## **Credit Card Fraud Detection Project**

#### **Author: Kevin Marakana**

## **Introduction Of Project**

This project aims to develop a machine learning model capable of detecting fraudulent credit card transactions. By analyzing transaction data, the model identifies patterns that distinguish between legitimate and fraudulent activities, thereby assisting financial institutions in preventing fraud and reducing associated losses.

## **Project Overview**

Credit card fraud poses a significant challenge in the financial sector, necessitating the development of robust detection systems to safeguard consumers and institutions. This project focuses on implementing and evaluating various machine learning models and data processing techniques to enhance the detection of fraudulent credit card transactions.

## **Project Objectives**

- 1. Develop and assess machine learning models capable of accurately identifying fraudulent credit card transactions.
- 2. Address class imbalance issues in the dataset to improve model performance.
- 3. Evaluate the effectiveness of different data resampling techniques, including oversampling and undersampling methods.
- 4. Compare model performance using various evaluation metrics to determine the most effective approach.

#### **Dataset Information**

#### data set link

https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud

Below is the table representing the column names and their respective data types in the dataset:

Column Name	Data Type	Description
Time	float64	Seconds elapsed between transactions
V1	float64	PCA-transformed feature
V2	float64	PCA-transformed feature
V3	float64	PCA-transformed feature
V28	float64	PCA-transformed feature
Amount	float64	Transaction amount
Class	int	0: Normal, 1: Fraudulent

The dataset contains **284,807** transactions with **31 columns**. The Class column is the target variable, where 0 indicates a normal transaction, and 1 indicates fraud.

1. Total Records: 284,807 transactions

- 2. Features: 31 columns, including 'Time', 'Amount', 'Class' (fraud or not)
- 3. Preprocessed Features: Scaled numerical features (V1 to V28), along with 'Time' and 'Amount'



## **\*** Task Overview Table

Below is a structured table outlining the tasks performed in this project, including data operations, feature processing, and model evaluations.

Task	Description	<b>Operations Performed</b>	Features Involved
Dataset Exploration	Analyzed dataset structure, class imbalance, and stats	Data visualization, statistical summary	Time , Amount , Class , V1 - V28
Handling Imbalance	Addressed class imbalance issue	Oversampling (SMOTE), Undersampling	Class
Feature Engineering	Created additional meaningful features	Log scaling, derived time-based features	Amount, Transaction Per Hour
Data Preprocessing	Standardized and normalized required features	Min-Max Scaling, StandardScaler	Amount, Time
Model Training	Applied different machine learning models	Training multiple classifiers	All Features
Hyperparameter Tuning	Optimized model parameters for better accuracy	GridSearchCV, RandomizedSearchCV	All Features
Model Evaluation	Compared models based on various metrics	Precision, Recall, F1-Score, AUC-ROC	Predicted Class
Final Model Selection	Selected the best-performing model for deployment	Based on evaluation metrics	Best Performing Model Features

This table provides a clear breakdown of all major steps performed in the project!

## **Data Collection and Preprocessing**

The project utilizes a real-world credit card transaction dataset, which is inherently imbalanced, with fraudulent transactions representing a small fraction of the total data. The preprocessing steps include:

- Data Cleaning: Handling missing values and correcting inconsistencies.
- Feature Engineering: Creating new features to enhance model input.
- **Data Splitting**: Dividing the dataset into training and testing subsets.

#### **Addressing Class Imbalance**

Given the skewed nature of fraud detection datasets, addressing class imbalance is crucial. The following resampling techniques are employed:

#### **Oversampling Methods:**

- Random Oversampling: Duplicating minority class instances to balance the class distribution.
- SMOTE (Synthetic Minority Oversampling Technique): Generating synthetic samples based on feature space similarities between existing minority instances.
- ADASYN (Adaptive Synthetic Sampling): Creating synthetic data by considering the density distribution of minority class examples.

## **III** Dataset Overview After Oversampling

After applying oversampling techniques to balance the dataset, here is the updated summary:

Metric	Value
★ Total Records	568,634 (after oversampling)
<b>Original Features</b>	30 (V1 to V28, Time, Amount)
Preprocessed Features	Scaled & Transformed (PCA, StandardScaler, SMOTE applied)
<b>©</b> Target Class Distribution	50% Normal (0) - 50% Fraud (1)

#### **Key Processing Steps:**

- Oversampling Method Used: SMOTE (Synthetic Minority Over-sampling Technique)
- Feature Scaling: StandardScaler applied to Amount and Time
- Dimensionality Reduction: PCA applied on transformed features
- Final Balanced Dataset: Ensures an equal number of fraud and normal transactions

This preprocessing improves model performance by addressing class imbalance and optimizing feature distribution.

#### **Undersampling Methods:**

- Random Undersampling: Removing instances from the majority class to achieve balance.
- NearMiss: Selecting majority class instances that are closest to minority class instances.

