

Node Injection for Class-specific Network Poisoning

Ansh Kumar Sharma^{a,1}, Rahul Kukreja^{a,1}, Mayank Kharbanda^{a,1}, Tanmoy Chakraborty^b

^a*Indraprastha Institute of Information Technology, Delhi, India*

^b*Indian Institute of Technology, Delhi, India*

Abstract

Graph Neural Networks (GNNs) are powerful in learning rich network representations that aid the performance of downstream tasks. However, recent studies showed that GNNs are vulnerable to adversarial attacks involving node injection and network perturbation. Among these, node injection attacks are more practical as they don't require manipulation in the existing network and can be performed more realistically. In this paper, we propose a novel problem statement – a *class-specific poison attack* on graphs in which the attacker aims to misclassify specific nodes in the target class into a different class using node injection. Additionally, nodes are injected in such a way that they camouflage as benign nodes. We propose NICKI, a novel attacking strategy that utilises an optimization-based approach to sabotage the performance of GNN-based node classifiers. NICKI works in two phases – it first learns the node representation and then generates the features and edges of the injected nodes. Extensive experiments and ablation studies on four benchmark networks show that NICKI is consistently better than four baseline attacking strategies for misclassifying nodes in the target class. We also show that the injected nodes are properly camouflaged as benign, thus making the poisoned graph indistinguishable from its clean version w.r.t various topological properties.

Keywords: Adversarial attack, network poisoning, graph neural networks

¹First three authors contributed equally.

1. Introduction

The network is often used as a simple abstraction of a complex system to solve various real-world tasks such as node classification (Kipf and Welling, 2017; Hamilton et al., 2017), link prediction (Kipf and Welling, 2016) and community detection (Jin et al., 2019). Recent studies have mainly focused on obtaining a rich representation of networks using graph representation learning (GRL) that further improves the performance of the downstream tasks. Of late, graph neural networks (GNNs) have dominated the literature of GRL by showing remarkable performance in tasks such as node classification. With the increasing usage of GNNs, another body of literature (Zügner et al., 2018; Dai et al., 2018; Sun et al., 2020) alerted their limitations by showing how they are vulnerable to adversarial attacks on graphs.

Adversarial attacks on graphs. Adversarial attacks on graph embedding algorithms are challenging due to their discrete and combinatorial nature (Bojchevski and Günnemann, 2019). The majority of studies on adversarial graph attacks assume that the attacker gains access to an existing node and then manipulates its adjacent edges (Zügner et al., 2018; Dai et al., 2018). If the network is attributed, then the features of the corresponding nodes are also manipulated. The attack is generally carried out in a restricted setting, limiting the number of perturbations allowed. A few studies argued that such an attacking strategy is infeasible as it requires a high authority to compromise an existing node (e.g., a user account in a social network) (Zügner et al., 2018). Recent studies introduced a new setting wherein the attacker injects new nodes into the network and performs network perturbations restricted to them (Wang et al., 2018; Sun et al., 2020), which is more realistic as the attacker would have complete command on the injected nodes. For instance, injecting a new node into social networks is equivalent to creating a fake user profile, which is comparatively easier than compromising an existing user account.

Evasion vs poison attacks. Graph adversarial attacks can be classified into two broad categories – evasion attack and poison attack. In an evasion attack (Dai et al., 2018), the attack takes place on the test data; this is in contrast to the poison attack, which makes changes to the training data (Zügner and Günnemann, 2019; Bojchevski and Günnemann, 2019). Once poisoned, the graph is trained again, which requires solving a bi-level optimization function, making the poison attack more difficult than the evasion attack. Due to the dynamic nature of networks like social networks and

transaction networks, graph poisoning attack seems realistic and challenging since these networks require models to be retrained more frequently as new data appears over time.

Contribution I: *Class-specific network poisoning* – A novel problem statement. Our attack aims to misclassify specific *target class* nodes into a different class using node injection. We call this different class our *base class*. Many real-world scenarios back such an attack – (i) Misguiding any fraud detection algorithm in a two-class network of transactions with classes denoting ‘fraudulent’ and ‘legitimate’ users. The goal of the attacker would be to get fraudulent users labelled as legitimate. (ii) Misclassifying low credit users to high credit ones in a network where user relationships are leveraged for credit risk assessment. (iii) Injection of spies by a counter-terrorism squad into a network of ‘terrorists’, ‘civilians’ and ‘soldiers’ so as to baffle the node classifier in identifying the spies in the ‘terrorist’ target class.

Shafahi et al. (2018) defined a similar problem statement for the image domain which works only for continuous spaces. However, we overcome the challenge of discrete space in graph and introduce a novel optimisation based method. During network poisoning, it is equally important that an injected node and an existing node look alike so that the defender won’t trivially spot the attacker node. To address this, we also take measures to ensure that the injected nodes replicate both the topological structure as well as the features of the nodes present in the *base class*. This eventually helps in polarising the target nodes into the *base class*. Continuing with the above examples, in the transaction network, the attacking users will be disguised as legitimate users (base class) and make the model misclassify fraudulent users as legitimate.

Contribution II: NICKI – A novel graph poisoning model. We propose NICKI, an optimization-based node injection approach to hinder the *class-specific network poisoning* task within a budget constraint. We further provide a method to hide the attacker nodes under the hood of base class so as to confuse the annotator and the node classifier

Contribution III: Extensive evaluation. We demonstrate the efficacy of NICKI on different datasets – discrete and continuously attributed, small and medium scale. NICKI outperforms four network poisoning methods with significant margins with up to 42% drop in node classification accuracy even with tighter constraints. We empirically show the working of our hiding module by studying the generated features and edges, and evaluating them under robust models and homophily defenders.

Reproducibility. The source codes and datasets and other execution-related information can be found at <https://github.com/rahulk207/nicki>.

2. Related Work

In general, deep learning algorithms have been shown to be vulnerable to adversarial attacks (Goodfellow et al., 2015; Jia and Liang, 2017). In graph learning, GNNs have achieved fascinating results in a variety of graph-based tasks, such as node classification (Kipf and Welling, 2017; Hamilton et al., 2017; Veličković et al., 2017; Xu et al., 2019b,a), drug design (Jiang et al., 2020) and social recommendation (Ying et al., 2018; Fan et al., 2019). Of late, various studies (Zügner et al., 2018; Dai et al., 2018; Zügner and Günnemann, 2019; Bojchevski and Günnemann, 2019; Sun et al., 2020; Wang et al., 2020) pointed out that GNNs too are susceptible to adversarial attacks, such as their counterparts in the vision and language domains.

Adversarial attack on graphs. Adversarial attack on graphs has recently been explored to show the robustness of GNNs. Dai et al. (2018) proposed different attacking strategies on both node and graph classification tasks. These strategies include a greedy attack based on gradients (GradArgmax), a genetic algorithm-based attack (GeneticAlg), and a hierarchical Q-learning attack (RL-S2V). Zügner et al. (2018) introduced adversarial attacks on attributed graphs by following a greedy approximation scheme while making unnoticeable perturbations. Zügner and Günnemann (2019) used a meta-learning-based attack and showed a significant performance drop in classification accuracy even with small perturbations. In contrast to the above methods which attack node classifiers, Bojchevski and Günnemann (2019) attempted to attack unsupervised embedding algorithms like Deepwalk Perozzi et al. (2014) and node2vec Grover and Leskovec (2016) by perturbing the adjacency matrix. All the aforementioned attacks make a common assumption that the attacker has access to the existing nodes in the graph; therefore, the attacker can perturb edges/features associated with those nodes. Gaining complete access to already existing nodes is a hard assumption and impractical. Therefore, there have been recent studies (Wang et al., 2018; Sun et al., 2020; Wang et al., 2020), which attack via *node injection*.

Node injection attacks. Owing to the more realistic setting of attack, node injection attacks have gained popularity in the adversarial learning paradigm. Similar to RL-S2V, NIPA (Sun et al., 2020) used hierarchical Q-

learning to generate labels and edges of the injected nodes. AFGSM (Wang et al., 2020) provides an approximate closed-form solution for generating new features and edges of the injected nodes. AFGSM is a scalable method, which works well on both small and large-scale graphs and has shown performance at par, if not improved, against methods proposed in Zügner et al. (2018); Zügner and Günnemann (2019) when changed according to node injection strategies. G-NIA (Tao et al., 2021) introduces an attack with the extreme setting of a single-node injection using an optimization-based method and can even work with continuous attributed graphs, unlike previous methods. Later, TDGIA (Zou et al., 2021) utilizes a topological edge selection strategy and a smooth feature optimization objective to generate new edges and features of the injected nodes, respectively. Recently, there has been a study in node injection attacks which focuses on homophily preservation. Fang et al. (2022) is a genetic algorithm based attack, trying to make unnoticeable perturbations in both structural and feature domains. The methods proposed by Chen et al. (2022) and Tao et al. (2022) are developed as plug-ins on already existing graph injection attacks to improve their unnoticeability. With respect to the time of perturbation, an adversarial attack can be classified into two types – evasion and poison attacks. In an evasion attack (Dai et al., 2018), the attack takes place on the test graph after the model has been trained. On the other hand, poison attack (Zügner and Günnemann, 2019; Bojchevski and Günnemann, 2019) occurs on the training graph, and therefore, is comparatively harder to solve because the graph learning protocol is compromised due to the attack. Both poison (Wang et al., 2020; Sun et al., 2020) and evasion (Tao et al., 2021; Zou et al., 2021; Fang et al., 2022; Chen et al., 2022; Tao et al., 2022) attacks have been studied for node injection.

