CERTIFIED DEFENSE ON THE FAIRNESS OF GRAPH NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph Neural Networks (GNNs) have emerged as a prominent graph learning model in various graph-based tasks over the years. Nevertheless, due to the vulnerabilities of GNNs, it has been empirically proved that malicious attackers could easily corrupt the fairness level of their predictions by adding perturbations to the input graph data. In this paper, we take crucial steps to study a novel problem of certifiable defense on the fairness level of GNNs. Specifically, we propose a principled framework named ELEGANT and present a detailed theoretical certification analysis for the fairness of GNNs. ELEGANT takes any GNNs as its backbone, and the fairness level of such a backbone is theoretically impossible to be corrupted under certain perturbation budgets for attackers. Notably, ELEGANT does not have any assumption over the GNN structure or parameters, and does not require re-training the GNNs to realize certification. Hence it can serve as a plug-and-play framework for any optimized GNNs ready to be deployed. We verify the satisfactory effectiveness of ELEGANT in practice through extensive experiments on real-world datasets across different backbones of GNNs, where ELEGANT is also demonstrated to be beneficial for GNN debiasing.

025 026 027

003

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

Graph Neural Networks (GNNs) have emerged among the most popular models to handle learning tasks on graphs (Kipf & Welling, 2017; Veličković et al., 2018) and made remarkable achievements in various domains (Feng et al., 2022; Li et al., 2022). Nevertheless, as GNNs are increasingly deployed in real-world decision-making scenarios, there has been an increasing societal concern on the fairness of GNN predictions. A primary reason is that most traditional GNNs do not consider fairness, and thus could exhibit bias against certain demographic subgroups. Here the demographic subgroups are usually divided by certain sensitive attributes, such as gender and race. To prevent GNNs from biased predictions, multiple recent studies have proposed fairness-aware GNNs (Agarwal et al., 2021; Dai & Wang, 2021; Kang et al., 2022a) such that potential bias could be mitigated.

Unfortunately, despite existing efforts towards fair GNNs, it remains difficult to prevent the corruption of their fairness level due to their common vulnerability of lacking adversarial robustness. In fact, 040 malicious attackers can easily corrupt the fairness level of GNNs by perturbing the node attributes 041 (i.e., changing the values of node attributes) and/or the graph structure (i.e., adding and deleting 042 edges) (Hussain et al., 2022), which could lead to serious consequences in the test phase (Dai 043 & Wang, 2021; Hussain et al., 2022). For example, GNNs have been leveraged to perform bail 044 decision-making on the graph of defendants, where an edge between two defendants represents high 045 profile similarity (Agarwal et al., 2021). Yet, by simply injecting adversarial links in the graph data, 046 attackers can make GNNs deliver advantaged predictions for a subgroup (e.g., individuals with a 047 certain nationality) while damaging the interest of others (Hussain et al., 2022). Hence achieving 048 defense over the fairness of GNNs is crucial for the purpose of safe deployment.

It is worth noting that despite the abundant empirical defense strategies for GNNs (Zhang & Zitnik, 2020; Entezari et al., 2020; Jin & Zhang, 2019; Jin et al., 2020c; Wu et al., 2019b), they are always subsequently defeated by novel attacking techniques (Schuchardt et al., 2020; Carlini & Wagner, 2017), and the defense over the fairness of GNNs also faces the same problem. Therefore, an ideal way is to achieve certifiable defense on fairness (i.e., certified fairness defense). A few recent

054 works aim to certify the fairness for traditional deep learning models (Khedr & Shoukry, 2022; 055 Kang et al., 2022b; Jin et al., 2022; Mangold et al., 2022; Borca-Tasciuc et al., 2022; Ruoss et al., 056 2020). Nevertheless, most of them require specially designed training strategies (Khedr & Shoukry, 057 2022; Jin et al., 2022; Ruoss et al., 2020) and thus cannot be directly applied to optimized GNNs 058 ready to be deployed. More importantly, they mostly rely on assumptions on the optimization results (Khedr & Shoukry, 2022; Jin et al., 2022; Borca-Tasciuc et al., 2022; Ruoss et al., 2020) or 060 data distributions (Kang et al., 2022b; Mangold et al., 2022) over a continuous input space. Hence 061 they can hardly be generalized to GNNs due to the binary nature of the input graph topology. Several 062 other works propose certifiable GNN defense approaches to achieve theoretical guarantee (Wang et al., 2021; Bojchevski & Günnemann, 2019; Bojchevski et al., 2020; Jin et al., 2020a; Zügner & 063 Günnemann, 2019; 2020). However, they mainly focus on securing the GNN prediction for a certain 064 individual node to ensure model utility, ignoring the fairness defense over the entire population. 065 Therefore, despite the significance, the study in this field still remains in its infancy. 066

067 It is worth noting that achieving certifiable defense on the fairness of GNNs is a daunting task 068 due to the following key challenges: (1) Generality: different types of GNNs could be designed 069 and optimized for different real-world applications (Zhou et al., 2020). Correspondingly, our first 070 challenge is to design a plug-and-play framework that can achieve certified defense on fairness for 071 any optimized GNN models that are ready to be deployed. (2) Vulnerability: a plethora of existing 072 studies have empirically verified that most GNNs are sensitive to input data perturbations (Zhang & 073 Zitnik, 2020; Zügner et al., 2020; Xu et al., 2019). In other words, small input perturbations may cause significant changes in the GNN output. Hence our second challenge is to properly mitigate the 074 common vulnerabilities of GNNs without changing its structure or re-training. (3) Multi-Modality: 075 the input data of GNNs naturally bears multiple modalities. For example, there are node attributes 076 and graph topology in the widely studied attributed networks. In practice, both data modalities may 077 be perturbed by malicious attackers. Therefore, our third challenge is to achieve certified defenses of 078 fairness on both data modalities for GNNs. 079

080 As an early attempt to address the aforementioned challenges, in this paper, we propose a principled 081 framework named ELEGANT (cEtifiabLE GNNs over the fAirNess of PredicTions). Specifically, 082 we focus on the widely studied task of node classification and formulate a novel research problem 083 of Certifying GNN Classifiers on Fairness. To handle the first challenge, we propose to develop 084 ELEGANT on top of an optimized GNN model without any assumptions over its structure or parameters. Hence ELEGANT is able to serve as a plug-and-play framework for any optimized GNN 085 model ready to be deployed. To handle the second challenge, we propose to leverage randomized smoothing (Wang et al., 2021; Cohen et al., 2019) to defend against malicious attacks, where most 087 GNNs can then be more robust over the prediction fairness level. To handle the third challenge, we propose two different strategies working in a concurrent manner, such that certified defense against the attacks on both the node attributes (i.e., add and subtract attribute values) and graph topology 090 (i.e., flip the existence of edges) can be realized. Finally, we evaluate the effectiveness of ELEGANT 091 on multiple real-world network datasets. In summary, our contributions are three-fold: (1) Problem 092 **Formulation.** We formulate and make an initial investigation on a novel research problem of Certifying GNN Classifiers on Fairness. (2) Algorithm Design. We propose a framework ELEGANT 094 to achieve certified fairness defense against attacks on both node attributes and graph structure 095 without relying on assumptions about any specific GNNs. (3) Experimental Evaluation. We 096 perform comprehensive experiments on real-world datasets to verify the effectiveness of ELEGANT. 097

098 099

100

2 PROBLEM DEFINITION

101 **Preliminaries.** Let $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ be an undirected attributed network, where $\mathcal{V} = \{v_1, ..., v_n\}$ is the 102 set of *n* nodes; $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges. Let $\mathbf{A} \in \{0, 1\}^{n \times n}$ and $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the adjacency 103 matrix and attribute matrix of \mathcal{G} , respectively. Assume each node in \mathcal{G} represents an individual, 104 and sensitive attribute *s* divides the population into different demographic subgroups. We follow a 105 widely studied setting (Agarwal et al., 2021; Dai & Wang, 2021) to assume the sensitive attribute is 106 binary, i.e., $s \in \{0, 1\}$. We use s_i to denote the value of the sensitive attribute for node v_i . In node 107 classification tasks, we use \mathcal{V}_{trn} and \mathcal{V}_{tst} ($\mathcal{V}_{trn}, \mathcal{V}_{tst} \in \mathcal{V}$) to represent the training and test node set, respectively. We denote the GNN node classifier as f_{θ} parameterized by θ . f_{θ} takes A and X as input, and outputs \hat{Y} as the predictions for the nodes in \mathcal{G} . Each row in \hat{Y} is a one-hot vector flagging the predicted class. We use f_{θ^*} to denote the GNN with optimal parameter θ^* .

Threat Model. We focus on the attacking scenario of model evasion, i.e., the attack happens 111 in the test phase. In particular, we assume that the victim model under attack is an optimized 112 GNN node classifier f_{θ^*} . We follow a widely adopted setting (Bojchevski & Günnemann, 2019; 113 Zügner & Günnemann, 2019; Ma et al., 2020; Mu et al., 2021) to assume that a subset of nodes 114 $\mathcal{V}_{vul} \in \mathcal{V}_{tst}$ are vulnerable to attacks. Specifically, attackers may perturb their links (i.e., flip the 115 edge existence) to other nodes and/or their node attributes (i.e., change their attribute values). We 116 denote the perturbations on adjacency matrix as $A \oplus \Delta_A$. Here \oplus denotes the element-wise XOR 117 operator; $\Delta_A \in \{0,1\}^{n \times n}$ is the matrix representing the perturbations made by the attacker, where 1 118 only appears in rows and columns associated with the vulnerable nodes while 0 appears elsewhere. 119 Correspondingly, in Δ_A , 1 entries represent edges that attackers intend to flip, while 0 entries are 120 associated with edges that are not attacked. Similarly, we denote the perturbations on node attribute 121 matrix as $X + \Delta_X$, where $\Delta_X \in \mathbb{R}^{n \times n}$ is the matrix representing the perturbations made by the 122 attacker. Usually, if the total magnitude of perturbations is within certain budgets (i.e., $\|\Delta_A\|_0 \le \epsilon_A$ 123 for A and $\|\Delta_X\|_2 \leq \epsilon_X$ for X), the perturbations are regarded as unnoticeable. The goal of an 124 attacker is to add unnoticeable perturbations to nodes in \mathcal{V}_{vul} , such that the GNN predictions for nodes 125 in \mathcal{V}_{tst} based on the perturbed graph exhibit as much bias as possible. In addition, we assume that the 126 attacker has access to any information about the victim GNN (i.e., a white-box setting). This is the 127 worst case in practice, which makes it even more challenging to achieve defense.

To defend against the aforementioned attacks, we aim to establish a node classifier on top of an optimized GNN backbone, such that this classifier, theoretically, will not exhibit more bias than a given threshold no matter what unnoticeable perturbations (i.e., perturbations within budgets) are added. We formally formulate the problem of *Certifying GNN Classifiers on Fairness* below.

Problem 1. Certifying GNN Classifiers on Fairness. Given an attributed network \mathcal{G} , a test node set \mathcal{V}_{tst} , a vulnerable node set $\mathcal{V}_{vul} \in \mathcal{V}_{tst}$, a threshold η for the exhibited bias, and an optimized GNN classifier f_{θ^*} , our goal is to achieve a classifier on top of f_{θ^*} associated with budgets ϵ_A and ϵ_X , such that this classifier will bear comparable utility with f_{θ^*} but provably not exhibit more bias than η on the nodes in \mathcal{V}_{tst} , no matter what unnoticeable node attributes and/or graph structure perturbations (i.e., perturbations within budgets) are made over the nodes in \mathcal{V}_{vul} .

139

3 Methodology

140 141

Here we first introduce the modeling of attack and defense on the fairness of GNNs, then discuss how
we achieve certified defense on node attributes. After that, we propose a strategy to achieve both types
of certified defense (i.e., defense on node attributes and graph structure) at the same time. Finally, we
introduce strategies to achieve the designed certified fairness defense for GNNs in practice.

146 147

148

3.1 BIAS INDICATOR FUNCTION

149 We first construct an indicator g to mathematically model the attack and defense on the fairness of 150 GNNs. Our rationale is to use g to indicate whether the predictions of f_{θ^*} exhibit a level of bias 151 exceeding a given threshold. We present the formal definition below.

Definition 1. (Bias Indicator Function) Given adjacency matrix A and node attribute matrix X, a test node set V_{tst} , a threshold η for the exhibited bias, and an optimized GNN model f_{θ^*} , the bias indicator function is defined as $g(f_{\theta^*}, A, X, \eta, V_{tst}) = \mathbb{1}(\pi(f_{\theta^*}(A, X), V_{tst}) < \eta)$, where $\mathbb{1}(\cdot)$ takes an event as input and outputs 1 if the event happens (otherwise 0); $\pi(\cdot, \cdot)$ denotes any bias metric for GNN predictions (taken as its first parameter) over a set of nodes (taken as its second parameter). Traditional bias metrics include Δ_{SP} (Dai & Wang, 2021; Dwork et al., 2012) and Δ_{EO} (Dai & Wang, 2021; Hardt et al., 2016).

159

160 Correspondingly, the goal of the attacker is to ensure that the indicator g outputs 0 for an η as large 161 as possible, while the goal of certified defense is to ensure for a given threshold η , the indicator gprovably yields 1 as long as the attacks are within certain budgets. Note that a reasonable η should ensure that g outputs 1 based on the clean graph data (i.e., graph data without any attacks). Below we first discuss the certified fairness defense over node attributes to maintain the output of g as 1.

164 165 166

3.2 Certified Fairness Defense over Node Attributes

We now introduce how we achieve certified defense over the node attributes for the fairness of the predictions yielded by f_{θ^*} . Specifically, we propose to construct a smoothed bias indicator function $\tilde{g}_{\boldsymbol{X}}(f_{\theta^*}, \boldsymbol{A}, \boldsymbol{X}, \mathcal{V}_{vul}, \eta)$ via adding Gaussian noise over the node attributes of vulnerable nodes in \mathcal{V}_{vul} . For simplicity, we use $\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A}, \boldsymbol{X})$ to represent the smoothed bias indicator function over node attributes by omitting \mathcal{V}_{vul} , f_{θ^*} and η . We formally define $\tilde{g}_{\boldsymbol{X}}$ below.

Definition 2. (Bias Indicator with Node Attribute Smoothing) We define the bias indicator with smoothed node attributes over the nodes in \mathcal{V}_{vul} as $\tilde{g}_{\mathbf{X}}(\mathbf{A}, \mathbf{X}) = \operatorname{argmax}_{c \in \{0,1\}} \Pr(g(f_{\theta^*}, \mathbf{A}, \mathbf{X} + \gamma_{\mathbf{X}}(\boldsymbol{\omega}_{\mathbf{X}}, \mathcal{V}_{vul}), \eta, \mathcal{V}_{tst}) = c)$. Here $\boldsymbol{\omega}_{\mathbf{X}}$ is a $(d \cdot |\mathcal{V}_{vul}|)$ -dimensional vector, where each entry is a random variable following a Gaussian Distribution $\mathcal{N}(0, \sigma^2)$; $\gamma_{\mathbf{X}}(\cdot, \cdot)$ maps a vector (its first parameter) to an $(n \times d)$ -dimensional matrix, where the vector values are assigned to rows with the indices indicated by a set of nodes (its second parameter) while other entries are zeros.

179 We denote $\Gamma_{\mathbf{X}} = \gamma_{\mathbf{X}}(\boldsymbol{\omega}_{\mathbf{X}}, \mathcal{V}_{vul})$ and $g(\mathbf{A}, \mathbf{X} + \Gamma_{\mathbf{X}}) = g(f_{\theta^*}, \mathbf{A}, \mathbf{X} + \gamma_{\mathbf{X}}(\boldsymbol{\omega}_{\mathbf{X}}, \mathcal{V}_{vul}), \eta, \mathcal{V}_{tst})$ below 180 for simplicity. We are then able to derive the theoretical certification for the defense on fairness with 181 the defined $\tilde{g}_{\mathbf{X}}$ in Definition 2. We now present the defense certification on fairness below.

Theorem 1. (*Certified Fairness Defense for Node Attributes*) Denote the probability for $g(\mathbf{A}, \mathbf{X} + \mathbf{\Gamma}_{\mathbf{X}})$ to return class $c (c \in \{0, 1\})$ as P(c). Then $\tilde{g}_{\mathbf{X}}(\mathbf{A}, \mathbf{X})$ will provably return $\operatorname{argmax}_{c \in \{0, 1\}} P(c)$ for any perturbations (over the attributes of vulnerable nodes) within an l_2 radius $\epsilon_{\mathbf{X}} = \frac{\sigma}{2} \left(\Phi^{-1}(\max_{c \in \{0,1\}} P(c)) - \Phi^{-1}(\min_{c \in \{0,1\}} P(c)) \right)$, where $\Phi^{-1}(\cdot)$ is the inverse of the standard Gaussian cumulative distribution function.

Correspondingly, for an η that enables $\max_{c \in \{0,1\}} P(c) = 1$, it is then safe to say that no matter 188 what perturbations Δ_X are made on vulnerable nodes, as long as $\|\Delta_X\|_2 \leq \tilde{\epsilon}_X$, the constructed 189 $\tilde{g}_{\mathbf{X}}$ will provably not yield predictions for \mathcal{V}_{tst} with a level of bias exceeding η . Nevertheless, it is 190 worth noting that, in GNNs, perturbations may also be made on the structure of the vulnerable nodes, 191 i.e., adding and/or deleting edges between these vulnerable nodes and any nodes in the graph. Hence 192 it is also necessary to achieve certified defense against such structural attacks. Here we propose to 193 also smooth the constructed $\tilde{g}_{\mathbf{X}}$ over the graph structure (of the vulnerable nodes) for the purpose of 194 certified fairness defense on the graph structure. However, the adjacency matrix describing the graph 195 structure is naturally binary, and thus should be smoothed in a different way. 196

197 198

187

3.3 CERTIFIED FAIRNESS DEFENSE OVER NODE ATTRIBUTES AND GRAPH STRUCTURE

We then introduce achieving certified fairness defense against attacks on both node attributes and graph structure. We propose a strategy to leverage noise following Bernoulli distribution to smooth \tilde{g}_X over the rows and columns (due to symmetricity) associated with the vulnerable nodes in A. In this way, we can smooth both the node attributes and graph structure for g in a randomized manner, and we denote the constructed function as $\tilde{g}_{A,X}$. We present the formal definition below.

