

Credit Card Fraud Detection

Agenda



Problem Definition & Societal Impact
Exploratory Data Analysis (EDA): Dataset Overview
EDA: Data Quality Checks
EDA: Feature Engineering



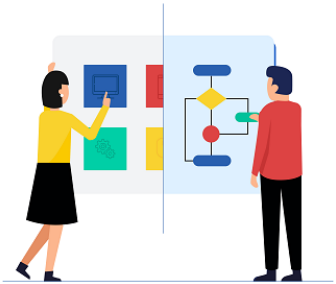
Model Selection and Rationale
Model Training, Optimization, & Validation
Model Performance Comparison
Insights from the Best Model



Deployment and Business Implications
Limitations and Future Improvements
Conclusion
References

Problem Definition & Societal Impact

Credit-card fraud costs about \$28 billion annually, erodes customer trust, and strains financial institutions—making automated, real-time detection essential.



Problem Overview

- Credit-card fraud is on the rise, causing roughly \$28 billion in annual losses.
- These incidents inflict both direct financial hits and extra operational costs on banks and merchants.
- Traditional, manual review processes can't keep pace with growing transaction volumes.
- Advanced, automated detection systems are needed for rapid identification and prevention of fraud.

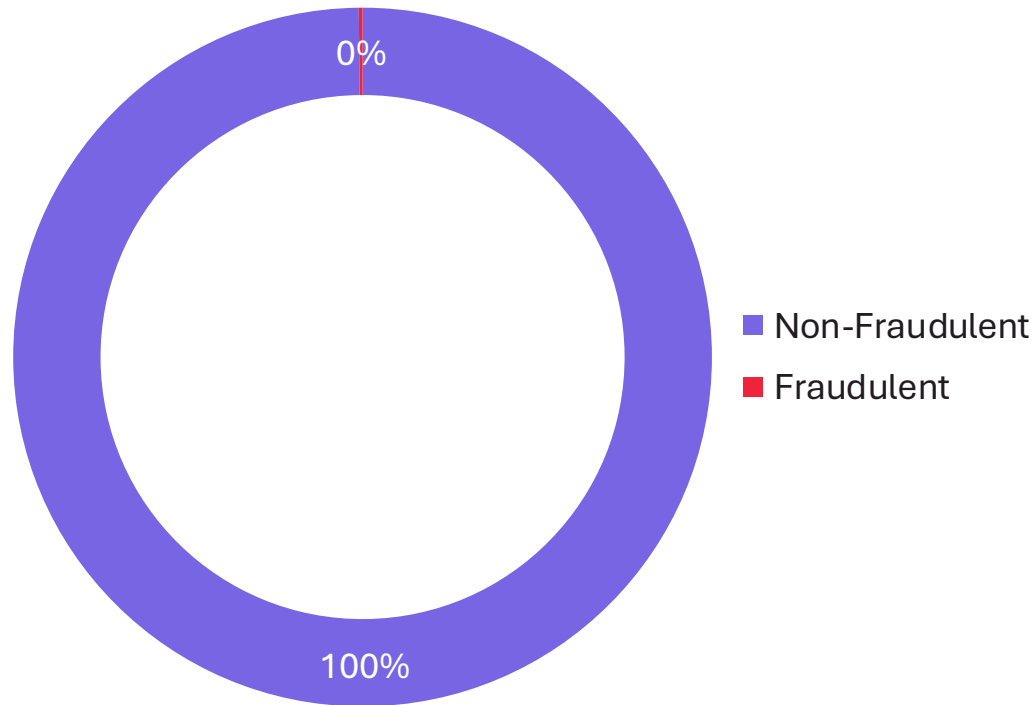


Societal and Economic Impact

- Erosion of customer trust in financial systems due to frequent fraud cases.
- Increased regulatory scrutiny and compliance costs for financial institutions.
- Higher transaction fees and insurance costs passed on to consumers.
- Potential for identity theft and broader financial crime beyond credit-card fraud.

