

# A Survey of Adversarial Learning on Graph

LIANG CHEN, Sun Yat-sen University of China  
 JINTANG LI, Sun Yat-sen University of China  
 JIAYING PENG, Sun Yat-sen University of China  
 TAO XIE, Sun Yat-sen University of China  
 ZENGXU CAO, Hangzhou Dianzi University of China  
 KUN XU, Sun Yat-sen University of China  
 XIANGNAN HE, University of Science and Technology of China  
 ZIBIN ZHENG\*, Sun Yat-sen University of China

---

Deep learning models on graphs have achieved remarkable performance in various graph analysis tasks, e.g., node classification, link prediction and graph clustering. However, they expose uncertainty and unreliability against the well-designed inputs, i.e., *adversarial examples*. Accordingly, a line of studies have emerged for both attack and defense addressed in different graph analysis tasks, leading to the arms race in graph adversarial learning.

Despite the booming works, there still lacks a unified problem definition and a comprehensive review. To bridge this gap, we investigate and summarize the existing works on graph adversarial learning tasks systemically. Specifically, we survey and unify the existing works w.r.t. attack and defense in graph analysis tasks, and give appropriate definitions and taxonomies at the same time. Besides, we emphasize the importance of related evaluation metrics, investigate and summarize them comprehensively. Hopefully, our works can provide a comprehensive overview and offer insights for the relevant researchers. More details of our works are available at <https://github.com/gitgiter/Graph-Adversarial-Learning>.

CCS Concepts: •**Computing methodologies** → **Semi-supervised learning settings**; **Neural networks**; •**Security and privacy** → **Software and application security**; •**Networks** → **Network reliability**;

Additional Key Words and Phrases: adversarial learning, graph neural networks, adversarial attack and defense, adversarial examples, network robustness

## ACM Reference format:

Liang Chen, Jintang Li, Jiaying Peng, Tao Xie, Zengxu Cao, Kun Xu, Xiangnan He, and Zibin Zheng. 2018. A Survey of Adversarial Learning on Graph. *J. ACM* 37, 4, Article 111 (August 2018), 28 pages.  
 DOI: 10.1145/1122445.1122456

## 1 INTRODUCTION

Over the past decade, deep learning has enjoyed the status of crown jewels in artificial intelligence, which shows an impressive performance in various applications, including speech and language processing [19, 74], face recognition [47] and object detection [35]. However, the frequently used deep learning models recently have been proved unstable and unreliable due to the vulnerability

---

\*Contact Author

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed 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).

© 2018 ACM. 0004-5411/2018/8-ART111 \$15.00  
 DOI: 10.1145/1122445.1122456

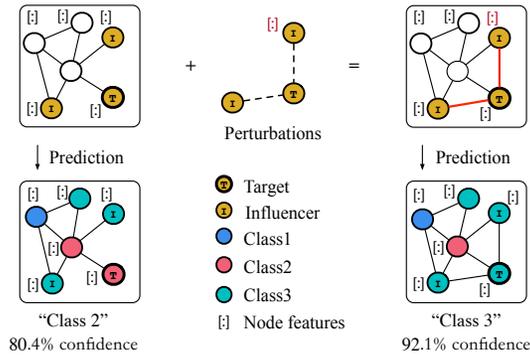


Fig. 1. A misclassification of the target caused by a small perturbations of the graph structure and node features.

against perturbations. For example, slight changes on several pixels of a picture, which appears imperceptible for human eyes but strongly affect the outputs of deep learning models [46]. As stated by Szegedy et al. [61], deep learning models that are well-defined and learned by backpropagation have intrinsic blind spots and non-intuitive characteristics, which should have been generalized to the data distribution in an obvious way.

On the other hand, deep learning on graphs has received significant research interest recently. As a powerful representation, graph plays an important role and has been widely applied in real world [27]. Naturally, deep learning's research on graph is also a hot topic and brings lots of refreshing implementations in different fields, such as social networks [48], e-commerce networks [66] and recommendation systems [15, 73]. Unfortunately, graph analysis domain, a crucial field of machine learning, has also exposed the vulnerability of deep learning models against well-designed attacks [85, 87]. For example, consider the task of node classification, attackers usually have control over several fake nodes, aim to fool the target classifiers, leading to misclassification by adding or removing edges between fake nodes and other benign ones. As shown in Figure 1, performing small perturbations (two added links and several changed features of nodes) on a clean graph can lead to the misclassification of deep graph learning models.

With rising concerns being paid on the security of graph models, there exists a surge of researches on graph adversarial learning, i.e., a field on studying the security and vulnerability of graph models. On one hand, from the perspective of attacking a graph learning model, Zügner et al. [85] first study adversarial attacks on graph data, with little perturbations on node features and graph structure, the target classifiers are easily fooled and misclassify specified nodes. On the other hand, Wang et al. [67] propose a modified Graph Convolution Networks (GCNs) model with an adversarial defense framework to improve its robustness. Moreover, Sun et al. [57] study the existing works of adversarial attack and defense strategies on graph data and discuss their corresponding contributions and limitations. However, they mainly focus on the aspect of the adversarial attack, leaving works on defense unexplored.

*Challenges.* Despite a surge of works on the graph adversarial learning, there still exists several problems to solve. i) *Unified and specified formulation.* Current studies consider the problem definition and assumptions in graph adversarial learning with their own mathematical formulations and mostly lack of detailed explanations, which effectively hinders the progress of follow-up studies. ii) *Related evaluation metrics.* While for various tasks, evaluation metrics on corresponding

performance are rather different, and even have diverse standardization on it. Besides, special metrics for the graph adversarial learning scenario are necessary and timely to explore, e.g., evaluation on the attack impacts.

For the problem of inconsistent formulations and definitions, we survey the existing attack and defense works, give unified definitions and categorize them from diverse perspectives. Although there have been some efforts [20, 39, 85] to generalize the definitions, most formulations still make customization for their own models. So far, only one article [57] outlines these concepts from a review perspective, which is not sufficient to summarize the existing works comprehensively. Based on the previous literature, we summarize the different types of graphs and introduce the three main tasks according to the level, and subsequently give the unified formulations of attack and defense in Section 3.1 and 4.1, respectively.

Different models have various metrics due to different emphasis. To provide guidance for researchers and better evaluate their adversarial models, we have a more detailed summarize and discussion on metrics in Section 5. In particular, we first introduce some common metrics for both attack and defense, and then present some special metrics provided in their respective works from three categories: effectiveness, efficiency, and imperceptibility. For instance, the Attack Success Rate (ASR) [10] and the Average Defense Rate (ADR) [11] are proposed to measure the effectiveness of attack and defense, respectively.

In summary, our contributions can be listed as bellow:

- We thoroughly investigate related works in this area, and subsequently give the unified problem formulations and the clear definitions for current inconsistent concepts of both attack and defense.
- We give a clear overview on existing works and classify them from different perspectives based on reasonable criteria systematically.
- We emphasize the importance of evaluation metrics and investigate them to make a comprehensively summary.
- For such an emerging research area, we point out the limitations of current researches and provide some open questions to solve.

The survey is organized as follow. In Section 2, we will give some basic notations of typical graphs. In Section 3 and Section 4, we will separately introduce the definitions, taxonomies of adversarial attack and defense on graph data, and further give a clear overview. We then summarize the related metrics in Section 5 and try to discuss some open research questions in Section 6. Finally, we draw our conclusion in Section 7.

## 2 PRELIMINARY

Focusing on the graph structure data, we first give notations of typical graphs for simplicity and further introduce the mainstream tasks in graph analysis fields. The most frequently used symbols are summarized in Table 1.

### 2.1 Notations

Generally, a graph is represented as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$  denotes the set of  $N$  nodes and  $\mathcal{E} = \{e_1, e_2, \dots, e_M\}$  is the set of  $M$  existing edges in the graph, and naturally  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ . The connections of the graph could be represented as an adjacency matrix  $A \in \mathbb{R}^{N \times N}$ , where  $A_{i,j} \neq 0$  if there is an edge from node  $v_i$  to node  $v_j$  and  $A_{i,j} = 0$  otherwise.

Table 1. Summary of notations.

| Symbol        | Description   | Symbol        | Description  | Symbol        | Description   |
|---------------|---|---------------|--|---------------|---|
| $N$           | number of nodes   | $M$           | number of edges  | $n$           | number of graphs  |
| $F_{(\cdot)}$ | number of features w.r.t. nodes $F_{(node)}$ or edges $F_{(edge)}$                        | $\hat{G}$     | a graph instance, denoting a clean graph $G$ or a modified graph $\hat{G}$                     | $C$           | set of pre-defined class labels for nodes, edges or graphs  |
| $\mathcal{V}$ | set of nodes $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$                                     | $\mathcal{E}$ | set of edges $\mathcal{E} = \{e_1, e_2, \dots, e_M\}$  | $\mathcal{G}$ | set of graphs $\mathcal{G} = \{G_1, G_2, \dots, G_n\}$  |
| $\mathcal{D}$ | the whole dataset, $\mathcal{D} = (\mathcal{G}, X, C)$                                    | $S$           | set of instances, could be $\mathcal{G}$ , $\mathcal{V}$ , or $\mathcal{E}$                    | $S_L$         | set of labeled instances, could be $\mathcal{G}_L$ , $\mathcal{V}_L$ , or $\mathcal{E}_L$ , $S_L \subset S$ |
| $\mathcal{T}$ | set of unlabeled instances where $\mathcal{T} \subset S - S_L$ or $\mathcal{T} = S - S_L$ | $\mathcal{K}$ | attackers' knowledge of the dataset  | $O$           | operation set w.r.t. attackers' manipulation  |
| $A$           | adjacency matrix $A \in \mathbb{R}^{N \times N}$  | $A'$          | modified adjacency matrix  | $X$           | feature matrix $X \in \{0, 1\}^{N \times F_{(\cdot)}}$ or $X \in \mathbb{R}^{N \times F_{(\cdot)}}$         |
| $X'$          | modified feature matrix   | $f$           | deep learning model w.r.t. inductive learning $f^{(ind)}$ or transductive learning $f^{(tra)}$ | $\hat{f}$     | surrogate model   |
| $\tilde{f}$   | well-designed model for defense   | $\theta$      | set of parameters w.r.t. a specific model  | $Z$           | output of the model   |
| $t$           | target instance, can be a node, an edge or a graph  | $\Delta$      | attack budgets   | $y$           | ground-truth label w.r.t. an instance, $y \in C$  |
| $\Psi$        | perturbation space of attacks   | $\mathcal{L}$ | loss function of deep learning model   | $Q$           | similarity function of graphs   |

## 2.2 Taxonomies of Graphs

Different scenarios correspond to various types of graphs, hence we will introduce them further in the following parts based on the basic graph definition in Section 2.1.

**Directed and Undirected Graph.** Directed graph, also called a digraph or a directed network, is a graph where all the edges are directed from one node to another but not backwards [7]. On the contrary, a graph where the edges are bidirectional is called an undirected graph. For undirected graphs, the convention for denoting the adjacency matrix doesn't matter, as all edges are bidirectional. Generally,  $A_{i,j} \neq A_{j,i}$  for directed graph and  $A_{i,j} = A_{j,i}$  for undirected graph.

**Weighted and Unweighted Graph.** Typically a weighted graph refers to an edge-weighted graph where each edge is associated with a real value [7]. An unweighted graph may be used if a relationship in terms of magnitude doesn't exist, i.e., the connections between edges are treated as the same.

**Attributed Graph.** An attributed graph refers to a graph where both node and edge are available to have its own attributes/features [23]. Specifically, the attributes of nodes and edges could be denoted as  $X_{node} \in \mathbb{R}^{N \times F_{node}}$  and  $X_{edge} \in \mathbb{R}^{M \times F_{edge}}$ , respectively. In most cases a graph usually have attributes associated with nodes only, we use  $X$  to denote the node attributes/features for brevity.