204 **Definition 3.** (Bias Indicator with Attribute-Structure Smoothing) We define the bias indicator 205 function with smoothed node attributes and graph structure over the nodes in \mathcal{V}_{vul} as $\tilde{g}_{A,X}(A,X) =$ 206 $\operatorname{argmax}_{c \in \{0,1\}} \Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \gamma_{\boldsymbol{A}}(\boldsymbol{\omega}_{\boldsymbol{A}}, \mathcal{V}_{vul}), \boldsymbol{X}) = c).$ Here $\boldsymbol{\omega}_{\boldsymbol{A}}$ is an $(n \cdot |\mathcal{V}_{vul}|)$ -dimensional random 207 variable, where each dimension takes 0 and 1 with the probability of β (0.5 < $\beta \leq 1$) and 1 - β , 208 respectively; function $\gamma_{\mathbf{A}}(\cdot, \cdot)$ maps a vector (its first parameter) to a symmetric $(n \times n)$ -dimensional 209 matrix, where the vector values are assigned to rows whose indices associated with the indices of 210 a set of nodes (its second parameter) and then mirrored to the corresponding columns, while other 211 values are left as zeros.

212

213 We let $\Gamma_A = \gamma_A(\omega_A, \mathcal{V}_{vul})$ below for simplicity. To better illustrate how classifier $\tilde{g}_{A,X}$ achieves 214 certified fairness defense over both data modalities of an attributed network, we provide an ex-215 emplary case in Figure 1. Here we assume node $v_i \in \mathcal{V}_{vul}$. Considering the high dimen-316 sionality of node attributes and adjacency matrix, we only analyze two entries $X_{i,j}$ and $A_{i,j}$ and omit other entries after noise for simplicity. Here the superscript (i, j) represents the *i*-th row and *j*-th column of a matrix. Under binary noise, entry $A_{i,j}$ only has two possible values, i.e., $A_{i,j} \oplus 0$ and $A_{i,j} \oplus 1$. We denote the two cases as Case (1) and Case (2), respectively. We assume that the area where *g* returns 1 in the span of the two input random entries of *g* (i.e., $X_{i,j}$ and $A_{i,j}$ under random noise) is an ellipse (marked out with green), where the decision boundary is marked out with deep green. In Case (1), $X_{i,j}$ under random noise follows a Gaussian distribution, whose probability density function is marked out as deep red.



223

224

225

226

227

228

229

230

231

232

233

235

236

237

238

239

241 242

243

244

245 246

257

258

259 260

261

262

263 264

265

266 267 268 (marked out with shallow red) is larger than 0.5. Correspondingly, according to Definition 2, \tilde{g}_X returns 1 in this case. In Case (2), we similarly mark out the probability density function and the area used for integral within the range of the ellipse. We assume that in this case, the integral is smaller than 0.5, and thus \tilde{g}_X returns 0. Note that to compute the output of $\tilde{g}_{A,X}$, we need to identify the output of \tilde{g}_X with the largest probability. Notice that $\beta > 0.5$, we have that $A_{i,j} \oplus 0$ happens with a larger probability than $A_{i,j} \oplus 1$. Therefore, $\tilde{g}_{A,X}$ outputs 1 in this example. In other words, the bias level of the predictions of f_{θ^*} is satisfying (i.e., smaller than η) based on $\tilde{g}_{A,X}$.

We assume that, in this case, the integral of the proba-

bility density function within the range of the ellipse

Figure 1: An example illustrating how ELE-GANT works in the input space.

Below we introduce a desirable property of $\tilde{g}_{A,X}$, i.e., certified fairness defense associated with tractable budgets over both node attributes and graph topology can be achieved.

Lemma 1. (*Perturbation-Invariant Budgets Existence*) There exist tractable budgets ϵ_A and ϵ_X , such that for any perturbations made over the node attributes and graph structure of the vulnerable nodes within ϵ_A and ϵ_X , $\tilde{g}_{A,X}$ provably maintains the same classification results.

Correspondingly, for an η that enables $\tilde{g}_{A,X}$ to return 1, we are then able to achieve certified fairness 247 defense over $\tilde{g}_{A,X}$ against perturbations on both node attributes and graph structure. Below we derive 248 the certified fairness defense budgets over the graph structure ϵ_A and node attributes ϵ_X for $\tilde{g}_{A,X}$. 249 We first introduce the derivation of ϵ_A . Here, our rationale is: considering that $\tilde{g}_{A,X}$ is a binary 250 classifier, we need to ensure that under structure attacks, the probability of \tilde{g}_{X} returning 1 (denoted 251 as $\Pr(\tilde{g}_{X}(A \oplus \Delta_{A} \oplus \Gamma_{A}, X) = 1))$ is provably greater than 0.5, such that $\tilde{g}_{A,X}$ will still return 1. To this end, we propose to derive a lower bound of $Pr(\tilde{g}_X(A \oplus \Delta_A \oplus \Gamma_A, X) = 1)$, which we 253 denote as $P_{\tilde{g}_{X}=1}$. Finally, we identify the largest perturbation size that keeps such a lower bound 254 larger than 0.5, and the identified perturbation size is then the graph structure perturbation budget. 255 We present the lower bound of $\Pr(\tilde{g}_X(A \oplus \Delta_A \oplus \Gamma_A, X) = 1)$ below. 256

Lemma 2. (Positive Probability Bound Under Noises) There exists a tractable $\underline{P}_{\tilde{g}_{X}=1} \in (0, 1)$, such that $\Pr(\tilde{g}_{X}(A \oplus \Delta_{A} \oplus \Gamma_{A}, X) = 1) \ge P_{\tilde{g}_{X}=1}$.

To derive the perturbation budget ϵ_A , we only need to find a Δ_A with the largest l_0 -norm that still enables $\underline{P}_{\tilde{g}_{\boldsymbol{X}}=1}$ to be greater than 0.5 (according to Definition 3). Correspondingly, we derive the theoretical perturbation-invariant budget ϵ_A in Theorem 2 below.

Theorem 2. (Certified Defense Budget for Structure Perturbations) The certified defense budget over the graph structure ϵ_A for $\tilde{g}_{A,X}$ is given as

$$\epsilon_{\boldsymbol{A}} = \max \epsilon_{\boldsymbol{A}}, \text{ s.t. } \underline{P}_{\boldsymbol{\tilde{g}}_{\boldsymbol{X}}=1} > 0.5, \forall \|\boldsymbol{\Delta}_{\boldsymbol{A}}\|_{0} \le \epsilon_{\boldsymbol{A}}.$$
(1)

To solve the optimization problem in Equation (1), we introduce Theorem 3 to compute $P_{\tilde{g}_{X}=1}$.

Theorem 3. (Positive Probability Lower Bound) We have $\underline{P}_{\tilde{g}_{X}=1} = \Pr(A \oplus \Delta_{A} \oplus \Gamma_{A} \in \mathcal{H})$. Here $\mathcal{H} = \bigcup_{i=\mu+1}^{n \cdot |V_{vul}|} \mathcal{H}_{i} \cup \mathcal{H}'_{\mu}$; \mathcal{H}_{i} is given by

$$\mathcal{H}_{i} = \left\{ \bar{A} : \frac{\Pr(A \oplus \Gamma_{A} = \bar{A})}{\Pr(A \oplus \Delta_{A} \oplus \Gamma_{A} = \bar{A})} = \left(\frac{\beta}{1-\beta}\right)^{i}, \forall v_{i} \in \mathcal{V} \setminus \mathcal{V}_{vul}, \|\bar{A}_{i} - A_{i}\|_{0} = 0 \right\};$$

and μ is defined over the optimization problem of $\operatorname{argmax}_{-n \cdot |\mathcal{V}_{vul}| \leq j \leq n \cdot |\mathcal{V}_{vul}|} j$, s.t. $\operatorname{Pr}(\tilde{g}_{\mathbf{X}}(\mathbf{A} \oplus \mathbf{\Gamma}_{\mathbf{A}}, \mathbf{X}) = 1) \leq \operatorname{Pr}\left(\mathbf{A} \oplus \mathbf{\Gamma}_{\mathbf{A}} \in \bigcup_{k=j}^{n \cdot |\mathcal{V}_{vul}|} \mathcal{H}_{j}\right)$. Here \mathcal{H}'_{μ} is any subregion of \mathcal{H}_{μ} that satisfies $\operatorname{Pr}(\mathbf{A} \oplus \mathbf{\Gamma}_{\mathbf{A}} \in \mathcal{H}'_{\mu}) = \operatorname{Pr}(\tilde{g}_{\mathbf{X}}(\mathbf{A} \oplus \mathbf{\Gamma}_{\mathbf{A}}, \mathbf{X}) = 1) - \operatorname{Pr}\left(\mathbf{A} \oplus \mathbf{\Gamma}_{\mathbf{A}} \in \bigcup_{k=j}^{n \cdot |\mathcal{V}_{vul}|} \mathcal{H}_{j}\right)$.

We provide detailed steps to solve the optimization problem given in Equation (1) in Appendix. Now we introduce the theoretical analysis of how to derive ϵ_X in Theorem 4.

Theorem 4. (Certified Defense Budget over Node Attributes) Denote \overline{A} as the set of all possible $(n \times n)$ -matrices, where entries in rows whose indices associate with those vulnerable nodes may take 1 or 0, while other entries are zeros. The certified defense budget $\epsilon_{\mathbf{X}}$ for $\tilde{g}_{\mathbf{A},\mathbf{X}}$ is given as $\epsilon_{\mathbf{X}} = \min\{\epsilon_{\mathbf{X}}^{\mathbf{X}} : \epsilon_{\mathbf{X}}^{\mathbf{X}} \text{ is derived with classifier } \tilde{g}_{\mathbf{X}}(\mathbf{A} \oplus \Gamma_{\mathbf{A}}, \mathbf{X}), \text{ where } \Gamma_{\mathbf{A}} \in \overline{A}\}.$

3.4 CERTIFICATION IN PRACTICE

289 Estimating the Predicted Label Probabilities. According to Definition 3, it is necessary to obtain 290 $\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X}) = c) \ (c \in \{0, 1\})$ to determine the output of classifier $\tilde{g}_{\boldsymbol{X}}$. We propose 291 to leverage a Monte Carlo method to estimate such a probability. Specifically, we first randomly 292 pick N samples of Γ_A as $\overline{\mathcal{A}}'(\overline{\mathcal{A}}' \subset \overline{\mathcal{A}})$. Considering the output of \tilde{g}_X is binary, we then follow 293 a common strategy (Cohen et al., 2019) to consider this problem as a parameter estimation of a 294 Binomial distribution: we first count the number of returned label 1 and 0 under noise as N_1 and 295 N_0 $(N_1 + N_0 = N)$; then we choose a confidence level $1 - \alpha$ and take the α -th quantile of the beta 296 distribution with parameters N_1 and N_0 as the estimated probability lower bound for returning label 297 c = 1. We proved that all theoretical analysis still holds true for such an estimation in Appendix. We 298 follow a similar strategy to estimate the probability lower bound of yielding 1 for $q(\mathbf{A}, \mathbf{X} + \Gamma_{\mathbf{X}})$.

299 **Obtaining Fair Classification Results.** After achieving certified fairness defense based on $\tilde{g}_{A,X}$, we 300 also need to obtain the corresponding node classification results (given by f_{θ^*}) over \mathcal{V}_{tst} . We propose 301 to collect all classification results associated with the sampled $\Gamma'_A \in \bar{\mathcal{A}}'$ that leads to an estimated 302 lower bound of $\Pr(\tilde{g}_{X}(A \oplus \Gamma'_{A}, X) = 1)$ to be larger than 0.5 as $\hat{\mathcal{Y}}'$. Here $\hat{\mathcal{Y}}'$ is a set of output 303 matrices of f_{θ^*} , where each matrix consists of the one-hot output classification results (as each row 304 in the matrix) for all nodes. We propose to take $\operatorname{argmin}_{\hat{Y}'} \pi(\hat{Y}', \mathcal{V}_{tst}), s.t. \hat{Y}' \in \hat{\mathcal{Y}}'$ as the final node 305 classification results. Correspondingly, consider $\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Gamma}'_{\boldsymbol{A}}, \boldsymbol{X}) = 1)$ falls into the confidence 306 interval characterized by $1 - \alpha$, we have a neat probabilistic theoretical guarantee below. 307

Proposition 1. (Probabilistic Guarantee for the Fairness Level of Node Classification). For $\hat{\mathbf{Y}} = \underset{\hat{\mathbf{Y}}'}{\operatorname{argmin}} \hat{\mathbf{Y}}' \pi(\hat{\mathbf{Y}}', \mathcal{V}_{tst}), s.t. \, \hat{\mathbf{Y}}' \in \hat{\mathcal{Y}}', we have \Pr(\pi(\hat{\mathbf{Y}}, \mathcal{V}_{tst}) > \eta) < 0.5^{|\hat{\mathcal{Y}}'|}.$

310

273 274 275

276

277 278

279

281

282

283

284

285

286 287

288

Note that for a large enough sample size N, the cardinality of $\hat{\mathcal{Y}}'$ also tends to be large in practice. Hence it is safe to argue that $\Pr(\pi(\hat{Y}, \mathcal{V}_{tst}) > \eta)$ tends to be small enough. In other words, we have a probability that is large enough to obtain results with a bias level lower than threshold η .

Calculation of Perturbation Budgets. We calculate ϵ_A by solving the optimization problem given in Equation (1), and we provide the completed procedure in Appendix. For ϵ_X , we utilize a Monte Carlo method to estimate its value. More specifically, we leverage min $\{\tilde{\epsilon}_X : \tilde{\epsilon}_X \text{ is derived with classifier } \tilde{g}_X(A \oplus \Gamma'_A, X)$, where $\Gamma'_A \in \tilde{\mathcal{A}}'\}$ to estimate the value of ϵ_X .

318 319

4 EXPERIMENTAL EVALUATIONS

320 321

In this section, we aim to answer three research questions: **RQ1**: How well does ELEGANT perform in achieving certified fairness defense? **RQ2**: How does ELEGANT perform under fairness attacks compared to other popular fairness-aware GNNs? **RQ3**: How does ELEGANT perform under

different settings of parameters? We present the main experimental settings and representative results 325 in this section due to space limits. Detailed settings and supplementary experiments are in Appendix. 326

4.1 EXPERIMENTAL SETTINGS

324

327

328

Downstream Task and Datasets. We focus on the widely studied node classification task, which is 330 one of the most representative tasks in the domain of learning on graphs. We adopt three real-world 331 network datasets that are widely used to perform studies on the fairness of GNNs, namely German 332 Credit (Agarwal et al., 2021; Asuncion & Newman, 2007), Recidivism (Agarwal et al., 2021; Jordan 333 & Freiburger, 2015), and Credit Defaulter (Agarwal et al., 2021; Yeh & Lien, 2009). We provide 334 their basic information, including how these datasets are built and their statistics, in Appendix.

335 **Evaluation Metrics.** We perform evaluation from three main perspectives, including model utility, 336 fairness, and certified defense. To evaluate utility, we adopt the node classification accuracy. To 337 evaluate fairness, we adopt the widely used metrics Δ_{SP} (measuring bias under *Statistical Parity*) and 338 $\Delta_{\rm EO}$ (measuring bias under *Equal Opportunity*). To evaluate certified defense, we extend a traditional 339 metric named Certified Accuracy (Wang et al., 2021; Cohen et al., 2019) in our experiments, and we 340 name it as Fairness Certification Rate (FCR). Specifically, existing GNN certification works mainly 341 focus on a certain individual node, and utilize certified accuracy to measure the ratio of nodes that are 342 correctly classified and also successfully certified out of all test nodes (Wang et al., 2021). In this 343 paper, however, we perform certified (fairness) defense for individuals over an entire test set (instead 344 of for any specific individual). Accordingly, we propose to sample multiple test sets out of nodes that 345 are not involved in the training and validation set. Then we perform certified fairness defense for 346 all sampled test sets, and utilize the ratio of test sets that are successfully certified over all sampled 347 sets as the metric of certified defense. The rationale of FCR is leveraging a Monte Carlo method to estimate the probability of being successfully certified for a randomly sampled test node set. 348

349 **GNN Backbones and Baselines.** Note that ELEGANT serves as a plug-and-play framework for 350 any optimized GNNs ready to be deployed. To evaluate the generality of ELEGANT across GNNs, 351 we adopt three of the most representative GNNs spanning across simple and complex ones, namely 352 Graph Sample and Aggregate Networks (Hamilton et al., 2017) (GraphSAGE), Graph Convolutional 353 Networks (Kipf & Welling, 2017) (GCN), and Jumping Knowledge Networks (JK). Note that to the 354 best of our knowledge, existing works on fairness certification cannot certify the attacks over two 355 data modalities (i.e., continuous node attributes and binary graph topology) at the same time, and 356 thus cannot be naively generalized onto GNNs. Hence we compare the usability of GNNs before and 357 after certification with ELEGANT. Moreover, we also adopt two popular fairness-aware GNNs as 358 baselines to evaluate bias mitigation, including FairGNN (Dai & Wang, 2021) and NIFTY (Agarwal 359 et al., 2021). Specifically, FairGNN utilizes adversarial learning to debias node embeddings, while NIFTY designs regularizations to debias node embeddings. 360

361 **Threat Models.** We propose to evaluate the performance of ELEGANT and other fairness-aware 362 GNN models under actual attacks on fairness. We first introduce the threat model over graph structure. 363 To the best of our knowledge, FA-GNN (Hussain et al., 2022) is the only work that performs graph 364 structure attacks targeting the fairness of GNNs. Hence we adopt FA-GNN to attack graph structure. 365 In terms of node attributes, to the best of our knowledge, no existing work has made any explorations. 366 Hence we directly utilize gradient ascend to perform attacks. Specifically, after structure attacks have 367 been performed, we identify the top-ranked node attribute elements (out of the node attribute matrix) 368 that positively influence the exhibited bias the most via gradient ascend. For any given budget (of attacks) on node attributes, we add perturbations to these elements in proportion to their gradients. 369

370 371

372

4.2 RQ1: FAIRNESS CERTIFICATION EFFECTIVENESS

To answer RQ1, we investigate the performance of different GNNs after certification across different 373 real-world attributed network datasets over FCR, utility, and fairness. We present the experimental 374 results across three GNN backbones and three real-world attributed network datasets in Table 1. 375 Here bias is measured with Δ_{SP} , and we have similar observations on Δ_{EO} . We summarize the main 376 observations as follows: (1) Fairness Certification Rate (FCR). We observe that ELEGANT realizes 377 values of FCR around or even higher than 90% for all three GNN backbones and three attributed

