

# **Fraud Detection System for Financial Transactions – Project Report**

## **1. Project Overview**

Financial institutions need to detect fraudulent transactions in real-time to protect users and maintain trust. In this project, we developed a machine learning-based fraud detection system using historical transaction data.

The final solution includes data analysis, model training, real-time prediction, and a simple user interface built with Streamlit.

## **2. Dataset Description**

- Source: Kaggle (Credit Card Fraud Detection dataset)
- Rows: ~100000 transactions
- Features:
  - V1 to V28: PCA-transformed features
  - Time: Seconds from first transaction
  - Amount: Transaction amount
  - Class: Target (0 = normal, 1 = fraud)
- Class Distribution:
  - Extremely imbalanced: 99.83% legitimate, 0.17% fraud

## **3. EDA Summary**

- Performed basic statistical analysis and visualization.
- Class imbalance visualized with bar plot.
- Histograms for Amount across fraud vs. normal.
- Heatmap showed:
  - V10, V14, V17 had strong correlation with fraud.
  - Time and Amount were weakly correlated.
- Bar plot showed top features correlated with Class.

## **4. Feature Engineering**

- Normalized Amount and Time using StandardScaler.
- Retained V1–V28 as-is (already PCA-transformed).
- Combined features for model input:
  - Time, Amount, V1–V28

## **5. Model Training & Evaluation**

Trained 3 different models:

- Logistic Regression: Baseline linear model
- Decision Tree: Simple non-linear model
- XGBoost: Best performance overall

Data split: 80% train, 20% test (with stratification).

## **6. Metrics & Results**

Evaluation done using:

- Precision
- Recall
- F1 Score

- AUC-ROC

XGBoost outperformed other models with:

- High F1-Score
- ROC-AUC > 0.98
- Good balance of precision and recall

## **7. Deployment (Streamlit App)**

- Final model saved as fraud\_model.pkl
- Streamlit app allows:
  - Manual input of transaction data
  - Real-time prediction (fraud/not fraud)
  - Clean and minimal interface

## **8. Risks & Limitations**

- Class Imbalance: Very few fraud cases; addressed via undersampling (5:1 ratio)
- False Positives: May inconvenience legitimate users
- False Negatives: Dangerous: actual fraud may go undetected
- Feature Meaning: PCA features (V1–V28) are not interpretable
- Real-Time Scaling: Streamlit demo doesn't include user profile or history context