**Homogeneous and Heterogeneous Information Graph.** As well as attributes, the type of nodes or edges is another important property. A graph  $\mathcal{G}$  is called heterogeneous information graph if there are two or more types of objects/nodes or relations/edges in it [30, 37]; otherwise, it is called a homogeneous information graph.

**Dynamic and Static Graph.** Intuitively, nodes, edges and attributes are possibly changing over time in a dynamic graph [41], which could be represented at  $\mathcal{G}^{(t)}$  at time  $t$ . For a static graph, which is simply defined as  $\mathcal{G}$ , in which nodes, edges and attributes remain the same as the time changes.

Each type of graph has different strengths and weaknesses. It's better to pick the appropriate kind of graph to model the problem. As existing works are mainly focus on a simple graph, i.e., undirected and unweighted. Besides, they assume that the graph is static and homogeneous for simplicity. By default, the mentioned "graph" refers to "simple graph" in this paper.

### 2.3 Graph Analysis Tasks

In this section, we will introduce the major tasks of graph analysis, in which deep learning models are commonly applied to. From the perspective of nodes, edges and graphs, we divide these tasks into three categories: node-, link- and graph-level task.

**Node-level Task.** Node classification is one of the most common node-level tasks, for instance, identifying a person in a social network. Given a graph  $\mathcal{G}$ , with partial labeled nodes  $\mathcal{V}_L \subseteq \mathcal{V}$  and other unlabeled ones, the goal of classifiers is to learn a mapping function  $\phi : \mathcal{V} \rightarrow \mathcal{C}$ , where  $\mathcal{C}$  is a set of pre-defined class labels. The learned mapping function is applied to effectively identify the class labels for the unlabeled nodes [34]. To this end, inductive learning and transductive learning settings are specified based on the characteristic of training and testing procedures.

- *Inductive Setting.* For the inductive learning setting, a classifier is usually trained on a set of nodes and tested on others that never seen during training.
- *Transductive Setting.* Different from inductive setting, test samples (i.e., the unlabeled nodes) can be seen (but not their labels!) during the training procedure in transductive learning setting.

To conclude, the objective function that optimized by the classifier could be formulated as follows:

$$\Omega = \frac{1}{|\mathcal{V}_L|} \sum_{v_i \in \mathcal{V}_L} \mathcal{L}(f_\theta(\mathcal{G}, X), y_i) \quad (1)$$

where  $y_i$  is the class label of node  $v_i$ , and  $\mathcal{L}$  could be the loss of either inductive learning or transductive learning, and the classifier  $f$  with parameters  $\theta$  is similarly defined with the two settings.

**Link-level Task.** Link-level task relates to the edge classification and link prediction. Among them, link prediction is a more challenging task and widely used in real-world, which aims to predict the connection strength of an edge, e.g., predicting the potential relationship between two specified persons in a social network, and even the new or dissolution relationship in the future. Link prediction tasks take the same input as node classification tasks, while the output is binary which indicates an edge will exist or not. Therefore the objective function of link prediction tasks could be similarly defined as Eq.(1) by changing  $y_i \in \{0, 1\}$ , and replacing  $\mathcal{V}_L$  with a set of labeled edges  $\mathcal{E}_L$ .

**Graph-level Task.** While treating the graph as a special form of node, the graph-level tasks are similar to node-level tasks. Take the most frequent application, graph classification, as an example, the graph  $\mathcal{G}$  could be represented as  $\mathcal{G} = \{G_1, G_2, \dots, G_n\}$ , where  $G_i = (\mathcal{V}_i, \mathcal{E}_i)$  is a subgraph of the entire graph. The core idea to solve the graph classification problem is to learn a mapping function  $\mathcal{G} \rightarrow \mathcal{C}$ , here  $\mathcal{C}$  represents the set of graph categories, so as to predict the class for an unseen graph more accurately. The objective function of graph classification tasks is similar with Eq.(1) as well, in which the labeled training set is  $\mathcal{G}_L$  instead of  $\mathcal{V}_L$ , and  $y_i$  is the category of the graph  $G_i$ . In real

world, graph classification task plays an important role in many crucial applications, such as social and biological graph classification.

### 3 ADVERSARIAL ATTACK

In this section, we will first introduce the definition of adversarial attack against deep learning methods on graphs. Then, we categorize these attack models from different perspectives. Finally, we will give a clear overview on existing works.

#### 3.1 Definition

According to existing works on graph adversarial attacks, we summarize and give a unified formulation for them.

*Attack on Graph Data.* Considering  $f$  a deep learning function designed to tackle related downstream tasks. Given a set of target instances  $\mathcal{T} \subseteq \mathcal{S} - \mathcal{S}_L$ , where  $\mathcal{S}$  could be  $\mathcal{V}$ ,  $\mathcal{E}$  or  $\mathcal{G}$  respectively for different levels of tasks, and  $\mathcal{S}_L$  denotes the instances with labels,<sup>1</sup> the attacker aims to maximize the loss of the target node on  $f$  as much as possible, resulting in the degradation of prediction performance. Generally, we can define the attack against deep learning models on graphs as:

$$\begin{aligned} & \underset{\hat{G} \in \Psi(G)}{\text{maximize}} \sum_{t_i \in \mathcal{T}} \mathcal{L}(f_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t. } \theta^* = \arg \min_{\theta} \sum_{v_j \in \mathcal{S}_L} \mathcal{L}(f_{\theta}(\tilde{G}^{v_j}, X, t_j), y_j) \end{aligned} \quad (2)$$

We denote a graph  $G \in \mathcal{G}$  with node  $t_i$  in it as  $G^{t_i}$ .  $\hat{G}$  denotes the perturbed graph and  $\Psi(G)$  indicates the perturbation space on  $G$ , and we use  $\tilde{G}$  to represent original graph  $G$  or a modified one  $\hat{G}$ , respectively.

To make the attack as imperceptible as possible, we set a metric for comparison between the graphs before and after the attack, such that:

$$\begin{aligned} & Q(\hat{G}^{t_i}, G^{t_i}) < \epsilon \\ & \text{s.t. } \hat{G}^{t_i} \in \Psi(G) \end{aligned} \quad (3)$$

where  $Q$  denotes a similarity function, and  $\epsilon$  is an allowed threshold of changes. Specifically, more metrics will be detailed in Section 5.

#### 3.2 Taxonomies of Attacks

As we focus on a simple graph, different types of attacks could be conducted on the target systems. In this section, we provide taxonomies for the attack types mainly proposed by Sun et al. [57] and subsequently extended in our works according to various criteria. For different scenarios, we give relevant instructions and summarize them in the following:

**3.2.1 Attacker's Knowledge.** To conduct attacks on the target system, usually an attacker will possess certain knowledge about the target models and the dataset, which helps them achieve the adversarial goal. Based on the knowledge of the target models [57], we characterize different threatening levels of attacks.

<sup>1</sup>Note that attackers only focus on attacking the target instances in the test set.

**White-box Attack.** This is the simplest attacks while attackers possess the entire information of the target models, including the model architecture, parameters and gradient information, i.e., the target models are fully exposed to the attackers. By utilizing such rich information, attackers can easily affect the target models and bring destructive effects. However, it is impracticable in real world since it is costly to possess such complete knowledge of target models. Therefore white-box attack is less dangerous but often used to approximate the worst performance of a system under attack.

**Gray-box Attack.** In this case, attackers are strict to possess excessive knowledge about the target models, this reflects real world scenarios much better since it is more likely that attackers have limited access to get the information, e.g., only familiar with the architecture of the target models. Therefore, it's harder than conducting the white-box attack but more dangerous for the target models.

**Black-box Attack.** In contrast to white-box attack, black-box attack assumes that the attacker does not know anything about the targeted systems. Under this setting, attackers are only allowed to do black-box queries on limited samples at most. However, it will be the most dangerous attack once it works, since attackers can attack any models with limited (or no) information.

Here also comes a term “**no-box attack**” [55] which refers to an attack on the surrogate model based on their (limited) understanding of the target model.<sup>2</sup> As attackers have complete knowledge of the surrogate model, the adversarial examples are generated under white-box setting. No-box attack will become a white-box attack once the attackers build a surrogate model successfully and the adversarial examples are transferable to fool other target models. Theoretically, this kind of attack can hardly affect the systems in most cases with these constraints. However, as if the surrogate model has strong transferability, the no-box attack is also destructive.

As attackers are strict with the knowledge of the target models, they also have less access to the information of the dataset  $\mathcal{D}$ . Based on the different levels of knowledge of the targeted system [17], we have:

**Perfect Knowledge.** In this case, the dataset  $\mathcal{D}$ , including the entire graph structure, node features, and even the ground-truth labels of objects, are completely exposed to attackers, i.e., the attacker is assumed to know *everything* about the dataset. This is impractical but the most common setting in previous studies. However, considering the damage that could be done by an attacker with perfect knowledge is critical, since it may expose the potential weaknesses of the target system in a graph.

**Moderate Knowledge.** The moderate knowledge case represents an attacker with less information about the dataset. This makes attackers focus more on the performance of their attacks, or even search for more information through legitimate or illegitimate methods.

**Minimal Knowledge.** This is the hardest one as attackers can only conduct attacks with the minimal knowledge of datasets, such as partial connections of the graph or the feature information of certain nodes. This is essential information for an attacker, otherwise it will fail even under the white-box attack setting. The minimal knowledge case represents the least sophisticated attacker and naturally with less threat level.

To conclude, we use  $\mathcal{K}$  to denote the attackers' knowledge in an abstract knowledge space. Here  $\mathcal{K}$  is consisting of three main parts, the dataset knowledge, the training algorithm and the learned parameters of target models. For instance,  $\mathcal{K} = (\mathcal{D}, f, \theta)$  denotes white-box attack with

<sup>2</sup>We must realize that no-box is a kind of gray- or black-box attack so we don't don't have a separate category for it.

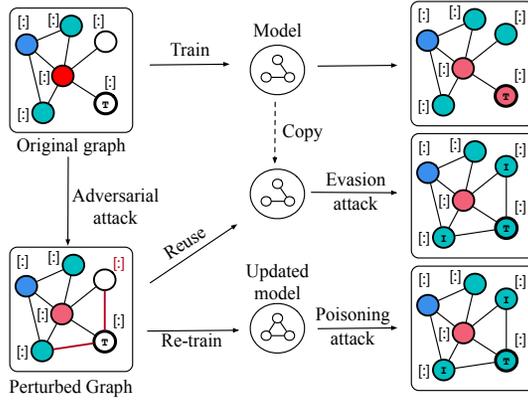


Fig. 2. An example of the evasion attack and the poisoning attack. (Image Credit: Zünger et al. [85])

perfect knowledge, and  $\mathcal{K} = (\mathcal{D}_{train}, \hat{f}, \hat{\theta})$  denotes gray- or black-box attack and with moderate knowledge, where  $\hat{f}$  and  $\hat{\theta}$  come from the surrogate model and  $\mathcal{D}_{train} \subset \mathcal{D}$  denotes the training dataset.

**3.2.2 Attacker's Goal.** The goal generally distinguishes three different types of attacks, i.e., *security violation*, *attack specificity* and *error specificity* [1]. They are not mutually exclusive in fact, and more details are discussed below.

**Security Violation.** Security violation can be categorized into *availability attack*, *integrity attack* and *others*. For availability attack, the attacker attempts to destroy the function of the system, thereby impairing its normal operation. The damage is global, that is, it attacks the overall performance of the whole system. For integrity attack, the attackers' purpose is to bypass or fool the detection of the system, which is different from availability attack in that it does not destroy the normal operation of the system. There are other goals, such as reverse-engineering model information to gain privacy knowledge.

**Error Specificity.** Take node classification task as an example, the error specific attack aims to misclassify the predictions as specific labels, while unspecific attack does not care what the prediction is, the attack is considered successful as long as the prediction is wrong.

**Attack Specificity.** This perspective focuses on the range of the attack, which can divide attacks into *targeted attack* and *non-targeted attack (general attack)*. The targeted attack focuses on a specific subset of nodes (usually a target node), while the non-targeted attack is undifferentiated and global. With reference to Eq.(2), the difference is the domain of  $\mathcal{T} \subset \mathcal{S} - \mathcal{S}_L$  or  $\mathcal{T} = \mathcal{S} - \mathcal{S}_L$ .