Here we propose a node-injection-based poison attacking strategy that is capable of learning both discrete and continuous attributes through optimization. Unlike existing methods which mostly attack specific nodes, our attack is *class-specific*, attempting to misclassify nodes present in a given target class. A similar type of attack has been studied in the vision domain (Shafahi et al., 2018).

3. Preliminaries

Let $G(V, A, X)$ be an undirected, unweighted and attributed network, where $A \in \{0, 1\}^{N \times N}$ is the adjacency matrix of the graph and $X \in \{0, 1\}^{N \times D}$ is the feature matrix such that $X^T = [x_1, x_2, \dots, x_N]$, where $x_i \in \{0, 1\}^D$ is

Table 1: Useful notations used throughout the paper.

Notation	Denotation
$G = (V, A, X)$	Original graph
$G' = (V', A', X')$	Poisoned graph
L	Set of class labels
V	Set of original nodes
V_L	Set of labelled nodes
V_A	Set of attacker nodes
V_b	Set of base class nodes (b denotes the base class)
X'	Poisoned featured matrix
k	Number of injected attacker nodes
Δ_e, Δ_x	Budgets for edge and feature perturbations
Z_A	Latent representation of attacker nodes
C_e, C_x	Candidate edges and features set
F_e, F_x	Multilayer Perceptrons
S_e, S_x	Scored edges and features

the D -dimensional feature vector of i^{th} node. Let V be the set of nodes with $|V|=N$ and V_L ($V_L \subset V$) be the set of labelled nodes where each node $u \in V_L$ is assigned a class l_i such that $l_i \in L = \{l_1, l_2, \dots, l_{|L|}\}$. In semi-supervised node classification, the goal of a classifier $g : V \rightarrow L$ is to correctly classify the nodes present in $V \setminus V_L$. The objective of the node classification task can then be defined as:

$$\theta' = \arg \min_{\theta} \sum_{v \in V_L} \mathcal{L}(g_{\theta}(G, v), l(v)) \tag{1}$$

Here θ is the set of parameters for the classifier g , $l(v)$ is the ground-truth class of node v , and \mathcal{L} is the node classification loss (typically a cross-entropy function).

Table 1 summarizes useful notations. We define a few important terminologies below:

Definition 3.1 (Attacker Nodes). *Given an original network G , a set of*

attacker nodes $V_{\mathcal{A}}$ are injected by an attacker \mathcal{A} to poison the original network so as to deteriorate the performance of the downstream node classifier g .

Definition 3.2 (Target Class). *Given the original network G and a set of labels L , the target class $l_t \in L$ is a class whose associated nodes are targeted by an attacker \mathcal{A} .*

Definition 3.3 (Class-specific Poison Attack). *Given an original network G , a target class l_t , attacker \mathcal{A} injects nodes $V_{\mathcal{A}}$ along with their vicious features and edges to generate a poisoned network G' . The node classifier g is trained on G' . The aim of the attacker is to let g misclassify the nodes in the target class as much as possible.*

4. Problem Formulation

In a network poison attack, we inject $V_{\mathcal{A}}$, a set of k attacker nodes. We then obtain a poisoned graph $G'(V', A', X')$ with $V' = V \cup V_{\mathcal{A}}$, $A' = \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}$ and $X' = \begin{bmatrix} X \\ X_{\mathcal{A}} \end{bmatrix}$. Here B is a submatrix containing the links between nodes in $V_{\mathcal{A}}$ and V , and C contains internal edges among nodes in $V_{\mathcal{A}}$. $X_{\mathcal{A}}^T = [x_{a_1}, x_{a_2}, \dots, x_{a_k}]$ denotes the feature matrix of the attacker nodes. Note that we start with $B = 0$, $C = I$, and $X_{\mathcal{A}} = 0$ and use the obtained G' as input to our framework pipeline. We iteratively update G' .

Problem 1 (Node misclassification). *The goal is to learn B , C and $X_{\mathcal{A}}$ such that g misclassifies nodes in $V_t \subset V$ having class l_t . The objective function of our poison attack can be defined as a bi-level optimization problem as follows:*

$$\max_{B, C, X_{\mathcal{A}}} \sum_{v \in V_t} \mathbb{I}(g_{\theta'}(G', v) \neq l_t) \quad (2)$$

$$: \|A'\|_0 - \|A\|_0 \leq 2\Delta_e, \|X'\|_0 - \|X\|_0 \leq \Delta_x \quad (3)$$

Here \mathbb{I} is the indicator function. Equation 2 maximizes the number of misclassified nodes belonging to the target class. θ' are the optimal parameters we receive from Equation 1, by tuning the classifier g on the poisoned graph G' . Equation 3 provides an upper bound to the total number of updates in B , C and $X_{\mathcal{A}}$. Note that we do not modify existing adjacency matrix A and features X . Since ours is a poisoning attack, the evaluation of the model is conducted after training g on the poisoned graph.

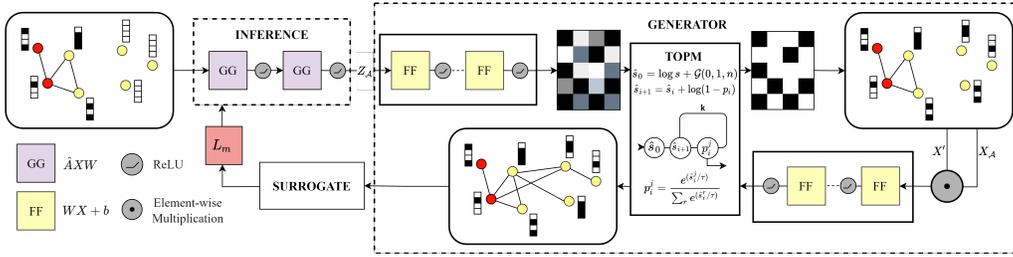


Figure 1: A Schematic diagram of NICKI. The two modules, **Inference** and **Generator**, run in series to produce the poisoned graph $G'(V', A', X')$. The **Surrogate** is trained at every iteration to solve the bi-level optimization, and the loss thus obtained is minimized using gradient descent.

The goal of a poison attack is to infect the training set, thereby making an assumption that all the attacker nodes in V_A also belong to the training set. However, an attacker does not have control over labelling, which poses a major challenge. If not accounted for, the attacker nodes have a high chance of being detected as an outlier and consequently, removed from the network. We address this challenge by making the adversaries belong to a benign class.

Problem 2 (Hiding attacker nodes). *The goal is to camouflage the attacker nodes V_A into a class l_b which we call our “base class”. The objective can be defined as follows:*

$$\arg \min_{\lambda} Im_{\lambda}(V_A, V_b) \quad (4)$$

Here V_b denotes the set of nodes in the base class. Im_{λ} measures the visual distance between V_A and V_b with λ denoting a set of parameters. Here visual distance is used to quantify the closeness of both feature and topological properties of the given two sets of nodes.

The above formulation restricts the nodes in V_A to look similar to those in V_b . This would help mislead the annotator in labelling the attacker nodes as the base class.

5. Proposed Framework

Here we describe our proposed attack framework, NICKI (**N**ode **I**njection for **C**lass-specific **N**etwo**RK** **P**o**I**soning). The attack takes place in a poison setting, i.e., the node classifier would misclassify nodes in the target class on

the poisoned graph. The proposed model, as shown in Figure 1, is a two-part network – (i) inference model, and (ii) generator model. The inference model is a graph encoder which learns a latent representation for each attacker node injected in the graph. The generator model comprises two modules that run sequentially to construct the poisoned graph – (i) a scoring module which scores each edge and feature such that their combination results in misclassification, and (ii) a top- m selector which chooses the best edges/features from scores generated by the scoring module.

The top- m module is used to ensure a differentiable method to impose two constraints on the number of perturbations – Δ_e and Δ_x , which are maximum number of permissible changes in edges and features, respectively.

Even though we impose a budget constraint, our injected attacker nodes could be exposed and subsequently be removed from the graph. Therefore, we also propose sophisticated hiding techniques, which camouflage’s attacker nodes with the nodes of a chosen base class. Once we have the attacker nodes labelled as l_b , we leverage it to misclassify our target nodes into l_b . We, therefore, update Equation 2 as:

$$\max_{B,C,X_A} \sum_{v \in V_t} \mathbb{I}(g_{\theta'}(G', v) = l_b) \quad (5)$$

Post hiding, we use Equation 16 to bring close the embeddings of attacker nodes and base class nodes so that when the classifier is retrained on the poisoned graph, the decision boundary in the embedding space would be expected to rotate such that the attacker nodes are labelled as base class. The target nodes being close to the attacker nodes would also be included in the base class, resulting in misclassification.

5.1. Inference Model

For our implementation, we introduce an intermediate graph $G_P(V', A_P, X_P)$, with $A_P \in \{0, 1\}^{(N+k) \times (N+k)}$ and $X_P \in \{0, 1\}^{(N+k) \times D}$ defined as:

$$(A_P)_{i,j} = \begin{cases} A'_{i,j}, & 0 \leq i, j < N \\ 1, & \text{otherwise} \end{cases} \quad (6)$$

$$(X_P)_{i,j} = \begin{cases} X'_{i,j}, & 0 \leq i < N, 0 \leq j < D \\ 1, & \text{otherwise} \end{cases} \quad (7)$$

We use A_P and X_P as two inputs to the inference model.

