# INDEX

DATA COLLECTION

DATA PREPROCESSING

FEATURE EXTRACTION AND SELECTION

MODEL BUILDING

CLASSIFICATION

RESULT EVALUATION

DASHBOARD

CONCLUSION

# PROBLEM STATEMENT

Our project aims to develop a robust machine learning model for identifying online payment fraud. We will utilize a comprehensive dataset from Kaggle, which contains historical transaction records encompassing both fraudulent and non-fraudulent payments. This rich dataset will serve as the foundation for training our model to effectively distinguish between legitimate transactions and potential fraud attempts.

# FEATURES

**STEP**
- Represents a unit of Time where 1 step equals 1 hour.

**TYPE**
- Type of Online Transaction.

**AMOUNT**
- The Amount of the Transaction.

**NAMEORG**
- Customer initiating the transaction.

**OLDBALORG**
- Balance of the initiator before the transaction.

**NEW BALORG**
- Balance of the initiator after the transaction.

**NAMEDEST**
- Recipient of the transaction.

**OLDBALDEST**
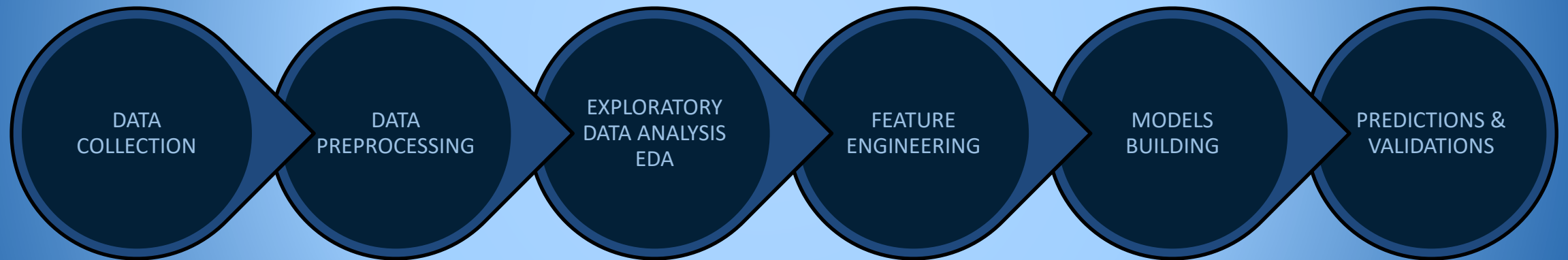- Balance of the recipient before the transaction.

**NEWBALANCE DEST**
- Balance of the recipient after the transaction.

**ISFRAUD**

• Target variable which indicates whether the transaction is Fraudulent or Not.

# METHODOLOGY

DATA COLLECTION

DATA PREPROCESSING

EXPLORATORY DATA ANALYSIS EDA

FEATURE ENGINEERING

MODELS BUILDING

PREDICTIONS & VALIDATIONS

# DATA COLLECTION

The Online Payment Fraud Detection data chosen for this project is taken from Kaggle. We used pandas read_csv() function to load the data to the notebook. This dataset comes with both categorical and numerical data which has been cleaned and processed to build the model.

