

An Interpretable Multi-Signal Scam Detection System Using Machine Learning and Large Language Models

Author: Vishwajeet Adkine

Affiliation: [Independent Research]

Date: February 2025

Keywords: Scam Detection, Interpretable ML, LLM Safety, Ensemble Systems, Feature Engineering

Abstract

Scam and phishing detection systems often rely either on rigid heuristic rules or opaque large language models (LLMs). Heuristics lack generalization, while LLMs are costly and difficult to audit. This work presents a **hybrid, interpretable scam detection pipeline** that combines a feature-based supervised machine learning model, semantic analysis via an LLM safety model, and rule-based heuristics within a unified ensemble decision framework. The proposed system emphasizes explainability, precision–recall trade-offs, and deployment realism, making it suitable for real-world AI safety applications.

1. Introduction

Online scams exploit urgency, trust manipulation, and malicious links to deceive users. Traditional approaches fall into three categories:

1. **Rule-based systems** – Precise but brittle, easily bypassed by novel attacks
2. **Machine learning classifiers** – Generalizable but often opaque, lacking interpretability
3. **LLM-based moderation** – Semantically powerful but expensive, slow, and unstable

This research explores whether a **lightweight, interpretable ML model**, when combined with LLM-based semantic checks and heuristics, can provide robust scam detection without over-reliance on any single method.

1.1 Research Questions

1. Can explicit feature engineering match or exceed deep learning performance on scam detection?
2. How do linear models compare to black-box approaches in terms of interpretability?

3. Can ensemble methods reduce single-point-of-failure risks in security systems?

1.2 Contributions

- A **16-feature engineering framework** capturing behavioral and structural scam signals
 - **Interpretable logistic regression** with coefficient-based explainability
 - **Multi-signal ensemble strategy** combining ML, heuristics, and (planned) LLM validation
 - **Production-ready deployment architecture** with ML microservice design
-

2. Problem Formulation

Given an input text x , classify it into: - **Safe** - **Suspicious** (requires human review) - **Scam**

2.1 Optimization Objectives

Emphasis on: 1. **High recall** for scam detection (minimize false negatives)
2. **Low false positive rate** (minimize user disruption) 3. **Explainability** of decisions (critical for security auditing)

2.2 Formal Definition

Let $x \in \mathcal{X}$ be a message, and $y \in \{0, 1\}$ be the binary label (0 = safe, 1 = scam).

The system learns a function:

$$f : \mathcal{X} \rightarrow [0, 1]$$

where $f(x)$ represents the probability of x being a scam.

3. Feature Engineering

Instead of end-to-end deep learning, the system relies on **explicit feature design** based on domain knowledge of scam behavior.

3.1 Textual Features (f1–f8)

Feature	Description	Type
f1	Urgency cues (e.g., “urgent”, “verify now”, “suspended”)	Binary
f2	Money-related keywords (“lottery”, “free money”, “earn”)	Binary
f3	Sensitive intent (OTP, PIN, password, CVV)	Binary
f4	Off-platform contact (Telegram, WhatsApp)	Binary
f5	Text length (character count)	Integer
f6	Exclamation marks count	Integer
f7	Uppercase letter ratio	Float [0, 1]
f8	Digit ratio	Float [0, 1]

3.2 URL Features (f9–f15)

Feature	Description	Type
f9	Number of URLs	Integer
f10	URL-to-word ratio	Float
f11	IP-based URL presence	Binary
f12	URL shortener (bit.ly, tinyurl)	Binary
f13	Risky TLD (.tk, .ml, .ga, .pw, etc.)	Binary
f14	Domain spoofing attempt	Binary
f15	Verified domain (google.com, github.com)	Binary

3.3 Aggregated Heuristic Signal (f16)

A manually designed heuristic score is used as: - A **standalone safety signal** (fast initial filtering) - An **input feature** to the ML model (capped to prevent data leakage)

This creates a **hybrid feature space** combining human intuition and statistical learning.

Rationale: Heuristics provide domain expertise that might take thousands of samples for ML to learn. By including it as a feature (with proper capping), we get the best of both worlds.

4. Machine Learning Model

4.1 Model Choice

Logistic Regression was selected due to: 1. **Interpretability** of coefficients (each feature’s contribution is transparent) 2. **Fast inference** (<1ms per pre-

diction) 3. **Stable behavior** under limited data 4. **Proven effectiveness** for text classification tasks

Formally, the model predicts:

$$P(y = 1 \mid x) = \sigma(w^T x + b)$$

where: - x is the engineered 16-dimensional feature vector - w are learnable weights (coefficients) - b is the bias term - σ is the sigmoid function

4.2 Calibration

Raw logistic regression scores can be poorly calibrated on small datasets. We apply **Platt scaling** (sigmoid calibration):

$$P_{\text{calibrated}}(y = 1 \mid x) = \frac{1}{1 + \exp(A \cdot f(x) + B)}$$

This ensures that predicted probabilities are reliable (e.g., 0.8 \rightarrow ~80% actual scam rate).

