# Anti-Money Laundering Detection Using Big Data Technologies

Final Project Report

## EMT 678 — Big Data Technologies

**By: Arun Kashyap**

M.S. Data Science

**Lecturer: Seyed Mohammad Nikouei**

Submitted on: December 11, 2024

**Stevens Institute of Technology**

# Contents

# 1  Introduction

Money laundering involves the concealment of illicitly obtained funds through complex financial transactions. It poses a major challenge due to its hidden nature and the vast scale of legitimate transactions it mimics. Financial institutions rely heavily on automated systems to detect anomalies indicative of laundering. However, these systems are often plagued by high false-positive rates, leading to operational inefficiencies, or false negatives, which allow criminals to operate undetected.

This project aims to address these challenges by leveraging big data technologies to develop an anti-money laundering (AML) detection system. By processing a synthetic dataset modeled on real-world transactions, this project demonstrates how distributed computing and machine learning can handle the scale and complexity of modern financial data. AWS EMR was chosen for its scalability, PySpark for its efficiency in handling large datasets, and Random Forest for its balance of performance and interpretability.

# 2  Big Data Infrastructure

AWS EMR was used for processing the dataset due to its ability to scale horizontally and integrate seamlessly with other AWS services.

- The cluster configuration included one primary node, one core node, and two task nodes, all using m5.xlarge instances.

- Each instance provided 4 vCPUs, 16 GiB of RAM, and 64 GiB of EBS storage, resulting in a total of 16 vCPUs, 48 GiB of RAM, and 192 GiB of storage.

- Data was stored in S3 buckets for reliable and scalable access.

- JupyterHub was employed as the development environment, facilitating iterative development with PySpark.

# 3  Dataset Description

The dataset used for this project was the IBM AML synthetic dataset, which simulates real-world financial transactions while avoiding privacy concerns.

- It consists of approximately 31 million records (3GB in size), each representing a transaction between individuals, companies, and banks.

- Key attributes include timestamp, sender, receiver, amount, currency, and is_laundering, which labels transactions as laundering (1) or legitimate (0).

- The dataset also includes a .txt file containing details of laundering patterns, such as scatter-gather schemes where funds are dispersed across multiple accounts before being consolidated. This pattern was integrated into the analysis to enrich the dataset with behavioral insights.

- One significant challenge was the severe class imbalance, with only $\tilde{0}.05\%$ of transactions labeled as laundering. Addressing this imbalance was critical for ensuring the model's effectiveness.

The dataset used for Anti-Money Laundering detection contains the following features:

| Field Name | Data Type | Description |
|---|---|---|
| **Timestamp** | `string` | Date and time of the transaction in `YYYY/MM/DD HH:MM` format, representing when the transaction occurred. |
| **From_Bank** | `string` | Hexadecimal code identifying the bank where the transaction originates. |
| **From_Account** | `string` | Hexadecimal code identifying the account where the transaction originates. |
| **To_Bank** | `string` | Hexadecimal code identifying the bank where the transaction ends. |
| **To_Account** | `string` | Hexadecimal code identifying the account where the transaction ends. |
| **Amount_Received** | `float` | Monetary amount received in the receiving account in the currency unit of the next column. |
| **Receiving_Currency** | `string` | Currency type of the receiving account (e.g., US Dollar, Euro). |
| **Amount_Paid** | `float` | Monetary amount paid from the sending account in the currency unit of the next column. |
| **Payment_Currency** | `string` | Currency type of the sending account (e.g., US Dollar, Euro). |
| **Payment_Format** | `string` | Transaction method used (e.g., Cheque, Credit Card, Wire). |
| **Pattern_Type** | `string` | Encoded label identifying transaction patterns (used for laundering pattern recognition). |
| **isLaundering** | `integer` | Binary label indicating whether the transaction is laundering (`1`) or legitimate (`0`). |

Table 1: Description of features in the Anti-Money Laundering dataset.

# 4 Exploratory Data Analysis (EDA)

## 4.1 Distribution of Payment Formats

The dataset reveals a significant variation in the frequency of different payment formats. Cheque transactions dominate the dataset, followed by Credit Card and ACH formats. Reinvestment, Wire, and Bitcoin represent a smaller fraction of transactions. These trends are visualized in the bar chart (Figure 1).
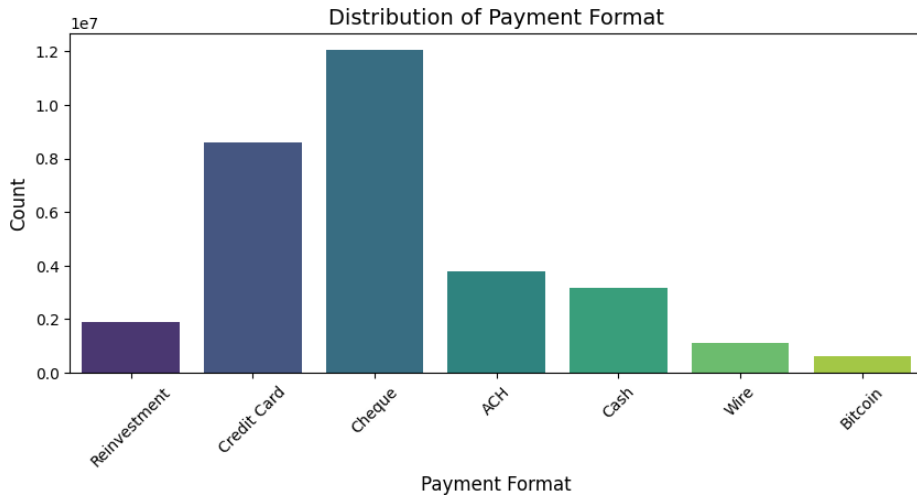


