

A Hard Label Black-box Adversarial Attack Against Graph Neural Networks

Jiaming Mu¹, Binghui Wang², Qi Li¹, Kun Sun³, Mingwei Xu¹, Zhuotao Liu¹

¹Institute for Network Sciences and Cyberspace, Department of Computer Science, and BNRist, Tsinghua University

²Illinois Institute of Technology ³George Mason University

{mujm19@mails, qli01@, xumw@, zhuotaoliu@}tsinghua.edu.cn, bwang70@iit.edu, ksun3@gmu.edu

ABSTRACT

Graph Neural Networks (GNNs) have achieved state-of-the-art performance in various graph structure related tasks such as node classification and graph classification. However, GNNs are vulnerable to adversarial attacks. Existing works mainly focus on attacking GNNs for node classification; nevertheless, the attacks against GNNs for graph classification have not been well explored.

In this work, we conduct a systematic study on adversarial attacks against GNNs for graph classification via perturbing the graph structure. In particular, we focus on the most challenging attack, i.e., *hard label black-box* attack, where an attacker has no knowledge about the target GNN model and can only obtain predicted labels through querying the target model. To achieve this goal, we formulate our attack as an optimization problem, whose objective is to minimize the number of edges to be perturbed in a graph while maintaining the high attack success rate. The original optimization problem is intractable to solve, and we relax the optimization problem to be a tractable one, which is solved with theoretical convergence guarantee. We also design a coarse-grained searching algorithm and a query-efficient gradient computation algorithm to decrease the number of queries to the target GNN model. Our experimental results on three real-world datasets demonstrate that our attack can effectively attack representative GNNs for graph classification with less queries and perturbations. We also evaluate the effectiveness of our attack under two defenses: one is well-designed adversarial graph detector and the other is that the target GNN model itself is equipped with a defense to prevent adversarial graph generation. Our experimental results show that such defenses are not effective enough, which highlights more advanced defenses.

CCS CONCEPTS

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

KEYWORDS

Black-box adversarial attack; structural perturbation; graph neural networks; graph classification

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.

CCS '21, November 15–19, 2021, Virtual Event, Republic of Korea

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8454-4/21/11...\$15.00

<https://doi.org/10.1145/3460120.3484796>

ACM Reference Format:

Jiaming Mu, Binghui Wang, Qi Li, Kun Sun, Mingwei Xu, Zhuotao Liu. 2021. A Hard Label Black-box Adversarial Attack Against Graph Neural Networks. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security (CCS '21), November 15–19, 2021, Virtual Event, Republic of Korea*. ACM, New York, NY, USA, 18 pages. <https://doi.org/10.1145/3460120.3484796>

1 INTRODUCTION

Graph neural networks (GNNs) have been widely applied to various graph structure related tasks, e.g., node classification [25], link prediction [59], and graph classification [55, 58, 62], and achieved state-of-the-art performance. For instance, in graph classification, given a set of graphs and each graph is associated with a label, a GNN learns the patterns of the graphs by minimizing the cross entropy between the predicted labels and the true labels of these graphs [55] and predicts a label for each graph. GNN has been used to perform graph classification in various applications such as malware detection [49], brain data analysis [31], superpixel graph classification [1], and protein pattern classification [43].

While GNNs significantly boost the performance of graph data processing, existing studies show that GNNs are vulnerable to adversarial attacks [13, 27, 41, 43, 43, 52, 63]. However, almost all the existing attacks focus on attacking GNNs for node classification, leaving attacks against GNNs for graph classification largely unexplored, though graph classification has been widely applied [1, 31, 43, 49]. Specifically, given a well-trained GNN model for graph classification and a target graph, an attacker aims to perturb the structure (e.g., delete existing edges, add new edges, or rewire edges [33]) of the target graph such that the GNN model will make a wrong prediction for the target graph. Such adversarial attacks could cause serious security issues. For instance, in malware detection [49], by intentionally perturbing a malware graph constructed by a certain malicious program, the malware detector could misclassify the malware to be benign. Therefore, we highlight that it is vital to explore the security of GNNs for graph classification under attack.

In this work, we investigate the most challenging and practical attack, termed *hard label and black-box adversarial attack*, against GNNs for graph classification. In this attack, an attacker cannot obtain any information about the target GNN model and can only obtain *hard labels* (i.e., no knowledge of the probabilities associated with the predicted labels) through querying the GNN model. In addition, we consider that the attacker performs the attack by perturbing the graph structure. The attacker's goal is then to fool the target GNN model by utilizing the hard label after querying the target model and with the minimal graph structural perturbations.

We formulate the attack as a discrete optimization problem, which aims to minimize the graph structure perturbations while maintaining high attack success rates. Note that our attack is harder than the existing black-box attacks (e.g., [11, 12]) that are continuous optimization problems. It is intractable to solve the formulated optimization problem due to the following reasons: (i) The objective function involves the L_0 norm, i.e., the number of perturbed edges in the target graph, and it is hard to be computed. (ii) The searching space for finding the edge perturbations increases exponentially as the number of nodes in a graph increases. That is, it is time-consuming and query-expensive to find appropriate initial perturbations. To address these challenges, we propose a three-phase method to construct our attack. First, we convert the intractable optimization problem to a tractable one via relaxing the L_0 norm to be the L_1 norm, where gradient descent can be applied. Second, we propose a coarse-grained searching algorithm to significantly reduce the search space and efficiently identify initial perturbations, i.e., a much smaller number of edges in the target graph to be perturbed. Note that this algorithm can effectively exploit the graph structural information. Third, we propose a query-efficient gradient computation (QEGC) algorithm to deal with hard labels and adopt the sign stochastic gradient descent (signSGD) algorithm to solve the reformulated attack problem. Note that our QEGC algorithm only needs one query each time to compute the sign of gradients. We also derive theoretical convergence guarantees of our attack.

We systematically evaluate our attack and compare it with two baseline attacks on three real-world datasets, i.e., COIL [37, 39], IMDB [57], and NCI1 [40, 45] from three different fields [35] and three representative GNN methods. Our experimental results demonstrate that our attack can effectively generate adversarial graphs with smaller perturbations and significantly outperforms the baseline attacks. For example, when assuming the same number of edges (e.g., 10% of the total edges in a graph) can be perturbed, our attack can successfully attack around 92% of the testing graphs in the NCI1 dataset, while the state-of-the-art RL-S2V attack [13] can only attack around 75% of the testing graphs. Moreover, only 4.33 edges on average are perturbed by our attack, while the random attack perturbs 10 times more of the edges. Furthermore, to show the effectiveness of our coarse-grained searching algorithm, we compare the performance of three different searching strategies. The results show that coarse-grained searching can significantly speed up the initial searching procedure, e.g., it can reduce 84.85% of the searching time on the NCI1 dataset. It can also help find initial perturbations that can achieve higher attack success rates, e.g., the success rate is improved by around 50%. We also evaluate the effectiveness of the proposed query-efficient gradient computation algorithm. Experimental results show that it decreases the number of queries dramatically. For instance, on the IMDB dataset, our attack with query-efficient gradient computation only needs 13.90% of the queries, compared with our attack without it.

We also explore the countermeasures against the adversarial graphs generated by our attack. Specifically, we propose two different defenses against our adversarial attack: one to detect adversarial graphs and the other to prevent adversarial graph generation. For the former defense, we train a binary GNN classifier, whose training dataset consists of both normal graphs and the corresponding adversarial graphs generated by our attack. Such a classifier aims

to distinguish the structural difference between adversarial graphs and normal graphs. Then the trained classifier is used to detect adversarial graphs generated by our attack on the testing graphs. Our experimental results indicate that such a detector is not effective enough to detect the adversarial graphs. For example, when applying the detector on the COIL dataset with 20% of the total edges are allowed to be perturbed, 47.50% of adversarial graphs can successfully evade the detector. For the latter one, we equip GNN methods with a defense strategy, in order to prevent the generation of adversarial graphs. Specifically, we generalize the low-rank based defense [14] for node classification to graph classification. The main idea is that only low-valued singular components of the adjacency matrix of a graph are affected by the adversarial attacks. Therefore, we propose to discard low-valued singular components to reduce the effects caused by attacks. Our experimental results show that such a defense achieves a clean accuracy-robustness tradeoff. Our contributions are summarized as follows:

- To our best knowledge, we develop the first optimization-based attack against GNNs for graph classification in the hard label and black-box setting.
- We formulate our attack as an optimization problem and solve the problem with convergence guarantee to implement the attack.
- We design a coarse-grained searching algorithm and query-efficient algorithm to significantly reduce the costs of our attack.
- We propose two different types of defenses against our attack.
- We systematically evaluate our attack and defenses on real-world datasets to demonstrate the effectiveness of our attack.

2 THREAT MODEL

Attack Goal. We consider adversarial attacks against GNNs for graph classification. Specifically, given a well-trained GNN model f for graph classification and a *target graph* G with a label y_0 , an attacker aims to perturb the target graph (e.g., delete existing edges, add new edges, or rewire edges in the graph) such that the perturbed target graph (denoted as G') is misclassified by the GNN model f . The attacks can be classified into *targeted attacks* and *non-targeted attacks*. In targeted attacks, an attacker will set a *target label*, e.g., y_c , for the target graph G . Then the attack succeeds, if the predicted label of the perturbed graph is y_c . In non-targeted attacks, the attack succeeds as long as the predicted label of the perturbed graph is different from y_0 . In this paper, we focus on non-targeted attacks and we will also show that our attack can be applied to targeted attacks in Section 4.2.

Attackers' Prior Knowledge. We consider the strictest hard label black-box setting. Specifically, we assume that the attacker can only query the target GNN model f with an input graph and obtain only the predicted hard label (instead of a confidence vector that indicates the probabilities that the graph belongs to each class) for the graph. All the other knowledge, e.g., training graphs, structures and parameters of target GNN model, is unavailable to the attacker.

Attacker's Capabilities. An attacker can perform an adversarial attack by perturbing one of three components in a graph: (i) perturbing nodes, i.e., adding new nodes or deleting existing nodes; (ii) perturbing node feature matrix, i.e., modifying nodes' feature vectors; and (iii) perturbing edges, i.e., adding new edges, deleting existing edges or rewiring edges, which ensures the total number

of edges is unchanged. In this paper, we focus on perturbing edges, which is practical in real-world scenarios. For example, in a social network, an attacker can influence the interactions between user accounts (i.e., modifying the edge status). However, it is hard for the attacker to close legitimate accounts (i.e., deleting nodes) or to modify the personal information of legitimate accounts (i.e., modifying the features). The attacker can also conduct adversarial attacks via adding new nodes to the target graph, which is called *fake node injection attack*. However, its attack performance is significantly impacted by locations of injected nodes, e.g., an attacker needs to add more edges if the injected nodes is on the boundary of a graph, which is however easily detected. What's worse, fake node injection only involves adding edges but it cannot delete edges. Thus, we focus on more generic cases, i.e., perturbing edges, in this paper. More specifically, we assume that the attacker can add new edges and delete existing edges to generate perturbations. To guarantee unnoticeable perturbations, we set a *budget* $b \in [0, 1]$ for perturbing each target graph. That is, perturbed graphs with a *perturbation rate* r , i.e., fraction of edges in the target graph is perturbed, exceeds the budget b are invalid.

As an attacker is often charged according to the number of queries, e.g., querying the model deployed by machine-learning-as-a-service platforms, we also assume that an attacker attempts to reduce the number of queries to save economic costs. In summary, the attacker aims to guarantee the attack success rate with as few queries as possible. Note that it is often a trade-off between the budget and the number of queries. For instance, with a smaller budget, the attacker needs to query the target GNN model more times. Our designed three-phase attack (see Section 4) will obtain a better trade-off.

3 PROBLEM FORMULATION

Given a target GNN model f and a target graph G with label y_0 (Please refer to Appendix A for more background on GNNs for graph classification, due to space limitation), the attacker attempts to generate an *untargeted adversarial graph* G' by perturbing the adjacency matrix A of G to be A' , such that the predicted label of G' will be different from y_0 . Let the *adversarial perturbation* be a binary matrix $\Theta \in \{0, 1\}^{N \times N}$. For ease of description, we fix the entries in the lower triangular part of Θ to be 0, i.e., $\Theta_{ij} = 0 \ \forall j \leq i$, and each entry in the upper triangular part indicates whether the corresponding edge is perturbed or not. Specifically, $\Theta_{ij} = 1, j > i$ means the attacker changes the edge status between nodes i and j , i.e., adding the new edge (i, j) if there is no edge between them in the original graph G or deleting the existing edge (i, j) from G . We keep the edge status between i and j unchanged if $\Theta_{ij} = 0, j > i$. Then the perturbed graph A' can be generated by a perturbation function h , i.e., $A' = h(A, \Theta)$, and h is defined as follows:

$$h(A, \Theta)_{ij} = h(A, \Theta)_{ji} = \begin{cases} A_{ij} & \Theta_{ij} = 0, j > i, \\ -A_{ij} & \Theta_{ij} = 1, j > i. \end{cases} \quad (1)$$

Moreover, the attacker ensures that the perturbation rate r will not exceed a given budget b . Formally, we formulate generating adversarial structural perturbations to a target graph (or called

Algorithm 1 Generating an adversarial graph for a target graph with a hard label black-box access

Input: A trained target GNN model f , a target graph A , perturbation budget b

Output: Adversarial graph A'

```

1: Search initial vector  $\Theta_0$  via coarse-grained searching;
2: for  $t = 1, 2, \dots, T$  do
3:   Randomly sample  $u_1, \dots, u_Q$  from a Gaussian distribution;
4:   Compute  $g(\Theta_t), g(\Theta_t + \mu u_q)$  for  $q = 1, \dots, Q$  via binary search;
5:   Compute  $p(\Theta_t), p(\Theta_t + \mu u_q)$  for  $q = 1, \dots, Q$  using Eq. (6);
6:   Estimate the gradient  $\nabla p(\Theta_t)$  using Eqs. (8) and (9);
7:   Update  $\Theta_{t+1} \leftarrow \Theta_t - \eta_t \nabla p(\Theta_t)$ ;
8: end for
9: Compute  $A' = h(A, \Theta_T)$ ,  $r = \|A' - A\|_0 / (N(N-1))$ ;
10: if  $r \leq b$  then return  $A'$  # succeed
11: else return  $A$  # failed
12: end if
```

adversarial graphs) as the following optimization problem:

$$\begin{aligned} \Theta^* = \arg \min_{\Theta} & \|A' - A\|_0, \\ \text{subject to} & \quad A' = h(A, \Theta), \\ & \quad f(A') \neq y_0, \\ & \quad r \leq b, \end{aligned} \quad (2)$$

where r is defined as $r = \|A' - A\|_0 / (N(N-1))$ and $\|M\|_0$ is the L_0 norm of M , which counts the number of nonzero entries in M .