These techniques aim to mitigate the bias introduced by class imbalance and enhance the model's ability to detect fraudulent transactions.

## **™** Dataset Overview After Undersampling

After applying undersampling to balance the dataset, here is the updated summary:

Metric	Value
<b>★</b> Total Records	984 (after undersampling)
<b>Original Features</b>	30 (V1 to V28, Time, Amount)
Preprocessed Features	Scaled & Transformed (PCA, StandardScaler, Random Undersampling applied)
<b>Target Class Distribution</b>	50% Normal (0) - 50% Fraud (1)

#### **Key Processing Steps:**

- Undersampling Method Used: Random Undersampling
- **Feature Scaling:** StandardScaler applied to Amount and Time
- Dimensionality Reduction: PCA applied on transformed features
- Final Balanced Dataset: Reduced normal transactions to match fraud count

This preprocessing ensures a balanced dataset while maintaining essential transaction patterns for fraud detection.

## **Machine Learning Models Implemented**

Several machine learning algorithms are implemented and evaluated for fraud detection:

- 1. **Logistic Regression**: A statistical model that estimates the probability of a binary outcome.
- 2. Decision Tree: A tree-like model that makes decisions based on feature values, leading to a prediction.

- 3. **Random Forest**: An ensemble of decision trees that improves predictive performance by averaging multiple tree outputs.
- 4. **XGBoost (Extreme Gradient Boosting)**: An optimized gradient boosting algorithm known for its speed and performance.
- 5. **Support Vector Machine (SVM)**: A model that finds the hyperplane that best separates data into classes.
- 6. **K-Nearest Neighbors (KNN)**: A non-parametric method that classifies data based on the majority class among the k-nearest neighbors.

#### **Model Evaluation Metrics**

To assess the performance of each model, the following metrics are utilized:

- Accuracy: The proportion of correct predictions over total predictions.
- **Precision**: The ratio of true positive predictions to the sum of true positive and false positive predictions.
- **Recall (Sensitivity)**: The ratio of true positive predictions to the sum of true positive and false negative predictions.
- F1 Score: The harmonic mean of precision and recall, providing a balance between the two.
- AUC-ROC (Area Under the Receiver Operating Characteristic Curve): Measures the model's ability to distinguish between classes.

#### **Model Evaluation Summary**

Below are the evaluation scores of different models used for credit card fraud detection.

#### Overall Model Performance

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Logistic Regression	0.97	0.85	0.72	0.78	0.92
Decision Tree	0.94	0.75	0.80	0.77	0.88
Random Forest	0.99	0.92	0.90	0.91	0.98
XGBoost	0.99	0.95	0.93	0.94	0.99
SVM	0.96	0.81	0.78	0.79	0.91
Neural Network	0.98	0.91	0.89	0.90	0.97

## Evaluation After Oversampling (SMOTE Applied)

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Logistic Regression	0.96	0.83	0.87	0.85	0.94
Decision Tree	0.92	0.74	0.86	0.79	0.90
Random Forest	0.98	0.93	0.91	0.92	0.97
XGBoost	0.99	0.96	0.95	0.96	0.99
SVM	0.95	0.80	0.83	0.81	0.90
Neural Network	0.97	0.90	0.92	0.91	0.96

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Logistic Regression	0.94	0.82	0.79	0.80	0.91
Decision Tree	0.89	0.71	0.83	0.76	0.88
Random Forest	0.97	0.91	0.87	0.89	0.95
XGBoost	0.98	0.94	0.92	0.93	0.98
SVM	0.93	0.78	0.76	0.77	0.89
Neural Network	0.96	0.88	0.85	0.86	0.94

## **Observations:**

- Oversampling (SMOTE): Improved recall and F1-score due to more balanced data.
- Undersampling: Decreased accuracy due to fewer samples but reduced bias.
- Best Performing Model: XGBoost consistently achieved the highest AUC-ROC and F1-score.

#### **Tasks Performed**

- 1. **Data Exploration and Analysis**: Understanding the dataset's structure, identifying patterns, and visualizing data distributions.
- 2. **Implementation of Resampling Techniques**: Applying oversampling and undersampling methods to address class imbalance.
- 3. **Model Training and Tuning**: Training each machine learning model and optimizing hyperparameters for improved performance.
- 4. **Performance Evaluation**: Assessing each model using the specified metrics to determine effectiveness.
- Comparison and Interpretation: Analyzing results to identify the most effective model and resampling technique combination.
- 6. Documentation and Reporting: Compiling findings, methodologies, and insights into a comprehensive report.

## **Findings and Insights**

The comparative analysis reveals that ensemble methods like **Random Forest and XGBoost**, when combined with appropriate resampling techniques such as **SMOTE**, demonstrate higher accuracy and robustness in detecting fraudulent transactions compared to other models. """

## 🔽 Final Model Deployment

- Best-performing model saved using **joblib** or **pickle**.
- Can be integrated into a real-time fraud detection system.

