# Credit Card Fraud Detection: A Comprehensive Report

## Table of Contents

## 1. Introduction

Credit card fraud is a serious issue affecting individuals and financial institutions worldwide. It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase. This report delves into the mechanisms of credit card fraud detection, focusing on the methodologies used to identify and prevent fraudulent activities.

## 2. Overview of Credit Card Fraud Project

The dataset which is being used here contains transactions made by credit cards in September 2013 by European cardholders. This dataset presents transactions that occurred in two days, where we have **492 frauds** out of **284,807** transactions. The dataset is highly unbalanced, the positive class (frauds) account for **0.172%** of all transactions.

**Data:**

The dataset for this project can be accessed by clicking the link provided below.

[creditcard.csv](creditcard.csv)

## 3. Importance of Fraud Detection

Detecting fraud is crucial to prevent financial losses, maintain customer trust, and uphold the integrity of financial systems. Effective fraud detection systems help mitigate risks and protect sensitive information.

# 4. Techniques for Fraud Detection

## 4.1 Data Preprocessing

Data preprocessing involves cleaning and transforming raw data into a suitable format for analysis. Key steps include:

- **Data Cleaning**: Removing duplicates and missing values and correcting errors.
- **Normalization**: Standardizing data to a uniform scale.
- **Feature Engineering**: Creating new features that can enhance model performance.

## 4.2  Model Selection and Model Training

We have built a variety of classification models in this section and then test them out individually to see which are the best suited for this task. Here's a structured approach to explaining each model:

### 1) Logistic Regression

**Overview**: Logistic Regression is a statistical method used for binary classification problems. Despite its name, it is a linear model used to estimate the probability that a given input belongs to a particular category.

**Key Points**:

- **Equation**: The logistic function (sigmoid) transforms the linear regression output into a probability between 0 and 1.
- **Interpretation**: Outputs can be interpreted as the likelihood of the instance belonging to the positive class.
- **Binary Classification**: Typically used for problems with two possible outcomes (e.g., spam vs. non-spam).

**Strengths**:

- **Simplicity**: Easy to implement and interpret.
- **Probabilistic Output**: Provides probability scores for each class.

- **Feature Importance**: Coefficients indicate the importance and effect of each feature.

**Weaknesses**:

- **Linearity**: Assumes a linear relationship between the independent variables and the log odds.
- **Performance**: May not perform well with complex, non-linear datasets without proper feature engineering.

**Use Cases**:

- Email spam detection.
- Credit scoring.
- Medical diagnosis (e.g., predicting the presence or absence of a disease).

## 2) Random Forest

**Overview**: Random Forest is an ensemble learning method that constructs multiple decision trees during training and merges their outputs to improve accuracy and control overfitting.

**Key Points**:

- **Ensemble Method**: Combines the predictions of multiple decision trees.
- **Bagging**: Uses bootstrap aggregating to create diverse trees.
- **Randomness**: Introduces randomness in feature selection and sample selection.

**Strengths**:

- **Accuracy**: Generally high accuracy due to ensemble nature.
- **Robustness**: Resistant to overfitting and noisy data.
- **Feature Importance**: Can assess the relative importance of different features.

**Weaknesses**:

- **Complexity**: Harder to interpret compared to single decision trees.
- **Computationally Intensive**: Requires more memory and computational power.

**Use Cases**:

  - Fraud detection.
  - Image classification.
  - Predictive maintenance.

## 3) K-Nearest Neighbors (KNN)

**Overview**: K-Nearest Neighbors is a simple, instance-based learning algorithm used for classification and regression. It classifies instances based on the majority label of their k-nearest neighbors in the feature space.

**Key Points**:

- **Lazy Learning**: No explicit training phase; the algorithm makes decisions based on the entire training dataset.
- **Distance Metrics**: Commonly uses Euclidean distance to find nearest neighbors.
- **Parameter**: The choice of k (number of neighbors) significantly affects performance.

**Strengths**:

- **Simplicity**: Easy to understand and implement.
- **Adaptability**: Can be used for both classification and regression tasks.
- **No Training Phase**: Fast prediction times once the data is stored.

**Weaknesses**:

- **Scalability**: Computationally expensive as it requires calculating the distance to all points in the training dataset for each prediction.
- **Sensitivity**: Performance is sensitive to the choice of k and the distance metric.
- **Curse of Dimensionality**: High-dimensional data can make distance metrics less effective.

**Use Cases**:

- Recommendation systems.
- Medical diagnosis.
- Pattern recognition.

## 4) Decision Tree

**Overview**: Decision Trees are non-parametric supervised learning methods used for classification and regression. They partition the data into subsets based on feature value tests.

**Key Points**:

- **Tree Structure**: Consists of nodes (features), branches (decision rules), and leaves (outcome).
- **Splitting Criteria**: Common criteria include Gini impurity, entropy, or variance reduction.
- **Interpretability**: Easy to visualize and interpret the decision-making process.

**Strengths**:

- **Interpretability**: Clear and easy to understand visual representation.
- **Non-Linear Relationships**: Can capture complex interactions between features.
- **Feature Importance**: Naturally performs feature selection during the tree-building process.

**Weaknesses**:

- **Overfitting**: Prone to overfitting, especially with deep trees.
- **Stability**: Small changes in the data can result in different tree structures.
- **Bias**: Tends to favor features with more levels.

**Use Cases**:

- Customer segmentation.
- Loan approval decisions.
- Risk assessment.

## Conclusion

Each machine learning model has its unique characteristics, strengths, and weaknesses, making them suitable for different types of problems and datasets. Understanding these aspects helps in selecting the appropriate model for a given task and in effectively communicating the results and rationale in a report.

## 5. Challenges in Fraud Detection

1. **Imbalanced Data**: Fraudulent transactions are rare compared to legitimate ones, leading to challenges in model training.
2. **Dynamic Fraud Patterns**: Fraud tactics continually evolve, necessitating adaptable models.
3. **False Positives**: Incorrectly flagging legitimate transactions can inconvenience customers and harm reputation.
4. **Data Privacy**: Ensuring that data used for detection does not violate privacy regulations.

# 6. Evaluation Metrics

To assess the performance of fraud detection models, several metrics are used:

- **Accuracy**: The proportion of correctly identified transactions (both fraud and non-fraud).
- **Precision**: The proportion of true positive frauds among all identified frauds.
- **Recall (Sensitivity)**: The proportion of actual frauds correctly identified.
- **F1 Score**: The harmonic mean of precision and recall, balancing both concerns.
- **ROC-AUC Curve**: Represents the model's ability to discriminate between fraud and non-fraud cases.

# 8. Conclusion

Each machine learning model discussed—Logistic Regression, Random Forest, KNN, and Decision Trees—brings unique strengths and considerations to the table. Logistic Regression is valuable for its simplicity and probabilistic interpretation, while Random Forest excels in accuracy and robustness against overfitting. KNN offers simplicity and flexibility, particularly in scenarios where interpretability is crucial. Decision Trees provide a straightforward, visual representation of decision-making but require careful management to prevent overfitting.

`Hence, we can say that Random Forest model performed best with an accuracy score of 98% for the detection of fraudulent and 93% for the detection of non-fraudulent transactions.`

## Regards,

**SHALINI BHATIA**

**DATA SCIENTIST**

https://github.com/09263

*SHALINI BHATIA*