

# TOPOLOGY MATTERS IN FAIR GRAPH LEARNING: A THEORETICAL PILOT STUDY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent advances in fair graph learning observe that graph neural networks (GNNs) further amplify prediction bias compared with multilayer perceptron (MLP), while the reason behind this is unknown. In this paper, as a pilot study, we provide a theoretical understanding of when and why the bias enhancement happens in GCN-like aggregation within contexture stochastic block model (CSBM)<sup>1</sup>. Specifically, we define bias enhancement as higher node representation bias after aggregation compared with that before aggregation. We provide a sufficient condition related to the statistical information of graph data to provably and comprehensively delineate instances of bias enhancement during aggregation. Additionally, the proposed sufficient condition goes beyond intuition, serving as a valuable tool to derive data-centric insights. Motivated by data-centric insights, we propose a fair graph rewiring algorithm, named *FairGR*, to reduce sensitive homophily coefficient while preserving useful graph topology. Experiments on node classification tasks demonstrate that *FairGR* can mitigate the prediction bias with comparable performance on three real-world datasets. Additionally, *FairGR* is compatible with many state-of-the-art methods, such as adding regularization, adversarial debiasing, and Fair mixup via refining graph topology, i.e., *FairGR* is a plug-in fairness method and can be adapted to improve existing fair graph learning strategies. The code is available in <https://anonymous.4open.science/r/FairGR-B65D/>.

## 1 INTRODUCTION

Graph neural networks (GNNs) (Kipf and Welling, 2017; Veličković et al., 2018; Wu et al., 2019) are widely adopted in various domains, such as social media mining (Hamilton et al., 2017), knowledge graph (Hamaguchi et al., 2017) and recommender system (Ying et al., 2018), due to remarkable performance in learning representations (Han et al., 2022a;b; Ling et al., 2023). Graph learning, a topic with growing popularity, aims to learn node representation containing both topological and attribute information in a given graph. Despite the outstanding performance in various tasks, GNNs still inherit or even amplify societal bias from input graph data (Dai and Wang, 2021).

In many real-world graphs, nodes with the same sensitive attribute (e.g., ages) are more likely to connect. We call this phenomenon “topology bias”. To be specific, the ratio of edges with the same sensitive attributes among all edges are 95.30%, 95.06%, and 72.37% for Pokec-n, Pokec-z, and NBA datasets. Even worse, in GNNs, the representation of each node is learned by aggregating the representations of its neighbors. Thus, nodes with the same sensitive attributes will be more similar after the aggregation. Such topology bias generates more similar node representation for those nodes with the same sensitive attribute and leads to higher representation bias.

Even though many GNNs empirically demonstrate higher bias compared with multilayer perceptron (MLP) (Dai and Wang, 2021), which aligns with the aforementioned intuition of topology bias, the comprehensive and rigorous theoretical analysis behind such observation from graph data perspective is missing. A natural question is raised:

*Can we theoretically understand when and why bias enhancement happens from a graph data perspective?*

<sup>1</sup>We leave more general analysis on other aggregation operations and random graph models in future work. See Appendix L for more discussions.

In this work, we move **the first step** to understand when and why bias enhancement happens. Specifically, we first define several graph statistical information (such as sensitive homophily, connection density, and the number of nodes) to generate graph data using contextual stochastic block model (CSBM)<sup>2</sup>. Under generated CSBM graph, we provide a sufficient condition that provably and comprehensively delineates instances of bias enhancement during aggregation. Besides, the sufficient condition goes beyond intuition, serving as a valuable tool to derive data-centric insights. Moreover, motivated by the derived data-centric insights, we develop a fair graph rewiring algorithm, named FairGR, to achieve fair GNN prediction via revising graph topology. More importantly, FairGR is a plug-in data rewiring method and compatible with many fair training strategies, such as regularization, adversarial debiasing, and Fair mixup. In short, the contributions can be summarized as follows:

- To the best of our knowledge, it is the first paper to theoretically investigate why and when bias enhancement happens in GNNs. We provide a sufficient condition that provably and comprehensively delineates instances of bias enhancement during GCN-like aggregation under CSBM. Several data-centric insights can be subsequently derived based on the proposed sufficient condition.
- Motivated by the derived data-centric insights, we propose a graph topology rewiring method, named FairGR, using topology-related fair loss to achieve fairness for node classification tasks.
- We experimentally show that the prediction bias of GNNs is larger than that of MLP on synthetic datasets, and validate the effectiveness of FairGR on three real-world datasets.

## 2 PRELIMINARIES

### 2.1 NOTATIONS

We adopt bold upper-case letters to denote matrices such as  $\mathbf{X}$ , bold lower-case letters such as  $\mathbf{x}$  to denote vectors or random variables, and calligraphic font such as  $\mathcal{X}$  to denote set. Given a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , the  $i$ -th row and  $j$ -th column are denoted as  $\mathbf{X}_i$  and  $\mathbf{X}_{\cdot,j}$ , and the element in  $i$ -th row and  $j$ -th column is  $\mathbf{X}_{i,j}$ . We use  $l_1$  norm of matrix  $\mathbf{X}$  as  $\|\mathbf{X}\|_1 = \sum_{ij} |\mathbf{X}_{ij}|$ . Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  be a graph with the node set  $\mathcal{V} = \{v_1, \dots, v_n\}$  and the undirected edge set  $\mathcal{E} = \{e_1, \dots, e_m\}$ , where  $n, m$  represent the number of node and edge, respectively. The graph structure  $\mathcal{G}$  can be represented as an adjacent matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , where  $\mathbf{A}_{ij} = 1$  if existing edge between node  $v_i$  and node  $v_j$ .  $\mathcal{N}(i)$  denotes the neighbors of node  $v_i$  and  $\tilde{\mathcal{N}}(i) = \mathcal{N}(i) \cup \{v_i\}$  denotes the self-inclusive neighbors. Suppose that each node is associated with a  $d$ -dimensional feature vector and a (binary) sensitive attribute, the feature for all nodes and sensitive attribute is denoted as  $\mathbf{X}_{ori} = \mathbb{R}^{n \times d}$  and  $\mathbf{s} \in \{-1, 1\}^n$ .  $I(\mathbf{s}, \mathbf{X})$  represents the mutual information between the sensitive attribute and node features.  $\mathbf{A} \odot \mathbf{B}$  represents Hadamard product for matrix element-wise multiplication.

### 2.2 HOMOPHILY COEFFICIENT IN GRAPHS

The behaviors of graph neural networks have been investigated in the context of label homophily for connected node pairs in graphs (Ma et al., 2022). Label homophily in graphs is typically defined to characterize the similarity of connected node labels in graphs. Here, similar node pair means that the connected nodes share the same label.

From the perspective of fairness, we also define the sensitive homophily coefficient to represent the sensitive attribute similarity among connected node pairs. Informally, the coefficients for label homophily and sensitive homophily are defined as the fraction of the edges connecting the nodes of the same class label and sensitive attributes in a graph (Zhu et al., 2020; Ma et al., 2022). We also provide the formal definition as follows:

**Definition 2.1 (Label and Sensitive Homophily Coefficient)** *Given a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  with node label vector  $\mathbf{y}$ , node sensitive attribute vector  $\mathbf{s}$ , the label and sensitive attribute homophily coefficients are defined as the fraction of edges that connect nodes with the same labels or sensitive attributes  $\epsilon_{label}(\mathcal{G}, \mathbf{y}) = \frac{1}{|\mathcal{E}|} \sum_{(i,j) \in \mathcal{E}} \mathbb{1}(\mathbf{y}_i = \mathbf{y}_j)$ , and  $\epsilon_{sens}(\mathcal{G}, \mathbf{s}) = \frac{1}{|\mathcal{E}|} \sum_{(i,j) \in \mathcal{E}} \mathbb{1}(\mathbf{s}_i = \mathbf{s}_j)$ , where  $|\mathcal{E}|$  is the number of edges and  $\mathbb{1}(\cdot)$  is the indicator function.*

<sup>2</sup>CSBM is a most popular graph generation model, which has been widely adopted in theoretical analysis on graphs (Ma et al., 2022; Baranwal et al., 2021).

Recent work (Ma et al., 2022; Chien et al., 2021; Zhu et al., 2020) aim to understand the relation between the message passing in GNNs and label homophily from the interactions between GNNs (model) and graph topology (data). For graph data with a high label homophily coefficient, existing work (Ma et al., 2022; Tang et al., 2020) have demonstrated, either provably or empirically, that the nodes with higher node degree obtain more prediction benefits in GCN, compared to the benefits that peripheral nodes obtain. As for graph data with a low label homophily coefficient, GNNs do not necessarily lead to better prediction performance compared with MLP since the node features of neighbors with different labels contaminate the node features during feature aggregation.

Although work (Dai and Wang, 2021) empirically points out that graph data with a large sensitive homophily coefficient may enhance bias in GNNs, the fundamental understanding of message passing and graph data properties, such as sensitive homophily coefficient, is still missing. We provide the pilot theoretical analysis in Section 3.

### 3 UNDERSTANDING WHY TOPOLOGY ENHANCES BIAS IN GNNs AGGREGATION

We aim to understand why topology enhances bias<sup>3</sup> in GNNs aggregation operation. For each node, GNNs aggregate their neighbors’ features to learn their representation. We observe a high sensitive homophily coefficient (which is even higher than the label homophily coefficient) in many real-world datasets. We provide an intuition why topology enhances bias happens. GNNs aggregation<sup>4</sup> smoothes node representations with their neighbors, i.e., node representation with the same sensitive attributes are more similar after aggregation.

However, this intuition is heuristic and the quantitative analysis between the topology bias and representation bias is missing. In this section, (1) we rectify such common belief and quantitatively reveal that **only sufficient high sensitive homophily coefficient** would lead to bias enhancement; (2) we analyze the **influence** of other graph’s statistical information, such as the number of nodes  $n$ , the connection density  $\rho_d = \mathbb{E}_{ij}[\mathbb{P}(\mathbf{A}_{ij} = 1)]$ , where adjacency matrix  $\mathbf{A}$  satisfies  $\mathbf{A}_{ij} \in \{0, 1\}$ , and sensitive group ratio  $c = \mathbb{E}_i[\mathbb{P}(s_i = 1)]$ , where the binary sensitive attribute satisfies  $s_i \in \{-1, 1\}$ , and sensitive homophily coefficient  $\epsilon_{sens} = \mathbb{P}(s_i = s_j | \mathbf{A}_{ij} = 1)$ , in term of bias enhancement<sup>5</sup>

#### 3.1 SYNTHETIC GRAPH

we consider the synthetic random graph generation using contextual stochastic block model (CSBM) (Fortunato and Hric, 2016), including graph topology and node features generation with Gaussian mixture distribution., which is well-studied in existing literature for theoretical investigation on graph data (Van Der Hofstad, 2016; Baranwal et al., 2023; 2021). This model allows for the analysis of aggregation for various dataset statistical information (such as graph size, label homophily coefficient, and sensitive homophily coefficient) in a parameterized manner.<sup>6</sup>

**Graph Topology.** We mainly focus on 4 characteristics of graph topology: the number of nodes  $n$ , the edge density  $\rho_d$ , sensitive homophily coefficient  $\epsilon_{sens}$ , and sensitive group ratio  $c$ . The synthetic random graph generation, including graph topology and node features generation are given by:

**Definition 3.1** ( $(n, \rho_d, \epsilon_{sens}, c)$ -**graph**) *The synthetic random graph  $\mathcal{G}$  sampled from  $(n, \rho_d, \epsilon_{sens}, c)$ -graph satisfies the following properties: 1) the graph node number is  $n$ ; 2) the connection density is given by  $\rho_d$ ; 3) sensitive homophily coefficient is  $\epsilon_{sens}$ ; 4) the sensitive attribute group ( $s_i = 1$ ) ratio is  $c$ ; 5) independent edge generation.*

**Node Features.** We assume that node attributes in synthetic graph follow Gaussian Mixture Model  $GMM(c, \mu_1, \Sigma_1, \mu_2, \Sigma_2)$ . For the node with sensitive attribute  $s_i = -1$  ( $s_i = 1$ ), the node attributes

<sup>3</sup>We use representation bias (See Section 3.2) to measure the representation bias of before/after GNNs aggregation.

