

Graph Mining Course Project Progress Report

Submission Date: [2026-01-07]
Course: Graph Mining [4041]
Instructor: Dr. Zeinab Maleki
Project Title: Graph-Based Illicit Transaction Detection Using Graph Neural Networks

Student Information

- **Student Name(s):** Rasoul Salehi
 - **Student ID(s):** 40125933
 - **Email(s):** [8raha6094@gmail.com]
-
- **Student Name(s):** Sayed Ali Amiri
 - **Student ID(s):** 40118393
 - **Email(s):** [amirisayedali@gmail.com]
-
- **Student Name(s):** Mojtaba Mollaei
 - **Student ID(s):** 40131383
 - **Email(s):** [m.mollaei@ec.iut.ac.ir]

Executive Summary

This report presents the current progress of the graph mining course project focused on detecting illicit transactions in the Bitcoin network using graph neural networks and ensemble learning strategies. Building upon the initial proposal, the Elliptic transaction dataset was fully preprocessed and represented as a large-scale directed graph. Baseline graph learning pipelines were successfully implemented and validated.

Beyond standard GNN models, the project has progressed toward advanced architectures and ensemble techniques. A Subgraph-based Graph Attention Network (SGAT) framework was designed to improve scalability and robustness by operating on sampled subgraphs rather than the full transaction graph. Additionally, bagging strategies were introduced by training multiple GNN models on different subgraph samples, enabling variance reduction and improved generalization. A stacking mechanism was further incorporated to aggregate predictions from multiple base learners using a meta-classifier.

Preliminary experiments confirm that ensemble-based graph learning improves classification stability compared to single-model baselines. The project remains aligned with its original objectives and is progressing toward a complete evaluation of SGAT-based ensemble methods for illicit transaction detection.

Progress on Objectives

Objective 1: Dataset Acquisition and Preprocessing

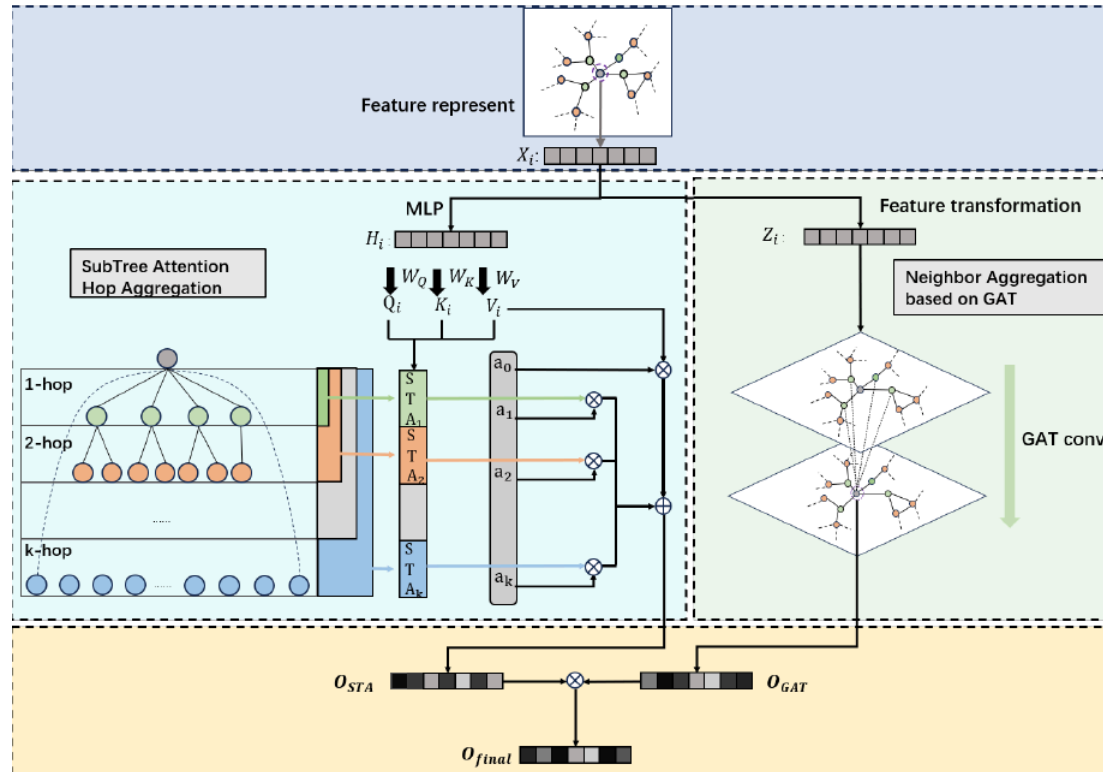
This objective has been fully completed. The Elliptic Bitcoin transaction dataset was successfully downloaded, explored, and transformed into graph-compatible structures. Node features, transaction edges, and class labels were aligned across data sources. Class imbalance and temporal segmentation were identified as key dataset characteristics. Masks for training, validation, and testing were generated to ensure reproducible and fair evaluation.

