

# Credit Card Fraud Detection Using Machine Learning

## Introduction

Credit card fraud is a significant global challenge, with billions of dollars lost each year. Machine learning offers a powerful tool for detecting fraudulent transactions by identifying patterns indicative of fraud. This report combines insights from two separate studies on credit card fraud detection to present a comprehensive analysis of machine learning techniques, their effectiveness, and the differences in their findings.

## Objectives

Both studies share a common goal: to identify fraudulent credit card transactions using machine learning models. They aim to evaluate multiple algorithms and determine the most effective model based on accuracy and other performance metrics.

## Dataset

Both studies used a publicly available dataset from Kaggle, containing:

- **284,807 transactions** recorded over two days in Europe.
- **492 transactions labeled as fraudulent** (highly imbalanced dataset).
- **31 attributes**, including 28 numerical variables transformed using Principal Component Analysis (PCA) for confidentiality, as well as **Time**, **Amount**, and **Class** (binary labels: 1 for fraud, 0 for non-fraud).

## Methodology

Both studies followed similar steps but differed slightly in their execution:

1. **Data Preprocessing:**
  - Both cleaned the data, handled missing values, and transformed the class labels for modeling.

- One study described additional transformations using tools like Sublime Text for integration with specific platforms (e.g., Weka).

2. **Algorithms Used:** Both studies evaluated the following machine learning models:

- K-Nearest Neighbors (KNN)
- Logistic Regression
- Support Vector Machine (SVM)
- Decision Tree (DT)
- Additionally, the first study also referenced Naïve Bayes in its review but did not implement it.

3. **Model Evaluation:**

- Both used confusion matrices, accuracy, precision, recall, and other metrics for comparison.
- Data splitting varied slightly, with a training-to-testing ratio of 70:30 in the first study, and unspecified proportions in the second study.

## Results

### First Study:

- **Best Model: Support Vector Machine (SVM)** achieved the highest accuracy of **99.94%**, demonstrating its strength in handling non-linear boundaries and high-dimensional data.
- **Other Results:**
  - Logistic Regression: **99.92% accuracy**
  - KNN: **99.89% accuracy**
  - Naïve Bayes: **97.76% accuracy**
  - Decision Tree: Not tested directly.
- **Observation:** SVM was identified as the most robust model for fraud detection due to its ability to maximize class separation with minimal misclassifications.

## Second Study:

- **Best Models:** KNN and **Decision Tree** both achieved **100% accuracy**, demonstrating perfect performance in detecting fraudulent transactions.
- **Other Results:**
  - SVM: **97.59% accuracy**
  - Logistic Regression: **93.51% accuracy (training), 91.88% (testing)**
- **Observation:** KNN and Decision Tree were identified as the top-performing models, attributed to their simplicity, interpretability, and strong performance on imbalanced datasets.

## Key Differences

### 1. Top-Performing Models:

- The **first study** identified SVM as the best-performing model, while the **second study** concluded that KNN and Decision Tree were superior.
- This discrepancy likely arises from differences in implementation, evaluation metrics, or dataset splits.

### 2. Model Rankings:

- Logistic Regression performed better in the first study (**99.92% accuracy**) compared to the second (**91.88% accuracy on test data**).
- KNN and Decision Tree were highlighted as top performers in the second study but were not emphasized in the first study.

### 3. Focus on Analysis:

- The **first study** focused on evaluating SVM's generalization capability and highlighted Naïve Bayes, which was absent in the second study.

- The **second study** emphasized KNN's simplicity and the Decision Tree's ability to handle diverse attributes.

#### 4. **Evaluation Metrics:**

- The first study detailed the use of tools like Weka and R for model implementation, while the second study relied on Jupyter Notebook for all models.

### **Recommendations**

#### **First Study:**

- Explore alternative datasets with diverse characteristics.
- Incorporate geolocation data to improve fraud detection by identifying mismatches in transaction locations.

#### **Second Study:**

- Test additional algorithms and hyperparameters.
- Experiment with different data-splitting ratios to assess their impact on model performance.

### **Conclusion**

Both studies demonstrate the effectiveness of machine learning in credit card fraud detection, though their results and interpretations differ:

- The **first study** positions SVM as the most effective model due to its accuracy and robust generalization capabilities.
- The **second study** highlights KNN and Decision Tree as the best performers, achieving perfect accuracy.

These differences underscore the importance of model selection, data preprocessing, and evaluation methods. Together, the findings provide valuable insights into leveraging machine

learning to combat credit card fraud effectively. Future work could focus on harmonizing methodologies to ensure consistent evaluations across studies.