



PROJECT REPORT: CREDIT CARD FRAUD DETECTION

1. Project Overview

Objective

The primary objective of this project is to build a **real-time fraud detection system** that can identify fraudulent credit card transactions with high accuracy. This project uses **Machine Learning (Logistic Regression)** and **Data Visualization** techniques to analyse transaction data and detect anomalies.

2. Dataset Details

Dataset Source

- Dataset: [Credit Card Fraud Detection Dataset](#)
- Source: **Kaggle**
- Number of Transactions: **284,807**
- Number of Fraudulent Transactions: **492**
- Fraudulent Transaction Percentage: **0.172% (highly imbalanced dataset)**

Features in the Dataset

- **Time**: Seconds elapsed between transactions
 - **V1-V28**: Anonymized numerical features obtained via PCA transformation
 - **Amount**: Transaction amount
 - **Class**: Target variable (0 = non-fraud, 1 = Fraud)
-

3. Tools & Technologies Used

Libraries & Frameworks

- **Python** for data analysis and machine learning
- **Pandas & NumPy** for data manipulation
- **Seaborn & Matplotlib** for visualization
- **Scikit-learn** for model training

- **Streamlit** for real-time interactive dashboard
 - **KaggleHub** for dataset retrieval
-

4. Data Preprocessing

✅ Steps Taken

1. Dataset Loading

- Retrieved dataset from Kaggle using kagglehub
- Loaded dataset using `pandas.read_csv()`

2. Exploratory Data Analysis (EDA)

- Checked the number of frauds vs. non-fraud transactions.
- Identified class imbalance (fraud transactions are only **0.172%**).
- Generated **fraud distribution plots** and **correlation heatmap**.

3. Feature Engineering

- Dropped '**Time**' feature as it is not relevant.
 - Kept '**Amount**' and PCA features (**V1-V28**) for model training.
 - Applied **Standard Scaling** (StandardScaler) to normalize numerical features.
-

5. Machine Learning Model

🚀 Model: Logistic Regression

• Why Logistic Regression?

- Works well for binary classification problems.
- Provides probability estimates for fraud detection.
- Efficient and interpretable.

🔧 Model Training

• Train-Test Split:

- **80% training set**

- **20% test set**
- Stratified sampling used to maintain class balance.
- **Performance Metrics:**
 - **Accuracy:** 0.9993
 - **Precision:** 0.857
 - **Recall:** 0.742
 - **F1-Score:** 0.795

Model Evaluation

Metric	Score
Accuracy	99.93%
Precision	85.7%
Recall	74.2%
F1-Score	79.5%

- The model performs **exceptionally well**, but recall needs improvement.
- **Confusion Matrix Analysis:**
 - **False Positives** (wrongly detected as fraud) are **very low**.
 - **False Negatives** (missed fraud cases) exist and need improvement.

6. Real-Time Fraud Detection System

Interactive Dashboard (Streamlit)

A **real-time dashboard** was built using **Streamlit** to:

- **Visualize dataset insights** (fraud vs. non-fraud transactions)
- **Train the model dynamically** and show **accuracy results**
- **Allow users to enter transaction details** and predict **fraud risk**

Fraud Prediction Feature

- Users input **transaction features** (Amount, V1-V28 values).
 - Model predicts **fraud likelihood** in real-time.
 - If fraud probability is high, it **raises an alert**.
-

7. Visualizations & Insights

Fraud vs. Non-Fraud Distribution

- **Highly imbalanced data:** Only **0.172% transactions** are fraudulent.
- **Fraudulent transactions** have different **distribution patterns**.

Correlation Heatmap

- **High correlation** between PCA features (V1-V28).
- **Amount feature** has **weak correlation** with fraud.

Confusion Matrix

- **Most fraud cases** are **detected correctly**.
 - **False negatives** need **improvement** (can be reduced using advanced models).
-

8. Challenges & Solutions

Class Imbalance

- **Issue:** Fraudulent transactions are **rare**, making it hard to train the model.
- **Solution:** Use **Stratified Sampling** to maintain fraud ratio in train-test split.

Improving Fraud Detection

- **Issue:** Recall is **74.2%**, meaning some fraud cases go undetected.
 - **Solution:** Use **SMOTE (Synthetic Minority Over-sampling Technique)**.
-

9. Future Improvements

Model Enhancements

- Train advanced models like **Random Forest, XGBoost, or Neural Networks**.

- Use **SMOTE** to oversample fraud cases for better recall.
- Deploy the **model as an API (Flask/FastAPI)** for real-world use.

Deployment & Scalability

- Deploy on **AWS, GCP, or Azure** for real-time fraud detection.
 - Integrate with a **banking system** to flag fraudulent transactions.
 - Implement **alert notifications** for real-time fraud alerts.
-

10. Conclusion

- ✓ Successfully built a **Credit Card Fraud Detection System**.
 - ✓ Developed an **interactive dashboard** for real-time predictions.
 - ✓ Achieved **high model accuracy (99.93%)**.
 - ✓ **Next Step:** Improve fraud detection **recall** using advanced techniques.
-

11. References

- Dataset: [Kaggle - Credit Card Fraud Detection](#)
 - Libraries Used: **Pandas, NumPy, Seaborn, Scikit-learn, Streamlit**
 - SMOTE Technique: **Handling Class Imbalance in Fraud Detection**
-

12. Appendix

Key Python Libraries Used



`pip install streamlit pandas numpy seaborn matplotlib scikit-learn kagglehub`

Steps to Run the Dashboard

`streamlit run app.py`

Final Thoughts

This project successfully demonstrates **real-time fraud detection** with **high accuracy** and an **interactive user-friendly interface**. Future enhancements can further improve recall, making it an **industry-ready solution**.

 **Let's make banking safer with AI-powered fraud detection!** 

 **Need More Enhancements?**

Let me know if you need any **custom modifications, advanced models, or deployment suggestions!**



SANGAM S BHAMARE