Objective 2: Graph Construction and Baseline GNN Modeling

This objective is largely completed. The transaction network was modeled as a directed graph where nodes represent transactions and edges represent fund transfers. Node features were normalized and adjacency structures were built for graph-based learning. Baseline GNN models were trained to validate the correctness of the graph pipeline and establish reference performance levels.

Objective 3: SGAT Architecture Design and Subgraph Sampling

This objective is partially completed. A Subgraph-based Graph Attention Network (SGAT) framework was introduced to address scalability and over-smoothing issues present in large graphs. Instead of operating on the full transaction graph, multiple subgraphs were generated using stochastic node and edge sampling strategies. Each subgraph preserves local neighborhood structure while significantly reducing computational cost. Attention mechanisms were applied within each subgraph to focus on influential neighboring transactions. Initial experiments show stable training behavior and improved robustness to noise.

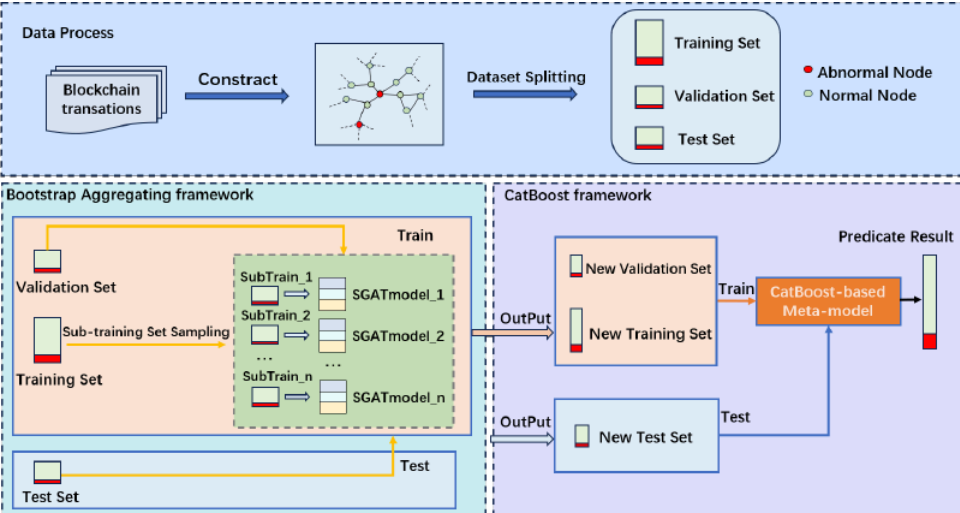


Objective 4: Bagging with Multiple GNN Models

This objective is partially completed. A bagging strategy was implemented by training multiple SGAT/GNN models on different sampled subgraphs. Each model learns complementary structural patterns from the transaction network. The diversity introduced by subgraph sampling helps reduce model variance and sensitivity to noisy or sparse regions of the graph. Predictions from individual models were aggregated using soft-voting and probability averaging. Early results indicate improved prediction stability compared to single-model approaches.

