# CIS5560 Term Project Tutorial

**Authors: Tanvi Gawade, Swarnim Jambhule, Sushant Burde**

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

**Date: 05/19/2019**

# Lab Tutorial

Tanvi ([tgawade@calstatela.edu](mailto:tgawade@calstatela.edu))

Swarnim ([sjambhu@calstatela.edu](mailto:sjambhu@calstatela.edu))

Sushant ([sburde@calstatela.edu](mailto:sburde@calstatela.edu))

05/19/2019

# Classification for Floating or Non-Floating Items from the Library Inventory On

# Data Bricks in Spark Machine Leaning

## Objectives

**List what your objectives are.** In this hands-on lab, you will learn how to:

- Get data manually

- Create Spark cluster

- Writing PySpark codes to develop a predictive model.

- Classification of Floating and Non-Floating values using Random Forest Classification and Gradient Boosting Tree Classification.

- Visualization

## Platform Spec

- Data Bricks PySpark

- Databricks Runtime Version: 5.2(Incl. Apache Spark 2.4.0, Scala 2.11)

- Execution: Single Node

- Memory: 6GB Capacity

# Step 1: Creating a Cluster in Data Bricks
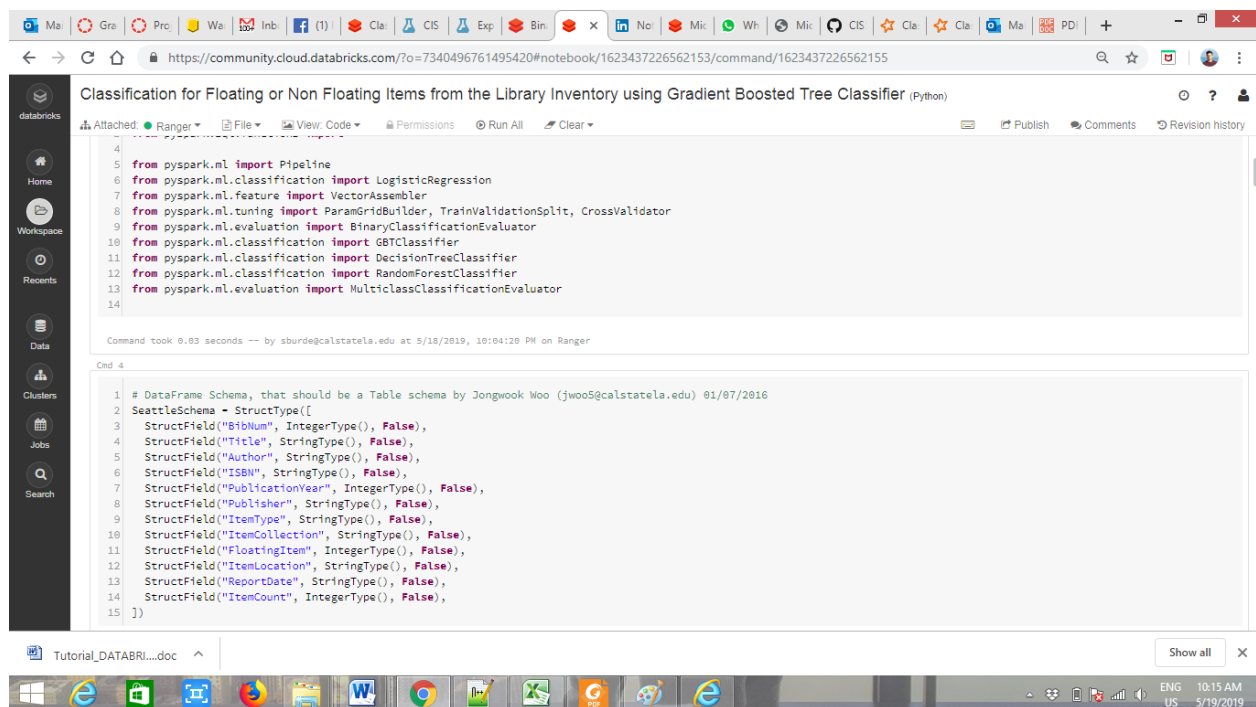
1. This step is to create a cluster for the execution of the codes.

