# Using predictive modeling to combat money laundering

**Ivan Lopez Torres,** *University of Charleston, ivanlopeztorres@ucwv.edu*

## Abstract

Anti-Money Laundering (AML) is crucial for preventing financial crimes. Predictive modeling techniques can help identify fraudulent transactions. This paper introduces the IBM Synthetic Financial Data Money Laundering dataset containing millions of records with legitimate money laundering transactions. Objectives include exploring effectiveness, improving AML compliance, reducing false positives, and understanding consequences.

## Introduction

Anti-Money Laundering (AML) refers to a set of laws, regulations, and procedures that financial institutions, governments, and other entities use to prevent, detect, and report financial crimes such as money laundering, terrorism financing, and other illicit activities (Tommie, 2023). Financial transactions, client interactions, and other activities that might be used to launder money or finance criminal operations are among the risks that AML systems seek to identify and reduce. For financial institutions and other businesses engaged in financial transactions, ensuring that AML standards are followed is essential since noncompliance carries serious repercussions to someone's reputation, finances, and legal grounds.

The landscape of anti-money laundering (AML) compliance and fraud detection in the financial industry has become increasingly complex and challenging due to factors such as the rapid evolution of financial systems, advancements in technology, the emergence of digital currencies, and the interconnectedness of global economies. According to estimates from the United Nations Office on Drugs and Crime, a significant portion of global GDP, ranging from 2% to 5%, is laundered annually, amounting to approximately $800 billion to $2 trillion USD (UNODC, n.d.). In response to these challenges, government departments and agencies collaborate to safeguard the integrity of financial systems, combat organized crime, and promote economic stability.

With the increase of sophisticated criminal activities and the volume of financial transactions, anti-money laundering (AML) has become a major problem for businesses and financial institutions worldwide. Finding and stopping money laundering has therefore become a top priority. The implications of money laundering extend beyond individual businesses and financial institutions, reaching into the broader context of global financial systems. Money laundering not only facilitates criminal activities but also undermines trust in financial institutions and erodes public confidence in the economy. Therefore, the implementation of effective AML measures is crucial to preserve the integrity and stability of the global financial system (Deloitte, 2020).

Traditional AML compliance methods, such as manual monitoring and rule-based systems, have proven to be insufficient in combating the ever-evolving techniques employed by money launderers. Manual examination of large volumes of data was a cumbersome and time-consuming process, limiting the effectiveness of detection. Financial institutions relied on these traditional methods to identify suspicious activities, such as unusual transaction patterns or high-risk customers (Tommie, 2023). However, machine learning, a subfield of artificial intelligence (AI), has the potential to revolutionize AML compliance for financial institutions and businesses. By leveraging advanced algorithms and vast amounts of data, machine learning models can learn patterns, detect anomalies, and make predictions, leading to more accurate and timely identification of suspicious activities. The utilization of machine learning algorithms has significantly enhanced AML compliance efforts by automating processes, increasing efficiency, and enabling more effective risk management.

By conducting this case analysis, the IBM Synthetic Financial Data Money Laundering dataset was selected based on its characteristics, which include being synthetic, having a large number of records, and containing a balance of legitimate and money laundering transactions.

This dataset contains information on financial transactions between banks and accounts, and is intended for use in anti-money laundering (AML) compliance and fraud detection efforts, it includes details on the amount sent and received in different currencies, the sending and receiving banks and accounts, the payment format used, and whether the transaction is suspected of being involved in money laundering.

The dataset is available in CSV format and has a single file containing over 45 million records.

