

Scalable Attack on Graph Data by Injecting Vicious Nodes

Jihong Wang · Minnan Luo · Fnu Suya ·
Jundong Li · Zijiang Yang · Qinghua Zheng

Received: date / Accepted: date

Abstract Recent studies have shown that graph convolution networks (GCNs) are vulnerable to carefully designed attacks, which aim to cause misclassification of a specific node on the graph with unnoticeable perturbations. However, a vast majority of existing works cannot handle large-scale graphs because of their high time complexity. Additionally, existing works mainly focus on manipulating existing nodes on the graph, while in practice, attackers usually do not have the privilege to modify information of existing nodes. In this paper, we develop a more scalable framework named Approximate Fast Gradient Sign Method (AFGSM) which considers a more practical attack scenario where adversaries can only inject new vicious nodes to the graph while having no control over the original graph. Methodologically, we provide an approximation strategy to linearize the model we attack and then derive an approximate closed-form solution with a lower time cost. To have a fair comparison with existing attack methods that manipulate the original graph, we adapt them to the new attack scenario by injecting vicious nodes. Empirical experimental results show that our proposed attack method can significantly reduce the classification accuracy of GCNs and is much faster than existing methods without jeopardizing the attack performance.

Keywords Graph Convolution Networks · Vicious Nodes · Scalable Attack

1 Introduction

Graphs are widely used to model various types of real-world data and many canonical learning tasks such as classification, clustering, and anomaly detection have been widely investigated for the graph-structured data (Bhagat et al., 2011; Tian et al., 2014; Perozzi et al., 2014a; Tang et al., 2016). In this paper, we focus on the task of node classification. Recently, to solve the node classification problem,

Jihong Wang · Minnan Luo · Qinghua Zheng
Xian Jiaotong University, China
E-mail: wang1946456505@stu.xjtu.edu.cn, minnluo@xjtu.edu.cn, tjlu@mail.xjtu.edu.cn

Fnu Suya · Jundong Li Zijiang Yang
University of Virginia, USA Western Michigan University, USA
E-mail: {suya, jundong}@virginia.edu E-mail: zijiang.yang@wmich.edu

graph convolution networks (GCNs) have gained a surge of research interests in the data mining and machine learning community because of their superior prediction performance (Pham et al., 2017; Cai et al., 2018; Monti et al., 2017). However, recent research efforts showed that various graph mining algorithms (e.g., GCNs) are vulnerable to carefully crafted adversarial examples, which are “unnoticeable” to humans but can cause the learning models to misclassify some target nodes (Zügner et al., 2018; Dai et al., 2018; Bojchevski and Günnemann, 2018). The vulnerabilities of these learning algorithms can lead to severe consequences in security-sensitive applications. For example, GCNs are often used in the risk management area to evaluate the credit level of users (Dai et al., 2018; Akoglu et al., 2015; Bolton et al., 2001), as user-user information is often used in this context, it provides ample opportunities for criminals to increase their credit score by connecting to high-credit users. In this paper, we focus on assessing the robustness of graph convolution networks against adversarial attacks. Different from existing efforts that directly manipulate the original graph, we investigate a more realistic attack scenario of injecting vicious nodes and develop a scalable solution to tackle the problem.

Limitations of Current Approaches. There are two common issues of existing works on attacking GCNs (Zügner et al., 2018; Dai et al., 2018; Bojcheski and Günnemann, 2018): the scalability issue and the explicit assumption that adversaries can easily manipulate existing nodes on the graph. First, GCNs are usually applied in large-scale graphs in various domains (Ying et al., 2018; Hamilton et al., 2017; Chen et al., 2018), which puts a high demand on the scalability of the underlying attack models. However, existing efforts (Zügner et al., 2018; Dai et al., 2018; Bojcheski and Günnemann, 2018) cannot be easily generalized to handle large-scale graphs. Second, in actual life, attackers may not be able to manipulate existing nodes in a graph. For example, GCNs are often used to classify users on social websites like Twitter or Weibo for content recommendation by exploring the friendship graph of users. But attackers usually have no ability to manipulate existing users on these websites. To achieve the purpose of attack, a simple way is to register some new accounts on these websites and enable these accounts to establish connections with existing users, *e.g.*, following other users or making comments on the same posts. Despite that, existing attack models seldomly consider this new attack scenario. The aforementioned two limitations motivate us to investigate the following research problem: *how to effectively and efficiently manipulate the prediction results of GCNs on a specific node by injecting vicious nodes to a large graph?*

Challenges. There are three challenges in the new attack scenario, where the first two are general challenges for devising efficient attacks on GCNs while the third one is a unique challenge of the new attack scenario we consider.

- **Discreteness.** Different from images which can be approximately regarded as a continuous field as the value of pixels can be any integers between 0 and 255, graph-structured data is often depicted in the discrete domains. Thus existing attack (Dong et al., 2018; Szegedy et al., 2013) strategies that are widely used in other domains (*e.g.*, computer vision) cannot be directly applied. The proposed solution needs to well handle the potential combinatorial problem efficiently.
- **Poisoning Attack.** Unlike the clear separation between training and test data in other domains, the node classification task is often conducted in a

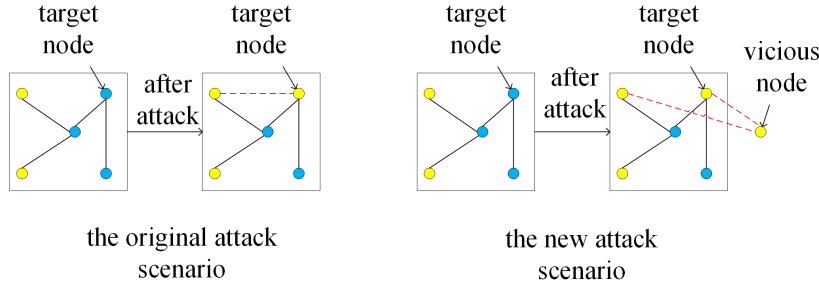


Fig. 1: Comparison between the attack scenario in existing literature and the new attack scenario considered in this paper.

transductive setting, where the test data (without ground-truth labels) is also considered in the training phase. Because of this, when the test data is manipulated, the graph is also dynamically updated. Therefore, it is important to propose attack strategies that remain effective after the model is retrained on the manipulated data. This leads to a bi-level optimization problem which is often computationally expensive to solve.

- **High complexity of Existing methods.** We can adapt existing methods to the new attack scenario. However, they all fail to scale to large graphs as their time complexities are very high (the complexity analysis is shown in Section 4.3). Considering that the node classification task is usually conducted on large graphs in practice, it is important to develop an efficient method that can be scaled to real-world large graphs.

Contributions. With the above-mentioned challenges, we propose a novel Approximate Fast Gradient Sign Method (AFGSM), which can modify and inject vicious nodes efficiently in the new attack setting. Specifically, our contributions can be summarized as follows:

- **A New Attack Scenario.** We consider a more practical attack scenario where adversaries can only inject vicious nodes to the graph while the original nodes on the graph remain unchanged.
- **Adapting Existing Attacks to New Scenario.** We adapt and carefully tune the existing attacks to our new attack scenario and adopt these attacks as the baselines for comparison.
- **A More Efficient and Effective Algorithm.** We propose a new attack strategy named Approximate Fast Gradient Sign Method (AFGSM), which can generate adversarial perturbations much more efficiently than the baseline attacks while maintaining similar attack performance.
- **Extensive Evaluation.** We empirically illustrate the effectiveness of our method on five benchmark datasets and also test on two state-of-the-art graph neural networks and one unsupervised network embedding model.

2 Related Work

Adversarial examples are extensively studied in the image classification task and recently researchers also show its existence in graph-related problems. Therefore, we will discuss related works in both the image domain and the graph domain.

Attack on Images. Szegedy et al. (2013) first demonstrate the vulnerability of deep learning models to adversarial examples using L-BFGS method and attributes the existence of adversarial examples to high non-linearity of deep models. Later on, Goodfellow et al. (2014) propose an efficient Fast Gradient Sign Method (FGSM) and instead demonstrates that adversarial examples exist because deep learning models are linear in nature. Carlini et al. (Carlini and Wagner, 2017) propose a stronger C&W attack that breaks heuristic defenses that are effective against adversarial examples generated using L-BFGS method. Madry et al. (Madry et al., 2017) propose the Projected Gradient Attack (PGD) and successfully break defenses that are effective against FGSM attacks. C&W and PGD attacks are commonly adopted as benchmarks for evaluating the robustness of new defenses as non-certified defenses can be easily evaded by considering some variants of the two attacks (Athalye et al., 2018). However, these attacks cannot be applied to our setting because of the discrete nature of graph data.

Attack on Graphs. Some earlier attacks on graphs focus on modifying the graph structure. Chaoji et al. (2012) add edges to maximize the content spreading in social network platforms such as Twitter. Researchers also reveal that the shortest path of a graph can be changed by slightly perturbing its structure (Israeli and Wood, 2002; Phillips, 1993). Csji et al. (2014) aim to maximize the PageRank score of a target node in a network with structure manipulation. However, all these attacks are not designed for graph learning algorithms (*e.g.*, graph neural networks or node embedding algorithms such as DeepWalk).

