# CSCE 5300: Intro. to Big Data and Data Science

# Project

# Fraud Detection in Financial Transactions

**Team members**
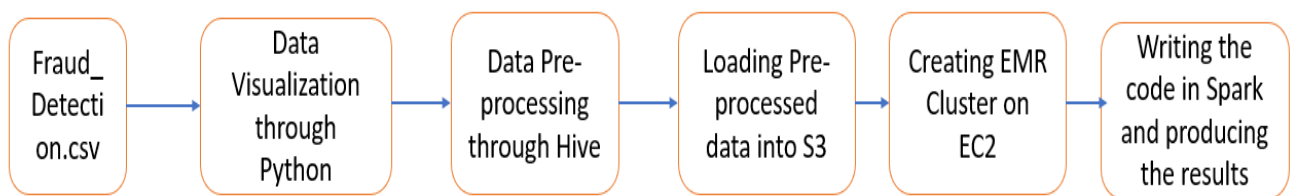
Hari Krishna Sai Rachuri

Vivek Nelluri

Sathwika Karingu

Krishna Anuhya Regalla

## Problem Statement:

Reducing the financial impact of credit card fraud, identity theft, and payment fraud requires constructing a robust fraud detection system. Big data analysis requires the use of sophisticated tools like Hadoop, Spark, and Hive because to the significant increase in digital transactions and the resulting expansion in data. A proficient system must incorporate data processing, machine learning, and real-time analytics to effectively prevent fraud.

## Workflow Diagram:



## Tools used in throughout the project :

- **AWS**
- **HIVE**

**Implementation Status Report :**

**Description :**

We are a team of 4 people working on the fraud transactions detection model. We have implemented Random Forest classifier to train our model in pyspark and Python.

**Responsibility:**

Sathwika – (Team Lead) Helped in building the EMR cluster and working throughout the Notebook.

Anuhya – Worked on the Data Visualization

Vivek – Data Preprocessing through Hive.

Hari - Writing the code for training the model and Hyper tuning parameters.

**Contributions:**

Sathwika took 25 percent of the project contribution in making the report and the presentation for the project. She is responsible with AWS and Cloudera setup.

Anuhya took 25 percent of the project contribution in making the report and the presentation for the project and also performed Data Visualization using python until the data Preprocessing steps

Vivek- took 25 percent of the project contribution in making the report and the presentation for the project and have written code in SQL.

Hari - took 25 percent of the project contribution in making the report and the presentation for the project and have written code in Python.

**Issues/ Concerns:**

We could have performed more on Hive. We have rather focused on Machine Learning and testing the model instead of working more on preprocessing part. That's our major concern in this project.

**Walking through the Data set :**

Step – It's a integer Data Type which refers to transactions numbering

Type – It's a String Data Type which contains the payment mode

Amount – Integer or Float value with amount

nameOrig -  This column likely represents the originating account or entity's name or identifier. It identifies the source of the transaction.

oldbalanceOrg – Integer value with amount before transaction

newbalanceOrig: Following the transaction, the balance in the originating account is shown in this column. It displays the account balance following a payment or transfer of funds.

nameDest: This field most likely contains the name or identification of the destination account or entity. It indicates who the transaction's recipient or destination is.

oldbalanceDest: The balance in the destination account prior to the transaction is shown in this column. It shows the amount in the recipient's account prior to the funds being sent.

newbalanceDest: Following the transaction, the balance in the destination account is shown in this column. It displays the recipient's account balance following the money transfer.