Objective 5: Stacking and Meta-Learning

This objective is in progress. A stacking framework was designed to combine outputs from multiple bagged SGAT models. Predictions from base learners were collected and used as input features for a higher-level meta-classifier. The meta-model learns how to weight different base models depending on their confidence and performance characteristics. This approach aims to exploit complementary strengths of individual learners and further improve detection accuracy on illicit transactions.



Work Accomplished

Dataset Preparation and Analysis

The Elliptic dataset was processed to ensure compatibility with graph neural networks. Feature matrices were constructed, edges were validated, and transaction IDs were consistently mapped. Exploratory analysis revealed a highly imbalanced class distribution and sparse graph connectivity. Temporal properties of the data were preserved to avoid information leakage across time steps.

Implementation Details

The project was implemented using Python with a modular pipeline that separates data processing, graph construction, and modeling. Key tools and techniques include:

- **Pandas / NumPy** for preprocessing and feature handling
- **Graph construction utilities** for adjacency and edge indexing
- **Graph Attention Networks (GAT/SGAT)** for node-level classification
- **Subgraph sampling** for scalability
- **Bagging and stacking** for ensemble learning

The notebook structure allows reproducibility and incremental experimentation.

Preliminary Results

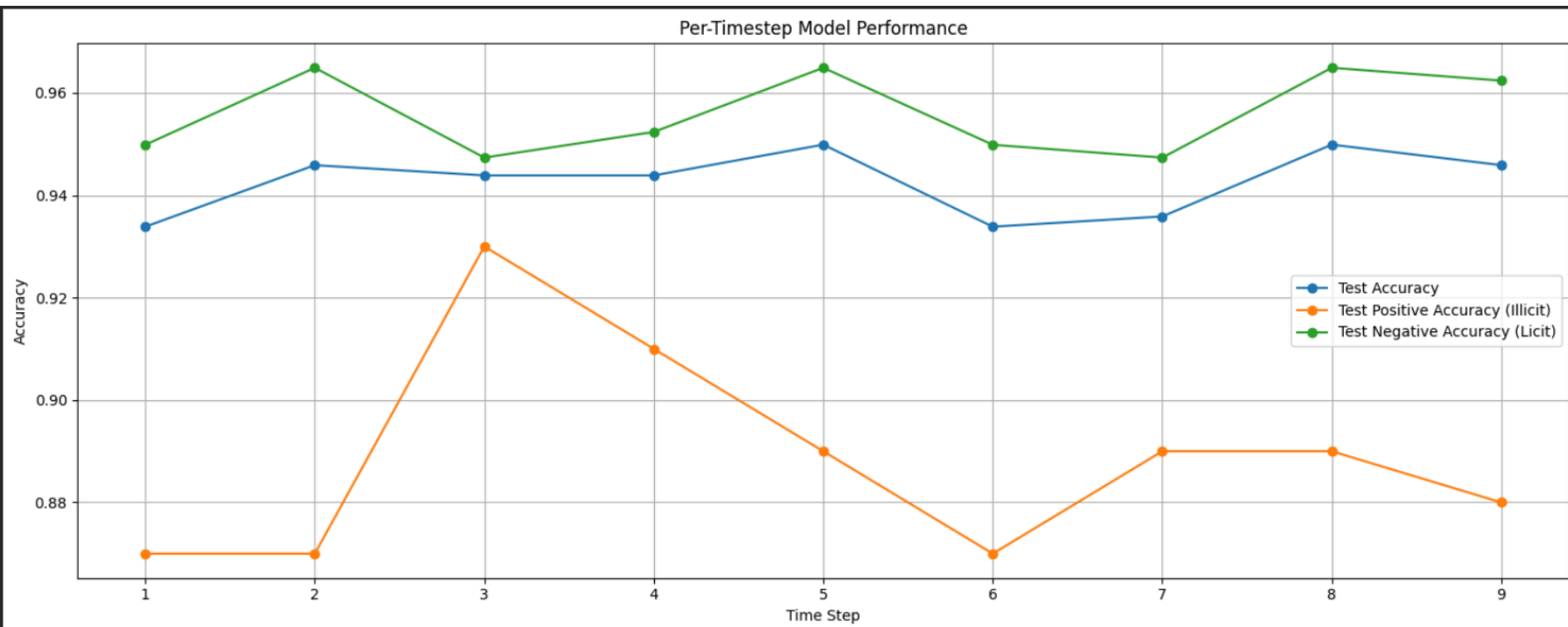
Preliminary experiments suggest that SGAT-based ensemble models outperform single-graph baselines in terms of stability and robustness.