It is worth noting that in some other fields (e.g., computer vision), *targeted attack* refers to the specific error attack, and *non-targeted attack* the unspecific error attack. In the field of graph adversarial learning, we suggest that distinguishing the targeted attack based on the attack range, and whether it is error specific based on the consequence.

**3.2.3 Attacker's Capability.** Attacks can be divided into poisoning attack and evasion attack according to adversaries' capabilities, which are occurred at different stages of the attacks.

**Poisoning Attack.** Poisoning attacks (a.k.a training-time attacks) try to affect the performance of the target models by modifying the dataset in the training stage, i.e., the target models are trained on the poisoned datasets. Since transductive learning is widely used in most graph analysis tasks,

the test samples (but not their labels) are participated in the training stage, which leads to the popularity of poisoning attacks. Under this scenario, the parameters of target models are retrained after the training dataset being modified, thus we can define poisoning attacks according to Eq.(2):

$$\begin{aligned} & \text{maximize } \sum_{\hat{G} \in \Psi(G)} \sum_{t_i \in \mathcal{T}} \mathcal{L}(f_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t. } \theta^* = \arg \min_{\theta} \sum_{v_j \in \mathcal{S}_L} \mathcal{L}(f_{\theta}(\hat{G}^{v_j}, X, v_j), y_j) \end{aligned} \quad (4)$$

**Evasion Attack.** While poisoning attacks focus on the training phase, evasion attacks (a.k.a test-time attacks) tend to add adversarial examples in test time. Evasion attacks occur after the target model is well trained on a clean graph, i.e., the learned parameters are fixed during evasion attacks. Therefore, we can define the formulation of evasion attacks by changing part of Eq.(4) slightly:

$$\theta^* = \arg \min_{\theta} \sum_{v_j \in \mathcal{S}_L} \mathcal{L}(f_{\theta}(G^{v_j}, X, v_j), y_j) \quad (5)$$

In Figure 2, we show a toy example of the poisoning attack and the evasion attack. In most cases the poisoning attack does not work well because the model is retrained to alleviate the adversarial impacts.

**3.2.4 Attack Strategy.** For attacking a target model on graph data, attackers may have a line of strategies to achieve their adversarial goals. In most instances, they will focus on the graph structure or node/edge features. Based on the strategy applied on a graph, we have:

**Topology Attack.** Attackers are mainly focus on the the topology of the graph, a.k.a, structure attack. For example, they are allowed to add or remove some edges legally between nodes in the graph to fool the target system. To specify this, we define an attack budget  $\Delta \in \mathbb{N}$ , thus the definition of the topology attack is:

$$\begin{aligned} & \text{maximize } \sum_{\hat{G} \in \Psi(G)} \sum_{t_i \in \mathcal{T}} \mathcal{L}(f_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t. } \sum_{u < v} |A_{u,v} - A'_{u,v}| \leq \Delta \end{aligned} \quad (6)$$

where  $A'$  is the adjacency matrix of the perturbed graph, and  $\theta^*$  is discussed in Section 3.2.3.

**Feature Attack.** Although the topology attack is more common, attackers can conduct feature attacks as well. Under this setting, the features of specified objects will be changed. However, unlike graph structure, the node/edge features could be either binary or continuous, i.e.,  $X \in \{0, 1\}^{N \times F(\cdot)}$  or  $X \in \mathbb{R}^{N \times F(\cdot)}$ . For binary features, attackers can flip them like edges, while for continuous features, attackers can add a small value on them, respectively. Thus the definition of feature attacks is:

$$\begin{aligned} & \text{maximize } \sum_{\hat{G} \in \Psi(G)} \sum_{t_i \in \mathcal{T}} \mathcal{L}(f_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t. } \sum_u \sum_j |X_{u,j} - X'_{u,j}| \leq \Delta \end{aligned} \quad (7)$$

where  $X'$  is similarly defined, and here  $\Delta \in \mathbb{N}$  for binary features and  $\Delta \in \mathbb{R}$  for continuous features.

**Hybrid.** Usually, attackers will conduct both attack strategies at the same time to exert more powerful impact. Besides, they could even add several fake nodes (with fake labels), which will have their own features and relationship with other benign instances. For example, some fake users will be added into the recommendation system and affect the results of the system [16, 29]. Therefore, we conclude a unified formulation as follows:

$$\begin{aligned} & \underset{\hat{G} \in \Psi(G)}{\text{maximize}} \sum_{t_i \in \mathcal{T}} \mathcal{L}(f_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t.} \sum_{u < v} |A_{u,v} - A'_{u,v}| + \sum_u \sum_j |X_{u,j} - X'_{u,j}| \leq \Delta \end{aligned} \quad (8)$$

$A'$  and  $X'$  may not have the same dimension with  $A$  and  $X$  respectively if some fake nodes are added into the graph.

**3.2.5 Attackerffs Manipulation.** Although attackers may have enough knowledge about the target system, they may not always have access to manipulation all of the dataset. Besides, different manipulations may have different budget cost. For example, in an e-commerce system [14], attackers could only purchase more items (add edges) but unable to delete the purchase records (remove edges). Based on different manipulation, we have:

**Add.** In a graph, attackers could add edges between different nodes, or add features on specified nodes/edges. This is the most simplest manipulation and naturally has lowest budget in most scenarios.

**Remove.** In a graph, attackers could remove edges between different nodes, or remove features on specified nodes/edges. This kind of manipulation will be harder than “add”, since attackers may not have enough access to execute the “remove” manipulation.

**Rewiring.** The manipulations mentioned above may become more noticeable for the target system if  $\Delta$  is larger. To address this problem, the rewiring operation is proposed in a less noticeable way than simple adding/removing manipulation [39, 70]. For instance, attackers can add an edge connected with a target node, while remove an edge connected with it at the same time. Each time of rewiring manipulation related to two single operations, which can preserve some basic graph properties (e.g., degree) and naturally more unnoticeable.

To conclude, we define an operation set  $O = \{O_{Add}, O_{Rem}, O_{Rew}\}$ , representing the above three manipulations of attackers, respectively.

**3.2.6 Attack Algorithm.** Generally speaking, current methods of adversarial attacks on generating adversarial examples are mainly based on the gradient information, either from the target model (white-box attack) or the surrogate model (black- or gray-box attack). Beyond that, there are several methods generating adversarial examples based on other algorithms. From the perspective of attack algorithm, we have:

**Gradient-based Algorithm.** Intuitively, gradient-based algorithm is simple but effective. The core idea is that: fix the parameters of a trained model, and regard the input as a hyperparameter to optimize. Similar to the training process, attackers could use the partial derivative of loss  $\mathcal{L}$  with respect to edges (topology attack) or features (feature attack), to decide how to manipulate the dataset. However, gradients could not applied directly into the input data due to the discreteness of the graph data, instead, attackers often choose the one with greatest absolute gradients, and manipulate it with a proper value. While most deep learning models are optimized by gradients, on the contrary, attackers could destroy them by gradients as well.

**Non-gradient-based Algorithm.** In addition to gradient information, attackers could generate adversarial examples in other ways. For example, from the perspective of the genetic algorithm, attackers can choose the population (adversarial examples) with highest fitness score (e.g., erroneous outputs of the target/surrogate model) generations by generations. Besides, reinforcement learning algorithms are also commonly used to solve this issue. Reinforcement learning based attack methods [20, 39] will learn the generalizable adversarial examples within an action space. Moreover, adversarial examples can even generated by a well-designed generative model.

**3.2.7 Target Task.** As discussed in Section 2.3, there exists three major tasks in graph domains. According to different levels of tasks, we can divide existing attack methods into the following three categories, respectively.

**Node-relevant Task.** Currently, there are several attack models against node-relevant tasks [2, 4, 10, 12, 20, 20, 58, 65, 68, 70, 72, 75, 77, 82, 85, 86]. Typically, Bojchevski et al. and Hou et al. [2, 29] utilize random walk [49] as an surrogate model to attack node embedding, Dai et al. [20] uses a reinforcement learning based framework to disturb node classification tasks. In general, most of the existing works address in node classification task due to its ubiquity in real world.

**Link-relevant Task.** Many relationships in real world can be represented by graph. For some of these graphs, such as social graph, are always dynamic in reality. Therefore, the link prediction on the graph comes into being, in order to predict the change of the edge. Link prediction is also the most common application of link-relevant tasks and there are many attack related studies have been proposed [10, 13, 16, 18, 58, 70, 82, 82].

**Graph-relevant Task.** Graph-relevant tasks treat graph as a basic unit. Compared to node- or link-relevant tasks, it is macro and large-scale. The application of graph-relevant methods are more inclined to the research community of biology, chemistry, environment, materials, etc. In this task, the model tries to extract features of nodes and spatial structures to represent a graph, so as to achieve downstream operations such as graph classification or clustering. Similar with Eq.(2), the targets of attack should be an graphs, while  $y$  is determined by the specific task. For graph-relevant tasks, some attack researches have also appeared [8, 17, 39].

### 3.3 Summary: Attack on Graph

In this part, we will discuss the main contributions of current works, and point out the limitations to be overcome, and propose several open questions in this area. Specifically, we cover these released papers and its characteristic in Table 2.

**3.3.1 Major Contributions.** So far, most of the existing works in the attack scenario are mainly based on the gradients, either the adjacency matrix or feature matrix, leading to topology attack and feature attack, respectively. However, the gradients' information of the target model is hard to acquire, instead, attackers will train a surrogate model to extract the gradients. In addition to gradient-based algorithms, several heuristic methods are proposed to achieve the goal of attack, such as genetic algorithm [71] and reinforcement learning [60] based algorithms. According to different tasks in graph analysis, we summarize the major contributions of current works in the following.

**Node-relevant Task.** Most of the researches focus on the node classification task. The work [85] is the first to study the adversarial attack on graph data, using an efficient greedy search method to perform perturbations on node features and graph structure and attacking traditional graph learning model. From the perspective of gradients, there are several works [12, 20, 68, 72, 75, 86] focus on

Table 2. Related works on attack in details.

