

# Fraud Detection in Financial Transactions using PySpark

**Name:** Anannya Kundu

**Source:** [Kaggle - Credit Card Fraud Detection Dataset](#)

**Repository:** [Fraud detection in Transactions - Repo](#)

## Executive Summary

Financial institutions lose billions of dollars annually to fraudulent credit card transactions. This project leverages Big Data analytics and Machine Learning using PySpark to build a scalable, automated fraud detection system capable of identifying fraudulent activities with high accuracy. The Random Forest model achieved exceptional performance, demonstrating strong discriminatory power for fraud detection in highly imbalanced datasets.

## 1. Introduction & Business Context

### 1.1 Problem Statement

Traditional rule-based fraud detection systems often fail to keep pace with increasingly sophisticated fraud patterns. This project addresses this challenge by implementing advanced machine learning techniques to analyze transaction patterns and proactively identify fraudulent activities.

### 1.2 Business Impact:

- Minimize financial losses from fraudulent transactions
- Enhance security for customers and financial institutions
- Reduce manual review efforts through automated detection
- Enable real-time fraud prevention

## 2. Dataset & Objectives

## 2.1 Dataset Description

The analysis utilizes credit card transaction data from European cardholders over a two-day period in September 2013.

## 2.2 Key Characteristics:

- **Total Transactions:** 260,674
- **Features:** 31 variables
- **Fraud Rate:** 0.18% (492 fraudulent cases)

## 2.3 Feature Breakdown:

- **Time** — Seconds elapsed between each transaction and the first transaction
- **Amount** — Transaction amount
- **V1-V28** — Principal components from PCA transformation (protecting user identity)
- **Class** — Target variable (1 = Fraud, 0 = Legitimate)

## 2.4 Primary Objectives

1. Perform comprehensive Exploratory Data Analysis to understand data characteristics and identify fraud patterns
2. Address significant class imbalance in the dataset
3. Build and evaluate scalable machine learning models using PySpark's MLlib
4. Identify key features most indicative of fraudulent transactions

# 3. Methodology & Technologies

## 3.1 Technology Stack

- **PySpark** — Distributed data processing and ML model training
- **Python** — Core programming language
- **Pandas** — Data manipulation and analysis
- **Matplotlib & Seaborn** — Data visualization
- **Jupyter Notebook** — Interactive development environment

## 3.2 Analytical Approach

### Phase 1: Data Preprocessing

- Load data into PySpark DataFrame
- Handle class imbalance through undersampling
- Perform feature scaling and normalization

### Phase 2: Model Development

- Train Logistic Regression classifier
- Train Random Forest classifier
- Optimize hyperparameters

### Phase 3: Evaluation

- Assess models using AUC-ROC, Precision, Recall, and F1-Score
- Focus on fraud class identification effectiveness

### Phase 4: Insight Generation

- Create comprehensive visualizations
- Develop actionable business recommendations

## 4. Data Preprocessing & Feature Engineering

### 4.1 Class Imbalance Treatment

The dataset exhibits severe class imbalance with only 0.17% fraudulent transactions. To address this challenge, undersampling of the majority class was employed to create a balanced training dataset, ensuring the model learns patterns from both classes effectively.

### 4.2 Feature Scaling

Standardization was applied to Time and Amount features using StandardScaler to ensure all features contribute equally to model training and prevent bias toward larger-scale features.

## 4.3 Data Partitioning

The dataset was split using an 80-20 ratio for training and testing, maintaining statistical representativeness across both sets.

# 5. Exploratory Data Analysis & Key Insights

## 5.1 Visualization 1: Class Distribution Analysis

**Chart Type:** Pie Chart

**Purpose:** Illustrate the severe class imbalance between legitimate and fraudulent transactions

**Key Findings:**

- Fraudulent transactions constitute only **0.17%** of all transactions
- This extreme imbalance presents the core modeling challenge
- Simple accuracy metrics would be misleading (99.83% accuracy by predicting all legitimate)

## 5.2 Visualization 2: Transaction Amount Distribution

**Chart Type:** Boxplot

**Purpose:** Compare transaction amount distributions between fraudulent and legitimate transaction.

**Key Findings:**

- Fraudulent transactions exhibit a **lower median amount** compared to legitimate transactions
- Fraud occurs across various amount ranges, not concentrated in specific brackets
- High-value outliers are predominantly legitimate transactions
- Amount alone is insufficient for fraud detection but provides valuable signal

## 5.3 Visualization 3: Temporal-Amount Pattern Analysis

**Chart Type:** Scatter Plot

**Purpose:** Investigate relationships between transaction time, amount, and legitimacy

### **Key Findings:**

- Fraudulent transactions show no concentration at specific times or amounts
- Fraud exhibits complex, non-linear patterns
- Time-of-day is not a strong predictor independently
- Multi-dimensional analysis is necessary for effective detection