	Ge	erman Cre	dit		Recidivisn	ı	Cr	edit Defau	lter
	ACC (†)	Bias (↓)	FCR (\uparrow)	ACC (†)	Bias (\downarrow)	FCR (†)	ACC (†)	Bias (\downarrow)	FCR (†)
SAGE	67.3 ±2.14	50.6 ±15.9	N/A	89.8 ±0.66	9.36 ±3.15	N/A	75.9 ±2.18	13.0 ±4.01	N/A
E-SAGE	71.0 $_{\pm 1.27}$	16.3 ± 10.9	$98.7_{\ \pm 1.89}$	89.9 ±0.90	$\textbf{6.39}_{\pm 2.85}$	$94.3_{\ \pm 6.65}$	$73.4{\scriptstyle~\pm 0.50}$	8.94 ±0.99	94.3 ±3.30
GCN	59.6 ±3.64	37.4 ±3.24	N/A	90.5 ±0.73	10.1 ±3.01	N/A	65.8 ±0.29	11.1 ±3.22	N/A
E-GCN	$58.2{\scriptstyle~\pm1.82}$	$\textbf{3.52} \scriptstyle \pm 3.77$	96.3 $_{\pm 1.89}$	$89.6{\scriptstyle~\pm 0.74}$	9.56 ±3.22	$96.0_{\ \pm 3.56}$	$65.2_{\pm 0.99}$	7.28 ±1.46	92.7 ±5.19
JK	63.3 ±4.11	$41.2{\scriptstyle~\pm18.1}$	N/A	91.9 $_{\pm 0.54}$	10.1 ±3.15	N/A	76.6 ±0.69	9.24 ±0.60	N/A
E-JK	62.3 ±4.07	22.4 ±1.95	97.0 ±3.00	89.3 ±0.33	6.26 ±2.78	89.5 ±10.5	77.7 ±0.27	3.37 ±2.64	99.3 ±0.47

Table 1: Comparison between vanilla GNNs and certified GNNs under ELEGANT over three popular
GNNs across three real-world datasets. Here ACC denotes node classification accuracy, and E- prefix
marks out the GNNs under ELEGANT with certification. ↑ denotes the larger, the better; ↓ represents
the opposite. Numerical values are in percentage, and the best ones are in bold.

network datasets, especially for the German Credit dataset, where vanilla GNNs tend to exhibit a high 396 level of bias. The corresponding intuition is that, for nodes in any randomly sampled test set, we have 397 a probability around or higher than 90% to successfully certify the fairness level of the predictions 398 yielded by the GNN model with our proposed framework ELEGANT. Hence ELEGANT achieves a 399 satisfying fairness certification rate across all adopted GNN backbones and datasets. (2) Utility. We 400 found that compared with those vanilla GNN backbones, certified GNNs with ELEGANT also exhibit 401 comparable and even higher node classification accuracy values in all cases. Hence we conclude that 402 our proposed framework ELEGANT does not significantly jeopardize the utility of the vanilla GNN 403 models, and those certified GNNs with ELEGANT still bear a high level of usability in terms of 404 node classification accuracy. (3) Fairness. Although the goal of ELEGANT is not debiasing GNNs, 405 we observe that certified GNNs with ELEGANT achieve better performances in all cases in terms 406 of algorithmic fairness compared with those vanilla GNNs. This demonstrates that the proposed 407 framework ELEGANT also contributes to bias mitigation. We conjecture that such an advantage of 408 debiasing could be a mixed result of (1) adding random noise on node attributes and graph topology 409 (as in Section 3.2 and Section 3.3) and (2) the proposed strategy of obtaining fair classification results 410 (as in Section 3.4). We provide a more detailed analysis in Appendix B.9.

411 412 413

414

382

391 392 393

4.3 RQ2: FAIRNESS CERTIFICATION UNDER ATTACKS

To answer RQ2, we perform attacks on the fairness of GCN, E-GCN, FairGNN (with a GCN backbone), and NIFTY (with a GCN backbone). Considering the large size of the quadratic space spanned by the sizes of perturbations Δ_A and Δ_X , we present the evaluation under four representative ($\|\Delta_A\|_0$, $\|\Delta_X\|_2$) pairs. We set the threshold for bias η to be 50% higher than the fairness level of the vanilla GCN model on clean data, since it empirically helps to achieve a high certification success rate under large perturbations.

421 We present the fairness levels of the four models in terms of $\Delta_{\rm FO}$ in Figure 2. Note that we utilize 422 a vanilla GCN to predict the labels for test nodes to obtain the classification results discussed in 423 Section 3.4, and we also have similar observations on other GNNs/datasets. (1) Fairness. We found 424 that the GCN model with the proposed framework ELEGANT achieves the lowest level of bias in all cases of fairness attacks. This observation is consistent with the superiority in fairness found 426 in Table 1, which demonstrates that the fairness superiority of ELEGANT maintains even under attacks within a wide range of attacking perturbation sizes. (2) Certification on Fairness. We now 427 compare the performance of E-GCN across different attacking perturbation sizes. We observed 428 that under relatively small attacking perturbation sizes, i.e., $(2^{0}, 10^{-1})$, $(2^{1}, 10^{0})$, and $(2^{2}, 10^{1})$, 429 ELEGANT successfully achieves certification over fairness, and the bias level increases slowly as the 430 size of attacks increases. Under relatively large attacking perturbation size, i.e., $(2^3, 10^2)$, although 431 the attacking budgets go beyond the certified budgets, GCN under ELEGANT still exhibits a fairness



439 440 441

446

447 448

449

(a) FCR of certification for σ over node attributes (b) FCR of certification for β over graph topology

Figure 3: Parameter study of σ over ϵ_X (a) and β over ϵ_A (b). Experimental results are presented based on GCN over German credit and Credit Defaulter for (a) and (b), respectively. Similar tendencies can also be observed based on other GNNs and datasets.

level far lower than the given bias threshold η , and the fairness superiority maintains. Hence the adopted estimation strategies are safe in achieving fairness certification.

4.4 RQ3: PARAMETER STUDY

To answer RQ3, we propose to perform pa-450 451 rameter study focusing on two most critical parameters, σ and β . To examine how σ 452 and β influence the effectiveness of ELE-453 GANT in terms of both FCR and certified 454 defense budgets, we set numerical ranges 455 for $\epsilon_{\boldsymbol{X}}$ (from 0 to 1e1) and $\epsilon_{\boldsymbol{A}}$ (from 0 to 456 2^4) and divide the two ranges into grids. 457 In both ranges, we consider the dividing 458 values of the grids as thresholds for certifi-459 cation budgets. In other words, under each 460 threshold, we only consider the test sets 461 with the corresponding certified defense 462 budget being larger than this threshold as 463 successfully certified ones, and the values 464 of FCR are re-computed accordingly. Our 465 rationale here is that with the thresholds 466 (for $\epsilon_{\mathbf{X}}$ and $\epsilon_{\mathbf{A}}$) increasing, if FCR reduces 467 slowly, this demonstrates that most success-



Figure 2: The bias levels of GCN, E-GCN, FairGNN, and NIFTY under fairness attacks on German Credit. The shaded bar indicates that certified budget $\epsilon_A \leq \|\Delta_A\|_0$ or $\epsilon_X \leq \|\Delta_X\|_2$. The y-axis is in logarithmic scale for better visualization purposes.

fully certified test sets are associated with large certified defense budgets. However, if FCR reduces
 fast, then most successfully certified test sets only bear small certified defense budgets.

470 Here we present the experimental results of σ and β with the most widely used GCN model based on 471 German Credit in Figure 3(a) and Credit Defaulter in Figure 3(b), respectively. We also have similar 472 observations on other GNNs and datasets. We summarize the main observations as follows: (1) 473 Analysis on σ . We observe that most cases with larger σ are associated with a larger FCR compared 474 with the cases where σ is relatively small. In other words, larger values of σ typically make FCR 475 reduce slower w.r.t. the increasing of $\epsilon_{\mathbf{X}}$ threshold. This indicates that increasing the value of σ 476 helps realize larger certified defense budgets on node attributes, i.e., the increase of σ dominates 477 the tendency of ϵ_X given in Theorem 4. Nevertheless, it is worth mentioning that if σ is too large, 478 the information encoded in the node attributes could be swamped by the Gaussian noise and finally 479 corrupt the classification accuracy. Hence moderately large values for σ , e.g., 5e-1 and 5e0, are recommended. (2) Analysis on β . We found that (1) for cases with relatively large β (e.g., 0.8 and 480 481 0.9), the FCR also tends to be larger (compared with cases where β is smaller) at ϵ_A threshold being 0. Such a tendency is reasonable, since in these cases, the expected magnitude of the added Bernoulli 482 noise is small. Correspondingly, GNNs under ELEGANT perform similarly to vanilla GNNs, and 483 thus an η larger than the bias level of vanilla GNNs is easier to be satisfied (compared with cases 484 under smaller values of β ; (2) for cases with relatively large β , the value of FCR reduces faster (w.r.t. 485 ϵ_A threshold) than cases where β is smaller. Therefore, we recommend that for any test set of nodes:

486 (1) if the primary goal is to achieve certification with a high probability, then larger values for β (e.g., 487 0.8 and 0.9) would be preferred; (2) if the goal is to achieve certification with larger certified defense 488 budgets on the graph topology, then smaller values for β (e.g., 0.6 and 0.7) should be selected.

489 490 491

5 **RELATED WORK**

492 493

495

496

497

498

499

500

501

Algorithmic Fairness in GNNs. Existing GNN works on fairness mainly focus on group fairness and 494 individual fairness (Dong et al., 2022b). Specifically, group fairness requires that each demographic subgroup (divided by sensitive attributes such as gender and race) in the graph should have their fair share of interest based on predictions (M. et al., 2021). Adversarial training is among the most popular strategies (Dai & Wang, 2021; Dong et al., 2022b). In addition, regularization (Agarwal et al., 2021; Fan et al., 2021; Zhang et al., 2021), topology modification (Dong et al., 2022a; Spinelli et al., 2021), and orthogonal projection (Palowitch & Perozzi, 2020) are also commonly used strategies. On the other hand, individual fairness it requires that similar individuals should be treated similarly (Dwork et al., 2012), where such similarity may be determined in different ways (Kang et al., 2020; Dong 502 et al., 2021). Designing optimization regularization terms to promote individual fairness for GNNs is a common strategy (Fan et al., 2021; Dong et al., 2021; Song et al., 2022). Nevertheless, despite the 504 research advancements in the field of algorithmic fairness on GNNs, the adversarial defense against 505 fairness attacks still remains in its infancy and has not been thoroughly explored. To the best of our 506 knowledge, our paper serves as the first comprehensive study dedicated to addressing this important 507 research problem, paving the way for future investigations in this under-explored area.

508 GNN Defense Against Attacks. Existing works on GNN defense are mainly categorized into four 509 mainstreams, namely adversarial training (Xu et al., 2019; Dai et al., 2019; Wang et al., 2019), 510 graph data purification (Entezari et al., 2020; Jin et al., 2020c; Wu et al., 2019a; Kipf & Welling, 511 2016), perturbation detection (Xu et al., 2018; Ioannidis et al., 2019; Jin et al., 2020b), and certified 512 defense (Schuchardt et al., 2020; Wang et al., 2021; Bojchevski & Günnemann, 2019; Zügner & 513 Günnemann, 2020; Jia et al., 2020). Adversarial training aims to inject adversarial examples (e.g., 514 edges) during training, such that the GNN tends to yield correct predictions for adversarial examples 515 during inference (Xu et al., 2019; Dai et al., 2019; Wang et al., 2019). Graph data purification 516 also works during training, where graph data is purified during learning to weaken the influence of 517 adversarial examples (Entezari et al., 2020; Jin et al., 2020c; Wu et al., 2019a; Kipf & Welling, 2016). 518 Perturbation detection is mostly applied in the pre-processing stage, where adversarial edges or nodes 519 can be identified before training (Xu et al., 2018; Ioannidis et al., 2019; Jin et al., 2020b). Different 520 from them, certified defense is the only approach that secures GNNs theoretically, such that attackers cannot find any adversary to fool the GNNs (Schuchardt et al., 2020; Wang et al., 2021; Bojchevski 521 & Günnemann, 2019; Zügner & Günnemann, 2020; Jia et al., 2020). Note that most certified defense 522 approaches only secure the prediction for a specific data point (e.g., a node in node classification). 523 Different from them, ELEGANT enables us to secure the fairness level for GNNs, which naturally 524 entwines with all predictions in the test set. 525

526 527

CONCLUSION 6

528 529

530 In this paper, we study a novel problem of certifying GNN node classifiers on fairness. To address 531 this problem, we propose ELEGANT, a framework designed to achieve certification on top of any 532 optimized GNN node classifier associated with certain perturbation budgets, ensuring that it is 533 impossible for attackers to degrade the fairness level of predictions within such budgets. Notably, 534 ELEGANT is designed to serve as a plug-and-play framework for any optimized GNNs and does not rely on any assumptions regarding GNN structures or re-training processes. Extensive experiments 535 verify the strong effectiveness and generalizability of ELEGANT across multiple GNN architectures 536 and real-world datasets. While this paper primarily focuses on the widely studied node classification 537 task, we also highlight the potential for extending this study to other graph-related tasks as a future 538 research direction. We expect positive broader impacts including deploying fairness-safe GNNs in applications, and no significant negative broader impact needs to be highlighted here.

540 REFERENCES

579

- 542 Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. Towards a unified framework for fair
 543 and stable graph representation learning. *UAI*, pp. 2114–2124, 2021.
- Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- Aleksandar Bojchevski and Stephan Günnemann. Certifiable robustness to graph perturbations.
 NeurIPS, 32, 2019.
- Aleksandar Bojchevski, Johannes Klicpera, and Stephan Günnemann. Efficient robustness certificates for discrete data: Sparsity-aware randomized smoothing for graphs, images and more. In *ICML*, 2020.
- Giorgian Borca-Tasciuc, Xingzhi Guo, Stanley Bak, and Steven Skiena. Provable fairness for neural network models using formal verification. *arXiv preprint arXiv:2212.08578*, 2022.
- Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *AISec*, pp. 3–14, 2017.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *ICML*, pp. 1725–1735, 2020.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *ICML*, 2019.
- Enyan Dai and Suhang Wang. Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information. In *WSDM*, pp. 680–688, 2021.
- Quanyu Dai, Xiao Shen, Liang Zhang, Qiang Li, and Dan Wang. Adversarial training methods for network embedding. In WWW, pp. 329–339, 2019.
- Yushun Dong, Jian Kang, Hanghang Tong, and Jundong Li. Individual fairness for graph neural networks: A ranking based approach. In *SIGKDD*, pp. 300–310, 2021.
- Yushun Dong, Ninghao Liu, Brian Jalaian, and Jundong Li. EDITS: modeling and mitigating data
 bias for graph neural networks. In WWW, pp. 1259–1269, 2022a.
- Yushun Dong, Jing Ma, Chen Chen, and Jundong Li. Fairness in graph mining: A survey. *arXiv* preprint arXiv:2204.09888, 2022b.
- 574 Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through
 575 awareness. In *ITCS*, pp. 214–226, 2012.
- Negin Entezari, Saba A Al-Sayouri, Amirali Darvishzadeh, and Evangelos E Papalexakis. All you need is low (rank) defending against adversarial attacks on graphs. In *WSDM*, 2020.
 - Wei Fan, Kunpeng Liu, Rui Xie, Hao Liu, Hui Xiong, and Yanjie Fu. Fair graph auto-encoder for unbiased graph representations with wasserstein distance. In *ICDM*, pp. 1054–1059, 2021.
- Shangbin Feng, Zhaoxuan Tan, Herun Wan, Ningnan Wang, Zilong Chen, Binchi Zhang, Qinghua Zheng, Wenqian Zhang, Zhenyu Lei, Shujie Yang, et al. Twibot-22: Towards graph-based twitter bot detection. *arXiv preprint arXiv:2206.04564*, 2022.
- Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *NeurIPS*, volume 30, 2017.
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *NeurIPS*, volume 29, 2016.
- Hussain Hussain, Meng Cao, Sandipan Sikdar, Denis Helic, Elisabeth Lex, Markus Strohmaier, and
 Roman Kern. Adversarial inter-group link injection degrades the fairness of graph neural networks. arXiv preprint arXiv:2209.05957, 2022.

594 595	Vassilis N Ioannidis, Dimitris Berberidis, and Georgios B Giannakis. Graphsac: Detecting anomalies in large-scale graphs. <i>arXiv preprint arXiv:1910.09589</i> , 2019.
596 597	Jinyuan Jia, Binghui Wang, Xiaoyu Cao, and Neil Zhenqiang Gong. Certified robustness of commu-
598	<i>Conf</i> pp. 2718–2724 2020
599	coly, pp. 2710-2721, 2020.
600	Hongwei Jin and Xinhua Zhang. Latent adversarial training of graph convolution networks. In <i>ICML</i>
601	workshop on learning and reasoning with graph-structured representations, 2019.
603	Hongwei Jin, Zhan Shi, Venkata Jaya Shankar Ashish Peruri, and Xinhua Zhang. Certified robustness
604	of graph convolution networks for graph classification under topological attacks. <i>NeurIPS</i> , 2020a.
605 606	Jiayin Jin, Zeru Zhang, Yang Zhou, and Lingfei Wu. Input-agnostic certified group fairness via gaussian parameter smoothing. In <i>ICML</i> , pp. 10340–10361. PMLR, 2022.
607 608 609 610	Wei Jin, Yaxin Li, Han Xu, Yiqi Wang, Shuiwang Ji, Charu Aggarwal, and Jiliang Tang. Adversarial attacks and defenses on graphs: A review, a tool and empirical studies. <i>arXiv preprint arXiv:2003.00653</i> , 2020b.
611 612	Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, and Jiliang Tang. Graph structure learning for robust graph neural networks. In <i>SIGKDD</i> , 2020c.
613	Wei Iin Tyler Derr Yigi Wang Yao Ma Zitao Liu and Iiliang Tang. Node similarity preserving
614	graph convolutional networks. In <i>Proceedings of the 14th ACM international conference on web</i>
615	search and data mining, pp. 148–156, 2021.
616	Kareem L Jordan and Tina L Freiburger. The effect of race/ethnicity on sentencing: Examining
617	sentence type, jail length, and prison length. J Ethn Crim Justice, 13(3):179–196, 2015.
618	Lion Kang, Lingrui He. Does Magisiawski, and Hanghang Tang. Informy Individual fairness on graph
620	mining. In <i>SIGKDD</i> , pp. 379–389, 2020.
621 622 623	Jian Kang, Yan Zhu, Jiebo Luo, Yinglong Xia, and Hanghang Tong. Rawlsgcn: Towards rawlsian difference principle on graph convolutional network. In <i>WWW</i> , pp. 1214–1225, 2022a.
624 625 626	Mintong Kang, Linyi Li, Maurice Weber, Yang Liu, Ce Zhang, and Bo Li. Certifying some distributional fairness with subpopulation decomposition. <i>arXiv preprint arXiv:2205.15494</i> , 2022b.
627 628	Haitham Khedr and Yasser Shoukry. Certifair: A framework for certified global fairness of neural networks, 2022.
629 630	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2015.
631	Thomas N Kipf and Max Welling Variational graph auto-encoders arXiv preprint arXiv:1611.07308
632	2016.
633	Themas N. Kinf and May Walling. Some supervised elessification with month convolutional networks
634	In ICLR, 2017.
635	
636	Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:
638	Representations ICLR 2019 New Orleans I.A. USA May 6-9 2019 2019
639	Representations, 1021(201), 1100 Orientis, 221, 0511, 1107 0 9, 2019, 2019.
640 641	Michelle M Li, Kexin Huang, and Marinka Zitnik. Graph representation learning in biomedicine and healthcare. <i>Nat. Biomed. Eng.</i> , pp. 1–17, 2022.
642 643	Ninareh M., Fred M., Nripsuta S., Kristina L., and Aram G. A survey on bias and fairness in machine learning. <i>CSUR</i> , 54(6):115:1–115:35, 2021.
644	Jiagi Ma Shuangrui Ding and Ojaozhu Mei. Towards more practical adversarial attacks on graph
645	neural networks. <i>NeurIPS</i> , 33:4756–4766, 2020.
646	Doub Manageld Michaël Domot Augilier Dellet and Man Terrary 'D'Constitution to the
647	impact on fairness in classification. <i>arXiv preprint arXiv:2210.16242</i> , 2022.