| Reference | Model  | Algorithm   | Target Task                              | Target Model                                      | Baseline   | Metric   | Dataset                                       |
|-----------|--|---|--|---|--|--|---|
| [6]       | GF-Attack  | Graph signal processing                           | Node classification                      | GCN, SGC, DeepWalk, LINE                          | Random, Degree, RL-S2V, $A_{class}$                        | Accuracy   | Cora, Citeseer, Pubmed                        |
| [72]      | IG-FGSM, IG-JSMA                                 | Gradient-based GCN                                | Node classification                      | GCN   | FGSM, JSMA, Nettack  | Classification margin, Accuracy                                  | Cora, Citeseer, PolBlogs                      |
| [9]       | EPA  | Genetic algorithm                                 | Community detection                      | GRE, INF, LOU                                     | $A_Q, A_S, A_B, A_D, D_S, D_W$                             | NMI, ARI   | Synthetic networks, Football, Email, Polblogs |
| [86]      | Meta-Self, Meta-Train                            | Gradient-Based GCN                                | Node classification                      | GCN, CLN, Deepwalk                                | DICE, Nettack, First-order                                 | Misclassification rate, Accuracy                                 | Cora-ML, Citeseer, PolBlogs, PubMed           |
| [2]       | $\mathcal{A}_{DW2}, \mathcal{A}_{DW3}$           | Gradient-based random walk                        | Node classification, Link prediction     | Deepwalk  | $\mathcal{B}_{rnd}, \mathcal{B}_{ctlg}, \mathcal{B}_{deg}$ | F1 score, Classification margin                                  | Cora, Citeseer, PolBlogs                      |
| [13]      | TGA-Tra, TGA-Gre                                 | Gradient-based DDNE                               | Link prediction                          | DDNE, cTRBM, GTRBM, dynAERNN                      | Random, DGA, CNA   | ASR, AML   | RADOSLAW, LKML, FB-WOSN                       |
| [39]      | ReWatt   | Reinforcement learning based on GCN               | Graph classification                     | GCN   | RL-S2V, Random   | ASR  | REDDIT-MULTI-12K, REDDIT-MULTI-5K, IMDB-MULTI |
| [75]      | PGD, Min-Max                                     | Gradient-based GCN                                | Node classification                      | GCN   | DICE, Meta-Self, Greedy                                    | Misclassification rate   | Cora, Citeseer                                |
| [77]      | EDA  | Genetic algorithm based on Deepwalk               | Node classification, Community detection | HOPE, LPA, EM, Deepwalk                           | Random, DICE, RLS, DBA                                     | NMI, Micro-F1, Macro-F1  | Karate, Game, Dolphin                         |
| [4]       | DAGAER   | Generative model based on VGAE                    | Node classification                      | GCN   | Nettack  | ASR  | Cora, Citeseer                                |
| [79]      | -  | Knowledge embedding                               | Fact plausibility prediction             | TransE, TransR, RESCAL                            | Random   | MRR, HR@K  | FB15k, WN18                                   |
| [65]      | -  | Based on LinLBP                                   | Node classification, Evasion             | LinLBP, JW, GCN LBP, RW, LINE, DeepWalk, Node2vec | Random, Nettack  | FNR, FPR   | Facebook, Enron, Epinions, Twitter, Google+   |
| [8]       | Q-Attack   | Genetic algorithm                                 | Community detection                      | FN, Lou, SOA, LPA, INF, Node2vec+KM               | Random, CDA, DBA   | NMI, Modularity Q  | Karate, Dolphins, Football, Polbooks          |
| [29]      | HG-attack  | Label propagation algorithm, Nodes injection      | Malware detection                        | Orig-HGC  | AN-Attack  | TP, TN, FP, FN, F1, Precision, Recall, Accuracy                  | Tencent Security Lab Dataset                  |
| [16]      | UNAttack   | Gradient-based similarity method, Nodes injection | Recommendation                           | Memory-based CF, BPRMF, NCF                       | Random, Average, Popular, Co-visitation                    | Hit@K  | Filmtrust, MovieLens, Amazon                  |
| [18]      | -  | Gradient-based GAN and MF, Nodes injection        | Recommendation                           | MF  | -  | Attack difference, TVD, JS, Est., Rank loss @K, Adversarial loss | MovieLens 100k, MovieLens 1M                  |
| [68]      | Greedy GAN                                       | Gradient-based GCN and GAN                        | Node classification                      | GCN   | Random   | Accuracy, F1 score, ASR  | Cora, Citeseer                                |
| [70]      | CTR, OTC   | Neighbor score based on graph structure           | Link prediction                          | Traditional link prediction algorithms            | -  | AUC, AP  | WTC 9/11, ScaleFree, Facebook, Random network |
| [10]      | IGA  | Gradient-based GAE                                | Link prediction                          | GAE,LRW, DeepWalk, Node2vec, CN, Random, Katz     | RAN, DICE, GA  | ASR, AML   | NS, Yeast, FaceBook                           |
| [20]      | RL-S2V   | Reinforcement learning                            | Node/Graph Classification                | GCN, GNN  | Random sampling  | Accuracy   | Citeseer, Cora, Finance, Pubmed               |
| [85]      | Nettack  | Greedy search & gradient based on GCN             | Node classification                      | GCN, CLN, Deepwalk                                | Rnd, FGSM  | Classification margin, Accuracy                                  | Cora-ML, Citeseer, PolBlogs                   |
| [12]      | FGA  | Gradient-based GCN                                | Node classification, Community detection | GCN, GraRep, DeepWalk, Node2vec, LINE, GraphGAN   | Random, DICE, Nettack                                      | ASR, AML   | Cora, Citeseer, PolBlogs                      |
| [58]      | Opt-attack                                       | Gradient-based Deepwalk and                       | Link prediction                          | DeepWalk, LINE, SC, Node2vec, GAE                 | Random, PageRank, Degree sum, Shortest path                | AP, Similarity Score   | Facebook, Cora, Citeseer                      |
| [82]      | Approx-Local                                     | Similarity methods                                | Link prediction                          | Local & global similarity metrics                 | RandomDel, GreedyBase                                      | Katz similarity, ACT distance, Similarity score                  | Random network, Facebook                      |
| [17]      | Targeted noise injection, Small community attack | Noise injection                                   | Graph clustering, Community detection    | SVD, Node2vec, Community detection algorithms     | -  | ASR, FPR   | Reverse, Engineered, DGA, Domains, NXDOMAIN   |

the topology attack, adding/removing edges between nodes based on the gradients information from various surrogate models. Specially, Xu et al. [75] present a novel optimized-based attack method that uses the gradients of surrogate model and facilitates the difficulty of tackling discrete graph data; Zügner et al. [86] use meta-gradients to solve the bilevel problem of poisoning a graph; Wang et al. [68] propose a greedy method based on Generative Adversarial Network (GAN) [25] to generate adjacency and feature matrices of fake nodes, which will be injected to a graph to misclassify the target models; Chen et al. and Dai et al. [12, 20] both use GCN as a surrogate model to extract gradients information and thus generating an adversarial graph; Wu et al. [72] argue that integrated gradients can better reflect the effect of perturbing certain features or edges. Moreover, Dai et al. [20] consider evasion attacks on the task of node classification and graph classification, and proposes two effective attack methods based on the reinforcement learning and genetic algorithms, respectively. Taking Deepwalk [49] as a base method, Bojchevski et al. [2] and Xuan et al. [77] propose an eigen decomposition and genetic algorithm based strategy to attack the network embedding, respectively. Also, Bose et al. [4] design a unified encoder-decoder framework from the generative perspective, which can be employed to attack diverse domains (images, text and graphs), Wang et al. [65] propose a threat model to manipulate the graph structure to evade detection by solving a graph-based optimization problem efficiently. Considering real-world application scenarios, Hou et al. [29] propose an algorithm that allows malware to evade detection by injecting nodes (apps).

**Link-relevant Task.** Link prediction is another fundamental research problems in network analysis. In this scenario, Waniek et al. [70] study the link connections, and propose heuristic algorithms to evade the detection by rewiring operations. Furthermore, Chen et al. [10] put forward a novel iterative gradient attack method based on a graph auto-encoder framework. Similarly, Chen et al. [13] also exploit the gradients information of surrogate model, and firstly study the works about adversarial attacks on dynamic network link prediction (DNLP). Besides, Sun et al. [58] focus on poisoning attacks and propose a unified optimization framework based on projected gradient descent. In the recommendation scenario, considering the interactions between users and items as a graph and treat it as a link prediction task, Christakopoulou et al. [18] and Chen et al. [16] propose the method of injecting fake users to degrade the recommendation performance of the system.

**Graph-relevant Task.** Few works study the adversarial attacks on this scenario, Ma et al. [39] propose a rewiring operation based algorithm, which uses reinforcement learning to learn the attack strategy on the task of graph classification. Besides, Chen et al. [8] introduce the problem of community detection and proposes a genetic algorithm based method. Chen et al. [17] focus on graph clustering and community detection scenarios, devise generic attack algorithms based on noise injection and demonstrates their effectiveness against a real-world system.

**3.3.2 Current Limitations.** Despite the remarkable achievements on attacking graph learning models, there are several limitations remain to be overcome:

- **Unnoticeability.** Most works are unaware of preserving the adversarial attacks from noticeable impacts, they simply consider the lower attack budgets but far from enough instead.
- **Scalability.** Existing works are mainly focus on a relatively small-scale graph, however, million-scale or even larger graphs are commonly seen in real life, and efficient attacks on larger graph are leaving unexplored.<sup>3</sup>

<sup>3</sup>There has been some efforts on large-scale graph computation that would be useful for graph learning methods. [78, 84]

- **Knowledge.** It is common to assume that attackers have *perfect knowledge* about the dataset, but it is unreasonable and impractical due to the limited access of attackers. Nevertheless, very few works conduct attacks with moderate or even minimal knowledge.
- **Physical Attack.** Most of the existing works conduct attacks on ideal datasets. However, in real world, attacks need to consider more factors. For instance, in a physical attack the adversarial examples as input will be distorted accidentally and it often fails to achieve the desired results. Unlike conducting attacks on ideal datasets, this brings more challenges to attackers.

## 4 DEFENSE

The proposed attack methods have made researchers realize the importance of the robustness of deep learning models. Relatively, some defense methods have also been proposed. In this section, we will give some general definitions of defense models against adversarial attack methods on graph data and its related concepts. In addition, this section systematically classifies existing defense methods and details some typical defense algorithms.

### 4.1 Definition

Simply put, the purpose of defense is to make the performance of the model still maintain a certain stability on the data that is maliciously disturbed. Although some defense models have been proposed, there is no clear and unified definition of the defense problem. To facilitate the discussion of the following, we propose a unified formulation of the defense problem.

*Defense on Graph Data.* Most symbols are the same as mentioned in Section 3, and we define  $\tilde{f}$  as a deep learning function with the loss function  $\tilde{\mathcal{L}}$  designed for defense, it receives a graph either perturbed or not. Then the defense problem can be defined as:

$$\begin{aligned} & \underset{\hat{G} \in \Psi(\mathcal{G}) \cup \mathcal{G}}{\text{minimize}} \sum_{t_i \in \mathcal{T}} \tilde{\mathcal{L}}(\tilde{f}_{\theta^*}(\hat{G}^{t_i}, X, t_i), y_i) \\ & \text{s.t. } \theta^* = \arg \min_{\theta} \sum_{v_j \in \mathcal{S}_L} \tilde{\mathcal{L}}(\tilde{f}_{\theta}(\tilde{G}^{v_j}, X, v_j), y_j) \end{aligned} \quad (9)$$

where  $\tilde{G} \in \Psi(\mathcal{G}) \cup \mathcal{G}$  can be a well-designed graph  $\hat{G}$  for the purpose of defense, or a clean graph  $G$ , which depends on whether the model have been attacked.

### 4.2 Taxonomies of Defenses

In this section, we divide the existing defense methods into several categories according to the used algorithms and describe them by examples. To the best of our knowledge, it's the first time for those defense methods to be clearly classified and the first time for those types to be clearly defined. All taxonomies are listed below:

***Preprocessing-based Defense.*** As the most intuitive way, directly manipulating the raw data has a great impact on the model performance. Additionally, preprocessing raw data is independent to model structures and training methods, which gives considerable scalability and transferability. Some existing works[72] improve model robustness by conducting certain preprocessing steps before training, and we define this type of defense methods as *Preprocessing-based Defense*. In addition, Wu et al. [72] try to drop edges that connect nodes with low similarity score, which could reduce the risk for those edges of being attack and nearly does no harm to the model performance.

**Structure-based Defense.** In addition to raw data, model structure is also crucial to the model performance. Some existing works modify the model structure, such as GCN, to gain more robustness, and we define this type of defense methods as *Structure-based Defense*. Instead of GCN's graph convolutional layers, Zhu et al. [83] use Gaussian-based graph convolutional layers, which learns node embeddings from Gaussian distributions and assign attention weight according to their variances. Wang et al. [67] propose dual-stage aggregation as a Graph Neural Network (GNN) encoder layer to learn the information of neighborhoods, and adopt GAN [25] to conduct contrastive learning. Tang et al. [62] initialize the model by meta-optimization and penalize the aggregation process of GNN. And Ioannidis et al. [32] propose a novel GCN architecture, named Adaptive Graph Convolutional Network (AGCN), to conduct robust semi-supervised learning.

**Adversarial-based Defense.** Adversarial training has been widely used in deep learning due to its excellent performance. Some researchers successfully adopt adversarial training from other fields into graph domain to improve model robustness, and we define the defense methods which use adversarial training as *Adversarial-based Defense*. There are two types of adversarial training: i) *Training with adversarial goals*. Some adversarial training methods gradually optimize the model in a continuous min-max way, under the guide of two opposite (minimize and maximize) objective functions [24, 56, 75]; ii) *Training with adversarial examples*. Other adversarial-based models are fed with adversarial samples during training, which helps the model learn to adjust to adversarial samples and thus reduces the negative impacts of those potential attack samples [11, 22, 28, 69].