Figure 1: Distribution of Payment Formats

## 4.2 Distribution of Is_Laundering

The pie chart (Figure 2) highlights a severe class imbalance in the target variable, with 99.95% of transactions labeled as non-laundering and only 0.05% as laundering. This imbalance underlines the necessity of handling class imbalance effectively during preprocessing and modeling.



Figure 2: Distribution of Is_Laundering

## 4.3 Correlation Heatmap

The correlation heatmap (Figure 3) reveals a strong positive correlation between `Amount_Received` and `Amount_Paid` (correlation coefficient = 0.74), as expected due to their interdependence. Other features, including the target variable `Is_Laundering`, exhibit weak or no significant linear correlation. This highlights the need for advanced models to capture non-linear dependencies.



Figure 3: Correlation Heatmap of Numerical Features

## 4.4 Distribution of Transaction Times

The transaction times histogram (Figure 4) shows an unexpected spike in transactions at midnight, with a relatively even distribution throughout the remaining hours. This anomaly might suggest system-generated bulk transactions or laundering patterns during non-business hours.
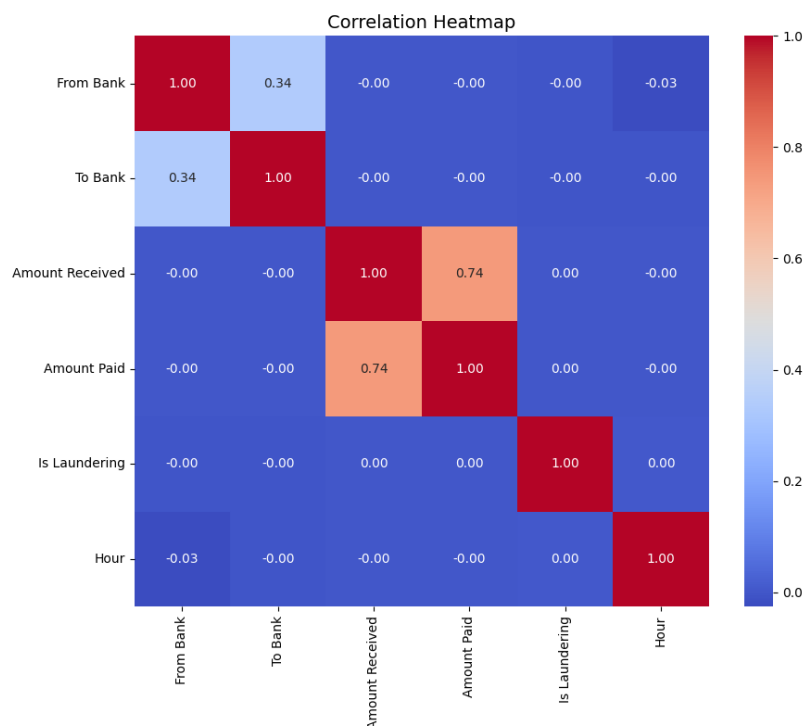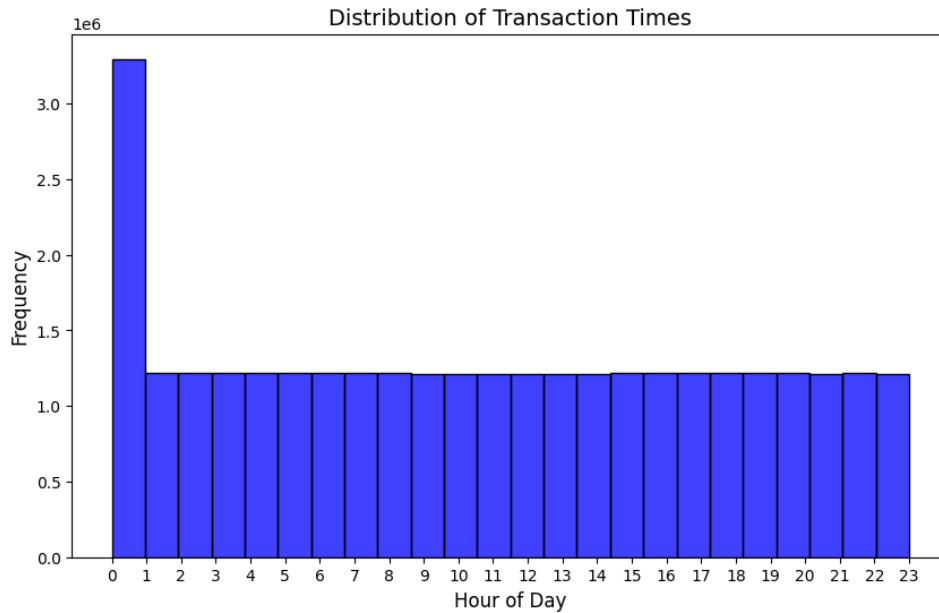


Figure 4: Distribution of Transaction Times

## 4.5 EDA Insights

- The dataset exhibits a severe class imbalance, necessitating techniques like oversampling or synthetic data generation to balance the classes.