<sup>4</sup>The sensitive attribute of each node is excluded for features aggregation due to the unallowable consideration in decision making (Barocas and Selbst, 2016).

<sup>5</sup>We discuss the difference with concentration property in GNNs, and persistent homology in Appendix F.

<sup>6</sup>We leave the analysis of other random graph models or feature distributions in future work.

$\mathbf{X}_i$  follow Gaussian distribution  $P_1 = \mathcal{N}(\mu_1, \Sigma_1)$  ( $P_2 = \mathcal{N}(\mu_2, \Sigma_2)$ ), where the node attributes with the same sensitive attribute are independent and identically distributed, and  $\mu_i, \Sigma_i$  ( $i = 1, 2$ ) represent the mean vector and covariance matrix.

### 3.2 REPRESENTATION BIAS MEASUREMENT

For the analysis of the bias difference of the representation before and after GNNs aggregation, a measurement of the node representations bias is required. In this paper, we adopt the mutual information between sensitive attribute and node attributes  $I(\mathbf{s}, \mathbf{X})$  as the measurement<sup>7</sup>. Note that the exact mutual information  $I(\mathbf{s}, \mathbf{X})$  is intractable to estimate, an upper bound on the exact mutual information is developed as a surrogate metric in the following Theorem 3.2:

**Theorem 3.2** *Suppose the synthetic graph node attribute  $\mathbf{X}$  is generated based on Gaussian Mixture Model  $GMM(c, \mu_1, \Sigma_1, \mu_2, \Sigma_2)$ , i.e., the probability density function of node attributes for the nodes of different sensitive attribute  $\mathbf{s} = \{-1, 1\}$  follows  $f_{\mathbf{X}}(\mathbf{X}_i = \mathbf{x} | \mathbf{s}_i = -1) \sim \mathcal{N}(\mu_1, \Sigma_1) \triangleq P_1$  and  $f_{\mathbf{X}}(\mathbf{X}_i = \mathbf{x} | \mathbf{s}_i = 1) \sim \mathcal{N}(\mu_2, \Sigma_2) \triangleq P_2$ , and the sensitive attribute ratio satisfies  $\mathbb{E}_i[\mathbb{P}(\mathbf{s}_i = 1)] = c$ , then the mutual information between sensitive attribute and node attributes  $I(\mathbf{s}, \mathbf{X})$  satisfies*

$$I(\mathbf{s}, \mathbf{X}) \leq -(1-c) \ln \left[ (1-c) + c \exp(-D_{KL}(P_1 || P_2)) \right] \\ - c \ln \left[ c + (1-c) \exp(-D_{KL}(P_2 || P_1)) \right] \triangleq I_{up}(\mathbf{s}, \mathbf{X}).$$

Based on Theorem 3.2, we can observe that lower distribution distance  $D_{KL}(P_1 || P_2)$  or  $D_{KL}(P_2 || P_1)$  is beneficial for reducing  $I_{up}(\mathbf{s}, \mathbf{X})$  and  $I(\mathbf{s}, \mathbf{X})$  since the sensitive attribute is less distinguishable based on node representations.

### 3.3 WHEN AND WHY AGGREGATION ENHANCES THE BIAS?

We focus on the role of message passing in terms of fairness. Suppose the graph adjacency matrix  $\mathbf{A}$  is sampled for  $(n, \rho_d, \epsilon_{sens}, c)$ -graph and we adopt the GCN-like message passing  $\tilde{\mathbf{X}} = \tilde{\mathbf{A}}\mathbf{X}$ , where  $\tilde{\mathbf{A}}$  is normalized adjacency matrix with self-loop. We define the bias difference for such message passing as  $\Delta I_{up} = I_{up}(\mathbf{s}, \tilde{\mathbf{X}}) - I_{up}(\mathbf{s}, \mathbf{X})$  to measure the role of graph topology. Before illustrating the main results, we first provide some related definitions based on  $(n, \rho_d, \epsilon_{sens}, c)$ -graph as follows:

**Definition 3.3 (intra-connect and inter-connect probability)** *For given  $(n, \rho_d, \epsilon_{sens}, c)$ -graph, the intra-connect probability is defined as the average connection probability given node pair with the same sensitive attributes, i.e.,  $p_{conn} = \mathbb{E}_{ij}[\mathbf{A}_{ij} | \mathbf{s}_i = \mathbf{s}_j]$ . Similarly, the inter-connect probability is defined as the average connection probability given node pair with different sensitive attributes, i.e.,  $q_{conn} = \mathbb{E}_{ij}[\mathbf{A}_{ij} | \mathbf{s}_i \mathbf{s}_j = -1]$ .*

**Definition 3.4 (Connection degree for demographic group)** *For given  $(n, \rho_d, \epsilon_{sens}, c)$ -graph, the average connection degree for the node with the same sensitive attribute is the same. Therefore, we define the connection degree for demographic group  $\zeta_l$  ( $l \in \{-1, 1\}$ ) as the average connection degree given sensitive attribute  $s = l$ , i.e.,  $\zeta_l \triangleq \sum_{j=0}^n \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_j = l]$ .*

We provide a sufficient condition on when graph topology enhances bias happens in Theorem 3.5.

**Theorem 3.5** *Suppose the synthetic node attribute  $\mathbf{X}$  is generated based on Gaussian Mixture Model  $GMM(c, \mu_1, \Sigma, \mu_2, \Sigma)$ , and adjacency matrix  $\mathbf{A}$  is generated from  $(n, \rho_d, \epsilon_{sens}, c)$ -graph.*

*If adopting GCN-like message passing  $\tilde{\mathbf{X}} = \tilde{\mathbf{A}}\mathbf{X}$ , representation bias will be enhanced, i.e.,  $\Delta I_{up} > 0$  if the bias-enhance condition holds:  $(\nu_{-1} - \nu_1)^2 \min\{\zeta_{-1}, \zeta_1\} > 1$ , where  $\nu_{-1}, \nu_1$  are given by*

$$\nu_{-1} = \frac{(n_{-1} - 1)p_{conn} + 1}{\zeta_{-1}}, \quad \nu_1 = \frac{n_{-1}q_{conn}}{\zeta_1}, \quad (1)$$

<sup>7</sup>Group fairness metrics, such as demographic parity and equal opportunities, can only adopt for classification problem. Here we consider node presentation bias and thus mutual information is a reasonable choice to measure the mutual dependence between node representation and sensitive attributes.

and connection degree for demographic group  $\zeta_{-1}, \zeta_1$  are given by Lemma D.1.

Based on Theorem 3.5, we can provide more discussion on the influence of 4 graph data information of graph: the number of nodes  $n$ , the edge density  $\rho_d$ , sensitive homophily coefficient  $\epsilon_{sens}$ , and sensitive group ratio  $c$  for bias enhancement as follows:

- **Node representation bias is enhanced by message passing for sufficiently large/small sensitive homophily  $\epsilon_{sens}$ .** According to Lemma D.1, for large sensitive homophily coefficient  $\epsilon_{sens} \rightarrow 1$ , the inter-connect probability  $q_{conn} \rightarrow 0$  and intra-connect probability  $p_{conn}$  keeps the maximal value. In this case, based on Theorem 3.5, it is easy to obtain that  $\nu_{-1} = 1, \nu_1 = 0$  and the distance for the mean aggregated node representation will keep the same, i.e.,  $\tilde{\mu}_1 - \tilde{\mu}_2 = (\nu_{-1} - \nu_1)(\mu_1 - \mu_2) = \mu_1 - \mu_2$ . In other words,  $\nu_{-1} - \nu_1 < 1$  represents the reduction coefficient of the distance between the mean node attributes of the two sensitive attributes groups. Additionally, the covariance will be diminished after aggregation since  $\zeta_{-1}$  and  $\zeta_1$  are strictly larger than 1. Therefore, for sufficient large sensitive homophily coefficient  $\epsilon_{sens}$ , the bias-enhancement condition  $(\nu_{-1} - \nu_1)^2 \min\{\zeta_{-1}, \zeta_1\} > 1$  holds. Similarly, when  $\epsilon_{sens} \rightarrow 0$ , the inter-connect probability  $q_{conn}$  keeps the maximal value and intra-connect probability  $p_{conn} \rightarrow 0$ . Hence, we can obtain  $\nu_{-1} - \nu_1 = 1 - \frac{1}{\zeta_{-1}} - \frac{1}{\zeta_1}$ . When  $\min\{\zeta_{-1}, \zeta_1\} = \min\{c, 1 - c\}nq_{conn} + 1 > 4$ , it is easy to check  $(\nu_{-1} - \nu_1)^2 \min\{\zeta_{-1}, \zeta_1\} > (1 - \frac{2}{\min\{\zeta_{-1}, \zeta_1\}})^2 \min\{\zeta_{-1}, \zeta_1\} > 1$ .
- **The bias enhancement implicitly depends on node representation geometric differentiation, including the distance between the mean node representation within the same sensitive attribute and the scale covariance matrix.** Theorem 3.2 implies that low mean representation distance and concentrated representation (low covariance matrix) lead to fair representation. However, GCN-like message passing renders the mean node representation distance reduction  $\nu_{-1} - \nu_1$  and concentrated for each sensitive attribute group, which is an "adversarial" effect for fairness and the mean distance and covariance reduction is controlled by sensitive homophily coefficient.
- **The bias is enlarged as node number  $n$  being increased.** For large node number  $n$ , the mean distance almost keeps constant since  $\nu_{-1} \approx \frac{(1-c)p_{conn}}{(1-c)p_{conn} + cq_{conn}}$ ,  $\nu_1 \approx \frac{(1-c)q_{conn}}{(1-c)q_{conn} + cp_{conn}}$ , and  $\zeta_{-1}, \zeta_1$  are almost proportional to node number  $n$ . Therefore, the bias-enhancement condition can be more easily satisfied, and  $\Delta I_{up}$  would be higher for large graph data. The intuition is that, given the graph density ratio, large graph data represents a higher average degree node. Hence, each aggregated node representation adopts more neighbor's information with the same sensitive attribute and thus leads to a lower covariance matrix value and higher representation bias.
- **The bias is enlarged as graph connection density  $\rho$  being increased.** Based on Lemma D.1, inter-connection probability and inter-connection probability are both proportional to graph connection density  $\rho_d$ . Therefore,  $\nu_{-1}$  and  $\nu_1$  almost keep constant and the distance of mean node representation is constant as well. As for the covariance matrix, message passing leads to more concentrated node representation since  $\zeta_{-1}$  and  $\zeta_1$  are larger for higher graph connection density  $\rho_d$ . The rationale is similar to the graph node number: given node number, higher graph connection density  $\rho_d$  means higher average node degree and each aggregated node representation adopts more neighbor's information with the same sensitive attribute.
- **More balanced sensitive attribute group, the larger bias the aggregation leads.** Based on Lemma D.1, given graph connection density  $\rho_d$  and graph node number  $n$ , the intra-connection probability  $p_{conn}$  is higher while inter-connection probability  $q_{conn}$  is lower for more balanced sensitive attribute group. In other words, intra-connection probability  $p_{conn}$  (inter-connection probability  $q_{conn}$ ) monotonically decreases (increases) with respect to  $|c - \frac{1}{2}|$ .

