

# Efficient, Direct, and Restricted Black-Box Graph Evasion Attacks to Any-Layer Graph Neural Networks via Influence Function

## ABSTRACT

Graph neural networks (GNNs) have achieved state-of-the-art performance in many graph-related tasks, e.g., node classification. However, recent works show that GNNs are vulnerable to graph evasion attacks, i.e., an attacker slightly perturbing the graph structure can fool trained GNN models. Existing evasion attacks to GNNs have one or more of the following drawbacks: 1) limited to only directly attack two-layer GNNs, while using (not effective enough) surrogate models to attack multi-layer GNNs; 2) inefficient, as they involve intensive computation; and 3) impractical, as they need to know full or part of GNN model parameters.

We address the above drawbacks and propose an influence-based *efficient, direct, and restricted black-box* evasion attack to *any-layer* GNNs, a completely different perspective from the existing works. Specifically, we first introduce two influence functions, i.e., feature-label influence and label influence, that are defined on GNNs and label propagation (LP), respectively. Then we observe that GNNs and LP are strongly connected in terms of our defined influences. Based on this, we can then reformulate the evasion attack to GNNs as calculating label influence on LP, which is *inherently* applicable to any-layer GNNs, while no need to know information about the internal GNN model. Finally, we propose an efficient algorithm to calculate label influence. We evaluate our influence-based attack on three benchmark graph datasets. Experimental results show that, compared to state-of-the-art white-box attacks, our attack can achieve comparable attack performance, but has a 5-50x speedup when attacking two-layer GNNs. Moreover, our attack is effective to attack multi-layer GNNs<sup>1</sup>.

## CCS CONCEPTS

• Security and privacy; • Computing methodologies → Machine learning;

## KEYWORDS

graph neural network, label propagation, attack, influence function

### ACM Reference Format:

. 2018. Efficient, Direct, and Restricted Black-Box Graph Evasion Attacks to Any-Layer Graph Neural Networks via Influence Function. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/1122445.1122456>

<sup>1</sup>Our source code and the full version are in the github: <https://shorturl.at/aeinV>

Unpublished working draft. Not for distribution.

Permission to make digital or hard copies of all or part of this work for personal or academic use, not for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

2023-06-02 03:25. Page 1 of 1–12.

## 1 INTRODUCTION

Learning with graph data, such as social networks, biological networks, financial networks, has drawn continuous attention recently. Moreover, graph neural network (GNN) has become the mainstream methodology for representation learning on graphs. GNN was first introduced in [27], which extended conventional neural network to process graph data. Then, various GNN methods have been proposed and achieved state-of-the-art performance in many graph-related tasks such as node classification [19, 35, 43], graph classification [14, 15], and link prediction [45]. However, recent works [7, 11, 23, 25, 30, 36, 38, 40, 42, 52, 53] show that GNNs are vulnerable to graph evasion attacks—Given a target node and a trained GNN model, an attacker slightly perturbing the graph structure<sup>2</sup> (e.g., add new edges to or delete existing edges from the graph) can make the GNN model misclassify the target node. Existing attacks to GNNs can be roughly classified as *optimization-based* attacks [38, 40, 42, 52] and *reinforcement learning (RL)*-based attacks [7, 11, 30].

In this paper, we focus on optimization-based attacks, as they are shown to be more effective [52]. Optimization-based attacks first formulate the graph evasion attack as a binary optimization problem, which is challenging to solve, and then design approximate algorithms to solve a tractable optimization problem. Although achieving promising attack performance, existing optimization-based attacks have one or more of the below key limitations:

- First, most of the existing attacks need to know the full/partial GNN model parameters, which is unrealistic in many real-world applications, e.g., when GNN models are confidential due to their commercial value and are deployed as an API. Thus, the practicability of the existing attacks are limited. Further, they are mainly designed to attack *two-layer GNNs*, while GNNs are multi-layer in essence. To attack multi-layer GNNs, they often first *indirectly* attack a surrogate two-layer GNN model, and then transfer the attack to the target multi-layer GNN. However, this strategy is not effective enough (See Figure 6 in Section 5).
- Second, they are not efficient, as they involve intensive computation, i.e., by multiplying GNN model parameters of different layers and with node feature matrix. If a GNN has many layers, such computation can be a bottleneck, especially for attackers who have limited computational resources or/and want to perform real-time attacks. For example, many fraud detection systems, such as detecting fake users in social networks and detecting anomalies from system logs, are updated frequently in order to reduce the loss caused by the evasion attacks' malicious activities. In these scenarios, efficiency is a major concern for the attack and an attacker performing efficient attacks is necessary, as otherwise the detection system may have already updated and identified the attack's malicious patterns before the attack is implemented.

<sup>2</sup>An attacker can also perturb node features to perform the attack. However, structure perturbation is shown to be much more effective than node feature perturbation [52].

**Our work:** We aim to address the above limitations in this paper. To this end, we propose an optimization-based evasion attack against any-layer GNNs based on influence function [20]—a completely different perspective from the existing works. Our influence-based attack is motivated by the strong connection between GNNs and label propagation (LP) [50]. Specifically, we first introduce two influence functions, i.e., feature-label influence and label influence, that are defined on GNNs and LP, respectively. Then, we prove that our label influence defined on LP is equivalent to feature-label influence on a particular well-known type of GNN, called Graph Convolutional Network (GCN) [19] (and its linearized version Simple Graph Convolutional (SGC) [39]). Based on this connection, we reformulate the evasion attack against GNNs to be related to calculating label influence on LP. As our influences are designed for any-layer GNNs, our attack is inherently applicable to attack any-layer GNNs. Note that label influence can be computed easily and we also design an efficient algorithm to compute it. Further, as our influence-based attack does not need to know any information about the GNN model (except the target node’s neighboring information), it is a more practical (restricted black-box) attack. Finally, we evaluate our attack against GCN/SGC on three benchmark graph datasets (i.e., Cora, Citeseer, and Pubmed). Compared to the state-of-the-art white-box attacks against two-layer GCN/SGC, our attack can achieve comparable attack performance but has a 5-50x speedup. Our attack is more effective to attack multi-layer GCN/SGC. For instance, our attack achieves a 93% attack success rate, when perturbing 4 edges per target node on Cora, while the surrogate model based attack only has 80% attack success rate. As a by product, our attack also shows promising transferability to attack other GNNs, and is more effective than existing black-box attacks. Our contributions can be summarized as follows:

- We propose graph evasion attacks to GNNs based on influence function, which is a completely new perspective.
- Our attack is effective, direct, efficient, and practical.
- Our attack has promising transferability.

## 2 RELATED WORK

**Attacks to graph neural networks.** Existing attacks to GNNs can be classified as graph *poisoning attacks* [7, 22, 30, 31, 41, 42, 47, 48, 52, 53] and *evasion attacks* [7, 23, 24, 40, 52]. In poisoning attacks, an attacker modifies the graph structure during the training process such that the trained GNN model has a low prediction accuracy on testing nodes. For instance, Zügner and Günnemann [53] proposed a poisoning attack, called Metattack, that perturbs the whole graph based on meta-learning. Xu et al. [42] developed a topology poisoning attack based on gradient-based optimization. Evasion attacks can be classified as untargeted attacks and targeted attacks, where the latter is more challenging. Given a target node and a trained GNN model, targeted attack means an attacker aims to perturb the graph structure such that the GNN model misclassifies the target node to be a target label, while untargeted attack misclassifies the target node to be an arbitrary label different from the target node’s label. For instance, Dai et al. [7] leveraged reinforcement learning techniques to design non-targeted evasion attacks to both graph classification and node classification. Zügner et al. [52] proposed a targeted evasion attack, called Nettack, against two-layer GCN

and achieved the state-of-the-art attack performance. Specifically, Nettack learns a surrogate linear model of GCN by removing the ReLU activation function and by defining a graph structure preserving perturbation that constrains the difference between the node degree distributions of the graph before and after attack. Our label influence-based attack is a targeted evasion attack.

Most of the existing GNN attacks are white/gray-box. Recently, two black-box attacks to GNNs [25, 38] have been proposed. For instance, Wang et al. [38] formulate the black-box attack to GNNs as an online optimization with bandit feedback. The original problem is NP-hard and they then propose an online attack based on (relaxed) bandit convex optimization which is proven to be sublinear to the query number. Our attack is a restricted black-box attack, where the attacker only needs to know the target node’s neighbors.

**Attacks to other graph-based methods.** Besides attacking GNNs, other adversarial attacks against graph data include attacking graph-based clustering [6], graph-based collective classification [34, 36], graph embedding [1, 4, 5, 8, 29], community detection [21], etc. For instance, Chen et al. [6] proposed a practical attack against spectral clustering, which is a well-known graph-based clustering method. Wang and Gong [36] designed an optimization-based attack against the collective classification method, called linearized belief propagation, by modifying the graph structure.

**Defending against graph perturbation attacks.** Existing defenses against the graph perturbation attacks can be classified as empirical defenses [10, 32, 33, 40, 42, 49] and provable defenses [2, 3, 17, 54]. The empirical defenses are shown to be easily broken by stronger/adaptive attacks [13, 26]. Provable defenses study certified robustness of GNNs against the worst-case graph perturbation attacks. For instance, Wang et al. [37] design a randomized smoothing-based provable defenses that achieves a tight certified robustness, when there are no assumptions about the GNN model. [37] achieves the state-of-the-art provable defense performance.

## 3 BACKGROUND

### 3.1 Graph Neural Network

Let  $G = (\mathcal{V}, \mathcal{E}, \mathbf{X})$  be a graph, where  $u \in \mathcal{V}$  is a node,  $(u, v) \in \mathcal{E}$  is an edge between  $u$  and  $v$ , and  $\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_n] \in \mathbb{R}^{n \times d}$  is the node feature matrix. We denote  $\mathbf{A} = [\mathbf{a}_1; \mathbf{a}_2; \dots; \mathbf{a}_n] \in \{0, 1\}^{n \times n}$  as the adjacency matrix, where  $A_{u,v} = 1$ , if  $(u, v) \in \mathcal{E}$  and  $A_{u,v} = 0$ , otherwise; Moreover, we denote  $d_u$  and  $\Gamma_u$  as  $u$ ’s node degree and the neighborhood set of  $u$  (including self-loop  $(u, u)$ ). We consider GNNs for node classification in this paper. In this context, each node  $u \in \mathcal{V}$  has a label  $y_u$  from a label set  $\mathcal{Y} = \{1, 2, \dots, C\}$ . Given a set of  $\mathcal{V}_L \subset \mathcal{V}$  labeled nodes  $\{(\mathbf{x}_u, y_u)\}_{u \in \mathcal{V}_L}$  as the training set, GNN for node classification is to take the graph  $G$  and labeled nodes as input and learn a node classifier that maps each node  $u \in \mathcal{V} \setminus \mathcal{V}_L$  to a class  $y \in \mathcal{Y}$ . In this paper, we focus on Graph Convolutional Network (GCN) [19], a widely used type of GNN, and its special case Simple Graph Convolution (SGC) [39].

