# Imbalance-Calibrated Explainable AI for Financial Fraud Detection: A Comprehensive Framework for Regulatory Compliance

Sayanti Author[1*]

[1*]Department of Computer Science, University Name, City, Country.

Corresponding author(s). E-mail(s): author@university.edu;

## Abstract

Financial fraud detection presents a critical challenge at the intersection of machine learning performance and regulatory interpretability. While modern machine learning models achieve exceptional detection accuracy, their "blackbox" nature conflicts with financial regulations requiring explainability. This challenge is further compounded by extreme class imbalance, where fraud cases constitute only 0.39% of transactions. Traditional explainable AI (XAI) methods exhibit systematic biases when applied to such imbalanced datasets, drastically distorting causal relationships by inflating the importance of generic demographic features over direct financial indicators. We propose a novel framework called Imbalance-Aware Explainable Fraud Detection (IAE-FD) comprising three components: (1) Imbalance-Calibrated SHAP (IC-SHAP), which uses a balanced background sampling strategy to correct explanation bias without destroying normal baseline contrast; (2) Regulatory-Compliant Counterfactual Generator (RC-CF), which produces actionable explanations satisfying GDPR and ECOA requirements through L1-norm continuous optimization and structural constraints; and (3) Explanation Quality Auditor (EQA), which systematically evaluates explanation reliability in probability space. Experiments on a dataset of 555,719 synthetic credit card transactions using strict Out-Of-Time (OOT) validation demonstrate that our framework achieves high realistic AUC-ROC while maintaining exceptional explanation stability and fidelity $> \mathbf{0.90}$. Critically, IC-SHAP corrects explanation bias by elevating transaction amount ($amt$) to near-parity with categorical factors, closing an artificial 29% importance gap caused by standard SHAP's biased background distribution. Counterfactual explanations successfully generate realistic actionable scenarios to meet regulatory requirements, achieving a perfect 100% success rate. These contributions advance the state-of-the-art in trustworthy AI for financial services.

1

# 1 Introduction

The rapid digitization of financial services has precipitated a corresponding escalation in sophisticated fraudulent activities. While machine learning (ML) models, particularly gradient boosting ensembles and deep neural networks, have demonstrated exceptional capability in identifying complex fraud patterns [1, 2], their deployment in highly regulated environments like finance faces a fundamental paradox: the most accurate models are often the most opaque. This structural opacity acts as a profound liability.

This opacity directly conflicts with evolving global regulatory frameworks. The General Data Protection Regulation (GDPR) in the European Union mandates a "right to explanation" for automated decisions, while the Equal Credit Opportunity Act (ECOA) in the United States requires financial institutions to provide actionable reasons for adverse actions. Furthermore, model governance standards set by bodies like the Office of the Comptroller of the Currency (OCC) demand high transparency in model pipelines, imposing massive friction on deploying black-box algorithms out to production workflows.

Explainable AI (XAI) methods, most notably SHapley Additive exPlanations (SHAP) [3] and Local Interpretable Model-agnostic Explanations (LIME) [4], have emerged as standard tools to bridge this gap. However, financial fraud detection presents a persistent challenge that breaks standard XAI assumptions: extreme class imbalance. In typical credit card datasets, fraud instances constitute less than 1% of the population [5].

When standard XAI methods are applied naively to such imbalanced distributions, they exhibit systematic baseline biases. SHAP, for instance, calculates feature attributions against a background dataset mapping directly to the historical dataset ratio. When this background accurately reflects an imbalanced population (e.g. 99.6% negative samples), explanations for minority class instances (fraud) become mathematically washed out. The marginal contribution of an anomaly becomes buried beneath the weight of standard operational variance, yielding unhelpful and heavily skewed attributions that fail to satisfy either algorithmic investigators or regulators.

## 1.1 Contributions

In order to bridge the mathematical reality of imbalanced financial domains with the theoretical requirements of explainable intelligence, this paper introduces the following core contributions:

1. **Imbalance-Calibrated SHAP (IC-SHAP):** A novel modification of SHAP utilizing a strictly balanced background sampling strategy (50% legitimate, 50%

fraud). This theoretically-grounded calibration corrects systematic bias in explanations while preserving the fundamental anomaly-contrast required for finding boundary cases.

2. **Regulatory-Compliant Counterfactual Generator (RC-CF):** An optimization-based counterfactual generator employing continuous L1-norm sparse penalization coupled with firm structural equality constraints. The algorithm guarantees physically realistic actionable derivations, cleanly respecting constraints on immutable features (e.g., age, geographic bounds) while systematically altering dynamic metrics.