### 3.4 MORE DISCUSSIONS

we highlight the novelty of this paper as follows:

- **Research problem:** We aim to theoretically investigate when and why GNNs enhance bias from a data-centric perspective. In contrast, (Ma et al., 2022) aims to characterize the "necessary" graph data property (particularly in heterophily) for GNNs, while (Baranwal et al., 2021) mainly focuses on out-of-distribution data. Specifically, our work mainly focuses on bias enhancement, which is defined as that the bias after aggregation is larger than that before aggregation, while (Ma et al., 2022; Baranwal et al., 2021) mainly focuses on the "linear separability" of the node presentation after aggregation. In other words, the research problems are significantly different.

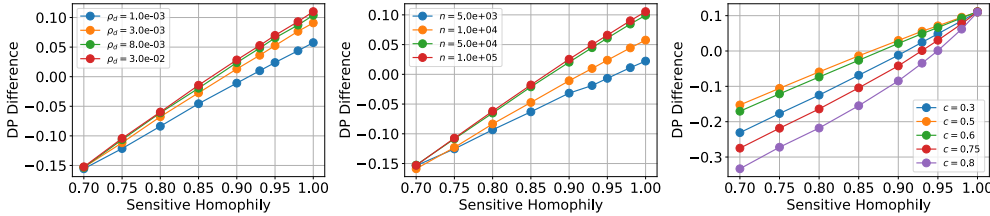


Figure 1: **Left:** DP difference for different connection density  $\rho_d$  with sensitive attribute ratio  $c = 0.5$  and number of nodes  $n = 10^4$ ; **Middle:** DP difference for different number of nodes  $n$  with sensitive attribute ratio  $c = 0.5$  and connection density  $\rho_d = 10^{-3}$ ; **Right:** DP difference for different sensitive attribute ratio  $c$  with connection density  $\rho_d = 10^{-3}$  and number of nodes  $n = 10^4$ ;

- **Technique:** We provide a sufficient condition to characterize when bias enhancement happens based on information theory. Conversely, (Ma et al., 2022; Baranwal et al., 2021) mainly relies on probability analysis (particularly involving concentration inequality) and the main theorems (such as Theorem.1 in (Ma et al., 2022), and Theorem.1 in (Baranwal et al., 2021)) hold with probability.

#### 4 LEARNING FAIR GRAPH TOPOLOGY

The above section provides a data-centric perspective (graph topology) to understand why message passing may enhance representation bias. However, searching for the optimal graph topology for GNNs is a non-trivial problem due to large-scale discrete optimization. In this section, motivated by the theoretical analysis, we propose **Fair Graph Rewiring** method, named **FairGR**<sup>8</sup>, to achieve fair prediction via learning the optimal graph topology. Specifically, we formulate the objective functions into three parts, including low sensitive homophily coefficient, high label homophily coefficient, and small topology perturbation. In this way, FairGR can explicitly reduce the topology bias while preserving useful topology information for prediction<sup>9</sup>.

**Objective loss design.** Considering binary sensitive attribute  $s_i \in \{-1, 1\}$  and binary classification task with label  $y_i \in \{-1, 1\}$  for  $i$ -th node, we aim to modify the graph topology to achieve low sensitive homophily coefficient  $\epsilon_{sens}$ , high label homophily coefficient  $\epsilon_{label}$ , and small topology perturbation. Note that only the sensitive attributes of training nodes are available during graph rewiring, we only modify the topology of the sub-graph used for training, while keeping the remaining edges intact. Therefore, we leave the subscript for choosing training nodes in the following loss formulation. For sensitive homophily coefficient with the binary sensitive attribute, the determination of any two nodes with the same sensitive attribute can be obtained via  $H(ss^T)$ , where  $H(\cdot)$  is Heaviside (unit) step function, sensitive attribute vector  $s \in \{-1, 1\}^{n \times 1}$ . Based on Definition 2.1, It is easy to rewrite sensitive homophily coefficient as  $\epsilon_{sens} = \frac{\|H(ss^T) \odot \mathbf{A}\|_1}{\|\mathbf{A}\|_1}$ , where  $\|\cdot\|_1$  denotes the entry-wise  $l_1$  norms (i.e., the summation over all absolute value of elements). Similarly, label homophily coefficient as  $\epsilon_{label} = \frac{\|H(yy^T) \odot \mathbf{A}\|_1}{\|\mathbf{A}\|_1}$ . As for graph topology perturbation, we can use the entry-wise  $l_1$  norms of the difference as the measurement. In a nutshell, the objective function to rewire graph connections can be formulated as:

$$\mathcal{L}(\hat{\mathbf{A}}|\mathbf{s}, \mathbf{y}, \mathbf{A}) = \max\{\epsilon_{sens}, 1 - \epsilon_{sens}\} - \alpha\epsilon_{label} + \beta\|\hat{\mathbf{A}} - \mathbf{A}\|_1, \tag{2}$$

Therefore, the rewired graph topology can be obtained via a constrained optimization problem

$$\min_{\hat{\mathbf{A}}} \mathcal{L}(\hat{\mathbf{A}}|\mathbf{s}, \mathbf{y}, \mathbf{A}) \quad s.t. \hat{\mathbf{A}}_{ij} \in \{0, 1\}, \tag{3}$$

where  $\alpha$  and  $\beta$  are the hyperparameters for label homophily coefficient and graph topology perturbation. Considering the formulated problem is a large-scale discrete optimization problem, we employ a heuristic optimization method to obtain the modified graph topology.

<sup>8</sup>Graph topology rewiring inevitably diminishes useful topology information and may decrease the performance. We provide more discussion on the follow-up works in Appendix L.

<sup>9</sup>The evaluation and computation complexity analysis can be found in Appendix G

Table 1: The performance on Node Classification (GR represents graph topology rewire). The boldface indicates that the prediction bias of our GR methods is lower than that without GR and MLP.

Models	Pokec-z			Pokec-n			NBA		
	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$
MLP	70.48 $\pm$ 0.77	1.61 $\pm$ 1.29	2.22 $\pm$ 1.01	72.48 $\pm$ 0.26	1.53 $\pm$ 0.89	3.39 $\pm$ 2.37	65.56 $\pm$ 1.62	22.37 $\pm$ 1.87	18.00 $\pm$ 3.52
GAT	69.76 $\pm$ 1.30	2.39 $\pm$ 0.62	2.91 $\pm$ 0.97	71.00 $\pm$ 0.48	3.71 $\pm$ 2.15	7.50 $\pm$ 2.88	57.78 $\pm$ 10.65	20.12 $\pm$ 16.18	13.00 $\pm$ 13.37
GAT-GR	56.75 $\pm$ 6.32	<b>1.04</b> $\pm$ 0.80	<b>1.14</b> $\pm$ 1.02	61.27 $\pm$ 9.34	<b>0.54</b> $\pm$ 0.51	<b>2.27</b> $\pm$ 1.55	53.65 $\pm$ 10.31	<b>4.16</b> $\pm$ 5.13	<b>3.67</b> $\pm$ 3.23
GCN	71.78 $\pm$ 0.37	3.25 $\pm$ 2.35	2.36 $\pm$ 2.09	73.09 $\pm$ 0.28	3.48 $\pm$ 0.47	5.16 $\pm$ 1.38	61.90 $\pm$ 1.00	23.70 $\pm$ 2.74	17.50 $\pm$ 2.63
GCN-GR	71.68 $\pm$ 0.58	1.94 $\pm$ 1.59	<b>1.27</b> $\pm$ 0.71	72.68 $\pm$ 0.44	<b>0.47</b> $\pm$ 0.39	<b>0.82</b> $\pm$ 0.78	61.59 $\pm$ 1.85	<b>20.24</b> $\pm$ 4.41	<b>9.50</b> $\pm$ 2.77
SGC	71.24 $\pm$ 0.46	4.81 $\pm$ 0.30	4.79 $\pm$ 2.27	71.46 $\pm$ 0.41	2.22 $\pm$ 0.29	3.85 $\pm$ 1.63	63.17 $\pm$ 0.63	22.56 $\pm$ 3.94	14.33 $\pm$ 2.16
SGC-GR	70.95 $\pm$ 0.91	3.32 $\pm$ 1.31	3.20 $\pm$ 1.90	71.91 $\pm$ 0.52	<b>0.71</b> $\pm$ 0.65	<b>2.39</b> $\pm$ 0.69	62.54 $\pm$ 1.62	<b>18.56</b> $\pm$ 2.81	<b>2.50</b> $\pm$ 1.66

**Optimization Strategy.** We first treat optimized variables as continuous variables to make Equation 3 differentiable and then adopt Proximal Gradient Descent (PGD) to optimize the formulated problem with constraint. Specifically, we first update the graph topology using gradient descent with  $\frac{\partial \mathcal{L}(\hat{\mathbf{A}}|s, y, \mathbf{A})}{\partial \hat{\mathbf{A}}}$ , then clip the graph topology  $\hat{\mathbf{A}}$  within  $\{0, 1\}$  in the projection operation of PGD. Such an operation is conducted multiple times to obtain the final graph topology. In practice, we only modify the graph topology within the training nodes. In other words, the three objectives, including sensitive homophily coefficient, label homophily coefficient, and graph topology perturbation, are calculated for the subgraph of training nodes. In this way, we can avoid information leakage from the test set and reduce the complexity of the optimization problem.

**Evaluation and Computation Complexity Analysis.** For algorithmic evaluation of pre-processing FairGR, we compared the prediction performance (including accuracy and fairness) using original graph topology and rewired graph topology across multiple GNN backbones. Denote the number of training nodes and update iterations as  $N$  and  $T$ , respectively. Then the computation complexity for gradient computation and projection of PGD are both  $O(n_{train}^2)$ . Therefore, the total computation complexity to obtain the final rewired graph topology is given by  $O(Tn_{train}^2)$ . The memory consumption is  $O(n_{train}^2)$  due to the storage of graph topology gradient.

## 5 EXPERIMENTS

In this section, we conduct experiments to validate the effectiveness (see Appendix K for more details.) of the proposed FairGR. For the node classification task, we adopt accuracy and two group fairness metrics, including *demographic parity* (DP) and *equal opportunity* (EO) to evaluate the classification performance and prediction bias (Louizos et al., 2015; Beutel et al., 2017). We firstly validate that GCN-like aggregation enhances representation bias for the graph data with large sensitive homophily via synthetic experiments. For real-world datasets, we conduct experiments to show that the prediction bias of GNN is larger than that of MLP. Moreover, we introduce the experimental settings and then evaluate our proposed FairGR compared with several baselines in terms of prediction performance and fairness metrics on real-world datasets.

### 5.1 SYNTHETIC EXPERIMENTS

In the synthetic experiments, we demonstrate the relation between DP difference across GCN-like aggregation operation and sensitive homophily coefficient. Specifically, we investigate the influence of graph node number  $n$ , graph connection density  $\rho_d$ , sensitive homophily  $\epsilon_{sens}$ , and sensitive attribute ration  $c$  for bias enhancement of GCN-like aggregation. **For evaluation metric**, we adopt the demographic parity (DP) difference during aggregation to measure the bias enhancement. **For node attribute generation**, we first generate node attribute with Gaussian distribution  $\mathcal{N}(\mu_1, \Sigma)$  and  $\mathcal{N}(\mu_2, \Sigma)$  for node with binary sensitive attribute, respectively, where  $\mu_1 = [0, 1]$ ,  $\mu_2 = [1, 0]$  and  $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ . **For adjacency matrix generation**, we randomly generate edges via a stochastic block model based on the intra-connection and inter-connection probability.