**GCN.** GCN is motivated by spectral graph convolution [9]. Suppose GCN has  $K$  layers. We denote node  $v$ ’s representation in the  $k$ -th layer as  $\mathbf{h}_v^{(k)}$ , where  $\mathbf{h}_v^{(0)} = \mathbf{x}_v$ . Then, GCN has the following form to update the node representation:

$$\mathbf{h}_v^{(k)} = \text{ReLU}\left(\mathbf{W}^{(k)} \left( \sum_{u \in \Gamma_v} d_u^{-1/2} d_v^{-1/2} \mathbf{h}_u^{(k-1)} \right)\right). \quad (1)$$

A node  $v$ 's final representation  $\mathbf{h}_v^{(K)} \in \mathbb{R}^{|\mathcal{Y}|}$  can capture the structural information of all nodes within  $v$ 's  $K$ -hop neighbors. Moreover, the final node representations of training nodes are used for training the node classifier. Specifically, let  $\Theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(K)}\}$  be the model parameters and  $v$ 's output be  $f_\Theta(\mathbf{A})_v = \text{softmax}(\mathbf{h}_v^{(K)}) \in \mathbb{R}^{|\mathcal{Y}|}$ , where  $f_\Theta(\mathbf{A})_{v,y}$  indicates the probability of node  $v$  being class  $y$ . Then,  $\Theta$  are learnt by minimizing the cross-entropy loss on the training nodes  $\mathcal{V}_L$ , i.e.,  $\Theta^* = \arg \min_{\Theta} - \sum_{v \in \mathcal{V}_L} \ln f_{\Theta^*}(\mathbf{A})_{v,y}$ . With the learnt  $\Theta^*$ , we can predict the label for each unlabeled nodes  $u \in \mathcal{V} \setminus \mathcal{V}_L$  as  $\hat{y}_u = \arg \max_y f_{\Theta^*}(\mathbf{A})_{u,y}$ .

**SGC.** SGC is a linearized version of GCN. Specifically, its node representation is updated as follows:

$$\mathbf{h}_v^{(k)} = \mathbf{W}^{(k)} \left( \sum_{u \in \Gamma_v} d_u^{-1/2} d_v^{-1/2} \mathbf{h}_u^{(k-1)} \right). \quad (2)$$

SGC has shown to have comparable node classification performance with GCN, but is much more efficient than GCN.

### 3.2 Label Propagation

Label Propagation (LP) is a conventional semi-supervised node classification method without training. The key idea behind LP is that two nodes having a high similarity (e.g., connected nodes in a graph) are likely to have the same label. Thus, LP iteratively propagates labels among the graph to unlabeled nodes based on node-pair similarity. Let  $\mathbf{y}_v \in \mathbb{R}^{|\mathcal{Y}|}$  be node  $v$ 's initial label vector (For notation reason, one should note that  $y_v$  is  $v$ 's categorical label). For instance,  $\mathbf{y}_v$  can be  $v$ 's one-hot label vector if  $v$  is a labeled node, and  $\mathbf{y}_v = \mathbf{0}$ , otherwise. Then, LP is formulated as follows:

$$\mathbf{y}_v^{(k)} = \sum_{u \in \Gamma_v} d_u^{-1/2} d_v^{-1/2} \mathbf{y}_u^{(k-1)}, \quad \mathbf{y}_v^{(0)} = \mathbf{y}_v. \quad (3)$$

With  $K$  iterations, an unlabeled node  $u$  is predicted to be class  $c$ , if  $c = \arg \max_i y_{u,i}^{(K)}$ .

**GNN vs. LP:** Viewing Equation (3) and Equations (1) and (2), we observe that LP and GNNs have similar iterative processes: LP propagates node labels  $\mathbf{y}_v$ , while GNNs propagate node features  $\mathbf{x}_v$ . The key difference is that LP does not involve model parameters, while GNN involves multiplying the parameter matrix  $\mathbf{W}^{(k)}$  in each  $k$ -th layer.

### 3.3 Problem Definition

We consider targeted evasion attacks<sup>3</sup> to GNNs. Suppose we are given a trained GNN model  $f_{\Theta^*}$  for node classification. We assume  $v$  is the *target node* and  $c$  is the *target label*. We consider an attacker can perturb the graph structure (i.e., add new edges to or delete existing edges from the graph) in order to make  $f_{\Theta^*}$  misclassify the target node  $v$  to be the target label  $c$ . We call the modified edges by the attacker as *attack edges*. In particular, we consider a practical *direct attack* [52], where an attacker can only modify the connection status between  $v$  and other nodes in the graph, while cannot modify the connection status among other nodes. We denote the perturbed graph as  $\tilde{G}$  (with the perturbed adjacency matrix  $\tilde{\mathbf{A}}$ ) after the attack and the attack budget as  $\Delta$ , i.e., at most  $\Delta$  edges can be perturbed for the target node. Then, the objective function of targeted evasion

attacks to GNNs is formally defined as:

$$\begin{aligned} \max_{\tilde{\mathbf{A}}_v} (f_{\Theta^*}(\tilde{\mathbf{A}})_{v,c} - f_{\Theta^*}(\tilde{\mathbf{A}})_{v,y_v}) &\Leftrightarrow \max_{\tilde{\mathbf{A}}_v} ([\tilde{\mathbf{h}}_v^{(K)}]_c - [\tilde{\mathbf{h}}_v^{(K)}]_{y_v}), \\ \text{s.t.}, \quad \sum_s |\tilde{A}_{v,s} - A_{v,s}| &\leq \Delta, \end{aligned} \quad (4)$$

where  $\tilde{\mathbf{h}}_v^{(K)}$  is  $v$ 's representation on the perturbed graph  $\tilde{G}$ .

A target node is called a success to attack the GNN model if the value of the attack's objective function is larger than 0, under the attack budget. Note that Equation (4) is a binary optimization problem and is challenging to solve in practice. Zügner et al. [52] proposed an optimization-based attack method, called *Nettack*, against two-layer GCN. Specifically, *Nettack* attacked a substitute GNN model (actually SGC) that removed the ReLU activation function in GCN. *Nettack* has achieved state-of-the-art attack performance. However, it is inefficient as it involves dense matrix multiplication (i.e., model parameters multiply node features); it also needs to know model parameters  $\Theta^*$ , and can only attack two-layer GNNs.

## 4 INFLUENCE-BASED EVASION ATTACK

In this section, we propose our evasion attack against GNNs via influence function. In contrast to existing optimization-based attacks that only focus on two-layer GNNs, our attack is applicable to any-layer GNNs. Specifically, we first define two influence functions associated with GNNs and LP, respectively, and build an equivalence relation between GNNs and LP with the defined influences. Next, we reformulate the attack objective function as relating to label influence defined on LP. Finally, we design an efficient algorithm to calculate label influence and realize our attack.

### 4.1 Equivalence between GNNs and LP in terms of Influence

**4.1.1 Motivation.** Due to GNN's complex network structure, existing optimization-based evasion attacks can only attack two-layer GNNs *directly*. However, we note that LP has a similar iterative process to GNNs, but it has good properties, e.g., LP does not involve model parameters. Motivated by this, we aim to discover an equivalence relation between LP and GNNs, such that the challenging problem of attacking multi-layer GNNs can be converted to a relatively easier problem by leveraging good properties of LP. We notice that influence function [20, 44] is an appropriate tool to bridge the gap, and our purpose is to explore equivalent influence functions defined on LP and on GNNs, respectively. As the attacker's goal is to change the target node's label, we thus need to define influences associated with the node label. As LP propagates node labels, we can naturally design the *label influence* function (see below Equation (6)). In addition, GNNs involve propagating node features. In order to also leverage node labels, we integrate both node features and node labels and design the *feature-label influence* function (see below Equation (5)). Next, we introduce our influence functions.

**4.1.2 Influence function.** Given two nodes  $u$  and  $v$ , an influence of  $u$  on  $v$  indicates how the output (e.g., final node representation in GNNs or estimated node label in LP) of  $v$  changes if the input of  $u$  is slightly perturbed. Inspired by [20, 44], we define the following feature-label influence on GNN and label influence on LP.

<sup>3</sup>As untargeted attacks are less powerful than targeted attacks, we only consider targeted attacks in this paper for simplicity.



**Definition 1** (Feature-label influence). We define the feature-label influence of node  $u$  on node  $v$  associated with  $u$ 's label on a  $K$ -layer GNN as follow:

$$I_{fl}(v, u; K) = \left\| \left[ \frac{\partial \mathbf{h}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_{y_u} \right\|_1 = \mathbf{1}_{y_u}^T \cdot \frac{\partial \mathbf{h}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)}, \quad (5)$$

where  $\mathbf{1}_{y_u} = [y_1, y_2, \dots, y_n]$  is an indicator vector where  $y_i = 1$  if  $i = u$  and  $y_i = 0$ , otherwise;  $\|\cdot\|_1$  is the vector  $\ell_1$ -norm;  $T$  is a transpose; and  $\mathbf{h}_u^{(0)} = \mathbf{x}_u$  is  $u$ 's node features.

**Definition 2** (Label influence). We define the label influence of node  $u$  on node  $v$  after  $K$  iterations of label propagation as follows:

$$I_l(v, u; K) = \frac{\partial y_v^{(K)}}{\partial y_u^{(0)}}. \quad (6)$$

Then, we have the following theorem showing the equivalence between GNN and LP in terms of influence.

**THEOREM 4.1.** If the GNN is a GCN/SGC, then:

$$I_{fl}(v, u; K) = C \cdot I_l(v, u; K), \quad (7)$$

where  $C = \rho \mathbf{1}_{y_u}^T [\prod_{l=1}^K \mathbf{W}^{(l)}] \mathbf{x}_u$  is constant related to GNN model parameters  $\Theta = \{\mathbf{W}^{(l)}\}$  and  $u$ 's node features  $\mathbf{x}_u$ .

**PROOF.** See the full version in <https://shorturl.at/aeinV>.  $\square$

Theorem 4.1 reveals that: given arbitrary node  $v$ , the feature-label influence defined on  $K$ -layer GCN/SGC of any other node  $u$  on the node  $v$  and the label influence defined on  $K$ -iteration LP of node  $u$  on the node  $v$  are equal (with a constant multiplier difference).

