

The Pattern of Probity vs. Plunder: Unmasking Fraud with Machine Learning

Introduction: The Hidden Battle Against Financial Fraud

In the digital age, financial fraud has evolved into a sophisticated cat-and-mouse game. Every transaction carries with it a story—a pattern that can reveal whether it's a legitimate purchase or a carefully disguised act of deception. This blog post takes you on a comprehensive journey through a fraud detection analysis, where we craft an AI investigator capable of sifting through routine transaction details to unmask fraudulent behavior.

The Challenge: Imagine being tasked with building a system that can examine mundane transaction attributes—cardholder age, transaction amount, transaction type—and identify the subtle behavioral signatures that betray fraudulent intent. This is not just a technical exercise; it's about protecting people's financial security and trust in digital transactions.

In machine learning terms, this is a supervised learning classification task focused on feature engineering and pattern recognition. Our goal? To predict whether a given transaction is fraudulent or legitimate based on its attributes.

The Dataset: Understanding Our Financial Landscape

Our analysis begins with a dataset of 50 transactions, each containing 6 key features:

- **TransactionID:** Unique identifier for each transaction
- **Amount:** The monetary value of the transaction
- **TransactionType:** Category of transaction (Online, ATM, or POS)
- **IsInternational:** Whether the transaction crossed borders (Yes/No)
- **CardHolderAge:** Age of the cardholder
- **FraudReported:** Our target variable (Yes/No)

A critical first observation: Our dataset contains 29 fraudulent transactions and 21 legitimate ones—a relatively balanced distribution, which is fortunate for model training.

Data Quality Assessment: Building on Solid Ground

Before diving into analysis, we performed essential data quality checks:

Missing Values: Zero missing values across all columns—our dataset is complete and ready for analysis.

Duplicate Rows: No duplicate entries found, ensuring each transaction represents a unique event.

This clean foundation allows us to proceed confidently to exploratory data analysis without the need for imputation or data cleansing.

Exploratory Data Analysis: Searching for Patterns

1. Two-Dimensional Scatter Plot Analysis

We began by examining relationships between individual features and fraud status:

Amount vs. Fraud: Fraudulent transactions occur across a wide range of amounts, with significant overlap with legitimate transactions. No clear separation emerges based solely on transaction amount.

CardHolderAge vs. Fraud: Both fraudulent and non-fraudulent transactions span across different age groups without distinct clustering. Age alone doesn't reveal a clear fraud pattern.

TransactionType vs. Fraud: Online transactions show a mix of both fraudulent and legitimate cases. While we see distribution differences across transaction types, no single type is exclusively associated with fraud.

IsInternational vs. Fraud: Both international and domestic transactions contain fraudulent instances, with no immediate strong separation.

Key Insight: Individual features don't provide linear separation between fraud and legitimate transactions, suggesting we'll need more sophisticated approaches or feature combinations.

2. Three-Dimensional Visualization: Adding Depth

To explore multi-feature relationships, we created a 3D scatter plot using Amount, CardHolderAge, and FraudReported. The visualization revealed that fraudulent (blue) and legitimate (red) transactions are largely intermingled in this 3D space. There's no clear plane or boundary that easily separates the two classes.

This reinforces our earlier observations: these two features alone are insufficient to linearly separate fraudulent from legitimate transactions. We need more complex relationships or additional features for effective classification.

3. Pair Plot Analysis: The Big Picture

Pair plots provided a comprehensive view of relationships between all numerical features:

Histograms: For both Amount and CardHolderAge, the distributions for fraudulent and non-fraudulent transactions show significant overlap. No clear separation in range or shape emerges.

Scatter Plots: Plots between feature pairs (Amount vs. CardHolderAge, etc.) show intermingled points for both classes without distinct clusters.

Conclusion: The pair plot reinforces that fraudulent and non-fraudulent transactions aren't easily separable based on simple linear relationships between these numerical features.