Figure 1 shows the DP difference during aggregation with respect to the sensitive homophily coefficient. We observe that a higher sensitive homophily coefficient generally leads to larger bias enhancement. Additionally, higher graph connection density  $\rho_d$ , larger node number  $n$ , and bal-

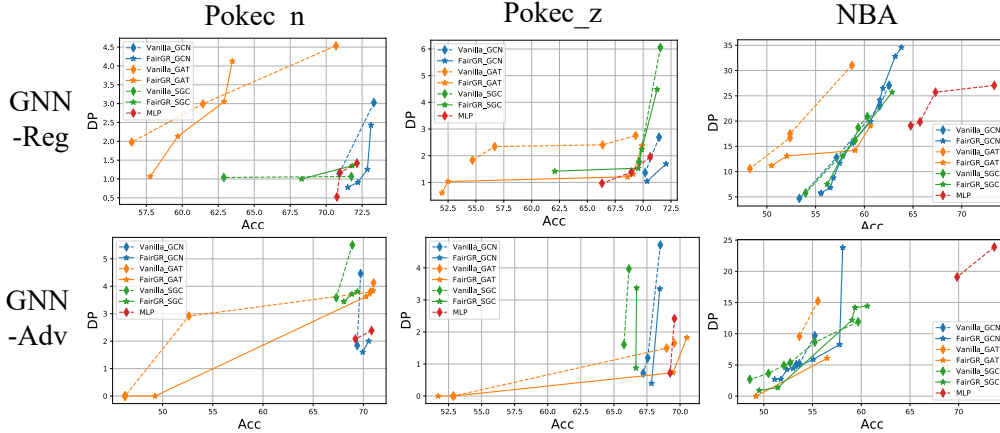


Figure 2: DP and Acc trade-off performance on three real-world datasets compared with adding regularization (Top) and adversarial debiasing (Bottom). The trade-off curve closer to the right bottom corner represents better trade-off performance.

anced sensitive attribute ratio  $c$  correspond to higher bias enhancement, which is consistent with our theoretical analysis in Theorem 3.5.

## 5.2 EXPERIMENTS ON REAL-WORLD DATASETS

We conduct experiments on three real-world datasets, including Pokec-z, Pokec-n, and NBA (Dai and Wang, 2021). We also adopt many representative GNNs, such as GCN (Kipf and Welling, 2017), GAT (Veličković et al., 2018), SGC (Wu et al., 2019). More details are in Appendix K.2.1.

**Baselines.** Considering that the proposed FairGR is a pre-processing method, we show that our proposed FairGR can improve many representative GNNs, such as GCN (Kipf and Welling, 2017), GAT (Veličković et al., 2018), SGC (Wu et al., 2019), in many fairness training strategies, such adding regularization, adversarial debiasing. For all models, we train 2 layers of neural networks with 64 hidden units for 300 epochs.

**Implementation Details.** For each experiment, we run 5 times and report the average performance for each method. We adopt Adam optimizer with 0.001 learning rate and  $1e^{-5}$  weight decay for all models training. In the adversarial debiasing setting, we train the classifier and adversary head with 70 and 30 epochs, respectively. The hyperparameters for adversarial debiasing are tuned in  $\{0.0, 0.5, 1.0, 2.0, 5.0, 8.0, 10.0, 50.0, 100.0\}$ . For adding regularization, we adopt the hyperparameter set  $\{0.0, 1.0, 1.5, 2.0, 5.0, 8.0, 10.0, 15.0, 25.0, 50.0, 80.0, 100.0\}$ . For FairGR, we first select the hyperparameters  $\alpha$  and  $\beta$  both in  $\{0.1, 0.5, 1.0, 5.0, 10.0\}$  to rewire graph topology with lowest demographic parity. Subsequently, the rewired graph can be adopted with several in-processing methods to further improve prediction and fairness tradeoff performance.

### 5.2.1 DOES GNNs HAVE A LARGER PREDICTION BIAS THAN MLP?

To validate the effect of bias enhancement of GNNs, we compare the performance of many representative GNNs over MLP on various real-world datasets in Table 1. We observe:

- Many representative GNNs have a higher prediction bias compared with MLP model on all three datasets in terms of demographic parity and equal opportunities. For demographic parity, the prediction bias of MLP is lower than that of GAT, GCN, and SGC by 32.64%, 50.46%, 66.53% and 58.72% on Pokec-z dataset. The higher prediction bias comes from the aggregation within the same-sensitive attribute nodes and topology bias in graph data.
- FairGR can mitigate bias for GCN and SGC backbone via rewiring graph topology in these three datasets. For GAT backbone, although the bias can be mitigated, the accuracy drop is significant due to the fact that GAT is more sensitive to graph topology rewiring.



### 5.2.2 DOES FAIRGR ACHIEVE BETTER TRADEOFF PERFORMANCE IN VARIOUS SETTINGS?

#### Comparison with GNN regularization (GNN-Reg) and GNN adversarial debiasing (GNN-Adv).

To validate that the proposed FairGR is compatible with these two strategies, we also show the prediction performance and fairness metric trade-off compared with adversarial debiasing (Bose and Hamilton, 2019; Dai and Wang, 2021) and add demographic parity regularization (Chuang and Mroueh, 2020). In adversarial debiasing, the output of GNNs is the input of the adversary, and the goal of the adversary is to predict the node sensitive attribute. For these two fair training strategies, we adopt GCN, GAT, and SGC as backbones. We randomly split 50%/25%/25% for training, validation, and test dataset. Figure 2 shows the Pareto optimality curve for all methods via adopting different hyperparameters in adding regularization and adversarial bias, where the right-bottom corner point represents the ideal performance (100% accuracy and 0% prediction bias). We observe that:

- For both adversarial debiasing and adding regularization training strategies, our proposed FairGR can achieve a better DP-Acc trade-off compared with that without any graph data rewiring for many GNNs. In other words, FairGR can effectively reduce training bias and is compatible with many existing fairness training strategies.
- Topology does matter in GNNs. For adding regularization or adversarial debiasing, FairGS embrace different tradeoff performance gain on top of different GNNs. Such observation implies that there is a complicated interaction between graph topology and aggregation algorithms. Additionally, FairGS provide the most tradeoff performance benefit in GAT compared with GCN and SGC. The high capacity of GAT may energize the aggregation algorithm to learn from data. Therefore, the tradeoff performance improvement is the highest in adding regularization and adversarial debiasing.

## 6 RELATED WORK

We briefly review the existing work on **graph neural networks** and **fairness-aware learning on graphs**. (Please refer to Appendix H for a more comprehensive discussion.). Existing GNNs can be roughly divided into spectral-based and spatial-based GNNs. Spectral-based GNNs provide graph convolution definition based on graph theory (Bruna et al., 2013). Spatial-based GNNs variants are popular due to explicit neighbors’ information aggregation, including Graph convolutional networks (GCN) (Kipf and Welling, 2017), graph attention network (GAT) (Veličković et al., 2018). As for fairness in graph data, many work have been developed to achieve fairness in machine learning community, including fair walk (Rahman et al., 2019), adversarial debiasing (Dai and Wang, 2021), Bayesian approach (Buyl and De Bie, 2020), and contrastive learning (Agarwal et al., 2021). Much literature empirically shows that GNN or aggregation may enhance prediction bias compared with MLP. However, the theoretical understanding of why such a phenomenon happens is still unclear. In this work, we take an initial step toward theoretically understanding why aggregation enhances bias from a data perspective. Based on this understanding, we develop a simple yet effective fair graph rewiring method to achieve better tradeoff performance. More importantly, the proposed FairGR is compatible with many fair training strategies.

## 7 CONCLUSION

In this work, we theoretically analyze when and why bias enhancement happens in GNNs, serving as a theoretical pilot study. We provide a sufficient condition provably and comprehensively delineates instances of bias enhancement for GCN-like aggregation under CSBM, and then derive several data-centric insights. Motivated by these insights, we develop a simple yet effective graph rewiring method, named FairGR, to reduce the sensitive homophily while preserving useful information. Experimental results on real-world datasets demonstrate the effectiveness of FairGR in many fair training strategies and GNNs backbones in node classification tasks.

## REFERENCES

Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. Towards a unified framework for fair and stable graph representation learning. *arXiv preprint arXiv:2102.13186*, 2021.

- Aseem Baranwal, Kimon Fountoulakis, and Aukosh Jagannath. Graph convolution for semi-supervised classification: Improved linear separability and out-of-distribution generalization. In *ICML*, 2021.
- Aseem Baranwal, Kimon Fountoulakis, and Aukosh Jagannath. Effects of graph convolutions in multi-layer networks. In *International Conference on Learning Representations*, 2023.
- Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *California law review*, pages 671–732, 2016.
- Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. Data decisions and theoretical implications when adversarially learning fair representations. *arXiv preprint arXiv:1707.00075*, 2017.
- Avishek Bose and William Hamilton. Compositional fairness constraints for graph embeddings. In *International Conference on Machine Learning*, pages 715–724. PMLR, 2019.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- Maarten Buyt and Tijl De Bie. Debayes: a bayesian method for debiasing network embeddings. In *International Conference on Machine Learning*, pages 1220–1229. PMLR, 2020.
- Mathieu Carriere, Frédéric Chazal, Marc Glisse, Yuichi Ike, Hariprasad Kannan, and Yuhei Umeda. Optimizing persistent homology based functions. In *International conference on machine learning*, pages 1294–1303. PMLR, 2021.
- Eli Chien, Jianhao Peng, Pan Li, and Olga Milenkovic. Adaptive universal generalized pagerank graph neural network. In *International Conference on Learning Representations*, 2021.
- Ching-Yao Chuang and Youssef Mroueh. Fair mixup: Fairness via interpolation. In *International Conference on Learning Representations*, 2020.
- Ching-Yao Chuang and Youssef Mroueh. Fair mixup: Fairness via interpolation. In *International Conference on Learning Representations*, 2021.
- Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In *International conference on machine learning*, pages 1436–1445. PMLR, 2019.
- Enyan Dai and Suhang Wang. Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 680–688, 2021.
- Yushun Dong, Jian Kang, Hanghang Tong, and Jundong Li. Individual fairness for graph neural networks: A ranking based approach. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 300–310, 2021.
- Yushun Dong, Ninghao Liu, Brian Jalaian, and Jundong Li. Edits: Modeling and mitigating data bias for graph neural networks. In *Proceedings of the ACM Web Conference 2022*, pages 1259–1269, 2022a.
- Yushun Dong, Song Wang, Jing Ma, Ninghao Liu, and Jundong Li. Interpreting unfairness in graph neural networks via training node attribution. *arXiv preprint arXiv:2211.14383*, 2022b.
- Yushun Dong, Binchi Zhang, Yiling Yuan, Na Zou, Qi Wang, and Jundong Li. Reliant: Fair knowledge distillation for graph neural networks. *arXiv preprint arXiv:2301.01150*, 2023.
- Mengnan Du, Subhabrata Mukherjee, Guanchu Wang, Ruixiang Tang, Ahmed Hassan Awadallah, and Xia Hu. Fairness via representation neutralization. *arXiv preprint arXiv:2106.12674*, 2021.
- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 259–268, 2015.