648 649	Jiaming Mu, Binghui Wang, Qi Li, Kun Sun, Mingwei Xu, and Zhuotao Liu. A hard label black-box adversarial attack against graph neural networks. In <i>SIGSAC</i> , pp. 108–125, 2021.
650 651 652	John Palowitch and Bryan Perozzi. Debiasing graph representations via metadata-orthogonal training. ASONAM, 2020.
653 654 655	Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In <i>NeurIPS</i> , 2017.
656 657 658	Anian Ruoss, Mislav Balunovic, Marc Fischer, and Martin Vechev. Learning certified individually fair representations. <i>NeurIPS</i> , 33:7584–7596, 2020.
659 660	Jan Schuchardt, Aleksandar Bojchevski, Johannes Gasteiger, and Stephan Günnemann. Collective robustness certificates: Exploiting interdependence in graph neural networks. In <i>ICLR</i> , 2020.
661 662	Weihao Song, Yushun Dong, Ninghao Liu, and Jundong Li. Guide: Group equality informed individual fairness in graph neural networks. In <i>SIGKDD</i> , pp. 1625–1634, 2022.
663 664 665	Indro Spinelli, Simone Scardapane, Amir Hussain, and Aurelio Uncini. Biased edge dropout for enhancing fairness in graph representation learning. <i>TAI</i> , 3(3):344–354, 2021.
666 667	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. In <i>ICLR</i> , 2018.
668 669 670	Binghui Wang, Jinyuan Jia, Xiaoyu Cao, and Neil Zhenqiang Gong. Certified robustness of graph neural networks against adversarial structural perturbation. In <i>SIGKDD</i> , 2021.
671 672	Xiaoyun Wang, Xuanqing Liu, and Cho-Jui Hsieh. Graphdefense: Towards robust graph convolutional networks. <i>arXiv preprint arXiv:1911.04429</i> , 2019.
673 674	Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. Adversarial examples for graph data: Deep insights into attack and defense. In <i>IJCAI</i> , pp. 4816–4823, 2019a.
675 676 677 678	Huijun Wu, Chen Wang, Yuriy Tyshetskiy, Andrew Docherty, Kai Lu, and Liming Zhu. Adversarial examples on graph data: Deep insights into attack and defense. <i>arXiv preprint arXiv:1903.01610</i> , 2019b.
679 680 681	Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue Lin. Topology attack and defense for graph neural networks: An optimization perspective. <i>arXiv</i> preprint arXiv:1906.04214, 2019.
682 683	Xiaojun Xu, Yue Yu, Bo Li, Le Song, Chengfeng Liu, and Carl Gunter. Characterizing malicious edges targeting on graph neural networks. 2018.
684 685 686	I-Cheng Yeh and Che-hui Lien. The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. <i>Expert Syst. Appl.</i> , 36(2):2473–2480, 2009.
687 688	Xiang Zhang and Marinka Zitnik. Gnnguard: Defending graph neural networks against adversarial attacks. <i>NeurIPS</i> , 2020.
689 690	Xu Zhang, Liang Zhang, Bo Jin, and Xinjiang Lu. A multi-view confidence-calibrated framework for fair and stable graph representation learning. In <i>ICDM</i> , pp. 1493–1498, 2021.
692 693 694	Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. <i>AI Open</i> , 1:57–81, 2020.
695 696	Daniel Zügner and Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks. In <i>SIGKDD</i> , 2019.
697 698 699	Daniel Zügner and Stephan Günnemann. Certifiable robustness of graph convolutional networks under structure perturbations. In <i>SIGKDD</i> , 2020.
700 701	Daniel Zügner, Oliver Borchert, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on graph neural networks: Perturbations and their patterns. <i>TKDD</i> , 2020.

Appendix

Table of Contents

Α	Proofs
	A.1 Proof of Theorem 1
	A.2 Proof of Lemma 1
	A.3 Proof of Lemma 2
	A.4 Proof of Theorem 2
	A.5 Proof of Theorem 3
	A.6 Solving the Optimization Problem in Theorem 2
	A.7 Proof of Theorem 4
	A.8 Proof of Proposition 1
	A.9 Rationale of Each Theoretical Result
В	Reproducibility and Supplementary Analysis
	B.1 Datasets
	B.2 Detailed Experimental Settings
	B.3 Algorithmic Routine
	B.4 Evaluation of Model Utility
	B.5 Certification under Different Fairness Metrics
	B.6 Ordering the Inner and Outer Defense
	B.7 Certification with Estimated Probabilities
	B.8 Time Complexity Analysis
	B.9 Additional Results on Different GNN Backbones & Baselines
	B.10 Complementary Results
С	Additional Discussion
	C.1 Why Certify A Classifier on top of An Optimized GNN?
	C.2 What Is the Difference Between the Attacking Performance of GNNs and the Eatrness of GNNs?
	C 3 Certification Without Considering the Binary Sensitive Attribute
	C 4 How Do the Main Theoretical Findings Differ From Existing Works on Robustness
	Certification of GNNs on Regular Attacks?
	C.5 Discussion: Difference with Existing Similar Works
	C.6 Additional Experiments on Different Datasets
	*

Proofs А

For better clarity, for a matrix X, we use X[i, j] to denote the element at the *i*-th row and the *j*-th column; for a vector \boldsymbol{x} , we use $\boldsymbol{x}[i]$ to denote its *i*-th component.

A.1 PROOF OF THEOREM 1

To prove Theorem 1, we formulate the theoretical prerequisite that Theorem 1 relies on as Lemma A 1. Similarly, the proof of Lemma A 1 relies on the results in Lemma A 2, and the proof of Lemma A 2 is based on Lemma A 3.

Proof. For simplification, we reshape the matrix $X \in \mathbb{R}^{n \times d}$ to the vector $x \in \mathbb{R}^{nd}$. We denote $h_c(\mathbf{x}) = \Pr(g(f_{\theta^*}, \mathbf{A}, \mathbf{X} + \gamma_{\mathbf{X}}(\boldsymbol{\omega}_{\mathbf{X}}, \mathcal{V}_{vul}), \eta, \mathcal{V}_{tst}) = c)$ as the function that returns P(c). Without loss of generality, we assume $\operatorname{argmax}_{c \in \{0,1\}} P(c) = 1$, so we have $\operatorname{argmin}_{c \in \{0,1\}} P(c) = 0$ conse-quently. We use Δ_X to denote a perturbation on X that satisfies $\Delta_X \leq \tilde{\epsilon_X}$, and Δ_x to denote the reshaped vector of Δ_X . Denote $\Phi^{-1}(\cdot)$ as the inverse of the standard Gaussian cumulative distribution function. According to Lemma A 1, we have $\Phi^{-1}(h_c(x))$ as a Lipschitz continuous function with a Lipschitz constant of $\frac{1}{\sigma}$ where σ is the standard deviation of the Gaussian noise ω_X . Based on the property of Lipschitz continuous functions, we have

$$|\Phi^{-1}(h_c(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) - \Phi^{-1}(h_c(\boldsymbol{x}))| \le \frac{\epsilon_{\boldsymbol{X}}}{\sigma}$$

Correspondingly, we have the following bounds for the output probabilities of class 0 and 1

$$\Phi^{-1}(h_1(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) - \Phi^{-1}(h_1(\boldsymbol{x})) \ge -\frac{\epsilon_{\boldsymbol{X}}}{\sigma},$$
(2)

$$\Phi^{-1}(h_0(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) - \Phi^{-1}(h_0(\boldsymbol{x})) \le \frac{\epsilon_{\boldsymbol{X}}}{\sigma}.$$
(3)

Combine Equation (2) and Equation (3), and we have

$$\Phi^{-1}(h_1(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) - \Phi^{-1}(h_0(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) \ge \Phi^{-1}(h_1(\boldsymbol{x})) - \Phi^{-1}(h_0(\boldsymbol{x})) - \frac{2\tilde{\epsilon_{\boldsymbol{X}}}}{\sigma}.$$
 (4)

Recall that $\tilde{\epsilon_X} = \frac{\sigma}{2} (\Phi^{-1}(\max_{c \in \{0,1\}} P(c)) - \Phi^{-1}(\min_{c \in \{0,1\}} P(c))) = \frac{\sigma}{2} (\Phi^{-1}(P(1)) - \Phi^{-1}(\min_{c \in \{0,1\}} P(c)))$ $\Phi^{-1}(P(0))) = \frac{\sigma}{2}(\Phi^{-1}(h_1(\boldsymbol{x})) - \Phi^{-1}(h_0(\boldsymbol{x}))))$, combine this condition with Equation (4), we have

$$\Phi^{-1}(h_1(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) - \Phi^{-1}(h_0(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}})) \ge 0.$$
(5)

Based on the strictly non-decreasing property of $\Phi(\cdot)$, we have

$$h_1(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}}) \ge h_0(\boldsymbol{x} + \boldsymbol{\Delta}_{\boldsymbol{x}}).$$
(6)

In Equation (6), $h_1(x + \Delta_x)$ and $h_0(x + \Delta_x)$ stand for the output probabilities of class 1 and 0 after the perturbation, correspondingly. Hence, the output for \tilde{g}_{X} will not change after the perturbation (still class 1). Noting that the exact probabilities $\max_{c \in \{0,1\}} P(c)$ and $\min_{c \in \{0,1\}} P(c)$ are difficult to calculate in practice, we can use a tractable lower bound p_{\max} and upper bound $\overline{p_{\min}}$ such that $\max_{c \in \{0,1\}} P(c) \geq \underline{p_{\max}} \geq \overline{p_{\min}} \geq \min_{c \in \{0,1\}} P(c)$ to replace them in $\tilde{\epsilon_X}$ as $\tilde{\epsilon_X} = \frac{\sigma}{2} (\Phi^{-1}(p_{\text{max}}) - \Phi^{-1}(\overline{p_{\text{min}}}))$. Because the practical perturbation budget $\tilde{\epsilon_X}$ derived by tractable bounds is smaller than the true budget, we can still obtain the same result as Equation (6).

Lemma A 1. Denote $h_c(\mathbf{x}) = \Pr(g(f_{\theta^*}, \mathbf{A}, \mathbf{X} + \gamma_{\mathbf{X}}(\boldsymbol{\omega}_{\mathbf{X}}, \mathcal{V}_{vul}), \eta, \mathcal{V}_{tst}) = c)$ as the function that returns P(c). Then, the function $\Phi^{-1}(h_c(\mathbf{x}))$ is a Lipschitz continuous function with respect to \mathbf{x} with a Lipschitz constant $L_{\Phi} = \frac{1}{\sigma}$, where $\Phi^{-1}(\cdot)$ is the inverse of the standard Gaussian cumulative distribution function.

Proof. To prove the Lipschitz continuity of $\Phi^{-1}(h_c(\boldsymbol{x}))$, we should find an upper bound of the norm of the gradient $\|\nabla_{\boldsymbol{x}} \Phi^{-1}(h_c(\boldsymbol{x}))\|_2$, denoted as L_{Φ} . The gradient $\nabla_{\boldsymbol{x}} \Phi^{-1}(h_c(\boldsymbol{x}))$ is computed as

$$\nabla_{\boldsymbol{x}} \Phi^{-1}(h_c(\boldsymbol{x})) = \frac{\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})}{\Phi'(\Phi^{-1}(h_c(\boldsymbol{x})))}$$

$$= \sqrt{2\pi} \exp(\frac{1}{2} \Phi^{-1}(h_c(\boldsymbol{x}))^2) \nabla_{\boldsymbol{x}} h_c(\boldsymbol{x}).$$

Therefore, the norm $\|\nabla_{\boldsymbol{x}} \Phi^{-1}(h_c(\boldsymbol{x}))\|_2$ is computed as

$$\|\nabla_{\boldsymbol{x}} \Phi^{-1}(h_c(\boldsymbol{x}))\|_2 = \sqrt{2\pi} \exp(\frac{1}{2} \Phi^{-1}(h_c(\boldsymbol{x}))^2) \|\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})\|_2.$$

According to Lemma A 2, the upper bound of $\|\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})\|_2$ is $\frac{1}{\sqrt{2\pi\sigma}} \exp(-\frac{1}{2}\Phi^{-1}(h_c(\boldsymbol{x}))^2)$. Consequently, we have

$$\|\nabla_{\boldsymbol{x}}\Phi^{-1}(h_c(\boldsymbol{x}))\|_2 \leq \frac{1}{\sigma}.$$

Finally, we have obtained the Lipschitz constant of $\Phi^{-1}(h_c(x))$ as $L_{\Phi} = \frac{1}{\sigma}$ and verified its Lipschitz continuity.

Lemma A 2. Denote $h_c(\mathbf{x}) = \Pr(g(f_{\theta^*}, \mathbf{A}, \mathbf{X} + \gamma_{\mathbf{X}}(\omega_{\mathbf{X}}, \mathcal{V}_{vul}), \eta, \mathcal{V}_{tst}) = c)$ as the function that returns P(c). Then, the function $h_c(x)$ is a Lipschitz continuous function with respect to x with a Lipschitz constant $L_h = \frac{1}{\sqrt{2\pi\sigma}} \exp(-\frac{1}{2}\Phi^{-1}(h_c(\boldsymbol{x}))^2)$, where $\Phi^{-1}(\cdot)$ is the inverse of the standard Gaussian cumulative distribution function.

Proof. To prove the Lipschitz continuity of $h_c(\mathbf{x})$, we should prove that the norm of the gradient $\|\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})\|_2$ is bounded by some constant L_h , i.e. $L_h = \sup_h \|\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})\|_2$. Let $\boldsymbol{\omega}_{vul} =$ $\gamma_{\boldsymbol{X}}(\boldsymbol{\omega}_{\boldsymbol{X}}, \mathcal{V}_{\text{vul}}) \in \mathbb{R}^{nd}$ where $\boldsymbol{\omega}_{\text{vul}}[i] \sim \mathcal{N}(0, \sigma^2)$ when $i \in \mathcal{V}_{\text{vul}}$ and $\boldsymbol{\omega}_{\text{vul}}[i] = 0$ otherwise. Consequently, we have $h_c(\mathbf{x}) = \Pr(g(f_{\theta^*}, \mathbf{A}, \mathbf{x} + \boldsymbol{\omega}_{vul}, \eta, \mathcal{V}_{tst}) = 1)$. Then, we compute the gradient $\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})$ as follows.

$$\begin{aligned} \nabla_{\boldsymbol{x}} h_c(\boldsymbol{x}) &= \nabla_{\boldsymbol{x}} \Pr(g(f_{\boldsymbol{\theta}^*}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}}) = c) \\ &= \nabla_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}}[g(f_{\boldsymbol{\theta}^*}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}})] \end{aligned}$$

$$= \nabla_{\boldsymbol{x}} \int_{\mathbb{R}^{|\mathcal{V}_{\text{vul}}|}} g(f_{\boldsymbol{\theta}^*}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}}) (2\pi\sigma^2)^{-\frac{|\mathcal{V}_{\text{vul}}|}{2}} \exp(-\frac{\|\boldsymbol{\omega}_{\text{vul}}\|_2^2}{2\sigma^2}) d\boldsymbol{\omega}_{\text{vul}}.$$

Substituting $t = x + \omega_{\text{vul}}$ into the above integration, we have

$$\begin{split} \nabla_{\boldsymbol{x}} h_{c}(\boldsymbol{x}) &= \nabla_{\boldsymbol{x}} \int_{\mathbb{R}^{|\mathcal{V}_{\text{vul}}|}} g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{t}, \eta, \mathcal{V}_{\text{tst}}) (2\pi\sigma^{2})^{-\frac{|\mathcal{V}_{\text{vul}}|}{2}} \exp(-\frac{\|\boldsymbol{t}-\boldsymbol{x}\|_{2}^{2}}{2\sigma^{2}}) d\boldsymbol{t} \\ &= \int_{\mathbb{R}^{|\mathcal{V}_{\text{vul}}|}} g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{t}, \eta, \mathcal{V}_{\text{tst}}) (2\pi\sigma^{2})^{-\frac{|\mathcal{V}_{\text{vul}}|}{2}} \nabla_{\boldsymbol{x}} \exp(-\frac{\|\boldsymbol{t}-\boldsymbol{x}\|_{2}^{2}}{2\sigma^{2}}) d\boldsymbol{t} \\ &= \int_{\mathbb{R}^{|\mathcal{V}_{\text{vul}}|}} g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{t}, \eta, \mathcal{V}_{\text{tst}}) (2\pi\sigma^{2})^{-\frac{|\mathcal{V}_{\text{vul}}|}{2}} \exp(-\frac{\|\boldsymbol{t}-\boldsymbol{x}\|_{2}^{2}}{2\sigma^{2}}) \frac{\boldsymbol{t}-\boldsymbol{x}}{\sigma^{2}} d\boldsymbol{t} \\ &= \int_{\mathbb{R}^{|\mathcal{V}_{\text{vul}}|}} g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}}) (2\pi\sigma^{2})^{-\frac{|\mathcal{V}_{\text{vul}}|}{2}} \exp(-\frac{\|\boldsymbol{\omega}_{\text{vul}}\|_{2}^{2}}{2\sigma^{2}}) \frac{\boldsymbol{\omega}_{\text{vul}}}{\sigma^{2}} d\boldsymbol{\omega}_{\text{vul}} \\ &= \frac{1}{\sigma^{2}} \mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}} [\boldsymbol{\omega}_{\text{vul}} g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}})] \\ &= \frac{1}{\sigma} \mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}} [\boldsymbol{\omega}_{\text{vul}}' g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{x} + \boldsymbol{\omega}_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}})]. \end{split}$$