## 4.2 Reformulating Evasion Attacks to Any-Layer GNNs as Calculating Label Influence on LP

Based on our influence functions and Theorem 4.1, we can first restate the challenging problem of attacking  $K$ -layer GNNs in Equation (4) in the form of feature-label influence, and further convert it to an equivalent problem related to label influence on LP. Before going into details, we first introduce the following lemma:

**LEMMA 1** (XU ET AL.[44]). Given a  $K$ -layer GCN. Assume all paths in the computation graph of the GCN model are activated (i.e., via ReLU) with the same probability of success  $\rho$ . Then,

$$\frac{\partial \mathbf{h}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} = \rho \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 a_{v_p^l, v_p^{l-1}} \cdot \mathbf{W}^{(l)}, \quad (8)$$

where  $\Psi_{v \rightarrow u}$  is the total number of paths  $[v_p^K, v_p^{K-1}, \dots, v_p^1, v_p^0]$  of length  $K+1$  from node  $v$  to the node  $u$  with  $v_p^K = v$  and  $v_p^0 = u$ . For  $l = 1, \dots, K$ ,  $v_p^{l-1} \in N(v_p^l)$ ,  $a_{v_p^l, v_p^{l-1}} = d_{v_p^l}^{-\frac{1}{2}} d_{v_p^{l-1}}^{-\frac{1}{2}}$  is the normalized weight of the edge  $(v_p^l, v_p^{l-1})$  in the path  $p$ .  $\Theta = \{\mathbf{W}^{(l)}\}$  is the  $K$ -layer GCN model parameters.

Then, according to Equation (8) in Lemma 1, the target node  $v$ 's final node representation  $\tilde{\mathbf{h}}_v^{(K)}$  learnt on the perturbed graph  $\tilde{G}$  can be expressed as  $\tilde{\mathbf{h}}_v^{(K)} = \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)}$ , where  $\tilde{\Lambda}_v^{(K)}$  is the node set containing  $v$ 's neighbors within  $K$ -hop on the perturbed

## Algorithm 1 Efficient calculation of label influence via DFS

**Input:** Path length  $K+1$ , the node  $n$  on the path, cumulative sum of weights  $s$ , target node  $v$  and label  $y_v$ , target label  $c$ .

**Output:** Difference between the label influence of label- $y_v$  nodes and label- $c$  nodes.

```

1:  $I_{y_v} \leftarrow 0, I_c \leftarrow 0$ 
2: function LABELINFLUENCE( $K, p, s$ )
3:   if  $K = 0$  then
4:      $I_{y_v} \leftarrow I_{y_v} + s \cdot y_{v, y_v}; I_c \leftarrow I_c + s \cdot y_{v, c};$  return
5:   end if
6:   for  $u \in \Gamma_p$  do
7:      $w = d_p^{-\frac{1}{2}} d_u^{-\frac{1}{2}}; \text{LabelInfluence}(K-1, u, s \cdot w);$ 
8:   end for
9: end function
10: return  $I_c - I_{y_v}$ 

```

graph  $\tilde{G}$ , i.e., after modifying the connection status between the target node  $v$  and other nodes in the clean graph  $G$ .

Then, the attack's objective function in Equation (4) is equivalent to the following objective function:

$$\begin{aligned} \max_{\tilde{\Lambda}_v} & \left( \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c - \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_{y_v} \right) \\ \text{s.t.,} & \sum_s |\tilde{A}_{v,s} - A_{v,s}| \leq \Delta, \end{aligned} \quad (9)$$

Finally, based on the following Assumption 1 and Theorem 4.1, we reach Theorem 4.2 that reformulates the evasion attack's objective function via label influence. We also conduct experiments (See Section 5.2) to verify that Assumption 1 holds in practice.

**Assumption 1.** Given a target node  $v$  and a target label  $c$ . We assume that any node  $u$ , within the  $K$ -hop neighbor of  $v$ , has a negligible feature-label influence on  $v$  if  $u$  is not a label- $c$  node. Formally,

$$\left[ \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c \approx 0, \quad \forall u \in \tilde{\Lambda}_v^{(K)}, y_u \neq c. \quad (10)$$

**THEOREM 4.2.** Let  $\tilde{I}_l(v, u; K)$  be the label influence of node  $u$  on the target node  $v$  with  $K$  iterations of LP after the attack. Then, the attack's objective function in Equation (4) equals to the following objective function on label influence:

$$\begin{aligned} \max_{\tilde{\Lambda}_v} & \left( \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u = c} \tilde{I}_l(v, u; K) - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z = y_v} \tilde{I}_l(v, z; K) \right), \\ \text{s.t.,} & \sum_s |\tilde{A}_{v,s} - A_{v,s}| \leq \Delta, \end{aligned} \quad (11)$$

where  $\tilde{I}_l(v, u; K)$  is defined as:

$$\tilde{I}_l(v, u; K) = \sum_{p=1}^{\tilde{\Psi}_{v \rightarrow u}} \prod_{l=K}^1 \tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}}, \quad (12)$$

where  $\tilde{\Psi}_{v \rightarrow u}$  is the total number of paths  $[v_p^K, v_p^{K-1}, \dots, v_p^1, v_p^0]$  of length  $K+1$  from  $v$  to  $u$  on the perturbed graph  $\tilde{G}$ , where  $v_p^K = v$  and  $v_p^0 = u$ .  $\tilde{d}_u$  is  $u$ 's degree on the perturbed graph  $\tilde{G}$  and  $\tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}}$  is the normalized weight of the edge  $(v_p^l, v_p^{l-1})$  in path  $p$  in  $\tilde{G}$ .

**PROOF.** See the full version in <https://shorturl.at/aeinV>.  $\square$

**Algorithm 2** (Label) influence-based targeted evasion attack**Input:** Adj. mat.  $A$ , layer  $K$ , target node  $v$ , target label  $c$ , budget  $\Delta$ .**Output:** *success* (I.e.,  $v$  attacks successfully or not)

```

1:  $AE \leftarrow 0$ ,  $success \leftarrow False$ ,  $\tilde{A} \leftarrow A$ ;
2:  $N_A, N_B \leftarrow$  Find two candidate node sets from  $A$ .
3: for  $a \in N_A$  do  $\Delta I_A(a) = d_v^{-\frac{1}{2}} d_a^{-\frac{1}{2}} \cdot LabelInfluence(K-1, a, 1)$ 
4: end for
5: for  $b \in N_B$  do  $\Delta I_B(b) = d_v^{-\frac{1}{2}} d_b^{-\frac{1}{2}} \cdot LabelInfluence(K-1, b, 1)$ 
6: end for
7: while  $AE < \Delta$  do
8:    $C_A \leftarrow LabelInfluence(K, v, 1)$ ; //  $d_v \leftarrow d_v + 1$ 
9:    $C_B \leftarrow LabelInfluence(K, v, 1)$ ; //  $d_v \leftarrow d_v - 1$ 
10:   $u^* \leftarrow \text{argsort}(\{C_A + \Delta I_A\} \cup \{C_B - \Delta I_B\})$ ;
11:  if  $u^* \in N_A$  then  $\tilde{A} \leftarrow A + (v, u^*)$ ;
12:  else  $\tilde{A} \leftarrow A - (v, u^*)$ ;
13:  end if
14:   $AE \leftarrow AE + 1$ 
15: end while
16: if  $(f_{\tilde{A}}(\tilde{A})_{v,c} - f_{\tilde{A}}(\tilde{A})_{v,y_v}) > 0$  then  $success \leftarrow True$ ;
17: end if
18: return  $success$ 

```

We have the following observations from Theorem 4.2.

- Our attack does not need to operate on model parameters  $\Theta^*$ , different from existing attacks that involve dense multiplication on  $\Theta^*$ . Thus, our attack is more efficient.
- Our attack can be applied to any layer GNN, as the label influence is defined for general  $K$ -iteration LP. However, most of the existing attacks can only directly attack two-layer GNNs. Thus, our attack is more practical.
- The only information our attack needs to know is the target node  $v$ 's within  $K$ -hop neighbors, whose labels are  $y_v$  or  $c$ . In practice, if the labels of these node are unknown, we can estimate them via querying the GNN model, and treat the estimated labels as the true labels. Thus, our attack can be seen a restricted black-box attack.

Next, we show how to fast calculate the label influence and design our influence-based targeted evasion attack.

### 4.3 Efficient Calculation of Label Influence

According to Theorem 4.2, the attack's goal is to select the minimum set of nodes such that when changing the connection status between the target node  $v$  and these selected nodes, the difference between the two label influence terms will be maximized. Observing Equation (11), we note that the two label influence terms are defined on two sets of nodes: a set of nodes having the same label as the target label  $c$ , and a set of nodes having the same label as the target node's label  $y_v$ . Intuitively, if we add an edge between  $v$  and a label- $c$  node, we can make  $v$  be close to label  $c$ ; and if we remove an edge between  $v$  and a label- $y_v$  node, we can make  $v$  away from label  $y_v$ . Thus, our idea to solve Equation (11) is as follows:

- First, we define a candidate set  $N_A \subset \{y_u = c, u \in \Lambda_v^{(K)}\}$  which contains label- $c$  nodes that are *not* connected with  $v$  in the clean graph, as well as a candidate set  $N_B \subset \{y_z = y_v, z \in \Lambda_v^{(K)}\}$  which contains label- $y_v$  nodes that are connected with  $v$  in the clean graph. We denote  $S$  as the final selected nodes from  $N_A$  and  $N_B$ , and initialize  $S = \{\}$ . For each node  $u \in N_A \cup N_B \setminus S$ ,

**Table 1: Dataset statistics.**

Dataset	#Nodes	#Edges	#Features	#Classes
Cora	2,708	5,429	1,433	7
Citeseer	3,327	4,732	3,703	6
Pubmed	19,717	44,338	500	3
OGB-arxiv	169,343	1,166,243	128	40

we change the connection status between  $v$  and  $u$  and compute the gap between two label influence terms.

- Next, we record the node  $u^*$  that obtains the largest positive gap. Then, we modify the connection status between  $v$  and  $u^*$ , calculate the value of the attack's objective function, and update  $S = S \cup \{u^*\}$ .
- We repeat above steps at most  $\Delta$  times and break if the value of attack's objective function is bigger than 0. Finally, we have the attack edges  $\{(v, u^*), u^* \in S\}$ .

However, note that when modifying the connection status between  $v$  and  $u^*$ , the normalized weight for all edges containing  $u^*$  in all paths  $\tilde{y}_{v \rightarrow u}$  in Equation (12) should be recalculated. When the candidate set has a large size or/and the number of recalculated edge weights is large, calculating the exact label influence will have a large computational complexity. To solve the problem, we propose an approximate algorithm to efficiently compute the label influences. More details are in the full version <https://shorturl.at/aeinV>.

Algorithm 1 illustrates how we efficiently calculate the label influences via depth first search (DFS), and Algorithm 2 shows the details of implementing our attack.

## 5 EVALUATION

### 5.1 Experimental Setup

**Datasets.** Following the existing works [12, 42, 52], we use three benchmark citation graphs (i.e., Cora, Citeseer, and Pubmed) [16, 28] to evaluate our attack. In these graphs, each node represents a documents and each edge indicates a citation between two documents. Each document treats the bag-of-words feature as the node feature vector, and has a label. Table 1 shows basic statistics of these graphs. **Training nodes and target nodes.** We use the public training nodes to train GNN models, and target nodes to evaluate attacks against the trained GNN models. For the target nodes, we employ a random sampling technique to select 100 nodes that are correctly classified by each GNN model as the target nodes. Similar to Netttack [52], for each target node, we choose the predicted label by the GNN model with a second largest probability as the target label.