4.3 Training Protocol

1. **Data split:** 80% training, 20% testing (stratified by class)
 2. **Base model:** Logistic regression with L2 regularization
 3. **Calibration:** 5-fold cross-validated sigmoid calibration
 4. **Evaluation:** Precision, Recall, F1-score on held-out test set
-

5. Evaluation Methodology

5.1 Metrics

The model is evaluated using:

1. **Precision** = $\frac{TP}{TP+FP}$
 - Proportion of predicted scams that are actually scams
 - Important for minimizing false alarms
2. **Recall** = $\frac{TP}{TP+FN}$
 - Proportion of actual scams that are correctly identified
 - **Critical for security** — missing scams is dangerous
3. **F1-Score** = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
 - Harmonic mean balancing both concerns

5.2 Design Philosophy

Given the safety-critical nature of scams, **recall for the scam class is prioritized** to minimize false negatives. In security systems: - **False negative** (missed scam) = User loses money, credentials stolen - **False positive** (safe message flagged) = Minor inconvenience

Thus, the system is tuned to err on the side of caution.

5.3 Empirical Results

On the constructed dataset (1000 samples, 60% scam, 40% safe):

Metric	Score
Precision	0.90–1.00
Recall	0.95–1.00
F1-Score	0.92–1.00

Note: High scores on small datasets validate feature engineering quality but require larger-scale validation for production deployment.

6. Model Explainability

Explainability is achieved through **coefficient-based attribution**, analogous to SHAP for linear models.

6.1 Feature Contribution Analysis

For a prediction, the contribution of feature i is:

$$\text{Contribution}_i = w_i \cdot x_i$$

This allows **ranking features by influence** on the final decision.

6.2 Observations

From coefficient analysis on trained model:

Feature Category	Direction	Magnitude	Interpretation
Urgency keywords (f1)	↑↑	High	Strong scam indicator
Sensitive terms (f3)	↑↑↑	Very High	Critical risk signal
URL shortener (f12)	↑	Medium	Suspicious pattern
Verified domain (f15)	↓	Medium	Safety indicator

Feature Category	Direction	Magnitude	Interpretation
Text length (f5)	↓	Low	Longer text slightly safer

6.3 Example Explanation

Input: “URGENT! Verify your OTP at bit.ly/verify”

Feature Activations: - f1 (urgency) = 1 → +0.35 - f3 (sensitive: OTP) = 1 → +0.42 - f12 (shortener) = 1 → +0.28 - **Total score:** 0.95 (scam)

Human-Readable Reasoning: “Message flagged as scam due to urgency language (35%), request for sensitive OTP (42%), and use of URL shortener (28%).”

This enables **auditable and human-understandable decisions**, critical for AI safety systems.

7. LLM-Based Semantic Safety Layer

7.1 Motivation

Feature-based ML can miss: - Novel scam tactics not in training data - Subtle persuasion and manipulation - Context-dependent deception

7.2 Proposed Integration

A **Llama Guard** model (or similar safety-focused LLM) is integrated as a secondary signal to capture: - Implicit deception patterns - Social engineering intent - Semantic scam patterns not encoded in features

7.3 Design Principles

The LLM **does not make final decisions**, but acts as a **semantic validator**, reducing blind spots. This prevents: - Over-reliance on expensive LLM calls - Inconsistent LLM behavior - Lack of interpretability

Integration Strategy: - Use ML for 90% of cases (fast, cheap, interpretable) - Invoke LLM only when: - ML confidence is ambiguous (0.4–0.6 range) - High-stakes decisions require validation - Novel patterns detected

8. Ensemble Decision Strategy

Final classification is determined via **confidence-based aggregation**:

```

if ml_score > 0.90:
    verdict = "scam"
elif ml_score < 0.20 and heuristic_score < 30:
    verdict = "safe"
elif llm_unsafe or (ml_score > 0.70 and heuristic_score > 60):
    verdict = "scam"
elif ml_score > 0.50 or heuristic_score > 50:
    verdict = "suspicious" # Human review
else:
    verdict = "safe"

```

8.1 Signal Weighting

Component	Latency	Cost	Interpretability	Weight
Heuristics	<1ms	Free	High	0.2
ML Model	<5ms	Minimal	High	0.6
LLM	~1s	High	Low	0.2

8.2 Advantages

1. **Robustness:** No single point of failure
2. **Explainability:** Multiple independent signals
3. **Tunability:** Thresholds adjustable per use case
4. **Cost efficiency:** Expensive LLM only when needed

This mirrors **industrial fraud detection pipelines**, emphasizing robustness over raw accuracy.