isFraud- An single integer to state whether the transaction is Fraud or not

isFlaggedFraud – This is also similar to the above variable, It's a potential
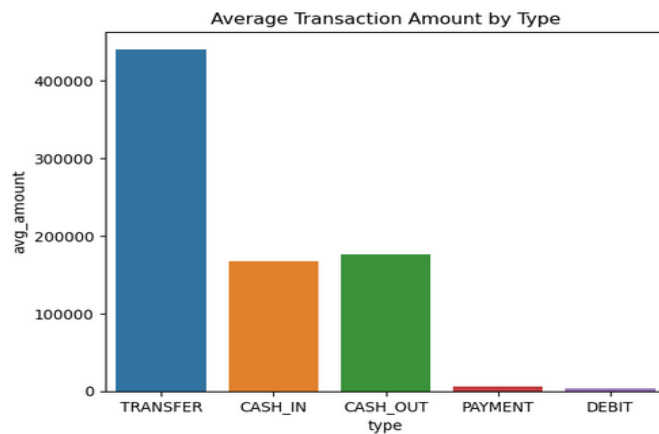
additional fraud detection measure.
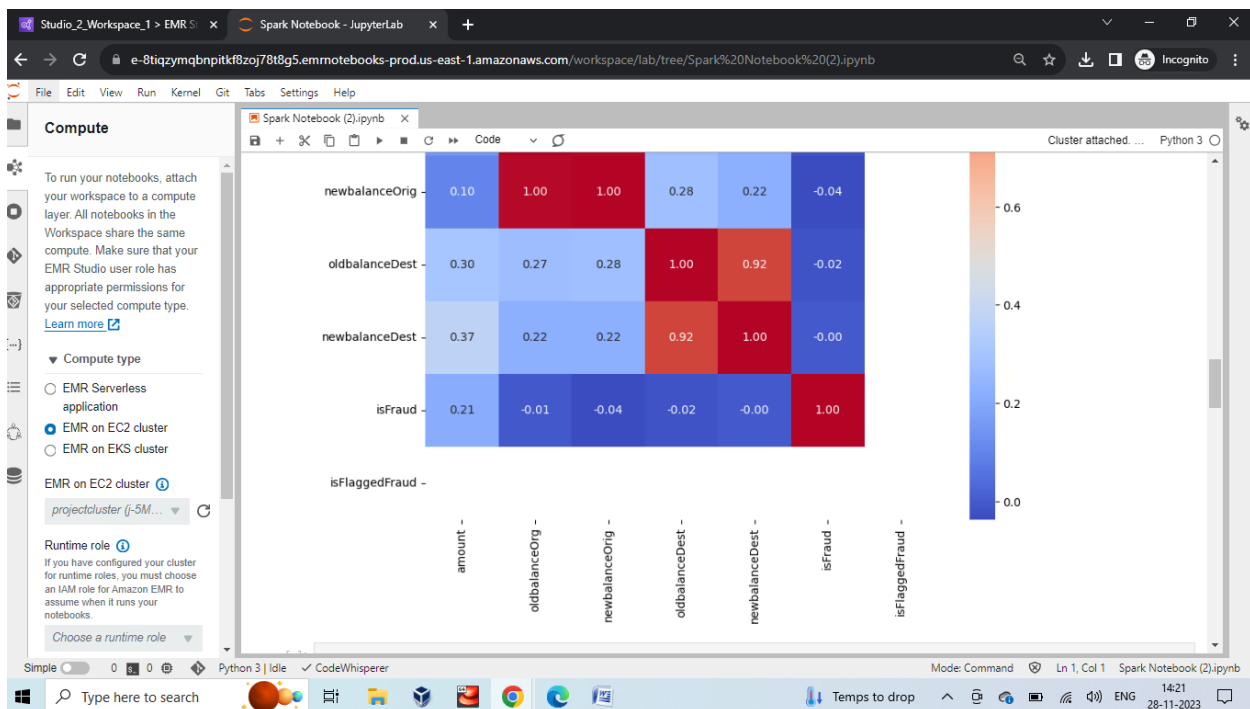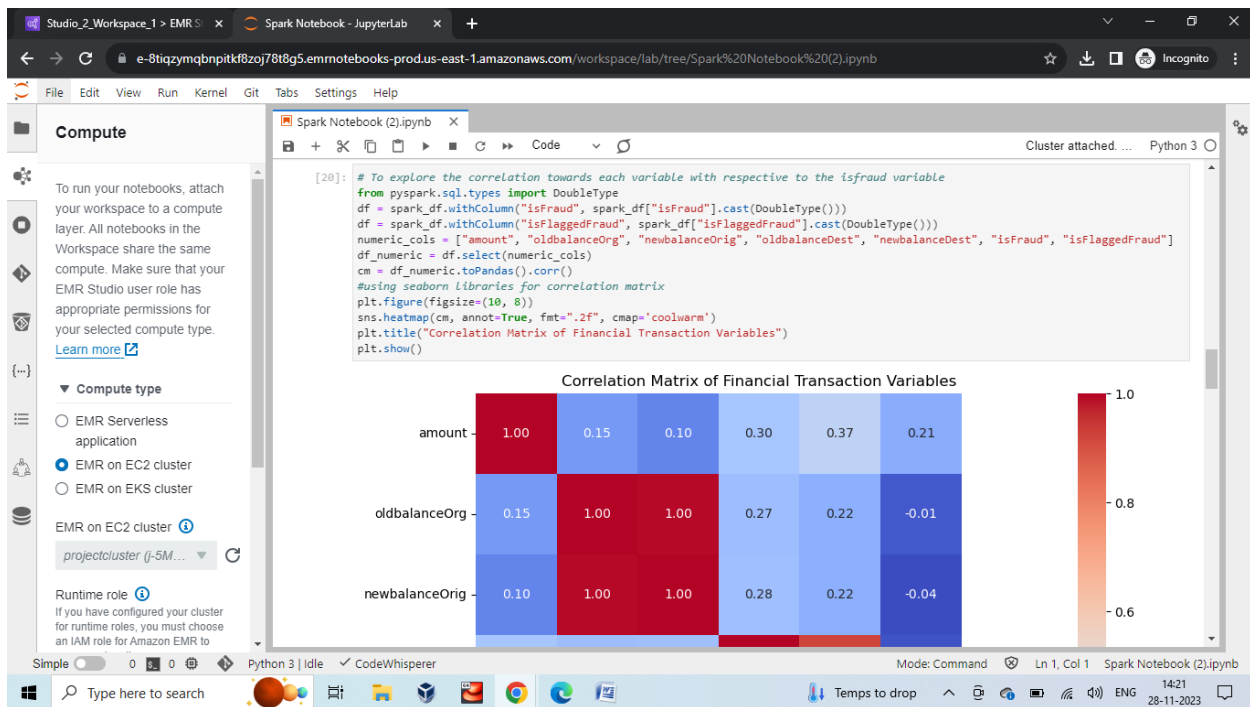
# Data Visualization in Python:

```
In [18]: #Data pre-processing
         selected_cols = ["amount", "oldbalanceOrg", "newbalanceOrig", "oldbalanceDest", "newbalanceDest", "isFraud"]
         data = spark_df.select(selected_cols)
         data = data.dropna()
         data = data.withColumn("label", lit(1.0).cast("double"))
```

```
In [19]: #Data Visualization to understand the relationship in between the variables.
         from pyspark.sql import functions as F

         type_vs_amount = spark_df.groupBy("type").agg(
             F.avg("amount").alias("avg_amount"),
             F.sum("amount").alias("total_amount")
         ).toPandas()
         #after converting to pandas dataframe i have used bar chart to compare the values
         import matplotlib.pyplot as plt
         import seaborn as sns

         sns.barplot(x='type', y='avg_amount', data=type_vs_amount)
         plt.title('Average Transaction Amount by Type')
         plt.show()
         #Looks like transcations with "transfers" are more in number.
```