1. **Timestamp**: The date and time of the transaction in YYYY/MM/DD HH:MM format. Analyzing the temporal aspect of transactions allows us to identify patterns, trends, and potential anomalies related to money laundering activities.
2. **From Bank**: The identifier of the bank that sent the funds. By examining the originating bank, we can investigate any connections to financial institutions that have been associated with money laundering or illicit activities in the past.
3. **Account**: The identifier of the account that sent the funds. Analyzing the transaction history of specific accounts can help identify suspicious behavior or links to known fraudulent activities.
4. **To Bank**: The identifier of the bank that received the funds. Exploring the recipient bank provides insights into potential networks involved in money laundering schemes.
5. **Account**: The identifier of the account that received the funds. Analyzing the receiving account can reveal connections to high-risk individuals or entities involved in illicit financial activities.
6. **Amount Received**: The amount received by the receiving account in the receiving currency. Unusually large or frequent amounts received may indicate suspicious transactions associated with money laundering.
7. **Receiving Currency**: The currency of the amount received. Analyzing transactions involving specific currencies can help identify regions or countries prone to money laundering activities.
8. **Amount Paid**: The amount paid by the sending account in the payment currency. Discrepancies between the amount paid and expected values can be indicative of fraudulent transactions.
9. **Payment Currency**: The currency of the amount paid. Analyzing the payment currency can help identify transactions involving high-risk currencies or currency conversion techniques commonly employed in money laundering.
10. **Payment Format**: The payment format used for the transaction, such as ACH (Automated Clearing House) or credit card. Certain payment formats may be more susceptible to fraudulent activities, and analyzing this variable can help identify patterns associated with money laundering.
11. **Is Laundering**: A binary variable indicating whether the transaction is suspected of being involved

in money laundering, with 1 indicating suspected money laundering and 0 indicating no suspected money laundering. This variable serves as the target variable for our predictive modeling efforts to classify transactions as potentially fraudulent or not.

In this project, we will explore the effectiveness of predictive modeling techniques in identifying potentially fraudulent financial transactions and flagging them for further investigation, some of the challenges and questions that we will try to answer are:

1. Can money laundering or other illegal acts be detected using predictive modeling approaches based on specific patterns or anomalies in financial transactions?
2. How can financial institutions and regulatory organizations use predictive modeling techniques to improve their AML compliance efforts and combat money laundering more effectively?
3. Can predictive modeling be used to reduce false positives in AML compliance efforts and minimize the impact on legitimate financial transactions?
4. What possible consequences could the of use predictive modeling methods have on AML compliance efforts?

There are several machine learning algorithms that can be used for fraud detection. Here are three popular algorithms that are commonly used for fraud detection:

Logistic Regression: This is a commonly used algorithm for binary classification problems like fraud detection. Logistic regression is a statistical model that uses a linear combination of input features to predict a binary output variable. It works well for problems with large datasets and can be used to model the probability of fraud based on the input features (Sqream Technologies, August 11, 2021).

1. **Random Forest**: Random forest is an ensemble learning algorithm that builds multiple decision trees and aggregates their outputs to make a final prediction. It is a popular algorithm for fraud detection due to its ability to handle large datasets with many features, and its ability to detect complex relationships between features (Sqream Technologies, August 11, 2021).

2. **Gradient Boosting**: Gradient boosting is another ensemble learning algorithm that works by building a series of decision trees, with each subsequent tree trained to correct the errors of the previous tree. It is a powerful algorithm for fraud detection because it can learn complex relationships between features and handle large datasets (Sqream Technologies, August 11, 2021).

It's important to keep in mind that the efficiency of these algorithms might change based on the specifics of the dataset and the situation. It's recommended to experiment with different algorithms and evaluate their performance using metrics such as precision, recall, and F1 score to determine the best approach for a given problem.

## Case Overview & Methodology

The case revolves around a comprehensive analysis of IBM's Synthetic Financial Data Money Laundering dataset. This dataset contains details about a series of financial transactions, encompassing variables such as the amount received, the amount paid, the receiving currency, payment currency, and payment format among others. The objective is to recognize potential instances of money laundering activities. The methodology of the study involves various stages:

1. **Data Preparation and Exploration:** As outlined by Garcia et al. (2015), is a crucial step aimed at

preparing the dataset for analysis. It involves data reading, data type specification, and treatment of null values. Our approach depends on Python's powerful libraries to manipulate our dataset effectively. We will begin by ingesting the CSV dataset via pandas, then cleanse the data by specifying accurate data types and handling missing values. Finally, to ensure compatibility with machine learning algorithms, categorical variables are converted into a numerical format using Label Encoding and One-Hot Encoding techniques. This systematic approach ensures our data is ready for subsequent machine learning model development and performance evaluation.

2. **Data Visualization:** In this section we proceed to visualize the data, aiming to uncover hidden patterns and trends. Different plots will be crafted to represent the count of transactions across various variables, segmented by their potential classification as money laundering activities.