4 CONSTRUCTING ADVERSARIAL GRAPHS

In this section, we design our hard label black-box adversarial attack to construct adversarial graphs by solving the optimization problem in Eq. (2).

4.1 Overview

Eq. (2) is an intractable optimization problem, and we cannot directly solve it. In order to address this issue, we convert the optimization problem into a tractable one and adopt a sign stochastic gradient descent (signSGD) algorithm to solve it with convergence guarantee. The signSGD algorithm computes gradients of graphs by iteratively querying the target GNN model. We also design two algorithms to reduce the number of queries: a coarse-grained searching algorithm by leveraging the graph structure and a query-efficient gradient computation algorithm that only requires one query in each time of computation. The overview of our attack framework is shown in Figure 1. The attack consists of three phases. First, we relax the intractable optimization problem to a new tractable one (Section 4.2). Second, we develop a coarse-grained searching algorithm to identify a better initial adversarial perturbation/graph (Section 4.3). Third, we propose a query-efficient gradient computation algorithm to deal with hard labels and construct the final adversarial graphs via signSGD (Section 4.4). The whole procedure of generating an adversarial graph for a given target graph is summarized in Algorithm 1.

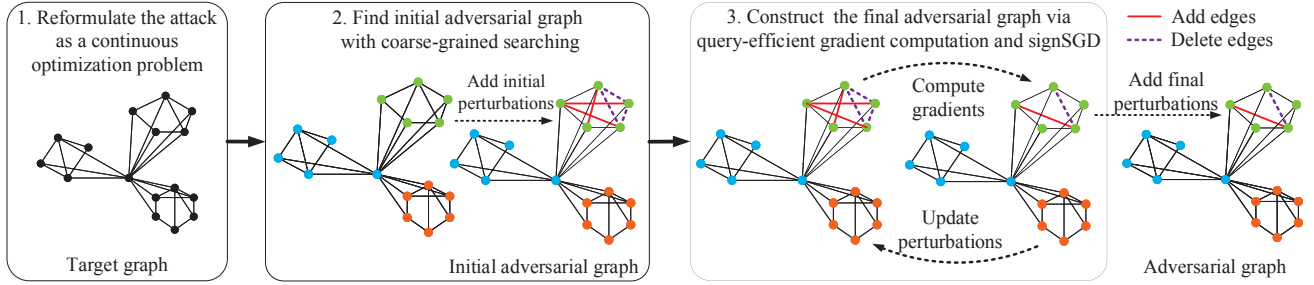


Figure 1: Overview of our hard label black-box attack: (1) We reformulate our attack as an continuous optimization problem that aims at minimizing perturbations on the target graph; (2) We design a coarse-grained searching algorithm to identify initial perturbations for efficient gradient descent computation; (3) We develop a query-efficient gradient computation algorithm, that only needs one query each time to compute the sign of gradients in signSGD. Finally, we obtain the adversarial graph via adding the final perturbations on the target graph.

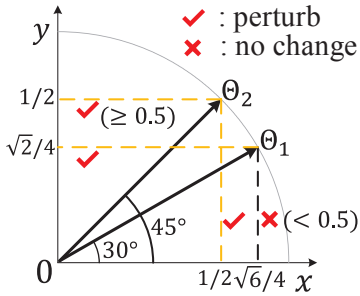


Figure 2: A toy sample: (i) Only one component of Θ_1 exceeds 0.5, which means we only need to perturb one edge in the graph in the direction of Θ_1 ; (ii) Both components of Θ_2 achieve 0.5 and thus we need to perturb both the two edges in the direction of Θ_2 .

4.2 Reformulating the Optimization Problem

The optimization problem defined in Eq. (2) is intractable to solve. This is because the objective function involves L_0 norm related to the variables of adversarial perturbation Θ , which is naturally NP-hard. To cope with this issue, we use the following steps to reformulate the original optimization problem.

Relaxing Θ to be continuous variables. We relax the binary entries $\{0, 1\}$ in Θ to be continuous variables ranging from 0 to 1, i.e., $\Theta_{ij} \in [0, 1] \forall j > i$, such that we can approximate the gradients of the objective function. Each relaxed entry can be treated as the probability that the corresponding edge between two nodes is changed. Specifically, we perturb the edge status between node i and node j if $\Theta_{ij} \geq 0.5$; otherwise not. Thus, the perturbation function h defined in Eq. (1) can be reformulated as follows:

$$h(A, \Theta)_{ij} = h(A, \Theta)_{ji} = \begin{cases} A_{ij} & \Theta_{ij} < 0.5, j > i, \\ -A_{ij} & \Theta_{ij} \geq 0.5, j > i. \end{cases} \quad (3)$$

Defining a new objective function. Next, we define a new objective function that replaces the L_0 norm with the L_1 norm. Similar to existing adversarial attacks against image classifiers [11], we can define a distance function $g(\Theta)$ for GNN to measure the distance

from the target graph to the classification boundary as follows:

$$g(\Theta) = \arg \min_{\lambda > 0} \{f(h(A, \lambda \Theta_{norm})) \neq y_0\}, \quad (4)$$

where Θ_{norm} is the normalized perturbation vector of the perturbation vector Θ^1 that satisfies $\|\Theta_{norm}\|_2 = 1$. $g(\Theta)$ measures the distance from the original graph A to the classification boundary, i.e., the minimal distance λ if we start at A and move to another class in the direction of Θ such that the predicted label of the perturbed graph $A' = h(A, \lambda \Theta_{norm})$ changes. We also denote $\hat{g}(\Theta)$ as a distance vector which starts from A and ends at classification boundary at the direction of Θ with a length of $g(\Theta)$, i.e., $\hat{g}(\Theta) = g(\Theta) \Theta_{norm}$.

A straightforward way of computing optimal Θ^* is to minimize $g(\Theta)$ because smaller $g(\Theta)$ may lead to less elements in Θ that exceed 0.5. Thus, we should change less edges in A for constructing the adversarial graphs. However, it is not effective enough as it does not consider the impact of the search direction, i.e., Θ , on the attack. Specifically, the metrics of our attack is the number of perturbed edges (i.e., the number of entries of Θ that exceed 0.5) instead of the L_2 norm distance (i.e., $g(\Theta)$). The perturbations with different Θ_1 and Θ_2 may be different even if they share the equal distance (i.e., $g(\Theta_1) = g(\Theta_2)$). We explain this via a toy sample shown in Figure 2. We assume that two distance vectors $\hat{g}(\Theta_1)$ and $\hat{g}(\Theta_2)$ with the dimension of 2 have the same length of $\sqrt{2}/2$ in the direction of Θ_1 and Θ_2 , respectively. The lengths of the two components of $\hat{g}(\Theta_1)$ along with x -axis and y -axis are $\sqrt{6}/4$ and $\sqrt{2}/4$, respectively. Thus, we only need to perturb one edge along with x -axis as only $\sqrt{6}/4 \geq 0.5$. However, the lengths of both two components of $\hat{g}(\Theta_2)$ are both 0.5, which means the attacker should perturb both edges because they both achieve the threshold 0.5. Motivated by this toy example and by considering both Θ and $g(\Theta)$, we define the following new objective function:

$$p(\Theta) = \|\text{clip}(\hat{g}(\Theta) - 0.5)\|_0, \quad (5)$$

where $\text{clip}(x)$ is a clip function which clips x into $[0, 1]$. $p(\Theta)$ denotes the number of elements of $\hat{g}(\Theta)$ that exceed 0.5. Thus, it can measure the desired perturbations in the direction of Θ . Here,

¹ Without loss of generality, we transform the triangle perturbation matrix Θ to the corresponding vector form and use perturbation vector and perturbation matrix interchangeably without otherwise mentioned.

in order to calculate the gradients, we also replace the L_0 norm in Eq. (5) with the L_1 norm as follows:

$$p(\Theta) = \|\text{clip}(\widehat{g}(\Theta) - 0.5)\|_1. \quad (6)$$

Converting the optimization problem. According to the definition of $p(\Theta)$, the attacker can find the optimal vector Θ^* by minimizing $p(\Theta)$. Finally, we convert the original optimization problem in Eq. (2) into a new one as follows:

$$\Theta^* = \arg \min_{\Theta} p(\Theta), \quad \text{subject to } r \leq b. \quad (7)$$

Note that, (i) this optimization problem is designed for non-targeted attacks. However, it can also be extended to targeted attacks via changing the condition in Eq. (4) to $f(h(A, \lambda\Theta_{\text{norm}})) = y_c$, where y_c is the target label; (ii) Eq. (7) approximates Eq. (2). Note that, we cannot guarantee that Eq. (2) and (7) have exactly the same optimal values. Nevertheless, our experimental results show that solving Eq. (7) can achieve promising attack performance.

4.3 Coarse-Grained Searching

In this section, we develop a coarse-grained searching algorithm to efficiently identify an initial perturbation vector Θ_0 that makes the corresponding adversarial graph have a different predicted label from y_0 in the direction specified by Θ_0 . We note that, it is difficult to find the valid Θ_0 because the searching space is extremely large when the number of nodes is large. Specifically, a graph with N nodes has $S = N(N-1)/2$ candidate edges. Each edge can be existent or nonexistent so that the search space has a volume of 2^S , which increases exponentially as N increases. The query and computation overhead of traversing all candidate graphs in the searching space is extremely large. Our coarse-grained searching algorithm aims to leverage the graph structure property to reduce the searching space.

We utilize the properties of a graph to reduce the number of the queries and to find a better initial Θ_0 that incurs small perturbations. Specifically, edges in a graph can reflect the similarity among nodes. For instance, there could be more number of edges within a set of nodes, but the number of edges between these nodes and other nodes is much smaller. This means this set of nodes are similar and we can group them into a node cluster. Inspired by graph partitioning [24], we split the original graph into several node clusters, where nodes within each cluster are more similar. We denote *supernode* as one node cluster and *superlink* as the set of links between nodes from two node clusters (See Figure 3). Then, there are three components of a graph for us to search, i.e., (i) supernode; (ii) superlink; and (iii) the whole graph. We take turns to traverse these three types of searching spaces. The reason why we target one specific type of components for perturbation each time is that we can ensure the perturbations follow the same direction, and thus can more effectively generate an adversarial graph. Otherwise, perturbations added in different components will interfere with each other. Thus, randomly selecting Θ_0 is not a good choice as it discards the structural information of the target graph. Our experimental results in Section 5.2 also support our idea.

Particularly, we first partition the graph into node clusters (or supernodes) using the popular and efficient Louvain algorithm [3]. After that, we traverse each supernode. In each supernode c , we

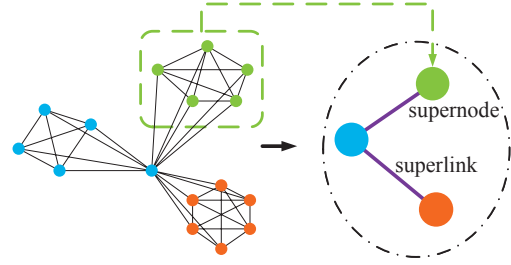


Figure 3: Coarse-grained searching. We partition the graph into several node clusters, denoted as *supernodes*, and links between two supernodes are denoted as a *superlink*.

uniformly choose a fraction $s \in [0, 1]$ at random to determine the number of perturbed edges n , i.e., $n = s \cdot N_c(N_c - 1)/2$, where N_c is the number of nodes in the supernode c . We then randomly select n edges to be perturbed and query the target GNN model to see whether the label of the target graph is changed. We repeat the above process, e.g., $5 \cdot N_c$ times used in our experiments, and always keep the initial perturbation vector with minimal number of perturbed edges. If we failed to find Θ_0 that can change the target label after searching all supernodes, we then search the space within each superlink and finally the whole graph if we still cannot find a successful Θ_0 . During the searching process, we can thus maintain the perturbation vector with the smallest number of edges to be perturbed.

Note that we can significantly reduce the searching overhead by searching supernodes, superlinks, and the whole graph in turn. First, we search supernodes before superlinks because the searching spaces defined by superlinks are larger than those of supernodes. It is not necessary to search within superlinks if we already find Θ_0 within supernodes. Second, the size of the searching space defined by the whole graph is 2^S , which is query and time expensive. As we search it at last, we can find successful Θ_0 in the former two phases (i.e., searching supernodes and superlinks) for most of the target graphs and only a few graphs need to search the whole space (see Section 5.2). Via performing coarse-grained searching, we can *exponentially* reduce the time and the number of queries when the number of supernodes is far more smaller than the number of nodes. The following theorem states the reduction in the time of space searching with our CGS algorithm:

THEOREM 4.1. *Given a graph G with N nodes, the reduction, denoted as β , in the time of space searching with coarse-grained searching satisfies $\beta \approx O(2^{\kappa^4})$, where κ is the number of node clusters and we assume $\kappa \ll N$.*

PROOF. See Appendix B. \square

4.4 Generating Adversarial Graphs via SignSGD

Now we present our sign stochastic gradient descent (signSGD) algorithm to solve the attack optimization problem in Eq (7). Before presenting signSGD, we first describe the method to compute $p(\Theta)$, where only hard label is returned when querying the GNN model; and then introduce a query-efficient gradient computation algorithm to compute the gradients of $p(\Theta)$.