File  Edit  View  Run  Kernel  Git  Tabs  Settings  Help

Compute

Spark Notebook (2).ipynb

Code    Cluster attached. ...    Python 3

```python
[20]:  # To explore the correlation towards each variable with respective to the isfraud variable
       from pyspark.sql.types import DoubleType
       df = spark_df.withColumn("isFraud", spark_df["isFraud"].cast(DoubleType()))
       df = spark_df.withColumn("isFlaggedFraud", spark_df["isFlaggedFraud"].cast(DoubleType()))
       numeric_cols = ["amount", "oldbalanceOrg", "newbalanceOrig", "oldbalanceDest", "newbalanceDest", "isFraud", "isFlaggedFraud"]
       df_numeric = df.select(numeric_cols)
       cm = df_numeric.toPandas().corr()
       #using seaborn libraries for correlation matrix
       plt.figure(figsize=(10, 8))
       sns.heatmap(cm, annot=True, fmt=".2f", cmap='coolwarm')
       plt.title("Correlation Matrix of Financial Transaction Variables")
       plt.show()
```

Correlation Matrix of Financial Transaction Variables

|               | amount | oldbalanceOrg | newbalanceOrig |      |      |       |
|---------------|--------|---------------|----------------|------|------|-------|
| amount        | 1.00   | 0.15          | 0.10           | 0.30 | 0.37 | 0.21  |
| oldbalanceOrg | 0.15   | 1.00          | 1.00           | 0.27 | 0.22 | -0.01 |
| newbalanceOrig| 0.10   | 1.00          | 1.00           | 0.28 | 0.22 | -0.04 |

To run your notebooks, attach your workspace to a compute layer. All notebooks in the Workspace share the same compute. Make sure that your EMR Studio user role has appropriate permissions for your selected compute type. Learn more

▼ Compute type

○ EMR Serverless application
● EMR on EC2 cluster
○ EMR on EKS cluster

EMR on EC2 cluster ⓘ
projectcluster (j-5M...  ▼

Runtime role ⓘ
If you have configured your cluster for runtime roles, you must choose an IAM role for Amazon EMR to

Simple    0    0    Python 3 | Idle    ✓ CodeWhisperer    Mode: Command    Ln 1, Col 1    Spark Notebook (2).ipynb

Temps to drop    ENG    14:21    28-11-2023

---

File  Edit  View  Run  Kernel  Git  Tabs  Settings  Help

Compute

Spark Notebook (2).ipynb

Code    Cluster attached. ...    Python 3

|                | amount | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newbalanceDest | isFraud | isFlaggedFraud |
|----------------|--------|---------------|----------------|----------------|----------------|---------|----------------|
| newbalanceOrig | 0.10   | 1.00          | 1.00           | 0.28           | 0.22           | -0.04   |                |
| oldbalanceDest | 0.30   | 0.27          | 0.28           | 1.00           | 0.92           | -0.02   |                |
| newbalanceDest | 0.37   | 0.22          | 0.22           | 0.92           | 1.00           | -0.00   |                |
| isFraud        | 0.21   | -0.01         | -0.04          | -0.02          | -0.00          | 1.00    |                |
| isFlaggedFraud |        |               |                |                |                |         |                |

To run your notebooks, attach your workspace to a compute layer. All notebooks in the Workspace share the same compute. Make sure that your EMR Studio user role has appropriate permissions for your selected compute type. Learn more

▼ Compute type

○ EMR Serverless application
● EMR on EC2 cluster
○ EMR on EKS cluster

EMR on EC2 cluster ⓘ
projectcluster (j-5M...  ▼

Runtime role ⓘ
If you have configured your cluster for runtime roles, you must choose an IAM role for Amazon EMR to assume when it runs your notebooks.
Choose a runtime role  ▼

Simple    0    0    Python 3 | Idle    ✓ CodeWhisperer    Mode: Command    Ln 1, Col 1    Spark Notebook (2).ipynb

Temps to drop    ENG    14:21    28-11-2023

**Data Pre-processing in Hive :**

Step -01 :

Setup Hive in the Cloudera Platform and create a empty table with the schema as it is present in the csv file.