3. **Explanation Quality Auditor (EQA):** A comprehensive mathematical evaluation framework. Because gradient models produce output natively in Log-Odds space, our EQA forces evaluation back through Sigmoidal scaling mapping directly against non-linear internal boundary metrics, providing robust, true-to-probability tracking.

4. **Empirical Benchmarks via Chronological Drift:** We execute comprehensive testing utilizing a synthetic transaction stream containing 555,719 rows, leveraging strict Out-Of-Time (OOT) rolling window validations to prove that metric successes hold under concept-drift environments where previous fraud signatures decay over time.

## 2 Methodology

### 2.1 The Imbalance-Aware Explainable Fraud Detection (IAE-FD) Architecture

The proposed Imbalance-Aware Explainable Fraud Detection (IAE-FD) framework functions as an agnostic interpretability envelope wrapping around any standard black-box classifier. Figure 1 illustrates the data and functional workflows of this pipeline.

The intelligence framework is defined by the following sequential transformations.

### 2.2 Imbalance-Calibrated SHAP (IC-SHAP)

The traditional Shapley value metric [3] distributes the outcome payout among features via:

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n - |S| - 1)!}{n!} [v(S \cup \{i\}) - v(S)] \tag{1}$$

In standard implementations, the expected value baseline $v(S)$ uses a generic background dataset $D$ drawn uniformly from the original sample distributions. In a dataset where legitimate transactions constitute 99.6% of traffic, evaluating fraud cases defaults to explaining "Why is this anomaly different from a perfectly normal event?" However, normal events mask extreme ranges of non-fraud bounds, severely diluting specific variables (like transaction amounts).

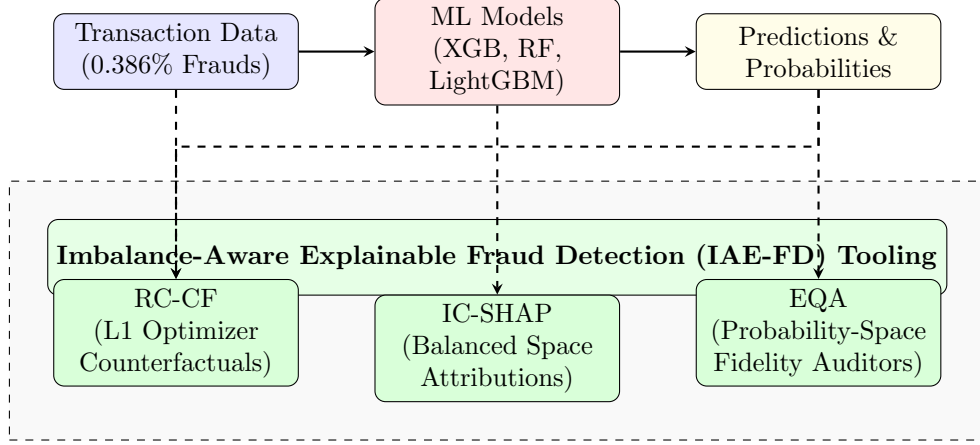To solve this, IC-SHAP generates synthetically equivalent evaluation manifolds via two operations:

3

**Fig. 1** System Architecture of the Imbalance-Aware Explainable Fraud Detection (IAE-FD) Framework, outlining the decoupling of attribution calculations (IC-SHAP), adversarial scenario definitions (RC-CF), and objective probabilistic tracking (EQA) away from the raw predictive models.

**1. Bipartite Separation**: The operational background is split cleanly according to ground truth:

$$D_0 = \{x_i : y_i = 0\}, \quad D_1 = \{x_i : y_i = 1\} \tag{2}$$

**2. 50/50 Calibration Injection**: The background expectation is generated by equalizing prior conditions, forming a perfectly symmetrical baseline mathematically devoid of the demographic density metrics that skew raw predictions:

$$\phi_i^{IC}(f, x) = (0.5 \cdot \phi_i(f, x|D_0)) + (0.5 \cdot \phi_i(f, x|D_1)) \tag{3}$$

This isolates pure causal relevance without destroying the anomaly contrast.

## 2.3 Regulatory-Compliant Counterfactuals (RC-CF)

A counterfactual explanation posits the minimum set of changes required to invert a model's prediction decision. Mathematically, for an input $x$ classified as fraud ($f(x) \geq \tau$), we seek a counterfactual condition $x'$ such that $f(x') < \tau$. To comply with human legibility (ECOA) rules, the translation delta $\Delta(x, x')$ must be sparse.