Sample view -

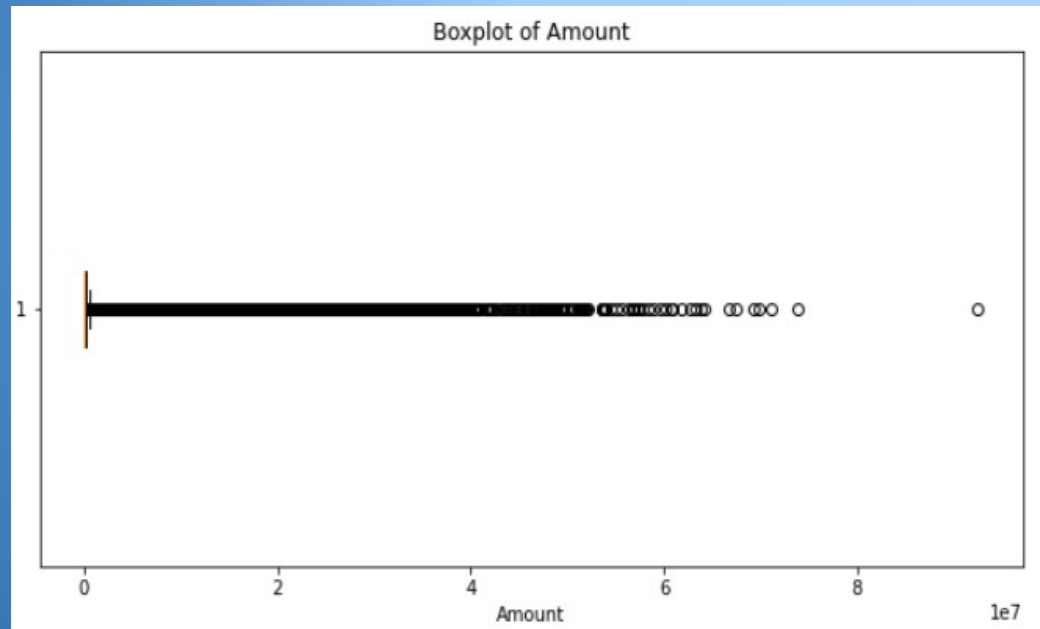| | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | nameDest | oldbalanceDest | newbalanceDest | isFraud | isFlaggedFraud |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 3200269 | 249 | CASH_IN | 73948.23 | C732203687 | 2708528.97 | 2782477.20 | C324416692 | 1940767.30 | 1866819.07 | 0 | 0 |
| 1347037 | 137 | CASH_OUT | 274511.15 | C73247380 | 0.00 | 0.00 | C1262048943 | 4544741.36 | 4819252.52 | 0 | 0 |
| 6211484 | 588 | PAYMENT | 2996.62 | C1329473655 | 287283.00 | 284286.38 | M394994211 | 0.00 | 0.00 | 0 | 0 |
| 954826 | 44 | TRANSFER | 186809.10 | C967821252 | 117181.00 | 0.00 | C1082076427 | 127791.20 | 314600.29 | 0 | 0 |
| 3702054 | 277 | PAYMENT | 6071.28 | C651810808 | 0.00 | 0.00 | M1667789905 | 0.00 | 0.00 | 0 | 0 |

Source  Data Source

# DATA PREPROCESSING

1. Implemented distinct imputation strategies for balance-related variables.
2. Addressed missing or inconsistent values in:
    -> oldbalanceOrg
    -> newbalanceOrig
    -> oldbalanceDest
    -> newbalanceDest
3. Ensured data integrity while preserving the unique characteristics of each transaction type.
4. Time-based Analysis
    -> Created 'stepbucket' feature for temporal aggregation.
    -> Visualized transaction patterns over time using line plots.
    -> Enabled identification of potential fraud trends or anomalies across different time periods.
5. Visualized amount distribution via box plots, identified outliers and applied appropriate handling techniques to ensure robust model performance.
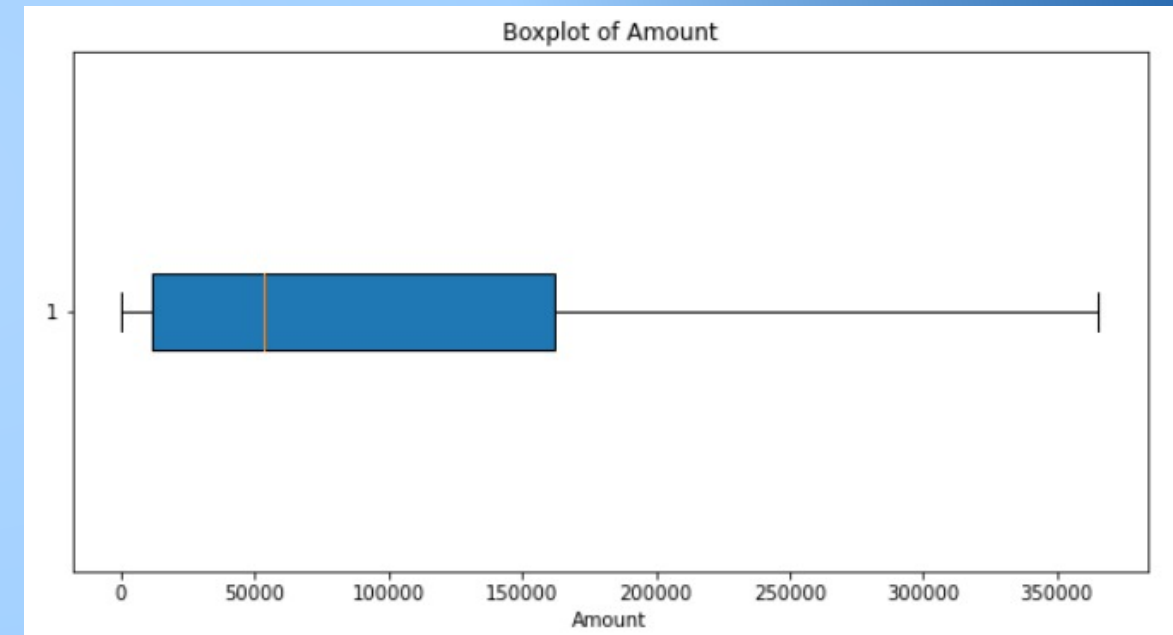
# EXPLORATORY DATA ANALYSIS

