

---

# FATE: Fairness Attacks on Graph Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study fairness attacks on graph learning to answer the following question: *How*  
2 *can we achieve poisoning attacks on a graph learning model to exacerbate the*  
3 *bias?* We answer this question via a bi-level optimization problem and propose a  
4 meta learning-based attacking framework named FATE. The proposed framework  
5 is broadly applicable with respect to various fairness definitions and graph learning  
6 models, as well as arbitrary choices of manipulation operations. We further instanti-  
7 ate FATE to attack statistical parity and individual fairness on graph neural networks.  
8 We conduct extensive experimental evaluations on real-world datasets in the task  
9 of semi-supervised node classification. The experimental results demonstrate that  
10 FATE could amplify the bias of graph neural networks with or without fairness  
11 consideration while maintaining the utility on the downstream task. We hope this  
12 paper provides insights into the adversarial robustness of fair graph learning and  
13 can shed light on designing robust and fair graph learning in future studies.

## 14 1 Introduction

15 Algorithmic fairness in graph learning has received much research attention [5, 20, 24]. Despite  
16 its substantial progress, existing studies mostly assume the benevolence of input graphs and aim  
17 to ensure that the bias would not be perpetuated or amplified in the learning process. However,  
18 malicious activities in the real world are commonplace. For example, consider a financial fraud  
19 detection system which utilizes a transaction network to classify whether a bank account is fraudulent  
20 or not [49, 45]. An adversary may manipulate the transaction network (e.g., malicious banker with  
21 access to the transaction data, theft of bank accounts to make malicious transactions), so that the  
22 graph-based fraud detection model would exhibit unfair classification results with respect to people  
23 of different demographic groups. Consequently, a biased fraud detection model may infringe civil  
24 liberty to certain financial activities and impact the well-being of an individual negatively [6]. It  
25 would also make the graph learning model fail to provide the same quality of service to people of  
26 certain demographic groups, causing the financial institutions to lose business in the communities  
27 of the corresponding demographic groups. Thus, it is critical to understand how resilient a graph  
28 learning model is with respect to adversarial attacks on fairness, which we term as *fairness attacks*.

29 To date, fairness attack has not been well studied. Sporadic literature often follows two strategies:  
30 (1) adversarial data point injection, which is often designed for tabular data rather than graphs [38,  
31 33, 8, 44] or (2) adversarial edge injection, which only attacks the group fairness of a graph neural  
32 network [19]. It is thus crucial to study how to attack different fairness definitions for a variety of  
33 graph learning models.

34 To achieve this goal, we study the Fairness attacks on graph learning (FATE) problem. We formulate  
35 it as a bi-level optimization, where the lower-level problem optimizes a task-specific loss function  
36 to make the fairness attacks deceptive and the upper-level problem leverages the supervision signal  
37 to modify the input graph and maximize the bias function corresponding to a user-defined fairness  
38 definition. To solve the bi-level optimization problem, we propose a meta learning-based solver

39 (FATE), whose key idea is to compute the meta-gradient of the upper-level bias function with respect  
40 to the input graph to guide the fairness attacks. Compared with existing works, our proposed  
41 FATE framework has two major advantages. First, it is capable of attacking *any* fairness definition  
42 on *any* graph learning model, as long as the corresponding bias function and the task-specific loss  
43 function are differentiable. Second, it is equipped with the ability for either continuous or discretized  
44 poisoning attacks on the graph topology. We also briefly discuss its ability for poisoning attacks on  
45 node features in a later section.

46 The major contributions of this paper are summarized as follows.

- 47 • **Problem definition.** We formally define the problem of fairness attacks on graph learning (the  
48 FATE problem). Based on the definition, we formulate it as a bi-level optimization problem, whose  
49 key idea is to maximize a bias function in the upper level while minimizing a task-specific loss  
50 function for a graph learning task.
- 51 • **Attacking framework.** We propose an end-to-end attacking framework named FATE. It learns a  
52 perturbed graph topology via meta learning, such that the bias with respect to the learning results  
53 trained with the perturbed graph will be amplified.
- 54 • **Empirical evaluation.** We conduct experiments on three benchmark datasets to demonstrate the  
55 efficacy of our proposed FATE framework in amplifying the bias while being the most deceptive  
56 method (i.e., achieving the highest micro F1 score) on semi-supervised node classification.

## 57 2 Preliminaries and Problem Definition

58 **A – Notations.** Throughout the paper, we use bold upper-case letter for matrix (e.g.,  $\mathbf{A}$ ), bold  
59 lower-case letter for vector (e.g.,  $\mathbf{x}$ ) and calligraphic letter for set (e.g.,  $\mathcal{G}$ ). We use superscript  $T$   
60 to denote the transpose of a matrix/vector (e.g.,  $\mathbf{x}^T$  is the transpose of  $\mathbf{x}$ ). Regarding matrix/vector  
61 indexing, we use conventions similar to NumPy in Python. For example,  $\mathbf{A}[i, j]$  is the entry of  $\mathbf{A}$  at  
62 the  $i$ -th row and  $j$ -th column;  $\mathbf{x}[i]$  is the  $i$ -th entry of  $\mathbf{x}$ ;  $\mathbf{A}[i, :]$  and  $\mathbf{A}[j, :]$  are the  $i$ -th row and  $j$ -th  
63 column of  $\mathbf{A}$ , respectively.

64 **B – Algorithmic fairness.** The general principle of algorithmic fairness is to ensure the learning  
65 results would not favor one side or another.<sup>1</sup> Among several fairness definitions that follow this  
66 principle, group fairness [16, 18] and individual fairness [15] are the most widely studied ones. Group  
67 fairness splits the entire population into multiple demographic groups by a sensitive attribute (e.g.,  
68 gender) and ensure the parity of a statistical property among learning results of those groups. For  
69 example, statistical parity, a classic group fairness definition, guarantees the statistical independence  
70 between the learning results (e.g., predicted labels of a classification algorithm) and the sensitive  
71 attribute [16]. Individual fairness suggests that similar individuals should be treated similarly. It is  
72 often formulated as a Lipschitz inequality such that distance between the learning results of two data  
73 points should be no larger than the difference between these two data points [15].

74 **C – Problem definition.** Existing work [19] for fairness attacks on graphs randomly injects adversar-  
75 ial edges so that the disparity between the learning results of two different demographic groups would  
76 be amplified. However, it suffers from three major limitations. (1) First, it only attacks statistical  
77 parity while overlooking other fairness definitions (e.g., individual fairness [15]). (2) Second, it only  
78 considers adversarial edge injection, excluding other manipulations like edge deletion or reweighting.  
79 Hence, it is essential to investigate the possibility to attack other fairness definitions on real-world  
80 graphs with an arbitrary choice of manipulation operations. (3) Third, it does not consider the utility  
81 of graph learning models while achieving the fairness attacks, resulting in performance degradation  
82 in the downstream tasks. However, an institution that applies the graph learning models are often  
83 utility-maximizing [28, 2]. Thus, a performance degradation in the utility would make the fairness  
84 attacks not deceptive from the perspective of a utility-maximizing institution.

85 In this paper, we seek to overcome the aforementioned limitations. To be specific, given an input  
86 graph, an optimization-based graph learning model, and a user-defined fairness definition, we aim  
87 to learn a modified graph such that a bias function of the corresponding fairness definition would  
88 be maximized for *effective* fairness attacks, while minimizing the task-specific loss function with  
89 respect to the graph learning model for *deceptive* fairness attacks. Formally, we define the problem of  
90 fairness attacks on graph learning, which is referred to as the FATE problem.

---

<sup>1</sup><https://www.merriam-webster.com/dictionary/fairness>

91 **Problem 1** FATE: *Fairness Attacks on Graph Learning*

92 **Given:** (1) An undirected graph  $\mathcal{G} = \{\mathbf{A}, \mathbf{X}\}$ ; (2) a task-specific loss function  $l(\mathcal{G}, \mathcal{Y}, \Theta, \theta)$  where  $\mathcal{Y}$   
 93 is the graph learning results,  $\Theta$  is the set of learnable variables and  $\theta$  is the set of hyperparameters; (3)  
 94 a bias function  $b(\mathbf{Y}, \Theta^*, \mathbf{F}, \theta)$  where  $\Theta^* = \arg \min_{\Theta} l(\mathcal{G}, \mathcal{Y}, \Theta, \theta)$  and  $\mathbf{F}$  is the matrix that contains  
 95 auxiliary fairness-related information (e.g., sensitive attribute values of all nodes in  $\mathcal{G}$  for group  
 96 fairness, pairwise node similarity matrix for individual fairness); (4) an integer budget  $B$ .

97 **Find:** a poisoned graph  $\tilde{\mathcal{G}} = \{\tilde{\mathbf{A}}, \tilde{\mathbf{X}}\}$  which satisfies the following properties: (1)  $d(\mathcal{G}, \tilde{\mathcal{G}}) \leq B$   
 98 where  $d(\mathcal{G}, \tilde{\mathcal{G}})$  is the distance between the input graph  $\mathcal{G}$  and the poisoned graph  $\tilde{\mathcal{G}}$  (e.g.,  $\|\mathbf{A}, \tilde{\mathbf{A}}\|_{1,1}$ );  
 99 (2) the bias function  $b(\mathbf{Y}, \Theta^*, \mathbf{F})$  is maximized for effectiveness; (3) the task-specific loss function  
 100  $l(\tilde{\mathcal{G}}, \mathcal{Y}, \Theta, \theta)$  is minimized for deceptiveness.