Properties: -

- Python Version: 2

- Driver Type: Community Optimized

- Availability Zone: us-west-2c

# Step 2: Prepare the Data

# Step 3: Select data from the table



## Question 1 (5%)

Read your flights.csv file from its table at Databricks

```python
# TODO: use LibrarySchema in order to adopt shcmea to read csv data set in the schema.
# Jongwook Woo (jwoo5@calstatela.edu) 01/07/2016

# Load the source data
csv = spark.sql("SELECT * FROM seattle_library_optimized_file_cd44f_csv")
```

csv: pyspark.sql.dataframe.DataFrame = [BibNum: integer, Title: string ... 11 more fields]
Command took 0.15 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:04:28 PM on Ranger

```python
# Select features and label
data = csv.select("BibNum", "PublicationYear", "ItemCount", col("FloatingItem").alias("label"))

# Split the data
splits = data.randomSplit([0.7, 0.3])
train = splits[0]
test = splits[1].withColumnRenamed("label", "trueLabel")
```

data: pyspark.sql.dataframe.DataFrame = [BibNum: integer, PublicationYear: integer ... 2 more fields]
train: pyspark.sql.dataframe.DataFrame = [BibNum: integer, PublicationYear: integer ... 2 more fields]
test: pyspark.sql.dataframe.DataFrame = [BibNum: integer, PublicationYear: integer ... 2 more fields]
Command took 0.06 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:05:51 PM on Ranger

# Step 4: Define for Train and Test the Data. (0.7 – Train, 0.3 – Test)

# Step 5: Define Pipeline for the data



# Step 5: Tune the parameter

**Tune Parameters**

You can tune parameters to find the best model for your data. A simple way to do this is to use **TrainValidationSplit** to evaluate each combination of parameters defined in a **ParameterGrid** against a subset of the training data in order to find the best performing parameters.

**Regularization**

is a way of avoiding Imbalances in the way that the data is trained against the training data so that the model ends up being over fit to the training data. In other words It works really well with the training data but it doesn't generalize well with other data. That we can use a **regularization parameter** to vary the way that the model balances that way.

**Training ratio of 0.7**

it's going to use 70% of the the data that it's got in its training set to train the model and then the remaining 30% is going to use to validate the trained model.

In **ParamGridBuilder**, all possible combinations are generated from regParam, maxIter, threshold. So it is going to try each combination of the parameters with 70% of the the data to train the model and 30% to to validate it.

# Step 5a: Train Validation Split with Threshold parameters

# Step 5b: Train Validation Split with elastic-net parameters

The second combination of parameters with (regParam: [0.01, 0.5, 2.0]), (elasticNetParam: [0.0, 0.5, 1]), (maxIter: [1, 5])

# Step 5c: Cross Validator with elastic net parameters

The combination of parameters with (regParam: [0.01, 0.5, 2.0]), (elasticNetParam: [0.0, 0.5, 1]), (maxIter: [1, 5])

Classification for Floating or Non Floating Items from the Library Inventory using Gradient Boosted Tree Classifier (Python)

```
1  # TODO: params refered to the reference above
2  paramGridCV = (ParamGridBuilder() \
3              .addGrid(lr[2].regParam, [0.01, 0.5, 2.0]) \
4              .addGrid(lr[2].elasticNetParam, [0.0, 0.5, 1]) \
5              .addGrid(lr[2].maxIter, [1, 5]) \
6              .build())
```

Command took 0.02 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:08:51 PM on Ranger

Cmd 21

```
1  # TODO: K = 2 you may test it with 5, 10
2  # K=2, 3, 5,
3  # K= 10 takes too long
4  cv = CrossValidator(estimator=pipeline[2], evaluator=BinaryClassificationEvaluator(), \
5                      estimatorParamMaps=paramGridCV, numFolds=5)
6
7  # the third best model
8  model.insert(2, cv.fit(train))
```

▸ (46) Spark Jobs

Command took 4.08 minutes -- by sburde@calstatela.edu at 5/18/2019, 10:08:56 PM on Ranger

Cmd 22

## Test the Model

Now you're ready to apply the model to the test data.

