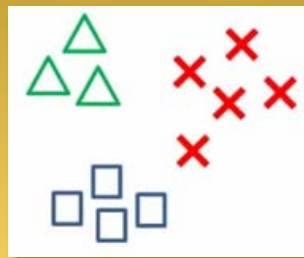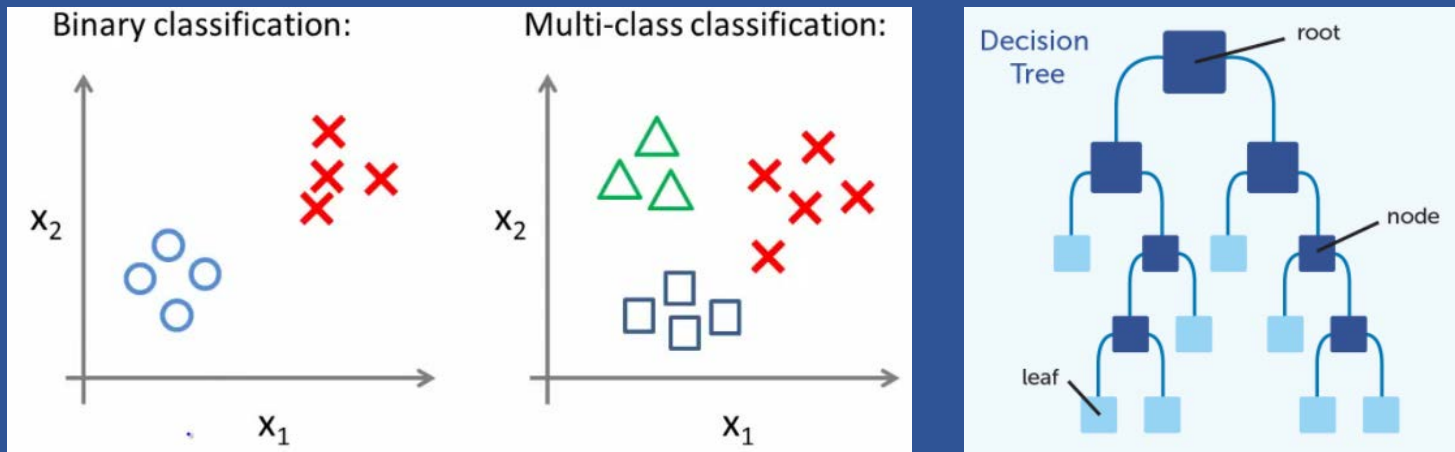# ALGORITHM MULTI CLASS

# Algorithm Multi Class

## In this session

- Multi Class Algorithms in Azure ML
- Data importing and engineering
- Feature engineering
- Modeling and evaluation
- Reuter Data set
- Edit Metadata
- Confusion Matrix