**Code to load the data from CSV to Hive Table:**

```
CREATE TABLE FraudTransactions (step INT,type STRING,amount DOUBLE,

    nameOrig STRING,

    oldbalanceOrg DOUBLE,

    newbalanceOrig DOUBLE,

    nameDest STRING,

    oldbalanceDest DOUBLE,

    newbalanceDest DOUBLE,

isFraud INT,

    isFlaggedFraud INT

)
ROW FORMAT DELIMITED

FIELDS TERMINATED BY ','

STORED AS TEXTFILE;
```

we can execute several Hive Queries to understand the data and perform necessary analysis in Hive: The following are the queries used to extract insights from the queries :

1) Total Number of transactions by type :

query: SELECT type, AVG(amount) AS avg_transaction_amount FROM Fraud GROUP BY type;

2) Money laundering transactions may involve multiple transfers with small amounts. You can identify such patterns by grouping transactions by nameOrig and calculating the total number of outgoing transactions and the total amount transferred for each account

SELECT nameOrig, COUNT(*) AS num_outgoing_transactions, SUM(amount) AS total_outgoing_amount

FROM your_table_name

WHERE type IN ('TRANSFER', 'CASH_OUT')

GROUP BY nameOrig

HAVING num_outgoing_transactions > 1 AND total_outgoing_amount <1000;

3) Based on your fraud detection criteria, you can flag transactions that meet certain conditions as potentially fraudulent. For example, transactions where isFraud is not flagged but the amount is significantly higher than average.

query: SELECT F.*,

     CASE

WHEN F.isFraud = 0 AND F.amount > A.avg_amount THEN 1

ELSE 0

END AS flagged_as_potential_fraud

FROM FraudTransactions F

CROSS JOIN (SELECT AVG(amount) as avg_amount FROM Fraud) A;

**Insights:**

- The last query will yield a new column which says potential fraud. This will flag transactions that meet certain criteria as potentially fraudulent.

- It will return all columns from the FraudTransactions table along with a new column flagged_as_potential_fraud. Transactions where isFraud is 0 (not flagged as fraud) and the amount is significantly higher than the average amount in the entire Fraud dataset will be flagged as 1; otherwise, they will be flagged as 0.

- Based on the isFraud and Potential_fraud transactions, there isn't much difference. So it's better to go with the same dataset we have taken.

