

**CSE 422**  
**FINAL PROJECT REPORT**

**Topic:** Credit Card Fraud Detection Analysis

**Submitted to:** MD. Mustakin Alam  
Labib Hasan Khan

**Submitted by:** Shuha Jahan [21301335]  
Ramisa Fariha Prova [20301001]

**Submission Date:** 27th April, 2024

## Table of Contents

Introduction.....	3
Dataset Description.....	3-5
Data pre-processing.....	5-6
Feature Scaling.....	6
Dataset Splitting.....	6
Model Training	
• KNN .....	6
• Decision Tree .....	7
• MLP.....	7
• NB.....	7
Model Selection / Comparison Analysis.....	7-8
Conclusion.....	8

## Introduction

In today's digital age, credit card usage has skyrocketed worldwide, with people increasingly relying on cashless transactions. However, this convenience comes with a significant downside, criminal credit card transactions result in substantial financial losses each year. According to the PwC global economic crime survey of 2017, nearly half of all organizations have experienced economic crime. This underscores the urgent need to address the issue of credit card fraud detection, especially with the emergence of new technologies providing additional avenues for criminals. Credit card fraud not only impacts financial institutions and merchants but also individual cardholders, damaging their trust and potentially tarnishing the reputation of businesses involved.

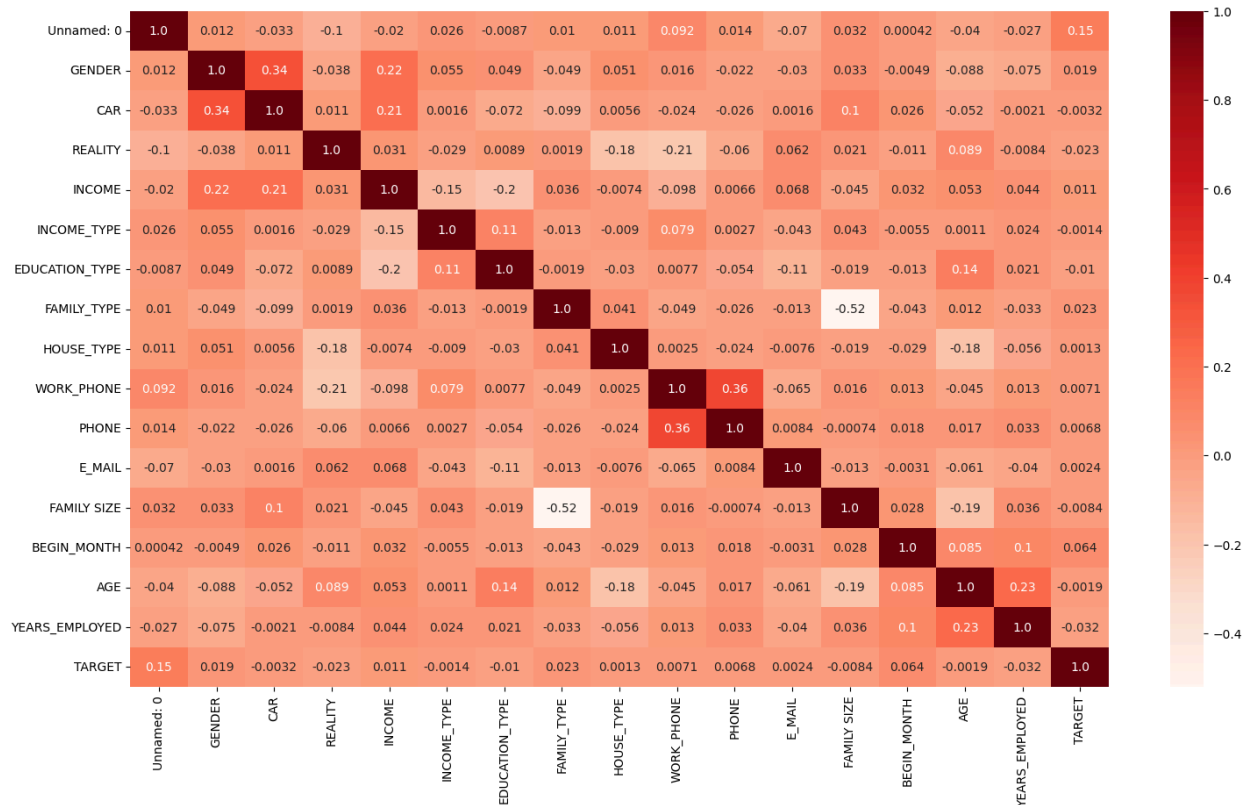
To combat this, we're implementing a systematic approach using three different machine learning methods — K-Nearest Neighbors (KNN), Decision Tree, and Multilayer Perceptron (MLP) to classify credit card datasets. In the Credit Card Fraud Detection project, we're essentially teaching a computer to learn from past credit card transactions, specifically those that were fraudulent. By understanding the patterns and characteristics of these fraudulent transactions, the computer can then predict whether a new transaction is likely to be fraudulent or not. It's like giving the computer a bunch of real-life examples of sneaky behavior with credit cards so it can become more vigilant and spot suspicious activity in the future. Through this project, we aim to enhance fraud detection and safeguard individuals and businesses from the devastating effects of credit card fraud.