# Algorithm Multi Class

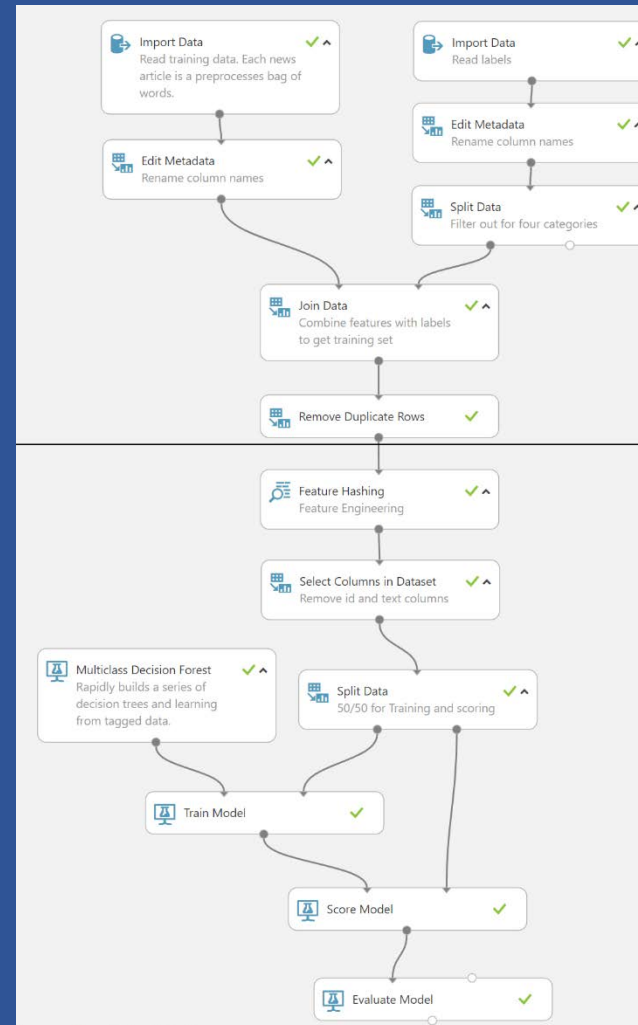## Multi Class Algorithms in Azure ML



## Multiclass Decision Forest

- Based on the decision forest algorithm
- Rapidly builds a series of decision trees
- learning from tagged data.
- Voting on the most popular output class
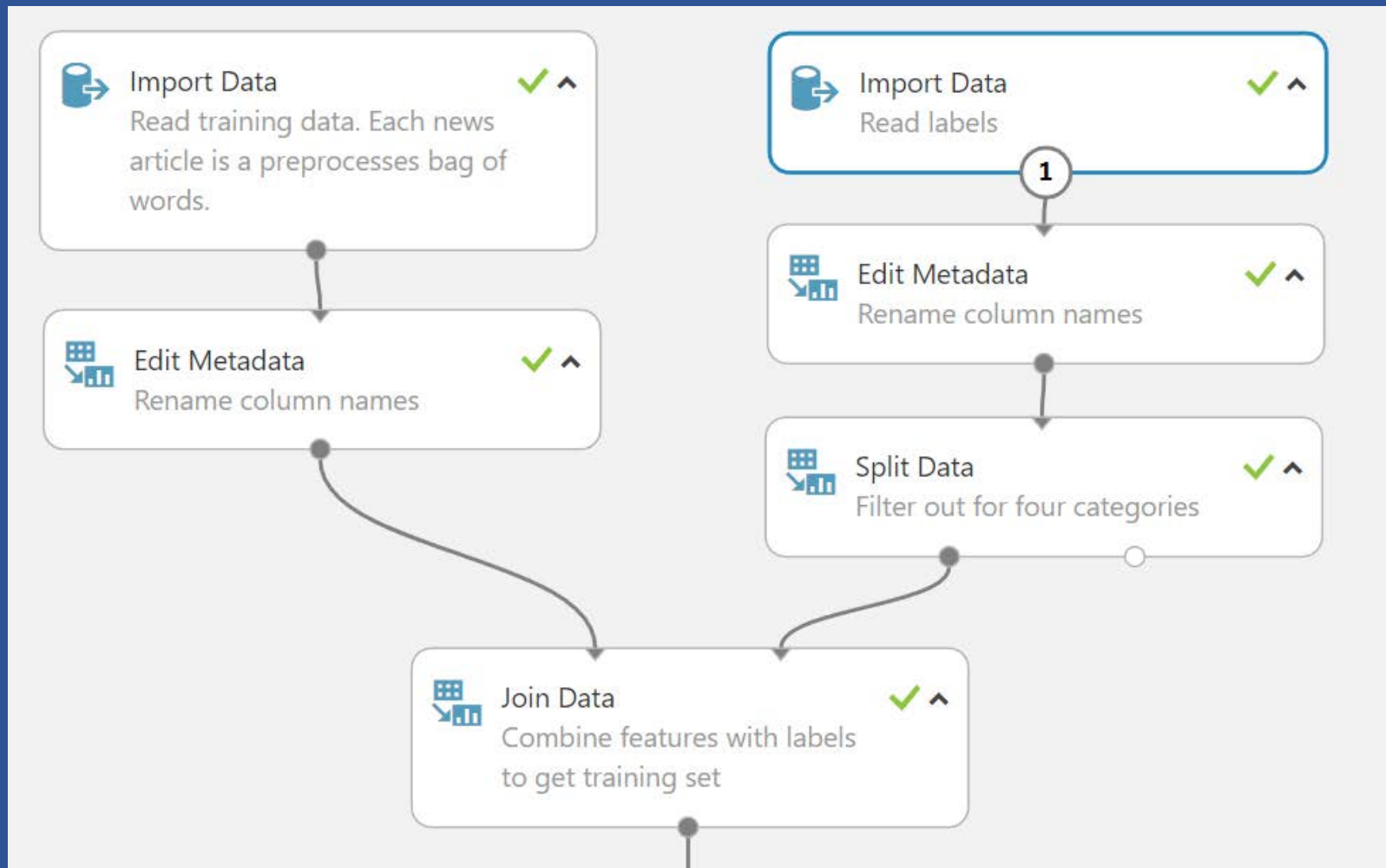- Voting is a form of aggregation