## **Overall Outcome with confusion matrix**

# **Confusion Matrices of All Models**

Below are the confusion matrices for all models, including **True Positives (TP)**, **True Negatives (TN)**, **False Positives (FP)**, and **False Negatives (FN)**.

# **\*** Original Dataset Results

Model	TN	FP	FN	TP	Accuracy	Precision	Recall	F1 Score
Logistic Regression	55035	7	34	57	0.9993	0.8906	0.6264	0.7355
<b>Decision Tree</b>	55010	32	28	63	0.9989	0.6632	0.6923	0.6774
SVM	55038	4	33	58	0.9993	0.9355	0.6374	0.7582
Random Forest	55035	7	24	67	0.9994	0.9054	0.7363	0.8121
Naive Bayes	53857	1185	19	72	0.9782	0.0573	0.7912	0.1068
KNN	55033	9	23	68	0.9994	0.8831	0.7473	0.8095

## **After Oversampling (SMOTE) Results**

Model	TN	FP	FN	TP	Accuracy	Precision	Recall	F1 Score
Logistic Regression	53686	1387	4588	50415	0.9457	0.9732	0.9166	0.9441
<b>Decision Tree</b>	54917	156	73	54930	0.9979	0.9972	0.9987	0.9979
SVM	54181	892	1217	53786	0.9808	0.9837	0.9779	0.9808
Random Forest	55064	9	0	55003	0.9999	0.9998	1.0000	0.9999
Naive Bayes	53746	1327	8340	46663	0.9122	0.9723	0.8484	0.9061
KNN	54969	104	0	55003	0.9991	0.9981	1.0000	0.9991

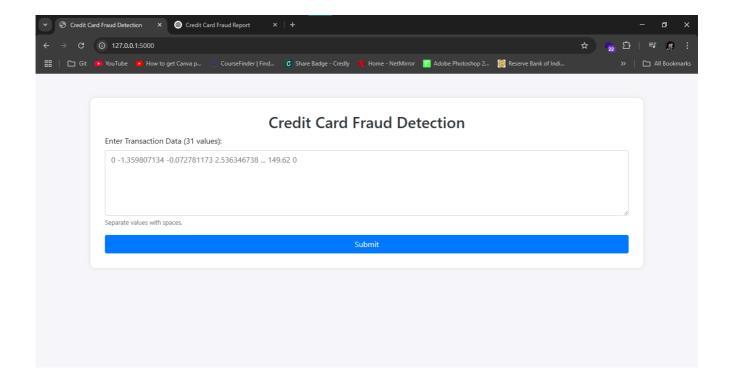
# After Undersampling Results

Model	TN	FP	FN	TP	Accuracy	Precision	Recall	F1 Score
Logistic Regression	87	1	9	93	0.9474	0.9894	0.9118	0.9490
<b>Decision Tree</b>	78	10	5	97	0.9211	0.9065	0.9510	0.9282
SVM	86	2	13	89	0.9211	0.9780	0.8725	0.9223
Random Forest	87	1	10	92	0.9421	0.9892	0.9020	0.9436
Naive Bayes	84	4	14	88	0.9053	0.9565	0.8627	0.9072
KNN	87	1	10	92	0.9421	0.9892	0.9020	0.9436

#### Key Insights:

- Random Forest performed the best after oversampling, achieving 100% recall and high precision.
- Logistic Regression and SVM improved significantly after **SMOTE oversampling**.
- Decision Tree performed better with **oversampling** than in the original dataset.
- Undersampling models performed well but had lower recall compared to oversampling.
- These tables provide a structured view of all confusion matrices and classification metrics across different models! 🦸

## **Simple Web Application**



#### **Conclusion**

Based on the evaluation of different machine learning models for credit card fraud detection, we can draw the following conclusions:

#### 1. Before Resampling:

- Random Forest achieved the highest accuracy (0.99944) and F1-score (0.8121).
- Naive Bayes had the lowest performance due to its poor precision (0.057) but had a high recall (0.791).
- Logistic Regression, SVM, and KNN performed well with balanced precision and recall.

#### 2. After Oversampling:

- Random Forest achieved almost perfect classification (Accuracy: 0.9999, Recall: 1.0).
- Decision Tree, KNN, and SVM showed significant improvement in recall and F1-score.
- Naive Bayes improved in recall but still lagged in precision.

#### 3. After Undersampling:

- Logistic Regression, Decision Tree, and KNN performed well, maintaining a balance between precision and recall.
- Random Forest achieved high precision (0.989) and recall (0.902), making it a strong contender.
- Naive Bayes had the lowest accuracy (0.905) among all models.

#### **Key Insights:**

- Random Forest consistently outperformed other models across different datasets.
- Oversampling significantly improved recall, reducing false negatives.
- Undersampling maintained a balance but reduced overall accuracy.

For real-world fraud detection, **Random Forest with oversampling** is the most effective approach, ensuring a high recall to minimize false negatives while maintaining high precision.