We use Graph Autoencoders (GAEs) Kipf and Welling (2016) to learn a joint representation of A_P and X_P given by the latent variable Z which is parameterized by a two-layer Graph Convolutional Network (GCN) as follows:

$$Z = \sigma_2(\hat{A}_P \sigma_1(\hat{A}_P X_P W^{(0)}) W^{(1)}) \quad (8)$$

Here \hat{A}_P represents the normalized adjacency matrix, described as $\hat{A}_P = \hat{D}^{-\frac{1}{2}}(A_P + I_n)\hat{D}^{-\frac{1}{2}}$, \hat{D} is a diagonal matrix with $\hat{D}_{ii} = \sum_j (A_P + I_n)_{ij}$, and I_n is the n th-order identity matrix. $W^{(l)}$ denote the l th-layer weights of the GCN. σ_1 and σ_2 are activation functions.

The inference model is essentially a GCN which propagates knowledge through edges. Therefore, making all the candidate edges and features as 1s in A_P and X_P , respectively (Equations 6 and 7) allows the attacker nodes to indirectly have access to the entire graph during encoding. This eventually results in more meaningful and useful latent embeddings for the attacker nodes.

5.2. Generator Model

Once we encode A_P and X_P into Z , our objective is to use the latent representation of the attacker nodes Z_A to generate both feature and edge perturbations. This is done by using two modules – *scoring* and *top- m selector*. The scoring module generates scores for each edge and feature in the candidate set, defined as C_e and C_x , respectively. The top- m module enforces the budget with a differentiable method and selects the edges and features with the highest scores.

5.2.1. Scoring Module

We derive our set of candidate edges from the submatrices B and C_U in A' , where C_U is the upper triangle of C not including the diagonal elements. We define $(B_{ij} = 0 \text{ or } C_{ij} = 0) \implies e_{ij} \in C_e$, where e_{ij} denotes an edge between nodes i and j . On the other hand, all the attributes present in X_A constitute our candidate features. Therefore, $x_{ij} \in X_A \iff x_{ij} \in C_x$, where x_{ij} denotes the j^{th} feature entry for node i . The scoring of features and edges is performed in series as described below:

Z_A is fed to a Multilayer Perceptron (MLP) F_x as:

$$S_x = F_x(Z_A), \text{ where } Z_A, S_x \in \mathbb{R}^{k \times D} \quad (9)$$

Here k is the number of attacker nodes, and F_x is a one-to-one mapping with S_x representing the score for each candidate feature. S_x is then passed through the top- m selector, which selects Δ_x candidate values to be switched to 1s. The following equation represents the same:

$$X_{\mathcal{A}} = TOPM_{\Delta_x}(S_x), \text{ where } X_{\mathcal{A}} \in \mathbb{R}^{k \times D} \quad (10)$$

where $TOPM_{\Delta_x}(S_x)$ denotes top- m selector for selecting the top Δ_x features according to their scores. Once $X_{\mathcal{A}}$ is generated, we use it to score the candidate edges as follows:

$$S_e = F_e(X' \odot X_{\mathcal{A}}), \text{ where } S_e \in \mathbb{R}^{(k+N) \times k} \quad (11)$$

$$S_e = S_e - \{s_e\} \quad \forall e \notin C_e \quad (12)$$

where $X' = \begin{bmatrix} X \\ X_{\mathcal{A}} \end{bmatrix}$, \odot denotes Vector-wise Dot-product for each pair of vectors, s_e denotes the obtained score for an edge e , and F_e is an MLP. Equation 11 is used to score all candidate edges by element-wise multiplying the features of the pair of nodes connecting each edge. Post scoring, top- m selector is used to select Δ_e edges as follows:

$$T_e = TOPM_{\Delta_e}(S_e), \text{ where } T_e \in \mathbb{R}^{k(k/2+N) \times 1}$$

T_e is further reshaped to matrix form A' by concatenating with original adjacency matrix A .

We use a GCN-based model to evaluate our attack’s performance. We call this our surrogate model which aims to mimic the original classifier g .

5.2.2. Top- m Selector

Given a vector $s = [s_1, s_2, s_3, \dots, s_n]$ of n scores, the goal of top- m selector is to provide a differentiable method, which samples a binary vector t of size n with exactly m 1s denoting the m -highest scores in s .

We adapt Relaxed Subset Sampling Xie and Ermon (2019) for our use case. The algorithm aims to iteratively select the i th maximum element in the i th iteration. In every iteration, first the score vector is updated as:

$$\hat{s}_{i+1} = \hat{s}_i + \log(1 - p_i); \text{ with } \hat{s}_0 = \log s + Gumbel(0, 1, n)$$

where $Gumbel(0, 1, n)$ is an n -dimensional vector sampled from the gumbel distribution. The maximum score is then selected from the updated

Algorithm 1: TOP- m Selector

Input : $s = [s_1, s_2, s_3 \dots s_n], m, \tau$
Output: Relaxed k -hot vector $t = [t_1, t_2, t_3, \dots t_n]$

- 1 $\hat{s} = \log s + \text{Gumbel}(0, 1, n)$
- 2 $p = [0, 0, \dots 0](n - \text{times})$
- 3 $t = [0, 0, \dots 0](n - \text{times})$
- 4 **for** $i \leftarrow 1$ **to** m **do**
- 5 $\hat{s} = \hat{s} + \log(1 - p)$
- 6 $p = \text{softmax}(\hat{s}/\tau)$
- 7 $t = t + p$
- 8 **end**

vector using *tempered softmax* to obtain a probability distribution $p_i = \{p_i^1, p_i^2, \dots p_i^n\}$ where, $p_i^j = \frac{e^{(\hat{s}_i^j/\tau)}}{\sum_r e^{(\hat{s}_i^r/\tau)}}$. Here p_i^j denotes the probability of the j th element of our input vector s being the i th maximum element. τ is the temperature parameter. After m iterations, we obtain our resultant binary vector $t = \sum_{i=1}^m p_i$. Algorithm 1 shows the pseudo-code of the Top- m module.

5.3. Hiding Attackers

We define a hiding module to solve Problem 2 (Section 4) by camouflaging the attacker nodes into a base class l_b , which allows them to be both included and labelled during the training process. We achieve the same by initialising the nodes in $V_{\mathcal{A}}$ with some *pretend edges* and *features*, assuring that they belong to the base class l_b . Our hiding framework is applied before the attack framework (Figure 2). Therefore, it can be understood as a pre-processing strategy.

5.3.1. Pretend Edges

Our goal is to make any node $i \in V_{\mathcal{A}}$ look like it belongs to V_b . To achieve this, we sample d_i , the degree of node i from the power-law degree distribution of nodes in V_b , which we obtain by fitting their Complementary Cumulative Distribution Function (CCDF). If $i \in V_{\mathcal{A}}$, we add a certain fraction of d_i edges from i to nodes in V_b . The connections are made based on the Barabási–Albert model (Barabási and Albert, 1999), where the probability of connection is directly proportional to the degree of nodes in V_b . We repeat

the process for every node $i \in V_A$ and add all such pairs (i, j) , where $j \in V_b$, to our final set of pretend edges, PE . Note that we sample d_i such that $d_i > d_\mu$, where d_μ is the average degree of the input graph. This is to ensure that we have enough edges to connect with the nodes in V_b .

5.3.2. Pretend Features

For pretend features, our objective is to initialise all $u \in V_A$ with features $x_u \in X_P$ such that u pretends to belong to V_b . All nodes belonging to the same class ought to have some feature commonalities that differentiate them from the nodes in other classes. The idea is to fit the features of nodes in l_b class to a distribution $P_b(X)$, from which we can sample $X_u, \forall u \in V_A$. We achieve the same by using the Conditional VAE (CVAE) (Sohn et al., 2015)) architecture. In contrast to VAEs, CVAEs have their output conditioned on l_i -class label which prevents repeating the training process for different base classes. The CVAE objective can be formulated as,

$$L(\phi, \theta, X, l_i) = -KL(q_\phi(z|X, l_i) \parallel p_\theta(z|l_i)) + E_{q_\phi(z|X, l_i)}[\log p_\theta(X|z, l_i)] \quad (13)$$

During testing or generative phase, we give our desired base class l_b as l_i and use the decoder to generate $x_u, \forall u \in V_A$. We stack all such x_u s to form our *pretend features* matrix PF such that $PF \in \{0, 1\}^{k \times D}$. We sample the features for each attacker node independently. Once we obtain our *pretend edges* PE and *pretend features* PF , we replace and re-initialise the input to our attack pipeline, A' and X' , as,

$$A'_{i,j} = \begin{cases} A'_{i,j}, & 0 \leq i, j < N \\ 1, & (i, j) \text{ or } (j, i) \in PE \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

$$X'_{i,j} = \begin{cases} X'_{i,j}, & 0 \leq i < N, 0 \leq j < D \\ PF_{(i-N),j}, & \text{otherwise} \end{cases} \quad (15)$$

We use these matrices, A' and X' , as inputs to our inference model, instead of what is defined in Section 4. Note that according to the definition of C_e in Section 5.2.1, this also changes our candidate edge set as the obtained *pretend edges* are to be masked from the attack pipeline as well.