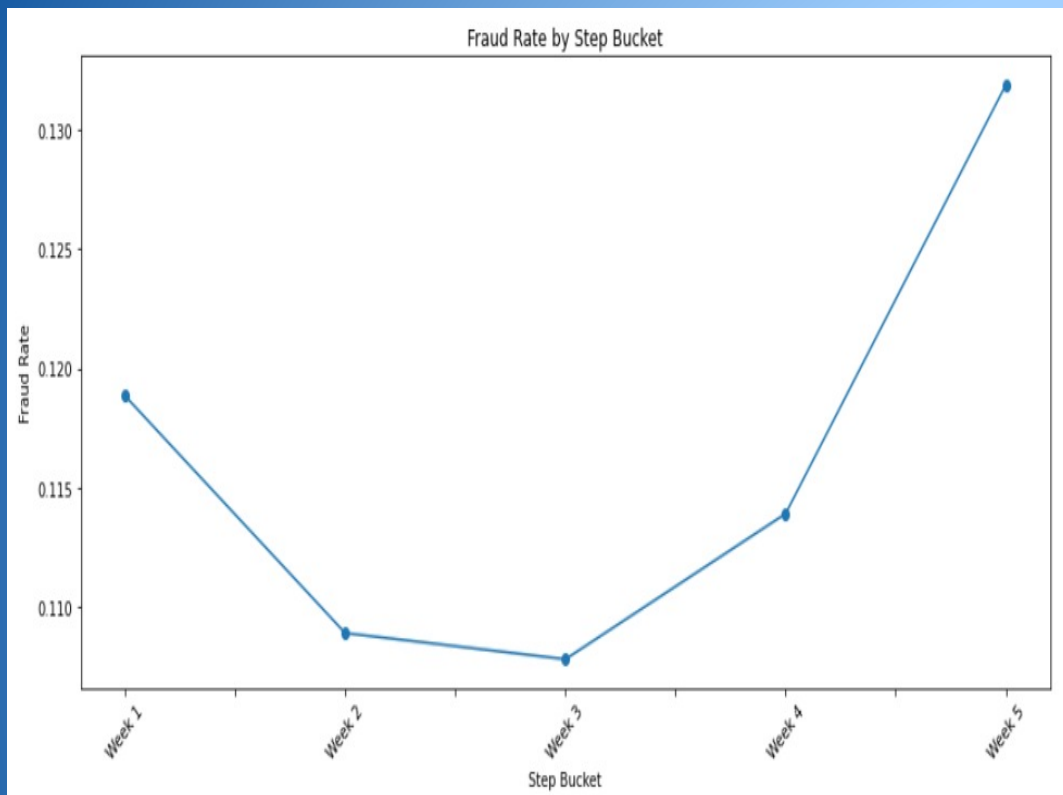Distribution of transaction amounts using Boxplot.
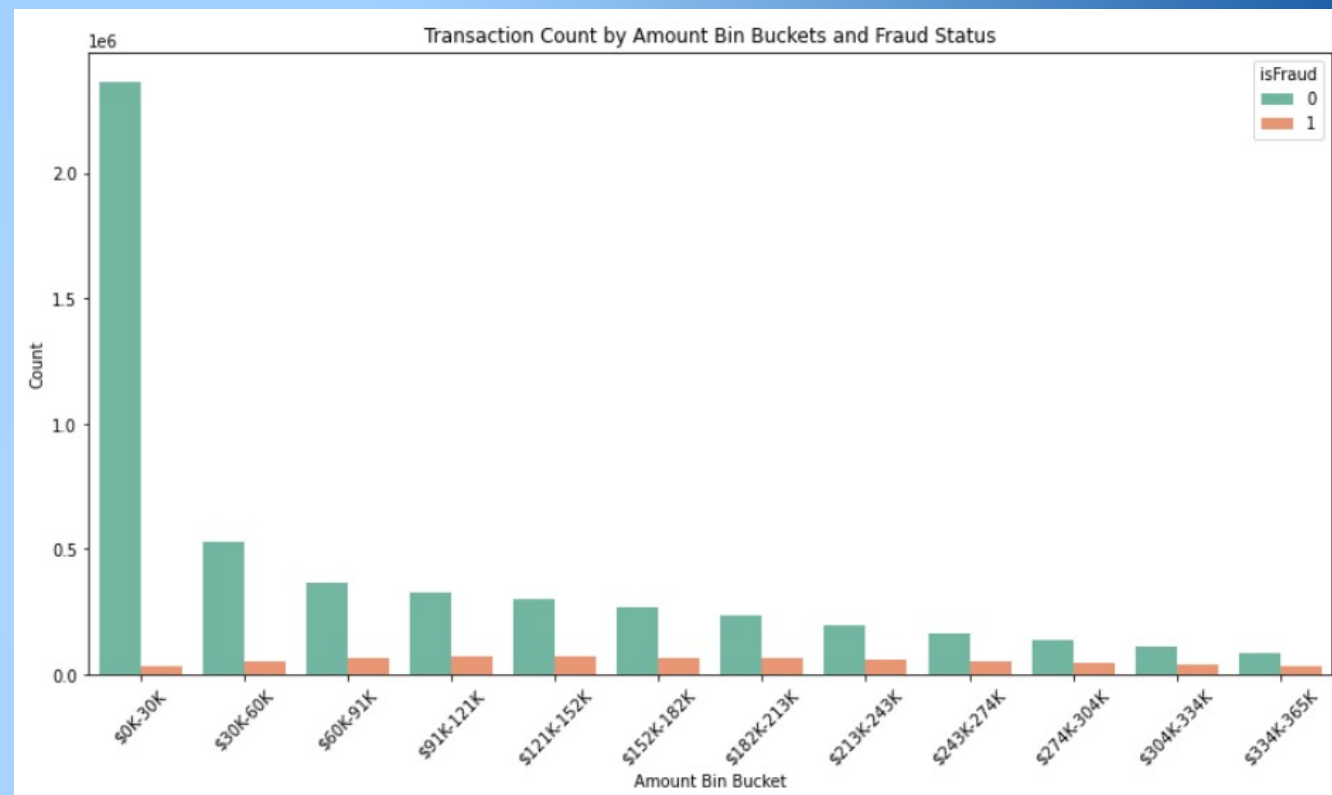
Before Outlier treatment -