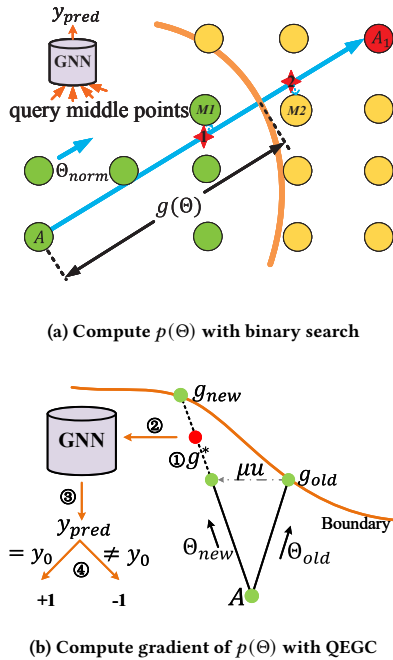


Figure 4: Constructing adversarial graphs. (a) Computing $g(\Theta)$ and $p(\Theta)$ by querying the target model until we find the classification boundary, which will incur many queries when computing gradients of $p(\Theta)$ by using the zeroth order oracle; (b) Query-efficient gradient computation (QEGC). To compare p_{new} and p_{old} (i.e., compute the sign of gradient of $p(\Theta)$), we find g^* in the direction of Θ_{new} such that $p^* = p_{old}$, and use the predicted label to judge if p_{new} is larger than p_{old} after querying $A^* = h(A, g^* \Theta_{new})$.

Computing $p(\Theta)$ via binary search. We describe computing $p(\Theta)$ with only hard label black-box access to the target model.

We first compute $g(\Theta)$ in Eq. (4) via repeatedly querying the target model and further obtain $p(\Theta)$ using Eq. (6). As shown in Figure 4 (a), each edge in the edge space of G can be either existent or nonexistent so that the searching space consists of lattice points (i.e., a lattice point is a symmetrical binary matrix $M \in \{0, 1\}^{N \times N}$) that have equal distance among each other. Suppose there is a classification boundary in the direction of Θ . $g(\Theta)$ is the length of direction vector $\hat{g}(\Theta)$ that begins at the target graph A and ends at the boundary. We can first find a graph A_1 with a different label from y_0 using our CGS algorithm. Since there will be a classification boundary between A and A_1 , we can then conduct a binary search between them, i.e., we query the middle point of the range $[A, A_1]$ (e.g., M_1 in Figure 4 (a)) and update the endpoints of the range based on the predicted label of the middle point in each iteration. The query process ends when the length of range decreases below a tolerance ϵ . With such query process, we can obtain $g(\Theta)$, and then we can compute $p(\Theta)$ easily.

Computing the gradient of $p(\Theta)$ via query-efficient gradient computation (QEGC). We now propose a query efficient algorithm to compute the sign of gradient of $p(\Theta)$, that aims at saving queries used in signSGD in the next part.

With zeroth order oracle, we can estimate the sign of gradient of $p(\Theta)$ via computing $\text{sign}((p(\Theta + \mu u) - p(\Theta))/\mu u)$, where u is a normalized i.i.d direction vector sampled randomly from a Gaussian distribution, and μ is a step constant. The sign can be acquired by computing $p(\Theta + \mu u)$ and $p(\Theta)$ separately. However, we need multiple queries to obtain the value of $p(\Theta)$. As we need to update Θ with many iterations, it is query expensive to compute all $p(\Theta_t + \mu u)$ and $p(\Theta_t)$ at each iteration during the signSGD. Fortunately, we only need to know which p is larger instead of the exact values of them. Thus, we propose a *query-efficient gradient computation (QEGC)* algorithm to compute the sign of gradient with only one query a time as shown in Figure 4 (b).

Suppose the current direction is Θ_{old} with $g(\Theta_{old}) = g_{old}$ and $p(\Theta_{old}) = p_{old}$. Now the direction steps forward with an increment of μu , i.e., $\Theta_{new} = \Theta_{old} + \mu u$. For simplicity of description, we assume Θ_{old} and Θ_{new} are both normalized vectors. We want to judge if p_{new} is larger than p_{old} or not. The idea is that we transfer p_{new} and p_{old} to g_{new} and g_{old} respectively and compare their values. Specifically, for p_{old} , we find g^* such that $p^* = \|clip(g^* \Theta_{new} - 0.5)\|_1 = p_{old}$. For p_{new} , the corresponding g_{new} is the distance from A to the classification boundary at the direction Θ_{new} . Then we query the target model f with graph $A^* = h(A, g^* \Theta_{new})$ to figure out whether g^* exceeds the boundary or not. We say that the classification boundary in the direction of Θ_{new} is closer than that of Θ_{old} if $f(A^*) \neq y_0$ because we cross the boundary with the same p_{old} at the direction of Θ_{new} , while we can only achieve the boundary (but not cross) at the direction of Θ_{old} . Thus, p_{new} is smaller than p_{old} and the sign of gradient is -1 . Similarly, $sign = +1$ if $f(A^*) = y_0$.

In summary, we compute the sign of a gradient as follows:

$$\text{sign}(p(\Theta + \mu u) - p(\Theta)) = \begin{cases} +1 & f(A^*) = y_0, \\ -1 & f(A^*) \neq y_0, \end{cases} \quad (8)$$

where A^* is the graph whose value of p equals to $p(\Theta)$ in the direction of $\Theta + \mu u$. We can use Eq. (8) to save the queries due to the following theorem.

THEOREM 4.2. *Given a normalized direction Θ_{old} with g_{old} and p_{old} , there is one and only one g^* at the direction of Θ_{new} that satisfies $p^* = \|\text{clip}(g^* \Theta_{new} - 0.5)\|_1 = p_{old}$.*

PROOF. See Appendix C.

Solving the converted attack problem via sign Stochastic Gradient Descent (signSGD). We utilize the sign stochastic gradient descent (signSGD) algorithm [2] to solve the converted optimization shown in Eq. (7). The reasons are twofold: (i) the sign operation that compresses the gradient into a binary value is suitable to the hard label scenario; (ii) the sign of the gradient can approximate the exact gradient, which can significantly reduce the query overhead.

Specifically, during the signSGD process, we use Eq. (8) to compute the sign of gradient of $p(\Theta)$ in the direction of u . To ease the noise of gradients, we average the signs of Q gradients in different directions to estimate the derivative of the vector $p(\Theta)$ as follows:

$$\nabla p(\Theta) = \frac{1}{Q} \sum_{q=1}^Q \text{sign} \left(\frac{p(\Theta + \mu u_q) - p(\Theta)}{\mu} u_q \right), \quad (9)$$

where $\nabla p(\Theta)$ is the estimated gradients of $p(\Theta)$, $u_q, q \in 1, 2, \dots, Q$ are normalized i.i.d direction vectors sampled randomly from a Gaussian distribution, and Q is the number of vectors. Recently, Maho et al. [34] proposed a black-box SurFree attack that also involves sampling the direction vector u from a Gaussian distribution. However, the purpose of using u is different from our method. Specifically, u in the SurFree attack is used to compute the distance from the original sample to the boundary, while u in our attack is used to approximate the gradients.

The sign calculated by Eq. (8) depends on a single direction vector u . In contrast, Eq. (9) computes the sign of the average of multiple directions, and can better approximate the real sign of gradient of $p(\Theta)$. Then, we use this gradient estimation to update the search vector Θ by computing $\Theta_{t+1} \leftarrow \Theta_t - \eta_t \nabla p(\Theta_t)$, where η_t is the learning rate in the t -th iteration. After T iterations, we can construct an adversarial graph $A' = h(A, \Theta_T)$ ². The following theorem shows the convergence guarantees of our signSGD for generating adversarial graphs.

Assumption 1. *At any time t , the gradient of the function $p(\Theta)$ is upper bounded by $\|\nabla p(\Theta_t)\|_2 \leq \sigma$, where σ is a non-negative constant.*

THEOREM 4.3. *Suppose that $p(\Theta)$ has L -Lipschitz continuous gradients and Assumption 1 holds. If we randomly pick Θ_R , whose dimensionality is d , from $\{\Theta_t\}_{t=0}^{T-1}$ with probability $P(R=t) = \frac{\eta_t}{\sum_{t=0}^{T-1} \eta_t}$, the convergence rate of our signSGD with $\eta_t = O\left(\frac{1}{\sqrt{dt}}\right)$ and $\mu = O\left(\frac{1}{\sqrt{dt}}\right)$ will give the following bound on $\mathbb{E}[\|\nabla p(\Theta)\|_2]$*

$$\mathbb{E}[\|\nabla p(\Theta)\|_2] = O\left(\frac{\sqrt{dL}}{\sqrt{T}} + \frac{\sqrt{d}}{\sqrt{Q}}\sqrt{Q+d}\right), \quad (10)$$

PROOF. See Appendix D. \square

5 ATTACK RESULTS

In this section, we evaluate the effectiveness our hard label black-box attacks against GNNs for graph classification.

5.1 Experimental Setup

Datasets. We use three real-world graph datasets from three different fields to construct our adversarial attacks, i.e., COIL [37, 39] in the computer vision field, IMDB [57] in the social networks field, and NCI1 [40, 45] in the small molecule field. Detailed statistics of these datasets are in Table 1. By using datasets from different fields with different sizes, we can effectively evaluate the effectiveness of our attacks in different real-world scenarios. We randomly split each dataset into 10 equal parts, of which 9 parts are used to train the target GNN model and the other 1 part is used for testing.

Target GNN model. We choose three representative GNN models, i.e., GIN [55], SAG [26], and GUNet [15] as the target GNN model. We train these models based on the authors' public available source code. The clean training/testing accuracy (without attack) of the three GNN models on the three graph datasets are shown in Table 2. Note that these results are close to those reported in the original

²Our attack can be easily extended to attack directed graphs via only changing the adjacency matrix for directed graphs.

Table 1: Dataset statistics.

Dataset	IMDB	COIL	NCI1
Num. of Graphs	1000	3900	4110
Num. of Classes	2	100	2
Avg. Num. of Nodes	19.77	21.54	29.87
Avg. Num. of Edges	96.53	54.24	32.30

Table 2: Clean accuracy of the three GNN models.

GNN model	Dataset	Train acc	Test acc
GIN	COIL	82.17%	77.95%
	IMDB	69.44%	77.00%
	NCI1	73.59%	77.37%
SAG	COIL	40.85%	42.56%
	IMDB	64.78%	68.00%
	NCI1	73.18%	72.02%
GUNet	COIL	31.25%	31.03%
	IMDB	64.44%	70.00%
	NCI1	69.59%	76.16%

papers. We can see that GIN achieves the best testing accuracy. Thus, we use GIN as the default target model in this paper, unless otherwise mentioned. We also observe that SAG and GUNet perform bad on COIL, and we thus do not conduct attacks on COIL for SAG and GUNet.

Target graphs. We focus on generating untargeted adversarial graphs, i.e., an attacker tries to deceive the target GNN model to output each testing graph a wrong label different from its original label. In our experiments, we select all testing graphs that are correctly classified by the target GNN model as the target graph. For example, the number of target graphs for GIN are 304 on COIL, 77 on IMDB, and 318 on NCI1, respectively.

Metrics. We use four metrics to evaluate the effectiveness of our attacks: (i) Success Rate (SR), i.e., the fraction of successful adversarial graphs over all the target graphs. (ii) Average Perturbation (AP), i.e., the average number of perturbed edges across the successful adversarial graphs. (iii) Average Queries (AQ), i.e., the average number of queries used in the whole attack. (iv) Average Time (AT), i.e., the average time used in the whole attack. We count queries and time for all target graphs even if the attack fails. Note that an attack has better attack performance if it achieves a larger SR or/and a smaller AP, AQ and AT.

Baselines. We compare our attack with state-of-the-art RL-S2V attack [13]. We also choose random attack as a baseline.

- **RL-S2V attack.** RL-S2V is a reinforcement learning based adversarial attack that models the attack as a Finite Horizon Markov Decision Process. To attack each target graph, it first decomposes the action of choosing one perturbed edge in the target graph into two hierarchical actions of choosing two nodes separately. Then it uses Q-learning to learn the Markov decision process. In the RL-S2V attack, the attacker needs to set a maximum number of perturbed edges before the attack. Thus, in our experiments, we first conduct our attack to obtain the perturbation rate and then we set the perturbation rate of RL-S2V attack the same as ours. Thus, the RL-S2V attack and our attack will have the same

Table 3: AQ and AT on three datasets.

Dataset	Metric	Our	RL-S2V	Random
COIL	AQ	1621	1728	1621
	AT (s)	121	245	97
IMDB	AQ	1800	1740	1800
	AT (s)	109	7291	88
NCI1	AQ	1822	1809	1822
	AT (s)	163	3071	104

APs (see Figure 7 and 8). For ease of comparison, we also tune RL-S2V to have a close number of queries as our attack. Then, we compare our attack with RL-S2V in terms of SR and AT.

- Random attack. The attacker first chooses a perturbation ratio uniformly at random. Then, given a target graph, the attacker randomly perturbs the corresponding number of edges in the target graph. For ease of comparison, the attacker will repeat this process and has the same number of queries as our attack, and choose the successful adversarial graph with a *minimal* perturbation as the final adversarial graph. Note that, the random attack we consider is the strongest, as the attacker always chooses the successful adversarial graph with a minimal perturbation.

Parameter setting. All the four metrics are impacted by the pre-set budget b . Unless otherwise mentioned, we set a default $b = 0.2$. Note that we also study the impact of b in our experiments. For other parameters such as Q and μ in signSGD, we set $Q = 100$ and $\mu = 0.1$ by default. In each experiment, we repeat the trail 10 times and use the average results of these trails as the final results to ease the influence of randomness.

5.2 Effectiveness of Our Attack

We conduct experiments to evaluate our hard-label black-box attacks. Specifically, we study the impact of the attack budget, the impact of our coarse-grained searching algorithm, and the impact of query-efficient gradient computation.