- Joseph Fisher, Arpit Mittal, Dave Palfrey, and Christos Christodoulopoulos. Debiasing knowledge graph embeddings. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7332–7345, 2020.
- Santo Fortunato and Darko Hric. Community detection in networks: A user guide. *Physics reports*, 659:1–44, 2016.
- Hongyang Gao, Zhengyang Wang, and Shuiwang Ji. Large-scale learnable graph convolutional networks. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1416–1424, 2018.
- Takuo Hamaguchi, Hidekazu Oiwa, Masashi Shimbo, and Yuji Matsumoto. Knowledge transfer for out-of-knowledge-base entities: A graph neural network approach. *arXiv preprint arXiv:1706.05674*, 2017.
- William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1025–1035, 2017.
- Xiaotian Han, Zhimeng Jiang, Ninghao Liu, and Xia Hu. G-mixup: Graph data augmentation for graph classification. In *International Conference on Machine Learning*, pages 1–9. PMLR, 2022a.
- Xiaotian Han, Zhimeng Jiang, Ninghao Liu, Qingquan Song, Jundong Li, and Xia Hu. Geometric graph representation learning via maximizing rate reduction. In *Proceedings of the ACM Web Conference 2022*, pages 1226–1237, 2022b.
- Xiaotian Han, Zhimeng Jiang, Hongye Jin, Zirui Liu, Na Zou, Qifan Wang, and Xia Hu. Retiring  $\Delta$ DP: New distribution-level metrics for demographic parity. *arXiv preprint arXiv:2301.13443*, 2023.
- Zhimeng Jiang, Xiaotian Han, Chao Fan, Fan Yang, Ali Mostafavi, and Xia Hu. Generalized demographic parity for group fairness. In *International Conference on Learning Representations*, 2022.
- Zhimeng Jiang, Xiaotian Han, Hongye Jin, Guanchu Wang, Na Zou, and Xia Hu. Weight perturbation can help fairness under distribution shift. *arXiv preprint arXiv:2303.03300*, 2023.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- Jon Kleinberg and Eva Tardos. *Algorithm design*. Pearson Education India, 2006.
- Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations*, 2019.
- Öykü Deniz Köse and Yanning Shen. Fairness-aware node representation learning. *arXiv preprint arXiv:2106.05391*, 2021.
- Solomon Kullback. *Information theory and statistics*. Courier Corporation, 1997.
- Charlotte Laclau, Ievgen Redko, Manvi Choudhary, and Christine Largeron. All of the fairness for edge prediction with optimal transport. In *International Conference on Artificial Intelligence and Statistics*, pages 1774–1782. PMLR, 2021.
- Peizhao Li, Yifei Wang, Han Zhao, Pengyu Hong, and Hongfu Liu. On dyadic fairness: Exploring and mitigating bias in graph connections. In *International Conference on Learning Representations*, 2021.
- Xiao Lin, Jian Kang, Weilin Cong, and Hanghang Tong. Bemap: Balanced message passing for fair graph neural network. *arXiv preprint arXiv:2306.04107*, 2023.
- Hongyi Ling, Zhimeng Jiang, Youzhi Luo, Shuiwang Ji, and Na Zou. Learning fair graph representations via automated data augmentations. In *International Conference on Learning Representations*, 2023.

- Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. Is homophily a necessity for graph neural networks? In *International Conference on Learning Representations*, 2022.
- Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5115–5124, 2017.
- Hoang Nt and Takanori Maehara. Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550*, 2019.
- Tahleen Rahman, Bartłomiej Surma, Michael Backes, and Yang Zhang. Fairwalk: Towards fair graph embedding. 2019.
- Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In *International scientific conference and international workshop present day trends of innovations*, volume 1, 2012.
- Xianfeng Tang, Huaxiu Yao, Yiwei Sun, Yiqi Wang, Jiliang Tang, Charu Aggarwal, Prasenjit Mitra, and Suhang Wang. Investigating and mitigating degree-related biases in graph convolutional networks. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 1435–1444, 2020.
- Remco Van Der Hofstad. *Random graphs and complex networks*, volume 43. Cambridge university press, 2016.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 974–983, 2018.
- Mikhail Yurochkin and Yuekai Sun. Sensei: Sensitive set invariance for enforcing individual fairness. In *International Conference on Learning Representations*, 2020.
- Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 335–340, 2018.
- Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond homophily in graph neural networks: Current limitations and effective designs. *Advances in neural information processing systems*, 2020.

## A PROOF OF THEOREM 3.2

We provide a more general proof for categorical sensitive attribute  $\mathbf{s} \in \{1, 2, \dots, K\}$  and the prior probability is given by  $\mathbb{P}(\mathbf{s} = i) = c_i$ . Suppose the conditional node attribute  $\mathbf{x}$  distribution given node sensitive attribute  $\mathbf{s} = i$  satisfies normal distribution  $P_i(\mathbf{x}) \triangleq \mathcal{N}(\mu_i, \Sigma_i)$ , the distribution of node sensitive attribute is the mixed Gaussian distribution  $f(\mathbf{x}) = \sum_{i=1}^K c_i P_i(\mathbf{x})$ . Based on the definition of mutual information, we have

$$I(\mathbf{s}, \mathbf{x}) = H(\mathbf{x}) - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i); \quad (4)$$

where  $H(\cdot)$  represents Shannon entropy for random variable. Subsequently, we focus on the entropy of the mixed Gaussian distribution  $H(\mathbf{x})$ . We show that such entropy can be upper bounded by the pairwise Kullback-Leibler (KL) divergence as follows:

$$\begin{aligned} I(\mathbf{s}, \mathbf{x}) &= - \sum_{i=1}^K c_i \mathbb{E}_{P_i} \left[ \ln \sum_{j=1}^K c_j P_j(\mathbf{x}) \right] - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i) \\ &\stackrel{(a)}{\leq} - \sum_{i=1}^K c_i \left[ \ln \sum_{j=1}^K c_j e^{\mathbb{E}_{P_i}[\ln P_j(\mathbf{x})]} \right] - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i) \\ &= - \sum_{i=1}^K c_i \left[ \ln \sum_{j=1}^K c_j e^{-H(P_i||P_j)} \right] - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i) \\ &\stackrel{(b)}{=} - \sum_{i=1}^K c_i \left[ \ln \sum_{j=1}^K c_j e^{-H(P_i) - D_{KL}(P_i||P_j)} \right] - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i) \\ &= - \sum_{i=1}^K c_i \left[ \ln \sum_{j=1}^K c_j e^{-D_{KL}(P_i||P_j)} \right] + \sum_{i=1}^K c_i H(P_i) - \sum_{i=1}^K c_i H(\mathbf{x}|\mathbf{s} = i), \\ &= - \sum_{i=1}^K c_i \left[ \ln \sum_{j=1}^K c_j e^{-D_{KL}(P_i||P_j)} \right], \end{aligned}$$

where KL divergence  $D_{KL}(P_i||P_j) \triangleq \int P_i(\mathbf{x}) \ln \frac{P_i(\mathbf{x})}{P_j(\mathbf{x})} d\mathbf{x}$  and cross entropy  $H(P_i||P_j) = - \int P_i(\mathbf{x}) \ln P_j(\mathbf{x}) d\mathbf{x}$ . The inequality (a) holds based on the variational lower bound on the expectation of a log-sum inequality  $\mathbb{E}[\ln \sum_i Z_i] \geq \ln [\sum_i e^{\mathbb{E}[\ln Z_i]}]$  (Kullback, 1997), and quality (2) holds based on  $H(P_i||P_j) = H(P_i) + D_{KL}(P_i||P_j)$ . As a special case for binary sensitive attribute, it is easy to obtain the following results:

$$I(\mathbf{s}, \mathbf{X}) \leq -(1-c) \ln \left[ (1-c) + c \exp(-D_{KL}(P_1||P_2)) \right] - c \ln \left[ c + (1-c) \exp(-D_{KL}(P_2||P_1)) \right].$$

## B PROOF OF LEMMA D.1

**Lemma B.1** Suppose the synthetic graph is generated from  $(n, \rho_d, \epsilon_{sens}, c)$ -graph, then we obtain the intra-connect and inter-connect probability as follows:

$$p_{conn} = \frac{\rho_d \epsilon_{sens}}{c^2 + (1-c)^2}, \quad q_{conn} = \frac{\rho_d(1-\epsilon_{sens})}{2c(1-c)}. \quad (5)$$

the connection degree for the demographic group (including the self-loop degree) is given by

$$\begin{aligned} \zeta_{-1} &= n_1 q_{conn} + (n_{-1} - 1) p_{conn} + 1, \\ \zeta_1 &= n_{-1} q_{conn} + (n_1 - 1) p_{conn} + 1, \end{aligned} \quad (6)$$

where the node number with for the demographic group  $s = -1$  and  $s = 1$  are given by  $n_{-1} = n(1-c)$ , and  $n_1 = nc$ , respectively.

**Proof:** Based on Bayes' rule, we have the intra-connect and inter-connect probability as follows

$$\begin{aligned} p_{conn} &= \mathbb{P}(\mathbf{A}_{ij} = 1 | \mathbf{s}_i = \mathbf{s}_j) = \frac{\mathbb{P}(\mathbf{A}_{ij} = 1)\mathbb{P}(\mathbf{s}_i = \mathbf{s}_j | \mathbf{A}_{ij} = 1)}{\mathbb{P}(\mathbf{s}_i = \mathbf{s}_j)} = \frac{\rho_d \epsilon_{sens}}{c^2 + (1-c)^2}, \\ q_{conn} &= \mathbb{P}(\mathbf{A}_{ij} = 1 | \mathbf{s}_i \mathbf{s}_j = -1) = \frac{\mathbb{P}(\mathbf{A}_{ij} = 1)\mathbb{P}(\mathbf{s}_i \mathbf{s}_j = -1 | \mathbf{A}_{ij} = 1)}{\mathbb{P}(\mathbf{s}_i \mathbf{s}_j = -1)} = \frac{\rho_d(1 - \epsilon_{sens})}{2c(1-c)}. \end{aligned} \quad (7)$$

It is easy to see that the node number for the demographic group  $s = -1$  and  $s = 1$  are given by  $n_{-1} = n(1-c)$ , and  $n_1 = nc$ , respectively. Consider there always exists a self-loop for each node, connection degree given sensitive attribute  $s = -1$  is given by

$$\begin{aligned} \zeta_{-1} &\triangleq \sum_{j=0}^n \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_j = -1] = 1 + \sum_{j=0, j \neq i}^n \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_j = -1] \\ &= 1 + (n_{-1} - 1)\mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_i = s_j = -1] + n_1 \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_i = 1, s_j = -1] \\ &= n_1 q_{conn} + (n_{-1} - 1)p_{conn} + 1, \end{aligned} \quad (8)$$

Similarly, we have

$$\begin{aligned} \zeta_1 &\triangleq \sum_{j=0}^n \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_j = 1] = 1 + \sum_{j=0, j \neq i}^n \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_j = 1] \\ &= 1 + (n_1 - 1)\mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_i = s_j = 1] + n_{-1} \mathbb{E}_{\mathbf{A}_{ij}}[\mathbf{A}_{ij} | s_i = -1, s_j = 1] \\ &= n_{-1} q_{conn} + (n_1 - 1)p_{conn} + 1, \end{aligned} \quad (9)$$

□

## C PROOF OF THEOREM 3.5

Before going deeper for our proof, we first introduce two useful lemmas on KL divergence and statistical information of graph.

**Lemma C.1** For two  $d$ -dimensional Gaussian distributions  $P = \mathcal{N}(\mu_p, \Sigma_p)$  and  $Q = \mathcal{N}(\mu_q, \Sigma_q)$ , the KL divergence  $D_{KL}(P||Q)$  is given by

$$D_{KL}(P||Q) = \frac{1}{2} \left[ \ln \frac{|\Sigma_q|}{|\Sigma_p|} - d + (\mu_p - \mu_q)^\top \Sigma_q^{-1} (\mu_p - \mu_q) + \text{Tr}(\Sigma_q^{-1} \Sigma_p) \right] \quad (10)$$

where  $\top$  is matrix transpose operation and  $\text{Tr}(\cdot)$  is trace of a square matrix.