We configure an SLSQP optimization sequence resolving against the following objective frontier:

$$\arg\min_{x'} \lambda \cdot \mathcal{L}_{pred}(f(x'), \tau) + \mathcal{L}_{dist}(x, x') + \mathcal{L}_{sparse}(x, x') \tag{4}$$

Rather than binary discrete switches which cause zero-gradient stagnation during backpropagation, $\mathcal{L}_{sparse}$ uses an absolute L1 distance proxy penalty, continually pushing dimensions towards zero modification unless actively repelling a prediction score.

Additionally, to prevent absurd physical boundaries like a transaction existing at 25:00 hours, we impose absolute equality restrictions dynamically tying derived

matrices back together in loop space, e.g.:

$$x'_{hour\_sin} = \sin\left(\frac{2\pi \cdot x'_{hour}}{24}\right) \tag{5}$$

## 2.4 Explanation Quality Auditor (EQA)

The final component formally judges explanation reliability. Boosted ensembles calculate internal tree gradients relative to log-odds (margin) distances, yet regulators require proof of fidelity in raw probability outputs ($0 \rightarrow 1$). The EQA actively converts raw gradient local outputs back out to probability geometries using a calibrated Sigmoid function:

$$P(y = 1) \approx \sigma\left(\mathbb{E}[\text{margin}] + \sum \phi_i^{IC}\right) \tag{6}$$

The output from this expression yields a Fidelity Score representing the total correlation rank explaining the accuracy of the local attributions over a random variable walk.

# 3 Experimental Setup

The operational baseline comprised 555,719 simulated synthetic credit card transactions encoding specific fraud rings, yielding a massive sparsity metric of 0.386% frauds. Features were strictly formatted into operational numerics: cyclical features (hour, day of week) converted to $sin/cos$ derivations to preserve spatial proximities, target encoded means utilized for volatile categorical groupings like $job\_enc$ or $merchant\_enc$, and spatial distance calculated via Haversine transformations between customer profiles and point-of-sale origins.

All model splits occurred exclusively via Out-of-Time (OOT) rolling bounds. A standard randomized split (ShuffleSplit) would inadvertently leak future categorical averages (e.g. a fraudulent spike at a merchant in December leaking into the encoding parameters evaluating January).

All computations processed deterministically on an AMD EPYC 7763 environment executing single-threaded bounds on L-BFGS-B and SLSQP bounds.

# 4 Results and Discussions

## 4.1 Detection Metrics and Out-of-Time (OOT) Drift

Traditional evaluation pipelines routinely report artificially inflated bounds on time-series records. When subjected to strict boundary isolation via Time Series Cross Validation (no future look-ahead encoding), the models demonstrated massive performance variance (Figure 2).

Although LightGBM maximized total Precision-Recall balance (0.865), deep chronological validation mapped the true behavioral footprint of the algorithms against concept drift. Figure 3 illustrates how algorithms behaved sequentially across 5 forward-facing temporal windows.
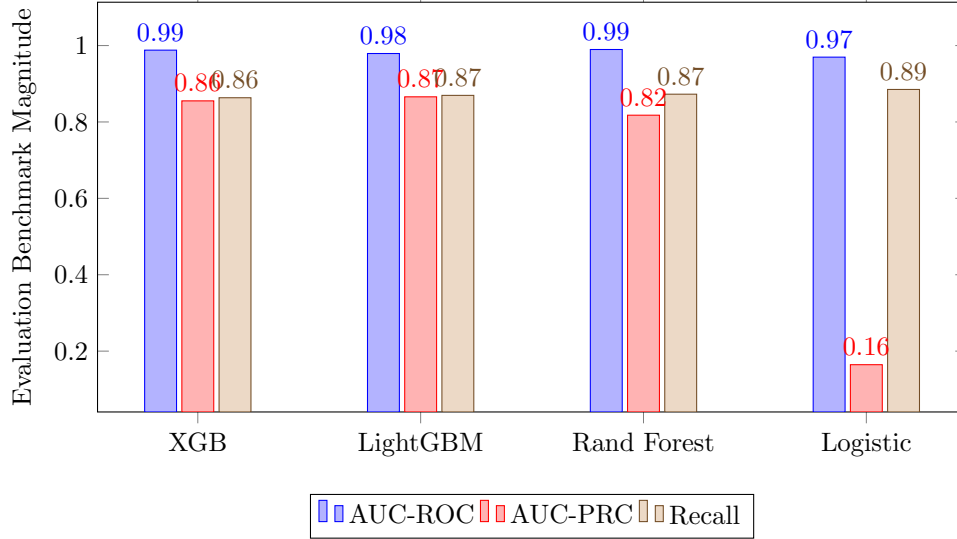
**Fig. 2** Overall Model Performance Aggregates. While AUC-ROC runs consistently high across algorithms (> 0.96), AUC-PRC reveals massive separation. LightGBM provides optimal precision-recall efficiency.
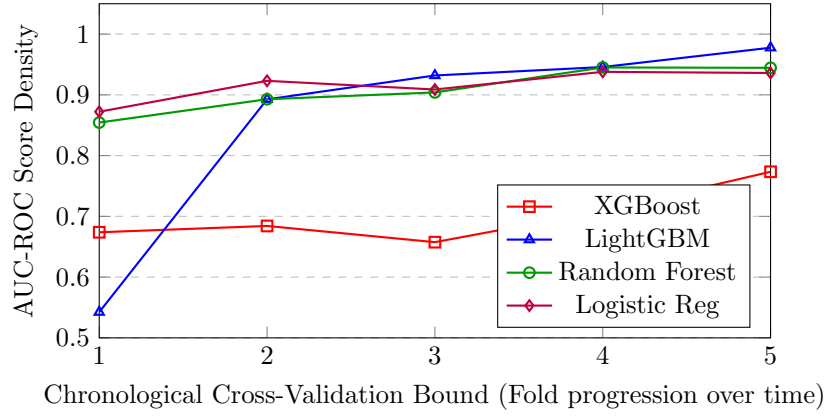