**Objective-based Defense.** As a simple and effective method, modifying objective function plays an important role in improving the model robustness. Many existing works attempt to train a robust model against adversarial attacks by optimizing the objective function, and we define this type of defense methods as *Objective-based Defense*. However, there is a little intersection between the definition of *Objective-based Defense* and *Adversarial-based Defense* mentioned above, because the min-max adversarial training (the first type of *Adversarial-based Defense*) is also in the range of objective optimization. Therefore, to accurately distinguish the definition boundary between *Adversarial-based Defense* and *Objective-based Defense*, we only consider those methods which are objective-based but not adversarial-based as the instances of *Objective-based Defense*. Zügner et al. [87] and Bojchevski et al. [3] combine hinge loss and cross entropy loss to perform a robust optimization. Jin et al. [33] and Chen et al. [11] regularize the training process by studying the characteristics of graph powering and smoothing the cross-entropy loss function, respectively.

**Detection-based Defense.** Some existing works focus on the detection of adversarial attacks, or the certification of model/node robustness, which we define as *Detection-based Defense*. Although those detection-based methods are unable to improve the model robustness directly, they can serve as supervisors who keep monitoring the model security and alarm for awareness when an attack is detected. Zügner et al. [87] and Bojchevski et al. [3] propose novel methods to certificate whether a node is robust, and a robust node means it won't be affected even under the worst-case/strongest attack. Pezeshkpour et al. [50] study the impacts of adding/removing edges by performing adversarial modifications on the graph. Xu et al. [76] adopt link prediction and its variants to detect the potential risks, and for example, link with low score could be a maliciously added edge. Zhang et al. [80] utilize perturbations to explore the model structure and propose a method to detect adversarial attack. Hou et al. [29] detect the target node by uncovering the poisoning nodes injected in the heterogeneous graph [59]. Ioannidis et al. [31] effectively detect anomalous nodes in large-scale graphs by applying a graph-based random sampling and consensus strategies.

**Hybrid Defense.** As the name suggested, we define *Hybrid Defense* to denote the defense method which consists of two or more types of different defense algorithms mentioned above. Many researches flexibly combine several types of defense algorithms to achieve better performance, thus alleviate the limitations of only using a single method. As mentioned above, Zügner et al. [87] and Bojchevski et al. [3] address the certification problem of node robustness and conduct robust training with objective-based optimization. Wang et al. [67] focus on improving the model structure and adversarially train the model by GAN. Chen et al. [11] adopt adversarial training and other regularization mechanisms (e.g., gradient smoothing, loss smoothing). Xu et al. [76] propose a novel graph generation method together with other detection mechanisms (e.g., link prediction, subgraph sampling, outlier detect) as preprocessing to detect potential malicious edges. Miller et al. [43] consider a combination of the features come from graph structure and origin node attributes and use well-designed training data selection methods to do a classification task. These are instances of the hybrid defense.

**Others.** Currently, the amount of researches for defense is far less than that of attack on the graph domain. To the best of our knowledge, most existing works for defense only focus on node classification tasks, and there are a lot of opportunities to study defense methods with different tasks on graph, which will enrich our defense types. For example, Pezeshkpour et al. [50] evaluate the robustness of several link prediction models and Zhou et al. [81] adopt the heuristic approach to reduce the damage of attacks.

### 4.3 Summary: Defense on Graph

In this section, we will introduce the contributions and limitations of current works about defense. Then, we will discuss the potential research opportunities in this area. The details and comparisons of various methods are placed in Table 3.

**4.3.1 Major Contributions.** Currently, most of the methods to improve the robustness of GNNs start from two aspects: a robust training method or a robust model structure. Among them, the training methods are mostly adversarial training, and many of the model structure improvements are made using the attention mechanism. In addition, there are some studies that do not directly improve the robustness of GNNs but try to verify the robustness or try to detect the data that is disturbed. In this part, considering the different ways in which existing methods contribute to the adversarial defense of graphs, we summarize the contributions of graph defense from the following three perspectives: adversarial learning, model improvement and others.

**Adversarial Learning.** As a successful method that has shown to be effective in defending the adversarial attacks of image data and test data, adversarial learning [26] augments the training data with adversarial examples during the training stage. Some researchers have also tried to apply this kind of idea to graph defense methods. Wang et al. [67] think the vulnerabilities of graph neural networks are related to the aggregation layer and the perceptron layer. To address these two disadvantages, they propose an adversarial training framework with a modified GNN model to improve the robustness of GNNs. Chen et al. [11] propose different defense strategies based on adversarial training for target and global adversarial attack with smoothing distillation and smoothing cross-entropy loss function. Feng et al. [24] propose a method of adversarial training for the attack on node features with a graph adversarial regularizer which encourages the model to generate similar predictions on the perturbed target node and its connected nodes. Sun et al. [56] successfully transfer the efficiency of Virtual Adversarial Training (VAT) to the semi-supervised node classification task on graphs and applies some regularization mechanisms on original GCN to refine its generalization performance. Wang et al. [69] point out that the values of perturbation

Table 3. Related works on defense in details. The type is divided according to the algorithm.

| Reference | Model                                   | Algorithm  | Type                  | Target Task               | Target Model                            | Baseline                            | Metric   | Dataset                                |
|-----------|---|--|-----------------------|---------------------------|---|-------------------------------------|--|--|
| [87]      | GNN (trained with RH-U)                 | Robustness certification, Objective based        | Hybrid                | Node classification       | GNN, GCN                                | GNN (trained with CE, RCE, RH)      | Accuracy, Averaged Worst-case Margin             | Citeseer, Cora-ML, Pubmed              |
| [75]      | -                                       | Adversarial training                             | Adversarial training  | Node classification       | GCN                                     | GCN                                 | Accuracy, Misclassification rate                 | Citeseer, Cora                         |
| [33]      | r-GCN, VPN                              | Graph powering                                   | Objective based       | Node classification       | GCN                                     | ManiReg, ICA, Vanilla GCN, ...      | Accuracy, Robustness merit, Attack deterioration | Citeseer, Cora, Pubmed                 |
| [72]      | -                                       | Drop edges                                       | Preprocessing         | Node classification       | GCN                                     | GCN                                 | Accuracy, Classification margin                  | Citeseer, Cora-ML, Polblogs            |
| [67]      | DefNet                                  | GAN, GER, ACL                                    | Hybrid                | Node classification       | GCN, GraphSAGE                          | GCN, GraphSAGE                      | Classification margin                            | Citeseer, Cora, Polblogs               |
| [50]      | CRIAGE                                  | Adversarial modification                         | Robustness evaluation | Link prediction           | Knowledge graph embeddings              | -                                   | Hits@K, MRR                                      | Nations, Kinship, WN18, YAGO3-10       |
| [83]      | RGCN                                    | Gaussian-based GCN                               | Structure based       | Node classification       | GCN                                     | GCN, GAT                            | Accuracy   | Citeseer, Cora, Pubmed                 |
| [11]      | Global-AT, Target-AT, SD, SCEL          | Adversarial training, Smooth defense             | Hybrid                | Node classification       | GNN                                     | AT                                  | ADR, ACD   | Citeseer, Cora, Polblogs               |
| [56]      | SVAT, DVAT                              | VAT  | Adversarial training  | Node classification       | GCN                                     | GCN                                 | Accuracy   | Citeseer, Cora, Pubmed                 |
| [80]      | -                                       | KL divergence                                    | Detection based       | Node classification       | GCN, GAT                                | -                                   | Classification margin, Accuracy, ROC, AUC        | Citeseer, Cora, Polblogs               |
| [24]      | GCN-GATV                                | GAT, VAT   | Adversarial training  | Node classification       | GCN                                     | DeepWalk, GCN, GraphSGAN, ...       | Accuracy   | Citeseer, Cora, NELL                   |
| [22]      | S-BVAT, O-BVAT                          | BVAT   | Adversarial training  | Node classification       | GCN                                     | ManiReg, GAT, GPNN, GCN, VAT, ...   | Accuracy   | Citeseer, Cora, Pubmed, NELL           |
| [62]      | PA-GNN                                  | Penalized aggregation, Meta learning             | Structure based       | Node classification       | GNN                                     | GCN, GAT, PreProcess, RGCN, VPN     | Accuracy   | Pubmed, Reddit, Yelp-Small, Yelp-Large |
| [69]      | GraphDefense                            | Adversarial training                             | Adversarial training  | Node/Graph classification | GCN                                     | Drop edges, Discrete AT             | Accuracy   | Cora, Citeseer, Reddit                 |
| [29]      | HG-HGC                                  | HG-Defense                                       | Detection based       | Malware detection         | Malware detection system                | Other malware detection systems     | Detection rate                                   | Tencent Security Lab Dataset           |
| [32]      | AGCN                                    | Adaptive GCN with edge dithering                 | Structure based       | Node classification       | GCN                                     | GCN                                 | Accuracy   | Citeseer, Polblogs, Cora, Pubmed       |
| [31]      | GraphSAC                                | Random sampling, Consensus                       | Detection based       | Anomaly detection         | Anomaly model                           | GAE, Degree, Cut ratio, ...         | AUC  | Citeseer, Polblogs, Cora, Pubmed       |
| [3]       | GNN (train with $L_{RCE}, L_{CEM}$ )    | Robustness certification, Objective based        | Hybrid                | Node classification       | GNN                                     | GNN                                 | Accuracy, Worst-case Margin                      | Cora-ML, Citeseer, Pubmed              |
| [81]      | IDOpt, IDRank                           | Integer Program, Edge Ranking                    | Heuristic algorithm   | Link prediction           | Similarity-based link prediction models | PPN                                 | DPR  | PA, PLD, TVShow, Gov                   |
| [43]      | SVM with a radial basis function kernel | Augmented feature, Edge selecting                | Hybrid                | Node classification       | SVM                                     | GCN                                 | Classification margin                            | Cora, Citeseer                         |
| [28]      | APR, AMF                                | MF-BPR based AT                                  | Adversarial training  | Recommendation            | MF-BPR                                  | ItemPop, MF-BPR, CD4E, NeuMF, IRGAN | HR, NDCG   | Yelp, Pinterest, Gowalla               |
| [76]      | SL, OD, P+GGD, ENS,GGD                  | Link prediction, Subsampling, Neighbour analysis | Hybrid                | Node classification       | GNN, GCN                                | LP                                  | AUC  | Citeseer, Cora                         |

in adversarial training could be continuous or even negative. And they propose an adversarial training method to improve the robustness of GCN. He et al. [28] adopt adversarial training to Bayesian Personalized Ranking (BPR) [52] on recommendation by adding adversarial perturbations on embedding vectors of the user and item.

**Model Improvement.** In addition to improvements in training methods, many studies have focused on improving the model itself. In recent years, the attention mechanism [63] has shown extraordinary performance in the field of natural language processing. Some studies of graph defense have borrowed the way that they can automatically give weight to different features to reduce the influence of perturbed edges on features. Zhu et al. [83] use Gaussian distributions to absorb the effects of adversarial attacks and introduce a variance-based attention mechanism to prevent the propagation of adversarial attacks. Tang et al. [62] propose a novel framework based on a penalized aggregation mechanism which restricts the negative impact of adversarial edges. Hou et al. [29] enhance the robustness of malware detection system in Android by their well-designed defense mechanism to uncover the possible injected poisoning nodes. Jin et al. [33] point out the basic flaws of the Laplacian operator in origin GCN and propose a variable power operator to alleviate the issues. Miller et al. [43] propose to use unsupervised methods to extract the features of the graph structure and use support vector machines to complete the task of node classification. In addition, they also propose two new training data partitioning strategies.