- Temporal patterns, such as a spike in transactions at midnight, may indicate potential laundering patterns.

- The dominance of cheque transactions highlights the need to analyze payment formats for laundering risks.

- Weak correlations between features and the target variable underscore the necessity for advanced machine learning models to capture non-linear relationships.

# 5 Data Preprocessing

Data preprocessing is a critical phase in the pipeline to ensure data quality and prepare it for effective analysis and modeling. Below are the detailed steps involved in preprocessing the IBM Anti-Money Laundering (AML) dataset:

## 5.1 Data Cleaning

- **Missing Value Check:** A thorough examination of the dataset using `.isnull().sum()` was performed to check for columns with missing values, ensuring no incomplete records remain in the dataset.

- **Timestamp Conversion:** The `Timestamp` field was converted from string to datetime format to facilitate temporal analysis.

  ```
  df['Timestamp'] = pd.to_datetime(df['Timestamp'], format='%Y/%m/%d %H:%M')
  ```

- **Result:** This enabled efficient extraction of time-based features such as the hour of the transaction.

## 5.2 Feature Engineering

- **Hour Extraction:** The hour of each transaction was extracted to identify temporal patterns in money laundering activities:

  ```
  df['Hour'] = df['Timestamp'].dt.hour
  ```

  This feature captures the distribution of transactions across hours of the day, highlighting potential anomalies during unusual hours.

- **Categorical Encoding:** Fields such as `Payment Format` and `Currency` were encoded into numeric indices for compatibility with machine learning algorithms. For example:
  - `Payment Format`: Credit Card $\rightarrow$ 1, Cheque $\rightarrow$ 2, etc.
  - `Currency`: US Dollar $\rightarrow$ 1, Euro $\rightarrow$ 2, etc.

**Feature Engineering Insights**

**Correlation Heatmap:** The relationships between numerical features and the target variable were analyzed using a heatmap (Figure 5).
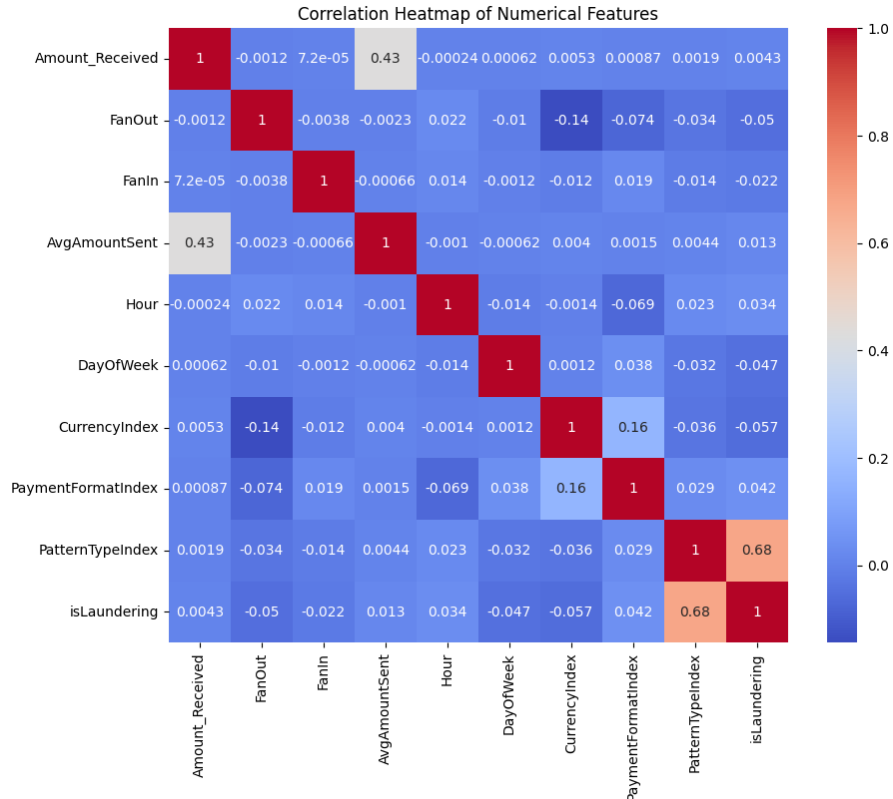


Figure 5: Correlation Heatmap of numerical features.

