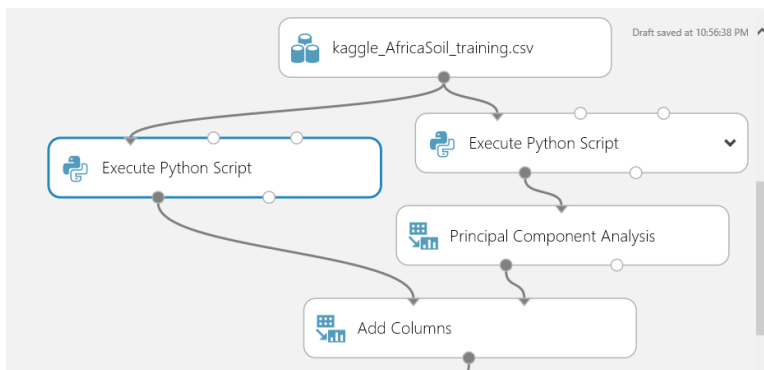
#### Execute Python Script

Python script

```

1 def azureml_main(dataframe1 = None, dataframe2 = None):
2     #Get all the columns
3     cols = dataframe1.columns.tolist()
4     #Select columns with name starting with letter 'm'
5     dataframe1=[col for col in cols if col.startswith('m')]
6     #Return the modified dataset
7     return dataframe1
8

```



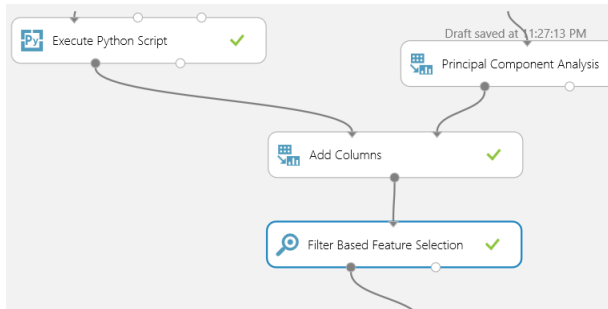
#### Execute Python Script

Python script

```

1 def azureml_main(dataframe1 = None, dataframe2 = None):
2     #Get all the columns
3     cols = dataframe1.columns.tolist()
4     #Select columns with name starting with letter 'm'
5     mCols=dataframe1[[col for col in cols if col.startswith('m')]]
6     #List of other columns to be excluded
7     exCols=['PIDN', 'Ca', 'pH', 'SOC', 'Sand']
8     #Drop all the columns with name starting with letter 'm'
9     dataframe1.drop(mCols, axis=1, inplace=True)
10    #Drop other columns - PIDN: 'Ca', 'pH', 'SOC', 'Sand'
11    dataframe1.drop(exCols,axis=1, inplace=True)
12    #Return the modified dataset
13    return dataframe1

```



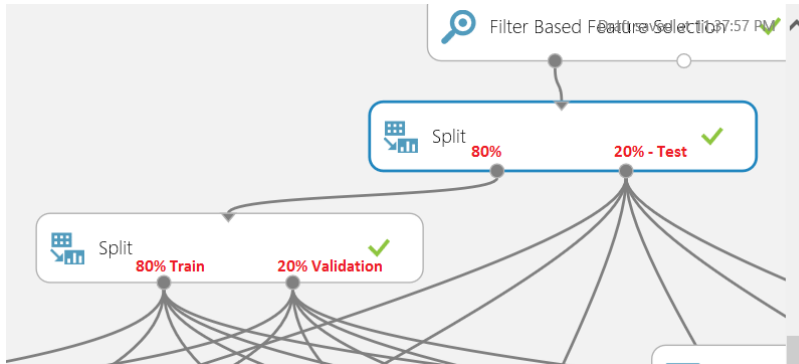
#### Filter Based Feature Selection

Feature scoring method

Operate on feature columns only

Target column  
 Selected columns:  
 Column names: P

Number of desired features



#### Split

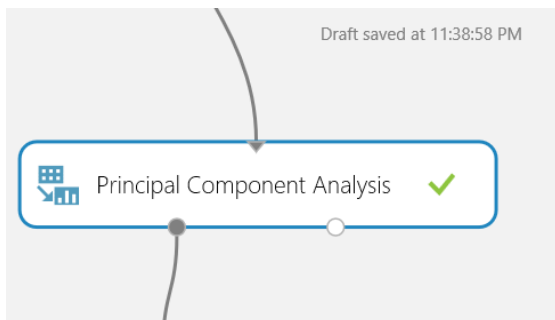
Splitting mode

Fraction of rows in the first output dataset

Randomized split

Random seed

Stratified split



#### Principal Component Analysis

Selected columns  
 Selected columns:  
 All columns

Number of dimensions to reduce to

Normalize dense dataset to zero mean