Metric	Value (Preliminary)	Notes
Graph Nodes	~200K	Bitcoin transactions
Graph Edges	~230K	Directed transaction flows
Number of Subgraphs	8	Used for bagging
Base Models	SGAT / GNN	Attention-based
Ensemble Stability	Improved	Compared to single-model baseline

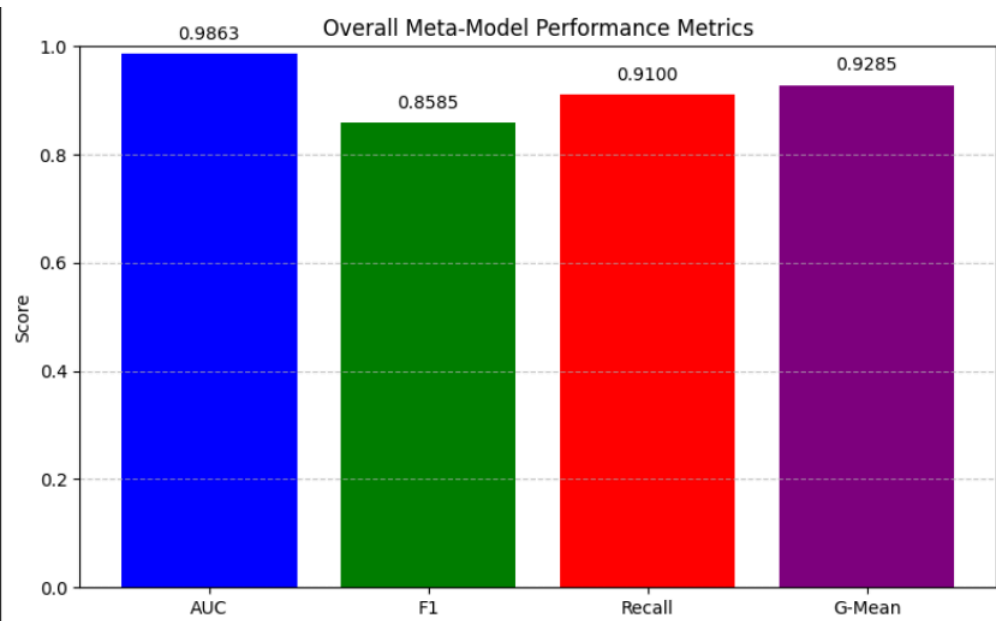
Result per each model

Time Step	Test Acc	Test Pos Acc	Test Neg Acc	Num Licit	Num Illicit
1	0.933868	0.87	0.949875	42019	4545
2	0.945892	0.87	0.964912	42019	4545
3	0.943888	0.93	0.947368	42019	4545
4	0.943888	0.91	0.952381	42019	4545
5	0.9499	0.89	0.964912	42019	4545
6	0.933868	0.87	0.949875	42019	4545
7	0.935872	0.89	0.947368	42019	4545
8	0.9499	0.89	0.964912	42019	4545
9	0.945892	0.88	0.962406	42019	4545

per-mode performance



metric results



Challenges Encountered and Resolutions

- **Challenge 1: Extreme Class Imbalance**
Illicit transactions form a small minority of labeled data.
Resolution: Ensemble learning and attention mechanisms were employed to mitigate bias toward the majority class.
 - **Challenge 2: Computational Constraints**
Training on the full graph is memory-intensive.
Resolution: Subgraph sampling and SGAT significantly reduced memory and runtime requirements.
 - **Challenge 3: Model Variance**
Single GNN models showed sensitivity to initialization and sampling.
Resolution: Bagging and stacking were introduced to improve robustness and generalization.
-
-

Student Signature: Rasoul Salehi , Sayed Ali Amiri, Mojtaba Mollaei
Date: [2026-01-07]