### 101 3 Methodology

102 In this section, we first formulate Problem 1 as a bi-level optimization problem, followed by a generic  
 103 meta learning-based solver named FATE.

#### 104 3.1 Problem Formulation

105 Given an input graph  $\mathcal{G} = \{\mathbf{A}, \mathbf{X}\}$  with adjacency matrix  $\mathbf{A}$  and node feature matrix  $\mathbf{X}$ , an attacker  
 106 aims to learn a poisoned graph  $\tilde{\mathcal{G}} = \{\tilde{\mathbf{A}}, \tilde{\mathbf{X}}\}$  such that the graph learning model will be maximally  
 107 biased when trained on  $\tilde{\mathcal{G}}$ . In this work, we consider the following settings for the attacker.

108 **The goal of the attacker.** The attacker aims to amplify the bias of the graph learning results output  
 109 by a victim graph learning model. And the bias to be amplified is a choice made by the attacker based  
 110 on which fairness definition the attacker aims to attack.

111 **The knowledge of the attacker.** Following similar settings in [19], we assume the attacker has  
 112 access to the adjacency matrix, the feature matrix of the input graph, and the sensitive attribute of  
 113 all nodes in the graph. For a (semi-)supervised learning problem, we assume that the ground-truth  
 114 labels of the training nodes are also available to the attacker. For example, for a graph-based financial  
 115 fraud detection problem, the malicious banker may have access to the demographic information (i.e.,  
 116 sensitive attribute) of the account holders and also know whether some bank accounts are fraudulent  
 117 or not, which are the ground-truth labels for training nodes. Similar to [51, 52, 19], the attacker has  
 118 no knowledge about the parameters of the victim model. Instead, the attacker will perform a gray-box  
 119 attack by attacking a surrogate graph learning model.

120 **The capability of the attacker.** The attacker is able to perturb up to  $B$  edges/features in the graph  
 121 (i.e.,  $\|\mathbf{A} - \tilde{\mathbf{A}}\|_{1,1} \leq B$  or  $\|\mathbf{X} - \tilde{\mathbf{X}}\|_{1,1} \leq B$ ).

122 Based on that, we formulate Problem 1 as a bi-level optimization problem as follows.

$$\begin{aligned} \tilde{\mathcal{G}} &= \arg \max_{\tilde{\mathcal{G}}} b(\mathbf{Y}, \Theta^*, \mathbf{F}) \\ \text{s.t. } \Theta^* &= \arg \min_{\Theta} l(\mathcal{G}, \mathbf{Y}, \Theta, \theta), d(\mathcal{G}, \tilde{\mathcal{G}}) \leq B \end{aligned} \quad (1)$$

123 where the lower-level problem learns an optimal surrogate graph learning model  $\Theta^*$  by minimizing  
 124  $l(\mathcal{G}, \mathbf{Y}, \Theta, \theta)$ , the upper-level problem finds a poisoned graph  $\tilde{\mathcal{G}}$  that could maximize a bias function  
 125  $b(\mathbf{Y}, \Theta^*, \mathbf{F})$  for the victim graph learning model and the distance between the input graph and the  
 126 poisoned graph  $d(\mathcal{G}, \tilde{\mathcal{G}})$  is constrained to satisfy the setting about the budgeted attack. Note that  
 127 Eq. (1) is applicable to attack *any* fairness definition on *any* graph learning model, as long as the bias  
 128 function  $b(\mathbf{Y}, \Theta^*, \mathbf{F})$  and the loss function  $l(\mathcal{G}, \mathbf{Y}, \Theta, \theta)$  are differentiable.

129 **A – Lower-level optimization problem.** A wide spectrum of graph learning models are essentially  
 130 solving an optimization problem. Take the graph convolutional network (GCN) [26] as an example.  
 131 It learns the node representation by aggregating information from its neighborhood, i.e., message  
 132 passing. Mathematically, for an  $L$ -layer GCN, the hidden representation at  $k$ -th layer can be  
 133 represented as  $\mathbf{E}^{(k)} = \sigma(\hat{\mathbf{A}}\mathbf{E}^{(k-1)}\mathbf{W}^{(k)})$  where  $\sigma$  is a nonlinear activation function (e.g., ReLU),  
 134  $\hat{\mathbf{A}} = \mathbf{D}^{-1/2}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-1/2}$  with  $\mathbf{D}$  being the degree matrix of  $(\mathbf{A} + \mathbf{I})$  and  $\mathbf{W}^{(k)}$  is the learnable  
 135 weight matrix of the  $k$ -th layer. Then the lower-level optimization problem aims to learn the set

136 of parameters  $\Theta^* = \{\mathbf{W}^{(k)} | k = 1, \dots, L\}$  that could minimize a task-specific loss function (e.g.,  
 137 cross-entropy loss for semi-supervised node classification). For more examples of graph learning  
 138 models from the optimization perspective, please refers to Appendix A.

139 **B – Upper-level optimization problem.** To attack the fairness aspect of a graph learning model, we  
 140 aim to maximize a differentiable bias function  $b(\mathbf{Y}, \Theta^*, \mathbf{F})$  with respect to a user-defined fairness  
 141 definition in the upper-level optimization problem. For example, for statistical parity [16], the fairness-  
 142 related auxiliary information matrix  $\mathbf{F}$  can be defined as the one-hot demographic membership matrix,  
 143 where  $\mathbf{F}[i, j] = 1$  if and only if node  $i$  belongs to  $j$ -th demographic group. Then the statistical parity  
 144 is equivalent to the statistical independence between the learning results  $\mathbf{Y}$  and  $\mathbf{F}$ . Based on that,  
 145 existing studies propose several differentiable measurements of the statistical dependence between  $\mathbf{Y}$   
 146 and  $\mathbf{F}$  as the bias function. For example, Bose et al. [5] use mutual information  $I(\mathbf{Y}; \mathbf{F})$  as the bias  
 147 function; Prost et al. [35] define the bias function as the Maximum Mean Discrepancy  $MMD(\mathcal{Y}_0, \mathcal{Y}_1)$   
 148 between the learning results of two different demographic groups  $\mathcal{Y}_0$  and  $\mathcal{Y}_1$ .

### 149 3.2 The FATE Framework

150 To solve Eq. (1), we propose a generic attacking framework named FATE to learn the poisoned graph.  
 151 The key idea is to view Eq. (1) as a meta learning problem, which aims to find suitable hyperparameter  
 152 settings for a learning task [3], and treat the graph  $\mathcal{G}$  as a hyperparameter. With that, we learn the  
 153 poisoned graph  $\tilde{\mathcal{G}}$  using the meta-gradient of the bias function  $b(\mathbf{Y}, \Theta^*, \mathbf{F})$  with respect to  $\mathcal{G}$ . In the  
 154 following, we introduce two key parts of FATE in details, including meta-gradient computation and  
 155 graph poisoning with meta-gradient.

156 **A – Meta-gradient computation.** The key term to learn the poisoned graph is the meta-gradient of  
 157 the bias function with respect to the graph  $\mathcal{G}$ . Before computing the meta-gradient, we assume that  
 158 the lower-level optimization problem converges in  $T$  epochs. Thus, we first pre-train the lower-level  
 159 optimization problem by  $T$  epochs to obtain the optimal model  $\Theta^* = \Theta^{(T)}$  before computing the  
 160 meta-gradient. The training of the lower-level optimization problem can also be viewed as a dynamic  
 161 system with the following updating rule

$$\Theta^{(t+1)} = \text{opt}^{(t+1)}(\mathcal{G}, \Theta^{(t)}, \theta, \mathbf{Y}), \forall t \in \{1, \dots, T\} \quad (2)$$

162 where  $\Theta^{(1)}$  refers to  $\Theta$  at initialization,  $\text{opt}^{(t+1)}(\cdot)$  is an optimizer that minimizes the lower-level  
 163 loss function  $l(\mathcal{G}, \mathbf{Y}, \Theta^{(t)}, \theta)$  at  $(t + 1)$ -th epoch. From the perspective of the dynamic system,  
 164 by applying the chain rule and unrolling the training of lower-level problem with Eq. (2), the  
 165 meta-gradient  $\nabla_{\mathcal{G}} b$  can be written as

$$\nabla_{\mathcal{G}} b = \nabla_{\mathcal{G}} b(\mathbf{Y}, \Theta^{(T)}, \mathbf{F}) + \sum_{t=0}^{T-2} A_t B_{t+1} \dots B_{T-1} \nabla_{\Theta^{(t)}} b(\mathbf{Y}, \Theta^{(T)}, \mathbf{F}) \quad (3)$$

166 where  $A_t = \nabla_{\mathcal{G}} \Theta^{(t+1)}$  and  $B_t = \nabla_{\Theta^{(t)}} \Theta^{(t+1)}$ . However, Eq. (3) is computationally expensive in  
 167 both time and space. To further speed up the computation, we adopt a first-order approximation of  
 168 the meta-gradient [17] and simplify the meta-gradient as

$$\nabla_{\mathcal{G}} b \approx \nabla_{\Theta^{(T)}} b(\mathbf{Y}, \Theta^{(T)}, \mathbf{F}) \cdot \nabla_{\mathcal{G}} \Theta^{(T)} \quad (4)$$

169 Since the input graph is undirected, the derivative of the symmetric adjacency matrix  $\mathbf{A}$  can be  
 170 computed as follows by applying the chain rule of a symmetric matrix [21].

$$\nabla_{\mathbf{A}} b \leftarrow \nabla_{\mathbf{A}} b + (\nabla_{\mathbf{A}} b)^T - \text{diag}(\nabla_{\mathbf{A}} b) \quad (5)$$

171 For the node feature matrix  $\mathbf{X}$ , its derivative is equal to the partial derivative  $\nabla_{\mathbf{X}} b$  since it is often an  
 172 asymmetric matrix.

