

Cryptocurrency AML Detection Using the Elliptic Bitcoin Dataset

Mario Paulin

Project Overview

This project focuses on detecting illicit Bitcoin transactions using the Elliptic Bitcoin Dataset. The dataset contains transaction-level data, including features, class labels (legitimate, illicit, or unknown), and an edgelist representing the transaction graph. The primary goal is to develop a random forest model that predicts whether a transaction is illicit, leveraging graph-based features.

Key Objectives

- **Graph-Based Feature Engineering:**
 - Extract meaningful features such as `betweenness centrality`, `PageRank`, and `clustering coefficient`.
 - Explore ego network features to capture local graph structures.
 - Analyze the evolution of illicit transactions and their graph structure over time.
- **Model Development:**
 - Train and hyper-parameter tune a Random Forest classifier to label transactions as legitimate or illicit.
 - Evaluate performance using precision, recall, and F1-score, focusing on the illicit class.
- **Graph Visualization:**
 - Visualize the transaction graph for specific timesteps to compare the graph structure of predictions against ground-truth labels.

Key Findings

1. Temporal and Graph-Based Insights

Early Timesteps (Before Timestep 43):

- Illicit transactions tend to appear at the periphery of the network.
- These transactions have lower centrality and connectivity.
- Features such as `betweenness centrality`, `PageRank`, and `average degree` are effective in this phase.

Later Timesteps (After Timestep 43):

- Illicit transactions form deeper chains and move toward the network's core.
- This makes them harder to detect unless the model is trained on these structural changes.

2. Model Performance

- **Precision (Class 1 - Illicit):** 0.91
91% of predicted illicit transactions were correctly identified.
- **Recall (Class 1 - Illicit):** 0.71
71% of actual illicit transactions were correctly identified.
- **F1-Score (Class 1 - Illicit):** 0.80
Provides a balanced measure of precision and recall.
- **Temporal Generalization:**
Performance degrades after timestep 43 due to structural drift.

3. Graph Visualization

- **Predicted Graph:** Highlights predicted illicit nodes, revealing suspicious clusters.
- **Actual Graph:** Visualizes ground truth illicit transactions to assess prediction gaps.
- **Insight:** Illicit nodes start on the periphery and move toward the core in later timesteps.

Challenges and Limitations

- **Structural Changes:** The model struggles to adapt to evolving graph structures of illicit behavior.
- **Graph Modeling:** Use larger ego network radii to capture structural drift.

Applications in AML Monitoring Systems

- **Real-Time Risk Scoring:** Assign a risk score based on the predicted probability of being illicit.
- **Graph-Based Insights:** Help compliance teams identify clusters and understand fund flows.
- **Explainability:** Use feature importances and partial dependence plots for compliance and audit purposes.