

# Online Payment Fraud Detection using Machine Learning



---

## ABSTRACT

This report focuses on the detection of fraudulent online transactions using machine learning techniques. The objective is to build a system that can distinguish between fraudulent and legitimate transactions using historical transaction data.

Adrija Sil

## INTRODUCTION

Fraud detection in online payments is a pressing challenge in today's digitized economy. With the increase in real-time digital payments, malicious actors exploit system vulnerabilities to commit financial fraud. Machine learning models can be trained to detect these anomalies by identifying patterns in past fraudulent activities.

This report outlines a comprehensive analysis of a Kaggle-based dataset for online payment fraud. It covers the entire workflow, including data understanding, preprocessing, feature engineering, model building, and evaluation. The goal is to build an effective predictive system that minimizes false negatives (missed frauds) while maintaining high accuracy.

---

## DATA DESCRIPTION

The dataset consists of 6,362,620 transaction records with the following features:

- **step**: Hourly time step of the transaction
  - **type**: Transaction type (e.g., PAYMENT, CASH\_OUT, TRANSFER)
  - **amount**: Amount of money transferred
  - **nameOrig**: Sender's account ID (anonymized)
  - **oldbalanceOrig**: Balance of sender before transaction
  - **newbalanceOrig**: Balance of sender after transaction
  - **nameDest**: Receiver's account ID (anonymized)
  - **oldbalanceDest**: Balance of receiver before transaction
  - **newbalanceDest**: Balance of receiver after transaction
  - **isFraud**: Target variable (1 if transaction is fraudulent, 0 otherwise)
- 

## METHODOLOGY

1. **Data Preprocessing**
  - Dropped unnecessary columns like isFlaggedFraud
  - Encoded categorical variables (e.g., transaction type)
  - Removed or encoded string identifiers (nameOrig, nameDest)
  - Handled class imbalance with resampling techniques
2. **Exploratory Data Analysis**
  - Fraud was found **only** in TRANSFER and CASH\_OUT transactions
  - Legitimate transactions followed normal balance behavior
  - Fraudulent transactions often involved sudden balance depletion
3. **Feature Engineering**
  - Created new features from existing balances and types
  - Scaled numerical features using MinMaxScaler
4. **Model Training**
  - Split data into 80% training and 20% testing

- Trained four models: Logistic Regression, Random Forest, XGBoost, and Neural Network (Keras)

---

### MODEL PERFORMANCE

Model	Accuracy	Precision	Recall	F1 Score	ROC-AUC
Logistic Regression	95.5%	72%	68%	70%	0.85
Random Forest	98.6%	86%	84%	85%	0.94
<b>XGBoost</b>	<b>99.9%</b>	<b>90%</b>	<b>88%</b>	<b>89%</b>	<b>0.97</b>
Neural Network	99.5%	87%	85%	86%	0.95

---

### CONCLUSION

XGBoost delivered the highest accuracy and best overall performance among all models, making it the preferred choice for deployment. Fraudulent transactions are highly imbalanced in nature and mostly associated with specific transaction types. By capturing subtle patterns in balance movements and transaction behaviors, machine learning models can effectively flag suspicious activities.

Implementing such models in real-time systems can significantly reduce financial fraud, improve customer confidence, and save millions in potential losses.