

# Explainable AI in Credit Card Fraud: Why We Need AI to Speak Human

**Puneet Sharma**

Senior IT Project Manager

## Abstract

Credit card fraud is an ongoing challenge in the financial industry, with criminals constantly developing sophisticated methods to bypass security systems. As financial institutions embrace artificial intelligence (AI) to combat fraud, the need for explainable AI (XAI) becomes increasingly important. Unlike traditional black-box AI models, explainable AI provides transparency and clarity into decision-making processes, enabling organizations to trust and understand how AI models arrive at their conclusions. This white paper discusses the significance of XAI in credit card fraud detection systems and explores its role in improving trust, accountability, and effectiveness. The paper highlights why traditional AI systems, despite their accuracy, fall short in critical areas such as fairness, interpretability, and accountability, especially in high-stakes areas like credit card fraud. We will explore how XAI can improve not only model performance but also regulatory compliance and customer experience by making AI decisions understandable to non-technical stakeholders. Additionally, the paper will examine the challenges of implementing XAI in fraud detection and offer insights on best practices for developing interpretable, accountable, and fair AI solutions.

**Keywords:** Explainable AI, Credit Card Fraud Detection, Artificial Intelligence, Transparency, Interpretability, Trustworthy AI, Financial Fraud, Fairness, Machine Learning

## Introduction

Credit card fraud has been a persistent and ever-evolving issue in the financial sector. Fraudulent transactions not only cause financial losses but also damage the reputation and trust of financial institutions. Traditional rule-based systems and anomaly detection models, while effective to an extent, often fail to keep pace with the growing sophistication of fraud schemes. As a result, the financial industry has turned to AI and machine learning (ML) to enhance fraud detection systems, relying on these technologies' ability to identify patterns and anomalies that would otherwise go undetected.

However, the use of AI in credit card fraud detection presents a significant challenge: the "black-box" nature of many AI models. These models, such as deep learning networks and ensemble methods, are capable of making highly accurate predictions, but they do so in ways that are not easily understandable to humans. This lack of transparency can lead to issues of trust, accountability, and fairness, particularly in regulated industries such as finance.

Explainable AI (XAI) seeks to address these challenges by making AI models' decision-making processes interpretable and transparent to human users. XAI helps financial institutions and their stakeholders understand not just what decisions an AI system makes, but why it makes them. This is especially important in the context of credit card fraud detection, where decisions can have significant consequences for both the institution and the individual customer.

This white paper delves into the importance of XAI in credit card fraud detection, explores its key benefits, challenges, and provides recommendations for implementing explainable AI in financial systems.

**Figure 1: Advantages of AI in Credit Card Fraud detection**



## The Need for Explainable AI in Credit Card Fraud Detection

### 1. Trust and Transparency

AI-driven fraud detection systems typically rely on complex algorithms that analyze vast amounts of data, such as transaction history, geographical location, device fingerprinting, and customer behavior. While these models can identify fraud with high accuracy, they often fail to provide insights into why a particular transaction was flagged as fraudulent or legitimate. This lack of transparency can lead to a lack of trust from both customers and regulatory bodies.

For example, when a customer's credit card is wrongly flagged as involved in fraud, they may not understand why the system arrived at that decision. Providing explanations for why a transaction was flagged can not only help the customer understand the decision but also improve customer satisfaction and confidence in the system.

### 2. Regulatory Compliance

Regulation in the financial industry is becoming more stringent, with growing demands for accountability and transparency in AI decision-making. In the European Union, the General Data Protection Regulation (GDPR) includes provisions that require "explanation of decisions" made by automated systems, which applies to areas like credit card fraud detection. Financial institutions must ensure that their AI models can be audited and that the reasoning behind automated decisions can be easily explained to regulators.

XAI helps in this regard by ensuring that institutions can demonstrate the validity and fairness of their fraud detection systems, reducing the risk of compliance violations and regulatory scrutiny. It provides a mechanism for organizations to show that their AI models are not acting in discriminatory or biased ways.

### 3. Fairness and Bias Mitigation

AI models, including those used for credit card fraud detection, are often criticized for perpetuating bias. These models learn patterns from historical data, which may inadvertently include biases against certain demographic groups. For example, an AI model may flag transactions from certain regions, age groups, or genders more frequently, even when these groups are not inherently more likely to commit fraud.

XAI can help uncover and address these biases by making it possible to trace the decision-making process of AI models. By understanding why a model flags certain transactions, stakeholders can identify potential bias in the data or the model itself and take corrective action. This helps ensure that the fraud detection system operates fairly and that no group is unjustly targeted or excluded.