Recently, researchers also demonstrate the vulnerability of graph learning algorithms. Chen et al. (2017) inject noise to a bipartite graph that represents DNS queries to mislead the result of graph clustering. However, their attack is generated through manual effort based on attacker’s domain knowledge. Zhao et al. (2018) study the poisoning attacks on multi-task relationship learning, but based on an assumption that the sampled nodes are i.i.d. within each task, which does not hold for the node classification task. Dai et al. (2018) study evasion attacks on node classification and graph classification problems. However, their perturbation is only limited to the edges of the graph, while attackers can benefit more by additionally manipulating features of nodes (shown in Section 5.4). Zügner et al. (2018) propose Nettack, which manipulates both edges and features of the graph with a greedy approach. In addition, the authors propose an efficient method to calculate the constraint condition on the perturbations to ensure the generated perturbations are “unnoticeable”. Bojcheski and Günnemann (2018) study poisoning attack on unsupervised node embeddings by borrowing ideas from matrix perturbation theory to maximize the loss of DeepWalk (Perozzi et al., 2014b) and change the embedding outcome. Sun et al. (2019) study the new attack scenario by injecting vicious nodes to perturb the graph using a reinforcement learning strategy. However, the complexity of reinforcement learning (Watkins and Dayan, 1992) is pretty high thus cannot scale to large graphs. Zügner and Günnemann (2019) propose a data poisoning attack named Meta-attack based on meta-learning, which can reduce the overall classification accuracy of the GCN by only perturbing small fraction of the training data.

Note that, all the aforementioned attacks except (Sun et al., 2019) on graph learning algorithms focus on changing edges or features of existing nodes and are not designed for injecting vicious nodes to the graph. As shown in Section 5,

directly adapting these attacks to the new attack scenario is not promising due to high complexity and we are motivated to devise a new attack strategy tailored for the vicious node setting with a low computational cost.

3 Problem Definition

3.1 Notation and Preliminary

In this section, we first introduce the notations used throughout this paper and preliminaries of GCNs. Here, following the standard notations in the literature (Kipf and Welling, 2016; Zhou et al., 2018; Wu et al., 2019), we assume that $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{F})$ denotes an undirected attributed network (graph) with n nodes (e.g., $v_i \in \mathcal{V}$), m edges (e.g., $e_{ij} = (v_i, v_j) \in \mathcal{E}$), and d attributes (features) (e.g., $f_i \in \mathcal{F}$). The features of these n nodes are given by $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \{0, 1\}^{n \times d}$, where $\mathbf{x}_i \in \{0, 1\}^d$ denotes the feature for the i -th node v_i . The adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ contains the information of node connections, where each component \mathbf{A}_{ij} denotes whether the edge e_{ij} exists in the graph. Here, we represent the graph as $\mathcal{G} = (\mathbf{A}, \mathbf{X})$ for simplicity. Note that only a limited number of nodes possess label information in many real-world scenarios and we denote these nodes as \mathcal{V}_L , where each node $v_i \in \mathcal{V}_L$ is affiliated with the label $c_v \in \mathcal{C}$.

Following the well-established work on node classification (Kipf and Welling, 2016), the probability of classification with GCN is formulated as:

$$\mathbf{Z} = f_{\theta}(\mathcal{G}) = \text{softmax} \left(\hat{\mathbf{A}} \sigma \left(\hat{\mathbf{A}} \mathbf{X} \mathbf{W}^{(1)} \right) \mathbf{W}^{(2)} \right), \quad (1)$$

where \mathbf{Z}_{vc} denotes the probability of assigning node v to the class c ; $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ is calculated by $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ and the diagonal matrix $\tilde{\mathbf{D}}$ with diagonal element $\tilde{D}_{ii} = \sum_j \tilde{\mathbf{A}}_{ij}$ for $i = 1, 2, \dots, n$; $\theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}\}$ collects all the trainable parameters; $\sigma(\cdot)$ is an activation function (ReLU is used in this paper). In the framework of semi-supervised classification, the optimal parameter θ^* is learned by minimizing the cross-entropy loss on the labeled nodes \mathcal{V}_L , i.e.,

$$\min_{\theta} \mathcal{L}_{train}(f_{\theta}(\mathcal{G})) = - \sum_{v \in \mathcal{V}_L} \ln \mathbf{Z}_{vc_v}. \quad (2)$$

3.2 Problem Definition and Methodology

Different from previous works (Dai et al., 2018; Zügner et al., 2018; Zügner and Günnemann, 2019), in this paper, we consider a new scenario: attacking a specific target node $v_0 \in \mathcal{V}$ to change its prediction by injecting n_{in} vicious nodes that are not on the original graph, denoted by \mathcal{V}_{in} with $|\mathcal{V}_{in}| = n_{in}$ and $\mathcal{V}_{in} \cap \mathcal{V} = \emptyset$. Formally, let $\mathcal{G}' = (\mathbf{A}', \mathbf{X}')$ be the new graph after performing small perturbations on the original graph \mathcal{G} , then we have

$$\mathbf{A}' = \begin{bmatrix} \mathbf{A} & \mathbf{E} \\ \mathbf{E}^\top & \mathbf{O} \end{bmatrix}, \quad \mathbf{X}' = \begin{bmatrix} \mathbf{X} \\ \mathbf{X}_{in} \end{bmatrix}. \quad (3)$$

Here $\mathbf{E} \in \{0, 1\}^{n \times n_{in}}$ denotes the relationship matrix between original nodes and vicious nodes. Specifically, $\mathbf{E}_{ij} = 1$ if the original node $v_i \in \mathcal{V}$ is connected to the vicious node $v_j \in \mathcal{V}_{in}$, and $\mathbf{E}_{ij} = 0$ otherwise. Symmetric matrix $\mathbf{O} \in \{0, 1\}^{n_{in} \times n_{in}}$ represents the relationships between vicious nodes in \mathcal{V}_{in} . The edge information

denoted by \mathbf{E} and \mathbf{O} are called vicious edges in this paper. $\mathbf{X}_{in} \in \{0, 1\}^{n_{in} \times d}$ is the feature matrix of vicious nodes. It is noteworthy that the perturbations can only be performed on \mathbf{E} , \mathbf{O} and \mathbf{X}_{in} while \mathbf{A} and \mathbf{X} remain unchanged.

Formally, the problem of adversarial attacks on graph \mathcal{G} in the new scenario is typically formulated as the following bi-level optimization problem ¹:

$$\begin{aligned} \min_{\{\mathbf{E}, \mathbf{O}, \mathbf{X}_{in}\} \in \Phi(\mathcal{G}')} \quad & \mathcal{L}_{atk}(f_{\theta^*}(\mathcal{G}')) = \mathbf{Z}'_{v_0 c_{v_0}} - \max_{c_{new} \neq c_{v_0}} \mathbf{Z}'_{v_0 c_{new}} \\ \text{s.t.} \quad & \theta^* = \arg \min_{\theta} \mathcal{L}_{train}(f_{\theta}(\mathcal{G}')). \end{aligned} \quad (4)$$

where $\mathbf{Z}' = f_{\theta^*}(\mathcal{G}')$ denotes the prediction confidence scores for all classes on the perturbed graph \mathcal{G}' , \mathcal{L}_{atk} is the loss function during attack, $\Phi(\mathcal{G}')$ is the constraints that \mathbf{E} , \mathbf{O} and \mathbf{X}_{in} should meet on the perturbed graph \mathcal{G}' . In the inner optimization problem, we get the optimal weights θ^* of the model f (i.e., GCN) on the current perturbed graph \mathcal{G}' and in the outer optimization problem, we get the optimal perturbation (i.e., \mathbf{E} , \mathbf{O} and \mathbf{X}_{in}) on the current model f_{θ^*} . Subsequently, we will introduce the loss function on attack and constraint conditions in details.

Loss function of attacks. The loss function of attacks \mathcal{L}_{atk} aims to find a perturbed graph \mathcal{G}' that classifies the target node v_0 as c_{new} and has maximal distance to the ground truth label c_{v_0} in terms of log-probabilities/logits. Thus, the smaller \mathcal{L}_{atk} is, the worse the classification performs on the target node v_0 . It is noteworthy that the surrogate model proposed in (Zügner et al., 2018) is typically used to generate perturbations instead of using Eq. (1) such that:

$$\mathbf{Z} = f_{\theta}(\mathcal{G}) = \hat{\mathbf{A}}^2 \mathbf{X} \mathbf{W}. \quad (5)$$

The simplification is necessary for efficiency since it enables us to derive an approximate optimal solution (shown in section 4). In this sense, the loss function on attack turns to

$$\mathcal{L}_{atk}(f_{\theta^*}(\mathcal{G}')) = [\hat{\mathbf{A}}'^2 \mathbf{X}' \mathbf{W}]_{v_0 c_{v_0}} - \max_{c_{new} \neq c_{v_0}} [\hat{\mathbf{A}}'^2 \mathbf{X}' \mathbf{W}]_{v_0 c_{new}}. \quad (6)$$

Constraint conditions. There should be some constraint conditions to ensure unnoticeable perturbations, i.e., the definition of $\Phi(\mathcal{G}')$. One of the most widely used constraints for the adversarial attack is ℓ_0 -norm constraint. Specifically, the number of vicious edges by a budget Δ_e should be sparse, i.e.,

$$\|\mathbf{E}\|_0 + 0.5\|\mathbf{O}\|_0 \leq \Delta_e, \quad (7)$$

where $\|\cdot\|_0$ denotes the ℓ_0 -norm of a matrix (i.e., the number of non-zero elements of a matrix). Additionally, if we inject vicious nodes with strange pairs of features (e.g., mutually exclusive features), it will be easily detected. For this issue, the work in (Zügner et al., 2018) proposes a statistical test based on the co-occurrence graph of features to decide that a feature is unnoticeable if it occurs together with a node's original features. However, note that there are no original features for vicious nodes in our new scenario. Indeed, we can design the features arbitrarily. In this paper, we take a more practical solution by modifying the constraint such

