# IBM TRANSACTIONS FOR ANTI MONEY LAUNDERING (AML)

**CIS 5560 INTRO TO BIG DATA SCIENCE**

Professor: Jongwook Woo

**Submitted by:**

Lekha Ajit Kumar

Sushmitha Dandu

Dauren Omarov

Navyasree Sriramoju

# AGENDA

| | | | |
|---|---|---|---|
| Introduction | Dataset Specifications | Technical Specifications | Prediction System Flowchart |
| Machine Learning Algorithms | Feature Importance | Algorithms Comparison | Classification |
| GitHub Link | References | Summary | |

# INTRODUCTION

- The IBM Transactions for Anti Money Laundering (AML) dataset is a synthetic dataset that contains financial transactions involving individuals, companies, and banks.

- The dataset contains both legitimate and laundering transactions that are labeled, making it ideal for training and testing Anti Money Laundering models.

- We chose "LI_Medium_Trans.csv," file from the Dataset which is of 2.98gb and has 11 Columns

# DATASET SPECIFICATIONS

**DATASET NAME:** IBM Transactions for Anti Money Laundering (AML)

**TOTAL DATASET SIZE:** 2.98 GB

**DATASET FORMAT:** CSV

**DATASET URL:**
https://www.kaggle.com/datasets/ealtman2019/ibm-transactions-for-anti-money-laundering-aml?select=LI-Medium_Trans.csv

# TECHNICAL SPECIFICATION

**Hadoop Version:** 3.1.2

**No. of CPUs:** 8

**Py Spark version**: 3.0.2
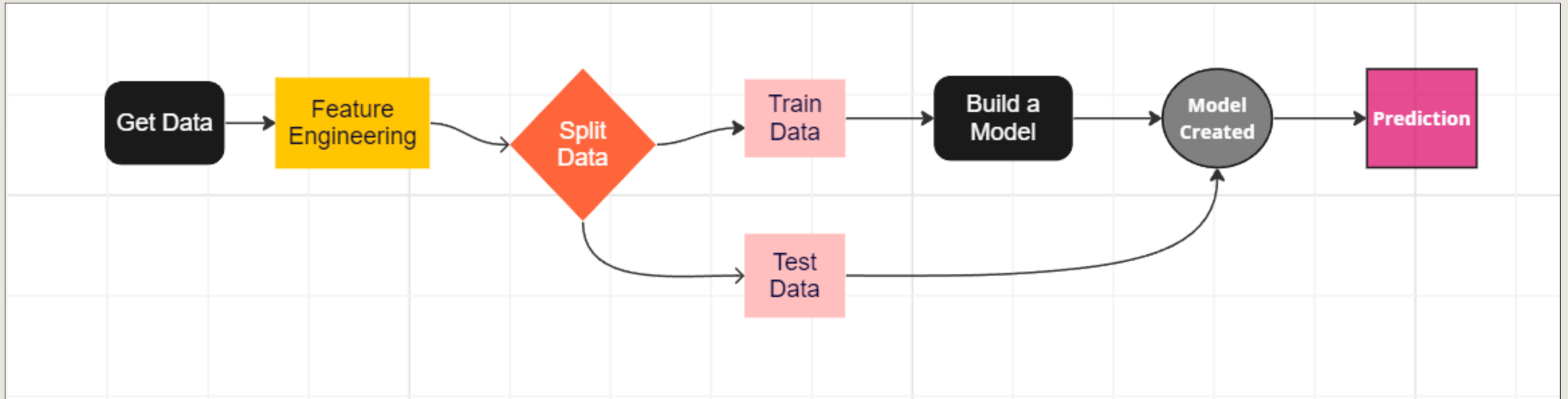
**Nodes:** 3

**Total Storage:**

**Databricks Community Version:** 10.4 LTS (includes Apache  Spark 3.1.1, Scala 2.12)

**File System:** DBFS (Data Bricks File System)

**Nodes**: 1

**Python Version**: 3.10.4

# PREDICTION SYSTEM FLOWCHART

# CLASSIFICATION

The Classification algorithm is a Supervised Learning technique that is used to identify the category of new observations based on training data.