# Algorithm Multi Class
## Over all Experiment

- multiclass classifiers
- Feature engineering using hashing
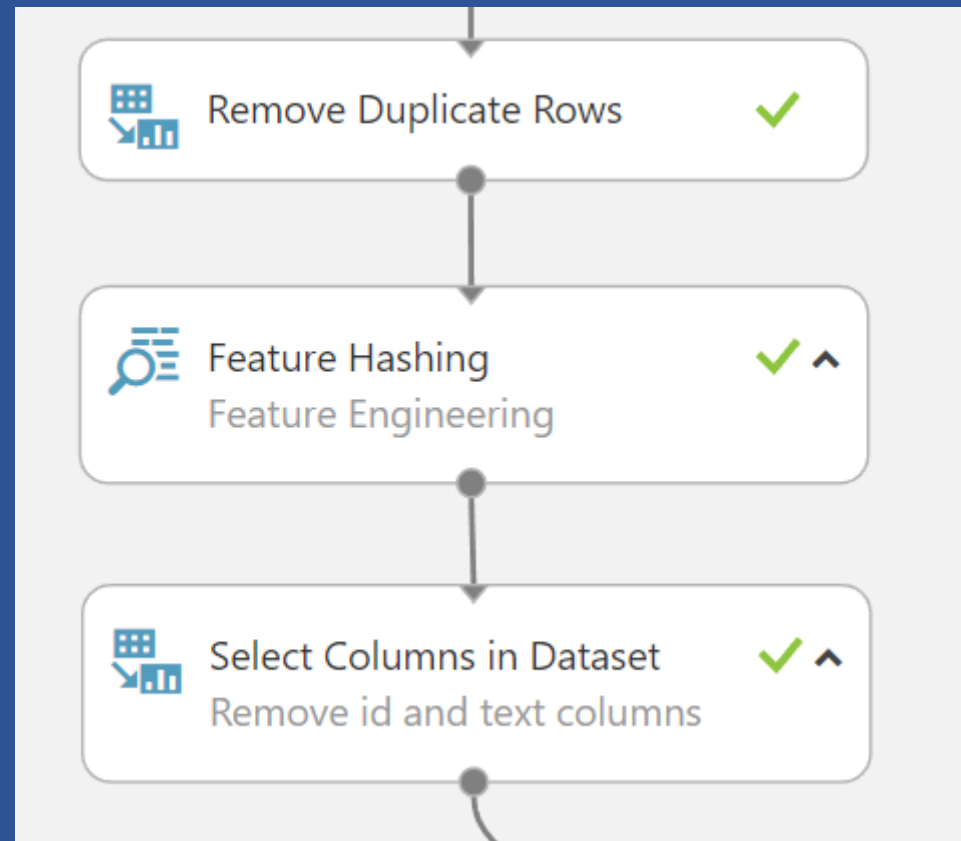- Classify news into four categories