> Here, ω'_{vul} is a normalized random vector that $\omega'_{vul}[i] \sim \mathcal{N}(0,1)$ when $i \in \mathcal{V}_{vul}$ and $\omega'_{vul}[i] = 0$ otherwise. Next, we compute the norm of the gradient $\|\nabla_{\boldsymbol{x}} h_c(\boldsymbol{x})\|_2$ as

862
863

$$\|\nabla_{\boldsymbol{x}}h_{c}(\boldsymbol{x})\|_{2} = \sup_{\|\boldsymbol{v}\|_{2}=1} \boldsymbol{v}^{\top} \nabla_{\boldsymbol{x}}h_{c}(\boldsymbol{x})$$

$$= \frac{1}{\sigma} \sup_{\|\boldsymbol{v}\|_{2}=1} \mathbb{E}_{\boldsymbol{\omega}'_{\text{vul}}}[\boldsymbol{v}^{\top}\boldsymbol{\omega}'_{\text{vul}}g(f_{\boldsymbol{\theta}^{*}}, \boldsymbol{A}, \boldsymbol{x} + \sigma\boldsymbol{\omega}'_{\text{vul}}, \eta, \mathcal{V}_{\text{tst}})].$$
(7)

L

To find L_h , we should consider the worst case (with the largest Lipschitz constant) among all possible classifiers. We let $\tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}'_{vul}) = g(f_{\boldsymbol{\theta}^*}, \boldsymbol{A}, \boldsymbol{x} + \sigma \boldsymbol{\omega}'_{vul}, \eta, \mathcal{V}_{tst})$. Then, we have the following optimization problem for solving L_h

868

879

883

885 886

887

888

889

890 891 892 $\sup_{\tilde{g}} \mathbb{E}_{\boldsymbol{\omega}_{vul}'}[\boldsymbol{v}^{\top} \boldsymbol{\omega}_{vul}' \tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}_{vul}')]$ s.t. $\tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}_{vul}') \in [0, 1], \|\boldsymbol{v}\|_2 = 1, \mathbb{E}_{\boldsymbol{\omega}_{uu}'}[\tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}_{vul}')] = h_c(\boldsymbol{x}).$ (8)

The rationale of this optimization problem is that we aim to find the function with the largest Lipschitz constant (objective) among all possible classifiers \tilde{g} with the same smoothing output $h_c(x)$ (constraints) when fixing the variable v. After solving this problem, we can find the largest objective among all possible v as L_h . To solve this problem, we have the following lemma.

Based on Lemma A 3, we have $\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}'_{vul}) = \mathbb{1}(\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul} \ge -\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x})))$ as the solution of the optimization problem in Equation (8). Next, we can compute the maximal objective when fixing the variable \boldsymbol{v} as

$$\mathbb{E}_{\boldsymbol{\omega}'_{\mathrm{vul}}}[\boldsymbol{v}^{\top}\boldsymbol{\omega}'_{\mathrm{vul}} \cdot \mathbb{I}(\boldsymbol{v}^{\top}\boldsymbol{\omega}'_{\mathrm{vul}} \geq -\varepsilon_{\boldsymbol{v}}\boldsymbol{\Phi}^{-1}(h_{c}(\boldsymbol{x})))] \\
= \mathbb{E}_{\boldsymbol{\omega}\sim\mathcal{N}(0,\varepsilon_{\boldsymbol{v}}^{2})}[\boldsymbol{\omega}\cdot\mathbb{I}(\boldsymbol{\omega}\geq -\varepsilon_{\boldsymbol{v}}\boldsymbol{\Phi}^{-1}(h_{c}(\boldsymbol{x})))] \\
= \mathbb{E}_{\boldsymbol{\omega}'\sim\mathcal{N}(0,1)}[\varepsilon_{\boldsymbol{v}}\boldsymbol{\omega}'\cdot\mathbb{I}(\boldsymbol{\omega}'\geq -\boldsymbol{\Phi}^{-1}(h_{c}(\boldsymbol{x})))] (\text{let }\boldsymbol{\omega}=\varepsilon_{\boldsymbol{v}}\boldsymbol{\omega}') \\
= \frac{\varepsilon_{\boldsymbol{v}}}{\sqrt{2\pi}} \int_{-\boldsymbol{\Phi}^{-1}(h_{c}(\boldsymbol{x}))}^{+\infty} \boldsymbol{\omega}\exp(-\frac{\boldsymbol{\omega}}{2})d\boldsymbol{\omega} \\
= \frac{\varepsilon_{\boldsymbol{v}}}{\sqrt{2\pi}}\exp(-\frac{1}{2}\boldsymbol{\Phi}^{-1}(h_{c}(\boldsymbol{x}))^{2}).$$
(9)

Therefore, we have $\sup_{\tilde{g}} \mathbb{E}_{\omega'_{vul}}[v^{\top} \omega'_{vul} \tilde{g}(\boldsymbol{A}, \omega'_{vul})] = \frac{\varepsilon_{\boldsymbol{v}}}{\sqrt{2\pi}} \exp(-\frac{1}{2} \Phi^{-1} (h_c(\boldsymbol{x}))^2)$. Combining this result with Equation (7), we have

$$\begin{split} {}_{h} = & \sup_{h} \| \nabla_{\boldsymbol{x}} h_{c}(\boldsymbol{x}) \|_{2} \\ = & \sup_{\tilde{g}, \|\boldsymbol{v}\|_{2}=1} \frac{1}{\sigma} \mathbb{E}_{\boldsymbol{\omega}'_{\text{vul}}} [\boldsymbol{v}^{\top} \boldsymbol{\omega}'_{\text{vul}} \tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}'_{\text{vul}})] \\ = & \sup_{\|\boldsymbol{v}\|_{2}=1} \frac{\varepsilon_{\boldsymbol{v}}}{\sqrt{2\pi\sigma}} \exp(-\frac{1}{2} \Phi^{-1} (h_{c}(\boldsymbol{x}))^{2}) \\ = & \frac{1}{\sqrt{2\pi\sigma}} \exp(-\frac{1}{2} \Phi^{-1} (h_{c}(\boldsymbol{x}))^{2}). \end{split}$$

895 896

900

901

902 903

894

Finally, we have proved that the function $h_c(\boldsymbol{x}) = \Pr(g(f_{\boldsymbol{\theta}^*}, \boldsymbol{A}, \boldsymbol{X} + \gamma_{\boldsymbol{X}}(\boldsymbol{\omega}_{\boldsymbol{X}}, \mathcal{V}_{\text{vul}}), \eta, \mathcal{V}_{\text{tst}}) = c)$ is a Lipschitz continuous function with respect to variable \boldsymbol{x} .

Lemma A 3. The solution to the optimization problem in Equation (8) is $\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}'_{vul}) = \mathbb{1}(\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul} \geq -\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x})))$, where $\varepsilon_{\boldsymbol{v}}^2 = \sum_{i \in \mathcal{V}_{vul}} \boldsymbol{v}[i]^2$.

Proof. First, we clarify the rationale for solving this problem. We note that $v^{\top} \omega'_{vul} \sim \mathcal{N}(0, \varepsilon_v^2)$ (based on the property of independent and identically distributed Gaussian), this optimization problem can be regarded as the reweighting of a Gaussian distribution where the range of the weight function $\tilde{g}(\mathbf{A}, \omega'_{vul})$ is [0, 1] and the constraint of the weight function is given by $\mathbb{E}_{\omega'_{vul}}[\tilde{g}(\mathbf{A}, \omega'_{vul})] = h_c(\mathbf{x})$. A straightforward solution here is to let the weight function at a large value of $v^{\top} \omega'_{vul}$ as large as possible. We let $\tilde{g}(\mathbf{A}, \omega'_{vul}) = 1$ where $v^{\top} \omega'_{vul} \ge -\varepsilon_v \Phi^{-1}(h_c(\mathbf{x}))(\omega'_{vul})$ and $\tilde{g}(\mathbf{A}, \omega'_{vul}) = 0$ otherwise. Here $\Phi(\cdot)$ is the cumulative distribution function of $\mathcal{N}(0, 1)$.

912 Next, we prove that $\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}'_{vul}) = \mathbb{1}(\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul} \ge -\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x})))$ is the exact solution of the 913 optimization problem in Equation (8). We first verify that \tilde{g}^* is a solution. It is obvious that \tilde{g}^* 914 suffices the first two constraints because the range of the indicator function is $\{0, 1\}$. For the last 915 constraint, $\mathbb{E}_{\boldsymbol{\omega}'_{vul}}[\mathbb{1}(\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul} \ge -\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x})))]$ is actually the probability of $\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul}$ being larger 916 than $-\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x}))$, which equals to $(h_c(\boldsymbol{x})$ apparently because $\boldsymbol{v}^\top \boldsymbol{\omega}'_{vul}/\varepsilon_{\boldsymbol{v}} \sim \mathcal{N}(0, 1)$. Therefore, 917 \tilde{g}^* satisfies all three constraints. We then prove that \tilde{g}^* is the optimal solution. We assume $\tilde{g} \neq \tilde{g}^*$ is another classifier that also suffices the constraints in the optimization problem in Equation (8). We use S to denote the support set $\{s \mid \tilde{g}^*(A, \omega'_{vul}) \neq 0\}$. Based on the final constraint in the optimization problem in Equation (8), we have

$$\mathbb{E}_{\boldsymbol{\omega}_{\mathrm{vul}}'}[\tilde{g}^*(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}') - \tilde{g}(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}')] = 0.$$

We divide this equation into two parts as

$$\mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}}[(\tilde{g}^*(\boldsymbol{A},\boldsymbol{\omega}_{\text{vul}}') - \tilde{g}(\boldsymbol{A},\boldsymbol{\omega}_{\text{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\text{vul}}' \in \mathcal{S})] + \mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}'}[(\tilde{g}^*(\boldsymbol{A},\boldsymbol{\omega}_{\text{vul}}') - \tilde{g}(\boldsymbol{A},\boldsymbol{\omega}_{\text{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\text{vul}}' \in \mathcal{S}^{\complement})] = 0,$$
(10)

where \mathcal{S}^{U} denotes the complement set of \mathcal{S} . We know that $\tilde{g}^{*}(\mathbf{A}, \boldsymbol{\omega}_{vul}') \equiv 1$ for $\boldsymbol{\omega}_{vul}' \in \mathcal{S}$. We also know that $\tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}'_{\text{vul}}) \leq 1$ and $\tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}'_{\text{vul}})$ cannot always equal to 1 for $\boldsymbol{\omega}'_{\text{vul}} \in \mathcal{S}$ because $\tilde{g} \neq \tilde{g}^*$. Therefore, we have

$$\mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}}[(\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}_{\text{vul}}') - \tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}_{\text{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\text{vul}}' \in \mathcal{S})] > 0,$$

$$\mathbb{E}_{\boldsymbol{\omega}_{\text{vul}}}[(\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}_{\text{vul}}') - \tilde{g}(\boldsymbol{A}, \boldsymbol{\omega}_{\text{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\text{vul}}' \in \mathcal{S}^{\complement})] < 0.$$
(11)

Moreover, we have $v^{\top}\omega_1 > v^{\top}\omega_0$ for any $\omega_1 \in S$ and $\omega_0 \in S^{\complement}$. Finally, combine this result with Equation (10) and Equation (11), we have

$$\begin{split} \mathbb{E}_{\boldsymbol{\omega}_{\mathrm{vul}}'}[\boldsymbol{v}^{\top}\boldsymbol{\omega}_{\mathrm{vul}}'(\tilde{g}^{*}(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}')-\tilde{g}(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\mathrm{vul}}'\in\mathcal{S})]\\ +\mathbb{E}_{\boldsymbol{\omega}_{\mathrm{vul}}'}[\boldsymbol{v}^{\top}\boldsymbol{\omega}_{\mathrm{vul}}'(\tilde{g}^{*}(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}')-\tilde{g}(\boldsymbol{A},\boldsymbol{\omega}_{\mathrm{vul}}'))\mathbb{1}(\boldsymbol{\omega}_{\mathrm{vul}}'\in\mathcal{S}^{\complement})] > 0. \end{split}$$

Consequently, we have $\mathbb{E}_{\omega'_{vul}}[v^{\top}\omega'_{vul}\tilde{g}^{*}(A,\omega'_{vul})] - \mathbb{E}_{\omega'_{vul}}[v^{\top}\omega'_{vul}\tilde{g}(A,\omega'_{vul})] > 0$. Therefore, we have proved that $\tilde{g}^*(\boldsymbol{A}, \boldsymbol{\omega}'_{\text{vul}}) = \mathbb{1}(\boldsymbol{v}^\top \boldsymbol{\omega}'_{\text{vul}} \ge -\varepsilon_{\boldsymbol{v}} \Phi^{-1}(h_c(\boldsymbol{x})))$ is the exact optimal solution of the optimization problem in Equation (8).

A.2 PROOF OF LEMMA 1

Proof. The tractable perturbation budgets ϵ_A and ϵ_X can be obtained according to Theorem 2 and Theorem 4, correspondingly.

A.3 PROOF OF LEMMA 2

Proof. The tractable probability lower bound $P_{\bar{g}_X=1}$ can be obtained according to Theorem 3.

A.4 PROOF OF THEOREM 2

Proof. To certify the fairness level, we assume that $\tilde{g}_{A,X}(A, X) = 1$. Refer to Lemma 2, we have $\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Delta}_{\boldsymbol{A}} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X}) = 1) \geq P_{\tilde{g}_{\boldsymbol{X}}=1}$. For any structure perturbation $\|\boldsymbol{\Delta}_{\boldsymbol{A}}\|_{0} \leq \tilde{\epsilon_{\boldsymbol{A}}}$, we combine this result with Equation (1) and obtain that

$$\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Delta}_{\boldsymbol{A}} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X}) = 1) \ge P_{\tilde{g}_{\boldsymbol{X}}=1} > 0.5,$$
(12)

As a consequence, we have $\tilde{g}_{A,X}(A \oplus \Delta_A, X) = \tilde{g}_{A,X}(\overline{A,X})$ for any structure perturbation $\|\boldsymbol{\Delta}_{\boldsymbol{A}}\|_0 \leq \tilde{\epsilon}_{\boldsymbol{A}}.$

A.5 PROOF OF THEOREM 3

Proof. To prove Theorem 3, we formulate the theoretical prerequisite that Theorem 3 relies on as Lemma A 4. The following Lemma A 4 indicates the relation between $A \oplus \Gamma_A$ and $A \oplus \Delta_A \oplus \Gamma_A$.

Lemma A 4. Let X and Y be two random vectors in the discrete space $\{0,1\}^n$ with prabability distributions $Pr(\mathbf{X})$ and $Pr(\mathbf{Y})$, correspondingly. Let $h : \{0,1\}^n \to \{0,1\}$ be a random or deterministic function. Let $S_1 = \{ \boldsymbol{z} \in \{0,1\}^n : \frac{\Pr(\boldsymbol{X}=\boldsymbol{z})}{\Pr(\boldsymbol{Y}=\boldsymbol{z})} > r \}$ and $S_2 = \{ \boldsymbol{z} \in \{0,1\}^n : \frac{\Pr(\boldsymbol{X}=\boldsymbol{z})}{\Pr(\boldsymbol{Y}=\boldsymbol{z})} > r \}$ $\frac{\Pr(\boldsymbol{X}=\boldsymbol{z})}{\Pr(\boldsymbol{Y}=\boldsymbol{z})} = r \} \text{ for some } r > 0. \text{ Assume } \mathcal{S}_3 \subseteq \mathcal{S}_2 \text{ and } \mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_3. \text{ If } \Pr(h(\boldsymbol{X})=1) \ge \Pr(\boldsymbol{X} \in \mathcal{S}),$ then $\operatorname{Pr}(h(\boldsymbol{Y}) = 1) \geq \operatorname{Pr}(\boldsymbol{Y} \in \mathcal{S}).$

