# Capstone Project

# Credit Card Fraud Detection Final Report

1. **Define the Problem Statement:** Credit card fraud is a major financial crime that affects millions of individuals and institutions worldwide, causing billions of dollars in annual losses. The challenge lies in detecting fraudulent transactions quickly and accurately within a vast volume of legitimate data. This project aims to build an effective machine learning pipeline that can identify fraudulent transactions with high recall, thereby minimizing false negatives. The ultimate goal is to improve fraud detection systems to reduce financial risk and increase customer trust for credit card service providers.

2. **Model Outcomes or Predictions:** This is a **supervised classification** problem where the goal is to predict whether a transaction is fraudulent (Class = 1) or not (Class = 0). The expected output is a binary classification label for each transaction. We selected supervised learning because the dataset contains historical labeled data indicating whether transactions were fraud or not. The machine learning models aim to output class labels and probabilities for each prediction.

3. **Data Acquisition:** The dataset used in this project was obtained from a publicly available source on [Kaggle](Kaggle). It consists of anonymized credit card transaction data from European cardholders in September 2013.

    - **Transactions:** 284,807
    - **Fraudulent cases:** 492 ($\approx$0.17%)
    - **Features:** V1-V28 (anonymized PCA components), Time, Amount, and target Class

    Although the dataset comes from a single source, it is well-suited for training fraud detection models due to its realistic class imbalance and anonymized features. Visualizations such as class distribution plots, amount distributions, and fraud vs. non-fraud comparisons were conducted to assess its suitability.

4. **Data Preprocessing/Preparation:**

    a.  **Missing Values and Inconsistencies**
        - The dataset was checked for missing values, and none were found.
        - Feature scaling was applied to Amount and Time using StandardScaler.

    b.  **Train-Test Split**
        - After preprocessing, we applied SMOTE to oversample the minority class (fraud).
        - We then used train_test_split with a 70/30 ratio to create training and test sets.

Prepared by: Bedirhan Ulas

c.  **Encoding and Analysis**
- No categorical features existed due to anonymized PCA-transformed columns.
- Data visualizations were used to identify pattern and class distributions.

5. **Modeling:** The following machine learning models were implemented and evaluated:
- **Logistic Regression:** Used as a baseline due to its interpretability
- **Random Forest Classifier:** Ensemble method capable of handling class imbalance and feature interactions
- **XGBoost Classifier:** Gradient-boosted model known for superior accuracy and robustness

All models were trained on the resampled (balanced) dataset.

6. **Model Evaluation:** Since this is a classification task with a highly imbalanced dataset, we focused on **precision**, **recall**, and **F1-score**, especially for the fraud class (Class = 1).

| Model | Accuracy | Precision (Fraud) | Recall (Fraud) | F1 Score (Fraud) |
|---|---|---|---|---|
| Logistic Regression | 95% | 0.97 | 0.92 | 0.95 |
| Random Forest | 100% | 1.00 | 1.00 | 1.00 |
| XGBoost | 100% | 1.00 | 1.00 | 1.00 |

Evaluation Result: While Logistic Regression performed well, both ensemble models (Random Forest and XGBoost) achieved perfect scores on the balanced test data. This demonstrates that ensemble techniques are highly effective for this use case.

7. **Conclusion:** The final model successfully identified fraudulent transactions with extremely high precision and recall after addressing class imbalance through SMOTE. Ensemble methods proved most effective for this task. This solution, once integrated into a production system, could significantly reduce undetected fraud, improve customer safety, and lower financial risk.

Prepared by: Bedirhan Ulas