5.2.1 Impact of the budget on the attack. Figure 5 and 6 show the SR of the compared attacks with different budgets on the three datasets and three GNN models. We sample 20 different budgets ranging from 0.01 to 0.20 with a step of 0.01. We can observe that: (i) Our attack outperforms the baseline attacks significantly in most cases. For instance, with a budget b less than 0.05, random attack fails to work on the three datasets, while our attack achieves a SR at least 40%; With a budget $b = 0.15$, our attack against GIN achieves a SR of 72% on IMDB, while the SR of RL-S2V is less than 40%. The results show that our proposed optimization-based attack is far more advantageous than the baseline methods. (ii) All methods have a higher SR with a larger budget. This is because a larger budget allows an attacker to perturb more edges in a graph.

We further calculate AP of successful adversarial graphs with different budgets b , and show the results in Figure 7 and 8. Note that, due to algorithmic issue, RL-S2V is set to have the same AP as our attack. We have several observations. (i) The AP of our attack is smaller for achieving a higher SR, which shows that our attack outperforms random attack significantly, even when the considered random attack is the strongest. For example, on the COIL dataset, our attack can achieve a SR of 91.52% when $b = 0.20$ and the

Table 4: Coarse-grained searching with different strategies.

Dataset	Strategy	SR	AP	AQ	AT (s)
COIL	I	0.89	8.88	175	3.07
	II	0.86	9.15	337	7.60
	III	0.84	14.46	339	29.29
IMDB	I	0.79	17.27	293	6.46
	II	0.79	17.22	279	6.66
	III	0.57	17.62	308	18.80
NCI1	I	0.88	12.57	437	7.55
	II	0.89	13.42	725	12.55
	III	0.59	43.09	463	49.87

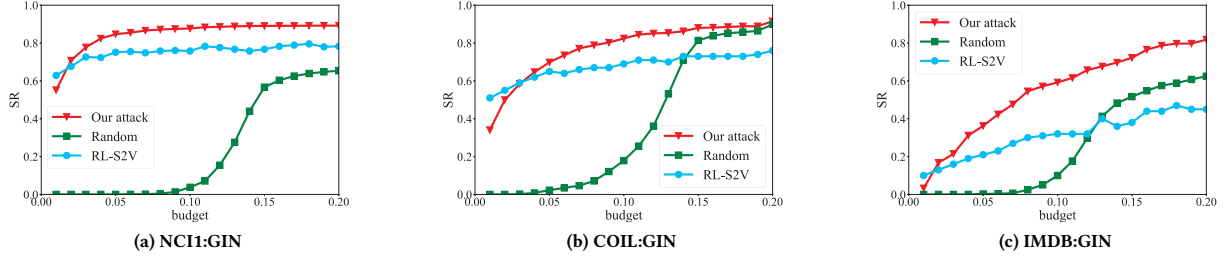
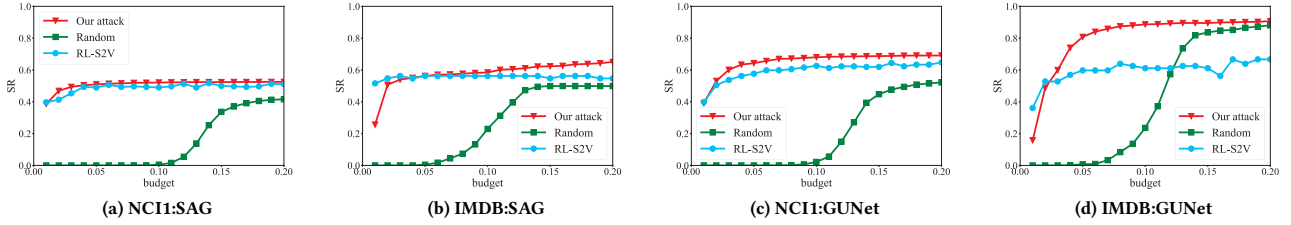
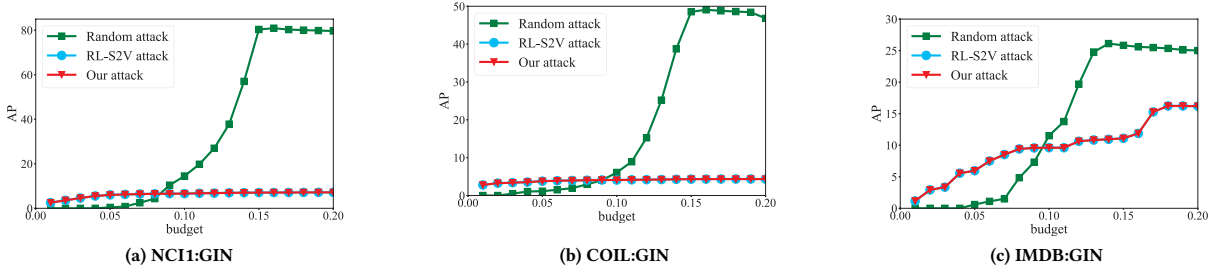
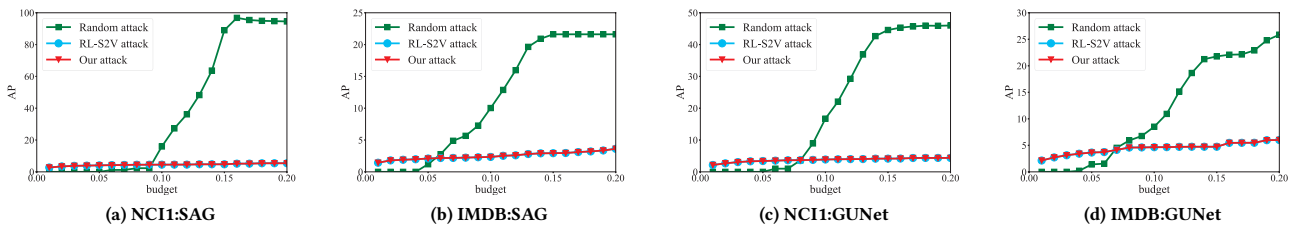
corresponding AP of adversarial graphs is 4.33. Under the same setting, random attack has a SR of only 9.25%. (ii) AP increases with budget b . It is obvious and reasonable because a larger budget means that the perturbed graphs with large perturbations have larger probabilities to generate successful adversarial graphs. (iii) The APs of our attack on three datasets are different. The reason is that these datasets have different average degree. Specifically, IMDB is the most dense graph while NCI1 is the least dense. This result demonstrates that it takes more effort to change the state of supernodes or superlinks of graphs in the dense graph, and thus we need to perturb more edges.

To evaluate the types of perturbations, we record the number of added edges and removed edges for each dataset in our attack. With the target GNN model as GIN and $b = 0.20$, the averaged (added edges, removed edges) on IMDB, COIL, and NCI1 are (8.46, 12.03), (2.51, 1.80), and (12.84, 1.12), respectively. Thus, we can see that we should remove more edges for denser datasets (e.g., IMDB) and add more edges for sparser datasets (e.g., COIL and NCI1).

We also record AQ and AT of the three attack methods on the three datasets, as shown in Table 3. Recall that the three methods are set to have very close number of queries. We observe that RL-S2V has far more AT than our attack and random attack. This is because the searching space of RL-S2V is exponential to the number of nodes of the target graph. Random attack has the smallest AT, as it does not need to compute gradients. Our attack has similar AT as random attack, although it needs to compute gradients.

5.2.2 Impact of coarse-grained searching (CGS) on the attack. In this experiment, we evaluate the impact of different strategies of CGS on the effectiveness of the attack. Specifically, we will validate the importance of initial search in our entire attack. We use three methods to search the initial perturbation vector Θ_0 : (i) Strategy-I (i.e., our strategy): supernode + superlink + whole graph, which means we search the space in the order of supernodes, superlinks and the whole graph (see Section 4.3); (ii) Strategy-II: superlink + supernode + whole graph; and (iii) Strategy-III: whole graph, which means we do not use CGS and search the whole space defined by the target graph directly. Note that, this strategy also means that we start our signSGD based on a randomly chosen Θ_0 .

Table 4 shows the attack results with different strategies against GIN. We have the following observations. (i) The SRs of strategy-I/-II are very close and both are much higher than that of strategy-III. For example, the SRs of strategy-I and -II are 0.88 and 0.89 on NCI1, while that of strategy-III is 0.59. (ii) The APs of strategy-I/-II are

Figure 5: Successful rate (SR) of our attack vs. budget b on the three datasets against GIN.Figure 6: Successful rate (SR) of our attack vs. budget b on IMDB and NCI1 against SAG and GUNet.Figure 7: Average perturbation (AP) of our attack vs. budget b on the three datasets against GIN.Figure 8: Average perturbation (AP) of our attack vs. budget b on IMDB and NCI1 against SAG and GUNet.

much less than that of strategy-III. For instance, AP of strategy-I on the NCI1 dataset is only 12.57, while that of strategy-III is 43.09, about 3.43 times more than the former. This result validates that CGS can find better initial vectors with less perturbations. (iii) Strategy-I has the least searching time and the least number of queries among the three strategies. For instance, it only requires 3.07 seconds to find Θ_0 for target graphs on COIL, while Strategy-II requires 2x time. (iv) The benefit of our CGS algorithm (e.g., Strategy III has a 1.94x AQ and 9.54x AT of our Strategy I on COIL) does

not reach the theoretical level as stated in Theorem 4.1 (i.e., $O(2^{\kappa^4})$). The reason is that from a practical perspective, we assume that the attacker only has maximum number of queries as $5N$, which is exponentially much less than 2^S . If we traverse the entire graph space with strategy-III as stated in Theorem 4.1, AQ and AT will be enlarged to $\frac{2^S}{5N}$ times, which also explains the gap between strategy-I/II and strategy-III. In summary, our proposed CGS algorithm can effectively find initial perturbation vectors with higher success rates, less perturbations, less queries, and shorter time.

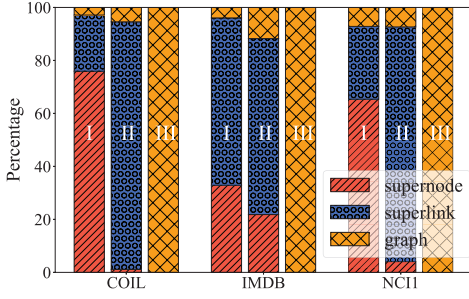


Figure 9: Percentage of the number of adversarial graphs with different initial perturbation vectors found in each component (i.e., supernode, superlink, and the whole graph) under three searching strategies.

Table 5: Impact of query-efficient gradient computation.

Dataset	QEGC	SR	AP	AQ	AT (s)
COIL	Yes	0.92	4.33	1,622	121.21
	No	0.92	4.44	9,859	808.76
IMDB	Yes	0.82	16.19	1,800	109.91
	No	0.82	16.05	12,943	787.67
NCI1	Yes	0.89	7.16	1,822	163.83
	No	0.89	7.60	10,305	1071.17

We further analyze the percentages of adversarial graphs whose initial perturbation vectors Θ_0 are found in each component (i.e., supernode, superlink, and graph). Figure 9 shows the results. Each bar illustrates the percentages of adversarial graphs whose Θ_0 is found by the three components. For instance, on COIL, we obtain 75.73% adversarial graphs whose initial Θ_0 is found in searching supernodes using strategy-I. On COIL and NCI1, we find effective initial vectors Θ_0 by using either strategy-I or strategy-II, i.e., either searching supernodes or superlinks first. Since the searching spaces of supernodes are often smaller than those of superlinks, strategy-I that searches within supernodes first is more suitable on these two datasets. Thus, strategy-I performs best among three strategies on these two datasets. However, on IMDB, we find Θ_0 for most target graphs within the superlinks in both strategy-I/-II. Thus, strategy-II that searches within superlinks first is a better strategy for initial search on IMDB. From Table 4, we can also see that strategy-II performs slightly better than strategy-I in terms of AP and AQ. Note that, the best searching strategies for different datasets are different. The possible reason is, due to the density of the datasets, i.e., strategy-II (i.e., searching superlinks first) is the best strategy for dense graphs (e.g., IMDB), while strategy-I (i.e., searching supernodes first) is for sparse graphs (e.g., NCI1), it is much harder to partition the graphs into supernodes in denser graphs than in sparser graphs.

5.2.3 Impact of query-efficient gradient computation (QEGC). We further conduct experiments to evaluate the impact of QEGC. Table 5 shows the results. We can observe that, under our attack, the number of required queries varies significantly, with and without QEGC. For instance, on IMDB, $AQ = 12,943$ when we do not apply QEGC, while the AQ is reduced to 1,800 when using QEGC, which

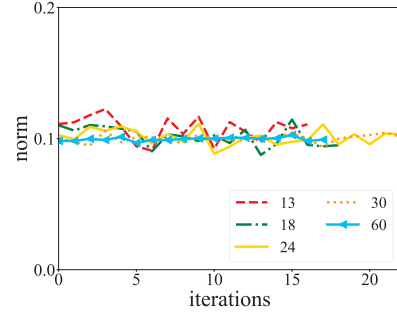


Figure 10: Gradient norms on the IMDB-BINARY dataset.

is only 13.91% of the former. The SR and AP vary slightly with and without QEGC. For example, the APs with and without QEGC on COIL are 4.33 and 4.44, respectively, and the difference is only 0.11. These results demonstrate that QEGC can significantly reduce the number of queries and thus the attack time in our attacks, while maintaining high success rate and incurring small perturbations.

5.2.4 Gradient norms in our attack. The convergence property of our optimization based hard label black-box attack is based on Assumption 1, which requires that the norm of gradient of $p(\Theta)$ should be bounded. Here, we conduct an experiment to verify whether this assumption is satisfied. Specifically, we randomly choose 5 target graphs from IMDB that are successfully attacked by our attack. The number of nodes of these graphs are 13, 18, 24, 30, and 60, respectively. From Figure 10, we can observe that gradient norms are relatively stable and are around 0.1 in all cases. Therefore, Assumption 1 is satisfied in our attacks.

6 DEFENDING AGAINST ADVERSARIAL GRAPHS

In this section, we propose two different defenses against adversarial graphs: one to detect adversarial graphs and the other to prevent adversarial graph generation.

6.1 Adversarial Graph Detection

We first train an adversarial graph detector and then use it to identify whether a testing graph is adversarially perturbed or not. We train GNN models as our detector. Next, we present our methods of generating the training and testing graphs for building the detector, and utilize three different structures of GNN models to construct our detectors. Finally we evaluate our attack under these detectors.