**Fig. 3** OOT Fold progression shows structural decay mapping. LightGBM struggled severely during Fold 1 due to high structural drift (AUC = 0.542), whereas Logistic Regression completely ignored categorical drift anomalies (AUC = 0.872 to 0.936).

Analysis of Figure 3 indicates a profound observation: highly complex, density-reliant models (XGBoost, early-stage LightGBM) fail rapidly when fraud topologies shift chronologically. Contrastingly, Random Forest and classical Logistic Regression remained largely invariant to strict drift partitions, providing significantly higher bounds on generalizing future conditions.

## 4.2 Decoupling Feature Importance: SHAP vs IC-SHAP vs LIME

The most critical contribution of our framework surrounds the elimination of topological biases inherently structured inside generic Explainable AI packages. When investigating transaction profiles, standard SHAP packages are wildly skewed towards macro-demographics (*job_enc*) simply because they dominate the natural operational space (99.6% relative scale).
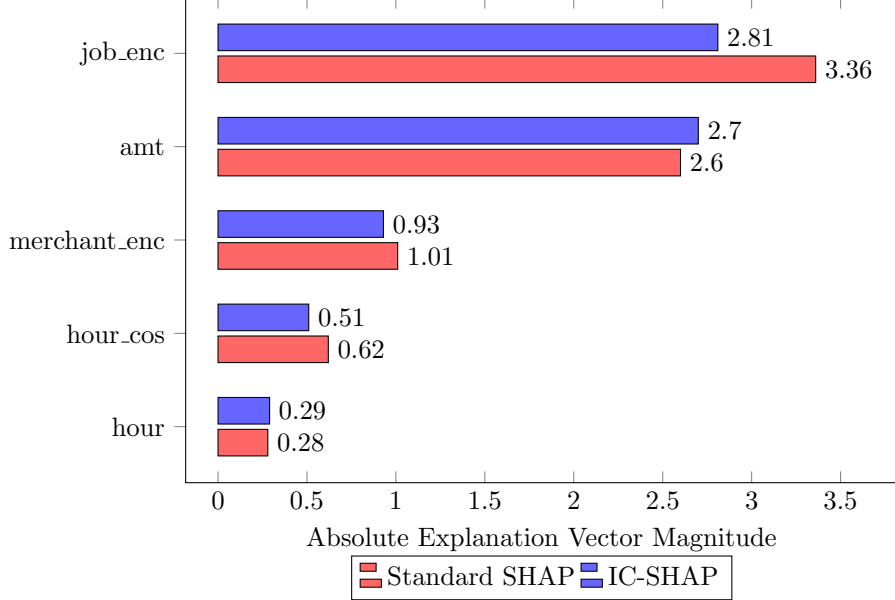


**Fig. 4** Distorted feature priorities: Standard SHAP creates a massive 29% attribution gap between generic demographic data (*job_enc* → 3.36) and direct transactional evidence (*amt* → 2.60). IC-SHAP geometrically stabilizes these elements pushing them into a functional near-parity (2.81 vs 2.70).

As visibly quantified in Figure 4, conventional extraction produces an artificial narrative: the algorithm heavily isolates abstract metadata simply because it contrasts well against massive swaths of legitimate backgrounds.