**Others.** In addition to improving the robustness of the model, there are also some studies that have made other contributions around robustness, such as detecting disturbed edges, certifying node robustness, analyses of current attack methods and so on. Zügner et al. [87] propose the first work on certifying the robustness of GNNs. The method can give robustness certification which states whether a node is robust under a certain space of perturbations. Zhang et al. [80] study the defense methods' performance under random attacks and Nettack concluded that graph defense models which use structure exploration are more robust in general. Otherwise, this paper proposes a method to detect the perturbed edges by calculating the mean of the KL divergences [36] between the softmax probabilities of the node and its neighbors. Pezeshkpour et al. [50] introduce a novel method to automatically detect the error for knowledge graphs by conducting adversarial modifications on knowledge graphs. In addition, the method can study the interpretability of knowledge graph representations by summarizing the most influential facts for each relation. Wu et al. [72] argue that the robustness issue is rooted in the local aggregation in GCN by analyzing attacks on GCN and propose an effective defense method based on preprocessing. Bojchevski et al. [3] propose a robustness certificate method that can certify the robustness of GNNs regarding perturbations of the graph structure and the label propagation. In addition, they also propose a new robust training method guided by their own certificate method. Zhou et al. [81] model the problem of robust link prediction as a Bayesian Stackelberg game [64] and propose two defense strategies to choose which link should be protected. Ioannidis et al. [32] propose a novel edge-dithering (ED) approach reconstructs the original neighborhood structure with high probability as the number of sampled graphs increases. Ioannidis et al. [31] introduce a graph-based random sampling and consensus approach to effectively detect anomalous nodes in large-scale graphs.

**4.3.2 Current Limitations.** As a new-rising branch that has not been sufficiently studied, current defense methods have several major limitations as follows:

- **Diversity.** At present, most works of defense mainly focus on node classification tasks only. From the perspective of defender, they need to improve their model robustness on different tasks.

- **Scalability.** Whether a defense model can be widely used in practice largely depends on the cost of model training. Most existing works lack the consideration of training costs.
- **Theoretical Proof.** Most of current methods only illustrate their effectiveness by showing experimental results and textual descriptions. It will be great if a new robust method's effectiveness can be proved theoretically.

## 5 METRICS

In this section, we first introduce the metrics that are common in graph analysis tasks, which are used in attack and defense scenarios as well. Next, we introduce some new metrics proposed in attack and defense works from three perspectives: effectiveness, efficiency, and imperceptibility.

### 5.1 Common Metrics

- **FNR and FPR.** In the classification or clustering tasks, here is a category of metrics based on False Positive (*FP*), False Negative (*FN*), True Positive (*TP*) and True Negative (*TN*). In attack and defense scenarios, commonly used to quantify are False Negative Rate (*FNR*) and False Positive Rate (*FPR*). Specifically, *FNR* is equal to  $FN / (TP + FN)$ , and *FPR* is equal to  $FP / (FP + TN)$ . For all negative samples, the former describes the proportion of false positive samples detected, while the latter describes the proportion of false negative samples detected. In other words, for negative samples, the former is the error rate, and the latter is the miss rate.
- **Accuracy (Acc).** The Accuracy is one of the most commonly used evaluation metrics, which measures the quality of results based on the percentage of correct predictions over total instances.
- **F1-score.** The F1 score [2] can be regarded as a harmonic average of model Precision and Recall score [51], with a maximum value of 1 and a minimum value of 0.
- **Area Under Curve (AUC).** The AUC is the area under the Receiver Operating Characteristic (ROC) curve [80]. Different models corresponding to different ROC curves, and it is difficult to compare which one is better if there is a crossover between the curves, thus comparison based on the AUC score is more reasonable. In general, AUC indicates the probability that the predicted positive examples rank in front of the negative examples.
- **Average Precision (AP).** As the area under the Precision-Recall curve, AP is described as one of the most robust metrics [5]. Precision-Recall curve depicts the relationships between Precision and Recall. A good model should improve the Recall while preserving the Precision a relatively high score. In contrast, weaker models may lose more Precision in order to improve Recall. Comparing with the Precision-Recall curve, AP can show the performance of the model more intuitively.
- **Mean Reciprocal Rank (MRR).** The MRR is a commonly used metric to measure a rank model. For a target query, if the first correct item is ranked  $n_{th}$ , then the MRR score is  $1/n$ , and once there is no match, the score is 0. The MRR of the model is the sum of the scores of all queries.
- **Hits@K.** By calculating the rank (e.g., MRR) of all the ground-truth triples, Hits@K is the proportion of correct entities ranked in top K.
- **Modularity.** The Modularity is an important measure based on the assortative mixing [44], which usually used to assess the quality of different division for a particular network, especially the community structure is unknown [45].
- **Normalized Mutual Information (NMI).** The NMI is another commonly used evaluation to measure the quality of clustering results, so as to analyze the network community

structure [21]. NMI further indicates the similarity between two partitions with mutual information and information entropy, while a larger value means a higher similarity.

- **Adjusted Rand Index (ARI).** The ARI is a measure of the similarity between two data clusterings: one given by the clustering process and the other defined by external criteria [53]. Such a correction establishes a baseline by using the expected similarity of all pair-wise comparisons between clusterings specified by a random model. ARI measures the relation between pairs of dataset elements without labels' information, which can cooperate with conventional performance measures to detect classification algorithm. In addition, ARI can also use labels information for feature selection [53].

## 5.2 Specific Metrics for Attack and Defense

In order to measure the performance, a number of specific metrics for attack and defense appear in literature. Next, we will organize these metrics from three perspectives: effectiveness, efficiency, and imperceptibility.

*5.2.1 Effectiveness-relevant Metrics.* Both attack and defense require metrics to measure the performance of the target model before and after the attack. Therefore, we have summarized some metrics for measuring effectiveness and list them below:

- **Attack Success Rate (ASR).** ASR is the ratio of targets which will be successfully attacked within a given fixed budget [10]. Correspondingly, we can conclude the formulation of ASR:

$$\text{ASR} = \frac{\text{Number of successful attacks}}{\text{Number of attacks}} \quad (10)$$

- **Classification Margin (CM).** CM is only designed for the classification task. Under this scenario, attackers aim to perturb a graph that misclassifies the target node and has maximal "distance" (in terms of log-probabilities/logits) to the correct class [85]. Based on this, CM is well formulated as:

$$\text{CM} = \max_{y'_t \neq y_t} Z_{t, y'_t} - Z_{t, y_t} \quad (11)$$

where  $Z_t$  is the output of the target model with respect to node  $t$ , and  $y_t \in C$  is the *correct* class label of  $v_t$  while  $y'_t$  is the *wrong* one.

- **Averaged Worst-case Margin (AWM).** Worst-case Margin (WM) is the minimum of the Classification Margin (CM) [87], and the average of WM is dynamically calculated over a mini-batch of nodes during training. The  $B_s$  in the following equation denotes batch size.

$$\text{AWM} = \frac{1}{B_s} \sum_{i=1}^{i=B_s} \text{CM}_{\text{worst}_i} \quad (12)$$

- **Robustness Merit (RM).** It evaluates the robustness of the model by calculating the difference between the post-attack accuracy and pre-attack accuracy on GNNs [33]. And  $\mathcal{V}^{\text{attacked}}$  denotes the set of nodes theoretically affected by the attack.

$$\text{RM} = \text{Acc}_{\mathcal{V}^{\text{attacked}}}^{\text{post-attacked}} - \text{Acc}_{\mathcal{V}^{\text{attacked}}}^{\text{pre-attacked}} \quad (13)$$

- **Attack Deterioration (AD).** It evaluates the attack effect of the model prediction. Note that any added/dropped edges can only affect the nodes within the spatial scope in the

origin network [33], due to the spatial scope limitations of the GNNs.

$$AD = 1 - \frac{\text{Acc}_{\mathcal{V}^{\text{attacked}}}^{\text{post-attacked}}}{\text{Acc}_{\mathcal{V}^{\text{attacked}}}^{\text{pre-attacked}}} \quad (14)$$

- **Average Defense Rate (ADR).** ADR is the ratio of the difference between the ASR of attack on GNNs with and without defense, versus the ASR of attack on GNNs without defense [11]. The higher ADR the better defense performance.

$$ADR = \frac{\text{ASR}_{\mathcal{V}^{\text{attacked}}}^{\text{with-defense}}}{\text{ASR}_{\mathcal{V}^{\text{attacked}}}^{\text{without-defense}}} - 1 \quad (15)$$

- **Average Confidence Different (ACD).** ACD is the average confidence different of nodes in  $N_s$  before and after attack [11], where  $N_s$  is the set of nodes which are classified correctly before attack in test set. Note that Confidence Different (CD) is equivalent to Classification Margin (CM).

$$ACD = \frac{1}{N_s} \sum_{v_i \in N_s} \text{CD}_i(\hat{A}_{v_i}) - \text{CD}_i(A) \quad (16)$$

- **Damage Prevention Ratio (DPR).** Damage prevention is defined to measure the amount of damage that can be prevented by defense [81]. Let  $L_0$  be the defender's loss without attack,  $L_A$  be the defender's loss under a certain attack strategy,  $A$  and  $L_D$  be the defender's loss with a certain defense strategy  $D$ . A better defense strategy leads to a larger DPR.

$$\text{DPR}_A^D = \frac{L_A - L_D}{L_A - L_0} \quad (17)$$

**5.2.2 Efficiency-relevant Metrics.** Here we introduce some efficiency metrics which measure the cost of the attack and defense. For example, the metric of quantifying how much perturbations are required for the same effect.

- **Average Modified Links (AML).** AML is designed for the topology attack, which indicates the average perturbation size leading to a successful attack [10]. Assume that attackers have limited budgets to attack a target model, the modified links (added or removed) are accumulated until attackers achieve their goal or run out of the budgets. Based on this, we can conclude the formulation of AML:

$$\text{AML} = \frac{\text{Number of modified links}}{\text{Number of attacks}} \quad (18)$$

**5.2.3 Imperceptibility-relevant Metrics.** Several metrics are proposed to measure the scale of the manipulations caused by attack and defense, which are summarized as follows:

- **Similarity Score.** Generally speaking, Similarity Score is a measure to infer the likelihoods of the existence of links, which usually applied in the link-level task [38, 82]. Specifically, suppose we want to figure out whether a particular link  $(u, v)$  exists in a network, Similarity Score can be used to quantify the topology properties of node  $u$  and  $v$  (e.g., common neighbors, degrees), and a higher Similarity Score indicates a greater likelihood of connection between this pair. Usually, Similarity Score could be measured by the cosine similarity matrix [58].
- **Test Statistic  $\Lambda$ .**  $\Lambda$  takes advantage of the distribution of node degrees to measure the similarity between graphs. Specifically, Zügner et al. [85] state that multiset of node degrees  $d_G$  in the graph  $G$  follow the power-law like distribution, and they provide a method to approximate the main parameter  $\alpha_G$  of the distribution. Through  $\alpha_G$  and  $d_G$ , we can

calculate  $G$ 's log-likelihood score  $l(d_G, \alpha_G)$ . For the pre-attack graph  $G_{pr}$  and post-attack graph  $G_{po}$ , we get  $(d_{G_{pr}}, \alpha_{G_{pr}})$  and  $(d_{G_{po}}, \alpha_{G_{po}})$ , respectively. Similarly, we can then define  $d_{G_q} = d_{G_{pr}} \cap d_{G_{po}}$  and estimate  $\alpha_{G_q}$ . The final test statistic of graphs can be formulated as:

$$\Lambda(G_{pr}, G_{po}) = 2 \cdot (-l(d_{G_q}, \alpha_{G_q}) + l(d_{G_{po}}, \alpha_{G_{po}}) + l(d_{G_{pr}}, \alpha_{G_{pr}})) \quad (19)$$

Finally, when the statistic  $\Lambda$  satisfies a specified constraint, the model considers the perturbations are unnoticeable.