173 **B – Graph poisoning with meta-gradient.** After computing the meta-gradient of the bias function  
 174  $\nabla_{\mathcal{G}} b$ , we aim to poison the input graph guided by  $\nabla_{\mathcal{G}} b$ . We introduce two poisoning strategies: (1)  
 175 continuous poisoning and (2) discretized poisoning.

176 *Continuous poisoning attack.* The continuous poisoning attack is straightforward by reweighting  
 177 edges in the graph. We first compute the meta-gradient of the bias function  $\nabla_{\mathbf{A}} b$ , then use it to poison  
 178 the input graph in a gradient descent-based updating rule as follows.

$$\mathbf{A} \leftarrow \mathbf{A} - \eta \nabla_{\mathbf{A}} b \quad (6)$$

179 where  $\eta$  is a learning rate to control the magnitude of the poisoning attack. The learning rate should  
 180 satisfy  $\eta \leq \frac{B}{\|\nabla_{\mathbf{A}}\|_{1,1}}$  to ensure that constraint on the budgeted attack.

181 *Discretized poisoning attack.* The discretized poisoning attack aims to select a set of edges to be  
 182 added/deleted. It is guided by a poisoning preference matrix defined as follows.

$$\nabla_{\mathbf{A}} = (\mathbf{1} - 2\mathbf{A}) \circ \nabla_{\mathbf{A}} b \quad (7)$$

183 where  $\mathbf{1}$  is an all-one matrix with the same dimension as  $\mathbf{A}$  and  $\circ$  denotes the Hadamard product.  
 184 A large positive  $\nabla_{\mathbf{A}}[i, j]$  indicates strong preference in adding an edge if nodes  $i$  and  $j$  are not  
 185 connected (i.e., positive  $\nabla_{\mathbf{A}} b[i, j]$ , positive  $(\mathbf{1} - 2\mathbf{A})[i, j]$ ) or deleting an edge if nodes  $i$  and  $j$  are  
 186 connected (i.e., negative  $\nabla_{\mathbf{A}} b[i, j]$ , negative  $(\mathbf{1} - 2\mathbf{A})[i, j]$ ). Then, a greedy selection strategy is  
 187 applied to find the set of edges  $\mathcal{E}_{\text{attack}}$  to be added/deleted.

$$\mathcal{E}_{\text{attack}} = \text{topk}(\nabla_{\mathbf{A}}, \delta) \quad (8)$$

188 where  $\text{topk}(\nabla_{\mathbf{A}}, \delta)$  selects  $\delta$  entries with highest preference score in  $\nabla_{\mathbf{A}}$ . Note that, if we only want  
 189 to add edges without any deletion, all negative entries in  $\nabla_{\mathbf{A}} b$  should be zeroed out before computing  
 190 Eq. (7). Likewise, if edges are only expected to be deleted, all positive entries should be zeroed out.

191 *Remarks.* Poisoning node feature matrix  $\mathbf{X}$  follows the same steps as poisoning adjacency matrix  $\mathbf{A}$   
 192 without applying Eq. (5).

193 **C – Overall framework.** FATE generally works as follows. (1) We first pre-train the surrogate graph  
 194 learning model and get the corresponding learning model  $\Theta^{(T)}$  as well as the learning results  $\mathbf{Y}^{(T)}$ .  
 195 (2) Then we compute the meta gradient of the bias function using Eqs. (4) and (5). (3) Finally, we  
 196 perform the discretized poisoning attack (Eqs. (7) and (8)) or continuous poisoning attack (Eq. (6)).  
 197 A detailed pseudo-code of FATE is provided in Appendix B.

198 **D – Limitations.** Since FATE leverages the meta-gradient to poison the input graph, it requires the  
 199 bias function  $b(\mathbf{Y}, \Theta^{(T)}, \mathbf{F})$  to be differentiable in order to calculate the meta-gradient  $\nabla_{\mathcal{G}} b$ . In  
 200 Sections 4 and 5, we present a carefully chosen bias function for FATE. And we leave it for future  
 201 work on exploring the ability of FATE in attacking other fairness definitions. Moreover, though the  
 202 meta-gradient can be efficiently computed via auto-differentiation in many deep learning packages  
 203 (e.g., PyTorch<sup>2</sup>, TensorFlow<sup>3</sup>), it requires  $O(n^2)$  space complexity to store the meta-gradient when  
 204 attacking fairness via edge flipping. It is still a challenging open problem on how to efficiently  
 205 compute the meta-gradient in terms of space. One possible remedy for discretized attack might be a  
 206 low-rank approximation on the perturbation matrix formed by  $\mathcal{E}_{\text{attack}}$ . Since the difference between  
 207 the benign graph and poisoned graph are often small and budgeted ( $d(\mathcal{G}, \tilde{\mathcal{G}}) \leq B$ ), it is likely that  
 208 the edge manipulations may be around a few set of nodes, which makes the perturbation matrix to be  
 209 an (approximately) low-rank matrix.

## 210 4 Instantiation #1: Statistical Parity on Graph Neural Networks

211 Here, we instantiate FATE framework by attacking statistical parity on graph neural networks in a  
 212 binary node classification problem with a binary sensitive attribute. We briefly discuss how to choose  
 213 (1) the surrogate graph learning model used by the attacker, (2) the task-specific loss function in the  
 214 lower-level optimization problem and (3) the bias function in the upper-level optimization problem.

215 **A – Surrogate graph learning model.** We assume that the surrogate model to be used by the attacker  
 216 is a 2-layer linear GCN [47] with different hidden dimensions and model parameters at initialization.

217 **B – Lower-level loss function.** We consider a semi-supervised node classification task for  
 218 the graph neural network to be attacked. Thus, the lower-level loss function is chosen as  
 219 the cross entropy between the ground-truth label and the predicted label:  $l(\mathcal{G}, \mathbf{Y}, \Theta, \theta) =$   
 220  $\frac{1}{|\mathcal{V}_{\text{train}}|} \sum_{i \in \mathcal{V}_{\text{train}}} \sum_{j=1}^c y_{i,j} \ln \hat{y}_{i,j}$ , where  $\mathcal{V}_{\text{train}}$  is the set of training nodes with ground-truth labels  
 221 with  $|\mathcal{V}_{\text{train}}|$  being its cardinality,  $c$  is the number of classes,  $y_{i,j}$  is a binary indicator of whether node  
 222  $i$  belongs to class  $j$  and  $\hat{y}_{i,j}$  is the prediction probability of node  $i$  belonging to class  $j$ .

223 **C – Upper-level bias function.** We aim to attack statistical parity in the upper-level problem, which  
 224 asks for  $\text{P}[\hat{y} = 1] = \text{P}[\hat{y} = 1 | s = 1]$ . Suppose  $p(\hat{y})$  is the probability density function (PDF) of  $\hat{y}_{i,1}$

<sup>2</sup><https://pytorch.org/>

<sup>3</sup><https://www.tensorflow.org/>

225 for any node  $i$  and  $p(\hat{y}|s=1)$  is the PDF of  $\hat{y}_{i,1}$  for any node  $i$  belong to the demographic group  
 226 with sensitive attribute value  $s=1$ . We observe that  $P[\hat{y}=1]$  and  $P[\hat{y}=1|s=1]$  are equivalent  
 227 to the cumulative distribution functions (CDF) of  $p(\hat{y} < \frac{1}{2})$  and  $p(\hat{y} < \frac{1}{2}|s=1)$ , respectively. To  
 228 estimate both  $P[\hat{y}=1]$  and  $P[\hat{y}=1|s=1]$  with a differentiable function, we first estimate their  
 229 probability density functions ( $p(\hat{y} < \frac{1}{2})$  and  $p(\hat{y} < \frac{1}{2}|s=1)$ ) with kernel density estimation (KDE,  
 230 Definition 1).

