# Hackathon Report: Credit Card Fraud Detection using Machine Learning

## Introduction

In this report, we present the workflow and results of a hackathon project focused on "Credit Card Fraud Detection" using Machine Learning techniques. The primary goal of this project is to develop a binary classification model that can accurately identify fraudulent credit card transactions. To achieve this, we used a Logistic Regression model and followed a structured workflow.

### Problem Statement

Credit card fraud is a significant concern for both financial institutions and consumers. Detecting fraudulent transactions is essential to protect customers and minimize financial losses. Machine learning can play a crucial role in automating the detection of such fraudulent activities.

### Workflow of the Model

Our approach to solving this problem can be broken down into the following steps:

- **Collection of Data**: We start by gathering the necessary dataset for training and testing our model.
- **Data Preprocessing**: Data preprocessing is crucial for ensuring the data's quality and suitability for training. This step includes handling missing values and scaling features if necessary.
- **Splitting Test and Training Data**: To evaluate the model's performance, we split the dataset into training and testing sets.
- **Model Training**: We use a Logistic Regression model for binary classification. The model learns from the training data to make predictions.
- **Model Evaluation**: We assess the model's performance using accuracy scores on both the training and testing data.
- **Prediction System**: Finally, we use the trained model to make predictions on new, unseen data to detect potential credit card fraud.

### Implementation

We implemented the above workflow using Python and popular libraries such as NumPy, Pandas, Matplotlib, Seaborn, and Scikit-Learn.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

## Data Analysis

Before proceeding with modeling, we conducted an exploratory data analysis (EDA) to understand the dataset's characteristics.

- The dataset contains 73,377 rows and 31 columns.
- We identified that the dataset has 28 missing values.
- The data is highly unbalanced, with a majority of transactions being legitimate (Class 0) and only a few being fraudulent (Class 1).

## Data Preprocessing

To address the class imbalance issue, we performed under-sampling, creating a balanced dataset with equal numbers of legitimate and fraudulent transactions.

## Data Visualization

We used data visualization techniques to gain insights into the data's features and relationships. A heatmap was created to visualize the correlation matrix between features.

## Model Training

We trained a Logistic Regression model on the preprocessed dataset. Logistic Regression is a suitable choice for binary classification tasks like fraud detection.

## Model Evaluation

The model's performance was evaluated using accuracy scores:

- Accuracy on the training data: `traning_data_accuracy`
- Accuracy on the testing data: `test_data_accuracy`

## Results

The model achieved an accuracy of approximately `traning_data_accuracy` on the training data and `test_data_accuracy` on the testing data. While accuracy is an essential metric, other metrics such as precision, recall, and F1-score should also be considered for a more comprehensive evaluation.

# Conclusion

In this hackathon project, we successfully developed a credit card fraud detection model using Logistic Regression. The model demonstrates promising accuracy in identifying fraudulent transactions. Further refinement and evaluation, including the consideration of additional metrics and techniques for handling imbalanced data, could lead to even better results. This project serves as a foundation for future work in credit card fraud detection using machine learning.