**EMR Cluster Creation on EC2 instance (AWS) :**

**EC2 instance Launch :**

# Cluster Creation :



# Workspace Creation :

# Loading the Notebook for writing the code in pyspark :



```
[16]: pip install pyspark
      sc.install_pypi_package("pandas==1.0.7") # Install a specific version of pandas
      sc.install_pypi_package("matplotlib", "https://pypi.org/simple") # Install the latest version of matplotlib from the given PyPI
      sc.install_pypi_package("seaborn", "https://pypi.org/simple")
      #sc.install_pypi_package("sklearn","https://github.com/scikit-Learn/sklearn-pypi-package")
      sc.install_pypi_package("scikit-learn")
```

```
[ ]: #Importing Libraries
     from pyspark.ml.classification import RandomForestClassifier
     from pyspark.ml.evaluation import BinaryClassificationEvaluator
     from pyspark import SparkContext
     from pyspark import SparkConf
     from pyspark.sql import SparkSession
     from pyspark.sql.functions import lit
     from pyspark.ml.feature import VectorAssembler, StandardScaler
     from pyspark.ml import Pipeline
     #For visualizations
     import seaborn as sns
     import matplotlib.pyplot as plt
     from sklearn.metrics import roc_curve, auc
```
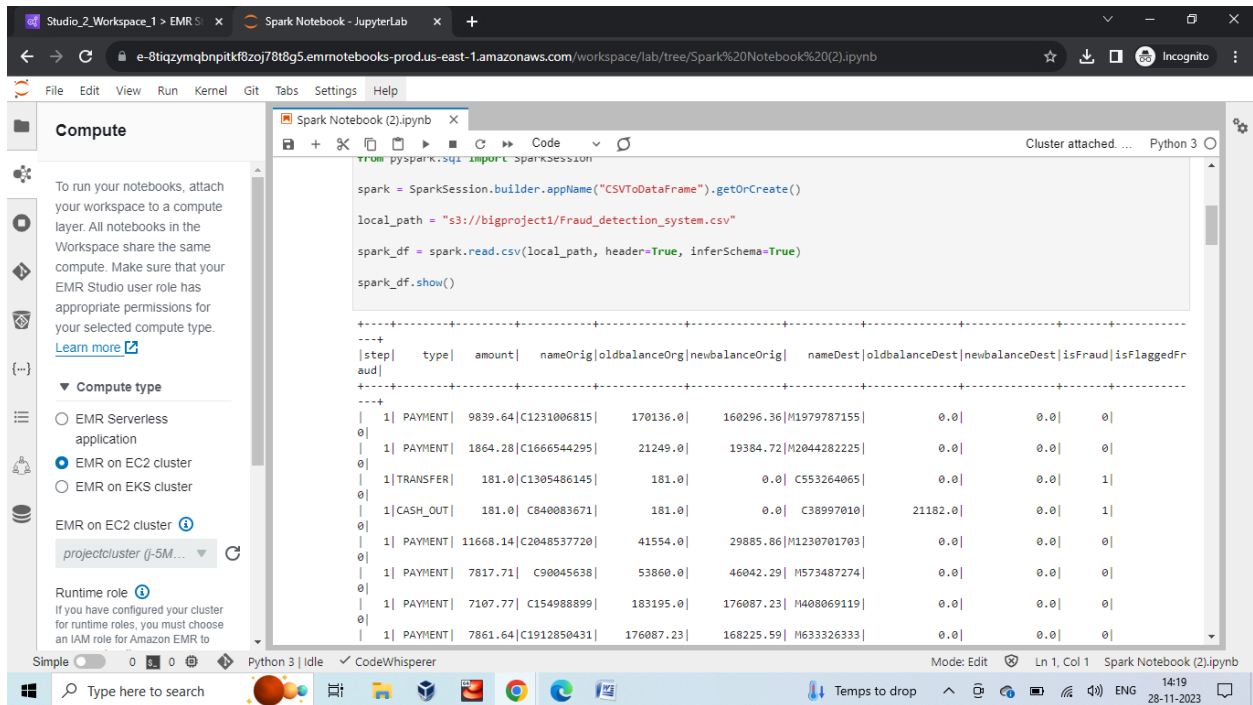
```
[17]: #Loading the dataset into Spark's Data frame
      from pyspark.sql import SparkSession

      spark = SparkSession.builder.appName("CSVToDataFrame").getOrCreate()

      local_path = "s3://bigproject1/Fraud_detection_system.csv"
```