6.1.1 Generating datasets for the detector. The detection process has two phases, i.e., training the detector and detecting testing (adversarial) graphs using the trained detector. Now, we describe how to generate the datasets for training and detection.

Testing dataset. The testing dataset includes all adversarial graphs generated by our attack in Section 5.2 and their corresponding normal (target) graphs. We set labels of adversarial graphs and normal graphs to be 1 and 0 respectively.

Training dataset. The training dataset contains normal graphs and adversarial graphs. Specifically, we first randomly select a number of normal graphs from the training dataset used to train the target GIN model (see Section 5.1). For each sampled normal graph, the

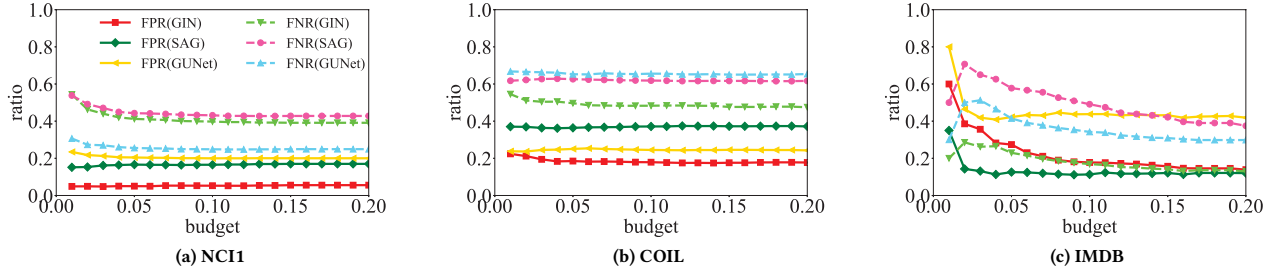


Figure 11: Detection performance vs. budget b on the testing dataset with the training dataset generated by our attack.

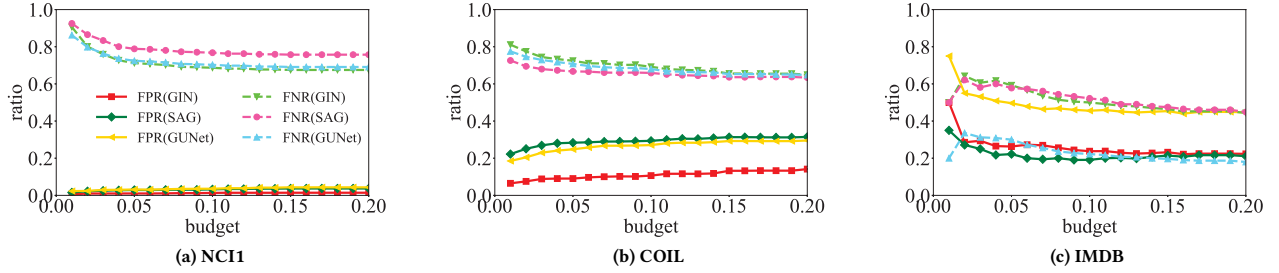


Figure 12: Detection performance vs. budget b on the testing dataset with the training dataset generated by PGD attack.

detector deploys an adversarial attack to generate the corresponding adversarial graph. We use all the sampled normal graphs and the corresponding adversarial graphs to form the training dataset. We consider that the detector uses two different attacks to generate the adversarial graphs: (i) the detector uses our attack, and (ii) the detector uses existing attacks. In our experiments, without loss of generality, we set the size of training dataset to be 3 times of the size of the testing dataset.

When using existing attacks, the detector chooses the projected gradient descent (PGD) attack [54]. PGD attack is a white-box adversarial attack against GCN model for node classification tasks. It first defines a perturbation budget as the maximum number of edges that the attacker can modify. Then it conducts projected gradient descent to minimize the attacker’s objective function. Specifically, in each iteration, the attacker computes the gradients of the objective function w.r.t the edge perturbation matrix. Then it updates the edge perturbation matrix in the opposite direction of the gradient and further projects the edge perturbation matrix into the constraint set such that the number of perturbations is within the pre-set budget. We extend PGD attack for graph classification. Finally, we choose the three aforementioned GNNs (i.e., GIN [55], SAG [26] and GUNet [15]) to train the binary detectors on the constructed training dataset. Note that, we do not use the aforementioned RL-S2V attack or the Random attack to generate training dataset because RL-S2V attack needs large AT and the random attack needs large AQ to produce a reasonable number of adversarial graphs (see Table 3). In contrast, the PGD attack is more efficient.

6.1.2 Detection results. In the detection process, we use *False Positive Rate (FPR)* and *False Negative Rate (FNR)* to evaluate the effectiveness of the detector, where FPR indicates the fraction of normal graphs which are falsely predicted as adversarial graphs, while FNR stands for the fraction of adversarial graphs that are falsely

predicted as normal graphs. Again, we repeat each experiment 10 trails and use the average results of them as the final results to ease the influence of randomness.

Detector results under our attack. The detection performance vs. budget b on the testing dataset with the training dataset generated by our attack is shown in Figure 11. Solid lines indicate FPRs and dashed lines indicate FNRs. We have several observations. (i) The detection performance increases as the budget is getting larger, which means that the detector can distinguish more adversarial graphs if the average perturbations of adversarial graphs are larger. For example, the FPR of GUNet detector on IMDB dataset decreases from 0.80 to 0.42 when the budget increases. This is because when more perturbations are added to the normal graphs, the structure difference between the corresponding adversarial graphs and the normal graphs is larger. Thus, it is easier to distinguish between them. (ii) The detection performances for a specific detector are different on the three datasets. For instance, when using GIN and the budget is 0.20, the FPRs on the three datasets are similar while the FNRs are quite different, i.e., 0.48 (COIL), 0.39 (NCI1) and 0.13 (IMDB), respectively. This is because the average perturbations of adversarial graphs for COIL are the smallest (i.e., 4.33), then for NCI1 (i.e., 7.16) and for IMDB are the largest (i.e., 16.19). (iii) The detector is not effective enough. For example, the smallest FNR on the COIL dataset is 0.48, which means that at least 48% of adversarial graphs cannot be identified by the detector. We have similar observation on the NCI1 dataset. We guess the reason is that the difference between adversarial graphs and normal graphs is too small on these two datasets, and the detector can hardly distinguish between them.

Detection results under the PGD attack. The detection performance vs. budget b on the testing dataset with the training dataset generated by PGD attack is shown in Figure 12. Similarly, we can

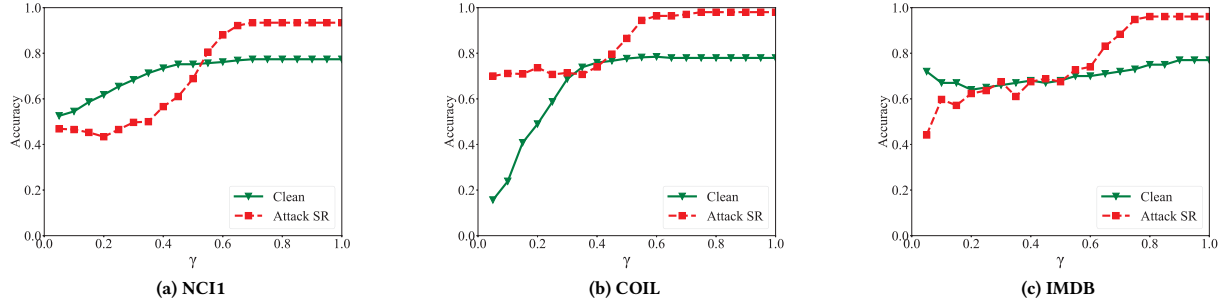


Figure 13: Defense performance against our attack for GIN vs. fraction γ of kept top largest singular values.

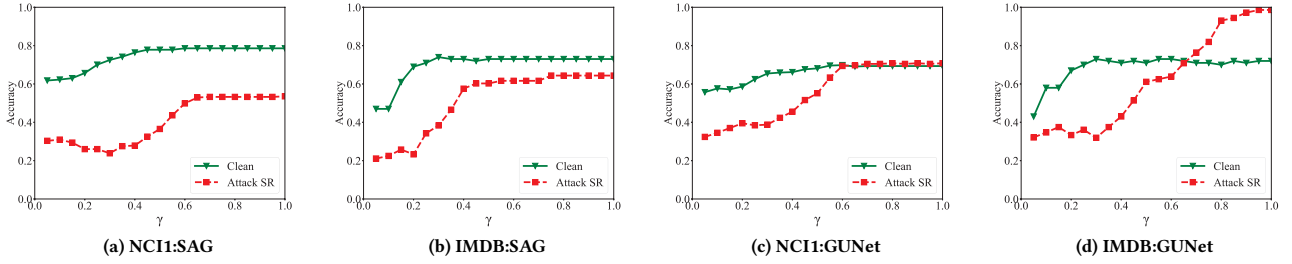


Figure 14: Defense performance against our attack for SAG and GUNets vs fraction γ of kept top largest singular values.

observe that the detection performance is better when the budget is larger. However, the detection performance with the PGD attack is worse than that with our attack. For instance, on the NCI1 dataset, FNRs are around 0.70 using the three detectors with the PGD attack, while the detector with our attack achieves as low as 0.25. One possible reason is that, when the adversarial graphs in the training set are generated by our attack, the detector trained on these graphs can be relatively easier to generalize to the adversarial graphs in the testing dataset that are also generated by our attack. On the other hand, when the adversarial graphs in the training set are generated by the PGD attack, it may be more difficult to generalize to the adversarial graphs in the testing dataset.

6.2 Preventing Adversarial Graph Generation

We propose to equip the GNN model with a defense strategy to prevent adversarial graph generation. Here, we generalize the state-of-the-art low-rank based defense [14] against GNN models for node classification to graph classification.

6.2.1 Low-rank based defense. The main idea is that only high-rank or low-valued singular components of the adjacency matrix of a graph are affected by the adversarial attacks. As these low-valued singular components contain little information of the graph structure, they can be discarded to reduce the effects caused by adversarial attacks, as well as maintaining testing performance of the GNN models. In the context of graph classification, we first conduct a singular value decomposition (SVD) to the adjacency matrix of each testing graph. Then, we keep the top largest singular values and discard the remaining ones. Based on the top largest singular values, we can obtain a new adjacency matrix, and the corresponding graph whose perturbations are removed.

6.2.2 Defense results. We calculate the SR of our attack and the clean testing accuracy after adopting the low-rank based defense. We use the target graphs described in Section 5.1 to calculate the SR. We use the original testing dataset without attack to compute the testing accuracy. For each target graph, we first generate a low-rank approximation of its adjacency matrix by removing small singular values and then feed the new graph into the target GNN model to see if the predicted label is correct or wrong. Figure 13 and 14 show the SR and clean testing accuracy with the low-rank based defense vs. fraction γ of kept top largest singular values for the three GNN models, respectively. γ ranges from 0.05 to 1.0 with a step of 0.05. We have several observations. (i) When γ is relatively small (e.g., ≤ 0.35), i.e., a small fraction of top singular values are kept, the clean accuracy decreases and even dramatically on NCI1 and COIL. One possible reason is that the testing graph structure is damaged. On the other hand, the SR also decreases, meaning adversarial perturbations in certain adversarial graphs are removed. (ii) When γ is relatively large (e.g., ≥ 0.35), i.e., a large fraction of top singular values are kept, the clean accuracy maintains and the SR is relatively high as well. This indicates that adversarial perturbations in a few graphs are removed. The above observations indicate that the fraction γ in low-rank based defense achieves an accuracy-robustness tradeoff. In practice, we need to carefully select γ in order to obtain high robustness against our attack, as well as promising clean testing performance. (iii) When γ is extremely small, e.g., $\gamma = 0.05$, which means that 95% singular values of a graph are removed, the clean accuracy does not decrease much (e.g., see Figure 14 (a) and (c)) or even slightly increases (e.g., see Figure 13 (c)). The results are similar to in [38]. We note that the labels of graphs may be different from the true labels when 95% of their smallest singular values are removed even they are

not perturbed. In our experiments, for ease of analysis, we assume that these graphs have true labels as our goal is to evaluate the impact of the low-rank based defense, i.e., measure if the attack SR significantly decreases with the maintained clean accuracy.

6.3 Discussion

The defense results in Section 6.1.2 and Section 6.2.2 indicate that our adversarial attack is still effective even under detection or prevention. For example, on the NCI1 dataset, the best detection performance is obtained when the budget is 0.20 with GUNet and our attack to train the detector. However, the FNR and FPR are still 0.25 and 0.20, even if the detector has a full knowledge of our attack. What's worse, it is even harder to detect adversarial graphs in real-world scenarios as the detector often does not know the true attack. Low-rank based defense can prevent adversarial graphs to some extent, while needing to sacrifice the testing performance on clean graphs, e.g., as large as 40% of adversarial graphs on NCI1 cannot be prevented even we remove 95% of the smallest singular values.

The proposed detector is a data-level defense strategy and it attempts to block the detected adversarial graphs before they query the target model. It has two key limitations: (i) it is heuristic and (ii) it needs substantial number of adversarial graphs to train the detector and the detection performance highly depends on the quality of the training dataset, i.e., the structure difference between adversarial graphs and normal graphs should be large. The proposed low-rank based defense is a model-level defense strategy and it equips the target GNN model with the smallest singular value removal such that the GNN model can accurately predict testing graphs even they are adversarially perturbed. There are several possible ways to empirically strengthen our defense: (i) Locating the vulnerable regions of graphs based on the feedback of our attacks; (ii) Designing attack-aware graph partitioning algorithm as the method used in our attack is generic and does not exploit the setting of adversarial attacks. (iii) *Adversarial training* [8, 20, 23]. It aims at training a robust GNN model by introducing a white-box adversarial attack and playing a min-max game when training the model. We do not adopt this method because there does not exist white-box attacks against the considered GNN models.

Another direction is to provide the *certified robustness* [4, 22, 47] of GNN models against adversarial structural perturbations. We will also leave those kind of defenses to defend our hard label black box adversarial attacks as the future work.