Proof. Note that we have $\frac{\Pr(\mathbf{X}=\mathbf{z})}{\Pr(\mathbf{Y}=\mathbf{z})} \ge r$ for any $\mathbf{z} \in S$ and $\frac{\Pr(\mathbf{X}=\mathbf{z})}{\Pr(\mathbf{Y}=\mathbf{z})} \le r$ for any $\mathbf{z} \in S^{\complement}$. Assuming h is random, we have $\Pr(h(\mathbf{Y}) = 1) - \Pr(\mathbf{Y} \in \mathcal{S})$ $= \sum_{\boldsymbol{z} \in \{0,1\}^n} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{Y} = \boldsymbol{z}) - \sum_{\boldsymbol{z} \in \mathcal{S}} \Pr(\boldsymbol{Y} = \boldsymbol{z})$ $= \sum_{\boldsymbol{z} \in \mathcal{S}} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{Y} = \boldsymbol{z}) + \sum_{\boldsymbol{z} \in \mathcal{S}^\complement} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{Y} = \boldsymbol{z})$ $-\sum_{\boldsymbol{z} \in S} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{Y} = \boldsymbol{z}) - \sum_{\boldsymbol{z} \in S} \Pr(h(\boldsymbol{z}) = 0) \Pr(\boldsymbol{Y} = \boldsymbol{z})$ $=\sum_{\boldsymbol{z}\in S} \Pr(h(\boldsymbol{z})=1)\Pr(\boldsymbol{Y}=\boldsymbol{z}) - \sum_{\boldsymbol{z}\in S} \Pr(h(\boldsymbol{z})=0)\Pr(\boldsymbol{Y}=\boldsymbol{z})$ $\geq \frac{1}{r} (\sum_{\boldsymbol{z} \in S^{\mathsf{f}}} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{X} = \boldsymbol{z}) - \sum_{\boldsymbol{z} \in S} \Pr(h(\boldsymbol{z}) = 0) \Pr(\boldsymbol{X} = \boldsymbol{z}))$ $= \frac{1}{r} (\sum_{\boldsymbol{z} \in \mathcal{S}^\complement} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{X} = \boldsymbol{z}) + \sum_{\boldsymbol{z} \in \mathcal{S}} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{X} = \boldsymbol{z})$ $-\sum_{\boldsymbol{z} \in S} \Pr(h(\boldsymbol{z}) = 0) \Pr(\boldsymbol{X} = \boldsymbol{z}) - \sum_{\boldsymbol{z} \in S} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{X} = \boldsymbol{z}))$ $= \frac{1}{r} \left(\sum_{\boldsymbol{z} \in \{0,1\}^n} \Pr(h(\boldsymbol{z}) = 1) \Pr(\boldsymbol{X} = \boldsymbol{z}) - \sum_{\boldsymbol{z} \in \mathcal{S}} \Pr(\boldsymbol{X} = \boldsymbol{z}) \right)$ $=\frac{1}{r}(\Pr(h(\boldsymbol{X})=1) - \Pr(\boldsymbol{X} \in \mathcal{S}))$ >0.Based on the definition of \mathcal{H} , we have $\Pr(\tilde{g}_{X}(A \oplus \Delta_{A}, X) = 1) \ge \Pr(A \oplus \Delta_{A} \in \mathcal{H})$ distinctly. According to Lemma A 4 (let $Y = A \oplus \Delta_A \oplus \Gamma_A$, $X = A \oplus \Delta_A$, $S = \mathcal{H}$, and $h = \tilde{g}_X$), we have $\Pr(\tilde{g}_X(A \oplus \Delta_A \oplus \Gamma_A, X) = 1) \ge \Pr(A \oplus \Delta_A \oplus \Gamma_A \in \mathcal{H})$ (X can be seen as a constant here), correspondingly. Noting that the exact probability $\Pr(\tilde{q}_{X}(A \oplus \Gamma_{A}, X) = 1)$ is difficult to calculate, we use a practical lower bound $P_{\tilde{g}_{X}=1}^{*} \leq \Pr(\tilde{g}_{X}(A \oplus \Gamma_{A}, X) = 1)$ to replace it in practice. Because we also have $\Pr(\tilde{g}_{X}(A \oplus \Delta_{A}, X) = 1) \ge \Pr(A \oplus \Delta_{A} \in \mathcal{H})$ for \mathcal{H} derived by $P_{\tilde{g}_{X}=1}^{*}$, the proof still holds. A.6 SOLVING THE OPTIMIZATION PROBLEM IN THEOREM 2 We can compute $\Pr(\mathbf{\Lambda} \oplus \mathbf{\Lambda} \oplus \mathbf{\Gamma} \oplus \mathbf{\Gamma})$

We can compute
$$\Pr(\mathbf{A} \oplus \Delta_{\mathbf{A}} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H})$$
 as
1021
1022 $\Pr(\mathbf{A} \oplus \Delta_{\mathbf{A}} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H}) = \sum_{j=\mu+1}^{n|\mathcal{V}_{vul}|} \Pr(\mathbf{A} \oplus \Delta_{\mathbf{A}} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H}_{j})$
1023
1024 $+ (\underline{P_{\tilde{g}_{\mathbf{X}}=1}}^{*} - \sum_{j=\mu+1}^{n|\mathcal{V}_{vul}|} \Pr(\mathbf{A} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H}_{j})) / (\frac{\beta}{1-\beta})^{\mu}.$
(13)

To compute Equation (13), we calculate the probability $\Pr(\mathbf{A} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H}_m)$ and $\Pr(\mathbf{A} \oplus \Delta_{\mathbf{A}} \oplus \Gamma_{\mathbf{A}} \in \mathcal{H}_m)$ as

$$\Pr(\boldsymbol{A} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}} \in \mathcal{H}_{m}) = \sum_{\substack{j=\max\{0,m\}\\m \in [n]\mathcal{V}_{vul}|, n|\mathcal{V}_{vul}|+m\}}}^{\min\{n|\mathcal{V}_{vul}|, n|\mathcal{V}_{vul}|+m\}} \beta_{n|\mathcal{V}_{vul}|-(j-m)} (1-\beta)^{(j-m)} t(m,j), \qquad (14)$$

1029 1030

1033

1036

1037

$$\Pr(\boldsymbol{A} \oplus \boldsymbol{\Delta}_{\boldsymbol{A}} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}} \in \mathcal{H}_m) = \sum_{j=\max\{0,m\}}^{\min\{n \mid \mathcal{V}_{\text{vul}}\mid, n \mid \mathcal{V}_{\text{vul}}\mid +m\}} \beta_{n \mid \mathcal{V}_{\text{vul}}\mid -j} (1-\beta)^j t(m,j),$$
(15)

where $-n|\mathcal{V}_{\text{vul}}| \le m \le n|\mathcal{V}_{\text{vul}}|, \|\Delta_{\boldsymbol{A}}\|_0 \le k$, and t(m, j) is defined as

$$t(m,j) = \begin{cases} 0, & \text{if } (m+k) \mod 2 \neq 0, \\ 0, & \text{if } 2j - m < k, \\ \binom{n|\mathcal{V}_{\text{vul}}|-k}{2} \binom{k}{2}, & \text{otherwise.} \end{cases}$$
(16)

1039 With Equation (13), we can traverse the perturbation budget $\|\Delta_A\|_0$ over 1, 2, ... until $\Pr(A \oplus \Delta_A \oplus \Gamma_A \in \mathcal{H}) < 0.5$.

1041

1048

1055

1057

1042 1043 A.7 Proof of Theorem 4

1044 1045 Proof. Recall that $\tilde{g}_{A,X}(A, X) = \operatorname{argmax}_{c \in \{0,1\}} \operatorname{Pr}(\tilde{g}_X(A \oplus \Gamma_A, X) = c)$ and $\tilde{g}_X(A, X) =$ 1046 $\operatorname{argmax}_{c \in \{0,1\}} \operatorname{Pr}(g(f_{\theta^*}, A, X + \Gamma_X, \eta, \mathcal{V}_{tst}) = c)$. To certify the fairness level, we assume that 1047 $\tilde{g}_{A,X}(A, X) = 1$, which means that

$$\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X}) = 1) > 0.5.$$
(17)

1049 For any perturbation $\|\Delta_X\|_2 \le \epsilon_X \le \tilde{\epsilon}_X$, we have $\tilde{g}_X(A \oplus \Gamma_A, X + \Delta_X) = \tilde{g}_X(A \oplus \Gamma_A, X)$ 1050 for any $\Gamma_A \in \bar{A}$ where $\tilde{\epsilon}_X$ is derived with classifier $\tilde{g}_X(A \oplus \Gamma_A, X)$. Regarding the randomness of 1051 Γ_A , we have

1052 $\Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X} + \boldsymbol{\Delta}_{\boldsymbol{X}}) = 1) = \Pr(\tilde{g}_{\boldsymbol{X}}(\boldsymbol{A} \oplus \boldsymbol{\Gamma}_{\boldsymbol{A}}, \boldsymbol{X}) = 1) > 0.5.$ (18) 1053 Hence we obtain that $\tilde{g}_{\boldsymbol{A},\boldsymbol{X}}(\boldsymbol{A}, \boldsymbol{X} + \boldsymbol{\Delta}_{\boldsymbol{X}}) = \tilde{g}_{\boldsymbol{A},\boldsymbol{X}}(\boldsymbol{A}, \boldsymbol{X})$ for any perturbation $\|\boldsymbol{\Delta}_{\boldsymbol{X}}\|_{2} \le \epsilon_{\boldsymbol{X}}$. \Box

1056 A.8 PROOF OF PROPOSITION 1

Proof. For the practical certification, we add perturbations within certified budgets and derive independent identically distributed output samples $\hat{\mathcal{Y}}'$ by Monte Carlo. For each sample $\hat{\mathbf{Y}}' \in \hat{\mathcal{Y}}'$, we have $\Pr(\pi(\hat{\mathbf{Y}}', \mathcal{V}_{tst}) > \eta) < 0.5$ according to Theorem 4 and Theorem 2. $\pi(\hat{\mathbf{Y}}, \mathcal{V}_{tst}) > \eta$ indicates that $\pi(\hat{\mathbf{Y}}', \mathcal{V}_{tst}) > \eta$ for any $\hat{\mathbf{Y}}' \in \hat{\mathcal{Y}}'$. Consequently, we have

$$\Pr(\pi(\hat{\boldsymbol{Y}}, \mathcal{V}_{tst}) > \eta) = \prod_{\hat{\boldsymbol{Y}}' \in \hat{\mathcal{Y}}'} \Pr(\pi(\hat{\boldsymbol{Y}}', \mathcal{V}_{tst}) > \eta) < 0.5^{|\mathcal{Y}'|}.$$
(19)

1064

1066

1062

A.9 RATIONALE OF EACH THEORETICAL RESULT

1067 1068 Here we provide an explanation below about the rationale of each theoretical result in this paper.

Theorem 1. (Certified Fairness Defense for Node Attributes): This theorem gives a way to compute the perturbation-invariant budget (i.e., the budget within which the fairness level will not reduce under a certain threshold) of node attributes. However, since we consider both input data modalities could be attacked, we still need to extend the analysis over the span of node attributes and graph topology (see Theorem 4).

Theorem 2. (Certified Defense Budget for Structure Perturbations): This theorem formulates an optimization problem, whose solution is the perturbation-invariant budget (i.e., the budget within which the fairness level will not reduce under a certain threshold) of graph topology under the smoothed node attributes. However, to solve this optimization problem, we need to explicitly compute $P_{\bar{g}_{X}=1}$ (see Theorem 3).

Theorem 3. (Positive Probability Lower Bound): This theorem provides a way to explicitly compute $\underline{P}_{\tilde{g}_{X}=1}$, which directly enables us to solve the optimization problem in Theorem 2.

Theorem 4. (Certified Defense Budget for Attribute Perturbations): This theorem is built upon
 Theorem 1 and provides a way to explicitly compute the perturbation-invariant budget of node
 attributes over the span of node attributes and graph topology.

1087 Lemma 1. (Perturbation-Invariant Budgets Existence): This lemma claims the existence and tractability of the perturbation-invariant budgets on both data modalities, which is further detailed by Theorem 2 and Theorem 4.

1090 1091 **Lemma 2.** (Positive Probability Bound Under Noises): This lemma claims the existence and 1092 tractability of $P_{\tilde{g}_{x}=1}$, which is further detailed by Theorem 3.

Proposition 1. (Probabilistic Guarantee for the Fairness Level of Node Classification): This
 proposition provides a neat probabilistic theoretical guarantee — we have a probability that is large
 enough to successfully achieve certified defense on fairness.

- 097
- 1099 B REPRODUCIBILITY AND SUPPLEMENTARY ANALYSIS
- 1100 1101

In this section, our primary emphasis is on ensuring the replicability of our experiments, which serves as an extension to Section 4. To begin with, we offer a comprehensive introduction of the three real-world datasets adopted in our experiments. Subsequently, we introduce the detailed experimental settings, as well as the implementation details of our proposed framework, ELEGANT, alongside GNNs and baseline models. Moreover, we outline those essential packages, including their versions, that were utilized in our experiments. Lastly, we elaborate on the supplementary analysis on the time complexity of ELEGANT.

- 1109
- 1110 1111

1112

B.1 DATASETS

In our experiments, we adopt three real-world network datasets that are widely used to perform
studies on the fairness of GNNs, namely German Credit (Asuncion & Newman, 2007; Agarwal et al.,
2021), Recidivism (Jordan & Freiburger, 2015; Agarwal et al., 2021), and Credit Defaulter (Yeh &
Lien, 2009; Agarwal et al., 2021). We introduce their basic information below.

(1) German Credit. Each node is a client in a German bank (Asuncion & Newman, 2007), while
each edge between any two clients represents that they bear similar credit accounts. Here the gender
of bank clients is considered as the sensitive attribute, and the task is to classify the credit risk of the
clients as high or low.

(2) Recidivism. Each node denotes a defendant released on bail at the U.S state courts during 1990-2009 (Jordan & Freiburger, 2015), and defendants are connected based on the similarity of their past criminal records and demographics. Here the race of defendants is considered as the sensitive attribute, and the task is to classify defendants into more likely vs. less likely to commit a violent crime after being released.

(3) Credit Defaulter. This dataset contains credit card users collected from financial agencies (Yeh & Lien, 2009). Specifically, each node in this network denotes a credit card user, and users are connected based on their spending and payment patterns. The sensitive attribute is the age period of users, and the task is to predict the future default of credit card for these users. We present the statistics pf the three datasets above in Table 2.

Dataset	German Credit	Recidivism	Credit Defaulter
# Nodes	1,000	18,876	30,000
# Edges	22,242	321,308	1,436,858
# Attributes	27	18	13
Avg. degree	44.5	34.0	95.8
Sens.	Gender	Race	Age
Label	Credit status	Bail decision	Future default

Table 2: The statistics and basic information about the six real-world datasets adopted for experimentalevaluation. Sens. represents the semantic meaning of sensitive attribute.

1148 1149

1134

For the three real-world datasets used in this paper, we adopt the split rate for the training set and validation set as 0.4 and 0.55, respectively. The input node features are normalized before they are fed into the GNNs and the corresponding explanation models. For the downstream task *node classification*, only the labels of the nodes in the training set is available for all models during the training process. The trained GNN models with the best performance on the validation set are preserved for test and explanation.

1158 B.2 DETAILED EXPERIMENTAL SETTINGS

Implementation of GNN Models. In our experiments, all GNN models are implemented in Py Torch (Paszke et al., 2017) with PyG (PyTorch Geometric) (Fey & Lenssen, 2019). For the corresponding hyper-parameters, we set the value of weight decay as 5e-4, with the hidden dimension number and dropout rate being 64 and 0.6, respectively. In addition, we set the learning rate and epoch number as 5e-2 and 200 for training.

1166 Implementation of ELEGANT. ELEGANT is implemented in PyTorch (Paszke et al., 2017) with 1167 MIT license and all GNNs under ELEGANT are optimized through Adam optimizer (Kingma & Ba, 1168 2015) on Nvidia A6000. In our experiments, the sampling sizes of Gaussian noise and Bernoulli 1169 noise are 150 and 200, respectively. All hyper-parameters for GNNs under ELEGANT are set as the 1170 same values as the hyper-parameters adopted for vanilla GNNs. We propose to add Gaussian and 1171 Bernoulli noise (to node attributes and graph topology) during training, which empirically leads to 1172 better certification performance, i.e., larger certification budgets over both node attributes and graph 1173 1174 topology. Specifically, we set the entry-wise probability of flipping the existence of an edge and the 1175 standard deviation of the added Gaussian noise as 2e-4 and 2e-5, respectively. In addition, we set the 1176 confidence level as 0.7 for estimation, since a lower confidence level helps exhibit a clearer tendency 1177 of the change of certified budgets w.r.t. other parameters under a limited number of sampling size, 1178 considering the computational costs. In the test phase, we set the sampled ratio for certification (from 1179 the nodes out of training and validation set) to be 0.9 to make the sampled size relatively large, in 1180 which way we include more nodes in the set of nodes to be certified. In each run, we sample 100 1181 times, and the value of FCR is averaged across three runs with different seeds. Finally, considering 1182 the sizes of the three datasets, we set the nodes that are vulnerable to be 5% for German Credit and 1183 1% for others. 1184

Selection of ϵ and β . There are two critical parameters, ϵ and β , that could affect the effectiveness of ELEGANT. These two parameters control the level of randomness for the added Gaussian and Bernoulli noise, respectively. Intuitively, larger ϵ and β will induce more randomness in the node

¹¹⁵⁷

1188 attributes and graph structure, which could make ELEGANT more robust to perturbations with larger 1189 sizes and thus achieve larger ϵ_X and ϵ_A . However, if ϵ_X and ϵ_A are too large, the randomness could 1190 go beyond what the GNN classifier can manage and could finally cause failure in certification. Hence 1191 it is necessary to first determine appropriate values of $\epsilon_{\mathbf{X}}$ and $\epsilon_{\mathbf{A}}$ for ELEGANT. Here we propose a 1192 strategy for parameter selection to realize as large certified defense budgets as possible. Specifically, 1193 we first set an empirical η to be 25% higher than the fairness level of the corresponding vanilla GNN 1194 model. Such a threshold calibrates across different GNNs and can be considered as a reasonable 1195 1196 threshold for the exhibited bias. Then we determine two wide search spaces for σ and β , respectively, 1197 and compute the averaged $\epsilon_{\mathbf{X}}$ and $\epsilon_{\mathbf{A}}$ from multiple runs over each pair of σ and β values. We now 1198 rank (σ, β) pairs based on the averaged $\epsilon_{\mathbf{X}}$ and $\epsilon_{\mathbf{A}}$ in a descending order, respectively. Finally, we 1199 truncate the obtained two rankings from their most top-ranked (σ, β) pair to the tail, until the two truncated rankings have the first overlapped (σ, β) pair. Such an identified (σ, β) pair can achieve 1201 large and balanced certification budgets over both A and X, and hence they are recommended. 1202

1203 Implementation of Baseline Models. In this paper, we include two fairness-aware GNNs as the 1204 baselines for comparison, namely FairGNN and NIFTY. We introduce the details below. (1) FairGNN. 1205 For FairGNN, we adopt the official implementations from (Dai & Wang, 2021). Hyper-parameters 1206 corresponding to the GNN model structure (such as the number of hidden dimensions) are ensured 1207 to be the same as the vanilla GNNs for a fair comparison. Other parameters are carefully tuned 1208 under the guidance of the recommended training settings. (2) NIFTY. For NIFTY, we use the official 1209 implementations provided from (Agarwal et al., 2021). We ensured that the parameters related to the 1210 GNN model structure stay the same as the original GNNs for a fair comparison. We also adjust other 1211 parameters based on the suggested training settings for better performance. 1212

Packages Required for Implementations. We list those key packages and their corresponding versions adopted in our experiments below.

- Python == 3.8.8
- torch == 1.10.1
- torch-geometric == 2.1.0
 - torch-scatter == 2.0.9
- torch-sparse == 0.6.13
 - cuda == 11.1
 - numpy == 1.20.1
 - tensorboard == 2.10.0
 - networkx == 2.5
 - scikit-learn == 0.24.2
 - pandas==1.2.4
 - scipy==1.6.2
- 1233 1234

1216

1217

1218 1219

1220

1222

1224 1225

1226

1227 1228

1229

1230 1231

1232

B.3 ALGORITHMIC ROUTINE

Now we introduce the pipeline of the proposed framework ELEGANT to obtain the node classification results in facing of the graph data that could have been perturbed by malicious attackers. We present the algorithmic routine in Algorithm 1. Note that ABSTAIN refers to the case where certification fails. Correspondingly, FCR measures the ratio of not returning ABSTAIN for the proposed framework ELEGANT, which generally reflects the usability of the certification defense.