¹ The problem is formulated as bi-level optimization because the perturbed test input is also used in the training procedure and the model weight is dependent on perturbed test data.

that the vicious nodes cannot import co-occurrence pairs of features that do not exist in the original graph. Formally:

$$\begin{aligned} \exists f_i, f_j \in \mathcal{F}, \exists u_{in} \in \mathcal{V}_{in}, [\mathbf{X}_{in}]_{u_{in}, i} = 1 \wedge [\mathbf{X}_{in}]_{u_{in}, j} = 1 \\ \text{only if } \exists u \in \mathcal{V}, \mathbf{X}_{u, i} = 1 \wedge \mathbf{X}_{u, j} = 1 \end{aligned} \quad (8)$$

This constraint means that if feature i and feature j occur together on the vicious node u_{in} , they must have occurred together in an original node u . In other words, no new co-occurrence pairs of features will be imported on the vicious nodes. Moreover, considering that vicious nodes with too few or too many features may be noticeable, we constrain the ℓ_0 -norm of vicious nodes to be equal to the mean of ℓ_0 -norm of original nodes, *i.e.*, $\forall u_{in} \in \mathcal{V}_{in}, \|\mathbf{X}_{in}\|_{u_{in}, \cdot} = \|\mathbf{X}\|_0/n$ where $[\mathbf{X}_{in}]_{u_{in}, \cdot}$ is the u_{in} -th row of matrix \mathbf{X}_{in} (*i.e.*, the feature vector of vicious node u_{in}).

4 The Proposed Framework

It is a very challenging problem to solve the proposed optimization problem in Eq. (4) due to the discreteness of the graph-structured data. In this section, we will first discuss how to adapt existing methods to the new attack scenario: injecting vicious nodes to perturb graphs. Moreover, considering the high computational complexity of existing methods, we propose a novel method (AFGSM) to speed up the computation and hence the developed method is scalable to large-scale graphs. Note that although we only focus on undirected graphs, all algorithms in this paper can be easily generalized to directed graphs.

4.1 Adapting Existing Methods for the New Scenario

In this paper, we consider three methods designed for the traditional scenario, including Nettack (Zügner et al., 2018), FGSM (Zügner et al., 2018) and Meta-attack (Zügner and Günnemann, 2019). To adapt them to the new attack scenario, we generate vicious nodes under the constraint conditions in Section 3.2 which is designed specifically for the new attack scenario instead of the constraints in (Zügner et al., 2018).

There are two strategies to adapt existing methods to the new scenario. First, we can inject all vicious nodes to the graph, and then consider the vicious nodes as special original ones. We call this strategy as **one-time injection**. Second, we can inject vicious nodes one by one and once we inject a vicious node, we optimize the corresponding edges and features. We call this strategy **sequential injection**.

Nettack for the new scenario. Nettack (Zügner et al., 2018) addresses the bi-level optimization problem following a greedy strategy. Specifically, we initialize the vicious edges in the graph \mathcal{G}' with \mathbf{E}, \mathbf{O} , and set the initial \mathbf{X}_{in} by randomly sampling from the original graph \mathcal{G} . Then we apply Nettack on the graph \mathcal{G}' and constrain that perturbations are performed on \mathbf{E}, \mathbf{O} and \mathbf{X}_{in} . In each iteration, Nettack assigns a score for each potential edge that satisfies the constraints and choose the edge with the maximum score to flip. The main cost of Nettack is the calculations of the scores, thus it derives an incremental update method that can get the new scores from the old scores after each iteration in constant time.

According to the analysis in (Zügner et al., 2018), the time complexity of Nettack in terms of n_{in} vicious nodes is $\mathcal{O}(\Delta_e \cdot n_{in} \cdot (n \cdot Nei_{v_0} + fd))$, where

Nei_{v_0} is the number of the one-order and second-order neighbors of the target node v_0 , f is the number of non-zero features in each vicious node's feature vector (*i.e.*, $\|X\|_0/n$). Specifically, there are Δ_e iterations. In each iteration, n scores for each vicious node should be calculated by considering all non-zero elements in $[\hat{A}'^2]_{v_0}$. (Nei_{v_0} edges average). And for each vicious node, once we select a feature (f at most) for it, we need to check whether the candidate features (d at most) of the node is co-occurring with it. Note that the strategies we adopt won't affect the computational complexity.

Meta-attack for the new scenario. Meta-attack (Zügner and Günnemann, 2019) utilizes the idea of meta-learning to optimize the perturbations generated on graphs. Note that Meta-attack is originally designed to attack the whole graph rather than attacking a specific target node v_0 . To apply Meta-attack to the new attack scenario, we replace the loss function \mathcal{L}_{atk} with Eq. (6) and update the meta-gradients by

$$\begin{aligned} \nabla_{\mathcal{G}'}^{meta} &= \nabla_{\mathcal{G}'} \mathcal{L}_{atk}(f_{\theta^*}(\mathcal{G}')) \\ s.t. \quad \theta^* &= \text{opt}_{\theta}(\mathcal{L}_{train}(f_{\theta}(\mathcal{G}')))) \end{aligned} \quad (9)$$

where $\text{opt}(\cdot)$ is a differentiable optimization procedure (*e.g.*, gradient descent or its variants). In each iteration, Meta-attack picks the edge and the feature value with the maximum meta-gradient to flip.

According to (Zügner and Günnemann, 2019), the time complexity of Meta-attack for the new scenario is $\mathcal{O}(\Delta_e (Tn^2 + n_{in}fd))$, where T is the number of iterations. For each iteration, the second-order gradient (Hessian matrix) should be calculated with a computational complexity of $\mathcal{O}(n^2)$. And for computation of features, it needs $\mathcal{O}(\Delta_e n_{in}fd)$ just like Nettack.

FGSM for the new scenario. FGSM (Zügner et al., 2018) computes the gradients of \mathcal{L}_{atk} w.r.t. edges and features and then it chooses the edge with the maximum revised gradient (*i.e.*, multiply -1 if the edge exists) to flip and optimize the features with the signs of gradients.

The time complexity of FGSM is $\mathcal{O}(\Delta_e (n^2 + n_{in}fd))$. For each iteration, the computation of gradients needs $\mathcal{O}(n^2)$. And like Nettack, the computation of features needs $\mathcal{O}(\Delta_e n_{in}fd)$.

4.2 Approximate Fast Gradient Sign Method (AFGSM)

Although existing methods can be adapted to the new attack scenario, their computational complexity is often too high to allow them to scale up in large-scale graphs in practice. To address this problem, we propose an efficient algorithm in this section, named Approximate Fast Gradient Sign Method (AFGSM) which injects vicious nodes one by one (*i.e.*, the sequential injection strategy).

Specifically, we inject vicious node v_{in} with edge information $\mathbf{e}_{in} \in \{0, 1\}^n$ and feature vector $\mathbf{x}_{in} \in \{0, 1\}^d$ to attack the target node v_0 in graph $\mathcal{G} = (\mathbf{A}, \mathbf{X})$, and denote the new graph by $\mathcal{G}' = (\mathbf{A}', \mathbf{X}')$, where

$$\mathbf{A}' = \begin{bmatrix} \mathbf{A} & \mathbf{e}_{in} \\ \mathbf{e}_{in}^\top & 0 \end{bmatrix} \in \{0, 1\}^{(n+1) \times (n+1)}, \mathbf{X}' = \begin{bmatrix} \mathbf{X} \\ \mathbf{x}_{in}^\top \end{bmatrix} \in \{0, 1\}^{(n+1) \times d}. \quad (10)$$

It is noteworthy that there are two possible values for the v_0 -th component of edge information \mathbf{e}_{in} , *i.e.*, $[\mathbf{e}_{in}]_{v_0} = 0$ or 1 , where $[\cdot]_k$ refers to the k -th component of

a vector. $[\mathbf{e}_{in}]_{v_0} = 1$ indicates that the vicious node v_{in} is allowed to connect to the target node v_0 directly, namely **direct attack**; otherwise, we call it **indirect attack**, which is more difficult to be detected by intelligent defense algorithms since the vicious node is not connected to the target node directly.