7 RELATED WORK

Existing studies have shown that GNNs are vulnerable to adversarial attacks [5, 6, 33, 46, 48, 61], which deceive a GNN to produce wrong labels for specific target graphs (in graph classification tasks) or target nodes (in node classification tasks). According to the stages when these attacks occur, they can be classified into training-time poisoning attacks [30, 46, 53, 63, 64] and testing time adversarial attacks [7, 9, 27, 32, 43, 48]. In this paper, we focus on testing time adversarial attacks against classification attacks.

Adversarial attacks against node classification. Existing adversarial attacks mainly attack GNN models for node classification. For node classification, the attacks can be divided into two categories, i.e., optimization based ones [32, 42, 44, 52] and heuristic

based ones that leverage greedy algorithms [9, 50] or reinforcement learning (RL) [13, 41]. In order to develop an optimization based method, the attacker formulates the attack as an optimization problem and solves it via typical techniques such as gradient descent. For example, Xu et al. [54] developed a CW-type loss as the attacker's objective function and utilized projected gradient descent to minimize the loss. As for heuristic based methods, an attacker can utilize a greedy based method, i.e., defining an objective function and traversing all candidate components (e.g., an edge or a node) for adding perturbations. The attacker can select the one that maximizes the objective function to perturb. This process will be repeated multiple times until the attacker finds an adversarial graph or the perturbations exceed the pre-set budget. For instance, Chen et al. [9] proposed an adversarial attack to GCN, which selects the edge of the maximal absolute link gradient and adds it to graph as the perturbation in each iteration.

Adversarial attacks against graph classification. Only a few attacks aim to interfere with graph classification tasks [13, 33, 43]. For instance, Ma et al. [33] proposed a RL based adversarial attack to GNN, which constructs the attack by perturbing the target graph via rewiring. Tang et al. [43] performed the attack against Hierarchical Graph Pooling (HGP) neural networks via a greedy based method. Different from these attacks that are either white-box or grey-box, we study the most challenging *hard label* and *black-box* attack against graph classification in this paper. Our attack is both time and query efficient and is also effective, i.e., high attack success rate with small perturbations.

8 CONCLUSION

We propose a black-box adversarial attack to fool graph neural networks for graph classification tasks in the hard label setting. We formulate the adversarial attack as an optimization problem, which is intractable to solve in its original form. We then relax our attack problem and design a sign stochastic gradient descent algorithm to solve it with convergence guarantee. We also propose two algorithms, i.e., coarse-grained searching and query-efficient gradient computation, to decrease the number of queries during the attack. We conduct our attack against three representative GNN models on real-world datasets from different fields. The experimental results show that our attack is more effective and efficient, when compared with the state-of-the-art attacks. Furthermore, we propose two defense methods to defend against our attack: one to detect adversarial graphs and the other to prevent adversarial graph generation. The evaluation results show that our attack is still effective, which highlights advanced defenses in future work.

ACKNOWLEDGMENTS

We would like to thank our shepherd Pin-Yu Chen and the anonymous reviewers for their comments. This work is supported in part by the National Key R&D Program of China under Grant 2018YFB1800304, NSFC under Grant 62132011, 61625203 and 61832013, U.S. ONR under Grant N00014-18-2893, and BNRist under Grant BNR2020 RC01013. Qi Li and Mingwei Xu are the corresponding authors of this paper.

REFERENCES

- [1] Pedro HC Avelar, Anderson R Tavares, Thiago LT da Silveira, Cláudio R Jung, and Luis C Lamb. 2020. Superpixel image classification with graph attention networks. In *2020 33rd SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI)*. IEEE, 203–209.
- [2] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. 2018. signSGD: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*. PMLR, 560–569.
- [3] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefevre. 2008. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment* 2008, 10 (2008), P10008.
- [4] Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. 2020. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In *International Conference on Machine Learning*. PMLR, 1003–1013.
- [5] Heng Chang, Yu Rong, Tingyang Xu, Wenbing Huang, Honglei Zhang, Peng Cui, Wenwu Zhu, and Junzhou Huang. 2020. A Restricted Black-Box Adversarial Framework Towards Attacking Graph Embedding Models. In *AAAI*. 3389–3396.
- [6] Jinyin Chen, Yixian Chen, Haibin Zheng, Shijing Shen, Shanjing Yu, Dan Zhang, and Qi Xuan. 2020. MGA: Momentum Gradient Attack on Network. *arXiv preprint arXiv:2002.11320* (2020).
- [7] Jinyin Chen, Xiang Lin, Ziqiang Shi, and Yi Liu. 2020. Link prediction adversarial attack via iterative gradient attack. *IEEE Transactions on Computational Social Systems* 7, 4 (2020), 1081–1094.
- [8] Jinyin Chen, Xiang Lin, Hui Xiong, Yangyang Wu, Haibin Zheng, and Qi Xuan. 2020. Smoothing Adversarial Training for GNN. *IEEE Transactions on Computational Social Systems* (2020).
- [9] 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).
- [10] Zhengdao Chen, Xiang Li, and Joan Bruna. 2017. Supervised community detection with line graph neural networks. *arXiv preprint arXiv:1705.08415* (2017).
- [11] Minhao Cheng, Thong Le, Pin-Yu Chen, Jinfeng Yi, Huan Zhang, and Cho-Jui Hsieh. 2018. Query-efficient hard-label black-box attack: An optimization-based approach. *arXiv preprint arXiv:1807.04457* (2018).
- [12] Minhao Cheng, Simranjit Singh, Patrick Chen, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. 2019. Sign-opt: A query-efficient hard-label adversarial attack. *arXiv preprint arXiv:1909.10773* (2019).
- [13] 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).
- [14] Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. 2020. All you need is low (rank) defending against adversarial attacks on graphs. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 169–177.
- [15] Hongyang Gao and Shuiwang Ji. 2019. Graph u-nets. *arXiv preprint arXiv:1905.05178* (2019).
- [16] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. 2017. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212* (2017).
- [17] Will Hamilton, Zhitaoying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *Advances in neural information processing systems*. 1024–1034.
- [18] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584* (2017).
- [19] Xinlei He, Jinyuan Jia, Michael Backes, Neil Zhenqiang Gong, and Yang Zhang. 2020. Stealing Links from Graph Neural Networks. *arXiv preprint arXiv:2005.02131* (2020).
- [20] Weibo Hu, Chuan Chen, Yaomin Chang, Zibin Zheng, and Yunfei Du. 2021. Robust graph convolutional networks with directional graph adversarial training. *Applied Intelligence* (2021), 1–15.
- [21] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. 2019. Strategies for Pre-training Graph Neural Networks. *arXiv preprint arXiv:1905.12265* (2019).
- [22] Hongwei Jin, Zhan Shi, Venkata Jaya Shankar Ashish Peruri, and Xinhua Zhang. 2020. Certified Robustness of Graph Convolution Networks for Graph Classification under Topological Attacks. *Advances in Neural Information Processing Systems* 33 (2020).
- [23] Hongwei Jin and Xinhua Zhang. 2021. Robust Training of Graph Convolutional Networks via Latent Perturbation. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part III*. Springer International Publishing, 394–411.
- [24] George Karypis and Vipin Kumar. 1998. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on scientific Computing* 20, 1 (1998), 359–392.
- [25] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [26] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. 2019. Self-attention graph pooling. *arXiv preprint arXiv:1904.08082* (2019).
- [27] Wanyu Lin, Shengxiang Ji, and Baohun Li. 2020. Adversarial Attacks on Link Prediction Algorithms Based on Graph Neural Networks. In *Proceedings of the 15th ACM Asia Conference on Computer and Communications Security*. 370–380.
- [28] Sijia Liu, Pin-Yu Chen, Xiangyi Chen, and Mingyi Hong. 2018. signSGD via zeroth-order oracle. In *International Conference on Learning Representations*.
- [29] Sijia Liu, Bhavya Kailkhura, Pin-Yu Chen, Paishun Ting, Shiyu Chang, and Lisa Amini. 2018. Zeroth-order stochastic variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems* 31 (2018), 3727–3737.
- [30] Xuanqing Liu, Si Si, Xiaojin Zhu, Yang Li, and Cho-Jui Hsieh. 2019. A unified framework for data poisoning attack to graph-based semi-supervised learning. *arXiv preprint arXiv:1910.14147* (2019).
- [31] Guixiang Ma, Nesreen K Ahmed, Theodore L Willke, Dipanjan Sengupta, Michael W Cole, Nicholas B Turk-Browne, and Philip S Yu. 2019. Deep graph similarity learning for brain data analysis. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 2743–2751.
- [32] Jiaqi Ma, Shuangrui Ding, and Qiaozhu Mei. 2020. Black-box adversarial attacks on graph neural networks with limited node access. *arXiv preprint arXiv:2006.05057* (2020).
- [33] Yao Ma, Suhang Wang, Tyler Derr, Lingfei Wu, and Jiliang Tang. 2019. Attacking graph convolutional networks via rewiring. *arXiv preprint arXiv:1906.03750* (2019).
- [34] Thibault Maho, Teddy Furon, and Erwan Le Merrer. 2021. SurFree: a fast surrogate-free black-box attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10430–10439.
- [35] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. 2020. TUDataset: A collection of benchmark datasets for learning with graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*. arXiv:2007.08663 www.graphlearning.io
- [36] Yurii Nesterov and Vladimir Spokoiny. 2017. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics* 17, 2 (2017), 527–566.
- [37] Kaspar Riesen and Horst Bunke. 2008. IAM graph database repository for graph based pattern recognition and machine learning. In *Joint IAPR International Workshops on Statistical Techniques in Pattern Recognition (SPR) and Structural and Syntactic Pattern Recognition (SSPR)*. Springer, 287–297.
- [38] Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. 2020. Dropedge: Towards deep graph convolutional networks on node classification. In *ICLR*.
- [39] S. K. Nayar S. A. Nene and H. Murase. 1996. Columbia Object Image Library. <https://www.cs.columbia.edu/CAVE/software/softlib/coil-100.php>. (1996).
- [40] Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. 2011. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research* 12, 9 (2011).
- [41] Yiwei Sun, Suhang Wang, Xianfeng Tang, Tsung-Yu Hsieh, and Vasant Honavar. 2019. Node injection attacks on graphs via reinforcement learning. *arXiv preprint arXiv:1909.06543* (2019).
- [42] Tsubasa Takahashi. 2019. Indirect Adversarial Attacks via Poisoning Neighbors for Graph Convolutional Networks. In *2019 IEEE International Conference on Big Data (Big Data)*. IEEE, 1395–1400.
- [43] Haoteng Tang, Guixiang Ma, Yurong Chen, Lei Guo, Wei Wang, Bo Zeng, and Liang Zhan. 2020. Adversarial Attack on Hierarchical Graph Pooling Neural Networks. *arXiv preprint arXiv:2005.11560* (2020).
- [44] Yunzhe Tian, Jiqiang Liu, Endong Tong, Wenjia Niu, Liang Chang, Qi Alfred Chen, Gang Li, and Wei Wang. 2021. Towards Revealing Parallel Adversarial Attack on Politician Socialnet of Graph Structure. *Security and Communication Networks* 2021 (2021).
- [45] Nikil Wale, Ian A Watson, and George Karypis. 2008. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and Information Systems* 14, 3 (2008), 347–375.
- [46] Binghui Wang and Neil Zhenqiang Gong. 2019. Attacking graph-based classification via manipulating the graph structure. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 2023–2040.
- [47] Binghui Wang, Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. 2021. Certified robustness of graph neural networks against adversarial structural perturbation. In *ACM SIGKDD*.
- [48] Binghui Wang, Tianxiang Zhou, Minhua Lin, Pan Zhou, Ang Li, Meng Pang, Cai Fu, Hai Li, and Yiran Chen. 2020. Evasion Attacks to Graph Neural Networks via Influence Function. *arXiv preprint arXiv:2009.00203* (2020).
- [49] Shen Wang, Zhengzhang Chen, Xiao Yu, Ding Li, Jingchao Ni, Lu-An Tang, Jiaping Gui, Zhichun Li, Haifeng Chen, and S Yu Philip. 2019. Heterogeneous Graph Matching Networks for Unknown Malware Detection. In *IJCAI*. 3762–3770.
- [50] Xiaoyun Wang, Minhao Cheng, Joe Eaton, Cho-Jui Hsieh, and Felix Wu. 2018. Attack graph convolutional networks by adding fake nodes. *arXiv preprint arXiv:1810.10751* (2018).
- [51] Felix Wu, Tianyi Zhang, Amauri Holanda de Souza Jr, Christopher Fifty, Tao Yu, and Kilian Q Weinberger. 2019. Simplifying graph convolutional networks. *arXiv preprint arXiv:1902.07153* (2019).

- [52] Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. 2019. Adversarial examples on graph data: Deep insights into attack and defense. *arXiv preprint arXiv:1903.01610* (2019).
- [53] Zhaohan Xi, Ren Pang, Shouling Ji, and Ting Wang. 2020. Graph backdoor. *arXiv preprint arXiv:2006.11890* (2020).
- [54] 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).
- [55] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826* (2018).
- [56] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. *arXiv preprint arXiv:1806.03536* (2018).
- [57] Pinar Yanardag and SVN Vishwanathan. 2015. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1365–1374.
- [58] Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. 2018. Hierarchical graph representation learning with differentiable pooling. In *Advances in neural information processing systems*. 4800–4810.
- [59] Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In *Advances in Neural Information Processing Systems*. 5165–5175.
- [60] Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. 2018. An end-to-end deep learning architecture for graph classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- [61] Zaixi Zhang, Jinyuan Jia, Binghui Wang, and Neil Zhenqiang Gong. 2020. Backdoor attacks to graph neural networks. *arXiv preprint arXiv:2006.11165* (2020).
- [62] Fan Zhou, Chengtai Cao, Kunpeng Zhang, Goce Trajcevski, Ting Zhong, and Ji Geng. 2019. Meta-GNN: On Few-shot Node Classification in Graph Meta-learning. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 2357–2360.
- [63] 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*. 2847–2856.
- [64] Daniel Zügner and Stephan Günnemann. 2019. Adversarial attacks on graph neural networks via meta learning. *arXiv preprint arXiv:1902.08412* (2019).