Algorithm 2 shows the pseudo-code for NICKI.

Algorithm 2: Attack Model

Input : Clean graph $G(A, X)$, label set L , budgets $\Delta_e, \Delta_x, \mathcal{A}$, base class l_b , target class l_t , pre labeled set V_L , HIDE

Output: Poisoned graph $G(A', X')$, and updated labeled set V'_L

- 1 **if** *HIDE* **then**
- 2 Connect nodes in $V_{\mathcal{A}}$ with base class nodes to hide.
- 3 Assign features to each node $i \in V_{\mathcal{A}}$ from the distribution of features in base class.
- 4 Update A and X with attacker nodes.
- 5 **repeat**
- 6 Compute embeddings of A, X using Equation 8.
- 7 Calculate scores for each candidate features as $S_x \leftarrow F_x(Z_{\mathcal{A}})$
- 8 Select top scorers with a constraint on budget using TOP-m Selector as described in Equation 10.
- 9 Update the feature matrix as $X' \leftarrow [X^T, X_{\mathcal{A}}^T]^T$
- 10 Use updated features to generate candidate edges as $S_e \leftarrow F_e(X' \odot X')$
- 11 Pass scores through TOPM module to get the updated adjacency matrix as $A' \leftarrow \text{reshape}(\text{TOPM}_{\Delta_e}(S_e))$
- 12 Train the surrogate model from the updated Adjacency and feature matrix $L_m \leftarrow \text{Surrogate}(A', X')$
- 13 **if** *HIDE* **then**
- 14 | Calculate change in feature using Equation 17
- 15 **until** *iterations* < *Max iterations*;

5.4. Loss Function

We use average cross-entropy for misclassification and hiding feature loss. For misclassification loss (L_m), we compare embeddings of attacker and mean target class nodes as:

$$L_m = \sum_{i=1}^k \frac{-\sum_{j=1}^D Q_{\mu t j} \log Q_{ij}}{k} \quad (16)$$

where $Q_{\mu t}$ is the mean embedding (surrogate model) of the nodes belonging to target class l_t , and Q_i is i th attacker node's embedding (surrogate). While in hiding feature loss (L_f), we minimize the distance between matrices PF

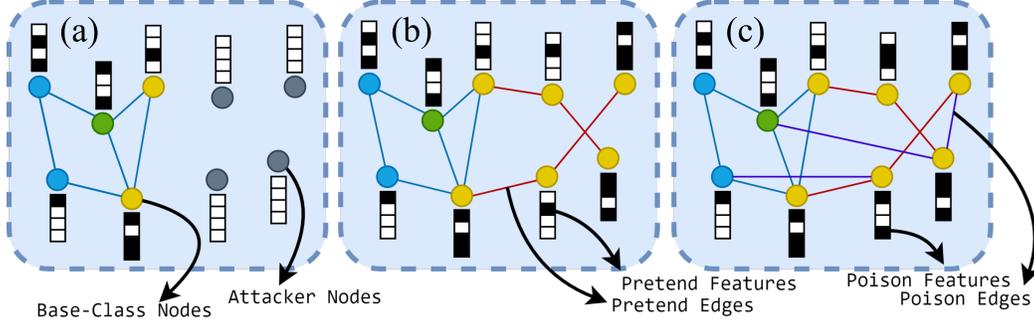


Figure 2: (a) Clean graph just after node injection. (b) Resultant graph post hiding injected attacker nodes. The attacker nodes pretend to be in base class (yellow) using pretend edges and features. (c) Graph after attack where new poison edges are introduced along with pretend edges. Note that the attributes of attacker nodes have changed a little owing to the regularization loss L_f .

and X_A using cross-entropy as follows:

$$L_f = - \sum_{i=1}^k \frac{\bar{X}_i \log \bar{P}F_i}{k} \quad (17)$$

where \bar{X}_i is the i th normalized feature vector of X_A we obtain from the generator module, and $\bar{P}F_i$ is the corresponding normalized *pretend* feature vector.

We train our adversarial attack network by minimising the following loss function L ,

$$L = L_m + \alpha L_f \quad (18)$$

Note that $\alpha \neq 0$ when we employ our hiding framework.

5.5. Time Complexity

The model architecture constitutes GCNs, MLPs and top- m modules with some time consumed by pre-processing in the hiding module. (i) For a single-layer GCN, we need to compute $Z = f(X, A) = \sigma(\hat{A}XW^{(0)})$. If $\hat{A} \in \mathbb{R}^{N \times N}$, $X \in \mathbb{R}^{N \times D}$ and $W^{(0)} \in \mathbb{R}^{D \times H}$, then $Z \in \mathbb{R}^{N \times H}$. The time complexity of the matrix multiplication before the activation function σ can be computed in $\mathcal{O}(ED + NDH)$, where E is the number of edges in matrix A . The time complexity of executing an activation function is linear in the size of input. Given an input $x \in \mathbb{R}^{m_1 \times m_2}$, a hidden layer $h \in \mathbb{R}^{m_2 \times m_3}$,

and an output layer $o \in \mathbb{R}^{m_3}$, the complexity for a forward pass in an MLP is simply the time to compute the matrix multiplication, i.e., $\mathcal{O}(m_1 m_2 m_3)$, with activation function executing in linear time on the input. The time complexity of the top- m selector is m times the size of input, as it executes tempered softmax at each iteration which runs in linear time. So given an input $x \in \mathbb{R}^n$, the time complexity would be $\mathcal{O}(mn)$. (ii) The complexity of a 2-layer GCN in the inference model will take $\mathcal{O}(2(ED + ND^2))$. There are two MLPs – F_x and F_e . The time complexity of F_x and F_e would be $\mathcal{O}(3kD^2 + 2kD)$ and $\mathcal{O}(5NkD^2 + 6NkD)$, respectively. The time complexity of 2 top- m selectors, $TOPM_{\Delta_x}$ and $TOPM_{\Delta_e}$ would be $\mathcal{O}(\Delta_x kD)$ and $\mathcal{O}(\Delta_e k(N + k))$, respectively.

Therefore, the overall complexity of NICKI is $\mathcal{O}(D(E + k(ND + \Delta_x)) + \Delta_e k(N + k))$.

6. Experiments

6.1. Datasets

We use four widely-used benchmark datasets for our experiments.

■ **Cora** (McCallum et al., 2000) is a citation network wherein nodes are scientific papers, and edges indicate citation relationships among them. Each node has a bag-of-word attribute vector obtained from processing the text in the corresponding document. The set of publications is divided into 7 different classes.

■ **Citeseer** (Giles et al., 1998) is also a citation network, similar to Cora. Here papers are divided into 6 classes.

■ **PolBlogs** (Adamic and Glance, 2005) is a network of political weblogs in the US with nodes representing blogs and edges referring to hyperlinks among blogs. Each blog has an associated political orientation, representing its class. There are two classes - Democratic and Republican.

■ The **Reddit** (Hamilton et al., 2017) network contains posts represented by nodes and edges showing post-to-post relationships. The nodes are divided into 41 different classes representing the subreddit the corresponding posts are part of. The attribute vector is obtained from the GloVe (Pennington et al., 2014) word embedding on the title and comments. Table 2 shows the statistics of the datasets.

We train GCN for node classification before and after graph poisoning to study the performance of NICKI. The citation networks have discrete attributes; therefore, showing our performance on them validates the working

Table 2: Network statistics. PolBlog is a featureless graph.

Dataset	N	$ E $	d_{avg}	D	$ L $
Cora	2708	6632	4.898	1433	7
Citeseer	2110	3694	3.501	3703	6
PolBlogs	1086	9502	17.499	-	2
Reddit	10004	31754	6.348	602	41

of our top- m selector module. We convert them into undirected networks. Polblogs is an attribute-less network, thus providing the efficacy of NICKI’s malicious edge generation. The Reddit dataset has continuous attributes, which shows NICKI’s adaptability to work on continuous attribute space. It also contains large number of nodes, which checks our model’s scalability.

6.2. Baseline Methods

NICKI is a node injection based poison attack and hence, we choose baseline attacking strategies based on two conditions – (i) poison, and (ii) node injection. We compare NICKI against two heuristic and two recent methods.

■ **Random.** This attack follows the Erdos-Rényi model (Erdos et al., 1960) for generating new edges. Each new edge is selected with probability $p = ||A||_0/|V|^2$, which is the normalized total degree of G . The process continues until the budget is exhausted. For feature generation, attributes of existing nodes are averaged, and the resultant vector defines the attribute of the attacker nodes.

■ **Preferential.** The attack is based on Barabási–Albert model (Barabási and Albert, 1999). Nodes are injected sequentially, and edges connecting the new node to the rest are selected through preferential attachment. We follow a feature generation procedure identical to that of a random attack.

■ **NIPA.** This is a node injection poison attack, which employs hierarchical Q-learning to generate adversarial edges (Sun et al., 2020). The adjacency matrix learnt by NIPA is used as it is. For attribute generation, it takes an average over the feature vectors of pre-existing nodes and adds some Gaussian noise to it. NIPA is a global attack, and therefore, we modify the test set to contain nodes only belonging to the target class.