Exploratory Data Analysis (EDA): Dataset Overview

- Class Distribution of Transactions



Dataset & Key Facts

Dataset contains 284,807 transactions collected over a 48-hour period.

Features include Time, Amount, and 28 anonymized PCA components (V1-V28).

Only 0.17% of transactions are labeled as fraudulent, indicating severe class imbalance.

This imbalance necessitates specialized techniques for effective fraud detection modeling.

EDA: Data Quality Checks

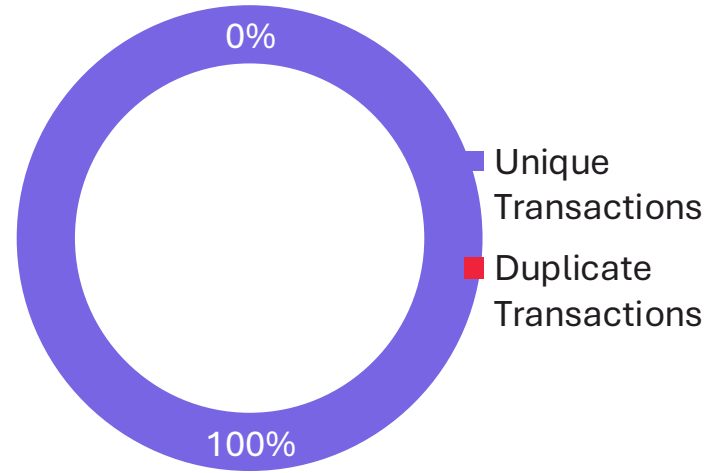
- Missing Values Check

■ Missing Values

Time	V2	V5	V8	V11	V14	V17	V20	V23	V26	Class
------	----	----	----	-----	-----	-----	-----	-----	-----	-------

No missing data detected across all features, ensuring data completeness for modeling.

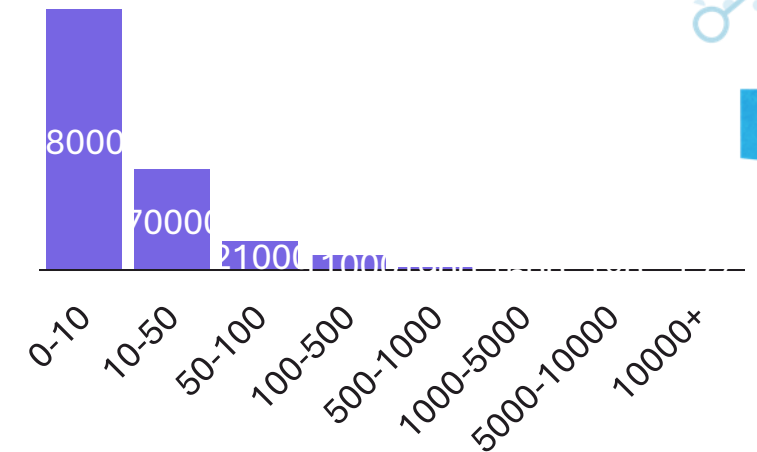
- Duplicate Transactions



No duplicate transactions found, confirming unique transaction records in the dataset.

- Outliers and Correlation Analysis

■ Frequency



Outliers in transaction amounts are retained as valid; PCA features show minimal correlation due to dimensionality reduction design.