3. **Model Construction:** Three models were selected for this study, namely Logistic Regression, Random Forest, and Gradient Boosting Classifier, which were chosen due to their suitability for binary classification problems (James et al., 2013). Logistic Regression is a simple yet robust model that can provide interpretable results and is particularly effective when features are scaled using techniques such as StandardScaler (Cramer, 2002). On the other hand, Random Forest is an ensemble learning method that operates by constructing multiple decision trees at training time and outputting the mode of the classes (Breiman, 2001). It's particularly known for its high accuracy, robustness, and ease of use. Lastly, the Gradient Boosting Classifier is another ensemble model, which constructs new predictors that aim to correct the residual errors of the preceding predictors (Friedman, 2001). It is renowned for its efficacy in minimizing bias and variance, thus improving overall model performance.

4. **Model Evaluation:** The efficacy of each machine learning model in our analysis will be gauged using a series of key performance indicators. These include accuracy, precision, recall, and the F1-score. Accuracy is the simplest measure and represents the proportion of correct predictions among the total number of cases (Japkowicz & Shah, 2011). Precision provides the proportion of true positive predictions among all positive predictions, useful in instances where the cost of false positives is high (Davis & Goadrich, 2006). Recall, also known as sensitivity, measures the proportion of actual positives that were correctly identified, and is particularly helpful in scenarios where the cost of false negatives is high (Davis & Goadrich, 2006). Lastly, the F1-score balances the trade-off between precision and recall by calculating their harmonic mean, providing a singular metric that considers both false positives and false negatives (Chicco & Jurman, 2020).

## Data Preparation and Exploration

It's important to understand the steps involved in preparing data for analysis, in this project we read the CSV file using Python's pandas library and stored it in a variable called aml. We specified the data type for the "Amount Received" and "Amount Paid" columns to be float to ensure that they are processed correctly. Additionally, we specified (null) values to be treated as missing data.

```python
aml = pd.read_csv(os.getcwd() + "\\Transaction Data.csv", dtype={"Amount Received": "float", "Amount Paid": "float"}, na_values=["(null)"])
```

Then, we converted the 'Timestamp' column to datetime format, which allows us to easily perform time-based analyses on the data.

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

These steps ensure that the data is in the correct format and will allow us to perform meaningful analysis to gain insights into the transactional data.

In data analysis, numeric variables play an important role in understanding trends and patterns within the IBM Synthetic Financial Data Money Laundering dataset, these variables can provide valuable insights into the financial transactions recorded. In this summary, we will explore the key statistical measures of the numeric variables, including their mean, standard deviation, minimum and maximum values, and quartile values.

```
        Amount Received    Amount Paid   Is Laundering
count    4.540402e+07     4.540642e+07    4.540693e+07
mean     3.140775e+04     3.071709e+04    1.220585e-03
std      1.106864e+06     1.095479e+06    3.491554e-02
min     -5.260059e+06    -5.260059e+06    0.000000e+00
25%      2.604000e+01     2.609000e+01    0.000000e+00
50%      3.301300e+02     3.314900e+02    0.000000e+00
75%      3.063790e+03     3.066670e+03    0.000000e+00
max      1.139233e+09     1.139233e+09    1.000000e+00
```

- The "Amount Received" and "Amount Paid" columns have a high standard deviation, indicating that the data is highly variable and spread out.
- The minimum "Amount Received" and "Amount Paid" values are both negative 5,260,059, while the maximum values for both columns are over 1 billion.
- The "Is Laundering" column has a low mean value of 0.0012, indicating that only a small percentage of the transactions are classified as potential money laundering activities.

## Data Visualization

Visualizations can be useful in identifying patterns or trends in the data, as well as variables that may be associated with potential money laundering activities. We will present three plots which provide a visual representation of the distribution of transactions across different variables in the IBM Synthetic Financial Data Money Laundering dataset, each plot shows the count of transactions for a specific variable, separated by the "Is Laundering" value, which indicates whether the transaction is classified as a potential money laundering activity or not. The first plot (Figure 1) shows the count of transactions for each "Receiving Currency" value, separated by the "Is Laundering" value.
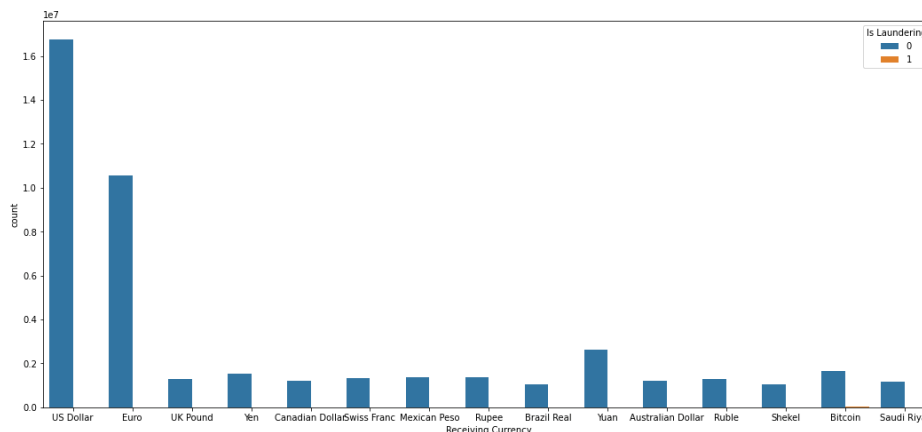