A BACKGROUND: GRAPH NEURAL NETWORK FOR GRAPH CLASSIFICATION

Graph Neural Networks (GNNs) has been proposed [15, 17, 21, 26, 55] to efficiently process graph data such as social networks, moleculars, financial networks, etc. [18, 19]. GNN learn embedding vectors for each node in the graph, which will be further used in various tasks, e.g., node classification [25], graph classification [55], community detection [10] and link prediction [59]. Specifically, in each hidden layer, the neural network iteratively computes an embedding vector for a node via aggregating the embedding vectors of the node's neighbors in the previous hidden layer [56], which is called *message passing* [16]. Normally, only the embedding vectors of the last hidden layer will be used for subsequent tasks. For example, in node classification, a logistic regression classifier can be used to classify the final embedding vectors to predict the labels of nodes [25]; In graph classifications, information of the embedding vectors in all hidden layers is utilized to jointly determine the graph's label [58, 60]. According to the strategies of message passing, various GNN methods have been designed for handling specific tasks. For instance, Graph Convolutional Network (GCN) [25], GraphSAGE [17], and Simplified Graph Convolution (SGC) [51] are mainly for node classification, while Graph Isomorphism Network (GIN) [55], SAG [26], and Graph U-Nets (GUNet) [15] are for graph classification. In this paper, we choose GIN [55], SAG [26], and GUNet as the target GNN models. Here, we briefly review GIN as it outperforms other GNN models for graph classification.

Graph Isomorphism Network (GIN). Suppose we are given a graph $G = (A, X)$ with label y_0 , where $A \in \{0, 1\}^{N \times N}$ is the symmetric adjacent matrix indicating the edge connections in G , i.e.,

$A_{ij} = 1$ if there is an edge between node i and node j and $A_{ij} = 0$ otherwise. N is the total number of nodes in the graph. $X \in \mathbb{R}^{N \times l}$ is the feature matrix for all nodes, where each row X_i denote the associated l -dimensional feature vector of node i . The process of message passing of an K -layer GIN can be formulated as follows [55]:

$$h_v^k = MLP^{(k)}((1 + \epsilon^{(k)}) \cdot h_v^{(k-1)} + \sum_{u \in \mathcal{N}_v} h_u^{(k-1)}), \quad (11)$$

where $h_v^k \in \mathbb{R}^k$ is the embedding vector of node v at the k -th layer and, for all nodes, $h_i^{(0)} = X_i$, MLP is a multi-layer perceptron whose parameters are trained together with the whole GIN model, $\epsilon^{(k)}$ is a learnable parameter at the k -th layer, and \mathcal{N}_v is the set of neighbor nodes of node v .

To fully utilize the structure information, GIN collects the information from all depth to predict the label of a graph in graph classification tasks. That is, the graph's embedding vector is obtained as follows:

$$h_G^{(k)} = READOUT(\{h_v^{(k)} | v \in G\}), \quad (12)$$

where $h_G^{(k)}$ is the embedding vector of the whole graph at the k -th layer and the $READOUT(\cdot)$ function aggregates node embedding vectors in this hidden layer. $READOUT(\cdot)$ can be a simple permutation invariant function (e.g., summation) or a more sophisticated graph pooling function. In this paper, we choose the graph add pooling function (i.e., adds node features of all nodes in a batch of graphs) as the $READOUT(\cdot)$ function. GIN finally enables a fully-connected layer to each $h_G^{(k)}$ and sum the results to predict the label of the graph, i.e.,

$$y_{pred} = softmax(\sum_{k=0}^{K-1} Linear(h_G^{(k)})), \quad (13)$$

where $Linear$ is a fully-connected layer and $softmax(\cdot)$ is a softmax layer that maps the logits of GIN to values in $[0, 1]$.

B PROOF OF THEOREM 4.1

We restate theorem 4.1:

THEOREM 4.1. *Given a graph G with N nodes, the reduction, denoted as β , in the time of space searching with coarse-grained searching satisfies $\beta \approx O(2^{\kappa^4})$, where κ is the number of node clusters and we assume $\kappa \ll N$.*

Suppose the graph G is partitioned into κ clusters and each cluster has $d_i, i = 1, 2, \dots, \kappa$ nodes. Note that $N = \sum_{i=1}^{\kappa} d_i$.

The searching space without coarse-grained searching (CGS) is:

$$S_{graph} = 2^{\frac{N(N-1)}{2}} \quad (14)$$

The total searching space of all supernodes is:

$$S_{node} = \sum_{i=1}^{\kappa} 2^{\frac{d_i(d_i-1)}{2}} \quad (15)$$

we define a convex function $f(x) = 2^{\frac{x(x-1)}{2}}$ and use Jensen's inequality:

$$S_{node} = \sum_{i=1}^{\kappa} f(d_i) \geq \kappa \cdot f\left(\frac{1}{\kappa} \cdot \sum_{i=1}^{\kappa} d_i\right) = \kappa \cdot f\left(\frac{N}{\kappa}\right) = \kappa \cdot 2^{\frac{d(d-1)}{2}} \quad (16)$$

The equal sign of the inequality holds when $d_1 = d_2 = \dots = d_\kappa = d = \frac{N}{\kappa}$, which means that κ clusters contain equal number of nodes. Similarly, the total searching space of superlinks is:

$$S_{link} = \frac{1}{2} \sum_{i=1}^{\kappa} \sum_{j=1, j \neq i}^{\kappa} 2^{d_i d_j} \quad (17)$$

We define a cluster of convex functions $f_i(x) = 2^{d_i x}$, $i = 1, 2, \dots, \kappa$ and again deploy Jensen's inequality to Eq.(17):

$$\begin{aligned} S_{link} &= \frac{1}{2} \sum_{i=1}^{\kappa} \sum_{j=1, j \neq i}^{\kappa} f_i(d_j) \geq \frac{1}{2} \sum_{i=1}^{\kappa} (\kappa - 1) \cdot f\left(\frac{1}{\kappa - 1} \sum_{j=1, j \neq i}^{\kappa} d_j\right) \\ &= \frac{\kappa - 1}{2} \sum_{i=1}^{\kappa} f\left(\frac{N - d_i}{\kappa - 1}\right) = \frac{\kappa - 1}{2} \sum_{i=1}^{\kappa} 2^{\frac{N d_i - d_i^2}{\kappa - 1}}, \end{aligned} \quad (18)$$

where the equal sign holds when $d_j = \frac{N - d_i}{\kappa - 1}$, $j = 1, 2, \dots, \kappa$, $j \neq i$. We further define a convex function $f_l(x) = 2^{\frac{N x - x^2}{\kappa - 1}}$ and use Jensen's inequality once again to the above equation, we have :

$$\begin{aligned} S_{link} &\geq \frac{\kappa - 1}{2} \sum_{i=1}^{\kappa} 2^{\frac{N d_i - d_i^2}{\kappa - 1}} = \frac{\kappa - 1}{2} \sum_{i=1}^{\kappa} f_l(d_i) \\ &\geq \frac{\kappa - 1}{2} \cdot \kappa \cdot f_l\left(\frac{N}{\kappa}\right) = \frac{\kappa(\kappa - 1)}{2} 2^{\frac{N^2}{\kappa^2}} = \frac{\kappa(\kappa - 1)}{2} 2^{d^2}, \end{aligned} \quad (19)$$

where the equal sign of the second inequality holds when $d_1 = d_2 = \dots = d_\kappa = d = \frac{N}{\kappa}$. In general situations, we can assume that this condition holds. Thus, if we first search within S_{node} and then search within S_{link} , β can be approximated as follows:

$$\begin{aligned} \beta &= \frac{S_{graph}}{S_{node} + S_{link}} \\ &\approx 2^{\frac{N(N-1)}{2}} \div [\kappa \cdot 2^{\frac{d(d-1)}{2}} + \frac{\kappa(\kappa-1)}{2} \cdot 2^{d^2}] \end{aligned} \quad (20)$$

Now suppose $d = t\kappa$ (thus $N = \kappa d = t\kappa^2$), where $t \gg 1$ often in practice. Then, we have:

$$\begin{aligned} \beta &\approx 2^{\frac{N(N-1)}{2}} \div [\kappa \cdot 2^{\frac{d(d-1)}{2}} + \frac{\kappa(\kappa-1)}{2} \cdot 2^{d^2}] \\ &= \frac{2^{\frac{t^2 \kappa^4 - t \kappa^2 + t^2 \kappa^2}{2}}}{\kappa \cdot 2^{\frac{t^2 \kappa^2 - t \kappa}{2}} + \frac{\kappa^2 - \kappa}{2} \cdot 2^{\frac{3t^2 \kappa^2}{2}}} \\ &> \frac{2^{\frac{t^2 \kappa^4 - t \kappa^2 + t^2 \kappa^2}{2}}}{\kappa^2 \cdot 2^{\frac{3t^2 \kappa^2}{2}} + \kappa^2 \cdot 2^{\frac{3t^2 \kappa^2}{2}}} \\ &> \frac{2^{\frac{t^2 \kappa^4}{2}}}{2\kappa^2 \cdot 2^{\frac{3t^2 \kappa^2}{2}}} \\ &= \frac{1}{2\kappa^2} \cdot 2^{\frac{t^2(\kappa^4 - 3\kappa^2)}{2}} \end{aligned} \quad (21)$$

Finally, β in general situations satisfies:

$$\beta \approx O(2^{\kappa^4}) \quad (22)$$

C PROOF OF THEOREM 4.2

We first restate theorem 4.2:

THEOREM 4.2. *Given a normalized direction Θ_{old} with g_{old} and p_{old} , there is one and only one g^* at the direction of Θ_{new} that satisfies $p^* = \|\text{clip}(g^* \Theta_{new} - 0.5)\|_1 = p_{old}$.*

We proof theorem 4.2 by showing that $p(\Theta)$ is a monotone increasing function of $g(\Theta)$. Without loss of generality, we assume two constants with $0 < g_1 < g_2$. They represent two points at the same direction Θ which have distances of g_1 and g_2 respectively from the original graph A . Then we have

$$\begin{aligned} p_1 &= \|\text{clip}(g_1 \Theta - 0.5)\|_1 \\ p_2 &= \|\text{clip}(g_2 \Theta - 0.5)\|_1 \end{aligned} \quad (23)$$

For simplicity, we assume that $\Theta = \{\Theta_1, \dots, \Theta_d\}$ here is a normalized direction vector. We denote I_+ as the set of indexes where the corresponding components of Θ are positive, i.e., $I_+ = \{i_1, i_2, \dots, i_l\}$ where $l = |I_+|$ and $\Theta_i > 0$ for $i \in I_+$. As the $\text{clip}(\cdot)$ function limits the inputs into $[0, 1]$ which will set all negative values as 0, we can rewrite p_1 and p_2 as follows

$$\begin{aligned} p_1 &= \sum_{i \in I_+} (\text{clip}(g_1 \Theta - 0.5))_i \\ p_2 &= \sum_{i \in I_+} (\text{clip}(g_2 \Theta - 0.5))_i \end{aligned} \quad (24)$$

Furthermore, the components of $g_1 \Theta - 0.5$ and $g_2 \Theta - 0.5$ may also be negative because of the -0.5 term. We thus further denote $I_+^{(1)}$ where $g_1 \Theta_i - 0.5 > 0 \forall i \in I_+^{(1)}$ and $I_+^{(2)}$ where $g_2 \Theta_j - 0.5 > 0 \forall j \in I_+^{(2)}$. It is obvious that $I_+^{(1)} \subseteq I_+^{(2)}$ as $0 < g_1 < g_2$ and $\Theta_k > 0 \forall k \in I_+^{(1)} \cup I_+^{(2)}$. Then we have

$$\begin{aligned} p_2 - p_1 &= \sum_{i \in I_+} (\text{clip}(g_2 \Theta - 0.5))_i - \sum_{i \in I_+} (\text{clip}(g_1 \Theta - 0.5))_i \\ &= \sum_{i \in I_+^{(2)}} (\text{clip}(g_2 \Theta - 0.5))_i - \sum_{i \in I_+^{(1)}} (\text{clip}(g_1 \Theta - 0.5))_i \\ &= \sum_{i \in I_+^{(1)}} (\text{clip}(g_2 \Theta - 0.5) - \text{clip}(g_1 \Theta - 0.5))_i \\ &\quad + \sum_{j \in I_+^{(2)} \setminus I_+^{(1)}} (\text{clip}(g_2 \Theta - 0.5))_j \\ &\geq \sum_{i \in I_+^{(1)}} (g_2 - g_1) \Theta_i + \sum_{j \in I_+^{(2)} \setminus I_+^{(1)}} (\text{clip}(g_2 \Theta - 0.5))_j \\ &\geq 0 \end{aligned} \quad (25)$$

The equal sign holds when one of the following two conditions satisfied:

- (i) $I_+^{(1)} = I_+^{(2)} = I_+$, which means that g_1 and g_2 are both large enough such that all positive components of $g_1 \Theta - 0.5$ and $g_2 \Theta - 0.5$ exceed 1.0. Under this condition, we will perturb all edges correspond to the positive components of $g_1 \Theta - 0.5$.
- (ii) $I_+^{(1)} = I_+^{(2)} = \emptyset$, which means that g_1 and g_2 are both small enough such that all positive components of $g_1 \Theta - 0.5$ and $g_2 \Theta - 0.5$ lower than 0. Under this condition, we do not perturb any edge.

Note that, the two conditions above will never be satisfied during our signSGD because we always start our gradient descent at an initial direction Θ_0 with a moderate g value. At each time t , when we step to a new direction, i.e., Θ_{t+1} , we may go into an extreme

condition where $I_+^{(t+1)} = I_+$ or $I_+^{(t+1)} = \emptyset$. However, both conditions will be rejected as the former leads to large perturbations and the later add no perturbations thus will never change the label of target graph. Therefore, during our signSGD, we will always have

$$p_1 < p_2, \quad \forall 0 < g_1 < g_2 \quad (26)$$

Then p is a monotonically increasing function of g thus p and g can be mutually uniquely determined.