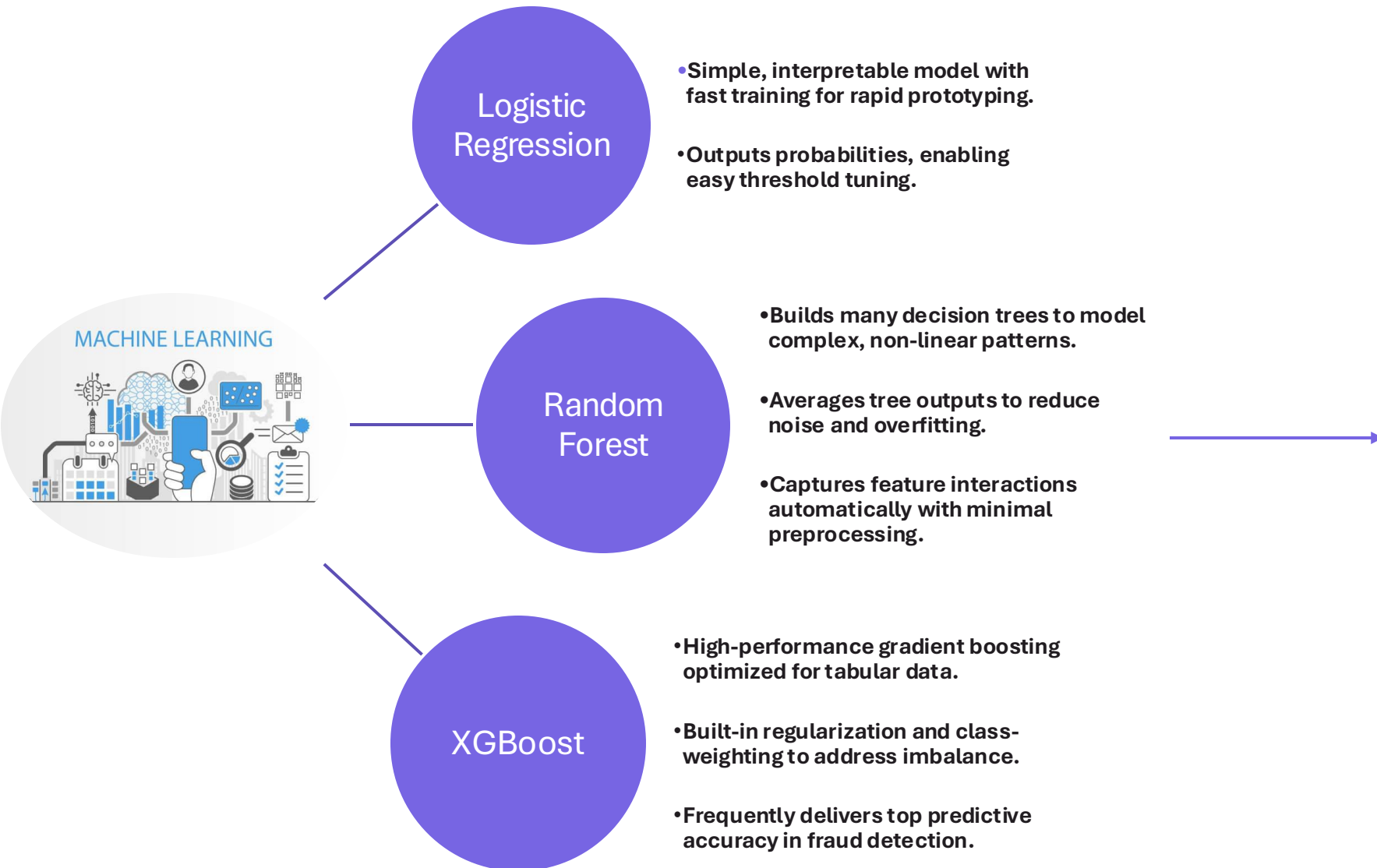
EDA: Feature Engineering

Applied log transformation to 'Amount' feature to reduce right skewness and improve normality of transaction values.

Scaled 'Time' and log-transformed 'Amount' to have zero mean and unit variance for consistent model input.

Included PCA components V1-V28 along with scaled Time and log-Amount features for model training.

Model Selection and Rationale



- We selected Logistic Regression for baseline interpretability.
- Random Forest for handling non-linearities and noise, and XGBoost for superior performance and imbalance handling.
- Evaluation focuses on precision, recall, F1-score, and ROC-AUC to balance fraud detection and false alarms.

Model Training, Optimization, & Validation

Effective model training and optimization using stratified splits and robust hyperparameter tuning methods significantly enhance fraud detection accuracy. Careful validation and threshold tuning ensure the model balances precision and recall for real-world deployment.