Here we assume that in large-scale graphs, the changes of degrees of original nodes can be ignored after injecting only one vicious node, thus we can approximate the self-loop degree matrix $\tilde{\mathbf{D}}'$ as

$$\tilde{\mathbf{D}}' \approx \begin{bmatrix} \tilde{\mathbf{D}} & \\ & \tilde{d}_{v_{in}} \end{bmatrix} \in \mathbb{R}^{(n+1) \times (n+1)}, \quad (11)$$

where $\tilde{d}_{v_{in}} = d_{v_{in}} + 1$; $d_{v_{in}}$ refers to the predefined degree of the vicious node v_{in} (i.e., how many connections it will build with existing nodes). In this sense, the Laplacian matrix of graph \mathcal{G}' is calculated by

$$\hat{\mathbf{A}}' = \tilde{\mathbf{D}}'^{-\frac{1}{2}} (\mathbf{A}' + \mathbf{I}) \tilde{\mathbf{D}}'^{-\frac{1}{2}} \approx \begin{bmatrix} \hat{\mathbf{A}} & \hat{\mathbf{e}}_{in} \\ \hat{\mathbf{e}}_{in}^\top & \tilde{d}_{v_{in}}^{-1} \end{bmatrix} \in \mathbb{R}^{(n+1) \times (n+1)},$$

where $\hat{\mathbf{e}}_{in} = (\tilde{d}_{v_{in}})^{-\frac{1}{2}} \tilde{\mathbf{D}}'^{-\frac{1}{2}} \mathbf{e}_{in}$. As a result, the probability of classification \mathbf{Z}' after perturbation is derived as

$$\mathbf{Z}' = \hat{\mathbf{A}}'^2 \mathbf{X}' \mathbf{W} \approx \begin{bmatrix} \hat{\mathbf{A}}^2 \mathbf{X} + \hat{\mathbf{e}}_{in} \hat{\mathbf{e}}_{in}^\top \mathbf{X} + \hat{\mathbf{A}} \hat{\mathbf{e}}_{in} \mathbf{x}_{in} + \tilde{d}_{v_{in}}^{-1} \hat{\mathbf{e}}_{in} \mathbf{x}_{in} \\ \hat{\mathbf{e}}_{in}^\top \hat{\mathbf{A}} \mathbf{X} + \tilde{d}_{v_{in}}^{-1} \hat{\mathbf{e}}_{in}^\top \mathbf{X} + \hat{\mathbf{e}}_{in}^\top \hat{\mathbf{e}}_{in} \mathbf{x}_{in} + \tilde{d}_{v_{in}}^{-2} \mathbf{x}_{in} \end{bmatrix} \mathbf{W}. \quad (12)$$

Specifically, the probability of classification for the target node v_0 turns to

$$\mathbf{Z}'_{v_0j} \approx [\hat{\mathbf{A}}^2 \mathbf{X} \mathbf{W}]_{v_0j} + [\hat{\mathbf{e}}_{in}]_{v_0} \hat{\mathbf{e}}_{in}^\top \mathbf{X} [\mathbf{W}]_{\cdot j} + \left(\hat{\mathbf{e}}_{in}^\top [\hat{\mathbf{A}}]_{\cdot v_0} + \tilde{d}_{v_{in}}^{-1} [\hat{\mathbf{e}}_{in}]_{v_0} \right) \mathbf{x}_{in} [\mathbf{W}]_{\cdot j} \quad (13)$$

for $j = 1, 2, \dots, |\mathcal{C}|$, where $[\cdot]_{\cdot i}$ and $[\cdot]_{\cdot j}$ denote the i -th row and j -th column of a matrix, respectively. $[\cdot]_{ij}$ refers to the i -th row and j -th component in a matrix.

Based on the approximation above, we observe that \mathbf{Z}'_{v_0j} turns out to be linear with feature vector $\mathbf{x}_{in} \in \{0, 1\}^d$ and edge information $\mathbf{e}_{in} \in \{0, 1\}^n$. Since the loss function \mathcal{L}_{atk} formulated in Eq. (6) is also linear with \mathbf{Z}'_{v_0j} , we can obtain the optimal closed-form solutions of \mathbf{x}_{in} and \mathbf{e}_{in} which minimize loss function \mathcal{L}_{atk} by their gradients.

An approximate closed-form solution of \mathbf{x}_{in} . From Eq. (13), the output \mathbf{Z}' is linear with respect to the variable \mathbf{x}_{in} . Therefore, a closed-form solution of \mathbf{x}_{in} w.r.t. the optimization problem in Eq. (4) can be obtained as follows

$$\mathbf{x}_{in}^* = -0.5 \text{sign} \left(\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{x}_{in}} \right) + 0.5, \quad (14)$$

where $\text{sign}(\cdot)$ is the element-wise function that takes the sign of a value. In other words, the features are set to 1 if the corresponding gradients in $\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{x}_{in}}$ are negative and 0 otherwise. The gradients of loss function \mathcal{L}_{atk} w.r.t. the variable \mathbf{x}_{in} can be calculated as follows

$$\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{x}_{in}} \approx \tilde{d}_{v_{in}}^{-1} \left(\mathbf{e}_{in}^\top \tilde{\mathbf{D}}'^{-\frac{1}{2}} [\hat{\mathbf{A}}]_{\cdot v_0} + \tilde{d}_{v_{in}}^{-1} \tilde{d}_{v_0}^{-\frac{1}{2}} [\mathbf{e}_{in}]_{v_0} \right) ([\mathbf{W}]_{\cdot c_{v_0}} - [\mathbf{W}]_{\cdot c_{new}}), \quad (15)$$

where $\tilde{d}_{v_{in}}^{-1} (\mathbf{e}_{in}^\top \tilde{\mathbf{D}}'^{-\frac{1}{2}} [\hat{\mathbf{A}}]_{\cdot v_0} + \tilde{d}_{v_{in}}^{-1} \tilde{d}_{v_0}^{-\frac{1}{2}} [\mathbf{e}_{in}]_{v_0})$ is always positive, and thus can be ignored since we only care about the signs of the gradients.

Algorithm 1 The proposed AFGSM algorithm.

Require: Graph $\mathcal{G} = (\mathbf{A}, \mathbf{X})$, the target node v_0 , the number of vicious nodes n_{in} , the budget of edge perturbations Δ_e ;

- 1: Train the surrogate model on the original graph \mathcal{G} and get its weight matrix \mathbf{W} ;
- 2: Randomly assign the degrees of vicious nodes $d_{v_{in}}^{(0)}, d_{v_{in}}^{(1)}, \dots, d_{v_{in}}^{(n_{in}-1)}$ to satisfy the budget constraint $\sum_{k=0}^{n_{in}-1} d_{v_{in}}^{(k)} = \Delta_e$;
- 3: $\mathcal{G}'^{(0)} \leftarrow \mathcal{G}, \mathbf{A}^{(0)} \leftarrow \mathbf{A}, \mathbf{X}^{(0)} \leftarrow \mathbf{X}$;
- 4: Calculate \mathbf{x}_{in}^* according to Eq. (14);
- 5: **for** $t = 0, \dots, n_{in} - 1$ **do**
- 6: Initialize $\mathbf{e}_{in}^{(t)} = \mathbf{0}, \mathbf{x}_{in}^{(t)} = \mathbf{0}$;
- 7: $\mathbf{x}_{in}^{(t)} \leftarrow$ Sample $\|\mathbf{X}\|_0/n$ features from \mathbf{x}_{in}^* under the constraint condition $\Phi(\mathcal{G})$;
- 8: $\mathbf{e}_{in}^{(t)} \leftarrow$ Calculate \mathbf{e}_{in}^* according to Eq. (16);
- 9: $\mathcal{G}'^{(t+1)} = (\mathbf{A}^{(t+1)}, \mathbf{X}^{(t+1)}) \leftarrow$ Update $\mathbf{A}^{(t)}$ and $\mathbf{X}^{(t)}$ according to Eq. (10);
- 10: **end for**
- 11: **return** $\mathcal{G}'^{(n_{in})} = (\mathbf{A}^{(n_{in})}, \mathbf{X}^{(n_{in})})$;

Algorithm 2 The proposed AFGSM-ada algorithm.

Require: Graph $\mathcal{G} = (\mathbf{A}, \mathbf{X})$, the target node v_0 , the number of vicious nodes n_{in} , the budget of edge perturbations Δ_e ;

- 1: Randomly assign the degrees of vicious nodes $d_{v_{in}}^{(0)}, d_{v_{in}}^{(1)}, \dots, d_{v_{in}}^{(n_{in}-1)}$ to satisfy the budget constraint $\sum_{k=0}^{n_{in}-1} d_{v_{in}}^{(k)} = \Delta_e$;
- 2: $\mathcal{G}'^{(0)} \leftarrow \mathcal{G}, \mathbf{A}^{(0)} \leftarrow \mathbf{A}, \mathbf{X}^{(0)} \leftarrow \mathbf{X}$;
- 3: **for** $t = 0, \dots, n_{in} - 1$ **do**
- 4: Train the surrogate model on the perturbed graph $\mathcal{G}'^{(t)}$ and get its weight matrix \mathbf{W} ;
- 5: Calculate \mathbf{x}_{in}^* according to Eq. (14);
- 6: Initialize $\mathbf{e}_{in}^{(t)} = \mathbf{0}, \mathbf{x}_{in}^{(t)} = \mathbf{0}$;
- 7: $\mathbf{x}_{in}^{(t)} \leftarrow$ Sample $\|\mathbf{X}\|_0/n$ features from \mathbf{x}_{in}^* under the constraint condition $\Phi(\mathcal{G})$;
- 8: $\mathbf{e}_{in}^{(t)} \leftarrow$ Calculate \mathbf{e}_{in}^* according to Eq. (16);
- 9: $\mathcal{G}'^{(t+1)} = (\mathbf{A}^{(t+1)}, \mathbf{X}^{(t+1)}) \leftarrow$ Update $\mathbf{A}^{(t)}$ and $\mathbf{X}^{(t)}$ according to Eq. (10);
- 10: **end for**
- 11: **return** $\mathcal{G}'^{(n_{in})} = (\mathbf{A}^{(n_{in})}, \mathbf{X}^{(n_{in})})$;