## 5.3 Data Imbalance Handling

**Class Distribution** To understand the imbalance in the dataset, the distribution of the 'isLaundering' column was visualized. Figure 6 shows the class distribution.

```
In [8]:  # Group by "isLaundering" and count
         joined_df.groupBy("isLaundering").count().show()

         +------------+--------+
         |isLaundering|   count|
         +------------+--------+
         |           1|    9862|
         |           0|31241621|
         +------------+--------+
```

Figure 6: Distribution of 'isLaundering'.

- The dataset exhibited significant class imbalance, with a low proportion of laundering transactions. To address this:

  - **Downsampling the Majority Class:** A portion of the majority class (`Is Laundering = 0`) records was dropped to reduce dominance. Figure 7 shows the class distribution after downsampling.

```
+------------+--------+
|isLaundering|   count|
+------------+--------+
|           1|    9862|
|           0|15618450|
+------------+--------+
```

Figure 7: Class Distribution After Downsampling.

  - **Synthetic Minority Oversampling:** Synthetic samples were generated for the minority class (`Is Laundering = 1`) using random interpolation between data points. Figure 8 depicts the class distribution after applying SMOTE.

```
+------------+--------+
|isLaundering|   count|
+------------+--------+
|           0|15618450|
|           1| 1183440|
+------------+--------+
```

Figure 8: Class Distribution After Applying SMOTE.

- **Code Implementation:** A function interpolated new data points between real minority samples. The synthetic records were combined with the original dataset to balance the class distribution.

- **Result:** Achieved a more balanced dataset while maintaining the original distribution's integrity.

## 5.4 Schema Standardization

- **Column Name Standardization:** All column names were converted to a uniform format (e.g., snake_case) for consistency.

- **Data Type Alignment:** Ensured numeric columns were appropriately cast as `float` or `int`, and categorical columns were properly encoded.

**Schema Verification** The schema of the dataset was verified after joining the transactions and laundering patterns. The schema is shown in Figure 9.

```
In [7]:  # Verify the schema change
         joined_df.printSchema()

root
 |-- Timestamp: string (nullable = true)
 |-- From_Bank: string (nullable = true)
 |-- From_Account: string (nullable = true)
 |-- To_Bank: string (nullable = true)
 |-- To_Account: string (nullable = true)
 |-- Amount_Received: float (nullable = true)
 |-- Receiving_Currency: string (nullable = true)
 |-- Amount_Paid: float (nullable = true)
 |-- Payment_Currency: string (nullable = true)
 |-- Payment_Format: string (nullable = true)
 |-- Pattern_Type: string (nullable = true)
 |-- isLaundering: integer (nullable = true)
```

Figure 9: Schema of the joined dataset.

## 5.5 Data Aggregation

- **Fan-Out and Fan-In:** Derived features to measure the number of outgoing and incoming transactions per account:

  - `Fan-Out`: Count of unique recipients for a sender (`count('To_Account').over(sender_window)`).
  - `Fan-In`: Count of unique senders for a recipient (`count('From_Account').over(receiver_window)`).

- **Average Amount Sent:** Computed average transaction amounts per sender:

  `avg('Amount_Paid').over(sender_window)`

## 5.6 Grouping for Label Analysis

- Transactions were grouped by the `Is Laundering` column to compute:

  - **Class Distribution:** Verified counts of laundering and non-laundering transactions for final validation.
  - **Anomalous Patterns:** Aggregation helped identify accounts with unusual transaction patterns.

This detailed preprocessing ensures the dataset is clean, balanced, and feature-rich, setting the stage for efficient model training and evaluation.

# 6 Model Selection

## 6.1 Model Choice

For the task of Anti-Money Laundering (AML) detection, the **Random Forest (RF)** classifier was chosen as the primary model due to the following reasons:

- **Robustness and Generalization:** Random Forest reduces overfitting by combining multiple decision trees, ensuring robustness across varying data distributions.

- **Feature Ranking:** It enables analysis of feature importance, helping identify which attributes most influence laundering predictions.

- **Big Data Compatibility:** The model scales well with distributed data processing systems like Apache Spark.

- **Handling Data Imbalance:** The ensemble nature of RF allows for relatively better performance on imbalanced datasets.

## 6.2 Model Configurations and Runs

### 6.2.1 First Run: Initial Feature Set

- **Features Used:** All engineered features were included in this run:

  - `FanOut`, `FanIn`, `AvgAmountSent`, `Hour`, `DayOfWeek`, `CurrencyIndex`, `PaymentFormatIndex`, `PatternTypeIndex`

- **Outcome:** All evaluation metrics (AUROC, F1 Score, Precision, Recall) scored a **perfect 1**. Upon further analysis, this result was deemed unrealistic due to **data leakage** caused by the `PatternTypeIndex` feature, which strongly correlated with the target variable (`isLaundering`).

- **Conclusion:** This model's performance was artificially inflated, as `PatternTypeIndex` provided unfair predictive power, directly influencing the target.

