

# Credit Card Fraud Detection Analysis

Objective is to build a predictive model using machine learning to classify credit card transactions as fraudulent or non-fraudulent. I employed an XGBoost model, a popular gradient boosting algorithm (combines weak model to create more accurate predictive model).

The dataset consists of transactions with 31 features and a target variable indicating whether the transaction was fraudulent (Class = 1) or not (Class = 0). The dataset is highly imbalanced, with non-fraudulent transactions significantly outnumbering fraudulent ones.

 by Sonu Gupta



# Exploratory Data Analysis

## Class Distribution

Severe class imbalance, with a large majority of non-fraudulent transactions. This imbalance can affect the model's performance and requires addressing.

## Feature Analysis

Summary statistics indicate that data preprocessing (e.g., scaling) is necessary. Outlier detection shows potential for extreme values.

## Missing Values

No missing values detected in the dataset, ensuring the data is ready for modeling.

## Correlation Analysis

A correlation heatmap was generated to identify potential relationships between features.

```
## Rows: 284,807
## Columns: 31
## $ Time    <dbl> 0, 0, 1, 1, 2, 2, 4, 7, 7, 9, 10, 10, 10, 11, 12, 12, 12, 13, 1...
## $ V1      <dbl> -1.3598071, 1.1918571, -1.3583541, -0.9662717, -1.1582331, -0.4...
## $ V2      <dbl> -0.07278117, 0.26615071, -1.34016307, -0.18522601, 0.87773675, ...
## $ V3      <dbl> 2.53634674, 0.16648011, 1.77320934, 1.79299334, 1.54871785, 1.1...
## $ V4      <dbl> 1.37815522, 0.44815408, 0.37977959, -0.86329128, 0.40303393, -0...
## $ V5      <dbl> -0.33832077, 0.06001765, -0.50319813, -0.01030888, -0.40719338,...
## $ V6      <dbl> 0.46238778, -0.08236081, 1.80049938, 1.24720317, 0.09592146, -0...
## $ V7      <dbl> 0.239598554, -0.078802983, 0.791460956, 0.237608940, 0.59294074...
## $ V8      <dbl> 0.098697901, 0.085101655, 0.247675787, 0.377435875, -0.27053267...
## $ V9      <dbl> 0.3637870, -0.2554251, -1.5146543, -1.3870241, 0.8177393, -0.56...
## $ V10     <dbl> 0.09079417, -0.16697441, 0.20764287, -0.05495192, 0.75307443, -...
## $ V11     <dbl> -0.55159953, 1.61272666, 0.62450146, -0.22648726, -0.82284288, ...
## $ V12     <dbl> -0.61780086, 1.06523531, 0.06608369, 0.17822823, 0.53819555, 0...
## $ V13     <dbl> -0.99138985, 0.48909502, 0.71729273, 0.50775687, 1.34585159, -0...
## $ V14     <dbl> -0.31116935, -0.14377230, -0.16594592, -0.28792375, -1.11966983...
## $ V15     <dbl> 1.468176972, 0.635558093, 2.345864949, -0.631418118, 0.17512113...
## $ V16     <dbl> -0.47040053, 0.46391704, -2.89008319, -1.05964725, -0.45144918,...
## $ V17     <dbl> 0.207971242, -0.114804663, 1.109969379, -0.684092786, -0.237033...
## $ V18     <dbl> 0.02579058, -0.18336127, -0.12135931, 1.96577500, -0.03819479, ...
## $ V19     <dbl> 0.40399296, -0.14578304, -2.26185710, -1.23262197, 0.80348692, ...
## $ V20     <dbl> 0.25141210, -0.06908314, 0.52497973, -0.20803778, 0.40854236, 0...
## $ V21     <dbl> -0.018306778, -0.225775248, 0.247998153, -0.108300452, -0.00943...
## $ V22     <dbl> 0.277837576, -0.638671953, 0.771679402, 0.005273597, 0.79827849...
## $ V23     <dbl> -0.110473910, 0.101288021, 0.909412262, -0.190320519, -0.137458...
## $ V24     <dbl> 0.06692807, -0.33984648, -0.68928096, -1.17557533, 0.14126698, ...
## $ V25     <dbl> 0.12853936, 0.16717040, -0.32764183, 0.64737603, -0.20600959, -...
## $ V26     <dbl> -0.18911484, 0.12589453, -0.13909657, -0.22192884, 0.50229222, ...
## $ V27     <dbl> 0.133558377, -0.008983099, -0.055352794, 0.062722849, 0.2194222...
## $ V28     <dbl> -0.021053053, 0.014724169, -0.059751841, 0.061457629, 0.2151531...
## $ Amount  <dbl> 149.62, 2.69, 378.66, 123.50, 69.99, 3.67, 4.99, 40.80, 93.20, ...
## $ Class   <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
```