# Loading the data into Spark Data Frame :



```python
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("CSVToDataFrame").getOrCreate()

local_path = "s3://bigproject1/Fraud_detection_system.csv"

spark_df = spark.read.csv(local_path, header=True, inferSchema=True)

spark_df.show()
```
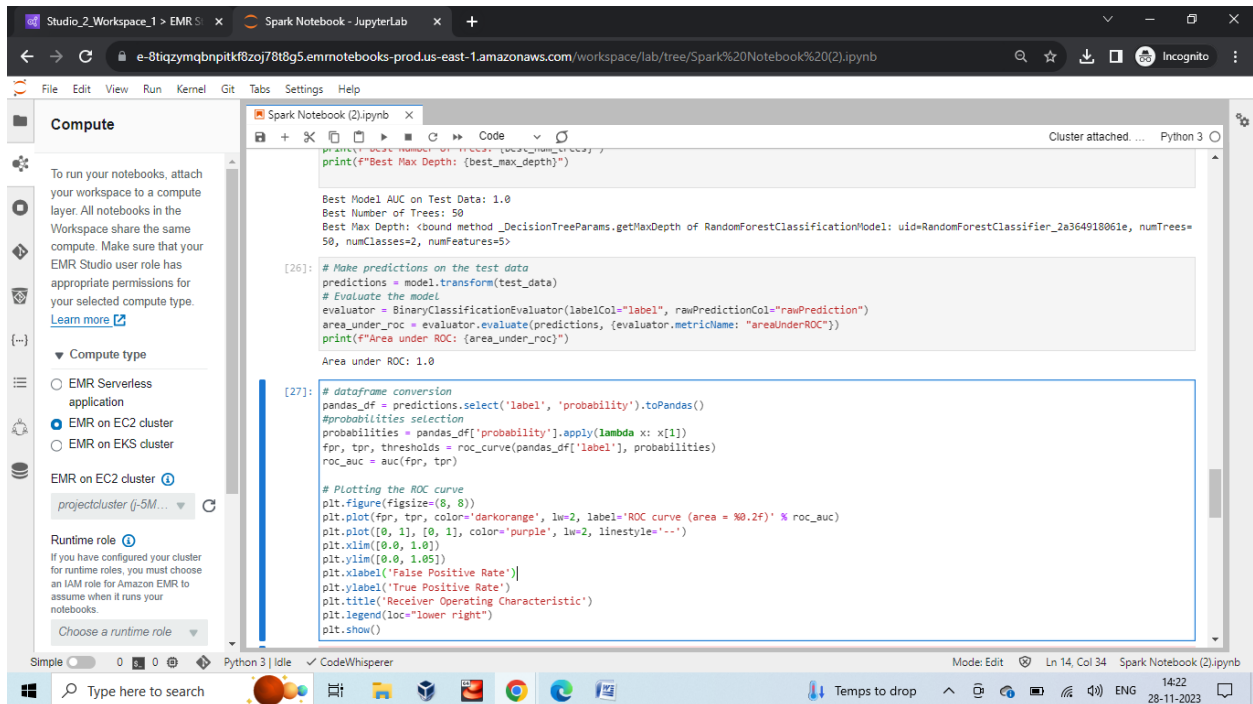
```
+----+--------+---------+-----------+-------------+--------------+-------------+--------------+--------------+-------+------------+
|step|    type|   amount|   nameOrig|oldbalanceOrg|newbalanceOrig|     nameDest|oldbalanceDest|newbalanceDest|isFraud|isFlaggedFraud|
+----+--------+---------+-----------+-------------+--------------+-------------+--------------+--------------+-------+------------+
|   1| PAYMENT|  9839.64|C1231006815|     170136.0|     160296.36|M1979787155|           0.0|           0.0|      0|           0|
|   1| PAYMENT|  1864.28|C1666544295|      21249.0|      19384.72|M2044282225|           0.0|           0.0|      0|           0|
|   1|TRANSFER|    181.0|C1305486145|        181.0|           0.0| C553264065|           0.0|           0.0|      1|           0|
|   1|CASH_OUT|    181.0| C840083671|        181.0|           0.0|  C38997010|       21182.0|           0.0|      1|           0|
|   1| PAYMENT| 11668.14|C2048537720|      41554.0|      29885.86|M1230701703|           0.0|           0.0|      0|           0|
|   1| PAYMENT|  7817.71|  C90045638|      53860.0|      46042.29| M573487274|           0.0|           0.0|      0|           0|
|   1| PAYMENT|  7107.77| C154988899|     183195.0|     176087.23|M408069119|           0.0|           0.0|      0|           0|
|   1| PAYMENT|  7861.64|C1912850431|     176087.23|    168225.59| M633326333|           0.0|           0.0|      0|           0|
```



```python
    print(f"Best Max Depth: {best_max_depth}")
```

```
Best Model AUC on Test Data: 1.0
Best Number of Trees: 50
Best Max Depth: <bound method _DecisionTreeParams.getMaxDepth of RandomForestClassificationModel: uid=RandomForestClassifier_2a364918061e, numTrees=50, numClasses=2, numFeatures=5>