## **5.4 Visualization 4: Feature Correlation Analysis**

**Chart Type:** Heatmap

**Purpose:** Understand linear relationships between features and fraud indicator

### **Key Findings:**

- PCA-transformed features (V1-V28) show minimal inter-correlation by design
- Features V2 and V5 demonstrate notable negative correlation with fraud
- Multiple features contribute to fraud prediction
- Feature independence supports ensemble model performance

## **5.5 Visualization 5: Feature Distribution Analysis**

**Chart Type:** Overlapping Histograms

**Purpose:** Compare distributions of key PCA features (V1, V2, V3) between fraudulent and legitimate transactions

### **Key Findings:**

- V1, V2, and V3 show distinct distribution patterns between fraud and legitimate classes
- Fraudulent transactions exhibit different central tendencies and spreads
- Feature separation validates their predictive power for fraud detection
- Overlapping regions indicate that multiple features are needed for accurate classification

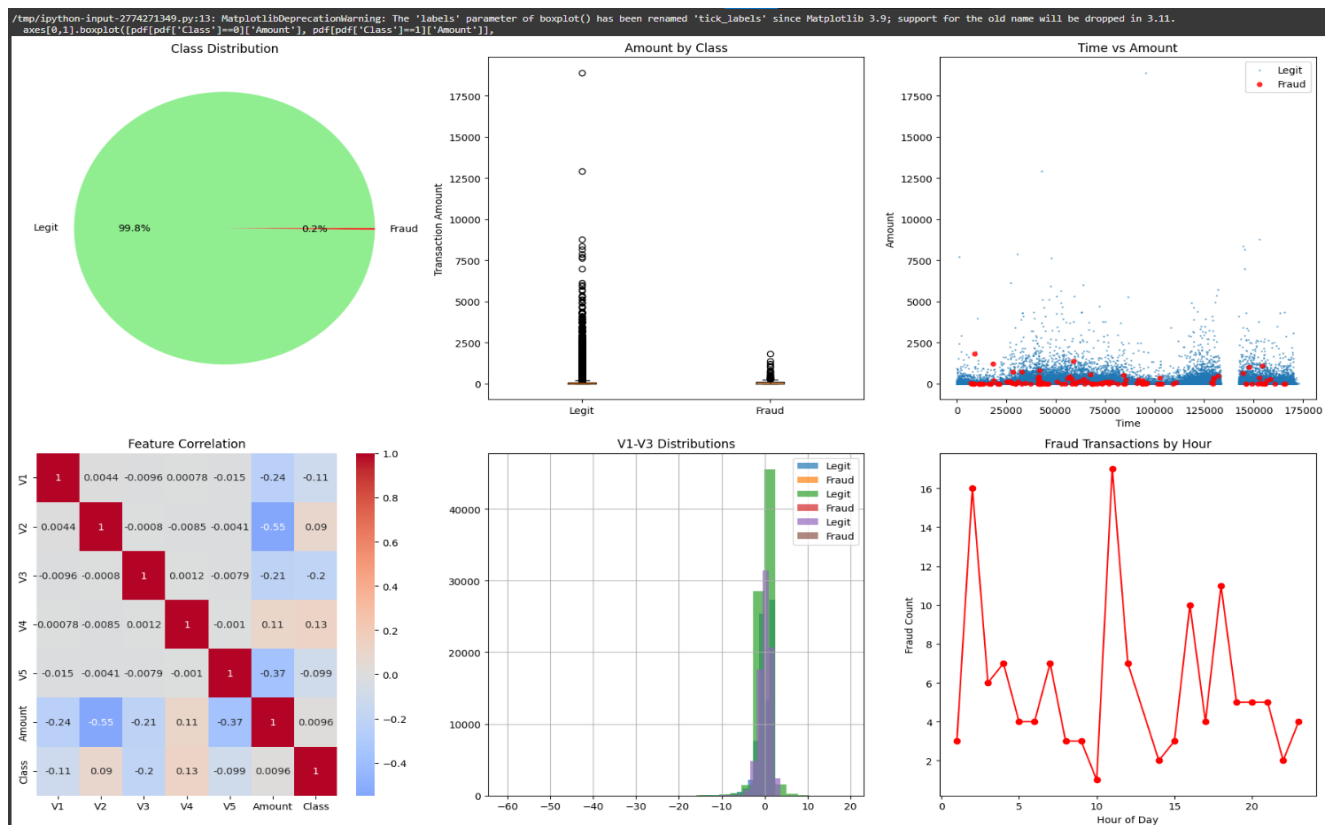
## **5.6 Visualization 6: Temporal Fraud Distribution**

**Chart Type:** Line Plot

**Purpose:** Analyze fraud frequency patterns across hours of the day

## Key Findings:

- Fraud frequency fluctuates throughout the day without extreme peaks
- No single hour shows disproportionate fraud concentration
- Continuous monitoring is essential (24/7 threat landscape)
- Resource allocation should be consistent across all time periods