We compare our influence-based attack with the state-of-the-art Netttack [52] for attacking two particular GNNs: GCN and SGC. Note that Netttack is mathematically designed to only attack two-layer GNNs and cannot directly attack multi-layer GNNs. To attack multi-layer GNNs, Netttack needs to be performed via an indirect way: It first attacks a surrogate two-layer GNN to generate the attack edges, and then transfers these attack edges to attack the target multi-layer GNNs. When computing the label influence, our attack needs to know the labels of unlabeled nodes in the graph. When our attack knows the true labels, we denote it as **Ours-KL**. When the true labels are unknown, our attack first queries the learnt GNN model to estimate labels for unlabeled nodes and then uses the estimated labels as the true labels. We denote this variant as **Ours-UL**. As a comparison, we also test our attack that is

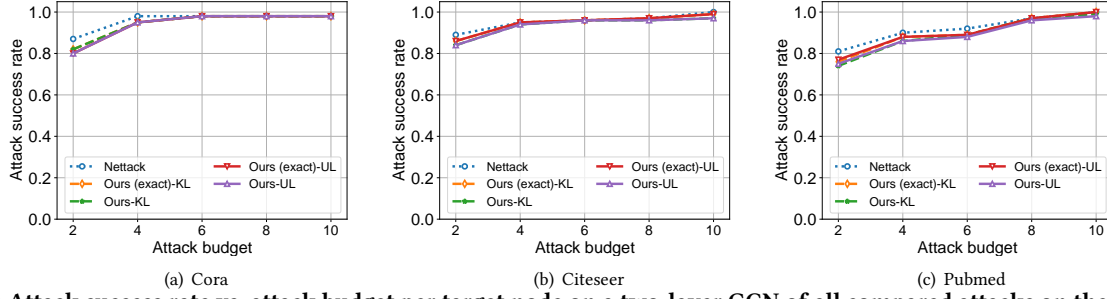


Figure 1: Attack success rate vs. attack budget per target node on a two-layer GCN of all compared attacks on the three graphs.

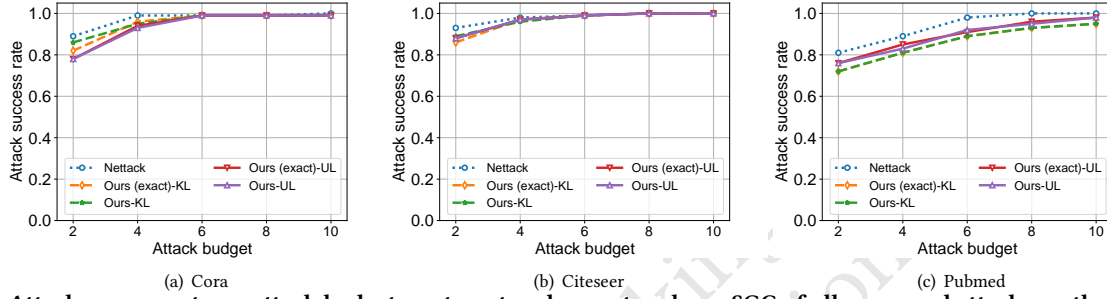


Figure 2: Attack success rate vs. attack budget per target node on a two-layer SGC of all compared attacks on the three graphs.

Table 2: Fraction of label influences of target nodes' within  $K$ -hop neighbors that have the same label as the target node.

Dataset	Cora			Citeseer			Pubmed		
K	2	3	4	2	3	4	2	3	4
Two-layer GCN	86%	83%	80%	89%	88%	87%	84%	81%	80%
Two-layer SGC	86%	82%	80%	92%	91%	90%	88%	84%	82%

implemented based on exact label influence calculation, and denote the corresponding two methods with known and unknowns labels as **Ours (exact)-KL** and **Ours (exact)-UL**, respectively.

Given a target GNN model, a set of target nodes, target label, and an attack budget  $\Delta$ , attack success rate is the fraction of target nodes that are misclassified by the target GNN to be the target label when the number of attack edges per target node is at most  $\Delta$ . Running time is reported on average across all the target nodes.

We train all GNNs using the public source code. We test Netack using the source code (<https://github.com/danielzuegner/netack>). We implement our attack in PyTorch. All experiments are conducted on an A6000 GPU with 48G memory.

## 5.2 Results on Attacking Two-layer GNNs

In this experiment, we aim to validate Assumption 1. Specifically, in the three datasets, for each of the 100 target nodes, we summarize the label influences of all its  $K$ -hop neighbors and select  $K = 2, 3, 4$  in our experiments (The shortest paths in these datasets are  $\leq 4$ ). The results on two-layer GCN/SGC are shown in Table 2. We observe that, in all datasets, a majority, i.e., more than 80%, of the total influences are from the nodes that have the same label as the target node. Such results validate Assumption 1 holds in general.

Next, we compare our attacks with Netack in terms of effectiveness (i.e., attack success rate) and efficiency (i.e., running time) against two-layer GCN/SGC. Figures 1 and 2 show the attack success rate against GCN and SGC on the three graphs, respectively. Moreover, Figures 3 and 4 show the running time of all attacks

against GCN and SGC on the three graphs, respectively. We have the following key observations.

- *Our attacks based on approximate label influence have similar performance with those based on exact label influence, but is much more efficient.* Specifically, the difference of the attack success rate between the two is less than 2% in all cases. This shows that our proposed efficient algorithm for label influence calculation is effective enough. Moreover, our attacks based on approximate label influence are 1–2 orders of magnitude more efficient than those based on exact label influence.
- *Our attacks with true labels and with estimated labels have similar performance.* Specifically, the difference of the attack success rate between Ours-KL and Ours-UL is negligible, i.e., less than 2% in all cases, and the running time of both Ours-KL and Ours-UL are almost the same. One reason is that the trained GNN model has accurate predictions on the unlabeled nodes, and thus most of the estimated labels match the true labels. One should note that Ours-UL knows very limited knowledge about the GNN model and thus it is a very practical attack.
- *Our attacks achieve comparable performance with Netack.* Netack achieves state-of-the-art attack performance against two-layer GCN. Our attacks have a slightly lower attack success rate than Netack when the attack budget is small, e.g., less than 4. This is possibly because our attack use some approximations on Assumption 1, and when the attack space is small, Assumption 1 negatively affects the attack effectiveness to some extent. However, when the attack budget is larger than 4, our attacks obtain almost the same performance with Netack.
- *Our attacks are much more efficient than Netack.* Specifically, our attacks have a 5–50x speedup over Netack across the three graphs. As the attack budget increases (from 2 to 10) or the graph size increases (from Cora to Pubmed), our attacks achieve better efficiencies. The reasons are two-fold. First, Netack needs

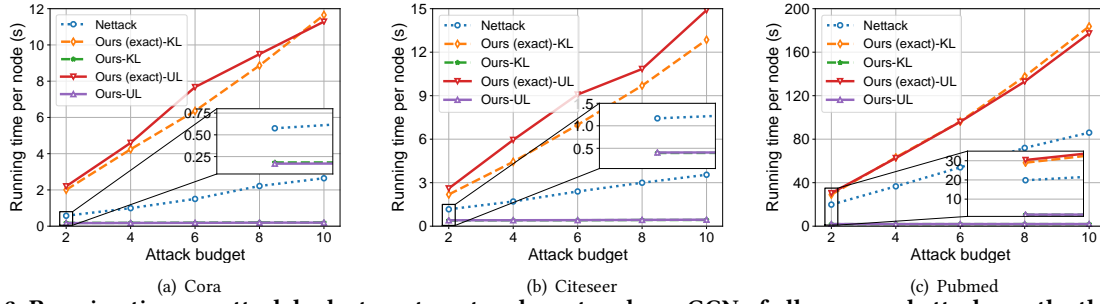


Figure 3: Running time vs. attack budget per target node on two-layer GCN of all compared attacks on the three graphs.

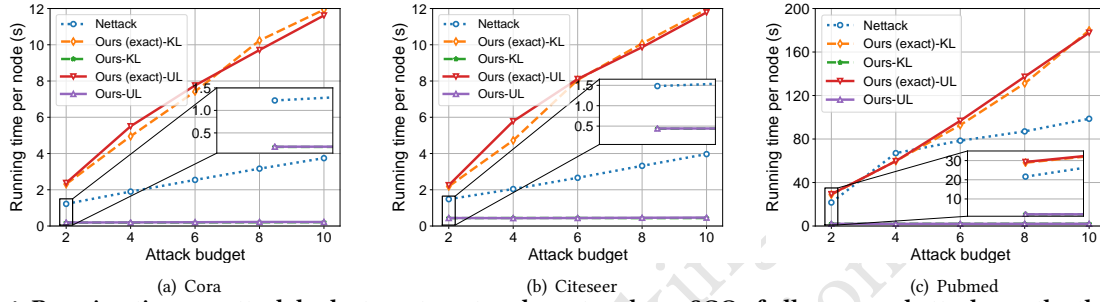


Figure 4: Running time vs. attack budget per target node on two-layer SGC of all compared attacks on the three graphs.

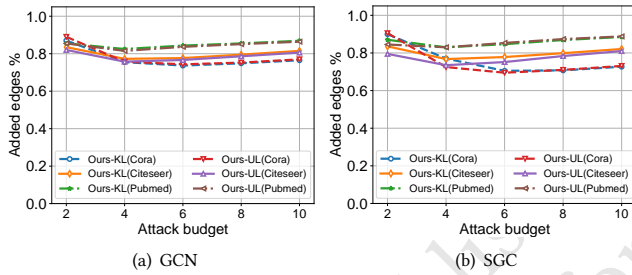


Figure 5: Fraction of the added edges among all attack edges generated by our attacks against (a) two-layer GCN and (b) two-layer SGC vs. attack budget per target node.

to multiply GNN model parameters in different layers, while our attacks do not. Second, Netstack involves multiplying the node hidden features, while ours is performed by calculating the label influence. Node hidden features are often high-dimensional, while label influence only needs scalar edge weights products.

We conduct experiments on the large-scale dataset OGB-arxiv to demonstrate the superior efficiency of our proposed attack method compared to baselines. We use a 2-layer GCN/SGC as the target GNN, achieving a clean accuracy of approximately 60% on test nodes. The attack method is set as Ours-KL. Note that time is denoted as the average running time of compared attack methods across the five attack budgets. We assess the attack performance on 100 target nodes that are accurately classified by GCN/SGC. The comparison results between Netstack and our attack are presented in Table 5. We observe that 1) Netstack encounters an out-of-memory (OOM) error on our platform due to the need for storing dense model weights and involving intensive matrix-matrix multiplication. 2) Our attacks achieve highly promising attack success rates while maintaining efficiency. Specifically, with an attack budget of 6, the attack success rates of Ours-KL against GCN and SGC are

73% and 95%, respectively. These results demonstrate the significant advantages of our attack method over baselines on large graphs.