In Classification, a program learns from the given dataset or observations and then classifies new observation into a number of classes or groups. Such as, **Yes or No, 0 or 1**

# SPLITTING THE DATASET

## Sample Data

```
1   #Finding the count of training and testing rows
2   train_rows = train.count()
3   test_rows = test.count()
4   print("Training Rows:", train_rows, " Testing Rows:", test_rows)
```

▸ (4) Spark Jobs

Training Rows: 3994   Testing Rows: 1714

## Full Data Set

```
>>> splits = data.randomSplit([0.7, 0.3])
>>> train = splits[0]
>>> test = splits[1].withColumnRenamed("label", "trueLabel")
>>> train_rows = train.count()
>>> test_rows = test.count()
>>> print("Training Rows:", train_rows, " Testing Rows:", test_rows)
Training Rows: 21882083   Testing Rows: 9373309
```

# MACHINE LEARNING ALGORITHMS

➢ Logistic Regression

➢ Gradient Boost Tree

➢ Decision Tree

➢ Random Forest

➢ Factorization Machine

➢ Support Vector Machine

# LOGISTIC REGRESSION

## CV & TVS and their respective AUC Values:

```python
1   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
2   auc_tvs_lr = evaluator.evaluate(prediction_lr_tvs)
3   print("AUC = ", auc_tvs_lr)
4
5   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
6   auc_cv_lr = evaluator.evaluate(prediction_lr_cv)
7   print("AUC = ", auc_cv_lr)
8
9   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
10  auc_lr = evaluator.evaluate(predicted)
11  print("AUC = ", auc_lr)
```

▶ (9) Spark Jobs

AUC =  0.5216146888078845
AUC =  0.5216146888078845
AUC =  0.5

Command took 1.33 seconds -- by lekha19202@gmail.com at 07/05/2023, 21:19:29 on Qs

# LOGISTIC REGRESSION

## Precision and recall Values:

```
1    # Precision and Recall
2    tp = float(prediction_lr_tvs.filter("prediction == 1.0 AND truelabel == 1").count())
3    fp = float(prediction_lr_tvs.filter("prediction == 1.0 AND truelabel == 0").count())
4    tn = float(prediction_lr_tvs.filter("prediction == 0.0 AND truelabel == 0").count())
5    fn = float(prediction_lr_tvs.filter("prediction == 0.0 AND truelabel == 1").count())
6    metrics2 = spark.createDataFrame([
7      ("TP", tp),
8      ("FP", fp),
9      ("TN", tn),
10     ("FN", fn),
11     ("Precision", tp / (tp + fp)),
12     ("Recall", tp / (tp + fn))],["metric", "value"])
13   metrics2.show()
```

▸ (11) Spark Jobs

▸ ▤  metrics2: pyspark.sql.dataframe.DataFrame = [metric: string, value: double]

```
+---------+-------------------+
|   metric|              value|
+---------+-------------------+
|       TP|                2.0|
|       FP|                0.0|
|       TN|             1669.0|
|       FN|               43.0|
|Precision|                1.0|
|   Recall|0.044444444444444446|
+---------+-------------------+
```

# GRADIENT BOOST TREE

CV & TVS and their respective AUC Values:

# DECISION TREE

CV & TVS and their respective AUC Values:

```python
1   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
2   auc_tvs_dt = evaluator.evaluate(predicted_dt_tvs)
3   print("AUC = ", auc_tvs_dt)
4
5   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
6   auc_cv_dt = evaluator.evaluate(predicted_dt_cv)
7   print("AUC = ", auc_cv_dt)
```

▸ (6) Spark Jobs

AUC =  0.5
AUC =  0.5

Command took 1.39 seconds -- by lekha19202@gmail.com at 07/05/2023, 22:14:37 on Qs

# RANDOM FOREST

## CV & TVS and their respective AUC Values:

```
Cmd 32
                                                                                    Python  ▶▾  ∨  —  ✕
1   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
2   auc_tvs1_rf = evaluator.evaluate(prediction)
3   print("AUC = ", auc_tvs1_rf)
4
5   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
6   auc_tvs_rf = evaluator.evaluate(predictiontvs)
7   print("AUC = ", auc_tvs_rf)
8
9   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
10  auc_cv_rf = evaluator.evaluate(predictionCV)
11  print("AUC = ", auc_cv_rf)
12

▶ (9) Spark Jobs

AUC =  0.917922235722965
AUC =  0.9269339813689754
AUC =  0.8963142972863507

Command took 1.77 seconds -- by lekha19202@gmail.com at 07/05/2023, 21:19:29 on Qs
```

# RANDOM FOREST

## Running time for the entire dataset

```
>>> #tvs = TrainValidationSplit(estimator=pipeline, evaluator=BinaryClassificationEvaluator(labelCol="label", rawPredictionCol="pr
UnderROC"), estimatorParamMaps=paramGrid, trainRatio=0.8)
...
>>> model = pipeline.fit(train)

>>> #model = tvs.fit(train)
...
>>> end = time()
>>> phrase = 'Random Forest tvs testing'
>>> print('{} takes {} seconds'.format(phrase, (end - start))) #round(end - start, 2)))
Random Forest tvs testing takes 1861.5901260375977 seconds
>>>
```

```
>>>
>>> end = time()
>>> phrase = 'Random Forest tvs2 testing'
>>> print('{} takes {} seconds'.format(phrase, (end - start))) #round(end - start, 2)))
Random Forest tvs2 testing takes 313.30773091316223 seconds
>>>
```

```
> end = time()
> phrase = 'Random Forest testing'
> print('{} takes {} seconds'.format(phrase, (end - start))) #round(end - start, 2)))
ndom Forest testing takes 301.2368438243866 seconds
>
```