**Figure 1: Count of Transactions by Receiving Currency**

Note: We can see that for most "Receiving Currency" values, the count of transactions classified as potential

money laundering activities ("Is Laundering" = 1) is almost zero, while the count of non-laundering transactions ("Is Laundering" = 0) is much higher. However, for the "Receiving Currency" value Bitcoin, we can see that there is a small quantity of transactions classified as potential money laundering activities. This information could be useful for businesses that deal with Bitcoin transactions.

Figure 2 shows the count of transactions by "Payment Currency and the color segments within each bar indicate the proportion of transactions that were flagged as "laundering".
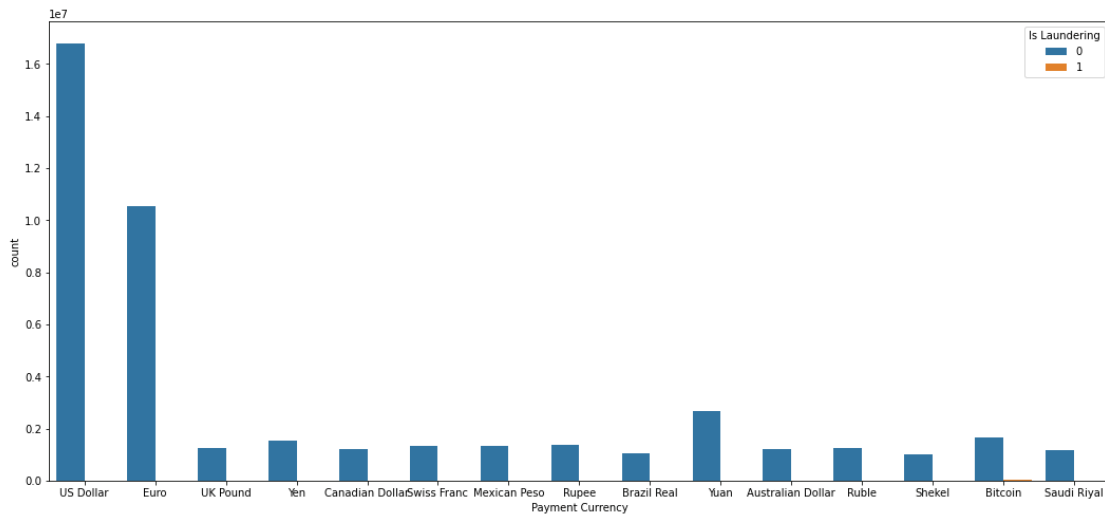
**Figure 2: Count of Transactions by Payment Currency**

Note: The Figure 2 plot highlights that transactions involving Bitcoin as the payment currency have a relatively small quantity of "laundering" transactions compared to other currencies.This plot shows the number of transactions by payment format and whether they are flagged as potential money laundering transactions.
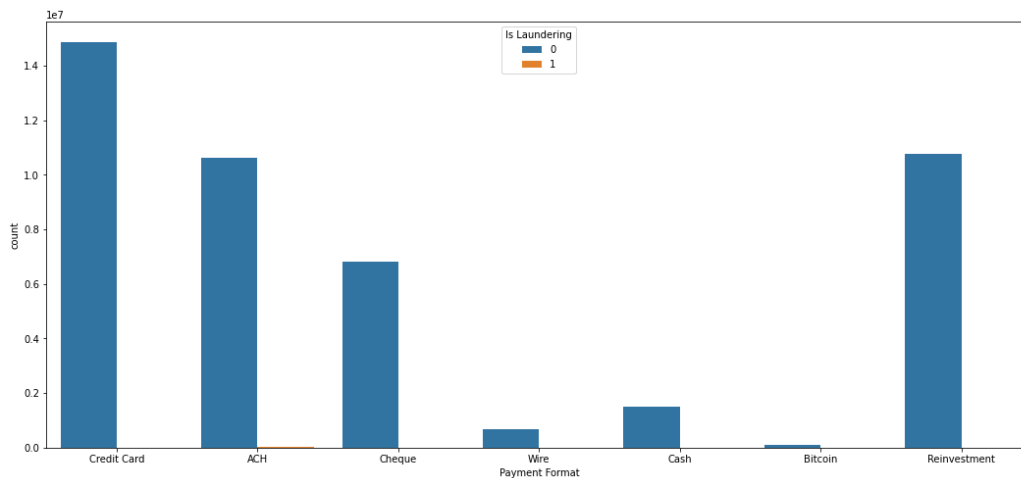
**Figure 3: Count of Transactions by Payment Currency**