- **Attack Budget  $\Delta$ .** To ensure unnoticeable perturbations in the attack, previous work [85] measures the perturbations in terms of the budget  $\Delta$  and the test statistics  $\Lambda$  w.r.t. log-likelihood score. More precisely, they accumulate the changes in node feature and the adjacency matrix, and limit it to a budget  $\Delta$  to constrain perturbations. However, it is not suitable to deal with complicated situations [85].
- **Attack Effect.** Generally, the attack effect evaluates the impacts in the community detection task. Given a confusion matrix  $J$ , where each element  $J_{i,j}$  represents the number of shared nodes between original community and a perturbed one, the attack effect is simply defined as the accumulation of normalized entropy [9]. Suppose all the communities keep exactly the same after the attack, the attack effect will be equal to 0.

## 6 OPEN PROBLEMS

Graph adversarial learning has many problems worth to be studied by observing existing researches. In this section, we try to introduce several major research issues and discuss the potential solutions. Generally, we will discuss the open problems in three parts: attack, defense and evaluation metric.

### 6.1 Attack

We mainly approach the problems from three perspectives, attacks' side effects, performance, and prerequisites. First, we hope the inevitable disturbances in the attack can be stealthy and unnoticeable. Second, now that data is huge, attacks need to be efficient and scalable. Finally, the premise of the attack should not be too ideal, that is, it should be practical and feasible. Based on the limitations of previous works detailed in Section 3.3.2 and the discussion above, we rise several open problems in the following:

- **Unnoticeable Attack.** Adversaries want to keep the attack unnoticeable to avoid censorship. To this end, Zügner et al. [86] argue that the main property of the graph should be marginally changed after attack, e.g., attackers are restricted to maintain the degree distribution while attacking a specified graph, which is evaluated via a test static. Nevertheless, most of existing works ignore such a constraint despite achieving outstanding performance. To be more concealed, we believe that attackers should focus more on their perturbation impacts on a target graph, with more constraints placed on the attacks.
- **Efficient and Effective Algorithms.** Seeing from Table 2, the most frequently used benchmark datasets are rather small-scale graphs, e.g., Cora [42] and Citeseer [54]. Currently, the majority of proposed methods are failed to attack a large-scale graph, for the reason that they need to store the entire information of the graph to compute the surrogate gradients or loss even with small changes, leading to the expensive computation and memory consumption. To address this problem, Zügner et al. [86] make an effort that derives an incremental update for all candidate set (potential perturbation edges) with a linear variant of GCN, which avoid the redundant computation while remain efficient attack performance. Unfortunately, the proposed method isn't suitable for a larger-scale graph

yet. With the shortcoming that unable to conduct attacks on a larger-scale graph, a more efficient and effective attack method should be studied to address this practical problem. For instance, given a target node, the message propagation are only exits within its  $k$ -hops neighborhood.<sup>4</sup> In other words, we believe that adversaries can explore a heuristic method to attack an entire graph (large) via its represented subgraph (small), and naturally the complexity of time and space are reduced.

- **Constraints on Attackers' Knowledge.** To conduct a real and practical black-box attack, attackers should not only be blind to the model knowledge, but also be strict to minimal data knowledge. This is a challenging setting from the perspective of attackers. Chen et al. [17] explore this setting by randomly sampling different size of subgraphs, showing that the attacks are failed with a small network size, and the performance increases as it gets larger. However, few works are aware of the constraints on attackers' knowledge, instead, they assume that perfect knowledge of the input data are available, and naturally perform well in attacking a fully "exposed" graph. Therefore, several works could be explored with such a strict constraint, e.g., attackers can train a substitute model on a certain part of the graph, learn the general patterns of conducting attacks on graph data, thus transfer to other models and entire graph with less prior knowledge.
- **Real-world Attack.** Attacks can't always be on an ideal assumption, i.e., studying attacks on a simple and static graphs. Such ideal, undistorted input is unrealistic in real cases. In other words, how can we improve existing models so that they can work in complex real-world environments?

## 6.2 Defense

Besides the open problems of attack mentioned above, there are still many interesting problems in defense on the graph, which deserve more attentions. Here we list some open problems that worth discussing:

- **Defense on Various Tasks.** From table 3 we can see that, most existing works of defense [33, 67, 72, 75, 87] are focusing on the node classification task on graph, achieving promising performance. Except for node classification, there are various tasks on graph domain are important and should be pay more attention to. However, only a few works [50, 81] try to investigate and improve model robustness on link prediction task on graph while a few works [28, 29, 31] make an effort to defense on some tasks (e.g., malware detection, anomaly detection, recommendation) on other graph domains. Therefore, it is valuable to think about how to improve model robustness on various tasks or transfer the current defense methods to other tasks.
- **Efficient Defense.** Training cost is critical important factor to be taken into consideration under the industrial scenes. Therefore, it's worth to be studied how to guarantee acceptable cost while improving robustness. However, the currently proposed defense methods have rarely discussed the space-time efficiency of their algorithms. Wang et al. [67] design a bottleneck mapping to reduce the dimension of input for a better efficiency, but their experiments absence the consideration of large-scale graph. Zügner et al. [87] test the number of coverage iterations of their certification methods under the different number of nodes, but the datasets they used are still small. Wang et al. [69] do not restrict the regularized adjacency matrix to be discrete during the training process, which improves the training efficiency.

<sup>4</sup>Here  $k$  depends on the specified method of message propagation and is usually of a smaller value, e.g.,  $k = 2$  with a two-layer GCN.

- **Certification on Robustness.** Robustness of the graph model is always the main issue among all the existing works including attack and defense mentioned above. However, most researchers only focus on improving model (e.g., GCN) robustness, and try to prove the model robustness via the performance results of their model. The experimental results can indeed prove the robustness of the model to a certain extent. However, considering that the model performance is easily influenced by the hyperparameters, implementation method, random initialization, and even some hardware constraints (e.g., cpu, gpu, memory), it is difficult to guarantee that the promising and stable performance can be reobtained on different scenarios (e.g., datasets, tasks, parameter settings). There is a novel and cheaper way to prove robustness via certification that should be taken more seriously into consideration, which certifies the node's absolute robustness under arbitrary perturbations, but currently only a few works have paid attention to the certification of GNNs' robustness [3, 87] which are also a valuable research direction.

### 6.3 Evaluation Metric

As seen in Section 5, we got so many evaluation metrics in graph adversarial learning field, however, the number of effectiveness-relevant metrics are far more than the other two, which reflects the research emphasis on model performance is unbalanced. Therefore, we propose several potential works can be further studied:

- **Measurement of Cost.** There are not many metrics to explore the model's efficiency which reflects the lack of research attention on it. The known methods roughly measure the cost via the number of modified edges which begs the question that is there another perspective to quantify the cost of attack and defense more precisely? In real world, the cost of adding an edge is rather different from removing one,<sup>5</sup> thus the cost between modified edges are unbalanced, which needs more reasonable evaluation metrics.
- **Measurement of Imperceptibility.** Different from the image data, where the perturbations are bounded in  $\ell_p$ -norm [26, 40, 57], and could be used as an evaluation metric on attack impacts, but it is ill-defined on graph data. Therefore, how to better measure the effect of perturbations and make the attack more unnoticeable?
- **Appropriate Metric Selection.** With so many metrics proposed to evaluate the performance of the attack and defense algorithms, here comes a problem that how to determine the evaluation metrics in different scenarios?

## 7 CONCLUSION

In this survey, we conduct a comprehensive review on graph adversarial learning, including attacks, defenses and corresponding evaluation metrics. To the best of our knowledge, this is the first work systemically summarize the extensive works in this field. Specifically, we present the recent developments of this area and introduce the arms race between attackers and defenders. Besides, we provide a reasonable taxonomy for both of them, and further give a unified problem formulation which makes it clear and understandable. Moreover, under the scenario of graph adversarial learning, we summarize and discuss the major contributions and limitations of current attack and defense methods respectively, along with the open problems of this area that worth exploring. Our works cover most of the relevant evaluation metrics in the graph adversarial learning field as well, aiming to provide a better understanding on these methods. In the future, we will measure the performance of proposed models with relevant metrics based on extensive empirical studies.

<sup>5</sup>Usually, removing an edge is more expensive than adding one.

Hopefully, our works will serve as a reference and give researchers a comprehensive and systematic understanding of the fundamental issues, thus become a well starting point to study in this field.

## ACKNOWLEDGMENTS

The paper was supported by the National Natural Science Foundation of China (61702568, U1711267), Key-Area Research and Development Program of Guangdong Province (2020B010165003), the Program for Guangdong Introducing Innovative and Entrepreneurial Teams (2017ZT07X355) and the Fundamental Research Funds for the Central Universities under Grant (17lgpy117).

## REFERENCES

- [1] Battista Biggio and Fabio Roli. 2018. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition* 84 (2018), 317–331.
- [2] Aleksandar Bojchevski and Stephan Günnemann. 2019. Adversarial Attacks on Node Embeddings via Graph Poisoning. In *International Conference on Machine Learning*. 695–704.
- [3] Aleksandar Bojchevski and Stephan Günnemann. 2019. Certifiable Robustness to Graph Perturbations. In *Advances in Neural Information Processing Systems*. 8317–8328.
- [4] Avishek Joey Bose, Andre Cianflone, and William Hamilton. 2019. Generalizable Adversarial Attacks Using Generative Models. *arXiv preprint arXiv:1905.10864* (2019).
- [5] Kendrick Boyd, Kevin H Eng, and C David Page. 2013. Area under the precision-recall curve: point estimates and confidence intervals. In *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 451–466.
- [6] Heng Chang, Yu Rong, Tingyang Xu, Wenbing Huang, Honglei Zhang, Peng Cui, Wenwu Zhu, and Junzhou Huang. 2019. A Restricted Black-box Adversarial Framework Towards Attacking Graph Embedding Models. *arXiv:cs.SI/1908.01297*
- [7] Gary Chartrand. 1977. *Introductory graph theory*. Courier Corporation.
- [8] Jinyin Chen, Lihong Chen, Yixian Chen, Minghao Zhao, Shanjing Yu, Qi Xuan, and Xiaoniu Yang. 2019. GA-Based Q-Attack on Community Detection. *IEEE Transactions on Computational Social Systems* 6, 3 (2019), 491–503.
- [9] Jinyin Chen, Yixian Chen, Lihong Chen, Minghao Zhao, and Qi Xuan. 2019. Multiscale Evolutionary Perturbation Attack on Community Detection. *arXiv preprint arXiv:1910.09741* (2019).
- [10] Jinyin Chen, Ziqiang Shi, Yangyang Wu, Xuanheng Xu, and Haibin Zheng. 2018. Link prediction adversarial attack. *arXiv preprint arXiv:1810.01110* (2018).
- [11] Jinyin Chen, Yangyang Wu, Xiang Lin, and Qi Xuan. 2019. Can Adversarial Network Attack be Defended? *arXiv preprint arXiv:1903.05994* (2019).
- [12] Jinyin Chen, Yangyang Wu, Xuanheng Xu, Yixian Chen, Haibin Zheng, and Qi Xuan. 2018. Fast gradient attack on network embedding. *arXiv preprint arXiv:1809.02797* (2018).
- [13] Jinyin Chen, Jian Zhang, Zhi Chen, Min Du, Feifei Li, and Qi Xuan. 2019. Time-aware Gradient Attack on Dynamic Network Link Prediction. *arXiv preprint arXiv:1911.10561* (2019).
- [14] Liang Chen, Yang Liu, Xiangnan He, Lianli Gao, and Zibin Zheng. 2019. Matching User with Item Set: Collaborative Bundle Recommendation with Deep Attention Network. In *IJCAI*. 2095–2101.
- [15] Liang Chen, Yang Liu, Zibin Zheng, and Philip Yu. 2018. Heterogeneous Neural Attentive Factorization Machine for Rating Prediction. In *CIKM*. ACM, 833–842.
- [16] Liang Chen, Yangjun Xu, Fenfang Xie, Min Huang, and Zibin Zheng. 2019. Data Poisoning Attacks on Neighborhood-based Recommender Systems. *arXiv:cs.IR/1912.04109*
- [17] Yizheng Chen, Yacin Nadji, Athanasios Kountouras, Fabian Monrose, Roberto Perdisci, Manos Antonakakis, and Nikolaos Vasiloglou. 2017. Practical attacks against graph-based clustering. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 1125–1142.
- [18] Konstantina Christakopoulou and Arindam Banerjee. 2018. Adversarial recommendation: Attack of the learned fake users. *arXiv preprint arXiv:1809.08336* (2018).
- [19] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *Journal of machine learning research* 12, Aug (2011), 2493–2537.
- [20] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. 2018. Adversarial attack on graph structured data. *arXiv preprint arXiv:1806.02371* (2018).
- [21] Leon Danon, Albert Diaz-Guilera, Jordi Duch, and Alex Arenas. 2005. Comparing community structure identification. *Journal of Statistical Mechanics: Theory and Experiment* 2005, 09 (2005), P09008.

