

Fraud & Risk Analytics: Credit Card Fraud Detection using Machine Learning

Business Context

Credit card fraud is a critical risk for financial institutions, resulting in **direct financial losses, increased operational costs, and erosion of customer trust**. A key challenge in fraud detection systems is the **extreme class imbalance** present in real-world transaction data—fraudulent transactions typically represent **less than 0.2%** of total volume.

In such environments, naïve accuracy-driven models fail, as they may appear “accurate” while missing fraud or unnecessarily blocking genuine customers. The real business challenge is to **detect fraud effectively while minimizing false positives**, thereby balancing **risk control with customer experience**.

Project Objective

The objective of this project was to design and evaluate machine learning models for **credit card fraud detection**, with a focus on **business-relevant performance metrics rather than raw accuracy**.

Specifically, the project aimed to:

- Accurately identify fraudulent transactions
- Compare an interpretable baseline model with a more powerful ensemble model
- Evaluate models using metrics aligned with real-world fraud operations
- Recommend a **production-ready solution** based on **business impact**, not theoretical performance

Dataset Overview

The analysis was conducted on a publicly available dataset containing approximately **284,000 anonymized credit card transactions**, with a fraud rate of **~0.17%**.

Key dataset characteristics:

- Highly imbalanced target variable (fraud vs genuine)
- Features transformed using **Principal Component Analysis (PCA)** to preserve customer privacy
- Transaction amounts retained in original scale to capture spending behavior

The PCA transformation ensures confidentiality while maintaining the **predictive structure required for modeling**.

Modeling Approach

Two machine learning models were implemented and compared:

Logistic Regression (Baseline Model)

- Selected for its **interpretability and transparency**
- Commonly preferred in regulated financial environments
- Serves as a benchmark to evaluate more complex models

Random Forest (Ensemble Model)

- Captures **non-linear relationships and interaction effects**
- Well-suited for complex fraud patterns that evolve over time
- More robust to noise and feature interactions

This comparison reflects a **real-world modeling trade-off** between explainability and predictive power.

Handling Class Imbalance & Evaluation Metrics

Given the extreme imbalance in fraud data, several techniques were applied to ensure meaningful evaluation:

- **Class-weighted learning** to penalize fraud misclassification
- **Stratified train-test split** to preserve fraud distribution
- Evaluation focused on:
 - Precision
 - Recall
 - F1-Score
 - ROC-AUC

Accuracy was intentionally deprioritized, as it is **misleading in imbalanced classification problems**.

Results & Business Insights

Logistic Regression

- Achieves **very high recall**, capturing most fraudulent transactions
- Generates a **large number of false positives**
- Would result in:
 - High customer friction
 - Increased manual review workload
 - Operational inefficiency

Random Forest

- Maintains strong fraud detection capability
- **Dramatically reduces false positives**
- Produces a better balance between:
 - Fraud prevention
 - Customer experience
 - Operational scalability

Confusion matrices included in this report visually demonstrate the trade-off between the two models.

Business Impact

From a business perspective, **false positives are often more damaging than missed fraud** due to:

- Customer dissatisfaction from declined transactions
- Increased call-center and investigation costs
- Loss of long-term customer trust

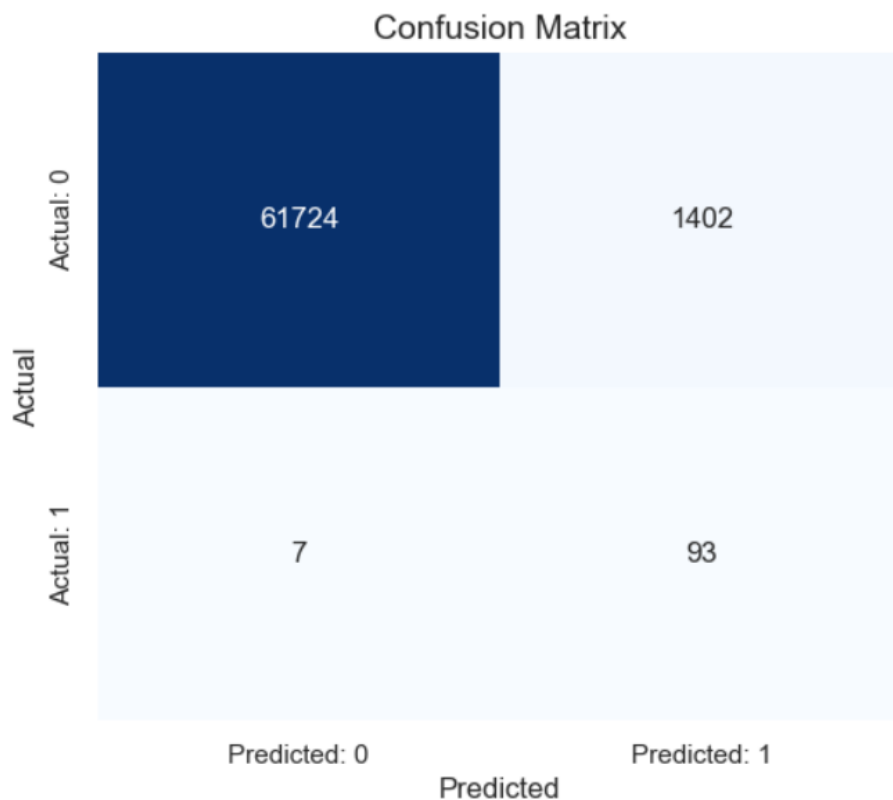
By significantly reducing false alarms while maintaining fraud detection strength, the Random Forest model offers:

- Lower customer friction
- Reduced operational cost
- Higher scalability for real-world deployment

Final Recommendation

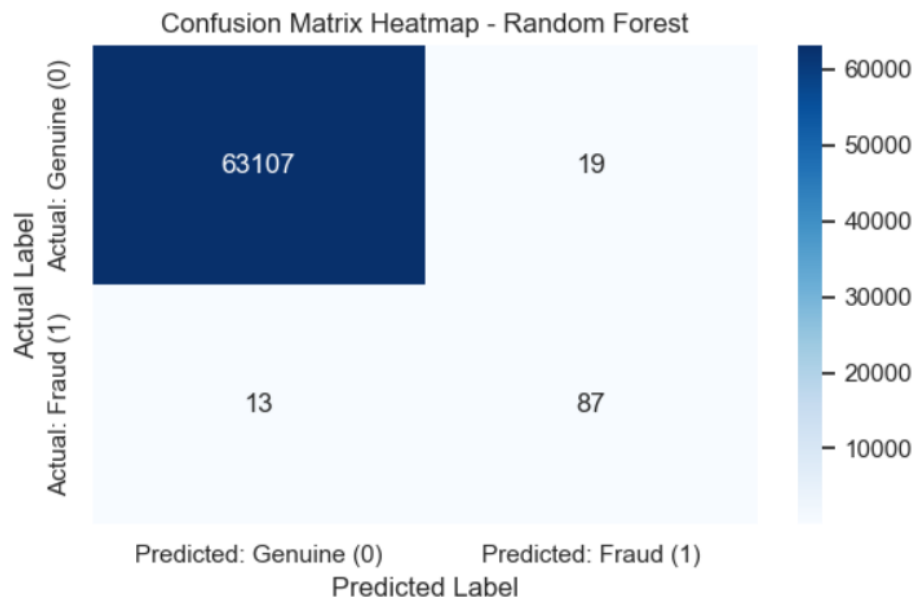
Based on both **technical performance and business considerations**, the **Random Forest model** is recommended for production deployment.

It provides the optimal trade-off between **fraud prevention effectiveness and customer experience**, making it suitable for real-world financial systems where trust, efficiency, and scalability are critical.



Classification Report:				
	precision	recall	f1-score	support
0	1.000	0.978	0.989	63126
1	0.062	0.930	0.117	100
accuracy			0.978	63226
macro avg	0.531	0.954	0.553	63226
weighted avg	0.998	0.978	0.987	63226

The Logistic Regression model achieves very high recall for fraudulent transactions, successfully identifying most fraud cases. However, it generates a **large number of false positives**, incorrectly flagging genuine transactions as fraud. While effective for fraud capture, this behaviour would lead to **high customer friction and increased manual review costs** in a real-world deployment.



The Random Forest model significantly reduces false positives while maintaining strong fraud detection capability. This results in fewer genuine customers being incorrectly blocked, improving customer experience while still preventing fraud. The model demonstrates a **better balance between fraud prevention and operational efficiency**, making it more suitable for production deployment.

Business Comparison Summary

Aspect	Logistic Regression	Random Forest
Fraud Detection	Very High	High
False Positives	Very High	Low
Customer Friction	High	Low
Operational Cost	High	Lower
Production Suitability	✗ No	✓ Yes