# FACTORIZATION MACHINE

## CV & TVS and their respective AUC Values:

```
1   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
2   auc_tvs_fm = evaluator.evaluate(predicted_fm_tvs)
3   print("AUC = ", auc_tvs_fm)
4
5   evaluator = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
6   auc_cv_fm = evaluator.evaluate(predicted_fm_cv)
7   print("AUC = ", auc_cv_fm)
```

▶ (6) Spark Jobs

AUC =  0.7959497772377667
AUC =  8.7959497772377667

Command took 1.13 seconds -- by lekha19202@gmail.com at 07/05/2023, 21:19:29 on Qs

# SUPPORT VECTOR MACHINE

## CV & TVS and their respective AUC Values:

```python
1   evaluatorSVM_tvs = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
2   auc_SVM_tvs = evaluatorSVM_tvs.evaluate(predictionSVM_tvs)
3   print("AUC = ", auc_SVM_tvs)
4
5   evaluatorSVM_cv = BinaryClassificationEvaluator(labelCol="trueLabel", rawPredictionCol="prediction", metricName="areaUnderROC")
6   auc_SVM_cv = evaluatorSVM_cv.evaluate(predictionSVM_cv)
7   print("AUC = ", auc_SVM_cv)
```

▶ (6) Spark Jobs

AUC =  0.5
AUC =  0.5

Command took 0.91 seconds -- by lekha19282@gmail.com at 07/05/2023, 21:19:29 on Qs

# FEATURE IMPORTANCE

```
1  import pandas as pd
2  featureImp = pd.DataFrame(list(zip(finalVect.getInputCols(),rfModel.featureImportances)),columns=["feature", "importance"])
3  featureImp.sort_values(by="importance", ascending=False)
```

Running command 27. Go to

| | feature | importance |
|---|---|---|
| 0 | TimestampIdx | 0.354191 |
| 1 | From Bank | 0.258357 |
| 3 | To Bank | 0.211514 |
| 2 | Account2 | 0.175938 |
| 0 | From Bank | 0.354191 |
| 1 | To Bank | 0.258357 |
| 3 | Amount Paid | 0.211514 |
| 2 | Amount Received | 0.175938 |

Feature Importance indicates how much each feature contributes to the model prediction. Basically, it determines the degree of usefulness of a specific variable for a current model and prediction.

# ALGORITHMS COMPARISON

| Algorithms | AUC | | Computation Time | |
|---|---|---|---|---|
| | CV | TVS | CV | TVS |
| Gradient Boost Tree | 0.895916 | 0.8959160 | 195.254653sec | 71.486958sec |
| Random Forest | 0.896314 | 0.926933 | 57.922800sec | 5.795582Sec(Tvs) 24.588072Sec(Tvs2) |
| Factorization Machine | 0.795949 | 0.795949 | 32.444634sec | 14.9555480 sec |
| Decision Tree | 0.5 | 0.5 | 448.105305sec | 140.117408 sec |
| Logistic Regression | 0.521614 | 0.521614 | 1566.377578sec | 497.938576 sec |
| Support Vector Machine | 0.5 | 0.5 | 138.414037sec | 57.775793sec |

# IMPLEMENTATION IN SPARK ML

# GITHUB LINK

https://github.com/Lekha19202/CIS-5560-big-data-science-project

# REFERENCES

- Li, Susan. "Machine Learning with Pyspark and MLlib - Solving a Binary Classification Problem." *Medium*, Towards Data Science, 7 May 2018, https://towardsdatascience.com/machine-learning-with-pyspark-and-mllib-solving-a-binary-classification-problem-96396065d2aa.

- "Regression Analysis." *Corporate Finance Institute*, 3 May 2023, https://corporatefinanceinstitute.com/resources/data-science/regression-analysis/.

- Jagdeesh. "PySpark Decision Tree – How to Build and Evaluate Decision Tree Model for Classification Using PySpark MLlib." *Machine Learning Plus*, 1 May 2023, https://www.machinelearningplus.com/pyspark/pyspark-decision-tree/.

**1** ——————— We learned how to compare the accuracy of different models based on AUC, Presession and Recall

**2** ——————— How to optimize computation time by changing the paramgrid parameters.

**3** ——————— How to increase the accuracy of the model.

**4** ——————— Building pipeline

# SUMMARY

# THANK YOU