Cmd 23

# Step 6: Test The model

# Step 7: Compute Confusion Matrix Metrics

Classifiers are typically evaluated by creating a confusion matrix, which indicates the number of:

- True Positives
- True Negatives
- False Positives
- False Negatives

From these core measures, other evaluation metrics such as precision and recall can be calculated.

Attached: ● Ranger ▾   📄 File ▾   🖥 View: Code ▾   🔒 Permissions   ⊙ Run All   ✐ Clear ▾   ⌨   📤 Publish   💬 Comments   ↺ Revision history

**Question 4 (5%): Complete Fill-in to calculate Recall**

Cmd 26

```python
# TODO: Complete the following [Fill-In] to calculate Recall
for i in range(3):
    tp = float(predicted[i].filter("prediction == 1.0 AND truelabel == 1").count())
    fp = float(predicted[i].filter("prediction == 1.0 AND truelabel == 0").count())
    tn = float(predicted[i].filter("prediction == 0.0 AND truelabel == 0").count())
    fn = float(predicted[i].filter("prediction == 0.0 AND truelabel == 1").count())
    metrics = spark.createDataFrame([
        ("TP", tp),
        ("FP", fp),
        ("TN", tn),
        ("FN", fn),
        ("Precision", tp / (tp + fp)),
        ("Recall", tp / (tp + fn))],["metric", "value"])
    metrics.show()
```

▸ (4) Spark Jobs

⊞ ZeroDivisionError: float division by zero

Command took 5.12 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:14:35 PM on Ranger

Cmd 27

## Review the Area Under ROC

Another way to assess the performance of a classification model is to measure the area under a ROC curve for the model. the spark.ml library includes a **BinaryClassificationEvaluator** class that you can use to compute this.
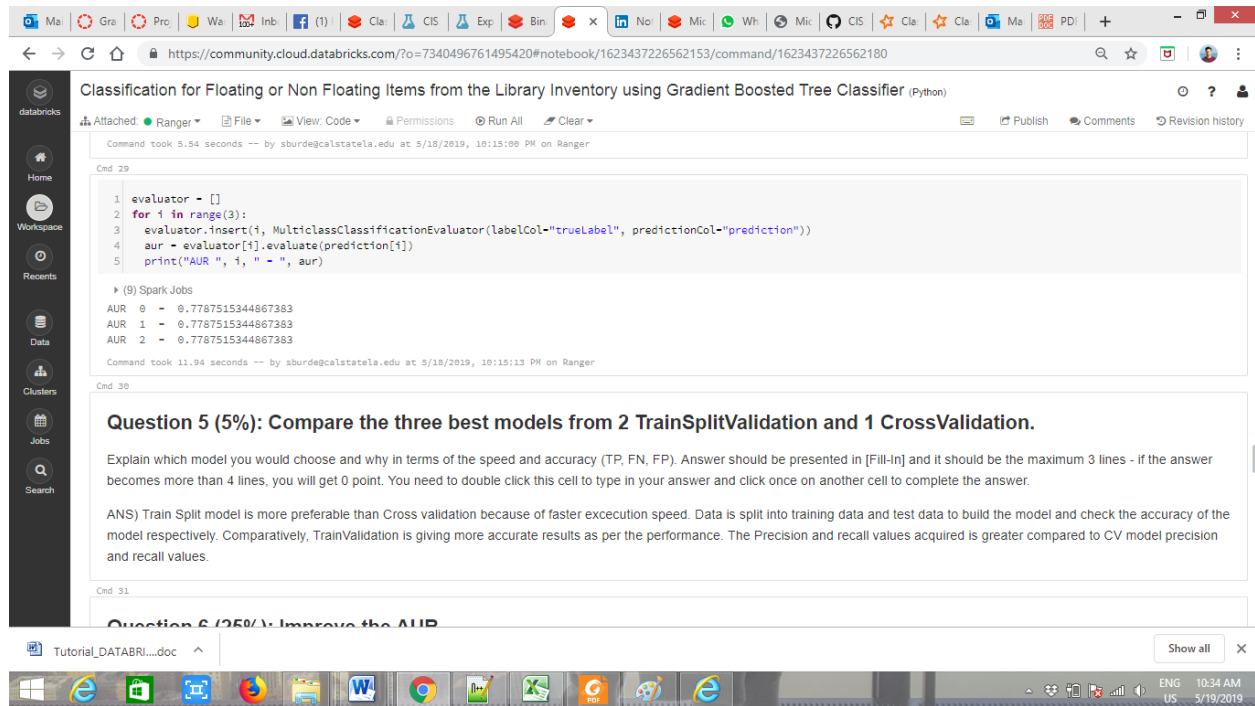
Cmd 28