1242 B.4 EVALUATION OF MODEL UTILITY

1244

1260

1266 1267



Figure 4: The utility of GCN, E-GCN, FairGNN, and NIFTY under fairness attacks on German Credit. The shaded bar indicates that certified budget $\epsilon_A \leq \|\Delta_A\|_0$ or $\epsilon_X \leq \|\Delta_X\|_2$.

In Section 4.3, we present the comparison between ELEGANT and baseline models over the fairness level under attacks. We now present the comparison over the utility under attacks. Specifically, we utilize node classification accuracy as the indicator of model utility, and we present the results in Figure 4. The fairnessaware GNNs are found to exhibit better utility compared with the vanilla GNNs, which is a common observation consistent with a series of existing works (Agarwal et al., 2021; Dong et al., 2022a). More importantly, we observe that the ELEGANT does not jeopardize the performance

of GNN compared with the utility of the vanilla GNN. This demonstrate a high level of usability for ELEGANT in real-world applications.

B.5 Certification under Different Fairness Metrics

1268 In Section 4.2, we present the experimental results based on the fairness metric of Δ_{SP} , which 1269 measures the exhibited bias under the fairness notion of *Statistical Parity*. We also perform the 1270 experiments based on $\Delta_{\rm EO}$, which measures the exhibited bias under the fairness notion of Equal 1271 *Opportunity.* We present the experimental results in Table 3. We summarize the observations below. 1272 (1) Fairness Certification Rate (FCR). We observe that ELEGANT realizes large values of FCR 1273 (larger than 80%) for all three GNN backbones and three attributed network datasets. Similar to 1274 our discussion in Section 4.2, this demonstrate that for nodes in any randomly sampled test set, 1275 1276 we have a probability around or larger than 80% to successfully certify the fairness level of the 1277 predictions yielded by the GNN model with our proposed framework ELEGANT. As a consequence, 1278 we argue that ELEGANT also achieves a satisfying fairness certification rate across all adopted GNN 1279 backbones and datasets on the basis of $\Delta_{\rm EO}$. In addition, we also observe that the German Credit 1280 dataset bears relatively larger values of FCR, while the values of FCR are relatively smaller with 1281 relatively larger standard deviation values on Recidivism and Credit Defaulter datasets. A possible 1282 reason is that we set the threshold (i.e., η) as a value 25% higher than the bias exhibited by the vanilla 1283 GNNs. Consequently, if the vanilla GNNs already exhibit a low level of bias, the threshold determined 1284 with such a strategy could be hard to satisfy under the added noise. This evidence indicates that the 1285 proposed framework ELEGANT tends to deliver better performance under scenarios where vanilla 1286 GNNs exhibit a high level of bias with the proposed strategy. (2) Utility. Compared with vanilla 1287 1288 GNNs, certified GNNs with ELEGANT exhibit comparable and even higher node classification accuracy values in all cases. Therefore, we argue that the proposed framework ELEGANT does 1290 not significantly jeopardize the utility of the vanilla GNN models in certifying the fairness level 1291 of node classification. (3) Fairness. We observe that certified GNNs with ELEGANT are able to achieve better performances in terms of algorithmic fairness compared with those vanilla GNNs. This 1293 evidence indicates that the proposed framework ELEGANT also helps to mitigate the exhibited bias 1294 (by the backbone GNN models). We conjecture that such bias mitigation should be attributed to the 1295 same reason discussed in Section 4.2.

Table 3: Comparison between vanilla GNNs and certified GNNs under ELEGANT over three popular GNNs across three real-world datasets. Here ACC is node classification accuracy, and E- prefix marks out the GNNs under ELEGANT with certification. \uparrow denotes the larger, the better; \downarrow denotes the opposite. Different from the table in Section 4.2 (where the bias is measured with Δ_{SP}), the bias is measured with Δ_{EO} here. Numerical values are in percentage, and the best ones are in bold.

	Ge	erman Cre	dit		Recidivism			Credit Defaulter		
	ACC (†)	Bias (\downarrow)	FCR (\uparrow)	ACC (†)	Bias (\downarrow)	FCR (†)	ACC (↑)	Bias (\downarrow)	FCR (†)	
SAGE	67.3 _{±2.14}	41.8 ±11.0	N/A	89.8 ±0.66	6.09 ±3.10	N/A	75.9 ±2.18	$10.4_{\pm 1.59}$	N/A	
E-SAG	E 72.2 ±1.26	$8.63{\scriptstyle~\pm6.15}$	$100{\scriptstyle~\pm 0.00}$	90.8 ±0.97	3.12 ±3.64	$81.0_{\ \pm 13.0}$	$73.4{\scriptstyle~\pm 0.61}$	7.18 ±1.06	88.7 ±6.02	
GCN	59.6 ±3.64	35.0 ±4.77	N/A	90.5 ±0.73	$6.35_{\pm 1.65}$	N/A	65.8 ±0.29	13.5 ±4.23	N/A	
E-GCN	58.8 +3 74	29.8 +6 82	93.3 +8 73	89.3 ± 0.92	3.93 +3 12	96.0 +4 97	63.5 ± 0.37	9.12 ±0.95	80.5 +14 5	
IK	63.3 +4.11	37.7 +15.9	N/A	91.9 ±0.54	5.26 +3.25	N/A	76.6 ±0.69	8.04 ±0.57	N/A	
E-JK	63.4 ±3.68	31.2 ±15.5	93.7 ±8.96	90.1 ±0.55	2.54 ±1.62	83.7 ±8.96	76.9 ±0.86	2.90 ±2.04	95.7 ±4.80	
laorit	hm 1 Carti	fied Defer	use on the	Fairmass	f CNNa					
Algorit	nin i Ceru	neu Deren	ise on the	raimess o	I GININS					
Input:					e					
G: g	graph data w	th potentia	al maliciou	is attacks; j	f_{θ^*} : an opt	imized GN	N node cla	ssifier; V_{train}	in, \mathcal{V} validation	
Vtest	$\in V$: the n	ode set for	training, v	alidation, a	ind test, res	pectively;	$V_{\text{vul}} \in V_{\text{test}}$	the set of	vulnerable	
node	es that may t	bear attacks	(on node a	$\frac{1}{1}$	id/or graph	topology);	N_1, N_2 : sat	mple size fo	or the set of	
Ber		aussian noi	ise, respect	ively; η : a	given thres	noid for the		bias; α : the	e parameter	
10 11			ever $(1 - c$	x) of the est	σ_{1}	: the std of	the added G	Jaussian no	bise; p: the	
prot	bability of re	turning zero	o of the add	ied Bernoul	lli noise;					
Output	•					A		c 1 1		
ϵ_{A} :	the certified	defense bi	idget over	the adjacen	icy matrix	$\mathbf{A}; \boldsymbol{\epsilon}_{\mathbf{X}}$: the	certified de	efense budg	get over the	
node	e attribute m	atrix X; Y	the outp	ut node clas	sification r	esults from	the certified	d classifier;		
1: Sam	ple a set of I	Bernoulli no	onse $Q_{\rm B}$ con	ntaining N_1	samples;					
2: Sam	ple a set of (Gaussian no	Dise $Q_{\rm G}$ col	ntaining N_2	samples;					
3: for ($\omega_A \in \mathcal{Q}_{\mathrm{B}}$ do)								
4: 10	or $\omega_{\mathbf{X}} \in \mathcal{Q}_{\mathbf{G}}$	do		C 6 1		c 1				
5:	Calculate a	ind collect t	the output of	of f_{θ^*} unde	r the noise	of ω_A and ω_A	$\omega_{\boldsymbol{X}};$			
6:	Calculate a	ind collect t	the output of	of g based of	on the outpu	it of <i>f</i> _{0*} ;				
/: e	nd for	11 (1	1 C		1.0		<i>c</i> 1			
8: U	Inder $Q_{\rm G}$, co	flect the nu	imber of g	returning I	and 0 as n_1	and n_0 , re	spectively;			
9: E	stimate the I	ower bound	1 of returni	ng c as $P_{g=}$	$\frac{1}{c}$ determin	ed by the la	rger one be	tween n_1 a	nd n_0 ;	
10: 1	$n_1 > n_0$ ar	$\frac{P_{g=1}}{5}$ is 1	larger than	0.5 with a c	confidence	level larger	than $1 - \alpha$	or $n_1 < n$	$_0$ and $\underline{P_{g=0}}$	
15	arger than	0.5 with a c	confidence	level larger	than $1 - \alpha$	then				
11:	Calculate a	and collect	the value of	$\epsilon_{\boldsymbol{X}};$						
12: e	ise									
13:	return AB	STAIN								
14: e	nd if									

- 1341
 14:
 end i

 1342
 15:
 end for
- 1343 16: Collect the number of cases where $n_1 > n_0$ and estimate the lower bound of returning 1 as $P_{\tilde{g}_{X}=1}$;
- 1344 17: if $P_{\tilde{g}_{X}=1}$ is larger than 0.5 with a confidence level larger than 1α then
- 1345 18: Calculate $\epsilon_{\mathbf{X}}$ (out of the collected $\tilde{\epsilon}_{\mathbf{X}}$) and $\epsilon_{\mathbf{A}}$ (based on the estimated $P_{\tilde{g}_{\mathbf{X}}=1}$);
- 1346 19: Find \mathbf{Y}' out of the collected output of f_{θ^*} ;
- **1347** 20: return Y', ϵ_X , and ϵ_A ;
- 1348 21: else
- 1349 22: return ABSTAIN
 - 23: end if

1350 B.6 ORDERING THE INNER AND OUTER DEFENSE 1351

1352 We first review the general pipeline to achieve certified fairness defense. Specifically, we first model 1353 the fairness attack and defense by formulating the bias indicator function q. Then, we achieve certified 1354 defense over the node attributes for g, which leads to classifier $\tilde{g}_{\mathbf{X}}$. Finally, we realize certified de-1355 fense for $\tilde{g}_{\mathbf{X}}$ over the graph topology, which leads to classifier $\tilde{g}_{\mathbf{A},\mathbf{X}}$. In general, we may consider the 1356 certified defense over node attributes and graph topology as the inner certified classifier and outer cer-1357 tified classifier, respectively. Now, a natural question is: is it possible to achieve certified defense in a 1358 different order, i.e., first achieve certified defense over the graph topology (as the inner classifier), and 1359 1360 then realize certified defense over the node attributes (as the outer classifier)? Note that this is not the 1361 research focus of this paper, but we will provide insights about this question. In fact, it is also feasible 1362 to achieve certified defense in the reversed order compared with the approach presented in our paper. 1363 We provide an illustration in Figure 5. We follow a similar setting to plot this figure as in Section 3.3. 1364 Specifically, in case (1), both $A_{i,j} \oplus 0$ and $A_{i,j} \oplus 1$ lead to a positive outcome for g; in case (2), both 1365 $A_{i,j} \oplus 0$ and $A_{i,j} \oplus 1$ lead to a negative outcome. However, considering the Gaussian distribution 1366 around $X_{i,j}$, samples will fall around case (1) with a much higher number compared with case (2). 1367

Hence, in this example, it would be reasonable to
assume that the classifier with Bernoulli noise over
graph topology (the inner certified classifier) will
return 1 with a higher probability. This example
thus illustrates how certification following a different
order returns 1.

However, such a formulation bears higher computa-1375 tional costs in calculating the certified budgets. The 1376 1377 reason is that we are able to utilize a closed-form solu-1378 tion to calculate $\epsilon_{\mathbf{X}}$ based on a set of Gaussian noise 1379 and the corresponding output from the bias indicator 1380 function. However, based on a set of Bernoulli noise 1381 and the corresponding output from the bias indica-1382 tor function, we will need to solve the optimization problem given in Theorem 2 to calculate ϵ_A , which 1384 bears a higher time complexity than calculating $\epsilon_{\mathbf{X}}$. 1385 If we follow the strategy provided in Section 3.4 to 1386 calculate the inner and outer certification budgets, the 1387 1388 certified budget of the inner certification will always



Figure 5: An example illustrating how ELE-GANT works with a different order to achieve certified defense.

be calculated multiple times, while the certified budget of the outer certification will only be calculated once. Considering the high computational cost of calculating ϵ_A , we thus argue that it is more efficient to realize the certification over graph topology as the outer certified classifier.

1393 1394

1395

B.7 CERTIFICATION WITH ESTIMATED PROBABILITIES

In Section 3.4, we proposed to utilize estimated lower bounds of the probabilities (including P(c) in Theorem 1 and $\Pr(\tilde{g}_{X}(A \oplus \Gamma_{A}, X) = 1)$ in Theorem 3) to perform certification in practice, considering the exact probability values are difficult to compute. In Appendix A.1 and Appendix A.5, we have discussed that both theorems hold no matter exact probability values or estimated lower bounds (of the probabilities above) are used. Now we present a brief review of other theoretical analysis to show that they also hold. (1) for Lemma 2, note that taking a lower bound estimation to replace the exact $\Pr(\tilde{g}_{X}(A \oplus \Gamma_{A}, X) = 1)$ reduces the total size of \mathcal{H} in Theorem 3. Correspondingly, the formulated $P_{\tilde{g}_{X}=1}$ based on the estimated $\Pr(\tilde{g}_{X}(A \oplus \Gamma_{A}, X) = 1)$ is smaller than that based on 1404 the exact $\Pr(\tilde{g}_X(A \oplus \Gamma_A, X) = 1)$. Hence Lemma 2 still holds when $P_{\tilde{g}_X=1}$ is replaced with one 1405 calculated based on the estimated $\Pr(\tilde{g}_{\mathbf{X}}(\mathbf{A} \oplus \boldsymbol{\Gamma}_{\mathbf{A}}, \mathbf{X}) = 1)$. (2) For Theorem 2, it holds no matter 1406 how P(c) in Theorem 1 and $\Pr(\tilde{g}_{X}(A \oplus \Gamma_{A}, X) = 1)$ in Theorem 3 are obtained. (3) For Theorem 1407 4, according to Appendix A.7, it still holds as long as Theorem 1 holds. (4) For Proposition 1, in 1408 all cases where $\Pr(g(A \oplus \Gamma'_A, X + \Gamma_X) = 1)$ is identified to be larger than 0.5 with an estimated 1409 lower bound, the (underlying) exact $\Pr(g(A \oplus \Gamma'_A, X + \Gamma_X) = 1)$ will be larger than the estimated 1410 probability value under the given confidence level, and thus will also be larger than 0.5. Here Γ'_A is a 1411 1412 sampled Bernoulli noise, i.e., $\Gamma'_A \in \overline{\mathcal{A}}'$. According to Appendix A.8, Proposition 1 still holds in this 1413 case. (5) Finally, we conclude that Lemma 1 still holds since Theorem 1, Lemma 2, Theorem 3, and 1414 Theorem 4 hold.

1415 1416

1417 B.8 TIME COMPLEXITY ANALYSIS

We now present a comprehensive analysis on the time complexity of ELEGANT. We present the analysis from both theoretical and experimental perspectives.

1421
1422**Theoretical.** The time complexity is linear w.r.t. the total number of the random perturbations N, i.e.,
 $\mathcal{O}(N)$. We perform 30,000 random perturbations over the span of node attributes and graph structure.
We note that the actual running time is acceptable since the certification does not require re-training
(which is the most costly process). In addition, all runnings do not rely on the prediction results from
each other. Hence they can be paralleled altogether theoretically to further reduce the running time.

1427 **Experimental.** We perform a study of running time, and we present the results in Table 4. Specifically, 1428 we compare the running time of a successful certification under 30,000 random noise samples and 1429 a regular training-inference cycle with vanilla GCN. We observe that (1) although ELEGANT 1430 improves the computational cost compared with the vanilla GNN backbones, the running time 1431 remains acceptable; and (2) ELEGANT has less running time growth rate on larger datasets. For 1432 example, E-SAGE has around 10x running time on German Credit (a smaller dataset) while only 1433 around 4x on Credit Default (a larger dataset) compared to vanilla SAGE. Hence we argue that 1434 ELEGANT bears a high level of usability in terms of complexity and running time. 1435

1436

1437 B.9 ADDITIONAL RESULTS ON DIFFERENT GNN BACKBONES & BASELINES 1438

We perform additional experiments over two popular GNNs, including APPNP (Klicpera et al., 2019)
and GCNII (Chen et al., 2020), to evaluate the generalization ability of ELEGANT onto different
backbones. We present all numerical results in Table 5 (in terms of accuracy), Table 6 (in terms of
fairness), and Table 7 (in terms of FCR). We observe that ELEGANT achieves comparable utility,
a superior level of fairness, and a large percentage of FCR. This verifies the satisfying usability of
ELEGANT, which remains consistent with the paper.

In addition, we provide a detailed fairness comparison between ELEGANT and robust GNNs from
(Jin et al., 2021) and (Wu et al., 2019b) in Table 8. We observe that the best performances still come
from the GNNs equipped with ELEGANT on all datasets. Hence we argue that ELEGANT exhibits
satisfying performance in usability, which remains consistent with the discussion in the paper.

Why ELEGANT Improves Fairness? We note that improving fairness is a byproduct of ELEGANT,
and our focus is to achieve certification over the fairness level of the prediction results. We now
provide a detailed discussion about why fairness is improved here. First, existing works found that
the distribution difference in the node attribute values and edge existence across different subgroups
is a significant source of bias (Dong et al., 2022a; Dai & Wang, 2021; Fan et al., 2021). However,
adding noise on both node attributes and graph topology may reduce such distributional divergence
and mitigate bias. Second, As mentioned in Section 3.4, the proposed strategy to obtain the output

predictions in ELEGANT is to select the fairest result among the output set $\hat{\mathcal{Y}}'$, where each output is derived based on a sample $\Gamma'_{A} \in \overline{\mathcal{A}}'$ (i.e., $\operatorname{argmin}_{\hat{\mathbf{Y}}'} \pi(\hat{\mathbf{Y}}', \mathcal{V}_{tst})$ s.t. $\hat{\mathbf{Y}}' \in \hat{\mathcal{Y}}'$). Such a strategy provides a large enough probability to achieve certification in light of Proposition 1. Meanwhile, we point out that such a strategy also helps to significantly improve fairness since highly biased outputs are excluded.