231 **Definition 1** (Kernel density estimation [7]) Given a set of  $n$  IID samples  $\{x_1, \dots, x_n\}$  drawn from  
 232 a distribution with an unknown probability density function  $f$ , the kernel density estimation of  $f$  at  
 233 point  $\tau$  is defined as follows.

$$\tilde{f}(\tau) = \frac{1}{na} \sum_{i=1}^n f_k\left(\frac{\tau - x_i}{a}\right) \quad (9)$$

234 where  $\tilde{f}$  is the estimated probability density function,  $f_k$  is the kernel function and  $a$  is a non-negative  
 235 bandwidth.

236 Moreover, we assume the kernel function in KDE is the Gaussian kernel  $f_k(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$ .  
 237 However, computing the CDF of a Gaussian distribution is non-trivial. Following [9], we leverage a  
 238 tractable approximation of the Gaussian Q-function as follows.

$$Q(\tau) = F_k(\tau) = \int_{\tau}^{\infty} f_k(x)dx \approx e^{-\alpha\tau^2 - \beta\tau - \gamma} \quad (10)$$

239 where  $f_k(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$  is a Gaussian distribution with zero mean,  $\alpha = 0.4920$ ,  $\beta = 0.2887$ ,  
 240  $\gamma = 1.1893$  [30]. The overall workflow of estimating  $P[\hat{y}=1]$  is as follows.

- 241 • For any node  $i$ , get its prediction probability  $\hat{y}_{i,1}$  with respect to class 1;
- 242 • Estimate the CDF  $P[\hat{y}=1]$  using a Gaussian KDE with bandwidth  $a$  by  $P[\hat{y}=1] =$   
 243  $\frac{1}{n} \sum_{i=1}^n \exp\left(-\alpha\left(\frac{0.5-\hat{y}_{i,1}}{a}\right)^2 - \beta\left(\frac{0.5-\hat{y}_{i,1}}{a}\right) - \gamma\right)$ , where  $\alpha = 0.4920$ ,  $\beta = 0.2887$ ,  $\gamma =$   
 244  $1.1893$  and  $\exp(x) = e^x$ .

245 Note that  $P[\hat{y}=1|s=1]$  can be estimated with a similar procedure with minor modifications. The  
 246 only modifications needed are: (1) get the prediction probability of nodes with  $s=1$  and (2) compute  
 247 the CDF using the Gaussian Q-function over nodes with  $s=1$  rather than all nodes in the graph.

## 248 5 Instantiation #2: Individual Fairness on Graph Neural Networks

249 We provide another instantiation of FATE framework by attacking individual fairness on graph neural  
 250 networks. Here, we consider the same surrogate graph learning model (i.e., 2-layer linear GCN)  
 251 and the same lower-level loss function (i.e., cross entropy) as described in Section 4. To attack  
 252 individual fairness, we define the upper-level bias function following the principles in [20]: the  
 253 fairness-related auxiliary information matrix  $\mathbf{F}$  is defined as the oracle symmetric pairwise node  
 254 similarity matrix  $\mathbf{S}$  (i.e.,  $\mathbf{F} = \mathbf{S}$ ), where  $\mathbf{S}[i, j]$  measures the similarity between node  $i$  and node  
 255  $j$ . Kang et al. [20] define that the overall individual bias to be  $\text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$ . Assuming that  $\mathbf{Y}$   
 256 is the output of an optimization-based graph learning model,  $\mathbf{Y}$  can be viewed as a function with  
 257 respect to the input graph  $\mathcal{G}$ , which makes  $\text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$  differentiable with respect to  $\mathcal{G}$ . Thus, the  
 258 bias function  $b(\cdot)$  can be naturally defined as the overall individual bias of the input graph  $\mathcal{G}$ , i.e.,  
 259  $b(\mathbf{Y}, \Theta^*, \mathbf{S}) = \text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$ .

## 260 6 Experiments

### 261 6.1 Attacking Statistical Parity on Graph Neural Networks

262 **Settings.** We compare FATE with 4 baseline methods, i.e., Random, DICE [46], FA-GNN [19], under  
 263 the same setting as in Section 4. That is, (1) the fairness definition to be attacked is statistical parity;  
 264 (2) the downstream task is binary semi-supervised node classification with binary sensitive attributes.  
 265 The experiments are conducted on 3 real-world datasets, i.e., Pokec-n, Pokec-z and Bail. Similar  
 266 to existing works, we use the 50%/25%/25% splits for train/validation/test sets. For all baseline

Table 1: Effectiveness of attacking group fairness on GCN. FATE poisons the graph via both edge flipping (FATE-flip) and edge addition (FATE-add) while all other baselines poison the graph via edge addition. Higher is better ( $\uparrow$ ) for micro F1 score (Micro F1) and  $\Delta_{SP}$ . Bold font indicates the success of fairness attack (i.e.,  $\Delta_{SP}$  is increased after fairness attack) with the highest micro F1 score. Underlined cell indicates the failure of fairness attack (i.e.,  $\Delta_{SP}$  is decreased after fairness attack).

Dataset	Ptb.	Random		DICE		FA-GNN		FATE-flip		FATE-add	
		Micro F1 ( $\uparrow$ )	$\Delta_{SP}$ ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	$\Delta_{SP}$ ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	$\Delta_{SP}$ ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	$\Delta_{SP}$ ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	$\Delta_{SP}$ ( $\uparrow$ )
Pokec-n	0.00	67.5 $\pm$ 0.3	7.1 $\pm$ 0.4	67.5 $\pm$ 0.3	7.1 $\pm$ 0.4						
	0.05	68.0 $\pm$ 0.3	6.2 $\pm$ 0.8	67.6 $\pm$ 0.2	6.8 $\pm$ 0.3	67.8 $\pm$ 0.1	3.3 $\pm$ 0.4	<b>67.9 <math>\pm</math> 0.4</b>	<b>9.3 <math>\pm</math> 1.2</b>	<b>67.9 <math>\pm</math> 0.4</b>	<b>9.3 <math>\pm</math> 1.2</b>
	0.10	66.8 $\pm$ 0.8	7.3 $\pm$ 0.7	<u>66.1 <math>\pm</math> 0.5</u>	<u>6.6 <math>\pm</math> 1.1</u>	66.0 $\pm$ 0.2	11.5 $\pm$ 0.6	<b>68.2 <math>\pm</math> 0.6</b>	<b>9.8 <math>\pm</math> 1.5</b>	<b>68.2 <math>\pm</math> 0.6</b>	<b>9.8 <math>\pm</math> 1.5</b>
	0.15	66.7 $\pm$ 0.4	8.1 $\pm$ 0.4	65.6 $\pm$ 0.4	7.7 $\pm$ 0.8	66.0 $\pm$ 0.4	15.6 $\pm$ 3.0	<b>68.0 <math>\pm</math> 0.3</b>	<b>11.5 <math>\pm</math> 1.0</b>	<b>68.0 <math>\pm</math> 0.3</b>	<b>11.5 <math>\pm</math> 1.0</b>
	0.20	66.3 $\pm$ 0.7	8.6 $\pm$ 1.8	<u>64.2 <math>\pm</math> 0.4</u>	<u>3.4 <math>\pm</math> 0.9</u>	65.8 $\pm$ 0.1	18.4 $\pm$ 0.7	<b>68.2 <math>\pm</math> 0.5</b>	<b>12.0 <math>\pm</math> 1.8</b>	<b>68.2 <math>\pm</math> 0.5</b>	<b>12.0 <math>\pm</math> 1.8</b>
	0.25	66.2 $\pm$ 0.6	8.5 $\pm$ 0.8	<u>63.4 <math>\pm</math> 0.2</u>	<u>6.3 <math>\pm</math> 0.8</u>	66.6 $\pm$ 0.2	23.3 $\pm$ 0.5	<b>68.3 <math>\pm</math> 0.4</b>	<b>12.1 <math>\pm</math> 2.1</b>	<b>68.3 <math>\pm</math> 0.4</b>	<b>12.1 <math>\pm</math> 2.1</b>
Pokec-z	0.00	68.4 $\pm$ 0.4	6.6 $\pm$ 0.9	68.4 $\pm$ 0.4	6.6 $\pm$ 0.9						
	0.05	68.8 $\pm$ 0.4	6.4 $\pm$ 0.6	67.4 $\pm$ 0.5	6.6 $\pm$ 0.3	68.1 $\pm$ 0.3	2.2 $\pm$ 0.4	<b>68.7 <math>\pm</math> 0.4</b>	<b>6.7 <math>\pm</math> 1.4</b>	<b>68.7 <math>\pm</math> 0.4</b>	<b>6.7 <math>\pm</math> 1.4</b>
	0.10	<b>68.7 <math>\pm</math> 0.3</b>	<b>8.0 <math>\pm</math> 0.6</b>	<u>66.5 <math>\pm</math> 0.2</u>	<u>6.3 <math>\pm</math> 0.8</u>	67.7 $\pm$ 0.4	13.5 $\pm$ 0.9	<b>68.7 <math>\pm</math> 0.6</b>	<b>7.5 <math>\pm</math> 0.7</b>	<b>68.7 <math>\pm</math> 0.6</b>	<b>7.5 <math>\pm</math> 0.7</b>
	0.15	67.9 $\pm$ 0.3	9.1 $\pm$ 0.8	<u>65.9 <math>\pm</math> 0.8</u>	<u>5.5 <math>\pm</math> 1.3</u>	66.6 $\pm$ 0.4	16.9 $\pm$ 2.6	<b>69.0 <math>\pm</math> 0.8</b>	<b>8.5 <math>\pm</math> 1.1</b>	<b>69.0 <math>\pm</math> 0.8</b>	<b>8.5 <math>\pm</math> 1.1</b>
	0.20	<b>68.5 <math>\pm</math> 0.4</b>	<b>9.3 <math>\pm</math> 1.0</b>	62.9 $\pm$ 0.7	8.7 $\pm$ 1.0	66.1 $\pm$ 0.2	25.4 $\pm$ 1.3	<b>68.5 <math>\pm</math> 0.6</b>	<b>8.8 <math>\pm</math> 1.1</b>	<b>68.5 <math>\pm</math> 0.6</b>	<b>8.8 <math>\pm</math> 1.1</b>
	0.25	68.3 $\pm$ 0.5	7.3 $\pm$ 0.5	<u>63.9 <math>\pm</math> 0.4</u>	<u>6.0 <math>\pm</math> 1.0</u>	65.5 $\pm$ 0.6	22.3 $\pm$ 2.8	<b>68.5 <math>\pm</math> 1.1</b>	<b>8.6 <math>\pm</math> 2.5</b>	<b>68.5 <math>\pm</math> 1.1</b>	<b>8.6 <math>\pm</math> 2.5</b>
Bail	0.00	93.1 $\pm$ 0.2	8.0 $\pm$ 0.2	93.1 $\pm$ 0.2	8.0 $\pm$ 0.2						
	0.05	<b>92.7 <math>\pm</math> 0.2</b>	<b>8.1 <math>\pm</math> 0.0</b>	91.6 $\pm$ 0.2	8.5 $\pm$ 0.1	91.7 $\pm$ 0.1	10.0 $\pm$ 0.4	92.6 $\pm$ 0.1	8.6 $\pm$ 0.1	92.5 $\pm$ 0.1	8.6 $\pm$ 0.1
	0.10	<u>92.2 <math>\pm</math> 0.2</u>	<u>7.8 <math>\pm</math> 0.2</u>	90.3 $\pm$ 0.1	8.5 $\pm$ 0.1	90.5 $\pm$ 0.0	10.3 $\pm$ 0.4	<b>92.4 <math>\pm</math> 0.1</b>	<b>8.9 <math>\pm</math> 0.1</b>	<b>92.4 <math>\pm</math> 0.1</b>	<b>8.9 <math>\pm</math> 0.1</b>
	0.15	<u>91.9 <math>\pm</math> 0.2</u>	<u>7.8 <math>\pm</math> 0.1</u>	89.2 $\pm$ 0.1	7.7 $\pm$ 0.1	90.0 $\pm$ 0.2	8.4 $\pm$ 0.2	92.2 $\pm$ 0.2	9.1 $\pm$ 0.1	<b>92.3 <math>\pm</math> 0.1</b>	<b>9.1 <math>\pm</math> 0.1</b>
	0.20	<u>91.6 <math>\pm</math> 0.2</u>	<u>7.8 <math>\pm</math> 0.1</u>	88.3 $\pm$ 0.1	8.3 $\pm$ 0.1	89.7 $\pm$ 0.1	7.4 $\pm$ 0.4	92.2 $\pm$ 0.2	9.3 $\pm$ 0.1	<b>92.3 <math>\pm</math> 0.1</b>	<b>9.3 <math>\pm</math> 0.2</b>
	0.25	91.4 $\pm$ 0.1	8.3 $\pm$ 0.1	<u>87.8 <math>\pm</math> 0.0</u>	<u>7.8 <math>\pm</math> 0.1</u>	<u>89.8 <math>\pm</math> 0.2</u>	<u>5.2 <math>\pm</math> 0.2</u>	<b>92.1 <math>\pm</math> 0.1</b>	<b>9.1 <math>\pm</math> 0.2</b>	<b>92.1 <math>\pm</math> 0.1</b>	<b>9.1 <math>\pm</math> 0.3</b>