Note: The Figure 3 plot shows that the payment format ACH has a small quantity of transactions flagged as money laundering compared to other payment formats, which have almost no transactions flagged as money laundering.

## Model

In this section, we'll be looking at the linear connections between the numerical variables in the dataset. Finding the correlations between various variables and how changes in one variable affect changes in another will be required for this. We can create models that can be utilized to generate predictions about the data by understanding these relationships.

The prediction model will incorporate these variables to capture the link between the variables and the outcome. To mathematically evaluate these relationships in a correlation matrix, we need to ensure all variables are in a numerical form. Thus, we employ the LabelEncoder and One-Hot Encoding techniques to convert categorical variables in our dataset into numerical ones.

First, we will extract the unique values of the "From Bank", "To Bank", "Account", "Account.1", "Receiving Currency", and "Payment Currency" columns in the dataset and merge them into a set of unique values. The LabelEncoder then fits and transforms the unique values of these columns into numerical labels, which is a common technique used to encode categorical variables into numerical labels.

```python
# Get the unique values
unique_values_banks = set(aml["From Bank"].unique()).union(set(aml["To Bank"].unique()))
unique_values_accounts = set(aml["Account"].unique()).union(set(aml["Account.1"].unique()))
unique_values_currencies = set(aml["Receiving Currency"].unique()).union(set(aml["Payment Currency"].unique()))

le = LabelEncoder()
le.fit(list(unique_values_banks) + list(unique_values_accounts) + list(unique_values_currencies))

aml["From Bank"] = le.transform(aml["From Bank"])
aml["To Bank"] = le.transform(aml["To Bank"])
aml["Account"] = le.transform(aml["Account"])
aml["Account.1"] = le.transform(aml["Account.1"])
aml["Receiving Currency"] = le.transform(aml["Receiving Currency"])
aml["Payment Currency"] = le.transform(aml["Payment Currency"])
```

Next, the "Payment Format" column is converted into a categorical column which will be transformed into a binary column using the pd.get_dummies function. This technique is called One-Hot Encoding and is used to convert categorical variables into binary columns where each unique value of the categorical variable becomes a separate binary column.

```python
aml["Payment Format"] = aml["Payment Format"].astype("category")

# Convert column to binary column
aml = pd.get_dummies(aml, columns=["Payment Format"])
```

Once we've prepared the dataset, we remove the "Timestamp" column and any rows containing missing values in the "Amount Received" and "Amount Paid" columns.

```python
final_aml = aml.copy()

final_aml = final_aml.drop(["Timestamp"], axis=1)
final_aml = final_aml.dropna(subset=["Amount Received", "Amount Paid"])
```