An approximate closed-form solution of \mathbf{e}_{in} . Without loss of generality, we assume that the vicious node v_{in} connects to a fixed number of nodes in the original graph \mathcal{G} , and therefore $d_{v_{in}} = \|\mathbf{e}_{in}\|_0$ holds. Since the output \mathbf{Z}' is linear w.r.t. the variable \mathbf{e}_{in} , we achieve a closed-form solution of \mathbf{e}_{in} with constraint in Eq. (7), *i.e.*,

$$\mathbf{e}_{in}^* = -0.5 \text{sign} \left(\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{e}_{in}} - g_{d_{in}} \mathbf{1} \right) + 0.5, \quad (16)$$

where $g_{d_{in}}$ is the d_{in} -th smallest element of $\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{e}_{in}}$. In other words, the i -th component of \mathbf{e}_{in}^* is set to 1 if the corresponding gradients are less than the d_{in} -th smallest gradient in $\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{e}_{in}}$. The gradients of loss function \mathcal{L}_{atk} w.r.t. the variable \mathbf{e}_{in} is derived as follows

$$\frac{\partial \mathcal{L}_{atk}}{\partial \mathbf{e}_{in}} \approx (\tilde{d}_{v_{in}} \tilde{d}_{v_0})^{-\frac{1}{2}} \left(\tilde{d}_{v_{in}}^{-\frac{1}{2}} [\mathbf{e}_{in}]_{v_0} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} + \tilde{\mathbf{D}}^{-1} [\tilde{\mathbf{A}}]_{v_0} \mathbf{x}_{in} \right) ([\mathbf{W}]_{\cdot c_{v_0}} - [\mathbf{W}]_{\cdot c_{new}}), \quad (17)$$

where $(\tilde{d}_{v_{in}} \tilde{d}_{v_0})^{-\frac{1}{2}}$ is always positive, and thus can be ignored.

We summarize the procedure of the proposed AFGSM in Algorithm 1. Note that the approximate closed-form solution of feature vector \mathbf{x}_{in}^* is calculated for

Table 1: The detailed statistics of the used datasets. N_{LCC} and E_{LCC} are the numbers of nodes and edges in the largest connected component, d is the dimension of features and C is the number of classes.

Dataset	N_{LCC}	E_{LCC}	d	C	frequency of classes
Citeseer	2,110	3,668	3,703	6	532, 463, 388, 308, 304, 115
Cora	2,485	5,069	1,433	7	726, 406, 379, 344, 285, 214, 131
DBLP	16,766	44,422	2,476	4	6935, 6532, 1777, 1522
Pubmed	19,717	44,324	500	3	7875, 7739, 4103
Reddit	149,177	5,215,380	602	41	28163, 15163, 13963, 13065, 12742, ... ²

once as it only depends on the weights \mathbf{W} . However, the approximate closed-form solution of edge information \mathbf{e}_{in}^* relies on variables \mathbf{A} and \mathbf{X} and thus should be updated as the vicious nodes are injected one by one in a sequential manner.

Note that Algorithm 1 treats the model weight matrix \mathbf{W} as static. Alternatively, we also develop an adaptive version of AFGSM, namely *AFGSM-ada* by re-training the surrogate model once we inject a vicious node using the AFGSM. We summarize the procedure of *AFGSM-ada* in Algorithm 2. Similarly, we also develop the adaptive version of Nettack and FGSM correspondingly, namely *Nettack-ada* and *FGSM-ada*. As for Meta-attack, it updates the model weights dynamically, so there is no need to develop its variant.

Complexity of AFGSM. The time complexity of the proposed AFGSM algorithm is $\mathcal{O}(n_{in}(n + fd))$. This is because that the gradients in Eq. (15) and Eq. (17) is concise enough such that they can be calculated by several vectors and no matrix multiplication is involved. And for each vicious node, the computation of features needs $\mathcal{O}(fd)$ just like analyzed above.

4.3 Comparison of Complexity

Comparing the time complexity of different methods, we observe that

$$\begin{aligned} \text{AFGSM: } \mathcal{O}(n_{in}(n + fd)) &< \text{Nettack: } \mathcal{O}(\Delta_e \cdot n_{inj} \cdot (n \cdot N_{ei_{v_0}} + fd)) < \\ \text{FGSM: } \mathcal{O}(\Delta_e n^2 + n_{in}fd) &< \text{Meta-attack: } \mathcal{O}(\Delta_e T n^2 + n_{in}fd) \end{aligned}$$

As a result, the proposed AFGSM algorithm is the most efficient one in terms of time complexity and then followed by Nettack, FGSM and Meta-attack. Experimental results also substantiate the conclusion in the next section.

5 Experiments

In this section, we evaluate the effectiveness and efficiency of our attack in the new attack scenario. We first provide the experimental setup (Section 5.1). Then we show the results of adapting existing methods to the new attack scenario following two different strategies: one-time injection and sequential injection (Section 5.2). Next, we explore the performance of our methods on large graphs and analyze the time cost (Section 5.3). Finally, we consider two stricter scenarios where attackers

² The frequency of classes are: 28163, 15163, 13963, 13065, 12742, 12041, 11149, 10239, 7915, 5863, 5087, 5048, 4937, 4898, 4849, 4668, 4547, 4212, 4188, 4184, 4161, 4040, 3930, 3588, 3538, 3422, 3279, 2970, 2960, 2792, 2687, 2630, 2304, 2232, 2115, 1696, 1645, 1588, 1554, 991, 328

Table 2: Accuracy of victim learning models against different attacks with one-time injection strategy w.r.t. different number of initial connections between vicious nodes and target node. *Clean* denotes the model without any attacks. *Random* denotes the model in which the vicious nodes connect to existing nodes randomly and their features are also sampled randomly according to those of existing nodes.

Method	Citeseer			Cora		
	GCN	GAT	Deepwalk	GCN	GAT	Deepwalk
Clean	0.892 ± 0.010	0.804 ± 0.008	0.736 ± 0.054	0.928 ± 0.016	0.788 ± 0.027	0.840 ± 0.028
Random	0.772 ± 0.020	0.708 ± 0.016	0.648 ± 0.041	0.868 ± 0.029	0.700 ± 0.022	0.728 ± 0.063
Nettack-0%	0.480 ± 0.025	0.460 ± 0.022	0.652 ± 0.041	0.424 ± 0.008	0.516 ± 0.023	0.744 ± 0.029
Nettack-50%	0.336 ± 0.029	0.304 ± 0.015	0.592 ± 0.081	0.352 ± 0.016	0.428 ± 0.016	0.656 ± 0.064
Nettack-100%	0.248 ± 0.010	0.244 ± 0.015	0.604 ± 0.102	0.324 ± 0.023	0.356 ± 0.022	0.651 ± 0.032
Meta-attack-0%	0.148 ± 0.020	0.156 ± 0.015	0.460 ± 0.046	0.204 ± 0.015	0.272 ± 0.020	0.484 ± 0.061
Meta-attack-50%	0.112 ± 0.027	0.140 ± 0.028	0.412 ± 0.059	0.224 ± 0.039	0.232 ± 0.010	0.520 ± 0.073
Meta-attack-100%	0.104 ± 0.008	0.164 ± 0.015	0.352 ± 0.081	0.188 ± 0.020	0.196 ± 0.023	0.452 ± 0.047
FGSM-0%	0.208 ± 0.020	0.300 ± 0.033	0.524 ± 0.041	0.336 ± 0.015	0.436 ± 0.057	0.596 ± 0.085
FGSM-50%	0.200 ± 0.013	0.216 ± 0.023	0.504 ± 0.034	0.260 ± 0.013	0.304 ± 0.015	0.536 ± 0.048
FGSM-100%	0.216 ± 0.023	0.216 ± 0.150	0.504 ± 0.066	0.292 ± 0.010	0.252 ± 0.016	0.488 ± 0.032

have limited capability and show the performance of the AFGSM method in these restricted cases (Section 5.4).

5.1 Experimental Setup

We conduct our experiments on five well-known public datasets: Citeseer (Sen et al., 2008), Cora (McCallum et al., 2000), Pubmed (Sen et al., 2008), DBLP (Zhang et al., 2019) and Reddit (Hamilton et al., 2017). The first four are citation networks where nodes are documents, edges are citation links and features are selected as the words in the document after filtering out the stop words. And the last one is a post-to-post graph where nodes are posts and edges denote these posts are from the same user. Due to the high cost of training models (*e.g.*, GCN and Deepwalk) on the original Reddit graph (around 230k nodes), we randomly sample a subgraph with nearly 150K nodes. The detailed statistics of these datasets are shown in Table 1. Following the same attack setting in (Zügner et al., 2018), we only consider the largest connected component for convenience.

In the experiments, we split the datasets into training set (10%), validation set (10%), and test set (80%). Note that in practice, attackers rarely can manipulate the training data. Therefore, we only inject vicious nodes to the test set (without labels). We first train a surrogate model on the training set, and then among all nodes that are correctly classified, we select (i) the 10 nodes with the highest margin of classification, (ii) the 10 nodes with the lowest margin of classification, (iii) 30 nodes selected randomly as our target nodes to be attacked. Since we focus on transductive classification in this paper, the model is then retrained on the mixture of clean training and perturbed test data. For each target node, we repeat the retraining process 5 times with different random seeds to stabilize the results and the average performance is reported.

5.2 Adapting Existing Methods to the New Scenario

As mentioned in the previous section, the node injection process can be performed in two ways: inject nodes all at once (one-time injection) and inject nodes sequentially (sequential injection). In this section, we compare the performance of different attacks with these two different node injection strategies.