D PROOF OF THEOREM 4.3

We first restate theorem 4.3:

THEOREM 4.3. Suppose that $p(\Theta)$ has L -Lipschitz continuous gradients and Assumption 1 holds. If we randomly pick Θ_R , whose dimensionality is d , from $\{\Theta_t\}_{t=0}^{T-1}$ with probability $P(R=t) = \frac{\eta_t}{\sum_{t=0}^{T-1} \eta_t}$, the convergence rate of our signSGD with $\eta_t = O\left(\frac{1}{\sqrt{d}T}\right)$ and $\mu = O\left(\frac{1}{\sqrt{d}T}\right)$ will give the following bound on $\mathbb{E}[\|\nabla p(\Theta)\|_2]$

$$\mathbb{E}[\|\nabla p(\Theta)\|_2] = O\left(\frac{\sqrt{d}L}{\sqrt{T}} + \frac{\sqrt{d}}{\sqrt{Q}}\sqrt{Q+d}\right), \quad (10)$$

Recall that η_t is the learning rate of signSGD in Algorithm 1 at the t -th iteration, $\mu > 0$ is the smoothing parameter and d is the dimension of Θ .

We first define some notations as follows:

$$\hat{\nabla} p(\Theta_t; u_q) = \frac{p(\Theta_t + \mu u_q) - p(\Theta_t)}{\mu} u_q \quad (27)$$

$$\hat{\nabla} p(\Theta_t; u_q) = \text{sign}\left(\frac{p(\Theta_t + \mu u_q) - p(\Theta_t)}{\mu} u_q\right) \quad (28)$$

$$p_\mu(\Theta) = \mathbb{E}_u[p(\Theta + \mu u)] \quad (29)$$

$$\delta_l = \sqrt{\mathbb{E}[(\hat{\nabla} p(\Theta_t; u_q) - \nabla p_\mu(\Theta_t))_l^2]} \quad (30)$$

where $p_\mu(\Theta)$ is the randomized smoothing function of $p(\Theta)$. We can observe that $\hat{\nabla} p(\Theta_t; u_q) = \text{sign}(\hat{\nabla} p(\Theta_t; u_q))$. Moreover, the corresponding estimated gradients are defined as:

$$\hat{p}_t \approx \frac{1}{Q} \sum_{q=1}^Q \frac{p(\Theta_t + \mu u_q) - p(\Theta_t)}{\mu} u_q = \frac{1}{Q} \sum_{q=1}^Q \hat{\nabla} p(\Theta_t; u_q) \quad (31)$$

$$\hat{p}_t \approx \frac{1}{Q} \sum_{q=1}^Q \text{sign}\left(\frac{p(\Theta_t + \mu u_q) - p(\Theta_t)}{\mu} u_q\right) = \frac{1}{Q} \sum_{q=1}^Q \hat{\nabla} p(\Theta_t; u_q) \quad (32)$$

Next, we introduce some lemmas.

LEMMA D.1. $|\langle \nabla p_\mu(\Theta_t), \hat{p}_t \rangle| \Pr[\text{sign}((\hat{p}_t)_l) \neq \text{sign}((\nabla p_\mu(\Theta_t))_l)] \leq \frac{\delta_l}{\sqrt{Q}}$.

PROOF. The proof can be found in [12] Lemma 2. \square

LEMMA D.2. $\mathbb{E}[\|\hat{\nabla} p(\Theta_t; u_q) - \nabla p_\mu(\Theta_t)\|_2^2] \leq \frac{4(Q+1)}{Q} \sigma^2 + \frac{2}{Q} C(d, \mu)$, where $C(d, \mu) = 2d\sigma^2 + \frac{\mu^2 L^2 d^2}{2}$.

PROOF. The proof can be found in [28] proposition 2 with $b = 1$, $q = Q$, $\alpha_b = 1$ and $\beta_b = 0$. As the number of objective function is just one in our optimization problem, so we can choose $b = 1$. Then α_b and β_b can be further fixed. \square

LEMMA D.3. $p_\mu(\Theta_1) - p_\mu(\Theta_T) \leq p_\mu(\Theta_1) - p^* + \mu^2 L$, where p^* is the minimal value of $p(\Theta)$.

PROOF. The proof can be found in [29] Lemma C. \square

LEMMA D.4. $\mathbb{E}[\|\nabla p(\Theta)\|_2] \leq \sqrt{2} \mathbb{E}[\|\nabla p_\mu(\Theta)\|_2] + \frac{\mu L d}{\sqrt{2}}$

PROOF. The proof can be found in [28]. \square

Now we prove our Theorem 4.3. As $p(\Theta)$ has an L -Lipschitz continuous gradient, it is known from [36] that $p_\mu(\Theta)$ also has L -Lipschitz continuous gradient. Based on the L -smoothness of $p_\mu(\Theta)$, we have

$$\begin{aligned} p_\mu(\Theta_{t+1}) &\leq p_\mu(\Theta_t) + \langle \nabla p_\mu(\Theta_t), \Theta_{t+1} - \Theta_t \rangle + \frac{L}{2} \|\Theta_{t+1} - \Theta_t\|_2^2 \\ &= p_\mu(\Theta_t) - \eta_t \langle \nabla p_\mu(\Theta_t), \hat{p}_t \rangle + \frac{L}{2} \eta_t^2 \|\hat{p}_t\|_2^2 \end{aligned} \quad (33)$$

Moreover, we define $(S_t)_l = \frac{1}{Q} \left| \sum_{q=1}^Q \hat{\nabla} p(\Theta_t; u_q)_l \right|$, and thus $\hat{p}_t = S_t \odot \text{sign}(\hat{p}_t)$ and $\|\hat{p}_t\|_2 = \|S_t\|_2$. We can also have

$$\begin{aligned} \langle \nabla p_\mu(\Theta_t), \hat{p}_t \rangle &= \|\nabla p_\mu(\Theta_t)\|_2 \|\hat{p}_t\|_2 \cos(\alpha_{1t}) \\ &= \|\nabla p_\mu(\Theta_t)\|_2 \|S_t\|_2 \cos(\alpha_{1t}) \frac{\cos(\alpha_{2t})}{\cos(\alpha_{2t})} \frac{\|\text{sign}(\hat{p}_t)\|_2}{\|\text{sign}(\hat{p}_t)\|_2} \\ &= \|\nabla p_\mu(\Theta_t)\|_2 \|\text{sign}(\hat{p}_t)\|_2 \cos(\alpha_{2t}) \cdot \frac{\cos(\alpha_{1t})}{\cos(\alpha_{2t})} \frac{\|S_t\|_2}{\sqrt{d}} \\ &= \langle \nabla p_\mu(\Theta_t), \text{sign}(\hat{p}_t) \rangle \cdot \frac{\cos(\alpha_{1t})}{\cos(\alpha_{2t})} \frac{\|S_t\|_2}{\sqrt{d}}, \end{aligned} \quad (34)$$

where α_{1t} is the angle between $\nabla p_\mu(\Theta_t)$ and \hat{p}_t and α_{2t} is the angle between $\nabla p_\mu(\Theta_t)$ and $\text{sign}(\hat{p}_t)$. Substituting Eq. (34) into Eq. (33), and defining $\hat{\eta}_t = \eta_t \cdot \frac{\cos(\alpha_{1t})}{\cos(\alpha_{2t})} \frac{\|S_t\|_2}{\sqrt{d}}$, we have

$$\begin{aligned} p_\mu(\Theta_{t+1}) &\leq p_\mu(\Theta_t) - \hat{\eta}_t \langle \nabla p_\mu(\Theta_t), \text{sign}(\hat{p}_t) \rangle + \frac{dL}{2} \hat{\eta}_t^2 \frac{\cos(\alpha_{2t})^2}{\cos(\alpha_{1t})^2} \\ &= p_\mu(\Theta_t) - \hat{\eta}_t \|\nabla p_\mu(\Theta_t)\|_1 + \frac{dL}{2} \hat{\eta}_t^2 \frac{\cos(\alpha_{2t})^2}{\cos(\alpha_{1t})^2} \\ &\quad + 2\hat{\eta}_t \sum_{l=1}^d |(\nabla p_\mu(\Theta_t))_l| I[\text{sign}((\hat{p}_t)_l) \neq \text{sign}((\nabla p_\mu(\Theta_t))_l)] \end{aligned} \quad (35)$$

Let $c_t = \frac{\cos(\alpha_{2t})}{\cos(\alpha_{1t})}$ and take expectation on both sides, we have

$$\begin{aligned} \mathbb{E}[p_\mu(\Theta_{t+1}) - p_\mu(\Theta_t)] &\leq -\hat{\eta}_t \|\nabla p_\mu(\Theta_t)\|_1 + \frac{dL}{2} \hat{\eta}_t^2 c_t^2 \\ &\quad + 2\hat{\eta}_t \sum_{l=1}^d |(\nabla p_\mu(\Theta_t))_l| \text{Prob}[(\hat{p}_t)_l \neq \text{sign}((\nabla p_\mu(\Theta_t))_l)] \end{aligned} \quad (36)$$

Applying Lemma D.1 into the inequality, we have

$$\mathbb{E}[p_\mu(\Theta_{t+1}) - p_\mu(\Theta_t)] \leq -\hat{\eta}_t \|\nabla p_\mu(\Theta_t)\|_1 + \frac{dL}{2} \hat{\eta}_t^2 c_t^2 + \frac{2\hat{\eta}_t}{\sqrt{Q}} \sum_{l=1}^d \delta_l \quad (37)$$

Note that

$$\begin{aligned} \sum_{l=1}^d \delta_l &\leq \|\delta\|_1 \leq \sqrt{d} \|\delta\|_2 \\ &= \sqrt{d} \sqrt{\mathbb{E}[\|\hat{\nabla} p(\Theta_t; u_q) - \nabla p_\mu(\Theta_t)\|_2^2]} \\ &\leq \sqrt{\frac{d}{Q}} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)}, \end{aligned} \quad (38)$$

where we apply Lemma D.2 in the last inequality in Equation (38). Substituting Equation (38) into Equation (37), we have

$$\begin{aligned} \hat{\eta}_t \|\nabla p_\mu(\Theta_t)\|_1 &\leq \mathbb{E}[p_\mu(\Theta_t) - p_\mu(\Theta_{t+1})] + \frac{dL}{2} \hat{\eta}_t^2 c_t^2 \\ &\quad + \frac{2\sqrt{d}\hat{\eta}_t}{Q} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)} \end{aligned} \quad (39)$$

By summing all inequalities for all t s we obtain

$$\begin{aligned} \sum_{t=1}^T \hat{\eta}_t \mathbb{E}[\|\nabla p_\mu(\Theta_t)\|_1] &\leq \mathbb{E}[p_\mu(\Theta_1) - p_\mu(\Theta_T)] + \frac{dL}{2} \sum_{t=1}^T \hat{\eta}_t^2 c_t^2 \\ &\quad + \sum_{t=1}^T \frac{2\sqrt{d}\hat{\eta}_t}{Q} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)} \end{aligned} \quad (40)$$

Further substituting Lemma D.3 into Inequality (40), we have

$$\begin{aligned} \sum_{t=1}^T \hat{\eta}_t \mathbb{E}[\|\nabla p_\mu(\Theta_t)\|_1] &\leq p_\mu(\Theta_1) - p^* + \mu^2 L + \frac{dL}{2} \sum_{t=1}^T \hat{\eta}_t^2 c_t^2 \\ &\quad + \sum_{t=1}^T \frac{2\sqrt{d}\hat{\eta}_t}{Q} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)} \end{aligned} \quad (41)$$

Dividing $\sum_{t=1}^T \hat{\eta}_t$ on both sides and use the property that $\|\nabla p_\mu(\Theta_t)\|_2 \leq \|\nabla p_\mu(\Theta_t)\|_1$, the inequality (41) can be changed into

$$\begin{aligned} \sum_{t=1}^T \frac{\hat{\eta}_t}{\sum_{t=1}^T \hat{\eta}_t} \mathbb{E}[\|\nabla p_\mu(\Theta_t)\|_2] &\leq \frac{p_\mu(\Theta_1) - p^* + \mu^2 L}{\sum_{t=1}^T \hat{\eta}_t} \\ &\quad + \frac{dL}{2} \frac{\sum_{t=1}^T \hat{\eta}_t^2 c_t^2}{\sum_{t=1}^T \hat{\eta}_t} + \frac{2\sqrt{d}}{Q} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)} \end{aligned} \quad (42)$$

If we randomly pick R from $\{1, \dots, T\}$ with probability $P(R = t) = \frac{\hat{\eta}_t}{\sum_{t=1}^T \hat{\eta}_t}$, we will have

$$\begin{aligned} \mathbb{E}[\|\nabla p_\mu(\Theta_R)\|_2] &= \mathbb{E}[\mathbb{E}_R[\|\nabla p_\mu(\Theta_R)\|_2]] \\ &= \mathbb{E}\left[\sum_{t=1}^T P(R = t) \|\nabla p_\mu(\Theta_t)\|_2\right] \end{aligned} \quad (43)$$

Applying Lemma D.4 into the Equation (43), we have

$$\begin{aligned} \mathbb{E}[\|\nabla p(\Theta)\|_2] &\leq \frac{\sqrt{2}(p_\mu(\Theta_1) - p^* + \mu^2 L)}{\sum_{t=1}^T \hat{\eta}_t} + \frac{dL}{\sqrt{2}} \frac{\sum_{t=1}^T \hat{\eta}_t^2 c_t^2}{\sum_{t=1}^T \hat{\eta}_t} \\ &\quad + \frac{\mu L d}{\sqrt{2}} + \frac{2\sqrt{2}d}{Q} \sqrt{4(Q+1)\sigma^2 + 2C(d, \mu)} \end{aligned} \quad (44)$$

By choosing $\mu = O(\frac{1}{\sqrt{dT}})$ and $\eta_t = \eta = O(\frac{1}{\sqrt{dT}})$, the convergence rate in (44) simplifies to

$$\mathbb{E}[\|\nabla p(\Theta)\|_2] \leq O\left(\frac{\sqrt{d}L}{\sqrt{T}} + \frac{\sqrt{d}}{\sqrt{Q}} \sqrt{Q+d}\right). \quad (45)$$