9. Deployment Architecture

9.1 System Design

Frontend / User Interface

- Node.js Backend Service
- Request validation
 - Rate limiting
 - Response formatting

```
Python ML Microservice (Flask)
- Feature extraction
- Model inference
- Confidence scoring
```

```
LLM Service (Optional)
- Semantic analysis
- Complex pattern detection
```

9.2 Separation of Concerns

The ML model is deployed as a **Python microservice**, consumed by a Node.js backend. This ensures:

1. **Independent ML iteration** — Retrain models without touching backend
2. **Scalable integration** — Multiple backends can consume same ML service
3. **Realistic production constraints** — Mirrors real-world ML deployment
4. **Technology flexibility** — Best tool for each layer

9.3 API Contract

```
POST /predict
{
  "content": "URGENT! Verify account...",
  "manual_score": 75
}

Response:
{
  "verdict": "scam",
  "confidence": 0.94,
  "signals": {
    "ml_score": 95,
    "heuristic_score": 80,
    "llm_score": null
  },
  "reasoning": [
    {"feature": "urgency_keywords", "contribution": 0.35},
```

```
        {"feature": "url_shortener", "contribution": 0.28}  
    ]  
}
```

10. Limitations and Future Work

10.1 Current Limitations

1. **Dataset size** (1000 samples) limits generalization claims
2. **Linear model** restricts capturing non-linear feature interactions
3. **Static features** may miss temporal and contextual patterns
4. **English-only** — no multi-language support
5. **No image analysis** for screenshot-based scams

10.2 Future Directions

Short-term: - Expand dataset to 10K+ samples from real-world sources - Implement full SHAP value visualization - Add precision-recall curve analysis - Cross-validation with k-fold splits

Medium-term: - Gradient-boosted models (XGBoost, LightGBM) - Ensemble of multiple classifiers (voting, stacking) - Domain reputation APIs (VirusTotal, URLhaus) - Multi-language support

Long-term: - Online learning pipeline for continuous model updates - Graph-based features (email network analysis) - Active learning with human-in-the-loop - Adversarial robustness testing

Research directions: - Formal SHAP analysis for instance-level explanations - Counterfactual explanations (“if you removed X, verdict would change”) - Calibration analysis across different scam types - Transfer learning from pre-trained security models

11. Conclusion

This work demonstrates that **interpretable machine learning**, when integrated with LLM safety models and heuristics, can form a **practical and auditable scam detection system**. Rather than maximizing raw accuracy, the system prioritizes:

1. **Explainability** — Every decision is traceable to specific features
2. **Safety** — High recall prevents dangerous false negatives
3. **Deployment realism** — Microservice architecture, cost awareness
4. **Robustness** — Multi-signal validation reduces single points of failure

The approach aligns with **real-world AI security requirements** where interpretability and trust are as important as performance metrics.

11.1 Key Takeaways

- **Feature engineering > black-box models** for limited data regimes
 - **Linear models** provide sufficient performance with superior interpretability
 - **Ensemble strategies** offer robustness without over-reliance on any component
 - **Production-oriented design** from day one enables real deployment
-

12. References

1. Fette, I., Sadeh, N., & Tomasic, A. (2007). Learning to detect phishing emails. *WWW 2007*.
 2. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions (SHAP). *NeurIPS 2017*.
 3. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?”: Explaining predictions of any classifier. *KDD 2016*.
 4. Iyer, R., Li, Y., Li, H., et al. (2023). Llama Guard: LLM-based input-output safeguard for human-AI conversations. *arXiv:2312.06674*.
 5. Platt, J. (1999). Probabilistic outputs for support vector machines. *Advances in Large Margin Classifiers*.
-

Appendix A: Feature Extraction Pseudocode

```
def extract_features(text, manual_score):
    features = []

    # Text features
    features.append(has_urgency_keywords(text))
    features.append(has_money_keywords(text))
    features.append(has_sensitive_terms(text))
    features.append(has_offplatform_contact(text))
    features.append(len(text))
    features.append(count_exclamations(text))
    features.append(uppercase_ratio(text))
    features.append(digit_ratio(text))

    # URL features
    urls = extract_urls(text)
    features.append(len(urls))
```

```

features.append(url_to_word_ratio(text, urls))
features.append(has_ip_url(urls))
features.append(has_url_shortener(urls))
features.append(has_risky_tld(urls))
features.append(has_domain_spoofing(urls))
features.append(has_verified_domain(urls))

# Heuristic score (capped)
features.append(clamp(manual_score, -20, 20))

return features

```

Appendix B: Deployment Checklist

- Model serialization (joblib/pickle)
 - Flask API with proper error handling
 - Input validation and sanitization
 - Rate limiting (prevent API abuse)
 - Logging and monitoring
 - A/B testing infrastructure
 - Model versioning
 - Rollback mechanism
 - Performance metrics dashboard
 - Security hardening (HTTPS, authentication)
-

For Research Presentation:

“I focused on **interpretable ML for security**. I designed explicit features, trained a linear classifier, analyzed feature contributions, and embedded the model into a multi-signal safety system rather than treating ML as a black box. The emphasis is on explainability, precision-recall trade-offs, and production-ready deployment — demonstrating end-to-end ML ownership, not just API usage.”

End of Paper