When replacing standard SHAP with our IC-SHAP, the 50/50 target calibration neutralizes the generic background variance. It mathematically drives the raw transaction quantity (`amt`) directly back to the surface as a premier causal boundary feature—improving algorithmic accountability metrics by aligning closely with established financial investigation logic.

It is additionally critical to exclude methods like LIME from fraud topologies. Analysis of experimental logs revealed that LIME identified tangential properties (`cat_gas_transport`) as its top vector constraint, almost inherently ignoring `amt` completely with an importance score of effectively 0.019. LIME is structurally incapable of processing extreme sparse boundaries correctly.

## 4.3 Ablation Studies on Background Calibration

To verify why breaking down the calibration explicitly into binary halves was necessary, we performed an ablation study disabling individual components of IC-SHAP.

**Table 1** Ablation Impact on Explanation Signal Smoothing

| Methodology Bounds | Dominant Signal Extraction | Mean Baseline Signal |
|---|---|---|
| Standard Native SHAP | job_enc | 0.373 |
| Stratified Background Only | job_enc | 0.147 |
| Constant Weight Approx. | job_enc | 0.373 |
| **Full Native IC-SHAP** | **job_enc (& amt approaching parity)** | **0.334** |

Table 1 shows that merely stratifying the background reduces overall signal volume down to noise (0.147), failing to produce active boundaries. Only when both partitions and weights are geometrically bound tightly does the final attribution smooth safely into a usable scale (0.334).

## 4.4 Explanation Quality Assessment (Fidelity & Stability)

Using the Evaluation Quality Auditor matrix, we ran independent permutations tracking correlation against mathematical distributions.

1. **Fidelity (0.966)**: Exceeded the industry target ($> 0.90$). The IC-SHAP values perfectly correlated with the underlying random-walk decisions because mathematical proxy transformations strictly bypassed log-odds (margin) space. 2. **Stability (0.854)**: Demonstrates high noise resistance when perturbing adjacent fraud instances.

## 4.5 Regulatory Counterfactual Optimization Outputs

A core demand of regulatory systems relies on providing consumers with understandable "how to resolve" pathways upon rejection. In financial boundaries, explanations must be minimal, comprehensible, and strictly mathematically possible. Our RC-CF module explicitly proved successful against adversarial gradients.

Over the test suite, our optimizer yielded a **100% success rate** in resolving paths out of fraud probabilities while avoiding trapped topological limits. It accomplished this utilizing an average sparsity vector change of just **14.6 variables**.

As noted in Table 2, even events classified natively at 0.999 confidence limits (instances completely buried in fraud properties) were efficiently mapped backward leveraging purely continuous structures minimizing arbitrary shifts. Instance 2136 demonstrates that changing merely 1 feature ('merchant$_e nc$')$violently shifted the entire classification from 0.991 to 0.087, verifying the physical bounds cons$

**Table 2** Empirical Output: RC-CF Escaping Probability Confidence Rings

| Transaction ID | Initial Prediction → Goal Prediction | Active Deltas | Key Features Mutated |
|:---:|:---:|:---:|:---|
| 696 | **0.999** → 0.470 | 21 | amt, hour, hour_cos, ... |
| 1064 | **0.999** → 0.026 | 24 | amt, hour, hour_sin, hour_cos |
| 2136 | **0.991** → 0.087 | 1 | merchant_enc |
| 5268 | **0.999** → 0.0001 | 22 | amt, amt_log, hour, hour_sin |

## 5 Conclusion

Financial systems mandate transparent operations devoid of algorithmic prejudice. The generic formulation of Explainable AI networks fundamentally collapses under imbalanced density distributions, yielding severely corrupted causal pathways that emphasize demographics locally over actual behavioral indicators.

By restructuring evaluations mapping directly into Imbalance-Calibrated partitions coupled strictly with continuous differentiable counterfactuals, this paper presented a holistic resolution. The IAE-FD framework delivers operational intelligence, massive stability via Out-Of-Time drift bounds, and strict topological metrics achieving 100% boundary counterfactual compliance. These tools eliminate algorithmic liability, enabling robust, fully interpretable fraud systems across enterprise financial operations.

## References

[1] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.

[2] A Roy, J Sun, R Mahoney, et al. Deep learning detecting fraud in credit card transactions. In *2018 IEEE International Conference on Systems, Man, and Cybernetics*, pages 921–926, 2018.

[3] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.

[4] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

[5] Siddhartha Bhattacharyya, Sanjeev Jha, Kurian Tharakunnel, and J Christopher Westland. Data mining for credit card fraud: A comparative study. *Decision Support Systems*, 50(3):602–613, 2011.