267 methods, the victim models are set to GCN [26]. For each dataset, we use a fixed random seed to  
 268 learn the poisoned graph corresponding to each baseline method. Then we train the victim model 5  
 269 times with different random seeds. For fair comparison, we only attack the adjacency matrix in all  
 270 experiments. Please refer to Appendix C for detailed experimental settings.

271 **Main results.** For FATE, we conduct fairness attacks via both edge flipping (FATE-flip in Table 5)  
 272 and edge addition (FATE-add in Table 1). For all other baseline methods, edges are only added.  
 273 The effectiveness of fairness attacks on GCN are presented in Tables 5. From both tables, we have  
 274 the following key observations: (1) FATE-flip and FATE-add are the only methods that consistently  
 275 succeeds in fairness attacks, while all other baseline methods might fail in some cases (indicated by  
 276 the underlined  $\Delta_{SP}$  in both tables) because of the decrease in  $\Delta_{SP}$ . (2) FATE-flip and FATE-add can  
 277 not only amplify  $\Delta_{SP}$  consistently, but also achieve the best micro F1 score on node classification,  
 278 which makes FATE-flip and FATE-add more deceptive than all baseline methods. Notably, FATE-flip  
 279 and FATE-add are able to even increase micro F1 score on all datasets, while other baseline methods  
 280 attack the graph neural networks at the expense of utility (micro F1 score). (3) Though FA-GNN  
 281 could make the model more biased in some cases, it cannot guarantee consistent success in fairness  
 282 attacks on all three datasets as shown by the underlined  $\Delta_{SP}$  in both tables. All in all, our proposed  
 283 FATE framework is the framework that consistently succeeds in fairness attacks while being the most  
 284 deceptive (i.e., highest micro F1 score).

285 **Effect of the perturbation rate.** From Table 1, we have the following observations. First,  $\Delta_{SP}$  tends  
 286 to increase when the perturbation rate increases, which demonstrates the effectiveness of FATE-flip  
 287 and FATE-add for attacking fairness. Though in some cases  $\Delta_{SP}$  might have a marginal decrease,  
 288 FATE-flip and FATE-add still successfully attack the fairness compared with GCN trained on the  
 289 benign graph by being larger to the  $\Delta_{SP}$  when perturbation rate (Ptb.) is 0. Second, FATE-flip and  
 290 FATE-add are deceptive, meaning that the micro F1 scores is close to or even higher than the micro  
 291 F1 scores on the benign graph compared with the corresponding metrics trained . In summary, across  
 292 different perturbation rates, FATE-flip and FATE-add are both effective, i.e., amplifying more bias  
 293 with higher perturbation rate, and deceptive, i.e., achieving similar or even higher micro F1 score.

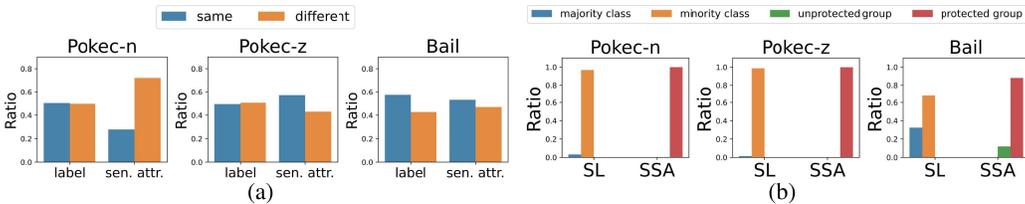


Figure 1: Attacking statistical parity with FATE-flip. (a) Ratios of flipped edges that connect two nodes with same/different label or sensitive attribute (sens. attr.). (b) SL (abbreviation for same label) refers to the ratios of flipped edges whose two endpoints are both from the same class. SSA (abbreviation for same sensitive attribute) refers to the ratios of manipulated edges whose two endpoints are both from the same demographic group. Majority/minority classes are determined by splitting the training nodes based on their class labels. The protected group is the demographic group with fewer nodes.

294 **Analysis on the manipulated edges.** Here, we aim to characterize the properties of edges that are  
 295 flipped by FATE (i.e., FATE-flip) in attacking statistical parity. The reason to only analyze FATE-flip is  
 296 that the majority of edges manipulated by FATE-flip on all three datasets is by addition (i.e., flipping  
 297 from non-existing to existing). Figure 1b suggests that, if the two endpoints of a manipulated edge  
 298 share the same class label or same sensitive attribute value, these two endpoints are most likely from  
 299 the minority class and protected group. Combining Figures 1a and 1b, FATE would significantly  
 300 increase the number of edges that are incident to nodes in the minority class and/or protected group.

301 **More experimental results.** Due to the space limitation, we defer more experimental results on  
 302 attacking statistical parity on graph neural networks in Appendix D. More specifically, we present the  
 303 performance evaluation under different metrics, i.e., Macro F1 and AUC, as well as the effectiveness  
 304 of FATE with a different victim model, i.e., FairGNN [11], which ensures statistical parity.