4. One-Dimensional Scatter Plots: Focused Analysis

Using strip plots, we examined the distribution of individual features:

Amount Distribution: Points for both fraud categories are intermingled across the range of amounts with no clear clustering at specific values.

Age Distribution: Fraudulent transactions spread across different age groups without concentration in distinct ranges.

These visualizations confirm that Amount and CardHolderAge individually lack clear linear separation power.

5. Distribution Analysis: PDF and CDF

We dove deeper into probability distributions:

CardHolderAge Analysis:

- PDF curves for fraudulent and non-fraudulent transactions show similar shapes with overlapping distributions
- CDF curves rise similarly for both groups
- Cumulative probability at any given age is comparable for both categories

Amount Analysis:

- PDF curves again show overlapping distributions
- No clear concentration of fraudulent transactions at specific amount ranges
- CDF curves show similar trends for both groups

Insight: Neither CardHolderAge nor Amount provides distinct distributional characteristics that clearly separate fraud from legitimate transactions.

6. Categorical Feature Analysis

Count plots revealed interesting patterns:

Transaction Type Distribution:

- Online transactions show the highest absolute count of both fraudulent and legitimate transactions
- ATM and POS transactions also contain both categories
- The relative proportions suggest potential relationships worth exploring

International Status:

- Both international and domestic transactions contain fraudulent instances
- The distribution doesn't suggest overwhelming fraud association with either category

7. Box Plot Analysis: Understanding Spread and Outliers

Amount Distribution:

- Interquartile ranges (IQR) overlap significantly between fraud categories
- Similar median values for both groups
- Some outliers present in both categories without clear differentiation

CardHolderAge Distribution:

- Overlapping IQRs and similar medians
- No distinct age group stands out as significantly more prone to fraud

Conclusion: Box plots confirm that Amount and CardHolderAge don't show significant distributional differences between fraud categories.

8. Violin Plot Analysis: Density Patterns

Violin plots combined box plot features with kernel density estimation:

Amount: Density distributions largely overlap, with similar shapes for both fraud statuses

CardHolderAge: Comparable density curves suggest similar age distributions for both categories

The violin plots reinforce that these features individually lack clear separation power.

Data Preprocessing: Preparing for Machine Learning

With our exploratory analysis complete, we prepared the data for modeling:

1. Encoding Categorical Variables

We applied one-hot encoding to TransactionType and IsInternational:

- TransactionType became: TransactionType_Online, TransactionType_POS
- IsInternational became: IsInternational_Yes
- We used drop_first=True to avoid multicollinearity

2. Feature and Target Definition

Features (X): Amount, CardHolderAge, and encoded categorical variables

Target (y): FraudReported (converted to 0 for No, 1 for Yes)

We excluded TransactionID as it's merely an identifier without predictive value.

3. Train-Test Split

We split the data using a 70-30 ratio:

- Training set: 35 transactions (5 features)
- Testing set: 15 transactions (5 features)
- Random state: 42 (for reproducibility)

This split allows us to train models and evaluate their performance on unseen data.

Model Training: The Battle of Algorithms

We trained and evaluated five classification algorithms:

1. Logistic Regression
2. Decision Tree
3. Random Forest
4. Support Vector Machine (SVM)
5. K-Nearest Neighbors (KNN)

Each model was trained on the training set and evaluated on the test set using four key metrics:

- Accuracy: Overall correctness
- Precision: Of predicted frauds, how many were actually fraud?
- Recall: Of actual frauds, how many did we catch?

- F1-Score: Harmonic mean of precision and recall

Model Performance: The Results

Here's how each model performed:

1. Logistic Regression

- Accuracy: 26.67%
- Precision: 50.00%
- Recall: 27.27%
- F1-Score: 35.29%

2. Decision Tree

- Accuracy: 46.67%
- Precision: 71.43%
- Recall: 45.45%
- F1-Score: 55.56%

3. Random Forest

- Accuracy: 46.67%
- Precision: 71.43%
- Recall: 45.45%
- F1-Score: 55.56%

4. Support Vector Machine

- Accuracy: 40.00%
- Precision: 66.67%
- Recall: 36.36%
- F1-Score: 47.06%

5. K-Nearest Neighbors

- Accuracy: 26.67%
- Precision: 50.00%
- Recall: 27.27%
- F1-Score: 35.29%

Best Performing Models: Decision Tree and Random Forest emerged as the top performers with identical results, achieving the highest F1-Score of 55.56% and the best Recall of 45.45%.