- [22] Zhijie Deng, Yinpeng Dong, and Jun Zhu. 2019. Batch Virtual Adversarial Training for Graph Convolutional Networks. *arXiv preprint arXiv:1902.09192* (2019).
- [23] Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. 2019. Deep Anomaly Detection on Attributed Networks. In *Proceedings of the 2019 SIAM International Conference on Data Mining*. SIAM, 594–602.
- [24] Fuli Feng, Xiangnan He, Jie Tang, and Tat-Seng Chua. 2019. Graph Adversarial Training: Dynamically Regularizing Based on Graph Structure. *arXiv preprint arXiv:1902.08226* (2019).
- [25] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Networks. *arXiv:stat.ML/1406.2661*
- [26] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [27] Palash Goyal and Emilio Ferrara. 2018. Graph embedding techniques, applications, and performance: A survey. *Knowledge-Based Systems* 151 (2018), 78–94.
- [28] Xiangnan He, Zhankui He, Xiaoyu Du, and Tat-Seng Chua. 2018. Adversarial personalized ranking for recommendation. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. ACM, 355–364.
- [29] Shifu Hou, Yujie Fan, Yiming Zhang, Yanfang Ye, Jingwei Lei, Wenqiang Wan, Jiabin Wang, Qi Xiong, and Fudong Shao. 2019.  $\alpha$ Cyber: Enhancing Robustness of Android Malware Detection System against Adversarial Attacks on Heterogeneous Graph based Model. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. ACM, 609–618.
- [30] Rana Hussein, Dingqi Yang, and Philippe Cudré-Mauroux. 2018. Are Meta-Paths Necessary?: Revisiting Heterogeneous Graph Embeddings. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 437–446.
- [31] Vassilis N Ioannidis, Dimitris Berberidis, and Georgios B Giannakis. 2019. GraphSAC: Detecting anomalies in large-scale graphs. *arXiv preprint arXiv:1910.09589* (2019).
- [32] Vassilis N Ioannidis and Georgios B Giannakis. 2019. Edge Dithering for Robust Adaptive Graph Convolutional Networks. *arXiv preprint arXiv:1910.09590* (2019).
- [33] Ming Jin, Heng Chang, Wenwu Zhu, and Somayeh Sojoudi. 2019. Power up! Robust Graph Convolutional Network against Evasion Attacks based on Graph Powering. *arXiv preprint arXiv:1905.10029* (2019).
- [34] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations (ICLR)*.
- [35] Balaji Krishnamurthy and Mausoom Sarkar. 2018. Deep-learning network architecture for object detection. US Patent 10,152,655.
- [36] Solomon Kullback and Richard A Leibler. 1951. On information and sufficiency. *The annals of mathematical statistics* 22, 1 (1951), 79–86.
- [37] Ziqi Liu, Chaochao Chen, Xinxing Yang, Jun Zhou, Xiaolong Li, and Le Song. 2018. Heterogeneous graph neural networks for malicious account detection. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 2077–2085.
- [38] Linyuan Lü and Tao Zhou. 2011. Link prediction in complex networks: A survey. *Physica A: statistical mechanics and its applications* 390, 6 (2011), 1150–1170.
- [39] Yao Ma, Suhang Wang, Lingfei Wu, and Jiliang Tang. 2019. Attacking Graph Convolutional Networks via Rewiring. *arXiv preprint arXiv:1906.03750* (2019).
- [40] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
- [41] Franco Manessi, Alessandro Rozza, and Mario Manzo. 2020. Dynamic graph convolutional networks. *Pattern Recognition* 97 (2020), 107000.
- [42] Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. 2000. Automating the construction of internet portals with machine learning. *Information Retrieval* 3, 2 (2000), 127–163.
- [43] Benjamin A Miller, Mustafa Çamurcu, Alexander J Gomez, Kevin Chan, and Tina Eliassi-Rad. 2019. Improving Robustness to Attacks Against Vertex Classification. (2019).
- [44] Mark EJ Newman. 2003. Mixing patterns in networks. *Physical Review E* 67, 2 (2003), 026126.
- [45] Mark EJ Newman and Michelle Girvan. 2004. Finding and evaluating community structure in networks. *Physical review E* 69, 2 (2004), 026113.
- [46] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*. ACM, 506–519.
- [47] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, et al. 2015. Deep face recognition. In *bmvc*, Vol. 1. 6.
- [48] Yulong Pei, Nilanjan Chakraborty, and Katia Sycara. 2015. Nonnegative matrix tri-factorization with graph regularization for community detection in social networks. In *Twenty-Fourth International Joint Conference on Artificial*

*Intelligence.*

- [49] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *KDD*. ACM, 701–710.
- [50] Pouya Pezeshkpour, CA Irvine, Yifan Tian, and Sameer Singh. 2019. Investigating Robustness and Interpretability of Link Prediction via Adversarial Modifications. In *Proceedings of NAACL-HLT*. 3336–3347.
- [51] David Martin Powers. 2011. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. (2011).
- [52] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of the twenty-fifth conference on uncertainty in artificial intelligence*. AUAI Press, 452–461.
- [53] Jorge M Santos and Mark Embrechts. 2009. On the use of the adjusted rand index as a metric for evaluating supervised classification. In *International conference on artificial neural networks*. Springer, 175–184.
- [54] 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.
- [55] Yash Sharma. 2018. *Gradient-based Adversarial Attacks to Deep Neural Networks in Limited Access Settings*. Ph.D. Dissertation. COOPER UNION.
- [56] Ke Sun, Zhouchen Lin, Hantao Guo, and Zhanxing Zhu. 2019. Virtual adversarial training on graph convolutional networks in node classification. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*. Springer, 431–443.
- [57] Lichao Sun, Ji Wang, Philip S Yu, and Bo Li. 2018. Adversarial attack and defense on graph data: A survey. *arXiv preprint arXiv:1812.10528* (2018).
- [58] 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 preprint arXiv:1810.12881* (2018).
- [59] Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S Yu, and Tianyi Wu. 2011. Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. *Proceedings of the VLDB Endowment* 4, 11 (2011), 992–1003.
- [60] Richard S Sutton and Andrew G Barto. 2018. *Reinforcement learning: An introduction*. MIT press.
- [61] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* (2013).
- [62] Xianfeng Tang, Yandong Li, Yiwei Sun, Huaxiu Yao, Prasenjit Mitra, and Suhang Wang. 2019. Robust graph neural network against poisoning attacks via transfer learning. *arXiv preprint arXiv:1908.07558* (2019).
- [63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [64] Heinrich Von Stackelberg. 2010. *Market structure and equilibrium*. Springer Science & Business Media.
- [65] Binghui Wang and Neil Zhenqiang Gong. 2019. Attacking Graph-based Classification via Manipulating the Graph Structure. *arXiv preprint arXiv:1903.00553* (2019).
- [66] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-scale commodity embedding for e-commerce recommendation in alibaba. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 839–848.
- [67] Shen Wang, Zhengzhang Chen, Jingchao Ni, Xiao Yu, Zhichun Li, Haifeng Chen, and Philip S Yu. 2019. Adversarial Defense Framework for Graph Neural Network. *arXiv preprint arXiv:1905.03679* (2019).
- [68] Xiaoyun Wang, Joe Eaton, Cho-Jui Hsieh, and Felix Wu. 2018. Attack graph convolutional networks by adding fake nodes. *arXiv preprint arXiv:1810.10751* (2018).
- [69] Xiaoyun Wang, Xuanqing Liu, and Cho-Jui Hsieh. 2019. GraphDefense: Towards Robust Graph Convolutional Networks. *arXiv preprint arXiv:1911.04429* (2019).
- [70] Marcin Waniek, Kai Zhou, Yevgeniy Vorobeychik, Esteban Moro, Tomasz P Michalak, and Talal Rahwan. 2018. Attack Tolerance of Link Prediction Algorithms: How to Hide Your Relations in a Social Network. *arXiv preprint arXiv:1809.00152* (2018).
- [71] Darrell Whitley. 1994. A genetic algorithm tutorial. *Statistics and computing* 4, 2 (1994), 65–85.
- [72] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. 2019. Adversarial examples for graph data: deep insights into attack and defense. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 4816–4823.
- [73] Fenfang Xie, Liang Chen, Yongjian Ye, Zibin Zheng, and Xiaola Lin. 2018. Factorization Machine Based Service Recommendation on Heterogeneous Information Networks. In *ICWS*. IEEE, 115–122.
- [74] Wayne Xiong, Jasha Droppo, Xuedong Huang, Frank Seide, Mike Seltzer, Andreas Stolcke, Dong Yu, and Geoffrey Zweig. 2016. Achieving human parity in conversational speech recognition. *arXiv preprint arXiv:1610.05256* (2016).
- [75] 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. *arXiv preprint arXiv:1906.04214* (2019).

- [76] Xiaojun Xu, Yue Yu, Bo Li, Le Song, Chengfeng Liu, and Carl Gunter. 2018. Characterizing Malicious Edges targeting on Graph Neural Networks. (2018).
- [77] Qi Xuan, Jun Zheng, Lihong Chen, Shanqing Yu, Jinyin Chen, Dan Zhang, and Qingpeng Zhang Member. 2019. Unsupervised Euclidean Distance Attack on Network Embedding. *arXiv preprint arXiv:1905.11015* (2019).
- [78] Ke Yang, MingXing Zhang, Kang Chen, Xiaosong Ma, Yang Bai, and Yong Jiang. 2019. KnightKing: a fast distributed graph random walk engine. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*. 524–537.
- [79] Hengtong Zhang, Tianhang Zheng, Jing Gao, Chenglin Miao, Lu Su, Yaliang Li, and Kui Ren. 2019. Data poisoning attack against knowledge graph embedding. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 4853–4859.
- [80] Yingxue Zhang, S Khan, and Mark Coates. 2019. Comparing and detecting adversarial attacks for graph deep learning. In *Proc. Representation Learning on Graphs and Manifolds Workshop, Int. Conf. Learning Representations, New Orleans, LA, USA*.
- [81] Kai Zhou, Tomasz P Michalak, and Yevgeniy Vorobeychik. 2019. Adversarial robustness of similarity-based link prediction. *arXiv preprint arXiv:1909.01432* (2019).
- [82] Kai Zhou, Tomasz P Michalak, Marcin Waniek, Talal Rahwan, and Yevgeniy Vorobeychik. 2018. Attacking Similarity-Based Link Prediction in Social Networks. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 305–313.
- [83] Dingyuan Zhu, Ziwei Zhang, Peng Cui, and Wenwu Zhu. 2019. Robust Graph Convolutional Networks Against Adversarial Attacks. (2019).
- [84] Xiaowei Zhu, Wenguang Chen, Weimin Zheng, and Xiaosong Ma. 2016. Gemini: A computation-centric distributed graph processing system. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. 301–316.
- [85] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. 2018. Adversarial attacks on neural networks for graph data. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 2847–2856.
- [86] Daniel Zügner and Stephan Günnemann. 2019. Adversarial attacks on graph neural networks via meta learning. *arXiv preprint arXiv:1902.08412* (2019).
- [87] Daniel Zügner and Stephan Günnemann. 2019. Certifiable robustness and robust training for graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 246–256.