## Data Description

The dataset used in this report is taken from the official Kaggle repository.

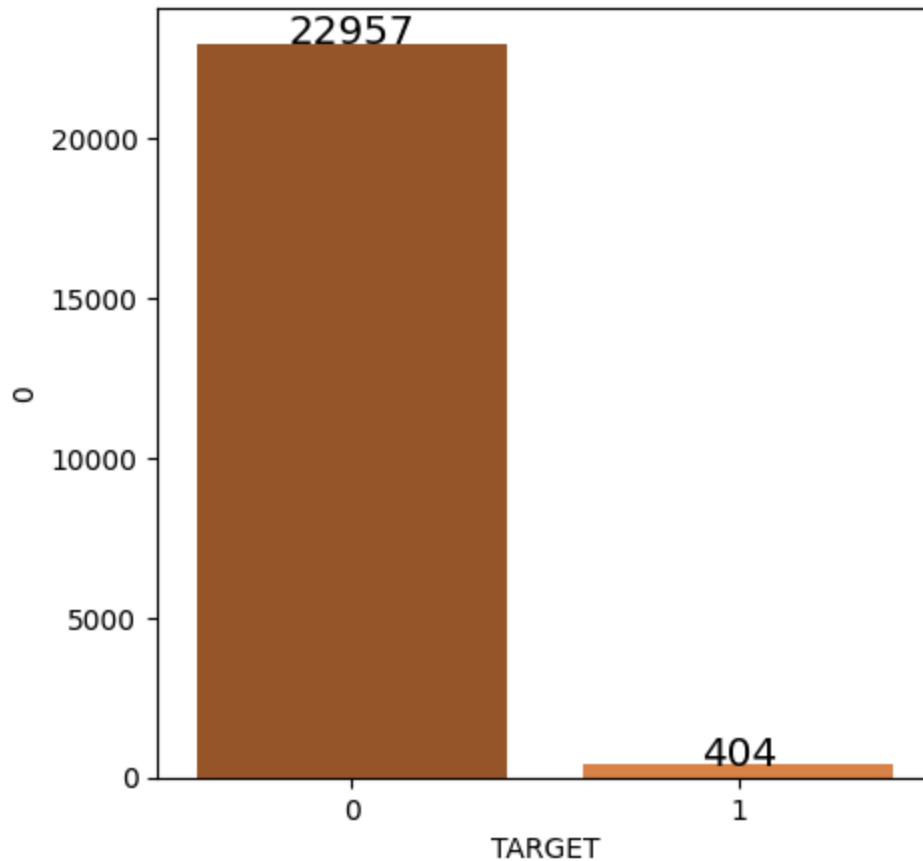
**Source:** <https://www.kaggle.com/datasets/dark06thunder/credit-card-dataset>

The dataset comprises 25,134 data points with 19 features. This problem is basically a classification task with a binary class dataset, as indicated by the “TARGET” feature, which categorizes individuals into those with payment difficulties (class 1) and those without (class 0). The features in the dataset consist of both quantitative and categorical types, including demographic information like age, income, and family size, as well as categorical attributes such as gender, education type, and family status. A correlation matrix was calculated to analyze the degree of linear relationship between each pair of features, including the output variable. A heatmap visualization of the correlation matrix provides insights into which features are most correlated with each other and with the target variable.



*Figure: The heatmap of the entire dataset*

Upon examining the dataset, it is evident that it suffers from class imbalance, as the distribution of instances across the two classes is highly skewed. Specifically, class 0 (indicating individuals without payment difficulties) dominates the dataset with 24,712 instances, while class 1 (representing individuals with payment difficulties) comprises only 422 instances. To visually represent the imbalance, a bar chart can be created where the x-axis denotes the class labels, and the y-axis represents the count of instances for each class. This chart will clearly illustrate the disproportionate distribution of instances across the classes, highlighting the challenge of working with imbalanced data.



*Figure: Bar chart of Imbalanced Data*

## Dataset Pre-processing

The project involves several data pre-processing steps to prepare the dataset for machine learning model training. Here,

### 1. Faults:

- NULL Values: The dataset was examined for missing values using the “isna().sum()” function, revealing no instances of missing data.
- Categorical Values: Categorical attributes such as “GENDER”, “CAR”, “REALITY”, “INCOME\_TYPE”, “EDUCATION\_TYPE”, “FAMILY\_TYPE”, and “HOUSE\_TYPE” were identified. These features required encoding to convert them into numerical format for compatibility with machine learning algorithms.

### 2. Solutions:

- NULL Values: Since no missing values were detected, no further steps were necessary for addressing missing data.