# Algorithm Multi Class
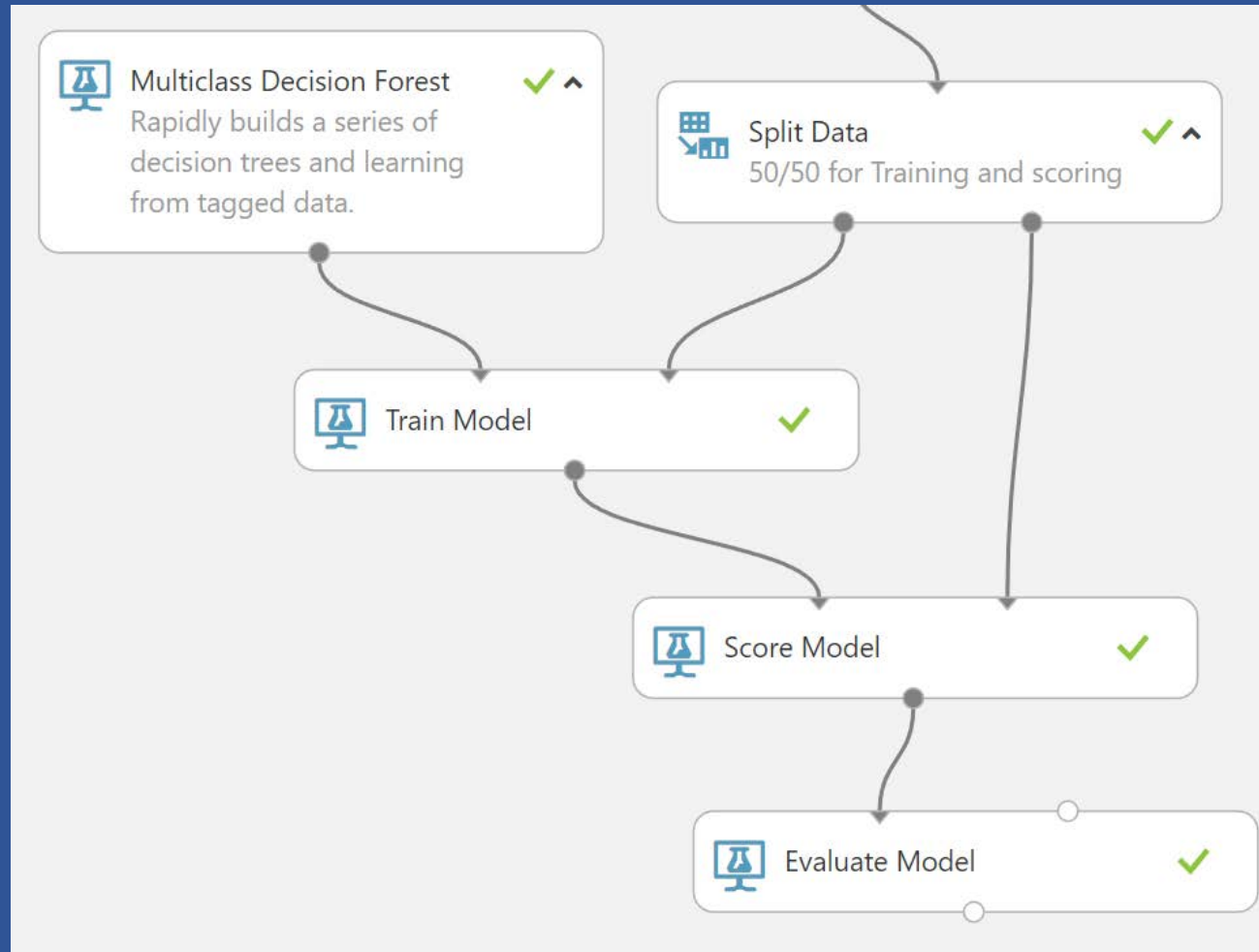
## Data importing and engineering

**GreatFriends.Biz**

# Algorithm Multi Class
## Feature engineering

# Algorithm Multi Class
## Modeling and evaluation

# Algorithm Multi Class
## Reuter Data set

- 2004 Reuters news dataset
- 10,000 News examples
- 5K Training / 5K Scoring

Data set has 103 categories that are organized into four hierarchies:

- Corporate-Industrial (CCAT)
- Government and Social (GCAT)
- Economics and Economic Indicators (ECAT)
- Securities and Commodities Trading and Market (MCAT)

# Algorithm Multi Class
## Import data set



A = https://raw.githubusercontent.com/laploy/ML/master/mul.csv

B = https://raw.githubusercontent.com/laploy/ML/master/mul_token.csv

# Algorithm Multi Class

## Edit Metadata

**GreatFriends.Biz**

# Algorithm Multi Class
## Splitting Data

Used only the rows already tagged with hierarchy names (CCAT,ECAT,GCAT,MCAT)

**Before splitting**          **After splitting**

# Algorithm Multi Class
## Feature & Clean



**Join Data**
Combine features with labels
to get training set

**Remove Duplicate Rows**
1

**Feature Hashing**

**Select Columns in Dataset**
Remove id and text columns

---

▲ **Join Data**

Join key columns for L

**Selected columns:**
**Column names**: id

Launch column selector

Join key columns for R

**Selected columns:**
**Column names**: id

Launch column selector

☑ Match case

Join type

Inner Join ▼

☐ Keep right key colu...

---

▲ **Remove Duplicate Rows**

Key column selection filter exp...

**Selected columns:**
**Column names**: id

Launch column selector