One-time injection. The one-time injection can be roughly treated as the special case of the attack scenario considered in (Zügner et al., 2018) as nodes are added in advance and perturbations are generated using existing attacks³. One-time injection proceeds by first connecting a fixed number of vicious nodes to the target node directly and then treat the newly added nodes as existing nodes and apply the attacks proposed in (Zügner et al., 2018; Zügner and Günnemann, 2019). We choose to connect the vicious nodes directly because connections to the target node usually lead to better attack performance (Zügner et al., 2018). However, we do not know the optimal number of vicious nodes that should be connected and testing all combinations of vicious nodes is not practical. Therefore, we randomly connect 0%, 50% and 100% of the vicious nodes to the target node. Although this simple approach cannot cover all the cases, as shown below, connecting all vicious nodes (to the target node) gives the best attack performance.

The results of different attacks under different number of initial connections (to the target node) are shown in Table 2. Note that we enforce the same perturbation budget (10 nodes and 20 edges) for all methods in Table 2 for a fair comparison. First, we observe that all attacks get the best performance when we connect 100% of vicious nodes to the target node, which is also consistent with our intuition and the results in (Zügner et al., 2018) that (more) connections with target node gives better attack results. Second, we find that Meta-attack performs the best in the one-time injection. The reason is that Meta-attack generates the perturbation on edges and features utilizing the second-order gradients which can provide more information by considering the model weights dynamically. However, the complexity of Meta-attack is often very high because of the higher-order gradients and hence, cannot scale to large graphs such as DBLP and Pumbed used in the paper. Third, interestingly, we observe that FGSM performs better than Netack in the new scenario which is quite different from the results in (Zügner et al., 2018). We hypothesize that it may be due to the limited search space of Netack in the new scenario. More specifically, the initial values of E , O are extremely sparse (*i.e.*, only a small number of connected edges) compared to the number of existing nodes and edges in the graph, even if we connect 100% of vicious nodes to the target node initially. Therefore, the search space for Netack is severely limited.

Sequential injection. Next, we explore the performance of different attacks using the sequential injection strategy. Following the same setting in the one-time injection, we still connect the vicious nodes to the target node, but in a sequential manner. The results are shown in Table 3.

We find that, in the sequential addition scenario, FGSM performs better while Meta-attack performs worse compared to their counterpart in the one-time injection scenario. For example, on the Cora dataset, FGSM-one-time only lowers the GCN accuracy from 92.8% to 29.2% while FGSM-sequential lowers the accuracy to 25.6%. Differently, still on the Cora dataset and with the GCN model, Meta-attack-one-time can lower the accuracy from 92.8% to 18.8% while Meta-attack-sequential can only lower it to 24.8%. As for Netack, there is no significant difference in the performance by following different node injection strategies. Netack-one-time can

³ There is a difference in the constraint for feature perturbations. As explained in Section 3.2, we do not have specific feature constraints (however, we do not allow the occurrence of pairs of features that do not exist in the original nodes) for the vicious nodes while in the original scenario, the number of feature perturbations cannot exceed a certain threshold.

Table 3: Accuracy of victim learning models against different attacks with two different strategies. Attacks with a postfix of *one-time* are attacks with one-time node injection strategy. Attacks with a postfix *sequential* are attacks with sequential node injection strategy.

Method	Citeseer			Cora		
	GCN	GAT	Deepwalk	GCN	GAT	Deepwalk
Clean	0.892 ± 0.010	0.804 ± 0.008	0.736 ± 0.054	0.928 ± 0.016	0.788 ± 0.027	0.840 ± 0.028
Nettack-one-time	0.248 ± 0.010	0.244 ± 0.015	0.604 ± 0.102	0.324 ± 0.023	0.356 ± 0.022	0.651 ± 0.032
Nettack-sequential	0.252 ± 0.024	0.228 ± 0.010	0.644 ± 0.066	0.316 ± 0.015	0.356 ± 0.020	0.784 ± 0.065
Meta-attack-one-time	0.104 ± 0.008	0.164 ± 0.015	0.352 ± 0.081	0.188 ± 0.020	0.196 ± 0.023	0.452 ± 0.047
Meta-attack-sequential	0.120 ± 0.012	0.228 ± 0.016	0.392 ± 0.071	0.248 ± 0.024	0.240 ± 0.013	0.420 ± 0.046
FGSM-one-time	0.216 ± 0.023	0.216 ± 0.150	0.504 ± 0.066	0.292 ± 0.010	0.252 ± 0.016	0.488 ± 0.032
FGSM-sequential	0.156 ± 0.020	0.144 ± 0.015	0.408 ± 0.061	0.256 ± 0.008	0.316 ± 0.008	0.456 ± 0.023

lower the accuracy on Cora to 32.4% while Nettack-sequential lowers it to 31.6%. For GCN on the Citeseer, Nettack-one-time lowers the accuracy to 24.8% while Nettack-sequential lowers it to 25.2%, which has a negligible difference. We hypothesize that different performance of Meta-attack and FGSM under the two injection strategies is related to the search space identified by the node injection strategies and the nature of the attacks.

The main difference between the two node injection strategies is the degree of freedom in their valid search space. For the sequential injection, attacks are limited to manipulate a limited number of edges for each vicious node as we have a constraint on the degree of each vicious node. Therefore, for attacks that utilize limited information (e.g., first-order gradient for FGSM), this limitation helps to prevent attacks from manipulating too many sub-optimal edges for a single node and hence avoids getting stuck into bad solutions. In contrast, for the one-time injection, attacks are only constrained by the total number of vicious edges and thus have higher freedom when perturbing edges. For attacks with limited information, a higher degree of freedom can easily lead to suboptimal solutions while with more information (e.g., second-order derivative for Meta-attack), a higher degree of freedom leads to better attack results. For Nettack, it still suffers from the limited search space with sequential injection (similar to the analysis in one-time injection) and hence, doesn't show significant difference under different node injection strategies. Moreover, we observe that FGSM with sequential injection performs relatively closely to the costly Meta-attack with the one-time injection and the gap can be further reduced by adaptively retraining the model during attack process and more details can be found in Table 4. However, in comparison to the complicated second order in Meta-attack, the first-order gradient in FGSM can be easily approximated and hence more efficient *AFGSM* is proposed in this paper.

For experiments presented in the rest part of the paper, when choosing the baseline methods, for each method, we select the best one under the two node injection strategies. Specifically, we select FGSM-sequential, Meta-attack-one-time, and Nettack-one-time. For convenience, we still denote them as FGSM, Meta-attack and Nettack, respectively.

5.3 Attack by AFGSM and its adaptive variant.

We conduct experiments on small and large-scale graphs to show the effectiveness and efficiency of our attack. For each attack in this section, we additionally consider

Table 4: Accuracy of victim learning models against different attacks and adaptive variants. Attacks with a postfix of *ada* are adaptive attacks with victim learning models retrained during the attack process.

Method	Citeseer			Cora		
	GCN	GAT	Deepwalk	GCN	GAT	Deepwalk
Clean	0.892 ± 0.010	0.804 ± 0.008	0.736 ± 0.054	0.928 ± 0.016	0.788 ± 0.027	0.840 ± 0.028
Random	0.772 ± 0.020	0.708 ± 0.016	0.648 ± 0.041	0.868 ± 0.029	0.700 ± 0.022	0.728 ± 0.063
Nettack	0.248 ± 0.010	0.244 ± 0.015	0.604 ± 0.102	0.324 ± 0.023	0.356 ± 0.022	0.651 ± 0.032
Nettack-ada	0.256 ± 0.008	0.224 ± 0.023	0.556 ± 0.069	0.304 ± 0.023	0.304 ± 0.019	0.656 ± 0.064
Meta-attack	0.104 ± 0.008	0.164 ± 0.015	0.352 ± 0.081	0.188 ± 0.020	0.196 ± 0.023	0.452 ± 0.047
FGSM	0.156 ± 0.020	0.144 ± 0.015	0.408 ± 0.061	0.256 ± 0.008	0.316 ± 0.008	0.456 ± 0.023
FGSM-ada	0.112 ± 0.010	0.136 ± 0.015	0.444 ± 0.057	0.208 ± 0.020	0.364 ± 0.015	0.472 ± 0.056
AFGSM	0.224 ± 0.023	0.192 ± 0.016	0.532 ± 0.060	0.304 ± 0.015	0.404 ± 0.041	0.580 ± 0.051
AFGSM-ada	0.104 ± 0.027	0.128 ± 0.020	0.484 ± 0.055	0.212 ± 0.008	0.388 ± 0.036	0.588 ± 0.056

an adaptive variant of the attack, which trains the model dynamically during the attack process. Therefore, both AFGSM and FGSM have two attack forms: one with fixed model during the attack and one with dynamically retrained model during the attack. Meta-attack does not have its adaptive variant because the original attack already retrains the model during the attack. We denote the attack with adaptive training by adding a postfix “ada”.

On small graphs. First, we conduct experiments on small graphs (*e.g.*, Citeseer and Cora). The results are shown in Table 4. We find that adaptively training models during the attack process helps FGSM and AFGSM achieve better performance. For example, AFGSM for GCN on the Citeseer dataset only lowers the accuracy from 89.2% to 22.4% while AFGSM-ada lowers it to 10.4%. Although Meta-attack still performs the best among all methods, FGSM-ada and AFGSM-ada get a pretty close performance to Meta-attack. AFGSM performs similar to FGSM because our approximation technique preserves the attack effectiveness while improves the efficiency significantly. To show that the good performance of AFGSM and AFGSM-ada does not depend on a specific set of hyperparameters, we also compare all the attacks using different constraint budgets. The results are shown in Figure 2. We can easily find that under different sets of attack hyperparameters, AFGSM is still very effective and AFGSM-ada performs close to the best performing Meta-attack. We emphasize that AFGSM achieves the comparably good performance with much lower computational complexity.