After this data cleaning and preprocessing step, we can create a correlation matrix, this matrix displays the correlation coefficient, which ranges from -1 to 1, for each variable pair. A correlation coefficient of 1 indicates a strong positive relationship, meaning that as one variable increases, the other variable also increases. A correlation coefficient of -1 indicates a strong negative relationship, meaning that as one variable increases, the other variable decreases. A correlation coefficient of 0 indicates no relationship between the two variables.
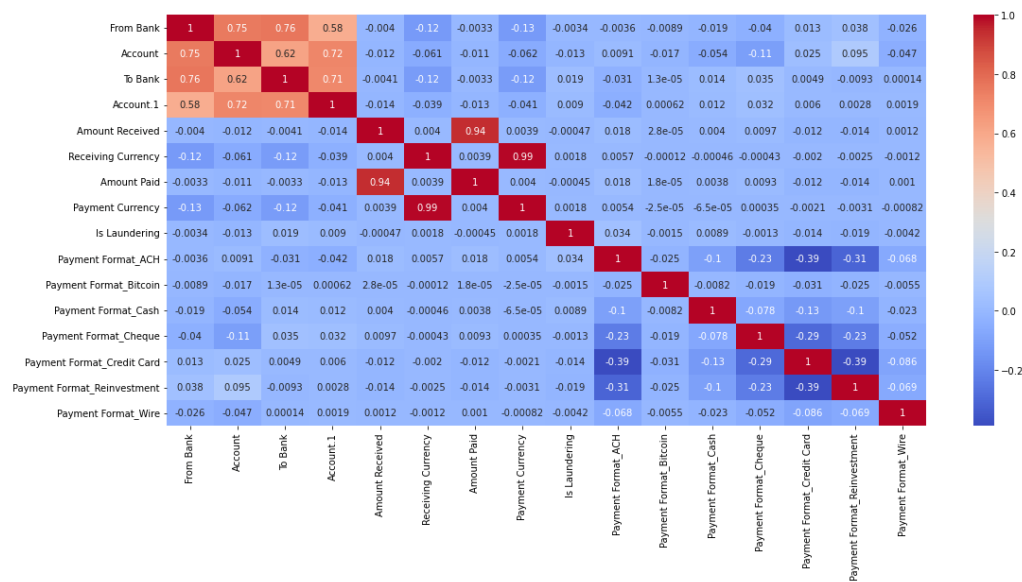
**Figure 4: Correlation Matrix**

Looking at the correlation matrix, we can interpret a few relationships:

1. **Strong Positive Correlation**: The variables 'From Bank', 'Account', 'To Bank', and 'Account.1' have strong positive correlations with each other. This means that if the value of one of these variables increases, the value of the others also tends to increase. The correlation between 'From Bank' and 'Account' is 0.75, between 'From Bank' and 'To Bank' is 0.76, and between 'Account' and 'To Bank' is 0.62.

2. **Positive Correlation with Payment Formats**: There are also some weak positive correlations between some of the payment format variables and the 'From Bank', 'Account', 'To Bank' and 'Account.1' variables. This suggests that there might be some relationships between these variables as well, though these are weaker.

3. **Weak Correlation with 'Is Laundering'**: The correlations between 'Is Laundering' and the other variables are relatively weak. The highest correlation is 0.02 with 'To Bank'. This suggests that these variables might not have a strong influence on whether a transaction is considered money laundering.

4. **High Correlation between 'Receiving Currency' and 'Payment Currency'**: The correlation between these two variables is very high (0.991055), suggesting that they often change together. This is expected, as the type of currency received is likely to be the same as the type of currency paid in most transactions.

5. **High Correlation between 'Amount Received' and 'Amount Paid'**: The correlation between these two variables is also very high (0.941515), which makes sense since the amount paid in a transaction is likely to be close to the amount received, minus any fees or other deductions.

6. **Correlations among Payment Formats**: Different payment formats show negative correlation with each other. This is because if one payment format is used, others are not. For example, if 'Payment Format_ACH' is used (has value 1), the other payment format variables ('Payment Format_Bitcoin', 'Payment Format_Cash', etc.) would likely have values of 0, and vice versa.

Remember that correlation does not imply causation. So, while these correlations can provide us with useful insights, they do not necessarily mean that one variable's change directly causes a change in another variable.

After preparing and understanding our data, we divide the dataset into a training set and a testing set, using the train_test_split function. This is crucial to avoid overfitting and to ensure our model generalizes well on unseen data. Our dataset, final_aml, containing all features, is split with 70% data used for training and 30% reserved for testing. After this step, we will have four datasets: X_train, y_train, X_test, and y_test, which we will use to train and test our machine learning models.

```python
# Feature selection
X = final_aml.drop(["Is Laundering"], axis=1)
y = final_aml["Is Laundering"]

# Split dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

As mentioned in the introduction, we will create a logistic regression model to predict whether a transaction is a case of money laundering or not. Before training the model, we first scaled the numerical features using the StandardScaler, this is important because logistic regression is sensitive to the scale of the features, and scaling the features can help improve the performance of the model.