After Outlier treatment -
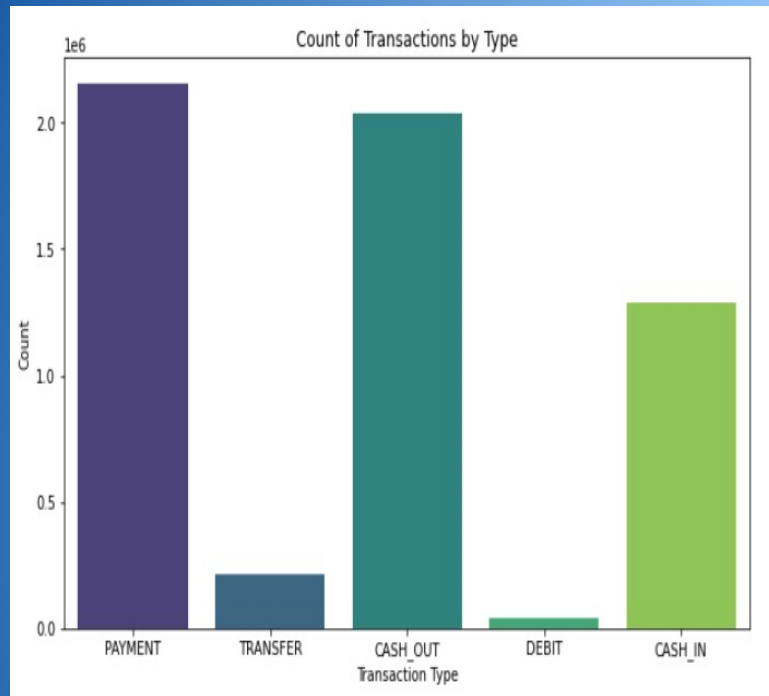
# Fraud Rate by Step Bucket