Key Insights and Learnings

1. Feature Complexity: Individual features don't provide linear separation between fraud and legitimate transactions. The overlapping distributions we observed in EDA explain why simple models struggle.

2. The Importance of Recall: In fraud detection, Recall is critically important—missing a fraudulent transaction (false negative) can be costly. Our best models caught 45.45% of fraudulent transactions.

3. Dataset Size Limitations: With only 50 transactions, our models have limited learning capacity. More data would likely improve performance significantly.

4. Feature Engineering Opportunities: The analysis suggests that derived features (combinations or interactions of existing features) might provide better separation.

5. Tree-Based Model Success: Decision Tree and Random Forest performed best, likely because they can capture non-linear relationships and feature interactions that linear models miss.

Practical Implications

For Financial Institutions:

- Multi-factor authentication and rule-based systems should complement ML models
- Focus on improving recall to catch more fraudulent transactions
- Continuously collect more diverse transaction data

For Data Scientists:

- Feature engineering is crucial when individual features lack separation power
- Consider ensemble methods and deep learning for complex pattern recognition
- Always validate models on test data to avoid overfitting

For Consumers:

- Fraudsters operate across all transaction types, amounts, and demographics
- No single factor guarantees safety; vigilance across all transactions matters
- Multi-layered security (2FA, alerts, monitoring) provides the best protection

Next Steps and Future Work

1. Data Expansion: Collect more transactions to improve model learning capacity

2. Advanced Feature Engineering:

- Time-based features (hour of day, day of week)
- Transaction velocity (frequency within time windows)

- Deviation from cardholder's typical behavior
- Geographic location patterns

3. Model Optimization:

- Hyperparameter tuning for Random Forest
- Try gradient boosting algorithms (XGBoost, LightGBM)
- Explore deep learning approaches (neural networks)

4. Ensemble Methods: Combine multiple models for improved performance

5. Real-time Implementation: Deploy models with streaming data pipelines

6. Explainability: Use SHAP or LIME to explain model predictions to stakeholders

Conclusion: The Ongoing Battle

Fraud detection is not a solved problem—it's an ongoing battle between fraudsters evolving their tactics and data scientists improving detection algorithms. This analysis demonstrates several critical truths:

- Fraud patterns are subtle and complex, not easily captured by individual features
- Machine learning offers powerful tools but requires careful feature engineering
- Model evaluation must prioritize the right metrics (Recall in fraud detection)
- Continuous improvement through more data and better features is essential

Our journey through this fraud detection analysis revealed that while Decision Tree and Random Forest models showed promise with 55.56% F1-Score, there's significant room for improvement. The overlapping distributions we observed throughout our exploratory analysis underscore the challenge: fraud looks remarkably similar to legitimate behavior in many dimensions.

The Pattern of Probit vs. Plunder is not about finding a silver bullet—it's about building layered defenses, continuously learning from new patterns, and staying one step ahead in an ever-evolving landscape.

As we move forward in the age of digital transactions, the tools and techniques explored here form the foundation for protecting financial systems and maintaining trust in the digital economy. The fight against fraud continues, armed with data, algorithms, and vigilance.

Technical Notes

Tools Used:

- Python 3
- Libraries: pandas, numpy, matplotlib, seaborn, scikit-learn
- Platform: Google Colaboratory
- Data: 50 transactions with 6 features

Code Availability: The complete analysis notebook is available on Google Colab with all visualizations and model implementations.

About This Analysis: This project was completed as part of exploring supervised learning applications in financial fraud detection, demonstrating the end-to-end machine learning workflow from data exploration to model evaluation.