B.10 COMPLEMENTARY RESULTS

We provide the results in terms of Δ_{EO} for Table 1 in Table 9, and we present the results of the baselines for Figure 2 in Table 10, Table 11, Table 12, and Table 13. For Table 9, we observe that ELEGANT does not constantly show a lower value of Δ_{EO} . This is because the certification goal in Table 1 is Δ_{SP} instead of Δ_{EO} . In addition, we note that debiasing existing GNN models is not the goal of this paper. In addition, we provide the corresponding results in terms of accuracy for Figure 3 in Table 14 and Table 15. We observe that although most performance remains stable, a stronger noise (i.e., larger σ and smaller β) generally leads to worse but still comparable performance. This is consistent with the discussion in Section 4.4, and this has been taken into consideration in the discussion of the parameter selection strategy in Section 4.4.

Table 4: Comparison of running time (in seconds) on different datasets using different methods.

German	Recidivism	Credit
5.27 ± 0.38	34.14 ± 1.08	40.11 ± 0.36
$\textbf{53.23} \pm \textbf{1.31}$	$\textbf{137.12} \pm \textbf{58.66}$	$\textbf{157.51} \pm \textbf{37.21}$
5.59 ± 0.37	34.94 ± 1.16	40.59 ± 0.32
$\textbf{53.79} \pm \textbf{30.19}$	$\textbf{212.94} \pm \textbf{10.38}$	$\textbf{214.11} \pm \textbf{10.31}$
5.78 ± 0.43	34.68 ± 0.88	39.44 ± 1.56
$\textbf{59.99} \pm \textbf{25.01}$	$\textbf{238.37} \pm \textbf{1.81}$	$\textbf{252.99} \pm \textbf{17.03}$
	German 5.27 ± 0.38 53.23 ± 1.31 5.59 ± 0.37 53.79 ± 30.19 5.78 ± 0.43 59.99 ± 25.01	GermanRecidivism 5.27 ± 0.38 34.14 ± 1.08 53.23 ± 1.31 137.12 ± 58.66 5.59 ± 0.37 34.94 ± 1.16 53.79 ± 30.19 212.94 ± 10.38 5.78 ± 0.43 34.68 ± 0.88 59.99 ± 25.01 238.37 ± 1.81

Table 5: Performance comparison of classification accuracy. Numbers are in percentage.

0.1				
91		German	Recidivism	Credit
93	SAGE	67.3 ± 2.14	89.8 ± 0.66	75.9 ± 2.18
14	E-SAGE	$\textbf{71.0} \pm \textbf{1.27}$	89.9 ± 0.90	$\textbf{73.4} \pm \textbf{0.50}$
5	GCN	59.6 ± 3.64	90.5 ± 0.73	65.8 ± 0.29
6	D CON	59.0 ± 5.04	90.5 ± 0.75	
7	E-GCN	$\textbf{58.2} \pm \textbf{1.82}$	89.6 ± 0.74	65.2 ± 0.99
8	JK	63.3 ± 4.11	91.9 ± 0.54	76.6 ± 0.69
9	E-JK	$\textbf{62.3} \pm \textbf{4.07}$	$\textbf{89.3} \pm \textbf{0.33}$	$\textbf{77.7} \pm \textbf{0.27}$
00	APPNP	69.9 ± 2.17	95.3 ± 0.78	74.4 ± 3.05
)1	E-APPNP	$\textbf{69.4} \pm \textbf{0.83}$	$\textbf{95.9} \pm \textbf{0.02}$	$\textbf{74.6} \pm \textbf{0.32}$
)2	GCNII	60.9 ± 1.00	90.4 ± 0.95	77.7 ± 0.22
13	E-GCNII	$\textbf{60.4} \pm \textbf{4.45}$	$\textbf{88.8} \pm \textbf{0.24}$	$\textbf{77.6} \pm \textbf{0.02}$

С ADDITIONAL DISCUSSION

C.1 WHY CERTIFY A CLASSIFIER ON TOP OF AN OPTIMIZED GNN?

We note that the rationale of certified defense is to provably maintain the classification results against attacks. Under this context, most existing works on certifying an existing deep learning model focus

	German	Recidivism	Credit
SAGE	50.6 ± 15.9	9.36 ± 3.15	13.0 ± 4.01
E-SAGE	$\textbf{16.3} \pm \textbf{10.9}$	$\textbf{6.39} \pm \textbf{2.85}$	8.94 ± 0.99
GCN	37.4 ± 3.24	10.1 ± 3.01	11.1 ± 3.22
E-GCN	$\textbf{3.52} \pm \textbf{3.77}$	$\textbf{9.56} \pm \textbf{3.22}$	7.28 ± 1.46
JK	41.2 ± 18.1	10.1 ± 3.15	9.24 ± 0.60
E-JK	$\textbf{22.4} \pm \textbf{1.95}$	$\textbf{6.26} \pm \textbf{2.78}$	3.37 ± 2.64
APPNP	27.4 ± 4.81	9.71 ± 3.57	12.3 ± 3.14
E-APPNP	$\textbf{13.1} \pm \textbf{5.97}$	$\textbf{2.23} \pm \textbf{0.04}$	10.8 ± 0.07
GCNII	51.4 ± 0.36	9.70 ± 3.37	7.62 ± 0.29
E-GCNII	$\textbf{24.9} \pm \textbf{0.47}$	$\textbf{3.78} \pm \textbf{0.93}$	1.72 ± 0.81

Table 6: Comparison of fairness (measured with Δ_{SP}). Numbers are in percentage.

Table 7: Performance in FCR on different datasets and backbone GNNs. Numbers are in percentage.

	German	Recidivism	Credit	
E-SAGE	98.7 ± 1.89	94.3 ± 6.65	94.3 ± 3.3	
E-GCN	96.3 ± 1.89	96.0 ± 3.56	92.7 ± 5.19	
E-JK	97.0 ± 3.00	89.5 ± 10.5	99.3 ± 0.47	
E-APPNP	97.8 ± 3.14	87.1 ± 3.79	95.5 ± 6.43	
E-GCNII	94.7 ± 5.27	92.9 ± 9.93	99.0 ± 1.41	

on certifying a specific predicted label over a given data point. Here, the prediction results to be certified are classification results. Correspondingly, these works are able to certify the model itself.

However, the strategy above is not feasible in our studied problem. This is because we seek to
certify the level of fairness of a group of nodes. The value of such a group-level property cannot be
directly considered as a classification result, and thus they are not feasible to be directly certified.
Therefore, we proposed to first formulate a classifier on top of an optimized GNN. As such, achieving
certification becomes feasible. In fact, this also serves as one of the contributions of our work.

 C.2 WHAT IS THE DIFFERENCE BETWEEN THE ATTACKING PERFORMANCE OF GNNS AND THE FAIRNESS OF GNNS?

In traditional attacks over the performance of GNNs, the objective of the attacker is simply formulated as having false predictions on as many nodes as possible, such that the overall performance is jeopardized. However, in attacks over the fairness of GNNs, whether the goal of the attacker can be achieved is jointly determined by the GNN predictions over all nodes. Such node-level dependency in achieving the attacking goal makes the defense over fairness attacks more difficult, since the defense cannot be directly performed at the node level but at the model level instead. Correspondingly, this necessitates (1) constructing an additional classifier as discussed in the previous reply, and (2) additional theoretical analysis over the constructed classifier as in Theorem 1, 2, and 3 to achieve certification.

1567				
1568 1569		German	Recidivism	Credit
1570	SAGE	50.6 ± 15.9	9.36 ± 3.15	13.0 ± 4.01
1571	E-SAGE	$\textbf{16.3} \pm \textbf{10.9}$	$\textbf{6.39} \pm \textbf{2.85}$	$\textbf{8.94} \pm \textbf{0.99}$
1572	GCN	37.4 ± 3.24	10.1 ± 3.01	11.1 ± 3.22
1573	E-GCN	$\textbf{3.52} \pm \textbf{3.77}$	$\textbf{9.56} \pm \textbf{3.22}$	$\textbf{7.28} \pm \textbf{1.46}$
1574	ЈК	41.2 ± 18.1	10.1 ± 3.15	9.24 ± 0.60
1575	E-JK	$\textbf{22.4} \pm \textbf{1.95}$	$\textbf{6.26} \pm \textbf{2.78}$	$\textbf{3.37} \pm \textbf{2.64}$
1576	(Jin et al., 2021)	14.8 ± 18.3	9.59 ± 0.65	3.84 ± 0.17
1577	(Wu et al., 2019b)	3.66 ± 0.52	8.04 ± 2.97	7.10 ± 5.10

Table 8: Comparison of fairness (measured with Δ_{SP}). Numbers are in percentage.

Table 9: The Δ_{EO} of Table 1 in the paper. All numerical numbers are in percentage.

	German	Recidivism	Credit
SAGE	30.43 ± 0.07	3.71 ± 0.01	5.56 ± 0.03
E-SAGE	$\textbf{12.21} \pm \textbf{0.04}$	$\textbf{6.95} \pm \textbf{0.02}$	$\textbf{7.18} \pm \textbf{0.01}$
GCN	35.19 ± 0.07	5.06 ± 0.01	11.9 ± 0.02
E-GCN	$\textbf{8.32} \pm \textbf{0.03}$	$\textbf{1.39} \pm \textbf{0.01}$	$\textbf{6.24} \pm \textbf{0.02}$
JK	18.10 ± 0.13	3.02 ± 0.01	9.47 ± 0.02
E-JK	$\textbf{23.68} \pm \textbf{0.02}$	$\textbf{2.74} \pm \textbf{0.01}$	$\textbf{2.55} \pm \textbf{0.01}$

C.3 CERTIFICATION WITHOUT CONSIDERING THE BINARY SENSITIVE ATTRIBUTE

We utilize the most widely studied setting to assume the sensitive attributes are binary. However, our certification approach is not designed to be tailored to the sensitive attributes. Therefore, our approach can be easily extended to scenarios where the sensitive attributes are multi-class and continuous by adopting the corresponding fairness metric as the function $\pi(\cdot)$ in Definition 1.

C.4 How Do the Main Theoretical Findings Differ From Existing Works on Robustness Certification of GNNs on Regular Attacks?

Most existing works for robustness certification can only defend against attacks on either node attributes or graph structure. Due to the multi-modal input data of GNNs, existing works usually fail to handle the attacks over node attributes and graph structure at the same time. However, ELEGANT is able to defend against attacks over both data modalities. This necessitates using both continuous and discrete noises for smoothing and the analysis for joint certification in the span of the two input data modalities (as shown in Figure 1).

Table 10: The results under $(2^0, 10^{-1})$ in terms of node classification accuracy, AUC score, F1 score, Δ_{SP} , and Δ_{EO} Figure 2. All numerical numbers are in percentage.

1615 1616	$(2^0, 10^{-1})$	Accuracy	AUC	F1 Score	Δ_{SP}	Δ_{EO}
1617	GCN	58.4%	66.4%	63.9%	41.4%	33.4%
1618	NIFTY	61.2%	68.1%	66.2%	33.9%	13.3%
1619	FairGNN	55.2%	62.2%	61.4%	16.4%	5.99%

Table 11: The results under $(2^1, 10^0)$ in terms of node classification accuracy, AUC score, F1 score, Δ_{SP} , and Δ_{EO} Figure 2. All numerical numbers are in percentage.

$(2^1, 10^0)$	Accuracy	AUC	F1 Score	Δ_{SP}	Δ_{EO}
GCN	58.4%	66.4%	63.9%	41.4%	33.4%
NIFTY	61.2%	68.2%	66.2%	36.1%	13.3%
FairGNN	55.2%	62.2%	61.4%	16.8%	7.77%

Table 12: The results under $(2^2, 10^1)$ in terms of node classification accuracy, AUC score, F1 score, Δ_{SP} , and Δ_{EO} for Figure 2. All numerical numbers are in percentage.

$(2^2, 10^1)$	Accuracy	AUC	F1 Score	Δ_{SP}	Δ_{EO}
GCN	58.0%	66.6%	63.7%	41.4%	37.8%
NIFTY	61.2%	68.1%	66.0%	42.1%	13.3%
FairGNN	55.6%	62.1%	61.9%	16.0%	9.56%

C.5 DISCUSSION: DIFFERENCE WITH EXISTING SIMILAR WORKS

Here we mainly focus on discussing the difference between this work and (Bojchevski et al., 2020).

We note that (1) the randomized smoothing technique adopted in (Bojchevski et al., 2020) is different from the proposed randomized smoothing approach on the graph topology in this paper and (2) the techniques in (Bojchevski et al., 2020) tackle a different problem from this paper. We elaborate on more details below.

The techniques in (Bojchevski et al., 2020) are different from this paper. Although both randomized smoothing approaches are able to handle binary data, we note that the randomized smoothing approach proposed in (Bojchevski et al., 2020) is data-dependent. However, the proposed randomized smoothing approach in this paper is data-independent. We note that in practice, a data-independent approach enables practitioners to pre-generate noises, which significantly improves usability.

The studied problem in (Bojchevski et al., 2020) is different from this paper. Although the authors claimed to achieve a joint certificate for graph topology and node attributes in (Bojchevski et al., 2020), all node attributes are assumed to be binary, which can only be applied to cases where these attributes are constructed as bag-of-words representations (as mentioned in the second last paragraph in the Introduction of (Bojchevski et al., 2020)). However, in this work, we follow a more realistic setting where only graph topology is assumed to be binary while node attributes are considered as continuous. This makes the problem more difficult to handle, since different strategies should be adopted for different data modalities. In summary, compared with (Bojchevski et al., 2020), the problem studied in this paper is more realistic and more suitable for GNNs.

Table 13: The results under $(2^3, 10^2)$ in terms of node classification accuracy, AUC score, F1 score, Δ_{SP} , and Δ_{EO} for Figure 2. All numerical numbers are in percentage.

1669 1670	$(2^3, 10^2)$	Accuracy	AUC	F1 Score	Δ_{SP}	Δ_{EO}
1671	GCN	58.0%	67.7%	63.7%	42.6%	45.7%
1672	NIFTY	58.8%	67.3%	63.1%	44.5%	19.4%
1673	FairGNN	54.4%	61.4%	61.0%	16.9%	23.8%

	5e-3	5e-2	5e-1	5e0
0	57.50 ± 1.50	57.51 ± 1.63	57.50 ± 1.58	55.67 ± 2.00
1e-3	57.50 ± 1.51	57.51 ± 1.63	57.50 ± 1.58	55.67 ± 2.0
5e-3	57.49 ± 1.52	57.50 ± 1.64	57.50 ± 1.58	55.67 ± 2.00
1e-2	57.55 ± 1.50	57.51 ± 1.65	57.50 ± 1.58	55.67 ± 2.0
5e-2	N/A	57.57 ± 1.59	57.50 ± 1.58	55.67 ± 2.0
1e-1	N/A	57.53 ± 1.57	57.50 ± 1.59	55.67 ± 2.0
5e-1	N/A	N/A	57.49 ± 1.60	55.67 ± 2.00
1e0	N/A	N/A	57.40 ± 1.58	55.67 ± 2.00
5e0	N/A	N/A	N/A	55.76 ± 1.86

Table 14: Classification accuracy in Figure 3(a) with different settings. Numbers are in percentage.

Table 15: Classification accuracy in Figure 3(b) with different settings. Numbers are in percentage.

	0.6	0.7	0.8	0.9
0	63.71 ± 0.64	64.03 ± 0.66	65.87 ± 0.49	64.88 ± 0.46
2^{0}	63.71 ± 0.64	64.03 ± 0.66	65.87 ± 0.49	64.88 ± 0.46
2^1	63.67 ± 0.64	64.04 ± 0.67	N/A	N/A
2^2	63.69 ± 0.67	N/A	N/A	N/A
2^3	N/A	N/A	N/A	N/A
2^4	N/A	N/A	N/A	N/A

1700

1687

1688

1674

Based on the discussion above, we would like to note that no existing work can be directly adopted to tackle the studied problem in this paper.

1701 1702

1704

1703 C.6 ADDITIONAL EXPERIMENTS ON DIFFERENT DATASETS

To further validate the performance of the proposed method, we also perform experiments with the 1705 same commonly used popular GNN backbone models (as Section 4.2) on two Pokec datasets, namely 1706 1707 Pokec-z and Pokec-n. We present the experimental results in Table 16, where all numerical numbers 1708 are in percentage. We observe that (1) the GNNs equipped with ELEGANT achieve comparable 1709 node classification accuracy; (2) the GNNs equipped with ELEGANT achieve consistently lower 1710 levels of bias; and (3) the values of the Fairness Certification Rate (FCR) for all GNNs equipped with 1711 ELEGANT exceed 90%, exhibiting satisfying usability. All three observations are consistent with 1712 the experimental results and the corresponding discussion presented in Section 4.2. Therefore, we 1713 argue that the effectiveness of the proposed approach is not determined by the dataset and is well 1714 generalizable over different graph datasets. 1715

1716

1718

1717 C.7 SCALABILITY OF ELEGANT

In this subsection, we discuss the scalability of ELEGANT. Specifically, we note that if the Gaussian 1719 1720 and Bernoulli noise is directly added over the whole graph, scaling to larger graphs would be difficult. 1721 However, the proposed approach can be easily extended to the batch training case, which has been 1722 widely adopted by existing scalable GNNs. Specifically, a commonly adopted batch training strategy 1723 of scalable GNNs is to only input a node and its surrounding subgraph into the GNN, since the 1724 prediction of GNNs only depends on the information of the node itself and its multi-hop neighbors, 1725 and the number of hops is determined by the layer number of GNNs. Since the approach proposed in 1726 our paper aligns with the basic pipeline of GNNs, the perturbation can also be performed for each 1727 specific batch of nodes. In this case, all theoretical analyses in this paper still hold, since they also do

		Pokec-z			Pokec-n	
	ACC (†)	Bias (↓)	FCR (†)	ACC (†)	Bias (↓)	FCR (†)
SAGE	63.13 ± 0.37	6.29 ± 0.20	-	57.60 ± 2.74	6.43 ± 1.08	-
E-SAGE	62.09 ± 2.22	4.18 ± 1.87	94.00 ± 5.66	60.74 ± 1.87	5.23 ± 0.13	91.50 ± 7.78
GCN	64.89 ± 0.93	3.44 ± 0.16	-	59.86 ± 0.09	4.26 ± 0.40	-
E-GCN	62.38 ± 0.26	1.52 ± 0.49	90.50 ± 0.71	59.83 ± 4.16	3.23 ± 1.20	94.00 ± 8.49
JK	63.06 ± 1.00	7.89 ± 3.05	-	57.70 ± 1.05	8.81 ± 2.46	-
E-JK	61.49 ± 2.55	3.63 ± 2.18	87.50 ± 2.12	61.19 ± 0.50	5.60 ± 0.01	93.00 ± 9.90

Table 16: Experimental results on Pokec-z and Pokec-n datasets.

not rely on the assumption of non-batch training. Therefore, we would like to argue that the proposed
approach can be easily scaled to large graphs.