**Proof:** Note that the probability density function of multivariate Normal distribution is given by:

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma_p|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right),$$

the KL divergence between distributions  $P$  and  $Q$  can be given by

$$\begin{aligned} D_{KL}(P||Q) &= \mathbb{E}_P[\ln(P) - \ln(Q)] \\ &= \mathbb{E}_P \left[ \frac{1}{2} \ln \frac{|\Sigma_q|}{|\Sigma_p|} - \frac{1}{2} (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) + \frac{1}{2} (\mathbf{x} - \mu_q)^\top \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \ln \frac{|\Sigma_q|}{|\Sigma_p|} - \underbrace{\frac{1}{2} \mathbb{E}_P \left[ (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right]}_{I_1} + \underbrace{\frac{1}{2} \mathbb{E}_P \left[ (\mathbf{x} - \mu_q)^\top \Sigma_q^{-1} (\mathbf{x} - \mu_q) \right]}_{I_2}. \end{aligned}$$

Using the commutative property of the trace operation, we have

$$\begin{aligned} I_1 &= \frac{1}{2} \mathbb{E}_P \left[ (\mathbf{x} - \mu_p)^\top \Sigma_p^{-1} (\mathbf{x} - \mu_p) \right] = \frac{1}{2} \text{Tr} \left( \mathbb{E}_P \left[ (\mathbf{x} - \mu_p)^\top (\mathbf{x} - \mu_p) \Sigma_p^{-1} \right] \right) \\ &= \frac{1}{2} \text{Tr} \left( \mathbb{E}_P \left[ (\mathbf{x} - \mu_p)^\top (\mathbf{x} - \mu_p) \right] \Sigma_p^{-1} \right) = \frac{1}{2} \text{Tr} \left( \Sigma_p \Sigma_p^{-1} \right) = \frac{d}{2}, \end{aligned} \quad (11)$$

As for the term  $I_2$ , note that  $\mathbf{x} - \mu_q = (\mathbf{x} - \mathbb{E}_P[\mathbf{x}]) + (\mathbb{E}_P[\mathbf{x}] - \mu_q)$ , we can obtain the following equation:

$$\begin{aligned} I_2 &= \frac{1}{2} \mathbb{E}_P \left[ (\mathbf{x} - \mu_q)^\top \boldsymbol{\Sigma}_q^{-1} (\mathbf{x} - \mu_q) \right] \\ &= \frac{1}{2} \left[ (\mu_p - \mu_q)^\top \boldsymbol{\Sigma}_q^{-1} (\mu_p - \mu_q) + \text{Tr}(\boldsymbol{\Sigma}_q^{-1} \boldsymbol{\Sigma}_p) \right], \end{aligned} \quad (12)$$

Therefore, the KL divergence  $D_{KL}(P||Q)$  is given by

$$D_{KL}(P||Q) = \frac{1}{2} \left[ \ln \frac{|\boldsymbol{\Sigma}_q|}{|\boldsymbol{\Sigma}_p|} - d + (\mu_p - \mu_q)^\top \boldsymbol{\Sigma}_q^{-1} (\mu_p - \mu_q) + \text{Tr}(\boldsymbol{\Sigma}_q^{-1} \boldsymbol{\Sigma}_p) \right]. \quad (13)$$

□

Note that the synthetic graph is generated from  $(n, \rho_d, \epsilon_{sens}, c)$ -graph. The sensitive attribute  $\mathbf{s}$  is generated to with ratio  $c$ , i.e., the number of node sensitive attribute  $\mathbf{s} = -1$  and  $\mathbf{s} = 1$  are  $n_{-1} = n(1 - c)$  and  $n_1 = nc$ . Based on the determined sensitive attribute  $\mathbf{s}$ , we randomly generate the edge based on parameters  $\rho_d$  and  $\epsilon_{sens}$  and Lemma D.1, i.e., the edges within and cross the same group are randomly generated based on intra-connect probability and inter-connect probability. Therefore, the adjacency matrix  $\mathbf{A}_{ij}$  is independent on node attribute  $\mathbf{X}_i$  and  $\mathbf{X}_j$  given sensitive attributes  $\mathbf{s}_i$  and  $\mathbf{s}_j$ , i.e.,  $\mathbf{A}_{ij} \perp\!\!\!\perp (\mathbf{X}_i, \mathbf{X}_j) | (\mathbf{s}_i, \mathbf{s}_j)$ . Similarly, the different node attributes and edges are also dependent on each other given sensitive attributes, i.e.,  $\mathbf{A}_{ij} \perp\!\!\!\perp \mathbf{A}_{ij} | (\mathbf{s}_i, \mathbf{s}_j, \mathbf{s}_k)$  and  $\mathbf{X}_i \perp\!\!\!\perp \mathbf{X}_j | (\mathbf{s}_i, \mathbf{s}_j)$ . Therefore, considering GCN-like aggregation  $\tilde{\mathbf{X}}_i = \sum_{j=1}^n \tilde{\mathbf{A}}_{ij} \mathbf{X}_j$ , we have the aggregated node attributes expectation given sensitive attribute as follows:

$$\begin{aligned} \tilde{\mu}_1 &= \mathbb{E}_{\tilde{\mathbf{X}}_i} [\tilde{\mathbf{X}}_i | \mathbf{s}_i = -1] = \sum_{j=1}^n \mathbb{E}_{\tilde{\mathbf{A}}_{ij}, \mathbf{X}_j} [\tilde{\mathbf{A}}_{ij} \mathbf{X}_j | \mathbf{s}_i = -1] \\ &= \sum_{j=1}^n \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = -1] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = -1] \\ &= (n_{-1} - 1) \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = -1, \mathbf{s}_j = -1] \mathbb{E}_{\mathbf{X}_j} [\tilde{\mathbf{X}}_j | \mathbf{s}_i = -1, \mathbf{s}_j = -1] \\ &\quad + \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = -1, i = j] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = -1] \\ &\quad + n_1 \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = -1, \mathbf{s}_j = 1] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = -1, \mathbf{s}_j = 1] \\ &= \frac{[(n_{-1} - 1)p_{conn} + 1]\mu_1 + n_1 q_{conn} \mu_2}{(n_{-1} - 1)p_{conn} + 1 + n_1 q_{conn}} \triangleq \nu_{-1} \mu_1 + (1 - \nu_{-1}) \mu_2. \end{aligned} \quad (14)$$

where  $\nu_{-1} = \frac{(n_{-1} - 1)p_{conn} + 1}{(n_{-1} - 1)p_{conn} + 1 + n_1 q_{conn}}$ . Similarly, for the node with sensitive attribute 1, we have

$$\begin{aligned} \tilde{\mu}_2 &= \mathbb{E}_{\tilde{\mathbf{X}}_i} [\tilde{\mathbf{X}}_i | \mathbf{s}_i = 1] = \sum_{j=1}^n \mathbb{E}_{\tilde{\mathbf{A}}_{ij}, \mathbf{X}_j} [\tilde{\mathbf{A}}_{ij} \mathbf{X}_j | \mathbf{s}_i = 1] \\ &= \sum_{j=1}^n \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = 1] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = 1] \\ &= n_{-1} \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = 1, \mathbf{s}_j = -1] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = 1, \mathbf{s}_j = -1] \\ &\quad + \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = 1, i = j] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = 1] \\ &\quad + (n_1 - 1) \mathbb{E}_{\tilde{\mathbf{A}}_{ij}} [\tilde{\mathbf{A}}_{ij} | \mathbf{s}_i = 1, \mathbf{s}_j = 1] \mathbb{E}_{\mathbf{X}_j} [\mathbf{X}_j | \mathbf{s}_i = 1, \mathbf{s}_j = 1] \\ &= \frac{n_{-1} q_{conn} \mu_1 + [(n_1 - 1)p_{conn} + 1]\mu_2}{n_{-1} q_{conn} + (n_1 - 1)p_{conn} + 1} \triangleq \nu_1 \mu_1 + (1 - \nu_1) \mu_2. \end{aligned} \quad (15)$$

where  $\nu_1 = \frac{n-1q_{conn}}{n-1q_{conn}+1+(n_1-1)p_{conn}}$ . As for the covariance matrix of aggregated node attributes  $\tilde{\mathbf{X}}$  given node sensitive attribute  $\mathbf{s} = -1$  and original sensitive attribute, note that we can obtain

$$\begin{aligned}\tilde{\Sigma}_1 &= \mathbb{D}_{\tilde{\mathbf{X}}_i}[\tilde{\mathbf{X}}_i|\mathbf{s}_i = -1] = \sum_{j=1}^n \mathbb{D}_{\tilde{\mathbf{A}}_{ij}, \mathbf{X}_j}[\tilde{\mathbf{A}}_{ij} \mathbf{X}_j|\mathbf{s}_i = -1] \\ &= \sum_{j=1}^n \mathbb{E}_{\tilde{\mathbf{A}}_{ij}}[\tilde{\mathbf{A}}_{ij}^2|\mathbf{s}_i = -1] \mathbf{D}_{\tilde{\mathbf{X}}_j}[\tilde{\mathbf{X}}_j|\mathbf{s}_i = -1] \\ &= \frac{(n-1)p_{conn}\Sigma + \Sigma + n_1q_{conn}\Sigma}{[(n-1)p_{conn} + 1 + n_1q_{conn}]^2} \\ &= \frac{\Sigma}{(n-1)p_{conn} + 1 + n_1q_{conn}} \triangleq \zeta_{-1}^{-1}\Sigma,\end{aligned}\quad (16)$$

where  $\zeta_{-1} = (n-1)p_{conn} + 1 + n_1q_{conn}$ . Similarly, given node sensitive attribute  $\mathbf{s} = 1$ , we have  $\tilde{\Sigma}_2 = \mathbb{D}_{\tilde{\mathbf{X}}_i}[\tilde{\mathbf{X}}_i|\mathbf{s}_i = 1] = \frac{\Sigma}{n-1q_{conn}+1+(n_1-1)p_{conn}} \triangleq \zeta_1^{-1}\Sigma$ , where  $\zeta_1 = n-1q_{conn} + 1 + (n_1 - 1)p_{conn}$ . In other words, the covariance matrix of the aggregated node attributes is lower than the original one since the ‘‘average’’ operation will make node representation more concentrated. Note that the summation over several Gaussian random variables is still Gaussian, we define the node attributes distribution for sensitive attribute  $\mathbf{s} = -1$  and  $\mathbf{s} = 1$  as  $P_1 = \mathcal{N}(\mu_1, \Sigma)$ ,  $P_2 = \mathcal{N}(\mu_2, \Sigma)$ , respectively. Similarly, the aggregated node representation distribution follows for sensitive attribute  $\mathbf{s} = -1$  and  $\mathbf{s} = 1$  as  $\tilde{P}_1 = \mathcal{N}(\tilde{\mu}_1, \tilde{\Sigma}_1)$ ,  $\tilde{P}_2 = \mathcal{N}(\tilde{\mu}_2, \tilde{\Sigma}_2)$ . Note that the sensitive attribute ratio keeps the same after the aggregation and larger KL divergence for these two sensitive attributes group distribution, the bias enhances  $\Delta Bias > 0$  if  $D_{KL}(\tilde{P}_1||\tilde{P}_2) > D_{KL}(P_1||P_2)$  and  $D_{KL}(\tilde{P}_2||\tilde{P}_1) > D_{KL}(P_2||P_1)$ . Therefore, we only focus on the KL divergence. According to Lemma C.1, it is easy to obtain KL divergence for original distribution as follows:

$$\begin{aligned}D_{KL}(P_1||P_2) &= \frac{1}{2} \left[ \ln \frac{|\Sigma|}{|\Sigma|} - d + (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2) + Tr(\Sigma^{-1}\Sigma) \right] \\ &= \frac{1}{2} (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2),\end{aligned}\quad (17)$$

As for KL divergence for aggregated distribution, similarly, we have

$$\begin{aligned}D_{KL}(\tilde{P}_1||\tilde{P}_2) &= \frac{1}{2} \left[ \ln \frac{|\tilde{\Sigma}_2|}{|\tilde{\Sigma}_1|} - d + (\tilde{\mu}_1 - \tilde{\mu}_2)^\top \tilde{\Sigma}_2^{-1}(\tilde{\mu}_1 - \tilde{\mu}_2) + Tr(\tilde{\Sigma}_2^{-1}\tilde{\Sigma}_1) \right] \\ &= \frac{1}{2} \left[ d \ln \frac{\zeta_{-1}}{\zeta_1} - d + (\nu_{-1} - \nu_1)^2 \zeta_1 (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2) + \frac{\zeta_1}{\zeta_{-1}} Tr(\mathbf{I}_d) \right] \\ &\stackrel{(c)}{\geq} \frac{1}{2} (\nu_{-1} - \nu_1)^2 \zeta_1 (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2),\end{aligned}\quad (18)$$

where inequality (c) holds since  $\ln x \leq x - 1$  for any  $x > 0$ . Compared with equations (17) and (18), it is seen that  $D_{KL}(\tilde{P}_1||\tilde{P}_2) > D_{KL}(P_1||P_2)$  if  $(\nu_{-1} - \nu_1)^2 \zeta_1 > 1$ . Similarly, we can have  $D_{KL}(\tilde{P}_2||\tilde{P}_1) > D_{KL}(P_2||P_1)$  if  $(\nu_{-1} - \nu_1)^2 \zeta_{-1} > 1$ . In a nutshell, the bias enhances  $\Delta Bias > 0$  after aggregation if  $(\nu_{-1} - \nu_1)^2 \min\{\zeta_{-1}, \zeta_1\} > 1$ .