```python
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_train_scaled = pd.DataFrame(X_train_scaled, columns=X_train.columns)
```

Next, we initialized a logistic regression model with a maximum of 1000 iterations and fit it to the scaled training data, finally, we used the trained model to make predictions on the scaled testing data using the predict method and store the predictions in the lr_preds variable.

```python
# Now we can train our logistic regression model on the scaled data
lr = LogisticRegression(max_iter=1000)
lr.fit(X=X_train_scaled, y=y_train)

X_test_scaled = scaler.transform(X_test)
X_test_scaled = pd.DataFrame(X_test_scaled, columns=X_test.columns)

lr_preds = lr.predict(X_test_scaled)
```

The next model to evaluate will be a random forest model, this model is trained using the RandomForestClassifier class, after the model is trained, the predict method is called on the fitted model object with the testing set X_test as input argument, the predicted values are stored in the rf_preds variable.

```python
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
rf_preds = rf.predict(X_test)
```

Finally, the last model to evaluate will be the Gradient Boosting Classifier, here, we are creating an instance of the GradientBoostingClassifier class and then fitting the model on our training data, we are then using the trained model to make predictions on our testing data using the predict method and storing the predictions in the gb_preds variable. Gradient Boosting is a type of ensemble learning algorithm that combines multiple weak learners to create a strong learner. It works by sequentially adding new models to the ensemble, with each new model trained to correct the errors of the previous models.

```
gb = GradientBoostingClassifier()
gb.fit(X_train, y_train)
gb_preds = gb.predict(X_test)
```

## Results & Findings

After evaluating three machine learning models: Logistic Regression, Random Forest, and Gradient Boosting Classifier. Each model was assessed based on its accuracy, precision, recall, and F1-score. The performance metrics of each model are summarized in the table below:

**Table 1: Performance Metrics of Machine Learning Models for Money Laundering Detection**

| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| Logistic Regression | 99.88% | 100% | 0% | 0% |
| Random Forest | 99.89% | 56.33% | 27.15% | 36.64% |
| Gradient Boosting | 99.88% | 35.82% | 0.14% | 0.29% |

The Logistic Regression model demonstrated exceptional accuracy, correctly predicting the nature of 99.88% of the transactions. This accuracy indicates the model's strength in general predictions, as it could correctly classify the majority of the transactions. However, its primary limitation was reflected in the recall score of 0.0, meaning it failed to identify any money laundering transactions. Despite its high overall accuracy, the lack of sensitivity towards money laundering instances poses a significant challenge, rendering this model unsuitable for our specific purpose.

```
Logistic Regression Results
Accuracy: 0.9987847487844551
Precision: 1.0
Recall: 0.0
F1 Score: 0.0
```

Maintaining a high accuracy rate, the Random Forest model displayed superior precision in identifying money laundering transactions. It managed to correctly label 56.33% of the actual money laundering instances. A key insight from the Random Forest model is the recognition of some money laundering transactions, although it failed to identify a significant proportion of them, as indicated by a recall score of 27.15%. Thus, it managed to spot 27.15% of all money laundering activities, neglecting 72.85%. The F1 score, which balances both precision and recall, stood at 36.64% suggesting a need for model refinement to improve its performance.

```
Random Forest Results
Accuracy: 0.9988588987106135
Precision: 0.5632990724492354
Recall: 0.2714915725246179
F1 Score: 0.3663935428641311
```

The Gradient Boosting Classifier model demonstrated high accuracy, similar to the Logistic Regression model. However, the precision and recall scores were drastically low compared to the previous two models, with precision at 35.82% and recall at 0.14%. This indicates that although the model accurately identified

a large proportion of the transactions, it was significantly lacking in the detection of money laundering instances. Its low F1 score of 0.29% further highlights the limitations of the Gradient Boosting Classifier model in detecting money laundering activities)

