Superviz Plenary Meeting WP4-TH4.1

Yufei Han@INRIA PIRAT

March 11, 2025

@Campus Cyber, Paris

- Why Graph-based Intrusion Detection ?
- Graph provides a structural representation of cyber attack behaviour



A survey on graph neural networks for intrusion detection systems: Methods, trends and challenges, Computer & Security, Vol.141, June 2024

• Target: Flow Graph based representation of network traffics





• A few words about Graph-Sage based GNN models



• What is an adversarial attack ?



• An adversarial attack against Graph Neural Network-based models



- What is the problem space constraint ?
- Any change to the input graph should not break the attack / normal traffic flows
 - An attack is still an attack
- Any change to the input graph should be made consistent with the deployed communication protocols
- In the end, any change to the input graph should be made compatiable with the profiles of real traffic flows between IP addresses.

• An example of problem-space adversarial attacks



Problem-space constrant:

Edge addition: An attacker can send traffic flows from a spoofed IP address to an IP address existing in the raw traffic records.

What edges to add: The added traffic flows should be added in a way consistent with the communication protocol, such as TCP-IP

A spoofed IP adress C2

• Mathematical nature of the edge-moification-based adversarial attacks



KDD '20, August 23–27, 2020, Virtual Event, USA KDD '20, August 23–27, 2020, Virtual Event, USA Problem Space-constrained Adversarial Attack against Graph-based Intrusion Detection

 Every GNN is vulnerable to edge modification-based adversarial attacks

> Let y be the target class label of the evasion attack. The goal is to make $f_y(\hat{\mathbf{x}})$ deviate from $f_y(\mathbf{x}) = 0$. In other words, we aim at maximizing $f_y(\hat{\mathbf{x}})$, so that x is modified to get y assigned to it. The evasion attack task can then be formulated as a process of set function optimization, defined as

$$S^* = \underset{|S| \le K}{\operatorname{arg max}} g(S)$$

where $g(S) = \underset{l \in S}{\max} f_y(\hat{\mathbf{x}}), \quad l = diff(\mathbf{b}, \hat{\mathbf{b}})$

T

This is in nature a weakly-submodular function maximization problem , which can be solved with Greedy Search with a polynormal complexity.

Algorithm 1: Unconstrained Adversarial Attack **Input:** Original flow graph $\mathcal{G} = \{A, X, E\}$, attack budget b, target model \mathcal{M} , loss function ℓ **Output:** Modified flow graph $\mathcal{G}' = \{A', X, E'\}$ 1 Initialize $A' \leftarrow A, E' \leftarrow E;$ 2 for $i \leftarrow 1$ to b do max loss $\leftarrow -\infty$; 3 for each controlled node u in \mathcal{G} do 4 for each possible destination node $v \neq u$ do 5 Simulate adding edge (u, v) to A', with 6 corresponding edge features $E'_{(u,v)}$; Compute loss: 7 $current_loss \leftarrow \ell(\mathcal{M}(\mathbf{A}', \mathbf{X}, \mathbf{E}'), \mathbf{Y});$ if $current_loss > max_loss$ then 8 $max_loss \leftarrow current_loss;$ 9 $best_edge \leftarrow (u, v);$ 10 Add *best_edge* to A', update E' with its features; 11

12 return $G' = \{A', X, E'\};$



- Questions remaining to answer
- We focus on activating an excessively large fase alarm rate as the attack goal, but can we also generate evasive samples (increasing the false negative rate ?)
 - Adding benign traffic flows targeting at one IP address in the dataset ?
- Divide the raw network traffics according to their temporal orders.
 - Attackers add adversarial manipulations to historical data, while aiming to deliver adversarial attack impacts over future traffic flows
- The impact of the degree of the target IP address
- Transferability of the attack
 - If an adversarial attack can mislead GNN-based intrusion detection systems, can it also mislead a flow-statistic based intrusion detection model (or not) ?
 - Think about other GNN models

- Next step
- Elevating both false alarm and false negative rate via attack
- Where to add edges and what edges to add ?
- Transferable attacks across both the graph-based IDS and statisticsbased IDS
- Publication plan:
 - Matthieu Mouzaoui at the second year of this thesis.
 - 2 papers submitted by the end of this year.

Thanks !