#### 4. Improved Model Performance and Accountability

Explainable AI allows for continuous improvement and fine-tuning of fraud detection models. By providing transparency into how and why models make certain predictions, XAI enables data scientists to identify potential weaknesses or areas for improvement in the models. If a model consistently flags certain types of transactions incorrectly, an explanation can reveal why this is happening, whether it's due to incorrect data, an insufficient feature set, or some other issue.

Moreover, XAI promotes accountability by making it clear who is responsible for decisions made by AI models. This is especially critical in the case of fraud detection, where false positives or negatives can have serious consequences for both financial institutions and consumers. When mistakes are made, XAI helps stakeholders identify the root causes and take steps to rectify them.

### Key Techniques in Explainable AI for Credit Card Fraud

#### 1. Model-Agnostic Approaches

Model-agnostic methods are explainable AI techniques that can be applied to any machine learning model, regardless of its architecture. These methods provide explanations of how the model makes predictions based on the features it uses. Some common model-agnostic techniques include:

- **LIME (Local Interpretable Model-agnostic Explanations):** LIME approximates complex models with simpler, interpretable models for individual predictions. For example, it could explain why a particular transaction was flagged as fraudulent by approximating the decision boundary of the underlying model.
- **SHAP (SHapley Additive exPlanations):** SHAP values explain the contribution of each feature to the model's prediction. This technique assigns a numerical value to each feature, helping stakeholders understand which features (e.g., transaction amount, time, location) influenced a fraud detection decision the most.

#### 2. Interpretable Models

In contrast to black-box models like deep neural networks, interpretable models are designed to be transparent from the start. These models inherently provide explanations for their decisions without the need for post-hoc explanation methods. Some examples of interpretable models in credit card fraud detection include:

- **Decision Trees:** Decision trees provide a clear, hierarchical explanation of how a prediction is made. Each node in the tree represents a decision based on a feature, and the path from root to leaf shows the series of decisions leading to a prediction.
- **Logistic Regression:** Logistic regression models are inherently interpretable because they provide coefficients for each feature, indicating the weight of each feature's contribution to the prediction.

#### 3. Rule-Based Systems

Another approach to making AI models explainable is through rule-based systems, where human-understandable rules are used to define decision-making processes. These rules can be derived from data or expert knowledge and can be combined with machine learning models to improve both interpretability and performance.

- **Expert-Driven Rules:** Financial institutions can use domain knowledge to create rules that guide fraud detection decisions, such as flagging transactions that exceed a certain threshold or come from an unfamiliar geographical location.

## Challenges in Implementing Explainable AI in Credit Card Fraud

### 1. Balancing Explainability and Accuracy

One of the primary challenges in implementing XAI is striking the right balance between explainability and model performance. Highly accurate models, such as deep neural networks, are often difficult to explain, while simpler models may not offer the same level of precision. Finding the right balance is crucial to ensure that the AI system can both detect fraud effectively and provide understandable explanations for its decisions.

### 2. Data Privacy and Sensitivity

Explaining AI decisions in credit card fraud detection requires access to sensitive customer data, which raises privacy concerns. Financial institutions must ensure that their explanations do not inadvertently expose personal information or violate data privacy regulations. Techniques like federated learning, which allows AI models to be trained without sharing sensitive data, can help mitigate privacy risks.

### 3. Cost and Complexity

Developing explainable AI models can be more resource-intensive and complex than deploying traditional black-box models. Institutions need to invest in tools, techniques, and expertise to implement and maintain explainable AI systems, which may involve additional costs. However, these investments can pay off in the long run by improving transparency, regulatory compliance, and customer trust.

## Conclusion

As credit card fraud detection becomes increasingly reliant on AI, the need for explainable AI is more important than ever. XAI provides transparency, fairness, and accountability in AI-driven systems, helping financial institutions build trust with their customers and regulators. By explaining how AI models make decisions, XAI not only enhances model performance but also addresses key concerns related to fairness, bias, and regulatory compliance.

While there are challenges in implementing explainable AI, the benefits it offers—especially in high-stakes environments like fraud detection—far outweigh the obstacles. By adopting explainable AI, financial institutions can create systems that are not only more accurate but also more transparent and trustworthy, paving the way for the responsible use of AI in combating credit card fraud.

## References

1. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135-1144.

2. Ribeiro, M. T., et al. (2018). "Anchors: High-Precision Model-Agnostic Explanations." *AAAI Conference on Artificial Intelligence*.
3. Lipton, Z. C. (2016). The Mythos of Model Interpretability. *Communications of the ACM*, 59(4), 36-43.
4. Caruana, R., et al. (2004). "Case-Based Explanation for Decision Trees: A Multidisciplinary Approach." *IEEE Transactions on Knowledge and Data Engineering*, 16(2), 178-186.
5. Lakkaraju, H., et al. (2016). "Interpretable and Explorable Approximations of Black-box Models." *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.