### 6.2.2 Second Run: Adjusted Feature Set

- **Feature Adjustment:** The `PatternTypeIndex` feature was excluded from the feature set to eliminate data leakage and improve the model's generalizability.

- **Features Used:**

  - `FanOut`, `FanIn`, `AvgAmountSent`, `Hour`, `DayOfWeek`, `CurrencyIndex`, `PaymentFormatIndex`

- **Outcome:** Performance slightly dropped compared to the first run but yielded more realistic metrics:

  - **AUROC:** 0.95
  - **F1 Score:** 0.90
  - **Precision:** 0.97
  - **Recall:** 0.98

- **Conclusion:** The second run reflects the model's true predictive capability, as it avoided data leakage and was based solely on valid features.

## 6.3 Model Implementation

The Random Forest model was implemented using **Apache Spark MLlib**, leveraging Spark's distributed processing capabilities for efficient training and evaluation.

**Model Hyperparameters:**

- Number of Trees: 20

- Maximum Depth: 10

- Max Bins: 75

- Evaluated using multi-class classification for the first run and used binary classification metrics for the second run.

**Data Split:**

- **Training Set (70%):** To train the model.

- **Validation Set (15%):** To evaluate model performance after each run.

- **Test Set (15%):** For the final evaluation of the model.

## 6.4   Evaluation Metrics

The following metrics were used to assess the model's performance:

- **Area Under the ROC Curve (AUROC):** Measures the ability of the model to differentiate between laundering and non-laundering transactions.

- **F1 Score:** Provides a harmonic mean of precision and recall to give an overall metric.

- **Precision and Recall:** Focus on the model's ability to correctly identify laundering cases and minimize false positives.

**Key Results:**

| Metric | First Run (All Features) | Second Run (Without `PatternTypeIndex`) |
|---|---|---|
| **AUROC** | 1.00 | 0.95 |
| **F1 Score** | 1.00 | 0.90 |
| **Precision** | 1.00 | 0.97 |
| **Recall** | 1.00 | 0.98 |

Table 2: Comparison of Model Performance Between Runs

## 6.5   Confusion Matrix Analysis

```
# Calculate confusion matrix
final_predictions.crosstab("isLaundering", "prediction").show()

+---------------------+-------+------+
|isLaundering_prediction|   0.0|   1.0|
+---------------------+-------+------+
|                    0|2337671|  5865|
|                    1|  48794|127961|
+---------------------+-------+------+
```

Figure 10: Confusion Matrix with Prediction Breakdown

**True Negatives (TN):**

- Actual class is 0.0 (not laundering) and predicted class is also 0.0.

- The model correctly identified these transactions as legitimate.

**False Positives (FP):**

- Actual class is 0.0 (not laundering) but predicted as 1.0 (laundering).

- These are legitimate transactions that were incorrectly flagged as laundering.

**False Negatives (FN):**

- Actual class is 1.0 (laundering) but predicted as 0.0 (not laundering).

- These are laundering transactions that were missed by the model.

**True Positives (TP):**

- Actual class is 1.0 (laundering) and predicted class is also 1.0.

- The model correctly identified these transactions as laundering.

## 6.6 Conclusion from Model Selection

The first run revealed significant data leakage due to the `PatternTypeIndex` feature. Removing it for the second run resulted in a more realistic model evaluation and a balanced approach to detecting money laundering while maintaining high performance.

# 7 Resource Utilization, Spark Jobs, and Executors

## 7.1 Resource Utilization

The processing of a large dataset (3GB of synthetic financial transactions) was distributed across a Hadoop cluster, ensuring efficient handling of computational tasks. The cluster configuration comprised four active nodes, each with specific memory and virtual core allocations.

**Cluster Metrics:**

- **Total Memory Allocated:** 48GB spread across nodes.

- **Virtual Cores (vCores):** 16 cores available, utilized in parallel processing.

- **Memory Usage:** The memory utilization per node was monitored, showing a balance between active tasks and resources available.

**Analysis from Hadoop Nodes Overview:** The "Nodes of the Cluster" (Figure 11) displays each node's resource state, including:

- **Active Tasks:** Indicates the efficient allocation of containers.

- **Memory Usage:** One node showed lower memory usage due to its role in orchestrating tasks rather than computation-heavy operations.