- Categorical Values: Label encoding was applied to the categorical columns using the “LabelEncoder” module from “sklearn.preprocessing”. This transformation converts categorical values into numerical representations, making them suitable for modeling while preserving ordinal information where applicable. Additionally, the "FLAG\_MOBIL" column, which contained only one unique value, was removed as it did not provide any useful information for modeling.

## Feature Scaling

In the Credit Card Fraud Detection project, Feature Scaling is applied as a preprocessing step to standardize the range of independent features or variables in the dataset. This process ensures that all features contribute equally to the analysis and model training, preventing features with larger scales from dominating the learning process. Specifically, the MinMaxScaler from sklearn.preprocessing is used to scale the features to a specified range (typically between 0 and 1) by subtracting the minimum value and dividing by the range of each feature. This normalization technique preserves the relative relationships between the values of each feature while ensuring consistency in their scales.

## Dataset Splitting

In the Credit Card Fraud Detection project, the dataset is split into training and testing sets using the “train\_test\_split function” from the “sklearn.model\_selection” module. The splitting process follows a random approach, where the data is divided randomly into two subsets: a training set and a testing set. By default, 80% of the data is allocated to the training set, while the remaining 20% is assigned to the testing set. This random splitting ensures that the model is trained on a diverse range of data samples and evaluated on unseen data, helping to assess its generalization performance.

## Model Training and Testing

In the Credit Card Fraud Detection project, we have used three different machine learning models, which are as follows:

1. **K-Nearest Neighbors (KNN):** KNN is a classification algorithm that assigns a class label to an input data point based on the majority class of its nearest neighbors. It calculates the distance between the input data point and all other data points in the training set, and then classifies the data point based on the class labels of its k nearest neighbors.

2. **Decision Tree (DT):** Decision Tree is a non-parametric supervised learning method used for classification and regression tasks. It partitions the feature space into regions and predicts the target variable based on the majority class or average value of the training instances within each region.
3. **Multilayer Perceptron (MLP):** MLP is a type of artificial neural network consisting of multiple layers of nodes, each connected to the next layer. It is a versatile architecture capable of learning complex relationships between input and output variables, making it suitable for various machine learning tasks, including classification.
4. **Naive Bayes (NB):** Naive Bayes is a probabilistic classifier based on Bayes' theorem with the "naive" assumption of feature independence. It calculates the probability of a class label given the observed feature values using conditional probability and assumes that features are conditionally independent given the class label. NB is computationally efficient, requires a small amount of training data, and performs well in text classification and other domains with high-dimensional feature spaces.

Each model is trained on the training data and evaluated on the testing data to assess its performance based on metrics such as accuracy, precision, recall, and F1 score. Additionally, confusion matrices are used to visualize the model's performance in predicting different classes.

## Model Selection and Comparison Analysis

After analyzing the model performance metrics in the Credit Card Fraud Detection project, we observe that:

1. **Bar Chart showcasing Prediction Accuracy:** The bar chart illustrates the prediction accuracy of each model. K-Nearest Neighbors (KNN) has the highest accuracy of 99.17%, followed closely by Decision Tree (DT) with 98.89%. Multi-Layer Perceptron (MLP) and Naive Bayes (NB) have slightly lower accuracies compared to KNN and DT.
2. **Precision and Recall Comparison:** Precision measures the accuracy of positive predictions, while recall measures the proportion of actual positives that were correctly identified. KNN has the highest precision (99.09%) and recall (99.17%), indicating that it has the best balance between precision and recall. Decision Tree and Naive Bayes have the same precision and recall scores, which are slightly lower than KNN. Multi-Layer Perceptron (MLP) has the lowest precision and recall scores among the models, indicating that it may be slightly less accurate in predicting positive cases.
3. **Confusion Matrix:** The confusion matrix provides a detailed breakdown of the model's performance in terms of true positives, false positives, true negatives, and false negatives. KNN has the fewest misclassifications, as evidenced by its high accuracy and precision. Decision Tree and Naive Bayes perform similarly in terms of confusion matrix metrics, while MLP shows slightly more misclassifications compared to the other models.

In summary, K-Nearest Neighbors (KNN) appears to be the best-performing model among the four, with the highest accuracy, precision, recall, and the fewest misclassifications. Decision Tree and Naive Bayes perform reasonably well but are slightly behind KNN. Multi-Layer Perceptron (MLP) shows slightly lower performance compared to the other models, indicating that it may not be as suitable for this particular task.

## **Conclusion**

In conclusion, the Credit Card Fraud Detection project utilized machine learning algorithms to identify fraudulent transactions. The dataset consisted of various demographic and financial attributes, with the target variable indicating payment difficulties. After preprocessing, feature engineering, and model training, K-Nearest Neighbors (KNN) and Decision Tree (DT) models emerged as the top performers, exhibiting high accuracy, precision, recall, and F1-score. These models demonstrated their effectiveness in distinguishing between fraudulent and non-fraudulent transactions, making them suitable for real-world application in credit card fraud detection systems. The project underscores the importance of leveraging machine learning techniques to enhance fraud detection mechanisms and protect consumers and financial institutions from fraudulent activities.