■ **AFGSM.** This is a greedy attack, which uses an approximate closed-form solution for generating edges and attributes making it scalable w.r.t the size of the network (Wang et al., 2020). Since AFGSM is a single-node target attack, we modify it to suit our setting by measuring scores corresponding to

each target node. We then sum the obtained scores corresponding to all the target nodes followed by adversarial selection. Therefore, the perturbation which has a higher total score is more likely to be chosen depending on the budget. We train the adaptive variant of AFGSM since it trains the surrogate dynamically during training, thus mimicking a poison attack.

Table 3: Accuracy of GCN-based node classification after adversarial attack, across varying budgets. r controls the budget (Section 6.3). Red (blue) color represents the best (second ranked) model. Lower value indicates better attack. Values within $[\cdot]$ in the first column indicate the classification accuracy on clean graph. See Table 4 for the same using GAT and GraphSAGE as node classifiers.

Dataset	Method	$r = 0.03$	$r = 0.07$	$r = 0.10$	$r = 0.15$
Cora [0.8189]	Random	0.7911	0.8050	0.8078	0.7994
	Preferential	0.7967	0.7911	0.8106	0.7939
	NIPA	0.7632	0.7465	0.7716	0.7994
	AFGSM	0.8161	0.7855	0.7827	0.7437
	NICKI	0.6741	0.5070	0.4345	0.3983
	NICKI (hide)	0.7409	0.7270	0.7103	0.7214
Citeseer [0.7395]	Random	0.7474	0.7500	0.7553	0.7316
	Preferential	0.7526	0.7605	0.7579	0.7342
	NIPA	0.7500	0.7474	0.7447	0.7579
	AFGSM	0.7421	0.7947	0.7421	0.7342
	NICKI	0.6632	0.5632	0.5500	0.4526
	NICKI (hide)	0.7131	0.7263	0.7316	0.7316
Polblogs [0.9316]	Random	0.9203	0.9227	0.9034	0.8913
	Preferential	0.9251	0.9275	0.9203	0.9300
	NIPA	0.9058	0.8986	0.8551	0.8937
	AFGSM	0.9420	0.9420	0.9203	0.9565
	NICKI	0.8164	0.7657	0.7657	0.6522
	NICKI (hide)	0.9109	0.9082	0.9034	0.8527
Reddit [0.9348]	Random	0.9144	0.9125	0.9163	0.9240
	Preferential	0.9183	0.9144	0.9144	0.9221
	NIPA	0.9144	0.9240	0.9202	0.9259
	AFGSM	0.8821	0.8460	0.8441	0.8536
	NICKI	0.8954	0.7529	0.8213	0.7795
	NICKI (hide)	0.9049	0.8935	0.9182	0.8631

6.3. Experimental Setup

To measure the performance of our model, we use the poisoned graph to obtain the classification accuracy w.r.t nodes in the target class, and compare the same with that of the clean graph. The lower the accuracy, the better the attack. Note that we report the accuracy of the node classifier w.r.t only

Table 4: We repeat the experiment shown in Table 3, but with GAT and GraphSAGE as our node classifiers. We do not show the results of random and preferential models as they are the worst among all. Values within $[x, y]$ in the first column indicate the classification accuracy on the clean graphs obtained from GAT (x) and GraphSAGE (y).

Dataset	Method	GAT				GraphSAGE			
		$r = 0.03$	$r = 0.07$	$r = 0.10$	$r = 0.15$	$r = 0.03$	$r = 0.07$	$r = 0.10$	$r = 0.15$
Cora [0.8134, 0.8078]	NIPA	0.7961	0.7701	0.7725	0.8144	0.7658	0.8006	0.7697	0.7950
	AFGSM	0.8056	0.7313	0.8324	0.7452	0.7750	0.7486	0.8464	0.7928
	NICKI	0.6500	0.6096	0.4266	0.5311	0.7806	0.6994	0.6441	0.6130
	NICKI (hide)	0.7000	0.7331	0.7345	0.7288	0.7556	0.7640	0.7006	0.7853
Citeseer [0.8000, 0.7947]	NIPA	0.8449	0.7506	0.7747	0.7755	0.7834	0.7224	0.7696	0.7807
	AFGSM	0.8080	0.7242	0.7519	0.7727	0.7314	0.7010	0.7114	0.7424
	NICKI	0.7325	0.5349	0.6211	0.3368	0.7377	0.6925	0.7010	0.6247
	NICKI (hide)	0.7844	0.8062	0.7655	0.7558	0.7766	0.7726	0.7758	0.7635
Polblogs [0.9444, 0.9444]	NIPA	0.8846	0.9502	0.8744	0.9329	0.8918	0.9005	0.8814	0.8472
	AFGSM	0.9231	0.9151	0.9372	0.9425	0.9591	0.9599	0.9721	0.9264
	NICKI	0.8155	0.8019	0.8406	0.7745	0.8762	0.7971	0.8092	0.7623
	NICKI (hide)	0.8835	0.9469	0.8188	0.9461	0.9417	0.8575	0.9058	0.9069
Reddit [0.9087, 0.9563]	NIPA	0.9032	0.8960	0.8807	0.8639	0.9279	0.9357	0.9545	0.9565
	AFGSM	0.8539	0.8623	0.8617	0.8480	0.9336	0.9302	0.9034	0.9212
	NICKI	0.8740	0.7699	0.7211	0.7324	0.8836	0.7965	0.7381	0.7154
	NICKI (hide)	0.9103	0.8918	0.8710	0.8368	0.9103	0.9013	0.7780	0.8273

the l_t -labelled nodes in the test set. We inject $k = r|V_t|$ nodes, where the controlling parameter r denotes the ratio of injected nodes to target class nodes.

Since ours is a poison attack, we evaluate the accuracy after training our surrogate on the poison graph. We split each dataset into 10%, 10%, 80% of the total nodes for training, validation, and testing, respectively, and use a two-layer GCN with a structure same as the one described in Section 5.1 as our surrogate model. Following the definition of poison attack, we include our attacker nodes V_A into the training set.

We trained our models for 50 epochs and used the learnt graph with minimum classification accuracy as our attack graph. There are 6 layers in F_e with $5D/4$, D , $D/2$, $D/2$, $D/2$, and $D/16$ being the respective number of hidden nodes in layers. In F_x , we have 3 layers with number of hidden nodes being $5D/4$ and $5D/4$. For NICKI (hide), we used 0.7 as parameter for internal degree. α is set to 0.1. We use Adaptive Moment Estimation (Adam) as the gradient descent optimizer. Our code is written in PyTorch. Due to the unavailability of NIPA’s original source code, we used DeepRobust (Li et al., 2020) library’s implementation.

Budget. We use two budgets Δ_e and Δ_x for constraining edge and feature perturbations, respectively. Formally, we set $\Delta_e = kd_{avg}$, where d_{avg}

denotes the average degree of nodes in the graph. Such a formulation helps in retaining the average degree of nodes in the clean graph to that in the poisoned graph, which is more indicative of a real-world attack. For discrete datasets, we set $\Delta_x = \frac{\sum_{i=1}^N \sum_{j=1}^D X_{ij}}{N}$, where X_{ij} denotes the j th attribute of node i . For continuous datasets like Reddit, we limit feature perturbations by enforcing the original features’ range over them. We do this by replacing $TOPM_{\Delta_x}$ from Equation 10 with a sigmoid function followed by mapping the values from $[0, 1]$ to the target range. PolBlogs, being a featureless graph, does not require feature generation. Therefore, we directly use the output of our inference model Z to score the edges in Equation 11. Note that for Polblogs, we use an identity matrix I of size $N + k$ as our feature matrix throughout the framework. To demonstrate our efficiency, we choose $r = \{0.03, 0.07, 0.1, 0.15\}$.

Default target and base classes. We select a target class, which has at least 100 nodes in the test set. This ensures that we have enough attacker nodes even in a parsimonious setting – very few attacker nodes w.r.t the target class nodes in test set. However, in Table 5, we show that our model consistently outperforms other baselines with varying choices of base and target classes.

Node classifier. Throughout the paper, we consider GCN as the default node classifier. We also consider two other neural node classifiers – GAT (Veličković et al., 2017) and GraphSAGE (Hamilton et al., 2017), and report their results in Table 4.

Hiding setup. We use our hiding framework as a pre-processing strategy. As described in Section 5.3, the added pretend edges connect the attacker nodes to nodes in V_b . Once added, we mask them in the attack module to prevent them from further modification. On the other hand, we use CVAE (Sohn et al., 2015) to generate pretend features. However, CVAE does not guarantee discrete features; therefore, for discrete datasets, we define a threshold δ such that $PF_{i,j} = \mathbb{I}(PF'_{i,j} > \delta)$, where PF' denotes the feature matrix obtained from CVAE. For both Cora and Citeseer, we set $\delta = 0.1$ as it best mimics the average number of 1s per vector in the original feature matrix X .

For Reddit, we do not need features to be discrete; hence, we skip the thresholding technique. In case of Polblogs, there is no feature matrix in the network; therefore, the pretend feature module is not required. For each network, we select the class containing the maximum number of nodes as

our base class l_b since injecting new nodes to the already large set (V_b) would make the attack unnoticeable.

6.4. Performance Comparison

We present results for two versions of our model – NICKI and NICKI (hide). The former represents our core attacking strategy described till Section 5.2, while the latter attaches the hiding framework on top of this. The experimental settings for both the versions are described Section 6.3. Table 3 shows the accuracy of the node classifier on different datasets across varying budgets.