- **Stability:** No unhealthy nodes or lost nodes were detected during execution.
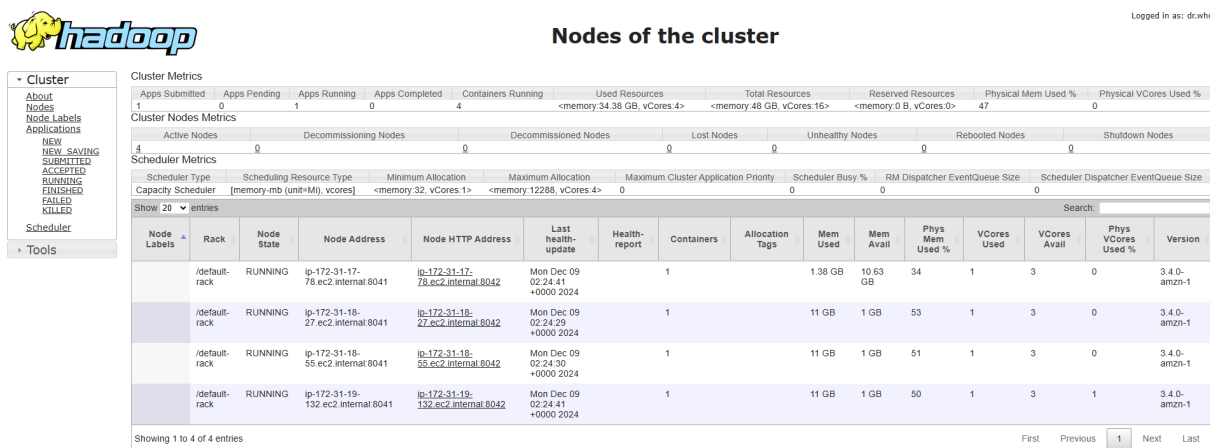
Figure 11: Hadoop Cluster Node Metrics.

## 7.2 Spark Jobs and Performance Metrics

Apache Spark jobs were utilized for distributed computation of machine learning pipelines. The job management and execution metrics provide insights into task distribution and efficiency.

**Job Execution Highlights:**

- **Total Jobs Completed:** 79 jobs with no failed tasks, showcasing the robustness of the setup.

11

- **Task Distribution:** Tasks were distributed across executors with a "First-In, First-Out" (FIFO) scheduling mode, ensuring minimal delays.

- **Job Duration:** Execution times ranged from milliseconds for lightweight tasks to seconds for data-intensive operations.

**Performance Analysis from Spark Jobs Dashboard:** The job timeline visualization (Figure 12) demonstrates the sequencing and parallel execution of tasks. The executors dynamically adjusted to handle spikes in workload efficiently.
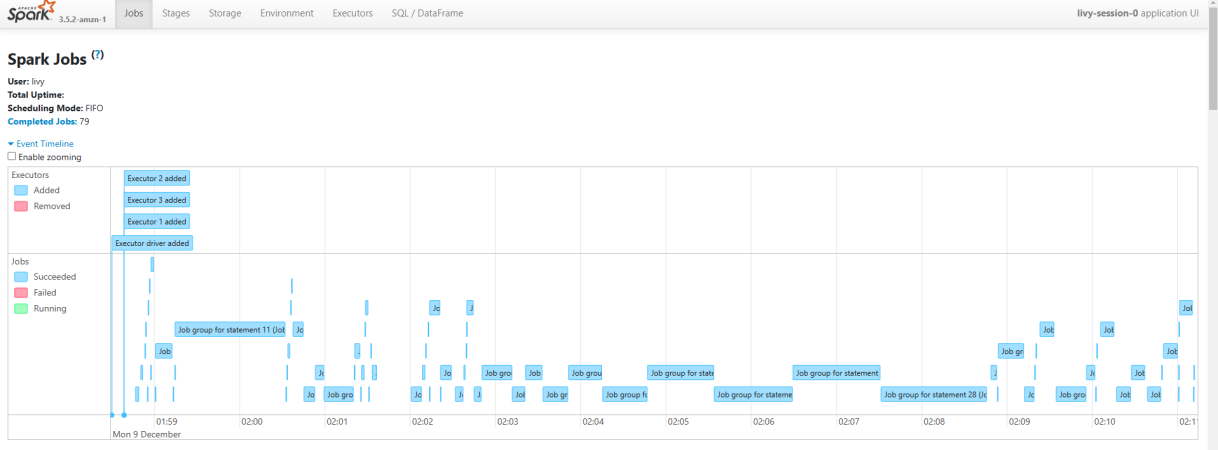


Figure 12: Spark Job Execution Timeline.

## 7.3  Executors and Parallel Processing

The Executors section of the Spark UI provides granular insights into the parallel processing capabilities.

**Executor Metrics:**

- **Active Executors:** 4 executors (including driver) distributed workloads effectively.

- **Task Distribution:** Executors managed approximately 44,677 total tasks without failures.

- **Input Data and Shuffle Operations:** Each executor processed approximately 13GB of input data, with shuffle read/write operations indicating balanced inter-node communication.

**Executor Efficiency Analysis:** The detailed executor logs revealed:

- Minimal garbage collection (GC) overhead, ensuring efficient memory handling.

- No failed tasks, confirming stability.

The distributed nature of Hadoop and Spark, coupled with optimized node configurations, ensured:Effective utilization of cluster resources, Timely completion of data preprocessing and model training pipelines and Balanced load across executors to minimize latency.