## D THE PARAMETERS RELATION IN $(n, \rho_d, \epsilon_{sens}, c)$ -GRAPH.

Given  $(n, \rho_d, \epsilon_{sens}, c)$ -graph, we derive intra-connect, inter-connect probability, and the connection degree for demographic group  $\zeta_i$  in Lemma D.1.

**Lemma D.1** *Suppose the synthetic graph is generated from  $(n, \rho_d, \epsilon_{sens}, c)$ -graph, then we obtain the intra-connect and inter-connect probability as follows:*

$$p_{conn} = \frac{\rho_d \epsilon_{sens}}{c^2 + (1-c)^2}, \quad q_{conn} = \frac{\rho_d(1 - \epsilon_{sens})}{2c(1-c)}.\quad (19)$$



the connection degree for the demographic group (including the self-loop degree) is given by

$$\begin{aligned}\zeta_{-1} &= n_1 q_{conn} + (n_{-1} - 1) p_{conn} + 1, \\ \zeta_1 &= n_{-1} q_{conn} + (n_1 - 1) p_{conn} + 1.\end{aligned}\quad (20)$$

where the node number with for the demographic group  $s = -1$  and  $s = 1$  are given by  $n_{-1} = n(1 - c)$ , and  $n_1 = nc$ , respectively.

## E TOPOLOGY AMPLIFIES BIAS IN ONE-LAYER GCN

Section 3 demonstrates that GCN-like aggregation operation amplifies node representation bias for graph data with large topology bias. However, it is still unclear whether such observation holds for general GNNs or not. Generally speaking, this problem is quite fundamental and challenging to understand the role of topology in fair graph learning. In this section, we try to move a step toward this problem by considering a one-layer GCN. Prior to comparing the prediction for one-layer GCN and one-layer MLP, we first provide the connection between demographic parity and mutual information of sensitive attributes and predictions. Then, we theoretically compare the prediction bias of one-layer GCN and one-layer MLP through the lens of mutual information.

### E.1 PREDICTION BIAS AND MUTUAL INFORMATION

Here, we only consider binary sensitive attributes  $\mathbf{s} \in \{-1, 1\}$  and binary labels  $\mathbf{y} \in \{-1, 1\}$ . Similarly, we can define  $\hat{\mathbf{y}} \in \{-1, 1\}$  as the binary model predictions. In the fairness community, demographic parity, defined as the average prediction difference among different sensitive attribute groups, is the most commonly used fairness metric, i.e.,  $\Delta DP = |\mathbb{P}(\hat{\mathbf{y}}|\mathbf{s} = 1) - \mathbb{P}(\hat{\mathbf{y}}|\mathbf{s} = -1)|$ . From the mutual information perspective, the correlation between sensitive attributes  $\mathbf{s}$  and prediction  $\hat{\mathbf{y}}$  can be measured by  $I(\mathbf{s}; \hat{\mathbf{y}})$ . In this subsection, we provide an inherent connection between demographic parity  $\Delta DP$  and mutual information  $I(\mathbf{s}; \hat{\mathbf{y}})$  as follows:

**Theorem E.1** *For binary sensitive attributes  $\mathbf{s} \in \{-1, 1\}$  and binary prediction  $\hat{\mathbf{y}} \in \{-1, 1\}$ , demographic parity  $\Delta DP$  and mutual information  $I(\mathbf{s}; \hat{\mathbf{y}})$  satisfies  $I(\mathbf{s}; \hat{\mathbf{y}}) \leq 2\Delta DP$*

**Proof:** *For simplicity, we defined the joint probability as  $\alpha_{i \cup j} = \mathbb{P}(\hat{\mathbf{y}} = i, \mathbf{s} = j)$  and condition probability as  $\alpha_{i|j} = \mathbb{P}(\hat{\mathbf{y}} = i|\mathbf{s} = j)$ . Considering the log ratio between joint distribution and margin product probability, we have*

$$\begin{aligned}\log \frac{\mathbb{P}(\hat{\mathbf{y}} = i, \mathbf{s} = j)}{\mathbb{P}(\hat{\mathbf{y}} = i)\mathbb{P}(\mathbf{s} = j)} &= \log \frac{\alpha_{i|j}}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)} = \log \left( 1 + \frac{(\alpha_{i|j} - \alpha_{i|-j})\mathbb{P}(\mathbf{s} = -j)}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)} \right) \\ &\stackrel{(d)}{\leq} (\alpha_{i|j} - \alpha_{i|-j}) \frac{\mathbb{P}(\mathbf{s} = -j)}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)} \leq \Delta DP \frac{\mathbb{P}(\mathbf{s} = -j)}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)}.\end{aligned}\quad (21)$$

where inequality (d) holds due to  $\log(1 + x) \leq x$  for any  $x > -1$ . According to the definition of mutual information, we have

$$\begin{aligned}I(\mathbf{s}; \hat{\mathbf{y}}) &= \sum_{i,j} \mathbb{P}(\hat{\mathbf{y}} = i, \mathbf{s} = j) \log \frac{\mathbb{P}(\hat{\mathbf{y}} = i, \mathbf{s} = j)}{\mathbb{P}(\hat{\mathbf{y}} = i)\mathbb{P}(\mathbf{s} = j)} = \sum_{i,j} \alpha_{i \cup j} \log \frac{\alpha_{i|j}}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)} \\ &\leq \Delta DP \sum_{i,j} \alpha_{i \cup j} \frac{\mathbb{P}(\mathbf{s} = -j)}{\sum_j \alpha_{i|j}\mathbb{P}(\mathbf{s} = j)} = \Delta DP \sum_{i,j} \mathbb{P}(\mathbf{s} = -j)\mathbb{P}(\mathbf{s} = j)|\hat{\mathbf{y}} = i) \\ &\stackrel{(f)}{\leq} \Delta DP \sum_i \left[ \sum_j \mathbb{P}(\mathbf{s} = -j) \right] \left[ \sum_j \mathbb{P}(\mathbf{s} = j|\hat{\mathbf{y}} = i) \right] = 2\Delta DP\end{aligned}\quad (22)$$

where inequality (f) holds due to  $\sum_i a_i b_i \leq \sum_i a_i \sum_i b_i$  for non-negative  $a_i$  and  $b_i$ .  $\square$

Theorem E.1 shows there is a strong connection between demographic parity and mutual information for binary sensitive attributes and binary labels, i.e., the mutual information is upper bounded by demographic parity multiplied by 2.

## E.2 PREDICTION BIAS COMPARISON BETWEEN ONE-LAYER GCN AND ONE-LAYER MLP

Considering the strong connection between mutual information and demographic parity, we investigate prediction bias comparison between one-layer GCN and one-layer MLP through the lens of mutual information. For one-layer MLP model, the prediction is given by  $\hat{y}_{MLP} = \sigma(\mathbf{X}\mathbf{W}_{MLP})$ , where  $\mathbf{W}_{MLP}$  is trainable parameter for MLP. Similarly, for one-layer GCN, the prediction is given by  $\hat{y}_{GCN} = \sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}_{GCN})$ , where  $\mathbf{W}_{GCN}$  is the trainable parameters for GCN. Define  $\tilde{\mathbf{X}} = \tilde{\mathbf{A}}\mathbf{X}$ , it is easy to see that one-layer MLP model and one-layer GCN model are almost the same except with different node features. Based on Theorem 3.5, the aggregated node features  $\tilde{\mathbf{X}}$  embrace higher presentation bias than that of  $\mathbf{X}$ . In other words, the bias of input data for one-layer GCN is higher than that of one-layer MLP.

For the prediction bias, note that sensitive attributes  $\mathbf{s}$ , node features  $\mathbf{X}$ , and prediction  $\hat{y}$  form a Markov chain  $\mathbf{X} \rightarrow \mathbf{s} \rightarrow \hat{y}$  since  $P(\hat{y}|\mathbf{s}, \mathbf{X}) = P(\hat{y}|\mathbf{X})$  for the model with vanilla training. Based on data processing inequality, it is easy to obtain

$$\begin{aligned} I(\mathbf{s}; \hat{y}_{MLP}) &= I(\mathbf{s}; \mathbf{X}) - I(\mathbf{s}; \mathbf{X}|\hat{y}_{MLP}), \\ I(\mathbf{s}; \hat{y}_{GCN}) &= I(\mathbf{s}; \tilde{\mathbf{X}}) - I(\mathbf{s}; \tilde{\mathbf{X}}|\hat{y}_{GCN}). \end{aligned} \quad (23)$$

In other words, when  $I(\mathbf{s}; \mathbf{X}|\hat{y}_{MLP}) = I(\mathbf{s}; \tilde{\mathbf{X}}|\hat{y}_{GCN})$ , the higher input data bias will lead to higher prediction bias. For one-layer MLP and one-layer GCN, if the condition in Theorem 3.5 are satisfied, the prediction bias of one-layer GCN is also larger than that of MLP.

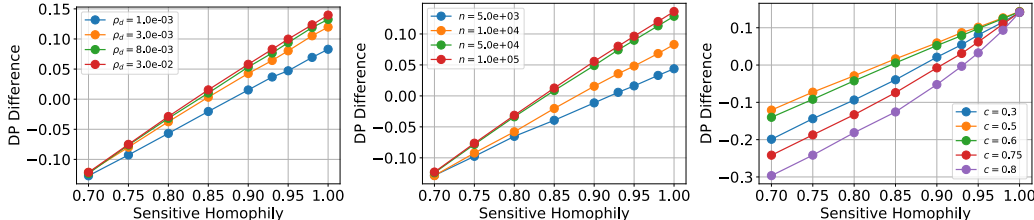


Figure 3: The difference of demographic parity for message passing with different initial covariance matrices. **Left:** DP difference for different graph connection density  $\rho_d$  with sensitive attribute ratio  $c = 0.5$  and number of nodes  $n = 10^4$ ; **Middle:** DP difference for different number of nodes  $n$  with sensitive attribute ratio  $c = 0.5$  and graph connection density  $\rho_d = 10^{-3}$ ; **Right:** DP difference for different sensitive attribute ratio  $c$  with graph connection density  $\rho_d = 10^{-3}$  and number of nodes  $n = 10^4$ .