# Model Building: Data Preprocessing

## Scaling and Centering

Applied scaling and centering to normalize the features before feeding them into the model. This ensures that the model treats all features equally, regardless of their scale.

## XGBoost Model

An initial XGBoost model was trained with the default hyperparameters, yielding a high accuracy of 99.87%. The model achieved high sensitivity and specificity, indicating good performance in distinguishing between fraudulent and non-fraudulent transactions.

# Model Performance Metrics

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    0    1
##           0 56804   19
##           1   55   83
##
##           Accuracy : 0.9987
##           95% CI : (0.9984, 0.999)
##       No Information Rate : 0.9982
##       P-Value [Acc > NIR] : 0.00221
##
##           Kappa : 0.691
##
## Mcnemar's Test P-Value : 4.728e-05
##
##           Sensitivity : 0.9990
##           Specificity : 0.8137
##           Pos Pred Value : 0.9997
##           Neg Pred Value : 0.6014
##           Prevalence : 0.9982
##           Detection Rate : 0.9972
##       Detection Prevalence : 0.9976
##           Balanced Accuracy : 0.9064
##
##           'Positive' Class : 0
##
```

```
roc_auc_xgb_rose <- roc_auc_vec(as.factor(y_test), xgb_preds_rose)
print(paste("ROSE XGBoost AUC:", roc_auc_xgb_rose))
```

99.90%

Sensitivity

True Positive Rate

81.37%

Specificity

True Negative Rate

99.97%

Positive Predictive Value

60.14%

Negative Predictive Value

The model achieved high AUC, indicating good discrimination between fraud and non-fraud cases.



# Model Improvement Techniques

1

## Hyperparameter Tuning

Performed using cross-validation, adjusting for class imbalance by setting the `scale_pos_weight` parameter.

```
roc_auc_xgb_tuned <- roc_auc_vec(as.factor(y_test), xgb_preds_tuned)
print(paste("Tuned XGBoost AUC:", roc_auc_xgb_tuned))
```

```
## [1] "Tuned XGBoost AUC: 0.0259478296674022"
```

2

## Handling Class Imbalance

Applied ROSE technique (Random Over-Sampling Examples) to balance the training data, improving the model's recall for detecting fraudulent transactions.

```
roc_auc_xgb_rose <- roc_auc_vec(as.factor(y_test), xgb_preds_rose)
print(paste("ROSE XGBoost AUC:", roc_auc_xgb_rose))
```

```
## [1] "ROSE XGBoost AUC: 0.0575820683362249"
```

3

## Advanced Resampling Techniques

Explore SMOTE or ADASYN to generate synthetic examples of fraudulent transactions.

# Further Improvement Strategies



## Ensemble Methods

Combine XGBoost with Random Forest or LightGBM to leverage strengths of multiple models.



## Feature Engineering

Create interaction terms or extract new features to capture hidden patterns in the data.



## Neural Networks

Explore ANNs and deep learning techniques to capture complex non-linear relationships.

# Conclusion and Next Steps

## Model Success

The model successfully identifies fraudulent transactions with high accuracy, but there's room for improvement in handling class imbalance and increasing specificity.

## Future Enhancements

Apply advanced techniques such as feature engineering, ensemble models, and real-time learning to further enhance the model's performance.

## Practical Application

Adapt and refine the model for practical fraud detection applications in real-world scenarios.