Cora, Citeseer and Polblogs. NICKI performs consistently better than every baseline by a considerable margin. The performance difference between NICKI and the best baseline is around 7–9% for the lowest budget and goes as high as 23–35% for a comparatively higher budget. AFGSM, as described in (Wang et al., 2020), uses a closed-form solution which is optimised to attack only a single target node. AFGSM scores each candidate perturbation w.r.t to misclassification of the target node. We adopt AFGSM to our multi-target setting by selecting perturbations that have a higher sum of scores for each node in the target class. The poor performance of AFGSM is indicative of its closed-form solution not being generalizable for attacking multiple target nodes. As expected, random and preferential models do not perform well; the misclassification increases as we increase the budget and in some cases gets better than AFGSM. NIPA, on the other hand, is a global attack, and hence, is better suited for our setting. The superior performance of NIPA compared to AFGSM on Cora and Polblogs can be attributed to the above fact. NICKI (hide), as expected, has an inferior performance than NICKI. However, for Cora and Citeseer, it still outperforms other baselines by a significant margin. We can observe that the baselines do not necessarily show improvement in performance with increasing budget. It is attributed to the fact that these algorithms are not optimized for low-budget settings like ours. NICKI, on the other hand, although not fully consistent due to probabilistic nature of gradient based algorithms, shows an upward trend in misclassification with increasing budget. The parameter α in Equation 18 controls the extent of hiding, and in turn, balances the trade-off between better misclassification and better hiding. For Cora, we also observe the reclassification of target nodes after the attack. From Figure 3, we can infer that majority of the nodes are reclassified as base class nodes. This stands as a testament to the working of NICKI towards fulfilling our objective defined in Equation 5.

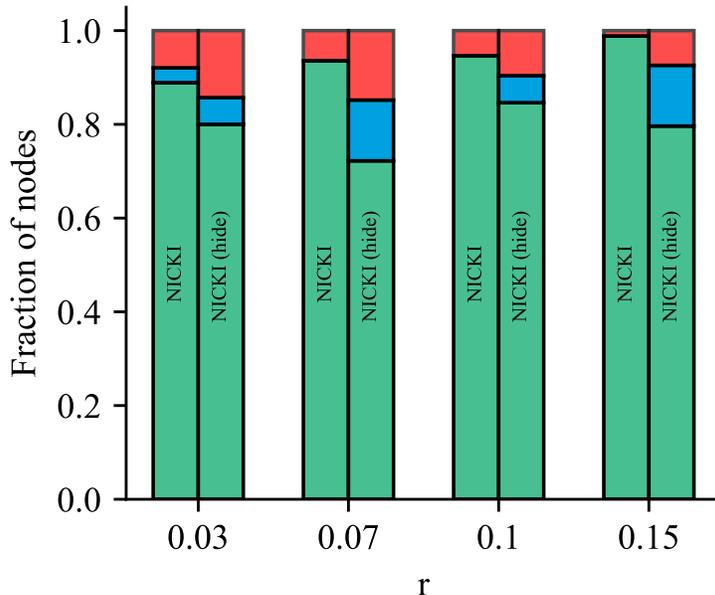


Figure 3: Fraction of target nodes in their respective predicted classes post network attack for Cora. **Green** indicates nodes misclassified as the base class (our goal). **Blue** indicates nodes correctly classified in target class, and **Red** indicates nodes misclassified into a class other than base class.

Reddit. For the lower budget at $r = 0.03$, AFGSM performs marginally better than NICKI. However, as the budget increases, NICKI outperforms AFGSM by a significant margin (around 7%). AFGSM only generates discrete features, whereas Reddit is a continuous dataset. We still evaluate on discrete features to allow AFGSM to reach its full-attack potential, which we think might be the reason for its superior performance in some cases. NIPA gives random features to attacker nodes which doesn't work well for Reddit, given its continuous features. NICKI (hide) performs better than random and preferential but is not able to outperform AFGSM. We feel that the restriction on feature generation given by the *hiding feature loss* does not allow our model to reach its attack potential in this case.

Varying base and target classes. Table 3 shows the results for varying budgets but fixed target and base classes. To demonstrate the performance of our model for any given pair of target and base classes, we randomly select four such pairs for Cora and present the node classification accuracy in Table 5 for $r = 0.07$. We observe that our model outperforms the best baseline AFGSM across all choices.

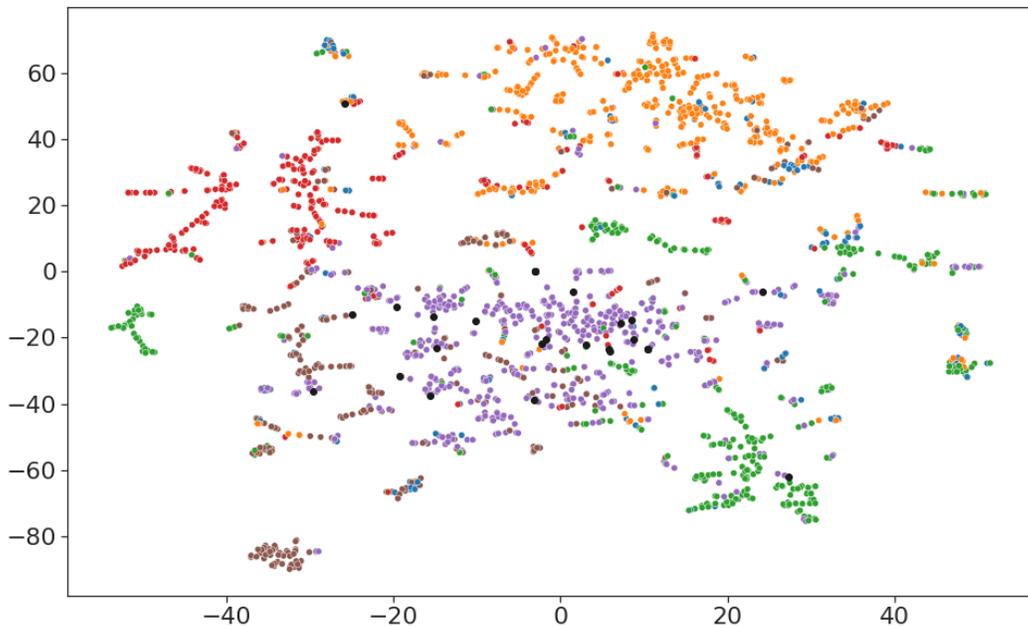


Figure 4: 2D t-SNE plot of Deepwalk embedding for Citeseer. Attacker nodes are black. Nodes in the base and target classes are violet and orange, respectively (see Figure 5 in Appendix for the feature embedding of Citeseer).

6.5. Hiding Efficiency

We test the effectiveness of our hiding strategy by checking the closeness of both features and topological structure assumed by attacker nodes to the nodes in the base class.

For features, we reduce dimensionality using t-SNE. Figure 5 displays the 2D t-SNE plot of representation of the nodes in Citeseer network. As visible, the attacker nodes (black) tend to be near the base class nodes (violet). Since the attack aims to bring the target and base class nodes together, the attacker nodes seem to bridge the nodes belonging to two classes. We can change the position of attacker nodes by tuning α in Equation 18.

For the topological structure, we use Deepwalk (Perozzi et al., 2014) to obtain node embeddings, and then pass it to t-SNE (Figure 4). As visible, our injected attacker nodes seem to lay towards the pre-existing base class nodes' distribution and are good at picking up the topological structure of base class nodes, and thus can be easily camouflaged with them.

Although removing the hiding framework increases the attack performance of NICKI in all cases, there is a trade-off with imperceptibility. With-

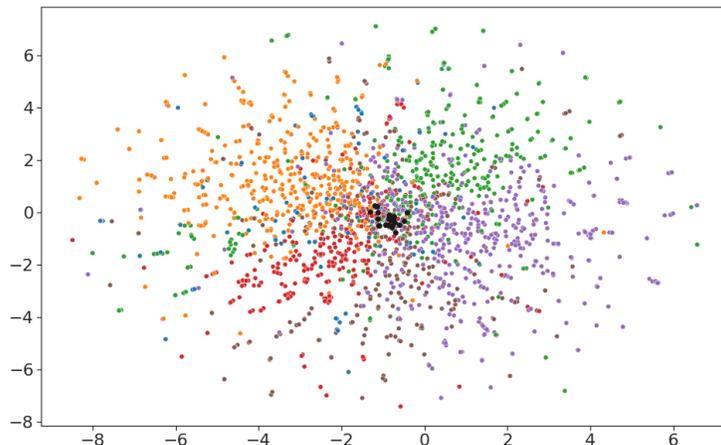


Figure 5: 2D t-SNE plot showing the feature embedding of Citeseer. Attacker nodes are black. Nodes in the base and target classes are violet and orange, respectively.

Table 5: Classification accuracy of target class nodes in test set across randomly selected base and target classes on Cora with $r = 0.07$. Lower value indicates better attack.

Base Class ID	Target Class ID	Clean	NICKI	NICKI AFGSM (hide)	NIPA	
5	0	0.818	0.507	0.727	0.785	0.746
0	2	0.841	0.722	0.837	0.819	0.770
1	6	0.821	0.805	0.793	0.818	0.805
3	4	0.728	0.585	0.685	0.667	0.650

out the hiding framework, the injected nodes might be detected easily since there is no mechanism to ensure their camouflage. Also, our assumption that the injected nodes will be labelled as base-class nodes will become more feasible with hiding since it essentially tries to camouflage injected nodes with base-class nodes only. Better hiding strategy can be explored in the future to improve NICKI (hide)’s performance.