Activities

• Training/Test Split

- Dataset split into 70% training and 30% testing sets with stratification to maintain class distribution.
- Ensures balanced representation of fraudulent and non-fraudulent transactions in both sets.
- Prevents data leakage and supports reliable model evaluation.

Deliverables

- Stratified training and testing datasets
- Class distribution report

• Hyperparameter Tuning

- RandomizedSearchCV applied to Random Forest to explore combinations of max depth, number of trees, and minimum samples split.
- Optimizes model complexity and generalization to avoid overfitting and underfitting.
- Speeds up the search process compared to exhaustive grid search.

- Optimized hyperparameters for Random Forest
- Tuning performance metrics

• Validation Strategy

- Implemented 3-fold stratified cross-validation to assess model stability across different subsets.
- Ensures consistent performance by averaging results over multiple folds.
- Maintains class balance in each fold for reliable fraud detection evaluation.

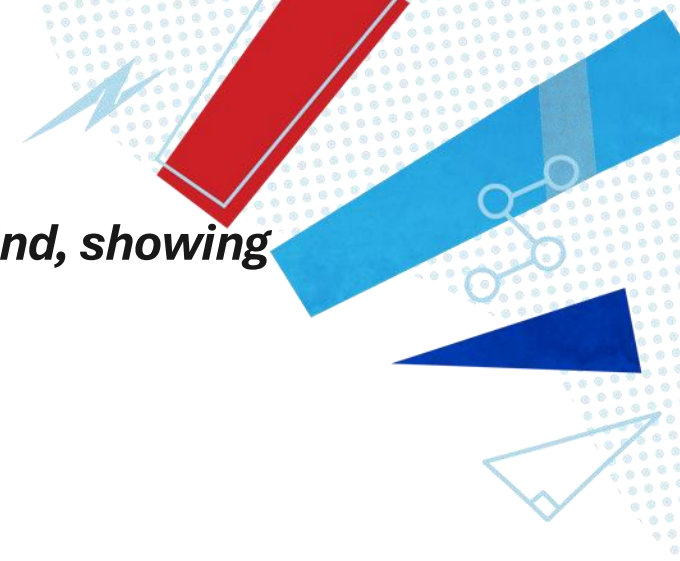
- Cross-validation results
- Performance consistency report

• Threshold Tuning

- Precision-Recall curve used to select optimal probability threshold rather than default 0.5.
- Threshold adjusted to maximize F1-score, balancing precision and recall effectively.
- Crucial for minimizing false positives and false negatives in fraud detection.