We further analyze the properties of the attack edges. Figure 5(a) and Figure 5(b) show the fraction of the added edges generated by our attacks against two-layer GCN and two-layer SGC, respectively. We have two key observations. First, Ours-KL and Ours-UL generate almost the same fraction of added edges in all attack budgets and all graphs. This again verifies the similar characteristics between Ours-KL and Ours-UL. Second, the fraction of added edges is larger than 0.5 in all cases. This indicates that when performing the targeted attack, adding new edges between the target node and the nodes with the target label could be more effective than removing existing edges between the target node and the nodes having the same label as the target node.

We consider the following three factors, i.e., *node degree*, *node centrality*, and *graph size*, that could affect the target node’s attack performance. Here, we adopt the normalized closeness centrality (NCC) as the metric to measure node centrality. Specifically, the NCC of a node is the average length of the shortest path between the node and all other nodes in the graph. We have the following conclusions: 1) *Nodes with smaller degrees are easier to attack.* Given a fixed attack budget (e.g., 4 in our experiment), we observe that in all the three datasets, 100% of the target nodes with degree  $\leq 4$  attack successfully, while at most 83% and 85% of the target nodes with degree  $> 4$  attack successfully against 2-layer GCN and 2-layer SGC, respectively. 2) *Nodes with larger centrality are easier to attack.* We assume the attack budget is 4 per node. Specifically, 100% and 98% of the 50 target nodes with the largest NCC successfully attack 2-layer GCN and 2-layer SGC in the three datasets, while  $< 90\%$  and  $< 80\%$  of the 50 target nodes with the smallest NCC can successfully perform the attack. 3) *No obvious relationship between graph size and attack success rate.* Specifically, graph size: Pubmed  $>$  Citeseer  $>$  Cora. When attacking 2-layer GCN and the attack budget is 6, we



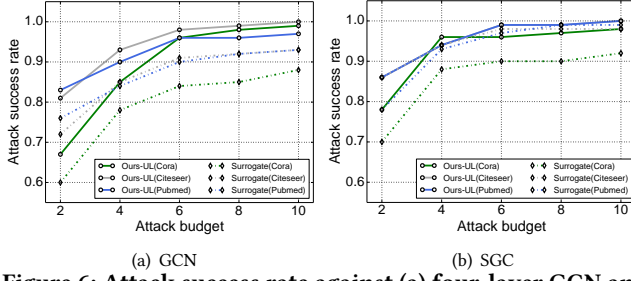


Figure 6: Attack success rate against (a) four-layer GCN and (b) four-layer SGC on the three graphs.

Table 3: Transferability of our attacks against two-layer GCN to other GNNs. Attack budget per target node is 6.

Dataset	Source	Target			
Cora	GCN	GCN	SGC	GAT	JK-Net
	No attack	0	0.01	0.03	0.02
	Ours-KL	0.98	0.82	0.66	0.67
	Ours-UL	0.98	0.84	0.65	0.70
Citeseer	GCN	GCN	SGC	GAT	JK-Net
	No attack	0	0.01	0.01	0.03
	Ours-KL	0.96	0.78	0.70	0.63
	Ours-UL	0.96	0.78	0.72	0.63
Pubmed	GCN	GCN	SGC	GAT	JK-Net
	No attack	0	0.03	0.04	0.05
	Ours-KL	0.89	0.80	0.80	0.79
	Ours-UL	0.88	0.80	0.80	0.77

have the attack success rate: Cora (0.98) > Citeseer (0.96) > Pubmed (0.87). However, when the attack budget is 10, we have the attack success rate: Pubmed (1.00) > Citeseer (0.99) > Cora (0.98).

### 5.3 Results on Attacking Multi-layer GNNs

In this experiment, we evaluate our attacks against multi-layer GCN/SGC. We denote Nettack that attacks a surrogate two-layer GNN model first and then transfers to attacking the target model as **Surrogate**. Figure 6 shows the attack success rate of our attack vs. attack budget against four-layer GCN and four-layer SGC on the three graphs, respectively. First, similarly, our attacks with both known label and unknown label are effective and achieve close attack performance, and we thus show results with unknown label for simplicity. When the attack budget is 6, our attacks achieve an attack success rate of  $\geq 90\%$  in all cases. Second, our attacks are more effective than the *indirect* surrogate model based attacks. Specifically, our attacks have more than 10% higher attack success rate than the surrogate model based attacks in almost all cases.

Moreover, Figure 7 shows the running time of our attacks vs. attack budget against four-layer GCN and four-layer SGC on the three graphs, respectively. Our attack is efficient. For instance, it takes our attacks less than 25s on average to attack a target node on the largest Pubmed in all cases. However, the surrogate model costs about 85s, validating that our attack is much more efficient.

In this experiment, we study the transferability of our attacks, i.e., whether the attack edges generated by our attacks against GCN/SGC can be also effective for other GNNs. Specifically, we use our attacks to generate the attack edges for each target node by attacking the source GNN (GCN or SGC), change the graph structure based on the attack edges, and adopt a target GNN to

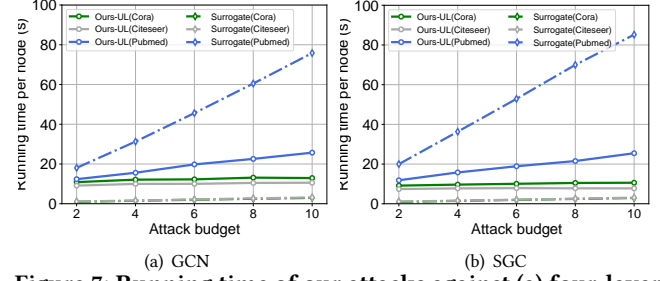


Figure 7: Running time of our attacks against (a) four-layer GCN and (b) four-layer SGC on the three graphs.

Table 4: Transferability of our attacks against two-layer SGC to other GNNs on the three graphs. Attack budget is 6.

Dataset	Source	Target			
Cora	SGC	SGC	GCN	GAT	JK-Net
	No attack	0	0.01	0.01	0.02
	Ours-KL	0.99	0.78	0.73	0.71
	Ours-UL	0.99	0.75	0.72	0.74
Citeseer	SGC	SGC	GCN	GAT	JK-Net
	No attack	0	0.03	0.03	0.04
	Ours-KL	0.99	0.81	0.79	0.69
	Ours-UL	0.99	0.84	0.77	0.70
Pubmed	SGC	SGC	GCN	GAT	JK-Net
	No attack	0	0.04	0.07	0.07
	Ours-KL	0.89	0.83	0.76	0.81
	Ours-UL	0.92	0.83	0.75	0.76

classify each target node on the perturbed graph. We select two additional representative GNNs, i.e., GAT [35] and JK-Net [44], as the target GNN. If a target node is also misclassified by the target GNN to be the target label, we say the attack edges generated by the source GNN are transferable.

Table 3 and Table 4 show the attack success rate of transferring of our attacks against two-layer GCN and two-layer SGC to attack other GNNs on the three graphs, where the attack budget per target node is 6. Note that we also show the attack performance for target GNNs without attack, i.e., the prediction error of target GNNs on the target nodes in the clean graph. We have the following observations. First, our attacks against GCN (or SGC) have the best transferability to SGC (or GCN). This is because SGC is a special case of GCN and they share similar model architectures. Second, our attacks are also effective against GAT and JK-Net. Specifically, on all the three graphs, our attacks can increase the classification errors by at least 60% when attacking GAT and JK-Net. This indicates that all the attack edges generated by our attack on the source GNN can be transferred to attack the target GNNs. Such good transferability further demonstrates the advantages of using (label) influence to perform the target evasion attacks.

## 6 DISCUSSIONS

**Evaluations on a larger dataset OGB-arXiv.** We follow the above attack setup and test our attack on a larger dataset OGB-arXiv [16]. The results against two-layer GCN/SGC are shown in Table 5. We can see that our attack is still very effective.

**Comparing with more attack baselines.** In the paper, we mainly choose the most state-of-the-art attack method, Nettack [52], as the attack baseline. The reasons are that: 1) Nettack and other evasion



**Table 5: Attack results of Ours-KL on OGB-arxiv.**

Attack budget	2	4	6	8	10	Time
Nettack	OOM	OOM	OOM	OOM	OOM	-
Ours-GCN	0.65	0.72	0.73	0.73	0.82	40.1s
Ours-SGC	0.90	0.93	0.95	0.95	0.96	40.7s

**Table 6: Comparing our attack vs. black-box attack [38].**

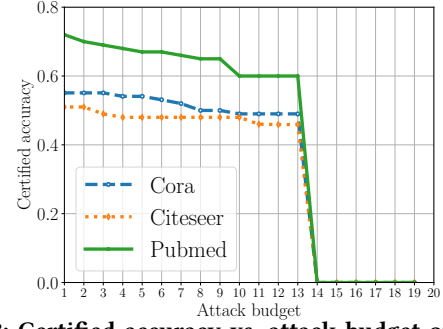
GCN	budget	2	4	6	GCN	budget	2	4	6
Cora	Ours	0.81	0.94	0.98	Citeseer	Ours	0.86	0.94	0.95
	[38]	0.68	0.79	0.80		[38]	0.84	0.91	0.92
SGC	budget	2	4	6	SGC	budget	2	4	6
Cora	Ours	0.81	0.92	0.98	Citeseer	Ours	0.90	0.96	0.98
	[38]	0.64	0.77	0.81		[38]	0.80	0.88	0.90

**Table 7: Comparing our attack vs. IG-FGSM [40] on Cora.**

Model	budget	2	4	6	8	10	Time
GCN	IG-FGSM	0.29	0.75	0.89	0.92	0.94	62s
	Nettack	0.85	0.96	0.97	0.97	0.97	1.5s
	Ours	0.90	0.93	0.95	0.95	0.96	0.1s
Model	budget	2	4	6	8	10	Time
SGC	IG-FGSM	0.32	0.74	0.89	0.95	0.96	55s
	Nettack	0.65	0.72	0.73	0.73	0.82	2.5s
	Ours	0.90	0.93	0.95	0.95	0.96	0.1s

attacks [5, 7, 40, 42] to GNNs have comparable performance. 2) Nettack explicitly keeps the graph structure property during the attack, making it difficult to defend against. To further demonstrate that our attack method is much more efficient than existing baselines, we choose to compare Ours-KL with a more recent attack IG-FGSM [40], where we use the same setting as that in Sec. 5.2 (i.e., 2-layer GCN/SGC as the target GNN, 100 target nodes). The comparison results on Cora are reported in Table 7. Note that time is denoted as the average running time of compared attack methods across the five attack budgets. We observe that: 1) Nettack not only outperforms IG-FGSM when the attack budget is small, but also is far more efficient than IG-FGSM. Specifically, when the attack budget is 2, the attack success rates of IG-FGSM against GCN and SGC are only 0.29 and 0.32, respectively, which are significantly lower than that of Nettack, 0.85 and 0.65. Moreover, the running times of IG-FGSM against GCN and SGC are 62s and 55s, respectively, which are much lower than that of Nettack, 1.5s and 2.5s, respectively. 2) Our method is even more efficient than Nettack. Specifically, when attacking GCN and SGC, the running time of our method is only 0.1s, which is lower than other two methods, further validating the efficiency of our method.