# Step 8: Review the Area Under ROC



# Step 9: Improve AUC: Gradient Boosted Tree Classifier

Gradient boosted tree classifier model using Multiclass Classification evaluator help us to increase the accuracy (AUR). I am using training split 0.8 with TrainValidationSplit model

# Step 9.a Define assembler and pipeline for GBT



# Step 9.b Train Validation Split with Bins parameters

Classification for Floating or Non Floating Items from the Library Inventory using Gradient Boosted Tree Classifier (Python)

Cmd 37

```
1  #Define paramGrid for GBT
2  paramGrid_GBT = (ParamGridBuilder() \
3              .addGrid(gbt.maxDepth, [2, 4, 6]) \
4              .addGrid(gbt.maxBins, [20, 60]) \
5              .addGrid(gbt.maxIter, [10, 20]) \
6              .build())
```

Command took 0.03 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:16:32 PM on Ranger

Cmd 38

## Test the Model

Now you're ready to apply the model to the test data.

Cmd 39

```
1  #TrainValidation split for BinaryClassification
2  GBT_tvs = TrainValidationSplit(estimator=pipeline, evaluator=BinaryClassificationEvaluator(), estimatorParamMaps=paramGrid_GBT, trainRatio=0.8)
3  model_gbt = GBT_tvs.fit(train)
```

▶ (56) Spark Jobs

Command took 5.28 minutes -- by sburde@calstatela.edu at 5/18/2019, 10:16:38 PM on Ranger

Cmd 40

## Transform the data with Testing Splits

# Step 9.b Test the Model

Classification for Floating or Non Floating Items from the Library Inventory using Gradient Boosted Tree Classifier (Python)

**Test the Model**

Now you're ready to apply the model to the test data.

Cmd 39

```
1  #TrainValidation split for BinaryClassification
2  GBT_tvs = TrainValidationSplit(estimator=pipeline, evaluator=BinaryClassificationEvaluator(), estimatorParamMaps=paramGrid_GBT, trainRatio=0.8)
3  model_gbt = GBT_tvs.fit(train)
```

▸ (56) Spark Jobs

Command took 5.28 minutes -- by sburde@calstatela.edu at 5/18/2019, 10:16:38 PM on Ranger