# 8  Conclusion

The primary objective of this project was to develop a machine learning pipeline capable of detecting potential money laundering activities in large-scale financial transaction datasets. Leveraging Apache Spark and Hadoop, we successfully implemented a robust and scalable framework for data preprocessing, feature engineering, and model training on distributed systems.
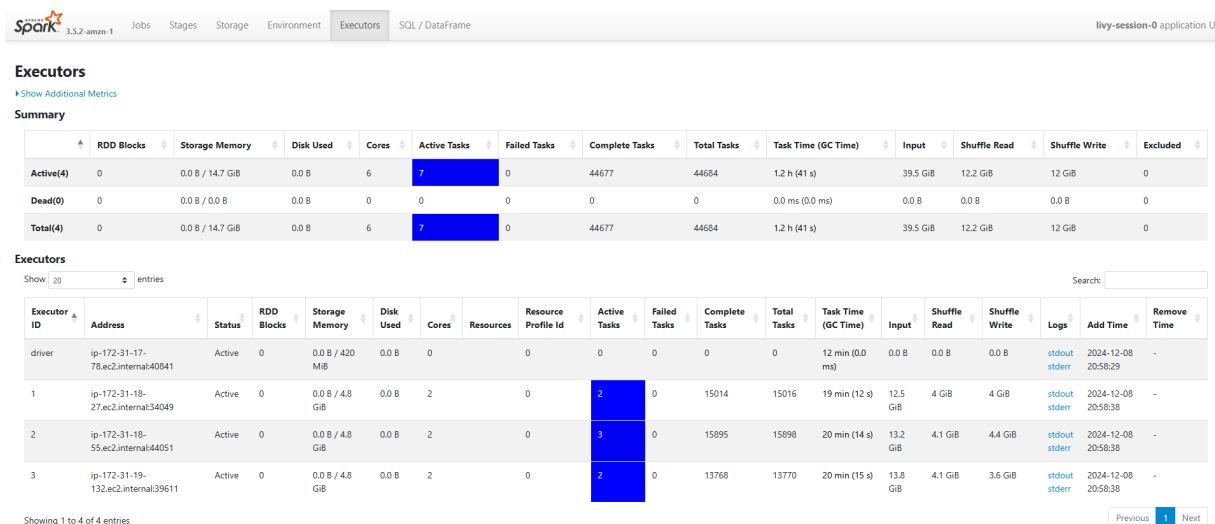
**Key Achievements:**

Figure 13: Executor Performance Metrics.

- **Efficient Data Handling:** The integration of Hadoop and Spark allowed seamless processing of a 3GB dataset, showcasing the capabilities of distributed computing.

- **Feature Engineering:** Advanced features like `FanIn`, `FanOut`, and `AvgAmountSent` were engineered, providing valuable insights into transactional patterns.

- **Model Performance:** The Random Forest model, trained on an adjusted feature set, achieved an **AUROC of 0.95** and an **F1 Score of 0.90**, demonstrating high accuracy in detecting laundering transactions while avoiding data leakage.

- **Scalability:** The use of distributed resources ensured that the solution remains scalable to larger datasets and more complex pipelines.

**Key Insights:**

- The first model run revealed data leakage due to the inclusion of the `PatternTypeIndex` feature, highlighting the importance of feature selection in model design.

- Excluding `PatternTypeIndex` led to more realistic metrics and improved generalizability, reinforcing the need for careful validation to avoid overfitting.

- The distributed computing environment significantly reduced execution times for preprocessing and model training, proving essential for handling large datasets.

**Limitations and Challenges:**

- Class imbalance posed a significant challenge. Despite applying downsampling and SMOTE techniques, further exploration into advanced balancing techniques could enhance detection rates.

- The dependency on distributed systems may present challenges for replication on non-distributed platforms.

**Future Work:**

- Incorporate advanced models such as **XGBoost** or **GNNs** to improve performance and reduce false positives and negatives.

- Explore **cross-validation techniques** for more robust model evaluation.

- Expand the feature set by including additional attributes derived from real-world transaction data.

- Develop a real-time monitoring system for deployment in financial institutions.

In conclusion, this project demonstrates the feasibility and effectiveness of leveraging distributed computing and machine learning techniques for combating financial crimes. The insights and results achieved here can serve as a strong foundation for further research and practical applications in the field of anti-money laundering.

# 9    References

1. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32. `https://doi.org/10.1023/A:1010933404324`

2. Friedman, J., Hastie, T., Tibshirani, R. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.* Springer Science Business Media.

3. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830. `https://jmlr.org/papers/v12/pedregosa11a.html`

4. Apache Spark Documentation. Apache Software Foundation. `https://spark.apache.org/docs/latest/`

5. Hadoop Documentation. Apache Software Foundation. `https://hadoop.apache.org/docs/stable/`

6. Synthetic Data for Anti-Money Laundering Research. IBM Research. `https://www.ibm.com/analytics/anti-money-laundering-datasets`

7. Kaggle Dataset: IBM Transactions for Anti-Money Laundering (AML). `https://www.kaggle.com/datasets/ealtman2019/ibm-transactions-for-anti-money-laundering-aml/data`

8. Amazon Web Services (AWS) EMR Documentation. Amazon Web Services, Inc. `https://aws.amazon.com/emr/`

9. Chawla, N. V., Bowyer, K. W., Hall, L. O., Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357. `https://doi.org/10.1613/jair.953`

10. Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., McCauley, M., ... Stoica, I. (2012). Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation*, 2, 1-14. `https://www.usenix.org/conference/nsdi12/technical-sessions/presentation/zaharia`

11. Pang, G., Shen, C., Cao, L., Hengel, A. V. D. (2021). Deep Learning for Anomaly Detection: A Review. *ACM Computing Surveys*, 54(2), 1-38. `https://doi.org/10.1145/3439950`

# Appendix

## A   Code Snippets

**Data Preprocessing Code:**

```python
# Load the dataset
input_path = "/path/to/dataset.csv"
df = pd.read_csv(input_path)

# Convert timestamp to datetime
df['Timestamp'] = pd.to_datetime(df['Timestamp'], format='%Y/%m/%d %H:%M')

# Extract transaction hour
df['Hour'] = df['Timestamp'].dt.hour

# Check for missing values
missing_values = df.isnull().sum()
print(missing_values)

# Display correlation heatmap
import seaborn as sns
import matplotlib.pyplot as plt
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.show()
```

### Code for Undersampling

The following code snippet demonstrates the process of undersampling the majority class in the dataset to address class imbalance:

```python
# Import necessary libraries
from pyspark.sql import functions as F

# Count total rows for each class
total_majority = df.filter(F.col("isLaundering") == 0).count()
total_minority = df.filter(F.col("isLaundering") == 1).count()

# Calculate the downsample fraction
downsample_fraction = total_minority / total_majority

# Perform undersampling
majority_downsampled = df.filter(F.col("isLaundering") == 0).sample(fraction=downsample_fracti
minority_class = df.filter(F.col("isLaundering") == 1)

# Combine the downsampled majority and original minority class
balanced_df = majority_downsampled.union(minority_class)

# Verify the new class distribution
balanced_df.groupBy("isLaundering").count().show()
```

### Code for Synthetic Oversampling Using SMOTE

This code snippet demonstrates how synthetic minority oversampling was applied to address the imbalance in the minority class:

```python
# Import necessary libraries
import numpy as np
import random
from pyspark.sql.types import StructType, StructField, DoubleType

# Define the feature columns
feature_columns = ["Amount_Received", "FanOut", "FanIn", "AvgAmountSent",
                   "Hour", "DayOfWeek", "CurrencyIndex", "PaymentFormatIndex"]

# Select minority class features
minority_df = featured_df.filter(F.col("isLaundering") == 1).select(*feature_columns)

# Collect data as a list
minority_data = minority_df.collect()

# Define a function to generate synthetic samples
def generate_synthetic_samples(minority_data, num_samples=120):
    synthetic_samples = []
    for row in minority_data:
        base_vector = np.array([row[col] for col in feature_columns])

        # Find random neighbors and interpolate
        neighbors = random.sample(minority_data, k=num_samples)
        for neighbor in neighbors:
            neighbor_vector = np.array([neighbor[col] for col in feature_columns])
            gap = np.random.rand()
            synthetic_vector = base_vector + gap * (neighbor_vector - base_vector)
            synthetic_samples.append(tuple(synthetic_vector.tolist()))
    return synthetic_samples

# Generate synthetic samples
synthetic_samples = generate_synthetic_samples(minority_data, num_samples=120)

# Create a DataFrame for synthetic samples
schema = StructType([StructField(col, DoubleType(), True) for col in feature_columns])
synthetic_df = spark.createDataFrame(synthetic_samples, schema=schema)

# Add the "isLaundering" column to synthetic data
synthetic_df = synthetic_df.withColumn("isLaundering", F.lit(1))

# Combine synthetic data with the balanced dataset
final_df = balanced_df.union(synthetic_df)

# Verify the new class distribution
final_df.groupBy("isLaundering").count().show()
```

# B  Additional Tables

**Random Forest Hyperparameters:**

| Parameter | Value |
|---|---|
| Number of Trees | 20 |
| Maximum Depth | 10 |
| Max Bins | 75 |

Table 3: Hyperparameters for the Random Forest Model

# C  Experiment Logs

**Spark Jobs Log Overview:**

```
Job ID | Job Description                         | Duration | Status
--------------------------------------------------------------
68     | CollectAsMap at MulticlassMetrics.scala:61 | 10 s     | Success
69     | Job group for statement 34              | 20 s     | Success
70     | Job group for statement 36              | 15 s     | Success
--------------------------------------------------------------
```

# D  Glossary

- **SMOTE:** Synthetic Minority Oversampling Technique used for handling class imbalance.

- **AUROC:** Area Under the Receiver Operating Characteristic Curve, a performance metric for classification models.

- **Fan-In/Fan-Out:** Derived features indicating the count of unique senders or recipients in transactions.

# E  Explanation

- **Undersampling:** The majority class is randomly downsampled to match the size of the minority class, reducing class imbalance.

- **SMOTE:** New synthetic samples are generated for the minority class by interpolating between existing samples, enhancing the model's ability to learn minority class patterns.

- **Final Distribution:** After combining undersampling and SMOTE, the dataset achieves a balanced distribution of the target variable.