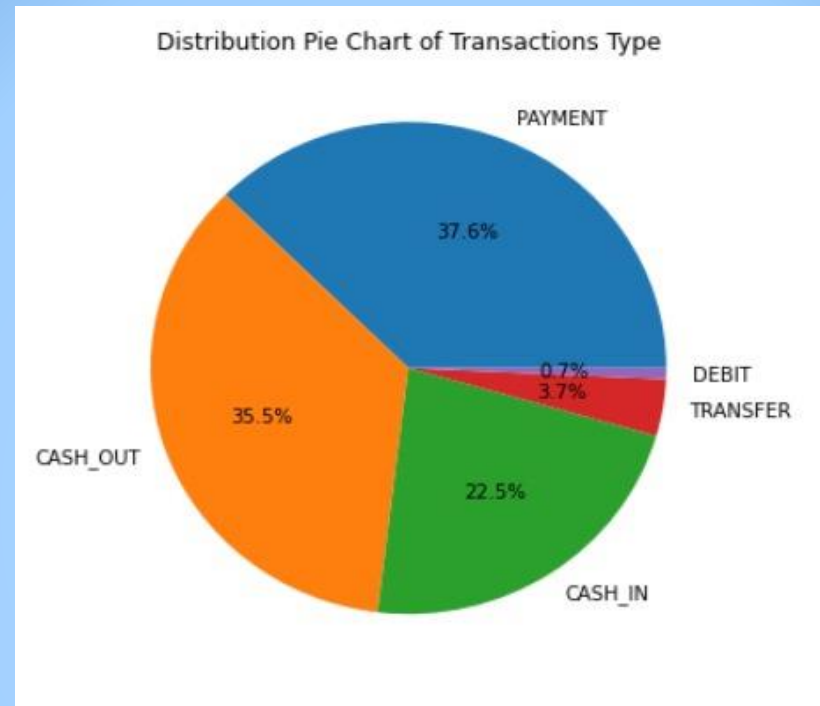# Transaction Count by Amount Bin and Fraud Status





- The sharp increase in fraud rate during Weeks 4 and 5 warrants immediate investigation, as it could indicate evolving fraud tactics or reduced security measures towards the end of the period.
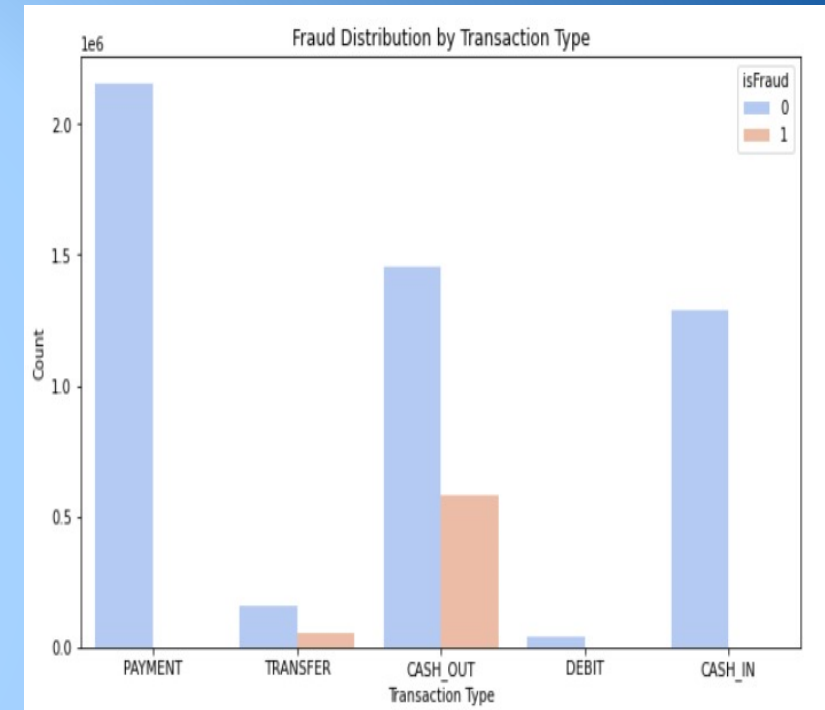- Lower transaction amounts (<$30K) show highest frequency, but fraud risk increases in higher amount brackets.

