

Title: Credit Card Fraud Detection Analysis

Author: John Pandian

Executive Summary

This report presents a comprehensive analysis aimed at detecting credit card fraud using advanced data processing and machine learning techniques. The analysis is based on a dataset of credit card transactions, with a focus on identifying fraudulent activities.

Data Overview

The data consists of two sets: a training dataset (`df_train`) and a testing dataset (`df_test`). Initial exploration using `df_train.head()` and `df_test.head()` provided an overview of the features available for analysis.

Exploratory Data Analysis (EDA)

- **Dataset Structure:** The shape of the datasets was examined (`df_train.shape`, `df_test.shape`), revealing the size and dimensionality of the data.
- **Data Types and Cleaning:** The `.info()` method identified the data types and presence of null values. An irrelevant 'Unnamed: 0' column was identified and removed from both datasets.
- **Date Conversion:** Key date columns were converted to datetime format for effective time-series analysis.

Data Preprocessing

Data preprocessing involved cleaning and formatting data correctly. The removal of irrelevant columns and conversion of date fields to datetime format set the stage for more detailed analysis.

Cost-Benefit Analysis

In this section, we evaluate the cost-benefit implications of implementing our credit card fraud detection model. This analysis considers the financial impact of correctly identifying fraudulent transactions (true positives) versus the costs associated with false positives (legitimate transactions incorrectly flagged as fraud) and false negatives (fraudulent transactions not detected).

Assumptions:

- **Cost of False Negatives:** The cost of a false negative is assumed to be the value of the fraudulent transaction, as this amount is typically lost in a fraud scenario.
- **Cost of False Positives:** The cost of false positives is harder to quantify but includes customer dissatisfaction, potential loss of clientele, and administrative costs involved in reviewing flagged transactions.
- **Savings from True Positives:** Successfully detected frauds prevent loss of transaction value and also save on potential legal and investigative costs.

Estimation:

- For this analysis, we use the model's performance metrics (precision, recall) to estimate the number of true positives, false positives, and false negatives.
- Based on the average transaction value in the dataset and estimated costs for false positives and negatives, we calculate the total cost/savings.

Results:

- **Total Savings from True Positives:** Calculated as the number of true positives multiplied by the average transaction value.
- **Total Cost from False Positives and Negatives:** Sum of costs from false positives and negatives.
- **Net Benefit:** The difference between total savings and total costs, providing a financial perspective on the model's effectiveness.

Conclusion

The EDA phase of this project lays the groundwork for subsequent in-depth analysis and model building. The next steps will involve applying statistical methods and machine learning algorithms to detect patterns indicative of fraud.