7. Comparison with Unnoticeable Evasion Attacks

We also compare our model with start-of-the-art unnoticeable evasion attacks – HAO Chen et al. (2022) and GANI Fang et al. (2022). We evaluate them in a poison setting by simply training a vanilla GCN once on the attacked graph before evaluating target class accuracy. The idea is to emulate

Table 6: Comparison with Unnoticeable Evasion attacks. Interpretation of results is same as that of Table 3

Dataset	Method	$r = 0.03$	$r = 0.07$	$r = 0.10$	$r = 0.15$
Cora [0.8189]	GANI	0.7966	0.7688	0.7966	0.7715
	AGIA + HAO	0.7855	0.7855	0.8245	0.7910
	NICKI (hide)	0.7409	0.7270	0.7103	0.7214
Citeseer [0.7395]	GANI	0.7368	0.7447	0.7236	0.7368
	TDGIA + HAO	0.7473	0.8184	0.8105	0.8289
	NICKI (hide)	0.7131	0.7263	0.7316	0.7316

baselines in a real-world poison setting in which perturbations happen prior to training. For this purpose, we use two datasets – Cora and Citeseer as both baselines scale well on them. Since HAO Chen et al. (2022) works as an unnoticeability inducing extension for other node-injection attacks, we use AGIA Chen et al. (2022) and TDGIA Zou et al. (2021) as base attacks for Cora and Citeseer, respectively, since they worked best in Chen et al. (2022) for the respective datasets. Additionally, we convert GANI Fang et al. (2022) into a targeted attack by optimising its genetic loss function only on target-class nodes. We present the results in Table 6. To conduct a fair comparison in terms of unnoticeable attacks, we only include NICKI (hide) results in Table 6.

7.1. Performance Comparison

We can observe from the Table 6 that NICKI (hide) outperforms in almost every case, showing that even in the presence of unnoticeability strategies, our model works better and compromises on accuracy much lesser than the other two baselines. GANI Fang et al. (2022) still performs at-par with our model; however HAO performs significantly worse. This can be attributed to two reasons - (i) HAO works on a significantly higher feature perturbation budget which we limited by taking the top-k scores, leading to a reduction in performance. (ii) Evasion attacks in general perform worse in a poison setting as they are not optimised for training like NICKI is (Equation 17).

8. Transferability of Attack

In our approach, we choose GCN (Kipf and Welling, 2017) as our surrogate model and use it to evaluate the poisoned graph, the results of which are presented in Table 3. However in a real-world scenario, an attacker has no

Table 7: Evaluation of attacks under robust defenses. Base, target class and the interpretation of results are same as that of Table 4.

Dataset	Method	G	NNGuard	E	GNNGuard
		$r = 0.03$	$r = 0.10$	$r = 0.03$	$r = 0.10$
Cora [0.4484, 0.8161]	GANI	0.4568	0.4540	0.8161	0.8077
	AGIA + HAO	0.3231	0.3454	0.7688	0.7883
	NICKI	0.3788	0.1838	0.7493	0.4038
	NICKI (hide)	0.2228	0.1559	0.5515	0.0863
Citeseer [0.8973, 0.7289]	GANI	0.8947	0.8789	0.7236	0.7078
	TDGIA + HAO	0.8973	0.9	0.7289	0.7421
	NICKI	0.6763	0.3421	0.6868	0.0473
	NICKI (hide)	0.3421	0	0.0473	0

prior knowledge of the classifier being used. Therefore, it only makes sense to verify the efficacy of our attack on other widely-used node classifiers, for which we select GAT (Veličković et al., 2017) and GraphSAGE (Hamilton et al., 2017). Note that since ours is a poison attack, we first train the obtained poisoned graph on the above mentioned classifiers and then evaluate the accuracy of target nodes. Results can be found in Table 4. We can observe that NICKI performs best in most cases followed by NICKI (hide); this shows that our attack can be effectively transferred to other GNN based classifiers.

9. NICKI Against Defense

In Table 4, we evaluate our attack on common GNN techniques - GAT and GraphSage. Realistically, institutions often employ more robust algorithms to prevent adversarial attacks from happening. Motivated by Chen et al. (2022) that homophily defenders work well against node-injection attacks, we also evaluate the efficacy of our attack against more robust GNN algorithms and homophily defenders - GNNGuard Zhang and Zitnik (2020) and EGNNGuard Chen et al. (2022).

■ **GNNGuard.** This is a general algorithm used to improve robustness of any GNN model. GNNGuard assigns higher weightage to edges between similar nodes and prunes edges between dissimilar nodes. The new edges provide robust message passing to mitigate the effect of attacks.

■ **EGNNGuard.** This is a more scalable version of GNNGuard based on similar ideology and tries to maintain homophily in the graph. It uses a different pruning method as it simply puts a threshold and removes edges between neighbour nodes with low similarity.

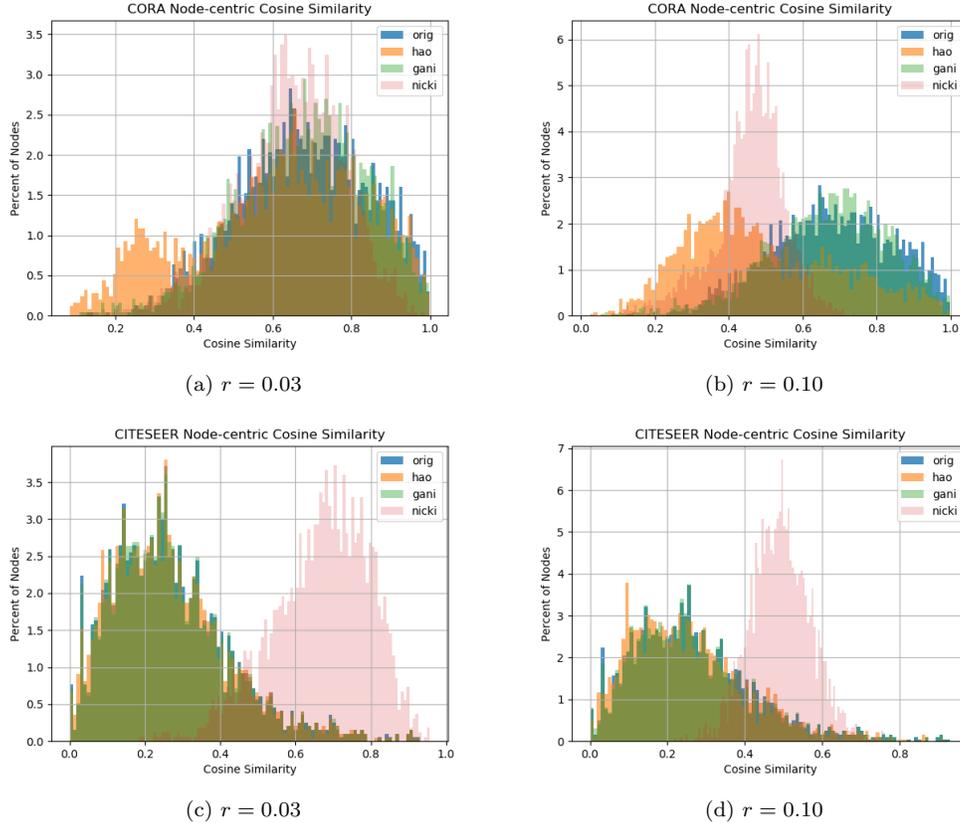


Figure 6: Plots to measure homophily change for different attacks on Cora and Citeseer. Note that for NICKI we presented results with hide. For Cora, HAO is used with AGIA and with TDGIA for Citeseer.

Table 7 shows that our attack is not affected by the above mentioned robust defenses. In fact, the accuracies are in most cases lower (better attack) than the ones in Table 6 in which we used a vanilla GCN. Since both these methods purge edges based on node similarities, we further plot pairwise node similarities and perform a homophily analysis to investigate further.