☑ Retain first duplicate r...

# Algorithm Multi Class
## Feature Engineering

◢ **Feature Hashing**

Target column(s)

> **Selected columns:**
> **Column names:** article

Launch column selector

Hashing bitsize     ☰

8

N-grams     ☰

1

◢ **Select Columns in Dataset**

Select columns

> **Selected columns:**
> **All columns**
> **Exclude column names:**
> id,article

Launch column selector

# Algorithm Multi Class
## Algorithm

**◢ Multiclass Decision Forest**

Resampling method ☰

Bagging ▼

Create trainer mode

Single Parameter ▼

Number of decision trees ☰

8

Maximum depth of the ... ☰

32

Number of random split... ☰

128

Minimum number of sa... ☰

1

☑ Allow unknown val... ☰

**◢ Train Model**

Label column

Selected columns:
Column indices: 1

Launch column selector

# Algorithm Multi Class
## Confusion Matrix

Test data = https://raw.githubusercontent.com/laploy/ML/master/mul-test.txt

◢ Metrics

| | |
|---|---|
| Overall accuracy | 0.813474 |
| Average accuracy | 0.906737 |
| Micro-averaged precision | 0.813474 |
| Macro-averaged precision | 0.802249 |
| Micro-averaged recall | 0.813474 |
| Macro-averaged recall | 0.749342 |

Predicted Class

| Actual Class | CCAT | ECAT | GCAT | MCAT |
|---|---|---|---|---|
| CCAT | 91.0% | 1.9% | 3.2% | 4.0% |
| ECAT | 27.8% | 52.1% | 8.8% | 11.2% |
| GCAT | 16.9% | 2.4% | 79.0% | 1.7% |
| MCAT | 17.0% | 4.1% | 1.1% | 77.7% |

# Algorithm Multi Class

## More information

Multiclass Decision Forest

https://msdn.microsoft.com/en-us/library/azure/dn906015.aspx

This Experiment

https://gallery.cortanaintelligence.com/Experiment/Multi-Class