## 305 6.2 Attacking Individual Fairness on Graph Neural Networks

306 **Settings.** To showcase the ability of FATE on attacking the individual fairness (Section 5), we further  
 307 compare FATE with the same set of baseline methods (Random, DICE [46], FA-GNN [19]) on the  
 308 same set of datasets (Pokec-n, Pokec-z, Bail). We follow the settings as in Section 5. We use the  
 309 50%/25%/25% splits for train/validation/test sets with GCN [26] being the victim model. For each  
 310 dataset, we use a fixed random seed to learn the poisoned graph corresponding to each baseline  
 311 method. Then we train the victim model 5 times with different random seeds. And each entry in the  
 312 oracle pairwise node similarity matrix is computed by the cosine similarity of the corresponding rows  
 313 in the adjacency matrix. That is,  $S[i, j] = \cos(\mathbf{A}[i, :], \mathbf{A}[j, :])$ , where  $\cos()$  is the function to compute  
 314 cosine similarity. For fair comparison, we only attack the adjacency matrix in all experiments. Please  
 315 refer to Appendix C for detailed experimental settings.

316 **Main results.** Similarly, we test FATE with both edge flipping (FATE-flip in Table 2) and edge  
 317 addition (FATE-add in Table 2), while all other baseline methods only add edges. From Table 2, we  
 318 have two key observations. (1) FATE-flip and FATE-add are effective: they are the only methods that  
 319 could consistently attack individual fairness whereas all other baseline methods mostly fail to attack  
 320 individual fairness. (2) FATE-flip and FATE-add are deceptive: they achieve comparable or even better  
 321 utility on all datasets compared with the utility on the benign graph. Hence, FATE framework is able  
 322 to achieve effective and deceptive attacks to exacerbate individual bias.

323 **Effect of the perturbation rate.** From Table 2, we obtain similar observations as in Section 6.1 for  
 324 Bail dataset. While for Pokec-n and Pokec-z, the correlation between the perturbation rate (Ptb.) and  
 325 the individual bias is weaker. One possible reason is that: for Pokec-n and Pokec-z, the discrepancy  
 326 between the oracle pairwise node similarity matrix and the benign graph is larger. Since the individual  
 327 bias is computed using the oracle pairwise node similarity matrix rather than the benign/poisoned  
 328 adjacency matrix, higher perturbation rate to poison the adjacency matrix may have less impact on  
 329 the computation of individual bias.

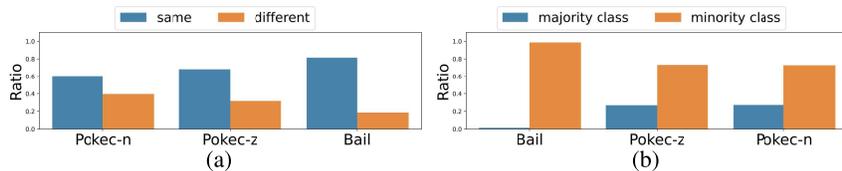


Figure 2: Attacking individual fairness with FATE-flip. (a) Ratios of flipped edges that connect two nodes with same/different label. (b) Ratios of flipped edges whose two endpoints are both from the majority/minority class. Majority/minority classes are formed by splitting the training nodes based on their class labels.

330 **Analysis on the manipulated edges.** Similarly, since the majority of edges manipulated by FATE-flip  
 331 is through addition, we only analyze FATE-flip here. From Figure 2, we can find out that FATE will  
 332 manipulate edges from the same class (especially from the minority class). In this way, FATE would  
 333 find edges that could increase individual bias and improve the utility of the minority class in order to  
 334 make the fairness attack deceptive.

335 **More experimental results.** Due to the space limitation, we defer more experimental results  
 336 on attacking individual fairness on graph neural networks in Appendix E. More specifically, we  
 337 present the performance evaluation under different metrics, i.e., Macro F1 and AUC, as well as

Table 2: Effectiveness of attacking individual fairness on GCN. FATE poisons the graph via both edge flipping (FATE-flip) and edge addition (FATE-add) while all other baselines poison the graph via edge addition. Higher is better ( $\uparrow$ ) for micro F1 score (Micro F1) and InFoRM bias (Bias). Bold font indicates the success of fairness attack (i.e., bias is increased after attack) with the highest micro F1 score. Underlined cell indicates the failure of fairness attack (i.e.,  $\Delta_{SP}$  is decreased after attack).

Dataset	Ptb.	Random		DICE		FA-GNN		FATE-flip		FATE-add	
		Micro F1 ( $\uparrow$ )	Bias ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	Bias ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	Bias ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	Bias ( $\uparrow$ )	Micro F1 ( $\uparrow$ )	Bias ( $\uparrow$ )
Pokey-n	0.00	67.5 $\pm$ 0.3	0.9 $\pm$ 0.2								
	0.05	67.6 $\pm$ 0.3	1.6 $\pm$ 0.3	66.9 $\pm$ 0.3	1.6 $\pm$ 0.2	<b>67.8 <math>\pm</math> 0.5</b>	<b>1.9 <math>\pm</math> 0.2</b>	<b>67.8 <math>\pm</math> 0.3</b>	<b>1.2 <math>\pm</math> 0.4</b>	67.6 $\pm$ 0.3	1.5 $\pm$ 0.6
	0.10	67.2 $\pm$ 0.5	1.4 $\pm$ 0.3	65.3 $\pm$ 0.7	1.1 $\pm$ 0.1	67.4 $\pm$ 0.4	1.2 $\pm$ 0.2	<b>67.9 <math>\pm</math> 0.4</b>	<b>1.3 <math>\pm</math> 0.3</b>	67.7 $\pm$ 0.4	1.6 $\pm$ 0.4
	0.15	67.2 $\pm$ 0.3	1.2 $\pm$ 0.4	63.9 $\pm$ 0.6	1.1 $\pm$ 0.2	66.1 $\pm$ 0.3	1.5 $\pm$ 0.3	<b>67.8 <math>\pm</math> 0.4</b>	<b>1.2 <math>\pm</math> 0.2</b>	67.6 $\pm$ 0.2	1.1 $\pm$ 0.3
	0.20	66.6 $\pm$ 0.3	1.1 $\pm$ 0.2	<u>63.8 <math>\pm</math> 0.1</u>	<u>0.8 <math>\pm</math> 0.1</u>	65.7 $\pm$ 0.6	1.5 $\pm$ 0.3	67.3 $\pm$ 0.4	1.1 $\pm$ 0.3	<b>68.2 <math>\pm</math> 1.0</b>	<b>1.7 <math>\pm</math> 0.8</b>
	0.25	66.7 $\pm$ 0.3	1.3 $\pm$ 0.4	62.5 $\pm$ 0.4	0.6 $\pm$ 0.0	65.2 $\pm$ 0.5	1.3 $\pm$ 0.4	<b>67.8 <math>\pm</math> 0.8</b>	<b>1.4 <math>\pm</math> 0.7</b>	<b>67.9 <math>\pm</math> 0.9</b>	<b>1.4 <math>\pm</math> 0.7</b>
Pokey-z	0.00	68.4 $\pm$ 0.4	2.6 $\pm$ 0.7								
	0.05	<b>69.0 <math>\pm</math> 0.4</b>	<b>3.4 <math>\pm</math> 0.5</b>	67.1 $\pm$ 0.5	2.7 $\pm$ 1.0	68.1 $\pm$ 0.4	2.9 $\pm$ 0.3	68.7 $\pm$ 0.5	2.9 $\pm$ 0.5	68.7 $\pm$ 0.4	3.1 $\pm$ 1.0
	0.10	68.7 $\pm$ 0.1	2.4 $\pm$ 0.5	66.3 $\pm$ 0.6	1.7 $\pm$ 0.6	68.2 $\pm$ 0.5	1.7 $\pm$ 0.5	<b>69.0 <math>\pm</math> 0.6</b>	<b>2.9 <math>\pm</math> 0.6</b>	<b>69.0 <math>\pm</math> 0.5</b>	<b>3.0 <math>\pm</math> 0.6</b>
	0.15	67.9 $\pm$ 0.3	2.8 $\pm$ 0.3	<u>65.5 <math>\pm</math> 0.3</u>	<u>1.4 <math>\pm</math> 0.3</u>	67.0 $\pm$ 0.5	1.3 $\pm$ 0.2	68.6 $\pm$ 0.5	2.9 $\pm$ 0.6	<b>69.0 <math>\pm</math> 0.7</b>	<b>2.7 <math>\pm</math> 0.4</b>
	0.20	67.9 $\pm$ 0.3	2.2 $\pm$ 0.6	<u>64.2 <math>\pm</math> 0.4</u>	<u>0.7 <math>\pm</math> 0.3</u>	66.1 $\pm$ 0.1	1.6 $\pm$ 0.5	68.8 $\pm$ 0.4	3.0 $\pm$ 0.4	<b>69.2 <math>\pm</math> 0.4</b>	<b>2.9 <math>\pm</math> 0.3</b>
	0.25	67.6 $\pm$ 0.3	1.9 $\pm$ 0.3	<u>64.2 <math>\pm</math> 0.3</u>	<u>0.5 <math>\pm</math> 0.1</u>	65.1 $\pm$ 0.3	1.9 $\pm$ 0.6	69.1 $\pm$ 0.3	2.9 $\pm$ 0.7	<b>69.3 <math>\pm</math> 0.3</b>	<b>2.7 <math>\pm</math> 0.6</b>
Bail	0.00	93.1 $\pm$ 0.2	7.2 $\pm$ 0.6								
	0.05	92.1 $\pm$ 0.3	8.0 $\pm$ 1.9	<u>91.8 <math>\pm</math> 0.1</u>	<u>7.1 <math>\pm</math> 1.1</u>	91.2 $\pm$ 0.2	5.6 $\pm$ 0.7	<b>93.0 <math>\pm</math> 0.3</b>	7.8 $\pm$ 1.0	92.9 $\pm$ 0.2	7.7 $\pm$ 1.0
	0.10	91.6 $\pm$ 0.1	7.3 $\pm$ 1.2	<u>90.3 <math>\pm</math> 0.1</u>	6.1 $\pm$ 0.6	<u>90.3 <math>\pm</math> 0.1</u>	5.1 $\pm$ 0.4	<b>93.0 <math>\pm</math> 0.1</b>	<b>8.0 <math>\pm</math> 0.7</b>	92.9 $\pm$ 0.2	7.9 $\pm$ 0.8
	0.15	<u>91.3 <math>\pm</math> 0.1</u>	6.5 $\pm$ 0.9	89.4 $\pm$ 0.0	4.8 $\pm$ 0.1	89.8 $\pm$ 0.1	5.2 $\pm$ 0.1	<b>93.1 <math>\pm</math> 0.1</b>	<b>8.2 <math>\pm</math> 0.6</b>	93.0 $\pm$ 0.2	7.8 $\pm$ 0.8
	0.20	<u>91.2 <math>\pm</math> 0.2</u>	6.6 $\pm$ 0.6	88.5 $\pm$ 0.1	4.0 $\pm$ 0.4	89.3 $\pm$ 0.1	5.3 $\pm$ 0.4	<b>93.1 <math>\pm</math> 0.1</b>	<b>7.9 <math>\pm</math> 0.6</b>	<b>93.1 <math>\pm</math> 0.1</b>	<b>8.2 <math>\pm</math> 0.6</b>
	0.25	<u>90.9 <math>\pm</math> 0.1</u>	6.8 $\pm$ 0.8	87.4 $\pm$ 0.3	3.6 $\pm$ 0.5	88.9 $\pm$ 0.1	5.4 $\pm$ 0.3	92.9 $\pm$ 0.1	7.6 $\pm$ 0.5	<b>93.0 <math>\pm</math> 0.2</b>	<b>7.8 <math>\pm</math> 0.7</b>