```
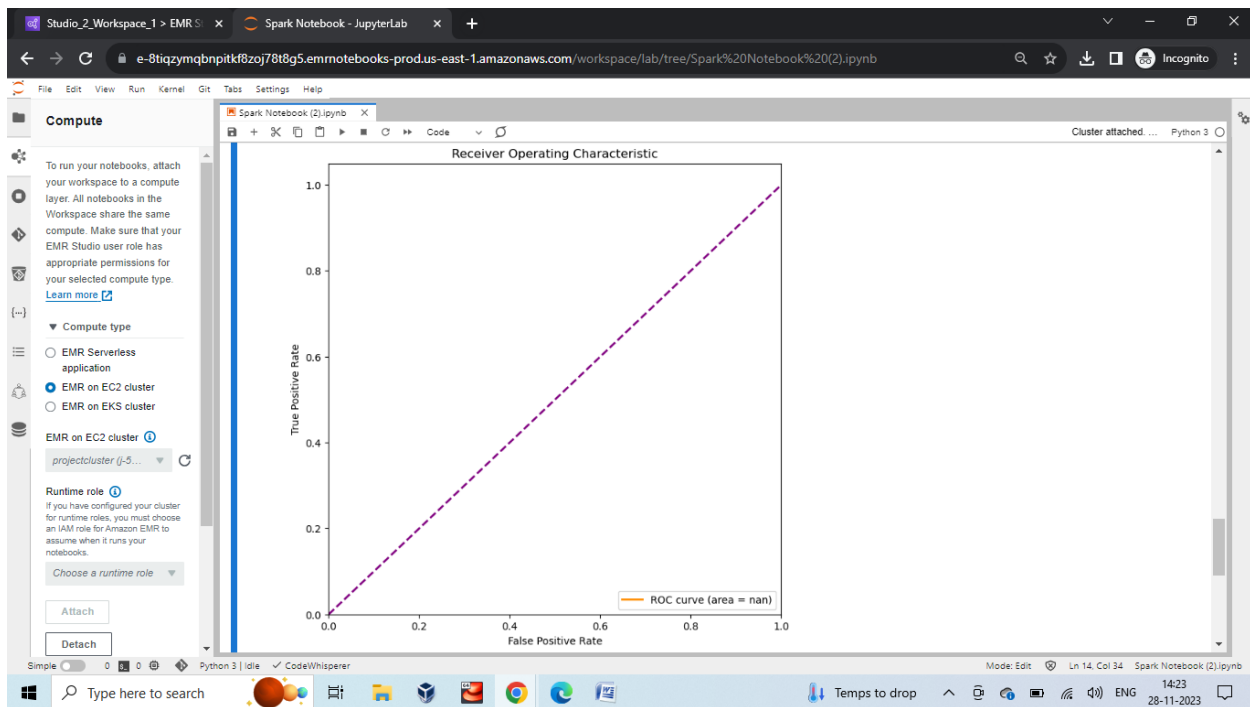
```python
[26]: # Make predictions on the test data
predictions = model.transform(test_data)
# Evaluate the model
evaluator = BinaryClassificationEvaluator(labelCol="label", rawPredictionCol="rawPrediction")
area_under_roc = evaluator.evaluate(predictions, {evaluator.metricName: "areaUnderROC"})
print(f"Area under ROC: {area_under_roc}")
```

```
Area under ROC: 1.0
```

```python
[27]: # dataframe conversion
pandas_df = predictions.select('label', 'probability').toPandas()
#probabilities selection
probabilities = pandas_df['probability'].apply(lambda x: x[1])
fpr, tpr, thresholds = roc_curve(pandas_df['label'], probabilities)
roc_auc = auc(fpr, tpr)

# Plotting the ROC curve
plt.figure(figsize=(8, 8))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='purple', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic')
plt.legend(loc="lower right")
plt.show()
```

- **Final Conclusion:**

- The Receiver Operating Characteristic (ROC) curve offers a graphical depiction of how well a model can differentiate between two categories, such as positive and negative instances. This curve enables the evaluation of the model's proficiency in distinguishing true positives (accurately identified positive cases) from false positives (wrongly identified positive cases), over a range of probability thresholds.

- The area under ROC curve indicates that 1 which indicates a perfect classifier

- The "evaluate method" of the evaluator computes the metric specified for the model predictions provided.

- In this case, {evaluator.metricName: "areaUnderROC"} sets the metric to the Area Under the ROC Curve (AUC - ROC).

- The AUC - ROC is a single scalar value that summarizes the performance of the binary classification model across all classification thresholds.
- It essentially measures the ability of the model to discriminate between the two classes.