**Comparing with black-box attacks.** In the paper, we mainly compare our attack with the strongest white-box attack. Note that our attack is restricted black-box, as it does not know any information about the internal GNN model, but the target node’s neighbors’ status. Here we also compare our attack with the stringent black-box attacks proposed in [38]. To best explore the attack capability, we do not restrict the number of queries in [38], and obtain the optimal attack successful rate for a given attack budget. Table 6 shows the comparison results on Cora and Citeseer (Note that [38] cannot run on Pubmed due to limited GPU memory) on attacking 2-layer GCN/SGC. We can see that our attack is more effective than [38], especially when the attack budget is small. One key reason is

**Figure 8: Certified accuracy vs. attack budget of GCN on three graphs. We use default parameters in [37], e.g., noise parameter is 0.7, confidence level is 99.9%, and number of samples is 100,000.**

that our attack utilizes the strong connection between GNN and label propagation, while [38] performs the attack based on the query feedback, i.e., the target node’s confidence score after querying the black-box GNN model.

**Defending against our attack.** As shown in Section 2, existing empirical defenses [10, 12, 18, 40, 46, 49, 51] are easy to be broken [26] when the adversary knows the defense mechanism. Hence, we propose to defend our attack via provable defenses and choose the state-of-the-art randomized smoothing-based provable defense [37]. Specifically, given a target node, a model is provably robust for the target node if the model correctly predicts the same label for the target node when the attacker *arbitrarily* modifies a bounded number of (e.g., at most  $R$ ) edges in the graph, where  $R$  is called *certified radius*. Hence, provably robust models can defend against the worst-case attack (including our attack). Accordingly, certified accuracy under  $R$  means the fraction of target nodes that are predicted accurately by modifying any  $R$  edges. That is, if a model achieves a larger certified accuracy at a given budget, it shows better provable robustness. We conduct experiments on defending two-layer GCN against the worst-cast attack via [37]. Results on the three datasets are shown in Figure 8. We can see that, when any 4 edges are allowed to be modified, the certified accuracy achieved by the method [37] on the three datasets are about 0.50, 0.55, and 0.70, respectively. However, the method [37] cannot provably defend against the worst-case attack when the attack budget is larger than 13, which implies the need to design more powerful provable defenses.

## 7 CONCLUSION

We propose an influence-based evasion attack against GNNs. Specifically, we first build the connection between GNNs and label propagation (LP) via carefully designed influence functions. Next, we reformulate the attack against GNNs to be related to label influence on LP. Then, we design an efficient algorithm to calculate label influences. Our attack is applicable to multi-layer GNNs and does not need to know the GNN model parameters. We evaluate our attack on multiple benchmark graph datasets. Experimental results demonstrate that our attack achieves comparable performance against state-of-the-art white-box attack, and has a 5-50x speedup when attacking two-layer GCNs. Our attack is also effective to attack multi-layer GNNs and is transferable to other GNNs. Finally, our attack is more effective than the state-of-the-art black-box attack.

## REFERENCES

- [1] Aleksandar Bojchevski and Stephan Günnemann. 2019. Adversarial Attacks on Node Embeddings via Graph Poisoning. In *ICML*.
- [2] Aleksandar Bojchevski and Stephan Günnemann. 2019. Certifiable Robustness to Graph Perturbations. In *NeurIPS*.
- [3] Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. 2020. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In *ICML*.
- [4] Heng Chang, Yu Rong, Tingyang Xu, Wenbing Huang, Honglei Zhang, Peng Cui, Wenwu Zhu, and Junzhou Huang. 2020. A Restricted Black-Box Adversarial Framework Towards Attacking Graph Embedding Models.. In *AAAI*.
- [5] Jinyin Chen, Yangyang Wu, Xuanheng Xu, Yixian Chen, Haibin Zheng, and Qi Xuan. 2018. Fast gradient attack on network embedding. *arXiv* (2018).
- [6] Yizheng Chen, Yacin Nadij, Athanasios Kountouras, Fabian Monrose, Roberto Perdisci, Manos Antonakakis, and Nikolaos Vasiloglou. 2017. Practical attacks against graph-based clustering. In *CCS*.
- [7] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. 2018. Adversarial attack on graph structured data. In *ICML*.
- [8] Quanyu Dai, Qiang Li, Jian Tang, and Dan Wang. 2018. Adversarial network embedding. In *AAAI*.
- [9] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *NIPS*.
- [10] Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. 2020. All you need is low (rank) defending against adversarial attacks on graphs. In *WSDM*.
- [11] Houxiang Fan, Binghui Wang, Pan Zhou, Ang Li, Zichuan Xu, Cai Fu, Hai Li, and Yiran Chen. 2021. Reinforcement learning-based black-box evasion attacks to link prediction in dynamic graphs. In *2021 IEEE 23rd Int Conf on High Performance Computing & Communications*. 933–940.
- [12] Simon Geisler, Tobias Schmidt, Hakan Şirin, Daniel Zügner, Aleksandar Bojchevski, and Stephan Günnemann. 2021. Robustness of graph neural networks at scale. *Advances in Neural Information Processing Systems* 34, 7637–7649.
- [13] Simon Geisler, Daniel Zügner, and Stephan Günnemann. 2020. Reliable graph neural networks via robust aggregation. In *Advances in Neural Information Processing Systems*.
- [14] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. In *ICML*.
- [15] Will Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NIPS*.
- [16] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems* 33, 22118–22133.
- [17] Hongwei Jin, Zhan Shi, Venkata Jaya Shankar Ashish Peruri, and Xinhua Zhang. 2020. Certified Robustness of Graph Convolution Networks for Graph Classification under Topological Attacks. In *NeurIPS*.
- [18] Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. 2020. Graph structure learning for robust graph neural networks. In *KDD*.
- [19] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR*.
- [20] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *ICML*.
- [21] Jia Li, Honglei Zhang, Zhichao Han, Yu Rong, Hong Cheng, and Junzhou Huang. 2020. Adversarial attack on community detection by hiding individuals. In *WWW*.
- [22] Xuanqing Liu, Si Si, Xiaojin Zhu, Yang Li, and Cho-Jui Hsieh. 2019. A unified framework for data poisoning attack to graph-based semi-supervised learning. *arXiv preprint arXiv:1910.14147* (2019).
- [23] Jiaqi Ma, Shuangrui Ding, and Qiaozhu Mei. 2020. Towards More Practical Adversarial Attacks on Graph Neural Networks. In *Advances in Neural Information Processing Systems*.
- [24] Yao Ma, Suhang Wang, Tyler Derr, Lingfei Wu, and Jiliang Tang. 2019. Attacking graph convolutional networks via rewiring. *arXiv preprint arXiv:1906.03750* (2019).
- [25] Jiaming Mu, Binghui Wang, Qi Li, Kun Sun, Mingwei Xu, and Zhuotao Liu. 2021. A Hard Label Black-box Adversarial Attack Against Graph Neural Networks. In *CCS*.
- [26] Felix Mujkanovic, Simon Geisler, Stephan Günnemann, and Aleksandar Bojchevski. 2022. Are Defenses for Graph Neural Networks Robust?. In *NeurIPS*.
- [27] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE Transactions on Neural Networks* (2008).
- [28] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI magazine* 29, 3 (2008), 93–93.
- [29] Mingjie Sun, Jian Tang, Huichen Li, Bo Li, Chaowei Xiao, Yao Chen, and Dawn Song. 2018. Data poisoning attack against unsupervised node embedding methods. *arXiv* (2018).
- [30] Yiwei Sun, Suhang Wang, Xianfeng Tang, Tsung-Yu Hsieh, and Vasant Honavar. 2020. Adversarial Attacks on Graph Neural Networks via Node Injections: A Hierarchical Reinforcement Learning Approach. In *The Web Conference*.
- [31] Tsubasa Takahashi. 2019. Indirect Adversarial Attacks via Poisoning Neighbors for Graph Convolutional Networks. In *2019 IEEE International Conference on Big Data (Big Data)*.
- [32] Xianfeng Tang, Yandong Li, Yiwei Sun, Huaxiu Yao, Prasenjit Mitra, and Suhang Wang. 2020. Transferring Robustness for Graph Neural Network Against Poisoning Attacks. In *WSDM*.
- [33] Shuchang Tao, Huawei Shen, Qi Cao, Liang Hou, and Xueqi Cheng. 2021. Adversarial Immunization for Certifiable Robustness on Graphs. In *WSDM*.
- [34] MohamadAli Torkamani and Daniel Lowd. 2013. Convex adversarial collective classification. In *ICML*.
- [35] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. In *ICLR*.
- [36] Binghui Wang and Neil Zhenqiang Gong. 2019. Attacking Graph-based Classification via Manipulating the Graph Structure. In *CCS*.
- [37] Binghui Wang, Jinyuan Jia, Xiaoyu Cao, and Neil Gong. 2021. Certified robustness of graph neural networks against adversarial structural perturbation. In *KDD*.
- [38] Binghui Wang, Youqi Li, and Pan Zhou. 2022. Bandits for Structure Perturbation-based Black-box Attacks to Graph Neural Networks with Theoretical Guarantees. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13379–13387.
- [39] Felix Wu, Tianyi Zhang, Amauri Holanda de Souza Jr, Christopher Fifty, Tao Yu, and Kilian Q Weinberger. 2019. Simplifying graph convolutional networks. In *ICML*.
- [40] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. 2019. Adversarial examples on graph data: Deep insights into attack and defense. In *IJCAI*.
- [41] Zhaohan Xi, Ren Pang, Shouling Ji, and Ting Wang. 2021. Graph backdoor. In *(USENIX) Security 21*.
- [42] Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin. 2019. Topology attack and defense for graph neural networks: An optimization perspective. In *IJCAI*.
- [43] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How powerful are graph neural networks?. In *ICLR*.
- [44] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. In *ICML*.
- [45] Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In *Advances in Neural Information Processing Systems*. 5165–5175.
- [46] Xiang Zhang and Marinka Zitnik. 2020. GnnGuard: Defending graph neural networks against adversarial attacks. In *NeurIPS*.
- [47] Zaixi Zhang, Jinyuan Jia, Binghui Wang, and Neil Zhenqiang Gong. 2021. Backdoor attacks to graph neural networks. (2021).
- [48] Zijie Zhang, Zeru Zhang, Yang Zhou, Yelong Shen, Ruoming Jin, and Dejing Dou. 2020. Adversarial Attacks on Deep Graph Matching. *Advances in Neural Information Processing Systems* 33.
- [49] Dingyuan Zhu, Ziwei Zhang, Peng Cui, and Wenwu Zhu. 2019. Robust graph convolutional networks against adversarial attacks. In *KDD*.
- [50] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. 2003. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*.
- [51] Jun Zhuang and Mohammad Al Hasan. 2022. Defending Graph Convolutional Networks against Dynamic Graph Perturbations via Bayesian Self-supervision. In *AAAI*.
- [52] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. 2018. Adversarial attacks on neural networks for graph data. In *KDD*.
- [53] Daniel Zügner and Stephan Günnemann. 2019. Adversarial attacks on graph neural networks via meta learning. *ICLR* (2019).
- [54] Daniel Zügner and Stephan Günnemann. 2020. Certifiable Robustness of Graph Convolutional Networks under Structure Perturbations. In *KDD*.