338 the effectiveness of FATE with a different victim model, i.e., InFoRM-GNN [20], which mitigates  
 339 individual bias.

## 340 7 Related Work

341 **Algorithmic fairness on graphs** aims to obtain debiased graph learning results such that a pre-  
 342 defined fairness definition can be satisfied with respect to the nodes/edges in the graph. Several  
 343 definitions of the fairness has been studied so far. Group fairness in graph embedding can be  
 344 ensured via several ways, including adversarial learning-based methods [5, 11], random walk-based  
 345 methods [36, 25] and dropout-based methods [39]. Individual fairness on graphs can be ensured  
 346 via Lipschitz regularization [20] and learning-to-rank [13]. Other than the aforementioned two  
 347 fairness definitions, several other fairness definitions are studied in the context of graph learning,  
 348 including counterfactual fairness [1, 31], degree fairness [42, 24, 29], dyadic fairness [32, 27] and  
 349 max-min fairness [37, 43]. For a comprehensive review of related works, please refer to existing  
 350 surveys [50, 10, 14] and tutorials [22, 23]. It should be noted that our work aims to attack fairness  
 351 (i.e., making the model more biased) rather than ensuring fairness as in the aforementioned literature.

352 **Adversarial attacks on graphs** aim to exacerbate the utility of graph learning models by perturbing  
 353 the input graph topology and/or node features. Several approaches have been proposed to attack graph  
 354 learning models, including reinforcement learning [12], bi-level optimization [51, 52], projected  
 355 gradient descent [40, 48] and edge rewiring/flipping [4, 31]. Other than adversarial attacks that  
 356 worsen the utility of a graph learning model, a few efforts have been made to attack the fairness of a  
 357 machine learning model for IID tabular data via label flipping [33], adversarial data injection [38, 8],  
 358 adversarial sampling [44]. Different from [38, 33, 8, 44], we aim to poison the input graph via  
 359 structural modifications on the topology rather than injecting adversarial data sample(s). The most  
 360 related work to our proposed method is by Hussain et al. [19], which degrade the group fairness of  
 361 graph neural networks by randomly injecting edges for nodes in different demographic groups and  
 362 with different class labels. In contrast, our proposed method could attack *any* fairness definition for  
 363 *any* graph learning models via arbitrary edge manipulation operations, as long as the bias function  
 364 and the utility loss are differentiable.

## 365 8 Conclusion

366 We study the problem of fairness attacks on graph learning models, whose goal is to amplify the bias  
 367 while maintaining the utility on the downstream task. We formally define the problem as a bi-level  
 368 optimization problem, where the upper-level optimization problem maximizes the bias function with  
 369 respect to a user-defined fairness definition and the lower-level optimization problem minimizes a  
 370 task-specific loss function. We then propose a meta learning-based framework named FATE to poison  
 371 the input graph using the meta-gradient of the bias function with respect to the input graph. We  
 372 instantiate FATE by attacking statistical parity on graph neural networks in a binary node classification  
 373 problem with binary sensitive attributes. Empirical evaluation demonstrates that FATE is effective  
 374 (consistently amplifying bias) and deceptive (achieving the highest micro F1 score).

## References

- 375
- 376 [1] Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. Towards a unified framework for  
377 fair and stable graph representation learning. In *Uncertainty in Artificial Intelligence*, pages  
378 2114–2124. PMLR, 2021.
- 379 [2] Joachim Baumann, Anikó Hannák, and Christoph Heitz. Enforcing group fairness in algorithmic  
380 decision making: Utility maximization under sufficiency. In *2022 ACM Conference on Fairness,  
381 Accountability, and Transparency*, pages 2315–2326, 2022.
- 382 [3] Yoshua Bengio. Gradient-based optimization of hyperparameters. *Neural computation*,  
383 12(8):1889–1900, 2000.
- 384 [4] Aleksandar Bojchevski and Stephan Günnemann. Adversarial attacks on node embeddings via  
385 graph poisoning. In *International Conference on Machine Learning*, pages 695–704. PMLR,  
386 2019.
- 387 [5] Avishek Bose and William Hamilton. Compositional fairness constraints for graph embeddings.  
388 In *International Conference on Machine Learning*, pages 715–724. PMLR, 2019.
- 389 [6] Consumer Financial Protection Bureau. CFPB targets unfair discrimination in  
390 consumer finance. [https://www.consumerfinance.gov/about-us/newsroom/  
391 cfpb-targets-unfair-discrimination-in-consumer-finance/](https://www.consumerfinance.gov/about-us/newsroom/cfpb-targets-unfair-discrimination-in-consumer-finance/), 2022. [Online;  
392 accessed 13-April-2023].
- 393 [7] Yen-Chi Chen. A tutorial on kernel density estimation and recent advances. *Biostatistics &  
394 Epidemiology*, 1(1):161–187, 2017.
- 395 [8] Anshuman Chhabra, Adish Singla, and Prasant Mohapatra. Fairness degrading adversarial  
396 attacks against clustering algorithms. *arXiv preprint arXiv:2110.12020*, 2021.
- 397 [9] Jaewoong Cho, Gyeongjo Hwang, and Changho Suh. A fair classifier using kernel density  
398 estimation. *Advances in neural information processing systems*, 33:15088–15099, 2020.
- 399 [10] Manvi Choudhary, Charlotte Laclau, and Christine Largeron. A survey on fairness for machine  
400 learning on graphs. *arXiv preprint arXiv:2205.05396*, 2022.
- 401 [11] Enyan Dai and Suhang Wang. Say no to the discrimination: Learning fair graph neural networks  
402 with limited sensitive attribute information. In *Proceedings of the 14th ACM International  
403 Conference on Web Search and Data Mining*, pages 680–688, 2021.
- 404 [12] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack  
405 on graph structured data. In *International conference on machine learning*, pages 1115–1124.  
406 PMLR, 2018.
- 407 [13] Yushun Dong, Jian Kang, Hanghang Tong, and Jundong Li. Individual fairness for graph neural  
408 networks: A ranking based approach. In *Proceedings of the 27th ACM SIGKDD Conference on  
409 Knowledge Discovery & Data Mining*, pages 300–310, 2021.
- 410 [14] Yushun Dong, Jing Ma, Chen Chen, and Jundong Li. Fairness in graph mining: A survey. *arXiv  
411 preprint arXiv:2204.09888*, 2022.
- 412 [15] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness  
413 through awareness. In *Proceedings of the 3rd innovations in theoretical computer science  
414 conference*, pages 214–226, 2012.
- 415 [16] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkata-  
416 subramanian. Certifying and removing disparate impact. In *proceedings of the 21th ACM  
417 SIGKDD international conference on knowledge discovery and data mining*, pages 259–268,  
418 2015.
- 419 [17] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-  
420 tation of deep networks. In *International conference on machine learning*, pages 1126–1135.  
421 PMLR, 2017.