- Optimal probability threshold
- Precision-Recall curve analysis

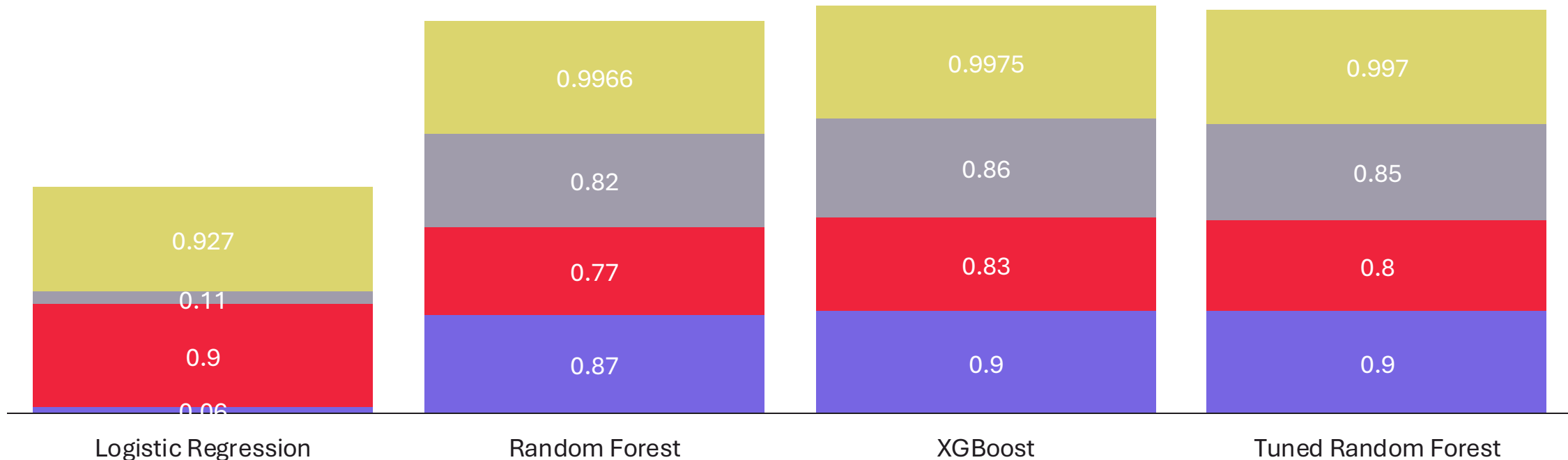
Model Performance Comparison



XGBoost achieves top performance, with tuned Random Forest close behind, showing the value of threshold optimization.

Performance Metrics of Fraud Detection Models

Precision Recall F1-Score ROC-AUC



Insights from the Best Model

- Model Insights & Threshold

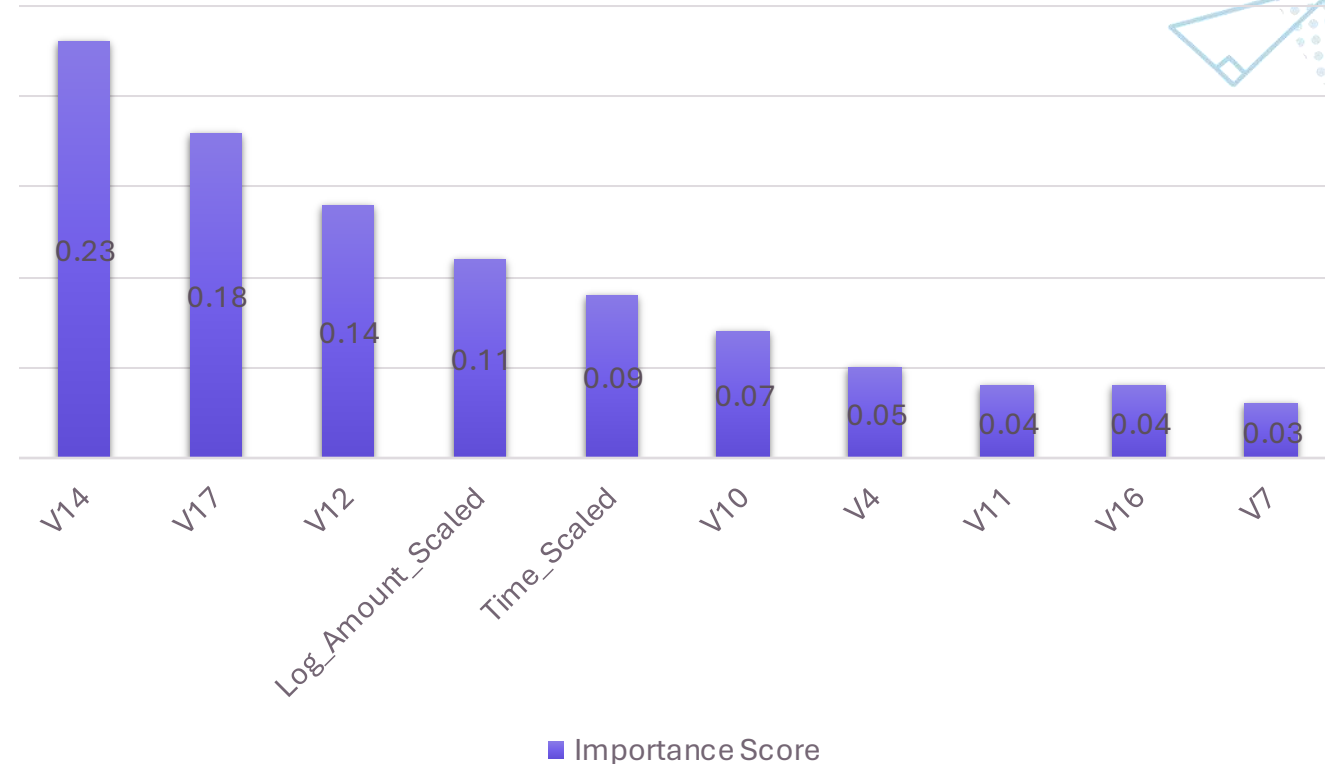
Top features influencing fraud prediction include PCA components V14, V17, V12, along with scaled Log Amount and Time variables.