9.1. Homophily Analysis

We plot node similarity for the graphs generated by all the unnoticeable attacks on Cora and Citeseer on two different budgets in Figure 6. As expected at lower budgets, each plot seems to move towards the clean graph. GANI does an excellent job at keeping the distributions almost overlapped

Table 8: Statistics – Gini coefficient (GC), distribution entropy (DE), power-law exponent (PE), and triangle count (TC) of different graphs after poisoning (NICKI;AFGSM). The values at $r = 0$ indicate the statistics of the clean graph.

	r	GC	DE	PE	TC
Cora	0	0.405	0.360	1.932	1630
	0.03	0.403;0.405	0.361;0.361	1.926;1.927	1633;1636
	0.07	0.402;0.405	0.364;0.364	1.918;1.921	1634;1642
	0.1	0.401;0.406	0.365;0.366	1.911;1.917	1635;1640
	0.15	0.399;0.408	0.367;0.369	1.902;1.911	1634;1640
Citeseer	0	0.427	0.352	2.058	1083
	0.03	0.424;0.428	0.354;0.352	2.047;2.062	1084;1083
	0.07	0.424;0.430	0.357;0.354	2.037;2.058	1086;1083
	0.10	0.424;0.431	0.360;0.355	2.031;2.054	1085;1086
	0.15	0.429;0.432	0.355;0.357	2.044;2.045	1114;1086
Polblogs	0	0.599	0.842	1.478	25649
	0.03	0.593;0.587	0.848;0.852	1.472;1.466	25811;25653
	0.07	0.588;0.576	0.856;0.858	1.465;1.458	25954;25661
	0.10	0.568;0.569	0.867;0.862	1.446;1.453	25742;25665
	0.15	0.558;0.560	0.873;0.869	1.437;1.446	25864;25681
Reddit	0	0.479	0.439	1.693	6850
	0.03	0.479;0.479	0.439;0.439	1.692;1.692	6850;6853
	0.07	0.477;0.478	0.440;0.439	1.689;1.691	6855;6855
	0.10	0.476;0.478	0.440;0.440	1.688;1.690	6856;6857
	0.15	0.475;0.478	0.441;0.440	1.686;1.689	6855;6862

across both graphs and budgets. However, in case of HAO, the graph is left shifted for Cora, i.e., node injection is destroying homophily. Interestingly, NICKI (hide) either creates a peak (in Cora) or shifts right (in Citeseer). Shifting right is analogous to an increase in homophily which is unexpected by graph injection attacks and thus fools the defense model. Similarly peak creation leads to most node pairs centred around the same similarity value and makes it difficult to pick the notorious edges. In fact, we believe that because of this nature, the GNN defenses purge benign edges which leads to even better attack efficacy and sometimes complete disruption of the graph. We can verify these findings from Table 7 as well.

10. Post Attack Statistics

Following Sun et al. (2020), we obtain the poisoned graphs from NICKI and AFGSM, and compare them with their corresponding clean graphs by presenting some key network statistics in Table 8. The change in the degree

distribution of a graph is considered as an important metric to measure *unnoticeability* of an attack (Zügner et al., 2018). It is interesting to observe that for our attack, the power-law exponent of the poisoned graph is similar to that of the clean graph. This shows that our pre- and post-attack graphs have similar degree distributions; indicating unnoticeability. As expected, the triangle count increases as more nodes are injected. However, note that the increment is minor in most cases, which can be attributed to our very small budget, leading to less number of new triangles being induced by the attacker nodes. AFGSM shows a similar trend, which is interesting since NICKI is able to achieve much better results (Table 8) while maintaining statistics similar to that of AFGSM. For the other statistics, we can observe that the numbers remain almost consistent even when the budget increases. Again, we impute this consistency to our minimal budget, which as we can see, helps in the overall unnoticeability of our attack.

11. Conclusion

In this paper, we introduced a novel problem – class-specific network poisoning, which unlike existing methods, aims to poison a network in such a way that the nodes in the target class get misclassified. We addressed this problem by introducing NICKI, a novel attacking strategy that leverages an optimization-based approach to deteriorate the performance of a node classifier. Extensive experiments on four real-world datasets showed significant performance gain over four baselines in terms of misclassifying the nodes in the target class. We also showed that the attack graph and clean graph look alike in terms of the several topological properties of the networks. We further empirically showed that the attacker nodes resemble benign nodes.

References

- Adamic, L.A., Glance, N., 2005. The political blogosphere and the 2004 us election: divided they blog, in: Proceedings of the 3rd international workshop on Link discovery, pp. 36–43.
- Barabási, A.L., Albert, R., 1999. Emergence of scaling in random networks. *science* 286, 509–512.
- Bojchevski, A., Günnemann, S., 2019. Adversarial attacks on node embeddings via graph poisoning, in: ICML, pp. 695–704.

- Chen, Y., Yang, H., Zhang, Y., Ma, K., Liu, T., Han, B., Cheng, J., 2022. Understanding and improving graph injection attack by promoting unnoticeability. arXiv preprint arXiv:2202.08057 .
- Dai, H., Li, H., Tian, T., Huang, X., Wang, L., Zhu, J., Song, L., 2018. Adversarial attack on graph structured data, in: ICML, pp. 1115–1124.
- Erdos, P., Rényi, A., et al., 1960. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci* 5, 17–60.
- Fan, W., Ma, Y., Li, Q., He, Y., Zhao, E., Tang, J., Yin, D., 2019. Graph neural networks for social recommendation, in: WWW, pp. 417–426.
- Fang, J., Wen, H., Wu, J., Xuan, Q., Zheng, Z., Tse, C.K., 2022. Gani: Global attacks on graph neural networks via imperceptible node injections. arXiv preprint arXiv:2210.12598 .
- Giles, C.L., Bollacker, K.D., Lawrence, S., 1998. Citeseer: An automatic citation indexing system, in: JCDL, p. 89–98.
- Goodfellow, I.J., Shlens, J., Szegedy, C., 2015. Explaining and harnessing adversarial examples, in: ICLR.
- Grover, A., Leskovec, J., 2016. node2vec: Scalable feature learning for networks, in: SIGKDD, pp. 855–864.
- Hamilton, W.L., Ying, R., Leskovec, J., 2017. Inductive representation learning on large graphs, in: NIPS, pp. 1025–1035.
- Jia, R., Liang, P., 2017. Adversarial examples for evaluating reading comprehension systems, in: EMNLP, pp. 2021–2031.
- Jiang, M., Li, Z., Zhang, S., Wang, S., Wang, X., Yuan, Q., Wei, Z., 2020. Drug–target affinity prediction using graph neural network and contact maps. *RSC Advances* 10, 20701–20712.
- Jin, D., Liu, Z., Li, W., He, D., Zhang, W., 2019. Graph convolutional networks meet markov random fields: Semi-supervised community detection in attribute networks, in: AAI, pp. 152–159.
- Kipf, T.N., Welling, M., 2016. Variational graph auto-encoders. *CoRR* abs/1611.07308.

- Kipf, T.N., Welling, M., 2017. Semi-supervised classification with graph convolutional networks, in: ICLR, pp. 1–14.
- Li, Y., Jin, W., Xu, H., Tang, J., 2020. Deeprobust: A pytorch library for adversarial attacks and defenses. arXiv preprint arXiv:2005.06149 .
- McCallum, A.K., Nigam, K., Rennie, J., Seymore, K., 2000. Automating the construction of internet portals with machine learning. *Information Retrieval* 3, 127–163.
- Pennington, J., Socher, R., Manning, C.D., 2014. Glove: Global vectors for word representation, in: EMNLP, pp. 1532–1543.
- Perozzi, B., Al-Rfou, R., Skiena, S., 2014. Deepwalk: Online learning of social representations, in: SIGKDD, pp. 701–710.
- Shafahi, A., Huang, W.R., Najibi, M., Suciu, O., Studer, C., Dumitras, T., Goldstein, T., 2018. Poison frogs! targeted clean-label poisoning attacks on neural networks, in: NIPS, p. 6106–6116.
- Sohn, K., Lee, H., Yan, X., 2015. Learning structured output representation using deep conditional generative models. *NIPS* 28, 3483–3491.
- Sun, Y., Wang, S., Tang, X., Hsieh, T.Y., Honavar, V., 2020. Adversarial attacks on graph neural networks via node injections: A hierarchical reinforcement learning approach, in: The WebConf, pp. 673–683.
- Tao, S., Cao, Q., Shen, H., Huang, J., Wu, Y., Cheng, X., 2021. Single node injection attack against graph neural networks, in: CIKM, pp. 1794–1803.
- Tao, S., Cao, Q., Shen, H., Wu, Y., Hou, L., Cheng, X., 2022. Adversarial camouflage for node injection attack on graphs. arXiv preprint arXiv:2208.01819 .
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y., 2017. Graph attention networks. arXiv preprint arXiv:1710.10903 .
- Wang, J., Luo, M., Suya, F., Li, J., Yang, Z., Zheng, Q., 2020. Scalable attack on graph data by injecting vicious nodes. *Data Mining and Knowledge Discovery* 34, 1363–1389.

- Wang, X., Cheng, M., Eaton, J., Hsieh, C.J., Wu, F., 2018. Attack graph convolutional networks by adding fake nodes. arXiv preprint arXiv:1810.10751 .
- Xie, S.M., Ermon, S., 2019. Reparameterizable subset sampling via continuous relaxations. arXiv preprint arXiv:1901.10517 .
- Xu, B., Shen, H., Cao, Q., Cen, K., Cheng, X., 2019a. Graph convolutional networks using heat kernel for semi-supervised learning, in: IJCAI, p. 1928–1934.
- Xu, B., Shen, H., Cao, Q., Qiu, Y., Cheng, X., 2019b. Graph wavelet neural network. arXiv preprint arXiv:1904.07785 .
- Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W.L., Leskovec, J., 2018. Graph convolutional neural networks for web-scale recommender systems, in: SIGKDD, pp. 974–983.
- Zhang, X., Zitnik, M., 2020. Gnn-guard: Defending graph neural networks against adversarial attacks. Advances in neural information processing systems 33, 9263–9275.
- Zou, X., Zheng, Q., Dong, Y., Guan, X., Kharlamov, E., Lu, J., Tang, J., 2021. Tdgia: Effective injection attacks on graph neural networks. arXiv preprint arXiv:2106.06663 .
- Zügner, D., Akbarnejad, A., Günnemann, S., 2018. Adversarial attacks on neural networks for graph data, in: SIGKDD, pp. 2847–2856.
- Zügner, D., Günnemann, S., 2019. Adversarial attacks on graph neural networks via meta learning. arXiv preprint arXiv:1902.08412 .