# Abstract

The objective of this project is to develop a machine learning model capable of detecting fraudulent transactions. Fraud detection is critical for financial institutions to prevent significant losses. By leveraging transaction data, this project aims to classify transactions as either fraudulent or non-fraudulent using various machine learning techniques. The project covers data collection, preprocessing, feature engineering, model selection, training, and evaluation. The findings demonstrate the effectiveness of machine learning models in identifying fraudulent transactions with high accuracy.

# Introduction

Fraudulent transactions are a major concern for financial institutions, resulting in substantial financial losses annually. Effective detection methods are crucial for mitigating these losses. Traditional rule-based fraud detection systems are becoming less effective as fraudsters develop more sophisticated techniques. Machine learning models offer a promising alternative by learning patterns from historical data and identifying anomalies indicative of fraud. This project explores the development of a machine learning-based fraud detection system using a publicly available transaction dataset.

# Methods

## Data Collection

For this project, we use the **Credit Card Fraud Detection** dataset from Kaggle, which contains anonymized features of credit card transactions over two days. The dataset includes 284,807 transactions, of which 492 are fraudulent.

## Data Preprocessing

Data preprocessing involves handling missing values, normalizing numerical features, and addressing class imbalance.

## Feature Engineering

Feature engineering is minimal since the features in the dataset are anonymized. We focus on splitting the data into training and testing sets.

## Model Selection

We experiment with several machine learning models including Logistic Regression, Decision Trees, and Random Forests. Due to class imbalance, we use precision, recall, and F1-score as evaluation metrics.

# Discussion

The results show that the Random Forest classifier outperforms Logistic Regression and Decision Tree models in detecting fraudulent transactions. This is attributed to the ensemble nature of Random Forest, which reduces overfitting and improves generalization. Precision and recall scores are crucial for fraud detection to minimize false positives and negatives. Further tuning of hyperparameters and testing with more sophisticated models like Gradient Boosting or Neural Networks could potentially improve the performance.

## Conclusion

In this project, we developed a machine learning model for fraud detection using the **Credit Card Fraud Detection** dataset. The Random Forest model demonstrated the highest accuracy, precision, and recall among the tested models. The findings suggest that machine learning techniques can effectively identify fraudulent transactions, offering a valuable tool for financial institutions. Future work includes exploring more advanced models and real-time fraud detection implementation.