# Transaction Count by Type
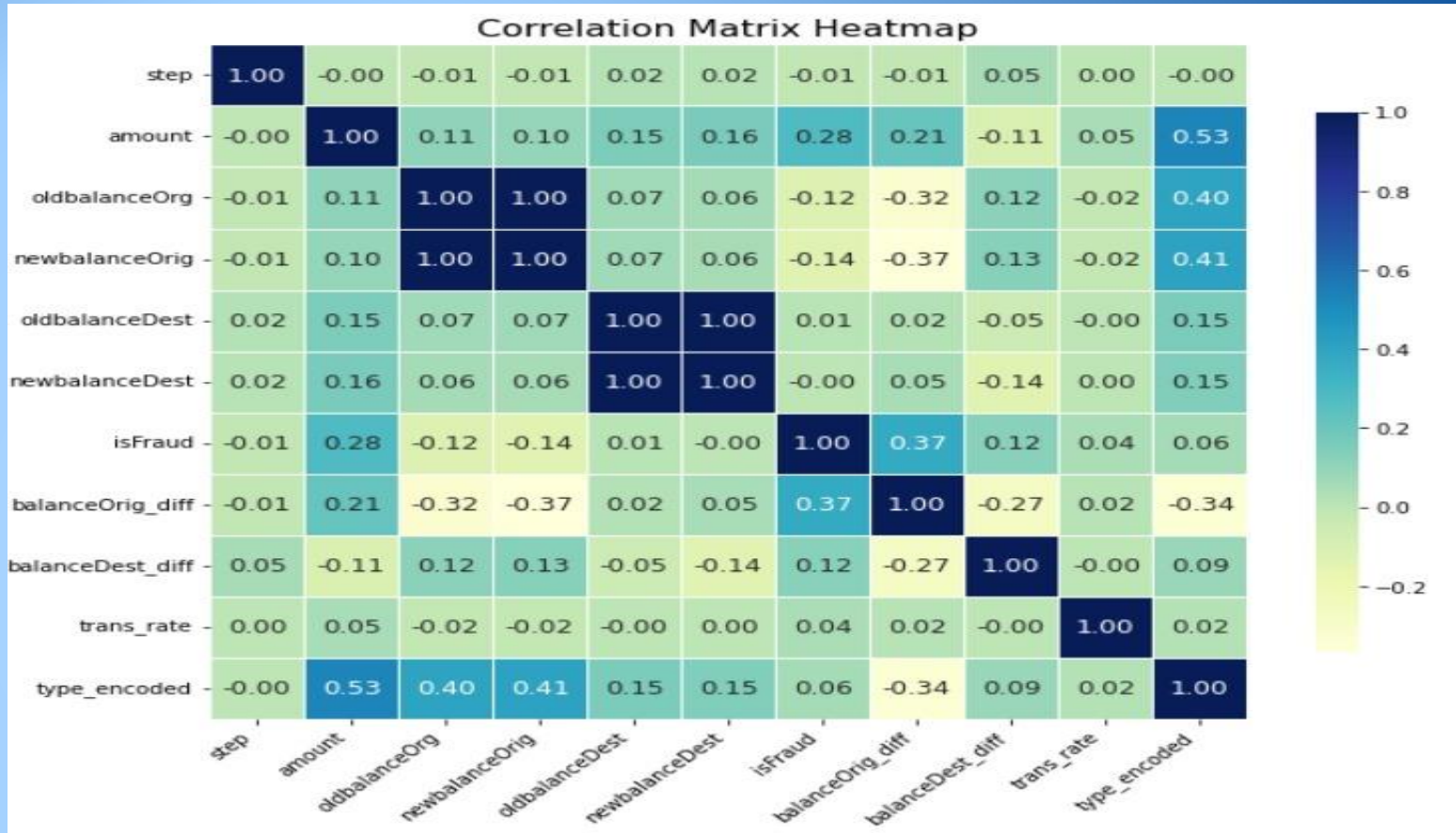
# Distribution of Transaction Type

# Fraud Dist. of Transaction Type



- Analysis reveals PAYMENT and CASH_OUT as the most frequent transaction types, while TRANSFER occurs least often.

- PAYMENT (37.6%) and CASH_OUT (35.5%) dominate, comprising over 70% of all transactions.

- CASH_OUT transactions show highest fraud incidence, indicating key area for model focus.

**CORRELATION MATRIX HEATMAP :-**



Correlation Matrix Heatmap

• The variables oldbalanceOrg with newbalanceOrig, and oldbalanceDest with newbalanceDest, exhibit perfect correlations (1.0). This suggests the potential for feature reduction by selecting one variable from each pair.

• Moderate positive correlations are observed between isFraud and both balanceOrig_diff (0.37) and amount (0.28), indicating that changes in the originator's balance and transaction amounts can be significant indicators of fraudulent activity.

• Type_encoded has a moderate correlation (0.53) with amount, suggesting larger sums are involved in certain transaction types. The step variable shows negligible correlation with other features, indicating fraud patterns are not strongly time-dependent.

# FEATURE ENGINEERING

- Created balanceOrig_diff and balanceDest_diff to capture differences between old and new balances for transaction originators and recipients, respectively.

- Introduced trans_rate to represent the ratio of the transaction amount to the originator's old balance, with safeguards against division by zero.

- Numerically encoded the type feature: PAYMENT (0), TRANSFER (1), CASH_OUT (2), DEBIT (3), CASH_IN (4).