```
Gradient Boosting Results
Accuracy: 0.9987833538848542
Precision: 0.3582089552238806
Recall: 0.001449888237781671
F1 Score: 0.0028880866425992774
```

The visualization of the data for Bitcoin transactions reveals an intriguing trend. The charts show that Bitcoin is a currency that is frequently used in possible money laundering operations. This suggests a pattern where Bitcoin transactions are more likely to be associated with fraudulent activity. As a result of Bitcoin transactions' more vulnerability to money laundering, this information could actually motivate companies and regulatory bodies to monitor them more closely and put strong anti-fraud measures in effect.

Also, the data analysis reveals that transactions conducted through the ACH payment format bear a comparatively higher incidence of flags for potential money laundering. This indicates that this particular payment method might be exploited more frequently for illicit activities.

Finally, the correlation analysis uncovers an intriguing trend among certain pairs of variables in the dataset. There's a strong positive correlation between "From Bank" and "Account", "From Bank" and "To Bank", and "Account" and "To Bank". This implies a strong interdependence between these pairs of variables; an increase in one variable is likely to correspond to an increase in the other. Since these variables move together, they could be used interchangeably to reduce the dataset's dimensionality without losing much information, or they might influence the creation of new, more informative features, this information can be crucial for feature selection during the training of predictive models. It's important to note that these correlations can be a sign of multicollinearity, which might affect the interpretability of certain models, such as logistic regression.

## Conclusions & Implications

This study explores the potential of predictive modeling techniques, particularly the Random Forest algorithm, in identifying money laundering and enhancing Anti-Money Laundering (AML) compliance. The techniques applied proved to be an effective tool in flagging potentially fraudulent financial transactions associated with money laundering. By identifying specific patterns and anomalies in financial transactions, these models provide early warnings for suspicious behavior, enabling proactive action.

Compared to existing literature and industry practices, our findings corroborate the efficacy of machine learning algorithms in detecting fraudulent activities, which aligns with studies like Bolton & Hand's (2002). However, the absence of research particularly addressing the connection between machine learning and the detection of money laundering increases the importance of our findings and suggestions.

Our research revealed significant findings, highlighting the susceptibility of Bitcoin transactions and the Automated Clearing House (ACH) payment format to potential money laundering activities. This revelation not only reinforces Reid & Harrigan's (2013) examination of potential misuse of cryptocurrencies, but also intensifies the call for stricter controls for transactions involving these mediums.

Despite these important findings, we encountered some challenges such as a significant imbalance in our dataset and weak correlations between certain variables. To enhance the accuracy and reliability of predictions, we recommend diversifying datasets and considering other classification models alongside the Random Forest algorithm.

The implementation of predictive modeling, specifically the Random Forest algorithm, in AML compliance efforts provides significant operational advantages. These include automated detection of suspicious transactions, efficient resource allocation, identification of high-risk individuals or entities, and a reduction in false positives. This approach can minimize disruption to legitimate financial activities while enhancing the precision in detecting potential money laundering cases. However, while applying these predictive modeling methods, organizations must ensure they adhere to standards of fairness, transparency, and regulatory compliance to maintain public trust and preserve the integrity of AML compliance efforts (Zerilli et al., 2019). Addressing potential concerns such as algorithmic bias, privacy issues, and the impact on legitimate financial transactions is integral to the successful implementation of these techniques.

In conclusion, our study reinforces the effectiveness of predictive modeling techniques, like the Random Forest algorithm, in enhancing AML compliance efforts. It aligns with existing practices while offering unique insights into areas susceptible to money laundering. Incorporating these findings into existing practices can significantly bolster efforts by financial institutions and regulatory bodies to combat money laundering.

## References

Breiman, L. (2001). Random Forests. Machine Learning.

Cramer, J. S. (2002). The origins of logistic regression. Tinbergen Institute.

Deloitte. (2020). Anti-Money Laundering Preparedness Survey Report 2020. Retrieved from
        https://www2.deloitte.com/content/dam/Deloitte/in/Documents/finance/Forensic/in-forensic-
        AML-Survey-report-2020-noexp.pdf

Garcia, S., Luengo, J., & Herrera, F. (2015). Data Preprocessing in Data Mining. Springer.

IBM. (2021). IBM AML Data. GitHub. Retrieved from https://github.com/IBM/AML-Data

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning.
        Springer.

Sqream Technologies. (2022, July 21). Best Machine Learning Algorithms for Fraud Detection. Retrieved
        from https://sqream.com/blog/fraud-detection-machine-learning/

Tommie,, C. (2023) From Data to *From Data to Solution: A Practical Machine Learning Notebook for
        Anit-Money Laundering Compliance*. Amazon Publishing.
        https://www.amazon.ca/gp/product/B0C4PVGB6T/ref=dbs_a_def_rwt_hsch_vapi_tkin_p1_i0

United Nations Office on Drugs and Crime. (n.d.). Money Laundering. Retrieved from
        https://www.unodc.org/unodc/en/money-laundering/overview.html