## PROOFS

### Proof of Theorem 4.1

We mainly focus on GCN, as SGC is a special case of GCN and the proof also applies. Our proof is based on the Lemma 1 in Section 4.2.

Based on Equation (8) in Lemma 1, the feature-label influence in Equation (5) can be expressed as follows:

$$\begin{aligned} I_{fl}(v, u; K) &= \mathbf{1}_{yu}^T \left[ \frac{\partial \mathbf{h}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \right] \mathbf{h}_u^{(0)} = \mathbf{1}_{yu}^T \left[ \prod_{l=K}^1 \mathbf{W}^{(l)} \cdot \rho \cdot \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 a_{v_p^l, v_p^{l-1}} \right] \mathbf{h}_u^{(0)} \\ &= \rho \cdot \mathbf{1}_{yu}^T \left[ \prod_{l=K}^1 \mathbf{W}^{(l)} \right] \mathbf{h}_u^{(0)} \cdot \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 a_{v_p^l, v_p^{l-1}} = C \cdot \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 a_{v_p^l, v_p^{l-1}}, \end{aligned} \quad (13)$$

where  $C = \rho \mathbf{1}_{yu}^T \left[ \prod_{l=1}^K \mathbf{W}^{(l)} \right] \mathbf{h}_u^{(0)}$  is a constant for a given GCN.

Comparing GCN with LP, we can find their iteration processes are similar, except that LP has no model parameters (which is constant for a trained GCN model). Specifically, we can calculate the label influence  $I_l$  as follows:

$$I_l(v, u; K) = \frac{\partial y_u^{(K)}}{\partial y_u^{(0)}} = \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 a_{v_p^l, v_p^{l-1}}. \quad (14)$$

Hence, we can build the following relationship between the feature-label influence in GNN and label influence in LP:

$$I_{fl}(v, u; K) = C \cdot I_l(v, u; K). \quad (15)$$

### Proof of Theorem 4.2

Let  $\tilde{\mathbf{h}}_v^{(K)}$  be  $u$ 's final node representation after the attack. Substituting Equation (??), the attack's objective function in Equation (4) is equivalent to the following form:

$$\begin{aligned} \max_{\tilde{\Lambda}_v} & \left( \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c - \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_{y_v} \right) \\ \text{s.t., } & \sum_s |\tilde{A}_{v,s} - A_{v,s}| \leq \Delta. \end{aligned}$$

Based on Assumption 1, we further deduce attack's objective function in Equation (9) as follows:

$$\begin{aligned} & \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c - \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_{y_v} \\ & \approx \left[ \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c - \left[ \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_z^{(0)}} \cdot \mathbf{h}_z^{(0)} \right]_{y_v} \\ & = \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \left[ \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]_c - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \left[ \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_z^{(0)}} \cdot \mathbf{h}_z^{(0)} \right]_{y_v} \\ & = \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \left[ \mathbf{1}_{yu}^T \cdot \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right] - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \left[ \mathbf{1}_{yz}^T \cdot \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_z^{(0)}} \cdot \mathbf{h}_z^{(0)} \right]. \end{aligned}$$

Note that the two terms  $\left[ \mathbf{1}_{yu}^T \cdot \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_u^{(0)}} \cdot \mathbf{h}_u^{(0)} \right]$  and  $\left[ \mathbf{1}_{yz}^T \cdot \frac{\partial \tilde{\mathbf{h}}_v^{(K)}}{\partial \mathbf{h}_z^{(0)}} \cdot \mathbf{h}_z^{(0)} \right]$  are exactly the feature-label influence after the attack, and we denote them as  $\tilde{I}_{fl}(v, u; K)$  and  $\tilde{I}_{fl}(v, z; K)$ , respectively. Based on the relationship between the feature-label influence and label influence in Equation (15), we thus have the following attack's objective function in terms of label influence:

$$\begin{aligned} & \max_{\tilde{\Lambda}_v} \left( \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \tilde{I}_l(v, u; K) - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \right) \\ & \Leftrightarrow \max_{\tilde{\Lambda}_v} \left( \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \tilde{I}_l(v, u; K) - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \right). \end{aligned}$$

## EFFICIENT LABEL INFLUENCE CALCULATION

First, we have the following two observations:

- When adding an edge between  $v$  and  $a \in \mathcal{N}_A$ , we will have new paths  $\{v, a, \dots, u\}$  from  $v$  to  $u$  passing through  $a$ . We denote the nodes in the new paths within  $v$ 's  $K$ -hop neighbors as  $\Delta \Lambda_{v,a}^{(K)}$ , and note that these nodes are within  $a$ 's  $(K-1)$ -hop neighbors in the clean graph, i.e.,  $\Delta \Lambda_{v,a}^{(K)} \subset \Lambda_a^{(K-1)}$ . Moreover,  $\tilde{\Lambda}_v^{(K)} = \Lambda_v^{(K)} \cup \Delta \Lambda_{v,a}^{(K)}$ .
- When deleting an edge between  $v$  and  $b \in \mathcal{N}_B$ , we will remove existing paths  $\{v, b, \dots, u\}$  from  $v$  to  $u$  passing through  $b$ . We denote the deleted nodes within  $v$ 's  $K$ -hop neighbors as  $\Delta \Lambda_{v,b}^{(K)}$ , and these nodes are within  $b$ 's  $(K-1)$ -hop neighbors in the clean graph, i.e.,  $\Delta \Lambda_{v,b}^{(K)} \subset \Lambda_b^{(K-1)}$ . Moreover,  $\tilde{\Lambda}_v^{(K)} = \Lambda_v^{(K)} \setminus \Delta \Lambda_{v,b}^{(K)}$ .

Based on the two observations, we can split each label influence in Equation (11) into two parts: an *approximate constant label influence* defined on the clean graph and an *approximate label influence* defined on the  $(K-1)$ -hop neighbors for each node in  $\mathcal{N}_A \cup \mathcal{N}_B$ .

We first consider **adding (+)** an edge between node  $u$  and node  $a \in \mathcal{N}_A$ . Specifically, we have:

$$\begin{aligned} & \sum_{u \in \tilde{\Lambda}_v^{(K)}, y_u=c} \tilde{I}_l(v, u; K) - \sum_{z \in \tilde{\Lambda}_v^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \\ & = \left[ \sum_{u \in \Lambda_v^{(K)}, y_u=c} \tilde{I}_l(v, u; K) + \sum_{u \in \Delta \Lambda_{v,a}^{(K)}, y_u=c} \tilde{I}_l(v, u; K) \right] \\ & \quad - \left[ \sum_{z \in \Lambda_v^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) + \sum_{z \in \Delta \Lambda_{v,a}^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \right] \\ & = \left[ \sum_{u \in \Lambda_v^{(K)}, y_u=c} \tilde{I}_l(v, u; K) - \sum_{z \in \Lambda_v^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \right] \\ & \quad + \left[ \sum_{u \in \Delta \Lambda_{v,a}^{(K)}, y_u=c} \tilde{I}_l(v, u; K) - \sum_{z \in \Delta \Lambda_{v,a}^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K) \right] \\ & = \left[ \sum_{u \in \Lambda_v^{(K)}, y_u=c} \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 \tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}} - \sum_{z \in \Lambda_v^{(K)}, y_z=y_v} \sum_{p=1}^{\Psi_{v \rightarrow z}} \prod_{l=K}^1 \tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}} \right] \\ & \quad + \left[ \sum_{u \in \Delta \Lambda_{v,a}^{(K)}, y_u=c} \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 \tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}} - \sum_{z \in \Delta \Lambda_{v,a}^{(K)}, y_z=y_v} \sum_{p=1}^{\Psi_{v \rightarrow z}} \prod_{l=K}^1 \tilde{d}_{v_p^l}^{-\frac{1}{2}} \tilde{d}_{v_p^{l-1}}^{-\frac{1}{2}} \right] \\ & \approx \left[ \sum_{u \in \Lambda_v^{(K)}, y_u=c} \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 d_{v_p^l}^{-\frac{1}{2}} d_{v_p^{l-1}}^{-\frac{1}{2}} - \sum_{z \in \Lambda_v^{(K)}, y_z=y_v} \sum_{p=1}^{\Psi_{v \rightarrow z}} \prod_{l=K}^1 d_{v_p^l}^{-\frac{1}{2}} d_{v_p^{l-1}}^{-\frac{1}{2}} \right] \\ & \quad + \left[ \sum_{u \in \Delta \Lambda_{v,a}^{(K)}, y_u=c} \sum_{p=1}^{\Psi_{v \rightarrow u}} \prod_{l=K}^1 d_{v_p^l}^{-\frac{1}{2}} d_{v_p^{l-1}}^{-\frac{1}{2}} - \sum_{z \in \Delta \Lambda_{v,a}^{(K)}, y_z=y_v} \sum_{p=1}^{\Psi_{v \rightarrow z}} \prod_{l=K}^1 d_{v_p^l}^{-\frac{1}{2}} d_{v_p^{l-1}}^{-\frac{1}{2}} \right] \\ & = C_A + \Delta I_A(a), \end{aligned}$$

where

- $\sum_{u \in \Delta \Lambda_{v,a}^{(K)}, y_u=c} \tilde{I}_l(v, u; K)$ ,  $\sum_{z \in \Delta \Lambda_{v,a}^{(K)}, y_z=y_v} \tilde{I}_l(v, z; K)$  represent the label influence of the label- $c$  nodes and label- $y_v$  nodes within the  $(K-1)$ -hop neighbors of  $a$  after adding the edge  $(v, a)$ , respectively. Note that the degree of most nodes in these paths do not change, except that the degree of the target node  $v$  and the node  $a$  increases by 1 due to the added edge. Here, for efficient computation, we assume all nodes in  $\Lambda_v^{(K)}$ , except the target node  $v$ , do not change the node degree. That means, we approximately set  $\tilde{d}_{v_p^l} = d_{v_p^l}$  and  $\tilde{d}_{v_p^{l-1}} = d_{v_p^{l-1}}$ , for  $v_p^l \neq v$  and  $v_p^{l-1} \neq v$ .
- $C_A$  is the approximate constant label influence defined on the *clean graph* and only needs to be calculated once.



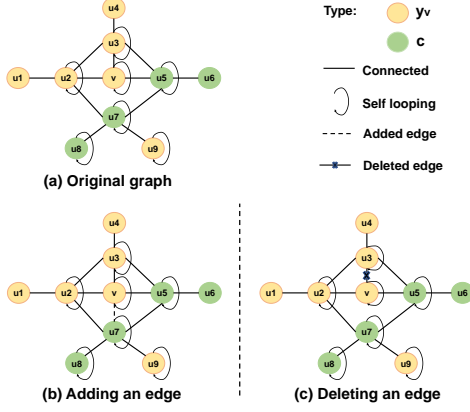


Figure 9: An example for calculating label influence.