An operational probability threshold of 0.70 maximizes the F1-score, balancing precision and recall effectively.

High precision reduces false positives, minimizing unnecessary alerts to customers and investigators.

High recall ensures the majority of fraudulent transactions are detected, enhancing security and reducing losses.

- Feature Importance from Tuned Random Forest Model



Deployment and Business Implications



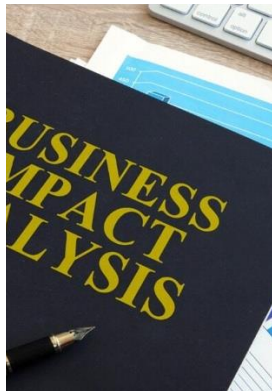
Real Time Fraud Detection

The fraud detection model scores transactions in milliseconds, enabling immediate assessment and action within the payment processing workflow.



Decision Framework

Transactions with a fraud probability above 0.85 are automatically blocked, while those between 0.70 and 0.85 undergo human review to balance security and customer convenience.



Business Impact

This approach significantly reduces manual review workload, accelerates fraud response times, and enhances customer experience by minimizing false positives and fraud losses.

Limitations and Future Improvements



The current credit-card fraud detection model faces challenges such as reduced interpretability due to PCA anonymization and the risk of model drift from evolving fraud patterns. Future improvements include integrating temporal features, advanced oversampling, and ensemble anomaly detection techniques to enhance accuracy and adaptability

- Data Limitations
 - **PCA anonymization obfuscates original feature meanings, limiting interpretability of model decisions.**
 - **Lack of raw feature transparency makes it challenging to understand specific fraud patterns directly.**
 - **Data collected over a short 48-hour period may not capture long-term trends and seasonal fraud variations.**
- Model Drift and Maintenance
 - Fraud tactics continuously evolve, causing potential degradation in model performance over time.
 - Periodic retraining and monitoring are essential to address changes in fraud behavior.
 - Without updates, the model risks increasing false negatives or false positives, reducing trust.
- Future Enhancements
 - Incorporate additional temporal features such as hour of day and day of week to capture time-based fraud patterns.
 - Apply advanced oversampling techniques like *SMOTE* and *ADASYN* to better address class imbalance.
 - Leverage ensemble methods combining multiple models and anomaly detection approaches like autoencoders for improved detection robustness.

Conclusion

The Tuned Random Forest model offers an effective, scalable solution for credit-card fraud detection with strong performance metrics and practical deployment guidelines

Best Model Performance

Tuned XGBoost emerged as the best performing model with an F1-score of approximately 0.85 and ROC-AUC near 0.997, indicating excellent balance between precision and recall.

Optimal Threshold

The optimized operational threshold of 0.70 maximizes fraud detection effectiveness while minimizing false alarms, ensuring practical usability.

Real Time Detection

Deployment strategy includes automatic blocking for high-probability fraud cases and human review for borderline cases, optimizing resource use and customer experience.

Deployment Strategy

The model supports real-time fraud detection, enabling rapid transaction scoring and timely decision-making.

References

- Kaggle Credit Card Fraud Detection dataset: <https://www.kaggle.com/mlg-ulb/creditcardfraud>
- Breiman, L. (2001). "Random Forests." Machine Learning, 45(1), 5-32.
- Chen, T., & Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System." Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). "SMOTE: Synthetic Minority Over-sampling Technique." Journal of Artificial Intelligence Research, 16, 321-357.



Thank You