On large graphs. To show the scalability of our algorithm, we first conduct experiments on large DBLP and Pubmed datasets (still with 10 vicious nodes and 20 edges). For graphs with 10K+ nodes, we can still obtain results for Nettack, FGSM, and AFGSM and their adaptive variants. Meta-attack cannot scale to these graphs and hence, we do not include the results of Meta-attack. Details of the experiments can be found in Table 5. We observe similar results as of small graphs: AFGSM performs closely to FGSM, and their adaptive variants provide better performance. One exception happens for Deepwalk on DBLP, which might be because of the poor transferability of the adaptive attack on GCN to Deepwalk (attacks on GAT and Deepwalk are all transferred from the attacks on GCN). Among all methods, FGSM-ada performs the best and AFGSM-ada is close to FGSM-ada. However, due to high complexity, FGSM or FGSM-ada cannot scale

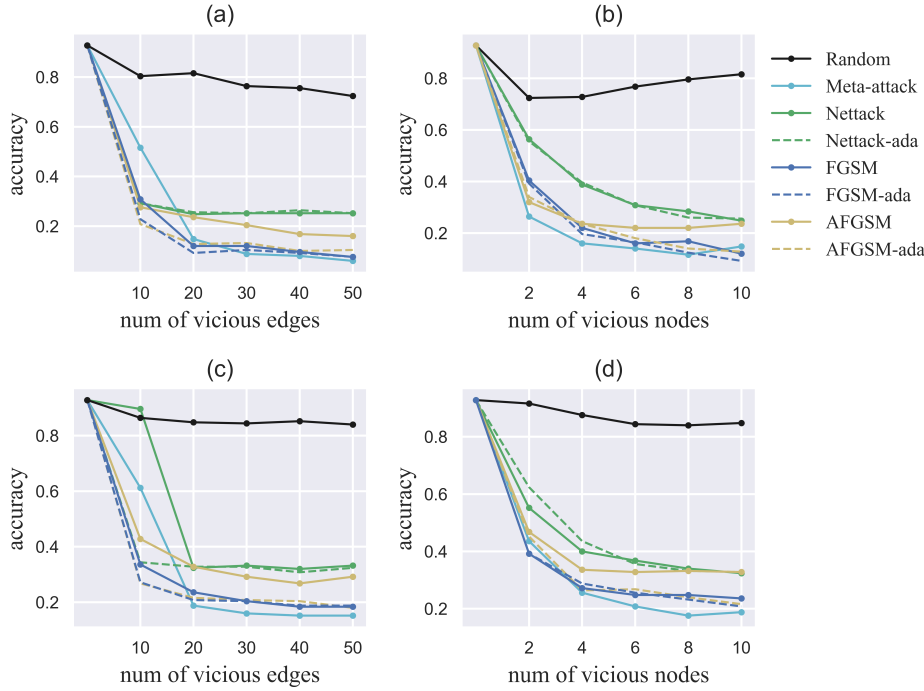


Fig. 2: Accuracy on GCN with different numbers of vicious nodes and edges. (a) and (c) denote the performance of GCN with different numbers of vicious edges and 10 nodes on Citeseer and Cora, respectively. (b) and (d) denote the performance of GCN with different numbers of vicious nodes and 10 nodes on Citeseer and Cora, respectively.

to graphs with more than 30K nodes (tested on our machine⁴ using sampled subgraphs from Reddit).

We further test the performance of our attack on larger graphs, where none of the baseline attacks can scale. We construct the large graph by subsampling from the Reddit dataset with 150K nodes. Note that we only perturb edges because the features of Reddit are preprocessed and it is impractical to directly manipulate the preprocessed features. Considering that the results of target nodes with higher degrees are harder to be mislead (Zügner et al., 2018), we attack the target nodes with $d_{v_0}/2$ vicious nodes and d_{v_0} edges. The results are shown in Table 6. Our AFGSM can still significantly reduce the accuracy on the GCN model. The transferability of the attack on GCN to Deepwalk is relatively low and we leave it as future work to improve the transferability of AFGSM on extremely large graphs. We also note that, when evaluating our attack on extremely large graphs, the main bottleneck will be in the model training instead of the attack process. Therefore, as long as there are efficient methods to train models on extremely large graphs, our attack can always scale and (highly probably) work.

⁴ We record the actual run time on the same machine with configuration: CPU (i9-7900X, 3.30GHz), 128GB RAM.

Table 5: Accuracy of victim learning models against different attacks on large graphs.

Method	DBLP			Pubmed		
	GCN	GAT	Deepwalk	GCN	GAT	Deepwalk
Clean	0.880 ± 0.025	0.808 ± 0.010	0.856 ± 0.023	0.904 ± 0.022	0.848 ± 0.020	0.832 ± 0.020
Random	0.712 ± 0.020	0.724 ± 0.015	0.832 ± 0.047	0.848 ± 0.016	0.796 ± 0.023	0.792 ± 0.032
Nettack	0.268 ± 0.016	0.420 ± 0.021	0.712 ± 0.035	0.152 ± 0.027	0.252 ± 0.010	0.796 ± 0.046
Nettack-ada	0.272 ± 0.016	0.436 ± 0.019	0.808 ± 0.016	0.156 ± 0.023	0.240 ± 0.000	0.824 ± 0.030
FGSM	0.260 ± 0.000	0.644 ± 0.015	0.652 ± 0.032	0.120 ± 0.013	0.220 ± 0.000	0.516 ± 0.037
FGSM-ada	0.240 ± 0.000	0.592 ± 0.016	0.664 ± 0.034	0.116 ± 0.015	0.220 ± 0.000	0.504 ± 0.029
AFGSM	0.252 ± 0.016	0.528 ± 0.016	0.604 ± 0.034	0.156 ± 0.008	0.224 ± 0.008	0.728 ± 0.030
AFGSM-ada	0.216 ± 0.013	0.460 ± 0.013	0.656 ± 0.028	0.136 ± 0.020	0.224 ± 0.023	0.648 ± 0.030

Table 6: Accuracy on Reddit

Method	Reddit	
	GCN	Deepwalk
Clean	0.860 ± 0.010	0.96 ± 0.000
Random	0.860 ± 0.042	0.944 ± 0.008
AFGSM	0.572 ± 0.027	0.892 ± 0.020

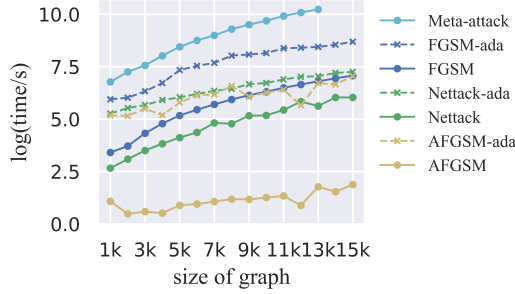


Fig. 3: Run time comparison on the GCN model.

Time cost analysis. To compare the computational cost of different methods, we conduct experiments (on the same machine mentioned above) on subgraphs of DBLP with different sizes (*i.e.*, varying from 1K to 15K nodes with a step size of 1K nodes). The results of time cost are shown in Figure 3. It is obvious that AFGSM is much more efficient than other baseline attacks, and the order of time cost aligns well with the complexity analysis in Section 4.3 (*i.e.*, AFGSM is the most efficient and followed by Nettack, FGSM, and Meta-attack). Furthermore, the execution time of AFGSM remains relatively stable when the size of the graph grows, while the run time of other attacks grow much faster when the graph size increases, especially for Meta-attack (time cost is 3,000 times higher than AFGSM and cannot scale to the DBLP dataset with more than 13K nodes).

5.4 Edge-only perturbation and Indirect perturbation

In this section, we explore two additional restricted perturbations: edges-only perturbation and indirect perturbation to verify the robustness of our attack algorithm in more practical and restricted settings.

Edge-only perturbation. From the practical point of view, manipulation of features can be hard since attackers may not have the knowledge of how features in the graph are selected and preprocessed. Therefore, we study the performance AFGSM when attackers are only allowed to change edges of vicious nodes. To do so, we inject vicious nodes with random features sampled from the original graph instead of some well-designed features to get rid of the impact of features and only focus on the edges. We denote our attack in the restricted setting as AFGSM-edges

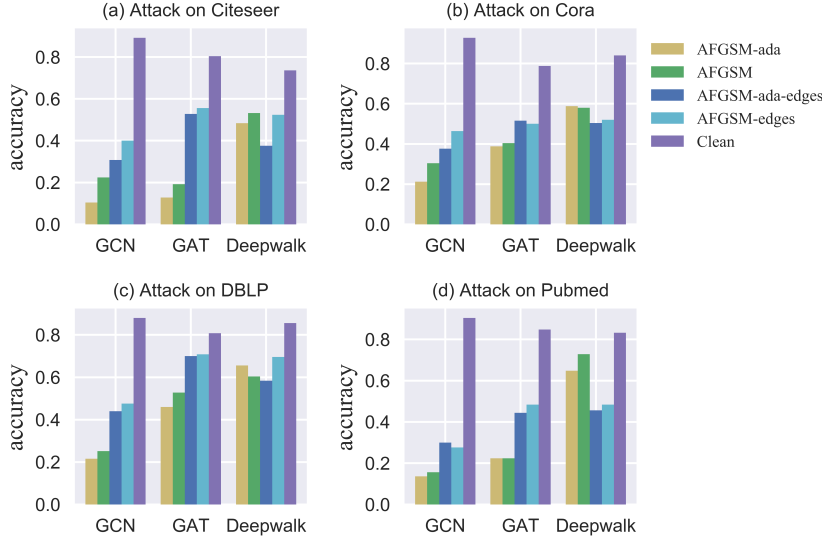


Fig. 4: Accuracy of victim learning model against edge-only attacks.