- Used a correlation heatmap to identify and retain relevant features while eliminating highly correlated ones to reduce redundancy and focus on those showing significant correlation with fraud indicators.

- Chose features including amount, oldbalanceOrg, oldbalanceDest, balanceOrig_diff, balanceDest_diff, trans_rate, and type_encoded, while excluding the highly correlated newbalanceOrig and newbalanceDest to streamline the model.
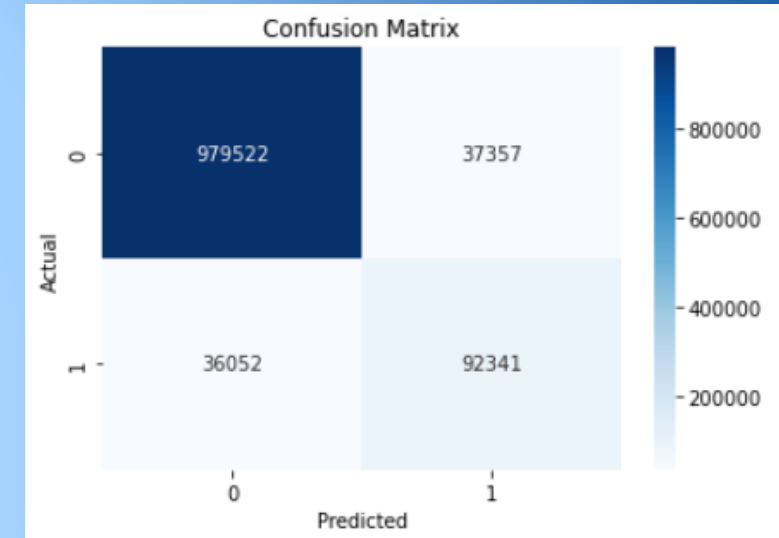
# MODEL BUILDING

- **Data Preparation**: Split data into train and test sets using train_test_split().

- **Model Implementation**: Utilized Decision Tree Classifier and LightGBM Classifier Models.

- **Evaluation Metrics**: Classification Report, Confusion Matrix & ROC AUC Score.

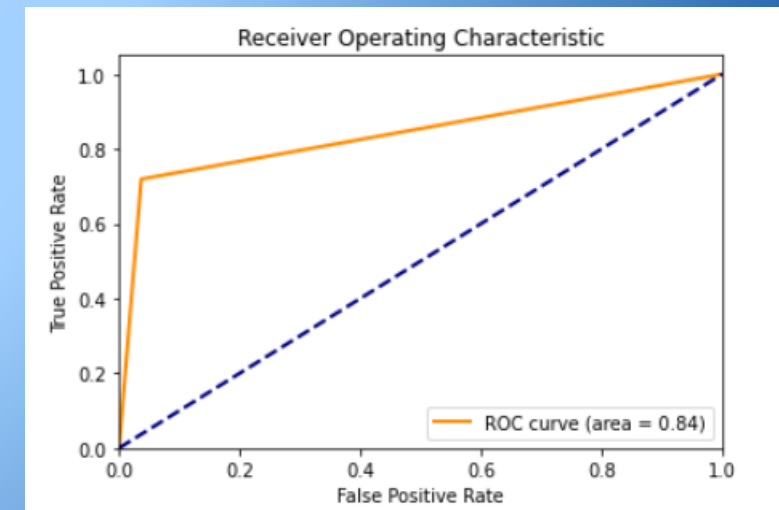- **Model Validation**: Plotted ROC AUC Curve and Visualized Confusion Matrix

# PREDICTIONS AND VALIDATION

**DECISION TREE CLASSIFIER :-**

```
Classification Report:
              precision    recall  f1-score   support

           0       0.96      0.96      0.96   1016879
           1       0.71      0.72      0.72    128393

    accuracy                           0.94   1145272
   macro avg       0.84      0.84      0.84   1145272
weighted avg       0.94      0.94      0.94   1145272
```