- 422 [18] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning.  
423 *Advances in neural information processing systems*, 29, 2016.
- 424 [19] Hussain Hussain, Meng Cao, Sandipan Sikdar, Denis Helic, Elisabeth Lex, Markus Strohmaier,  
425 and Roman Kern. Adversarial inter-group link injection degrades the fairness of graph neural  
426 networks. *arXiv preprint arXiv:2209.05957*, 2022.
- 427 [20] Jian Kang, Jingrui He, Ross Maciejewski, and Hanghang Tong. Inform: Individual fairness  
428 on graph mining. In *Proceedings of the 26th ACM SIGKDD International Conference on*  
429 *Knowledge Discovery & Data Mining*, pages 379–389, 2020.
- 430 [21] Jian Kang and Hanghang Tong. N2n: Network derivative mining. In *Proceedings of the 28th*  
431 *ACM International Conference on Information and Knowledge Management*, pages 861–870,  
432 2019.
- 433 [22] Jian Kang and Hanghang Tong. Fair graph mining. In *Proceedings of the 30th ACM International*  
434 *Conference on Information & Knowledge Management*, pages 4849–4852, 2021.
- 435 [23] Jian Kang and Hanghang Tong. Algorithmic fairness on graphs: Methods and trends. In  
436 *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,  
437 pages 4798–4799, 2022.
- 438 [24] Jian Kang, Yan Zhu, Yinglong Xia, Jiebo Luo, and Hanghang Tong. Rawlsgcn: Towards  
439 rawlsian difference principle on graph convolutional network. In *Proceedings of the ACM Web*  
440 *Conference 2022*, pages 1214–1225, 2022.
- 441 [25] Ahmad Khajehnejad, Moein Khajehnejad, Mahmoudreza Babaei, Krishna P Gummadi, Adrian  
442 Weller, and Baharan Mirzsoleiman. Crosswalk: Fairness-enhanced node representation learn-  
443 ing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages  
444 11963–11970, 2022.
- 445 [26] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional  
446 networks. In *International Conference on Learning Representations*, 2017.
- 447 [27] Peizhao Li, Yifei Wang, Han Zhao, Pengyu Hong, and Hongfu Liu. On dyadic fairness:  
448 Exploring and mitigating bias in graph connections. In *International Conference on Learning*  
449 *Representations*, 2021.
- 450 [28] Lydia T Liu, Sarah Dean, Esther Rolf, Max Simchowitz, and Moritz Hardt. Delayed impact of  
451 fair machine learning. In *International Conference on Machine Learning*, pages 3150–3158.  
452 PMLR, 2018.
- 453 [29] Zemin Liu, Trung-Kien Nguyen, and Yuan Fang. On generalized degree fairness in graph neural  
454 networks. *arXiv preprint arXiv:2302.03881*, 2023.
- 455 [30] Miguel López-Benítez and Fernando Casadevall. Versatile, accurate, and analytically tractable  
456 approximation for the gaussian q-function. *IEEE Transactions on Communications*, 59(4):917–  
457 922, 2011.
- 458 [31] Yao Ma, Suhang Wang, Tyler Derr, Lingfei Wu, and Jiliang Tang. Graph adversarial attack via  
459 rewiring. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery &*  
460 *Data Mining*, pages 1161–1169, 2021.
- 461 [32] Farzan Masrour, Tyler Wilson, Heng Yan, Pang-Ning Tan, and Abdol Esfahanian. Bursting the  
462 filter bubble: Fairness-aware network link prediction. In *Proceedings of the AAAI conference*  
463 *on artificial intelligence*, volume 34, pages 841–848, 2020.
- 464 [33] Ninareh Mehrabi, Muhammad Naveed, Fred Morstatter, and Aram Galstyan. Exacerbating  
465 algorithmic bias through fairness attacks. In *Proceedings of the AAAI Conference on Artificial*  
466 *Intelligence*, volume 35, pages 8930–8938, 2021.
- 467 [34] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation  
468 ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

- 469 [35] Flavien Prost, Hai Qian, Qiuwen Chen, Ed H Chi, Jilin Chen, and Alex Beutel. Toward a better  
470 trade-off between performance and fairness with kernel-based distribution matching. *arXiv*  
471 *preprint arXiv:1910.11779*, 2019.
- 472 [36] Tahleen Rahman, Bartłomiej Surma, Michael Backes, and Yang Zhang. Fairwalk: Towards  
473 fair graph embedding. In *Proceedings of the 28th International Joint Conference on Artificial*  
474 *Intelligence*, pages 3289–3295, 2019.
- 475 [37] Aida Rahmattalabi, Phebe Vayanos, Anthony Fulginiti, Eric Rice, Bryan Wilder, Amulya Yadav,  
476 and Milind Tambe. Exploring algorithmic fairness in robust graph covering problems. *Advances*  
477 *in Neural Information Processing Systems*, 32, 2019.
- 478 [38] David Solans, Battista Biggio, and Carlos Castillo. Poisoning attacks on algorithmic fairness.  
479 In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML*  
480 *PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part I*, pages 162–177.  
481 Springer, 2021.
- 482 [39] Indro Spinelli, Simone Scardapane, Amir Hussain, and Aurelio Uncini. Fairdrop: Biased edge  
483 dropout for enhancing fairness in graph representation learning. *IEEE Transactions on Artificial*  
484 *Intelligence*, 3(3):344–354, 2021.
- 485 [40] Mingjie Sun, Jian Tang, Huichen Li, Bo Li, Chaowei Xiao, Yao Chen, and Dawn Song. Data poi-  
486 soning attack against unsupervised node embedding methods. *arXiv preprint arXiv:1810.12881*,  
487 2018.
- 488 [41] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-  
489 scale information network embedding. In *Proceedings of the 24th international conference on*  
490 *world wide web*, pages 1067–1077, 2015.
- 491 [42] Xianfeng Tang, Huaxiu Yao, Yiwei Sun, Yiqi Wang, Jiliang Tang, Charu Aggarwal, Prasenjit  
492 Mitra, and Suhang Wang. Investigating and mitigating degree-related biases in graph convolu-  
493 tional networks. In *Proceedings of the 29th ACM International Conference on Information &*  
494 *Knowledge Management*, pages 1435–1444, 2020.
- 495 [43] Alan Tsang, Bryan Wilder, Eric Rice, Milind Tambe, and Yair Zick. Group-fairness in influ-  
496 ence maximization. In *Proceedings of the 28th International Joint Conference on Artificial*  
497 *Intelligence*, pages 5997–6005, 2019.
- 498 [44] Minh-Hao Van, Wei Du, Xintao Wu, and Aidong Lu. Poisoning attacks on fair machine learning.  
499 In *International Conference on Database Systems for Advanced Applications*, pages 370–386.  
500 Springer, 2022.
- 501 [45] Daixin Wang, Jianbin Lin, Peng Cui, Quanhui Jia, Zhen Wang, Yanming Fang, Quan Yu, Jun  
502 Zhou, Shuang Yang, and Yuan Qi. A semi-supervised graph attentive network for financial fraud  
503 detection. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 598–607.  
504 IEEE, 2019.
- 505 [46] Marcin Waniek, Tomasz P Michalak, Michael J Wooldridge, and Talal Rahwan. Hiding  
506 individuals and communities in a social network. *Nature Human Behaviour*, 2(2):139–147,  
507 2018.
- 508 [47] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.  
509 Simplifying graph convolutional networks. In *International conference on machine learning*,  
510 pages 6861–6871. PMLR, 2019.
- 511 [48] Kaidi Xu, Hongge Chen, Sijia Liu, Pin-Yu Chen, Tsui-Wei Weng, Mingyi Hong, and Xue  
512 Lin. Topology attack and defense for graph neural networks: An optimization perspective.  
513 In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages  
514 3961–3967, 2019.
- 515 [49] Si Zhang, Dawei Zhou, Mehmet Yigit Yildirim, Scott Alcorn, Jingrui He, Hasan Davulcu, and  
516 Hanghang Tong. Hidden: Hierarchical dense subgraph detection with application to financial  
517 fraud detection. In *Proceedings of the 2017 SIAM International Conference on Data Mining*,  
518 pages 570–578. SIAM, 2017.

- 519 [50] Wenbin Zhang, Jeremy C Weiss, Shuigeng Zhou, and Toby Walsh. Fairness amidst non-iid  
520 graph data: A literature review. *arXiv preprint arXiv:2202.07170*, 2022.
- 521 [51] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural  
522 networks for graph data. In *Proceedings of the 24th ACM SIGKDD international conference on*  
523 *knowledge discovery & data mining*, pages 2847–2856, 2018.
- 524 [52] Daniel Zügner and Stephan Günnemann. Adversarial attacks on graph neural networks via meta  
525 learning. In *International Conference on Learning Representations*, 2019.