Cmd 40

**Transform the data with Testing Splits**

Cmd 41

```
1  #Apply Transform to testing split and show 30 rows
2  prediction_gbt = model_gbt.transform(test)
3  predicted_gbt = prediction_gbt.select("features", "prediction", "probability", "trueLabel")
4  predicted_gbt.show(30)
```
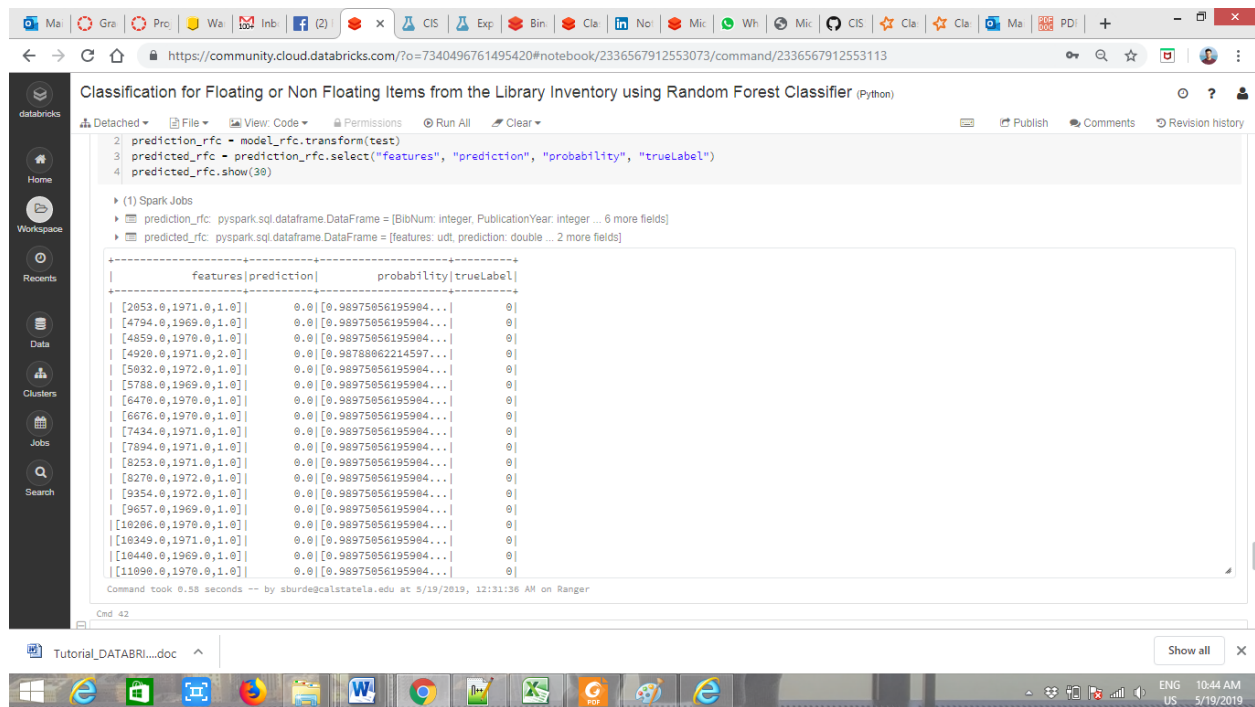
▸ (1) Spark Jobs

▸ ▦ prediction_gbt: pyspark.sql.dataframe.DataFrame = [BibNum: integer, PublicationYear: integer ... 6 more fields]

▸ ▦ predicted_gbt: pyspark.sql.dataframe.DataFrame = [features: udt, prediction: double ... 2 more fields]

# Step 9. c Review the Area under ROC



Classification for Floating or Non Floating Items from the Library Inventory using Gradient Boosted Tree Classifier (Python)

Cmd 42

**Review the Area Under ROC**

The performance of a classification model is to measure the area under a ROC curve using MulticlassClassificationEvaluator class to gain more accuracy compared with other classifier

Cmd 43

```
1  #print the AUR from the GBT model using Multiclass Classifier
2  evaluator_gbt = MulticlassClassificationEvaluator(labelCol="trueLabel", predictionCol="prediction")
3  aur_gbt = evaluator_gbt.evaluate(prediction_gbt)
4  print("AUR = ", aur_gbt)
```

▸ (3) Spark Jobs

AUR =  0.7873476526714038

Command took 3.35 seconds -- by sburde@calstatela.edu at 5/18/2019, 10:52:41 PM on Ranger

Cmd 44

**Result**

Area under ROC is much more accurate when using Multiclass Classification Evaluator than Binary Class Classifier using Gradient Boost Tree Classifier

Cmd 45

You may build an experiment in Azure ML Studio to get hints to add the function of data engineering or google to find out some hints how others have done. The following is the reference that you may look at as well

References to improve the accuracy

# Result

Area under ROC is much more accurate when using Multiclass Classification Evaluator than Binary Class Classifier using Gradient Boost Tree Classifier is 0.787.

# Step 10. Random Forest Classifier

Random Forest classifier model using Multiclass Classification evaluator help us to increase the accuracy (AUR). I am using training split 0.8 with TrainValidationSplit model

## Define assembler and pipline for RFC



## Train Validation Split with Bins parameters

# Test the Model

# Transform the data with Testing Splits



# Review the Area Under ROC by Random Forest Classifier

# Result

Area under ROC is much more accurate when using Multiclass Classification Evaluator than Binary Class Classifier using Random Forest Classifier is 0.788.

# References to improve the accuracy

1. Extracting, transforming and selecting features, https://spark.apache.org/docs/latest/ml-features.html
2. Basic data preparation in Pyspark — Capping, Normalizing and Scaling, http://bit.ly/2Ihs6Wa
3. Machine Learning with PySpark and MLlib — Solving a Binary Classification Problem, http://bit.ly/2Zb20tg