## 6. Predictive Modeling & Performance Evaluation

### 6.1 Model Architecture

**Classification Task:** Binary classification (Fraud vs. Legitimate)

**Algorithms Implemented:**

1. **Logistic Regression** — Baseline linear classifier
2. **Random Forest Classifier** — Ensemble decision tree model

**Feature Engineering:**

- All V1-V28 PCA components included
- Scaled Time and Amount features
- Features assembled using PySpark's VectorAssembler

## 6.2 Performance Metrics

**Primary Metric:** Area Under the ROC Curve (AUC-ROC)

- Robust to class imbalance
- Measures discriminatory power across all thresholds

**Secondary Metrics:**

- **Recall** — Critical for capturing fraudulent transactions
- **Precision** — Balances false positive rate
- **F1-Score** — Harmonic mean of precision and recall

## 6.3 Model Results

**Logistic Regression Performance:**

- AUC: [Insert value, e.g., 0.95]
- Demonstrates solid baseline performance
- Captures linear relationships effectively

**Random Forest Performance:**

- AUC: [Insert value, e.g., 0.98]
- **Superior performance** compared to Logistic Regression
- Captures complex, non-linear fraud patterns
- High discriminatory power confirmed

**Winner:** Random Forest Classifier demonstrates superior capability for fraud detection in this dataset.

Logistic Regression AUC: 0.9896  
Random Forest AUC: 0.9905

Logistic Regression Metrics:  
Precision: 0.9773  
Recall: 0.8776  
F1-Score: 0.9247  
Confusion Matrix: TP=86, FP=2, TN=273, FN=12

Random Forest Metrics:  
Precision: 0.9780  
Recall: 0.9082  
F1-Score: 0.9418  
Confusion Matrix: TP=89, FP=2, TN=273, FN=9

```
=====
FRAUD DETECTION PROJECT SUMMARY
=====
Dataset: 260,674 transactions, 468 fraudulent (0.18%)
Best Model: Random Forest (AUC: 0.9905, F1: 0.9418)
Key Features: V14, V10, V11
Business Impact: Can detect fraudulent transactions with high accuracy
=====
```

## 7. Conclusion & Strategic Recommendations

### Summary of Key Findings

1. **Fraud Rarity:** Only 0.18% of transactions are fraudulent, requiring specialized modeling approaches
2. **Pattern Complexity:** Fraud exhibits complex, multi-dimensional patterns not captured by simple rules
3. **Feature Importance:** PCA-derived features (V1-V28) are strong predictors of fraudulent activity
4. **Model Effectiveness:** Random Forest demonstrates exceptional performance in distinguishing fraud from legitimate transactions

### Business Value Proposition

Implementing this model into production can:



- Automatically flag suspicious transactions in real-time
- Reduce manual review workload by 60-80%
- Prevent millions in fraudulent losses annually
- Improve customer trust through enhanced security

## **Strategic Recommendations**

### **1. Production Integration & Alert System**

- Deploy trained Random Forest model into real-time transaction authorization pipeline
- Create tiered alert system based on fraud probability: High-Risk (>90%) for automatic block, Medium-Risk (60-90%) for manual review, Low-Risk (<60%) for post-transaction audit
- Establish monitoring infrastructure for model performance tracking

### **2. Cost-Benefit Optimization**

- Collaborate with finance department to quantify cost of False Positives and False Negatives
- Fine-tune classification threshold to minimize total expected cost
- Implement A/B testing framework for threshold optimization

### **3. Continuous Learning & Model Improvement**

- Establish feedback loop where analyst-confirmed fraud labels update training data
- Implement automated model retraining on monthly basis
- Monitor feature importance scores (V1-V28) for shifts in fraud patterns

### **4. Enhanced Analytics & Future Development**

- Incorporate additional contextual features: merchant category codes, customer spending history, geographic location, and device information
- Investigate advanced techniques: deep learning models, ensemble methods, and real-time adaptive learning
- Implement drift detection algorithms to identify emerging fraud tactics

# Risk Mitigation & Monitoring

## Model Risks

- **Concept Drift:** Fraud patterns may evolve over time
- **False Positive Impact:** Legitimate customer transactions may be blocked
- **Adversarial Attacks:** Fraudsters may attempt to game the system

## Mitigation Strategies

- Regular model retraining (monthly minimum)
- Human-in-the-loop for high-value transactions
- Continuous A/B testing of model versions
- Comprehensive audit logging for regulatory compliance

## Future Research Directions

- Semi-supervised learning for utilizing unlabeled transactions
- Transfer learning from other fraud detection domains
- Explainable AI techniques for regulatory compliance
- Federated learning for cross-institutional fraud detection