•  $\Delta I_A(a)$  is the approximate label influence caused by the perturbed edge  $(v, a)$ ,  $\forall a \in N_A$ . We can calculate  $\Delta I_A(a)$  for each attack edge  $(v, a)$  in advance.

When **removing** an edge between nodes  $u$  and  $b \in N_B$ , we replace  $+$  with  $-$  and assume all nodes' degree do not change, except the target node's degree decreases by 1. We denote the approximate constant label influence on the clean graph as  $C_B$  and the label influence per perturbed edge  $(v, b)$ ,  $\forall b \in N_B$ , as  $\Delta I_B(b)$ , respectively. Similarly,  $C_B$  only needs to be calculated once and  $\Delta I_B(b)$  can be calculated for each attack edge  $(v, b)$  in advance.

## EXAMPLE OF CALCULATING INFLUENCE

Now we show an example for exact and approximate label influence calculation. We take Figure 9 as an example to illustrate how label influence terms are calculated when adding a new edge or deleting an existing edge when  $K = 2$ . The original graph is shown in Figure 9(a). There are two classes of nodes, i.e., green color and yellow color.  $v$  is the target node with a label- $y_v$ . The candidate nodes for  $v$  to add edges are  $N_A = \{u_6, u_7, u_8\}$ , and the candidate nodes for  $v$  to delete edges are  $N_B = \{u_2, u_3\}$  (with label yellow).

In the clean graph, we have  $\Lambda_v^{(2)} = [v, u_1, u_2, u_3, u_4, u_5, u_6, u_7]$ ;  $\{u \in \Lambda_v^{(2)}, y_u = c\} = \{u_5, u_6, u_7\}$ ;  $\{z \in \Lambda_v^{(2)}, y_z = y_v\} = \{u_1, u_2, u_3, u_4\}$ .

Then, we have all the paths

- (a)  $\{\Psi_{v \rightarrow u} | u \in \Lambda_v^{(2)}, y_u = c\} = \{v - u_5 - u_5, v - v - u_5, v - u_3 - u_5, v - u_5 - u_6, v - u_2 - u_7, v - u_5 - u_7\}$ ;
- (b)  $\{\Psi_{v \rightarrow z} | z \in \Lambda_v^{(2)}, y_z = y_v\} = \{v - v - v, v - u_5 - v, v - u_2 - v, v - u_3 - v, v - u_2 - u_1, v - u_2 - u_2, v - v - u_2, v - u_3 - u_2, v - u_3 - u_3, v - v - u_3, v - u_2 - u_3, v - u_5 - u_3, v - u_3 - u_4\}$

**Adding an edge.** We first compute the exact label influence and then calculate the approximate label influence.

Figure 9(b) is the graph structure after adding an edge between  $u_7$  and  $v$ . In the perturbed graph, we have:  $\tilde{\Lambda}_v^{(2)} = [v, u_1, u_2, u_3, u_4, u_5, u_6, u_7, u_8, u_9]$ ;  $\Delta \Lambda_{v, u_7}^{(2)} = \{v, u_2, u_5, u_7, u_8, u_9\}$ ;  $\{u \in \Delta \Lambda_{v, u_7}^{(2)}, y_u = c\} = \{u_5, u_7, u_8\}$ ;  $\{z \in \Delta \Lambda_{v, u_7}^{(2)}, y_z = y_v\} = \{v, u_2, u_9\}$ .

The new added paths are passing through  $u_7$ . Specifically, all the new added paths are:

- (c)  $\{\Psi'_{v \rightarrow u} | u \in \Delta \Lambda_{v, u_7}^{(2)}, y_u = c\} = \{v - v - u_7, v - u_7 - u_5, v - u_7 - u_7, v - u_7 - u_8\}$ ;
- (d)  $\{\Psi'_{v \rightarrow z} | z \in \Delta \Lambda_{v, u_7}^{(2)}, y_z = y_v\} = \{v - u_7 - v, v - u_7 - u_2, v - u_7 - u_9\}$ .

**Exact label influence calculation.** We first precisely calculate the label influence as follows:

$$\begin{aligned} & \sum_{u \in \tilde{\Lambda}_v^{(2)}, y_u = c} \tilde{I}_l(v, u; 2) - \sum_{z \in \tilde{\Lambda}_v^{(2)}, y_z = y_v} \tilde{I}_l(v, z; 2) \\ &= \left[ (d_{u_5}^{-\frac{3}{2}} d_v'^{-\frac{1}{2}} + d_{u_5}^{-\frac{1}{2}} d_v'^{-\frac{3}{2}} + d_{u_5}^{-\frac{1}{2}} d_{u_3}^{-1} d_v'^{-\frac{1}{2}} \right. \\ & \quad \left. + d_{u_6}^{-\frac{1}{2}} d_{u_5}^{-1} d_v'^{-\frac{1}{2}} + d_{u_7}^{-\frac{1}{2}} d_{u_2}^{-1} d_v'^{-\frac{1}{2}} + d_{u_7}^{-\frac{1}{2}} d_{u_5}^{-1} d_v'^{-\frac{1}{2}} \right) \quad \text{(a)} \\ & \quad - (d_v'^{-2} + d_v'^{-1} d_{u_5}^{-1} + d_v'^{-1} d_{u_2}^{-1} + d_v'^{-1} d_{u_3}^{-1} + d_{u_1}^{-\frac{1}{2}} d_{u_2}^{-1} d_v'^{-\frac{1}{2}} \\ & \quad + d_{u_2}^{-\frac{3}{2}} d_v'^{-\frac{1}{2}} + d_{u_2}^{-\frac{1}{2}} d_v'^{-\frac{3}{2}} + d_{u_2}^{-\frac{1}{2}} d_{u_3}^{-1} d_v'^{-\frac{1}{2}} \\ & \quad + d_{u_3}^{-\frac{3}{2}} d_v'^{-\frac{1}{2}} + d_{u_3}^{-\frac{1}{2}} d_v'^{-\frac{3}{2}} + d_{u_3}^{-\frac{1}{2}} d_{u_2}^{-1} d_v'^{-\frac{1}{2}} + d_{u_3}^{-\frac{1}{2}} d_{u_5}^{-1} d_v'^{-\frac{1}{2}} + d_{u_4}^{-\frac{1}{2}} d_{u_3}^{-1} d_v'^{-\frac{1}{2}}) \quad \text{(b)} \\ & \quad + \left[ (d_{u_7}^{-\frac{1}{2}} d_v'^{-\frac{3}{2}} + d_{u_5}^{-\frac{1}{2}} d_{u_7}^{-1} d_v'^{-\frac{1}{2}} + d_{u_7}^{-\frac{3}{2}} d_v'^{-\frac{1}{2}} + d_{u_8}^{-\frac{1}{2}} d_{u_7}^{-1} d_v'^{-\frac{1}{2}}) \quad \text{(c)} \right. \\ & \quad \left. - (d_v'^{-1} d_{u_7}^{-1} + d_{u_2}^{-1} d_{u_7}^{-1} d_v'^{-\frac{1}{2}} + d_{u_5}^{-\frac{1}{2}} d_{u_7}^{-1} d_v'^{-\frac{1}{2}}) \right] \quad \text{(d)} \\ &= C_A(\text{exact}) + I_A(u_7) = (0.2563 - 0.5665) + (0.1530 - 0.1194) \\ &= -0.3102 + 0.0336 = -0.2766. \end{aligned}$$

Note that  $C_A(\text{exact}) = -0.3102$  is the exact constant label influence in the original paths, and  $I_A(u_7) = 0.0336$  is the exact label influence in the new added paths.

**Approximate label influence calculation.** When calculating the label influence in the original paths in the above Equation after adding the edge between  $v$  and  $u_7$ , we notice that only paths including  $u_7$  are affected, as  $u_7$ 's degree increases by 1. Here, for efficient calculation, we ignore the effect caused by  $u_7$ 's degree, i.e., we set  $d_{u_7}' = d_{u_7}$ , which means that we do not need to normalize the weights of edges associated with  $u_7$ . For the target node  $v$ , its degree increases by 1 and we still use  $d_v' = d_v + 1$ .

With this setting, we calculate the approximate constant label influence of the first term  $\sum_{u \in \tilde{\Lambda}_v^{(2)}, y_u = c} \tilde{I}_l(v, u; 2)$  is  $C_A(\text{approx}) = -0.3032$ , which is close to the precise value  $-0.3102$ . Moreover, we compute the approximate label influence for the second term  $\sum_{z \in \tilde{\Lambda}_v^{(2)}, y_z = y_v} \tilde{I}_l(v, z; 2)$ , whose value is  $I_A(\text{approx}, u_7) = 0.0299$  and close to the precise value 0.0336. Both the results verify the effectiveness of our efficient calculation of label influence

**Deleting an edge.** Figure 9(c) is the graph after deleting an edge between  $u_3$  and  $v$ . In the perturbed graph, we have  $\tilde{\Lambda}_v^{(2)} = [v, u_1, u_2, u_3, u_4, u_5, u_6, u_7]$ ;  $\Delta \Lambda_{v, u_3}^{(2)} = \{v, u_2, u_3, u_4, u_5\}$ ;  $\{u \in \Delta \Lambda_{v, u_3}^{(2)}, y_u = c\} = \{u_5\}$ ;  $\{z \in \Delta \Lambda_{v, u_3}^{(2)}, y_z = y_v\} = \{v, u_2, u_3, u_4\}$ . The new deleted paths are:  $\{\Psi'_{v \rightarrow u} | u \in \Delta \Lambda_{v, u_3}^{(2)}, y_u = c\} = \{v - u_3 - u_5\}$ .  $\{\Psi'_{v \rightarrow z} | z \in \Delta \Lambda_{v, u_3}^{(2)}, y_z = y_v\} = \{v - u_3 - v, v - u_3 - u_2, v - v - u_3, v - u_3 - u_3, v - u_3 - u_4\}$ .

Similarly, when calculating the label influence in the original paths after removing the edge between  $v$  and  $u_3$ , we notice that only paths including  $u_3$  are affected, as  $u_3$ 's degree decreases by 1. Here, for efficient calculation, we ignore the effect caused by  $u_3$ 's degree, i.e., we set  $d_{u_3}' = d_{u_3}$ , which means that we do not need to normalize the weights of edges associated with  $u_3$ . For the target node  $v$ , its degree decreases by 1 and we still use  $d_v' = d_v - 1$ .

With this setting, we calculate that the exact label influence of the first term is  $C_B(\text{exact}) = -0.5423$  and its approximate value is  $C_B(\text{approx}) = -0.5304$ . Moreover, the exact label influence of the second term is  $I_B(u_3) = -0.2860$ , and its approximate label influence is  $I_B(\text{approx}, u_3) = -0.2656$ . The two approximate values and precise values are also close.