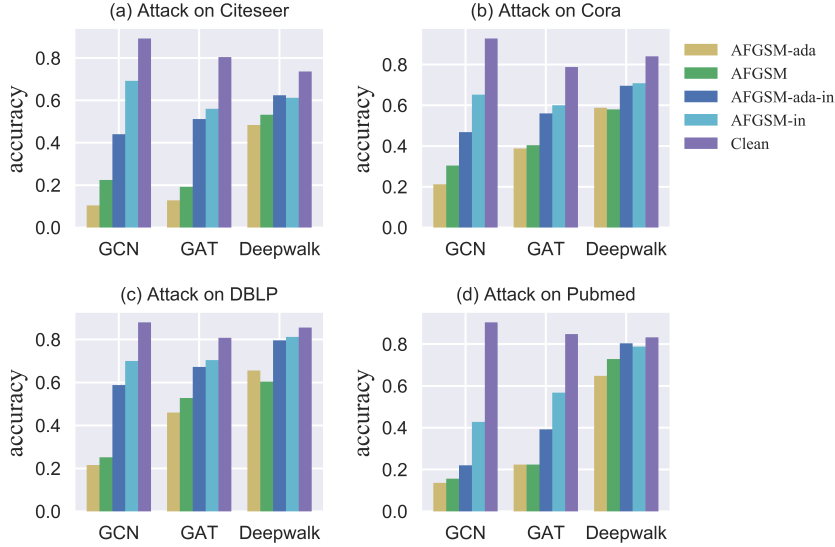


Fig. 5: Accuracy of victim learning model against indirect attacks.

since we only optimize edges during the attack process. The results are shown in Figure 4. We observe that AFGSM in the restricted setting still succeeds. We further verify whether features and edges are equally important in the new attack scenario. By comparing the performance of AFGSM and AFGSM-edges as well as AFGSM-ada and AFGSM-ada-edges, we observe that perturbing only the edges limits the attack performance significantly. This is in contrast to the findings in the attack scenario of (Zügner et al., 2018), where the authors observe that manip-

ulating features are not very important for successful attacks. Such contrast exists because, in the scenario of (Zügner et al., 2018), the attacker can only perturb features of the original nodes slightly (to remain unnoticeable) while in the new attack scenario, original features of vicious nodes can be rather arbitrary (but perturbations still follow the constraint in Eq. (7)). Hence, the new attack scenario grants the attacker more freedom in designing the perturbed features.

Indirect perturbation. In some cases, attackers may not be able to build connections with the target node directly. For example, some Facebook users can change their privacy settings and do not allow friend requests from unknown users who do share any mutual friends. Therefore, we evaluate the performance of our attack when attackers are not allowed to build direct connections to the target node. We denote the variant of our attack as AFGSM-in in the restricted setting. Results are shown in Figure 5. We observe that even in the restricted setting, the attack still succeeds in fooling the victim learning model. Unsurprisingly, we also observe that attacks in the current setting are much less successful than the attack in the case where vicious nodes can be directly connected to the target node. This observation also highlights the importance of building direct connections to the target node.

6 Conclusion

In this paper, we consider a practical attack scenario against GCN models, where attackers are only allowed to inject vicious nodes to the graph. We then propose a new attack algorithm named AFGSM to generate reliable adversarial perturbations efficiently. Through extensive experimental evaluations, we verify that GCN models are still vulnerable in the new attack setting and further show that our proposed AFGSM method outperforms the baselines significantly in terms of attack effectiveness and efficiency.

References

- Akoglu L, Tong H, Koutra D (2015) Graph based anomaly detection and description: a survey. *Data mining and knowledge discovery* 29(3):626–688
- Athalye A, Carlini N, Wagner D (2018) Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:180200420*
- Bhagat S, Cormode G, Muthukrishnan S (2011) Node classification in social networks. In: *Social network data analytics*, Springer, pp 115–148
- Bojcheski A, Günnemann S (2018) Adversarial attacks on node embeddings. *arXiv preprint arXiv:180901093*
- Bojchevski A, Günnemann S (2018) Adversarial attacks on node embeddings via graph poisoning. *arXiv preprint arXiv:180901093*
- Bolton RJ, Hand DJ, et al. (2001) Unsupervised profiling methods for fraud detection. *Credit scoring and credit control VII* pp 235–255
- Cai H, Zheng VW, Chang KCC (2018) A comprehensive survey of graph embedding: Problems, techniques, and applications. *IEEE Transactions on Knowledge and Data Engineering* 30(9):1616–1637
- Carlini N, Wagner D (2017) Towards evaluating the robustness of neural networks. In: *2017 IEEE Symposium on Security and Privacy (SP)*, IEEE, pp 39–57
- Chaoji V, Ranu S, Rastogi R, Bhatt R (2012) Recommendations to boost content spread in social networks
- Chen J, Ma T, Xiao C (2018) Fastgcn: fast learning with graph convolutional networks via importance sampling. *arXiv preprint arXiv:180110247*
- Chen Y, Nadji Y, Kountouras A, Monroe F, Vasiloglou N (2017) Practical attacks against graph-based clustering
- Csji BC, Jungers RM, Blondel VD (2014) Pagerank optimization by edge selection. *Discrete Applied Mathematics* 169(6):73–87

- Dai H, Li H, Tian T, Huang X, Wang L, Zhu J, Song L (2018) Adversarial attack on graph structured data. arXiv preprint arXiv:180602371
- Dong Y, Liao F, Pang T, Su H, Zhu J, Hu X, Li J (2018) Boosting adversarial attacks with momentum. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 9185–9193
- Goodfellow IJ, Shlens J, Szegedy C (2014) Explaining and harnessing adversarial examples. arXiv preprint arXiv:14126572
- Hamilton W, Ying Z, Leskovec J (2017) Inductive representation learning on large graphs. In: Advances in neural information processing systems, pp 1024–1034
- Israeli E, Wood RK (2002) Shortest-path network interdiction. *Networks: An International Journal* 40(2):97–111
- Kipf TN, Welling M (2016) Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:160902907
- Madry A, Makelov A, Schmidt L, Tsipras D, Vladu A (2017) Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations, 2018
- McCallum AK, Nigam K, Rennie J, Seymore K (2000) Automating the construction of internet portals with machine learning. *Information Retrieval* 3(2):127–163
- Monti F, Boscaini D, Masci J, Rodola E, Svoboda J, Bronstein MM (2017) Geometric deep learning on graphs and manifolds using mixture model cnns. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 5115–5124
- Perozzi B, Akoglu L, Iglesias Sánchez P, Müller E (2014a) Focused clustering and outlier detection in large attributed graphs. In: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, pp 1346–1355
- Perozzi B, Al-Rfou R, Skiena S (2014b) Deepwalk: Online learning of social representations. In: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, pp 701–710
- Pham T, Tran T, Phung D, Venkatesh S (2017) Column networks for collective classification. In: Thirty-First AAAI Conference on Artificial Intelligence
- Phillips CA (1993) The network inhibition problem. In: *Acm Symposium on Theory of Computing*
- Sen P, Namata G, Bilgic M, Getoor L, Galligher B, Eliassi-Rad T (2008) Collective classification in network data. *AI magazine* 29(3):93–93
- Sun Y, Wang S, Tang X, Hsieh TY, Honavar V (2019) Node injection attacks on graphs via reinforcement learning. arXiv preprint arXiv:190906543
- Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2013) Intriguing properties of neural networks. arXiv preprint arXiv:13126199
- Tang J, Aggarwal C, Liu H (2016) Node classification in signed social networks. In: Proceedings of the 2016 SIAM international conference on data mining, SIAM, pp 54–62
- Tian F, Gao B, Cui Q, Chen E, Liu TY (2014) Learning deep representations for graph clustering. In: Twenty-Eighth AAAI Conference on Artificial Intelligence
- Watkins CJ, Dayan P (1992) Q-learning. *Machine learning* 8(3-4):279–292
- Wu Z, Pan S, Chen F, Long G, Zhang C, Yu PS (2019) A comprehensive survey on graph neural networks. arXiv preprint arXiv:190100596
- Ying R, He R, Chen K, Eksombatchai P, Hamilton WL, Leskovec J (2018) Graph convolutional neural networks for web-scale recommender systems. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp 974–983
- Zhang D, Yin J, Zhu X, Zhang C (2019) Attributed network embedding via subspace discovery. arXiv preprint arXiv:190104095
- Zhao M, An B, Yu Y, Liu S, Pan SJ (2018) Data poisoning attacks on multi-task relationship learning. In: Thirty-Second AAAI Conference on Artificial Intelligence
- Zhou J, Cui G, Zhang Z, Yang C, Liu Z, Wang L, Li C, Sun M (2018) Graph neural networks: A review of methods and applications. arXiv preprint arXiv:181208434
- Zügner D, Günnemann S (2019) Adversarial attacks on graph neural networks via meta learning. arXiv preprint arXiv:190208412
- Zügner D, Akbarnejad A, Günnemann S (2018) Adversarial attacks on neural networks for graph data. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, ACM, pp 2847–2856