Confusion Matrix

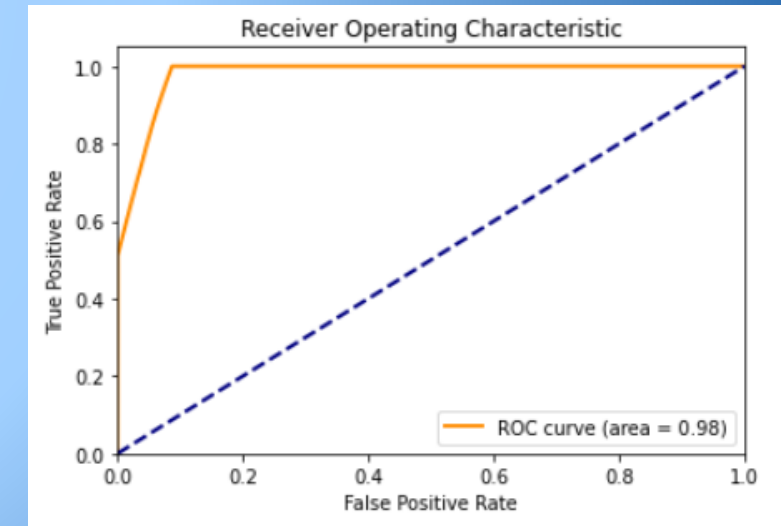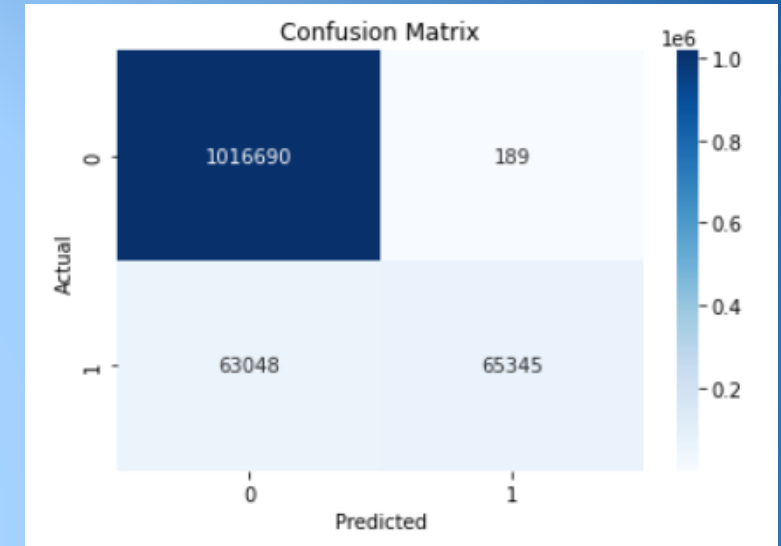

Receiver Operating Characteristic

- **Classification Report:** The model shows high accuracy (0.94) and performs exceptionally well for class 0 (0.96 precision, recall, and f1-score).

- **Confusion Matrix:** The model correctly predicts a large number of true positives and true negatives, with relatively few misclassifications.

- **ROC Curve:** The ROC curve shows strong model performance with an area under the curve (AUC) of 0.84, indicating good discrimination ability.

**LIGHT GBM CLASSIFIER :-**

```
Classification Report:
              precision    recall  f1-score   support

           0       0.94      1.00      0.97   1016879
           1       1.00      0.51      0.67    128393

    accuracy                           0.94   1145272
   macro avg       0.97      0.75      0.82   1145272
weighted avg       0.95      0.94      0.94   1145272
```



Confusion Matrix



Receiver Operating Characteristic

- **Classification Report:** The model shows high precision and recall for class 0, but lower recall for class 1, resulting in good overall accuracy of 0.94.

- **Confusion Matrix:** The model correctly predicts a large number of class 0 instances (1,016,690), but has more difficulty with class 1, showing some misclassifications.

- **ROC Curve:** With an area under the curve of 0.98, the model demonstrates excellent discriminative ability between the two classes.

# CONCLUSION

- ✓ 57,26,358 data points with 7 columns have been processed to build the classification models.

- ✓ Fraud Rates peaked in week 4 and week 5.

- ✓ Lower transaction amounts (<$30K) show highest frequency, but fraud risk increases in higher amount brackets.

- ✓ Analysis reveals PAYMENT and CASH_OUT as the most frequent transaction types, while TRANSFER occurs least often.

- ✓ PAYMENT (37.6%) and CASH_OUT (35.5%) dominate, comprising over 70% of all transactions.

- ✓ CASH_OUT transactions show highest fraud incidence, indicating key area for model focus.

- ✓ Feature reduction is achieved by addressing perfect correlations between balance pairs. Key fraud indicators are changes in the originator's balance and transaction amounts, showing moderate positive correlations with isFraud.

- ✓ Larger amounts are associated with specific transaction types, as indicated by the moderate correlation between type_encoded and amount. Fraud patterns exhibit no strong time dependency, indicating consistent occurrence throughout the observed period.

- ✓ With a ROC AUC score of 0.98 and overall accuracy of 94%, we can conclude that our model Light GBM Classifier demonstrates strong performance in detecting fraudulent transactions, though there's room for improvement in recall for the fraud class.

# POWER BI DASHBOARD

# THANK YOU