## F DISCUSSION ON TOPOLOGY ENHANCEMENT ANALYSIS

We provide a comprehensive comparison with CSBM, concentration property in GNNs, and persistent homology. CSBM is also adopted in GNNs literature (Baranwal et al., 2021; 2023). there are several significant differences from the following perspectives:

- **Research problem.** Our research problem (i.e., bias amplification in fair graph learning) is significantly different from that in [A, B]. Specifically, (Baranwal et al., 2021) theoretically investigates why graph convolution improves linear separability and out-of-distribution generalization, while (Baranwal et al., 2023) finds that graph convolution shows stronger classification ability than MLP given the distance of mean node features in CSBM. In other words, we investigate the node presentation bias difference before/after aggregation while (Baranwal et al., 2021; 2023) investigates classification performance comparison between GNNs and MLP.
- **Assumptions.** Firstly, we would like to clarify that, even though CSBM is considered in our work and (Baranwal et al., 2021; 2023), the parameters in CSBM are different. We highlight the sensitive homophily parameter in CSBM since this parameter is highly related to the sufficient condition for bias amplification. Secondly, there is no additional assumption in our analysis except the graph data generated by CSBM. However, the main findings in (Baranwal et al., 2021; 2023) are based on intra-class and inter-class connection probability  $p, q = \Omega(\frac{\log^2 n}{n})$ , where  $n$  is the number of nodes.

- Findings. (Baranwal et al., 2021; 2023) quantitatively measures the improvements of graph convolution compared with MLP in terms of prediction performance under assumption intra-class and inter-class connection probability  $p, q = \Omega(\frac{\log^2 n}{n})$ , while our work mainly identifies when and why bias amplification happens. It is shown that bias amplification conditionally happens, where the sufficient condition is related to many factors, including sensitive homophily, connection density, and the number of nodes. Additionally, we adopt an upper bound of mutual information as the node presentation bias (defined based on demographic group distribution distance) and then measure the bias difference before/after aggregation.

GNN’s concentration property represents all node presentation convergence after stacking of aggregations (Nt and Maehara, 2019; Ma et al., 2022; Baranwal et al., 2021). There are differences between bias enhancement and concentration property in GNN:

- Definition. Concentration property means that the node representation of all nodes convergence after GNN aggregation. Bias enhancement represents the node representation for different sensitive groups that are more distinguished. In fact, such two properties are somehow contrary since perfect concentration leads to zero bias.
- Aggregation. For concentration property, only the normal features and topology are involved in the analysis. The high-level interpretation of concentration is that aggregation acts like a low-frequency filter and such an “average” effect leads to node representation convergence. In the sensitive homophily coefficient, we would like to clarify that the sensitive attributes are not included for node feature aggregation due to law restrictions. Even though the sensitive attribute is included in GNN aggregation, and all node representations are the same, it does not represent the bias enhancement (actually zero bias.). The bias in GNN represents the highly different node representations or predictions among different sensitive groups (defined by sensitive attributes).
- Comparison. We also provide the connection between concentration property and bias enhancement in GNN. The intuition of why GNN enhances the bias, high sensitive homophily represents that the nodes with the same sensitive attributes are connected with high probability. Considering concentration property, the node representation for the same sensitive attribute is more similar after aggregation, therefore leading to highly different representations for different sensitive attribute groups. Notice that such behavior only happens for high sensitive homophily coefficients and shallow GNN. When the node with different sensitive attribute groups is connected randomly, the bias enhancement would not happen or be insignificant due to random concentration. When adopting deep GNNs, all node representations converge and have no bias enhancement. Unfortunately, due to concentration property, shallow GNNs are mainly used in practice. As for high sensitive homophily coefficient, such a condition is usually satisfied in practice due to natural graph data property.

Additionally, there are several differences between our proposed optimization scheme and work (Carriere et al., 2021), including definition, dependence, and optimization:

- Definition: Persistent homology is a method for calculating the importance of topological features in the simplicial complex. For example, giving a set of points in a point cloud corresponding to a chair, the task is to detect the object from the points. In this case, there is no connection between any pair of points. Persistent homology is a tool to identify the topological feature (or connection patterns) via gradually building up the connection between points. However, for graphs, they are 1-simplex with explicit connection patterns defined by the set of edges, thus many properties from persistent homology will degenerate to the field of graph theory. For example, applying the persistent homology on the graph is equivalent to building maximum spanning trees (MSTs) using Kruskal algorithm (Kleinberg and Tardos, 2006), which is irrelevant to our proposed optimization scheme.
- Dependence. Persistent homology is generally related to sample features, as shown in the example Point cloud optimization of (Carriere et al., 2021). In other words, persistent homology somehow represents the topological features of all samples. Differentially, in the graph data we focused on, there are node normal attributes, sensitive attributes, and adjacency matrix (graph topology). Based on the definition, the sensitive homophily coefficient is related to sensitive attributes and the adjacency matrix. However, the optimized adjacency matrix is generally dependent on sensitive attributes.
- Optimization. The main challenge for persistent homology-based optimization is generally undifferentiable except in some special cases. (Carriere et al., 2021) develops a general framework to study the differentiability of the persistence of parametrized families of filtrations. In this way, under

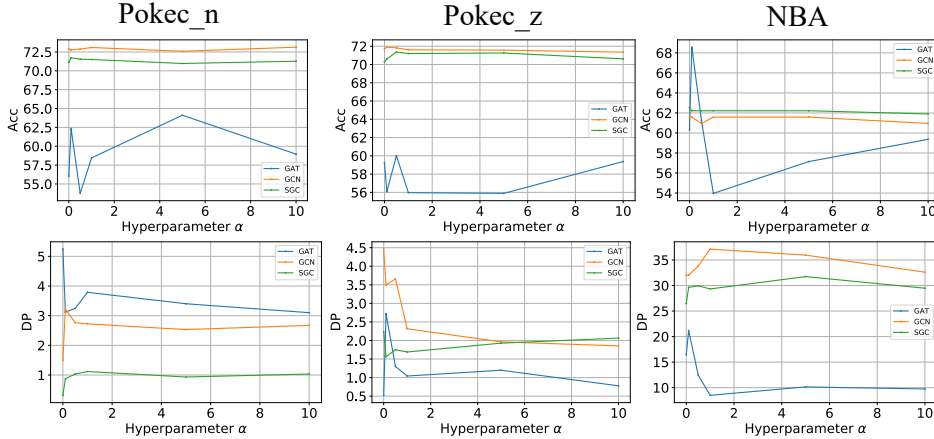


Figure 4: Ablation study on hyperparameter  $\alpha$  in terms of DP and Acc across different GNN backbones on three real-world datasets.

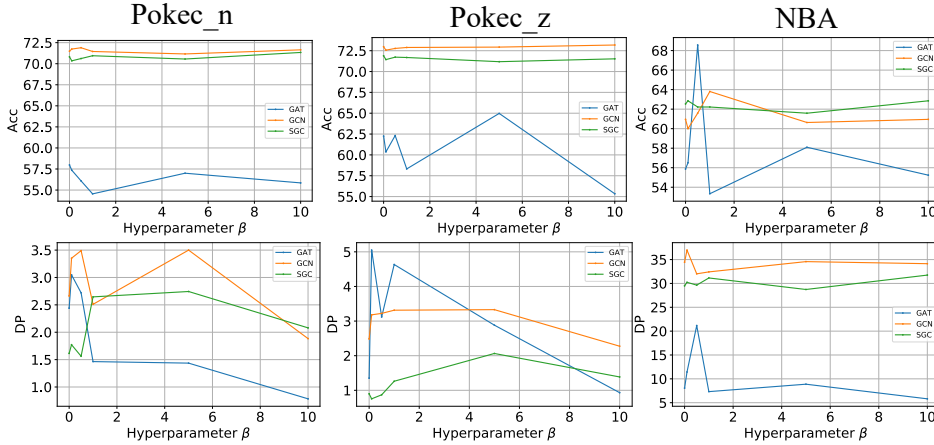


Figure 5: Ablation study on hyperparameter  $\beta$  in terms of DP and Acc across different GNN backbones on three real-world datasets.

mild assumptions, stochastic subgradient descent algorithms can be applied to such functions to converge almost surely to a critical point. For our problem (2), the gradient of loss over topology is differentiable in general. The challenge falls in the constraint of an element value. In our solution, we use a gradient-based method to update the adjacency matrix via variables relaxation and then adopt project operation to satisfy such constraint.

## G EVALUATION AND COMPUTATION COMPLEXITY ANALYSIS FOR FAIRGR.

For algorithmic evaluation of pre-processing FairGR, we compared the prediction performance (including accuracy and fairness) using original graph topology and rewired graph topology across multiple GNN backbones. Denote the number of training nodes and update iterations as  $N$  and  $T$ , respectively. Then the computation complexity for gradient computation and projection of PGD are both  $O(n_{train}^2)$ . Therefore, the total computation complexity to obtain the final rewired graph topology is given by  $O(Tn_{train}^2)$ . The memory consumption is  $O(n_{train}^2)$  due to the storage of graph topology gradient.

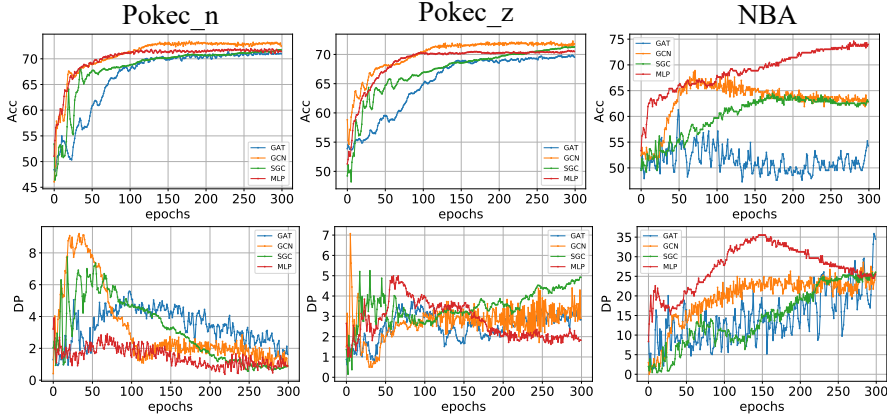


Figure 6: Accuracy (top) and DP (bottom) training curve w.r.t. epochs for different backbones, including GAT, GCN, SGC, and MLP, on three real-world datasets.

## H RELATED WORK

**Graph neural networks.** GNNs achieve state-of-the-art performance for various real-world applications. There are two categories in GNNs model backbones, including spectral-based and spatial-based GNNs. Spectral-based GNNs provide graph convolution operation together with feature transformation (Bruna et al., 2013). Many spatial-based GNNs are also proposed to aggregate the neighbors’ information, including graph attention network (GAT) (Veličković et al., 2018), GraphSAGE (Hamilton et al., 2017), SGC (Wu et al., 2019), APPNP (Klicpera et al., 2019), et al (Gao et al., 2018; Monti et al., 2017; Han et al., 2022b;a).

**Fairness-aware learning on graphs.** Fairness in machine learning has attracted many research efforts to mitigate prediction bias (Chuang and Mroueh, 2020; Zhang et al., 2018; Du et al., 2021; Yurochkin and Sun, 2020; Jiang et al., 2022; Creager et al., 2019; Feldman et al., 2015; Han et al., 2023; Jiang et al., 2023). Fair walk (Rahman et al., 2019) is a fair version of random walk to learn fair node representation via revising neighborhood sampling. From the bias mitigation perspective, adversarial debiasing and contrastive learning are also developed for graph data. For example, work (Dai and Wang, 2021; Bose and Hamilton, 2019; Fisher et al., 2020; Ling et al., 2023) also adopts the adversary to predict the sensitive attribute given the node representation. Fairness-aware representation learning is also developed via node feature masking, graph topology rewires (Agarwal et al., 2021; Köse and Shen, 2021; Dong et al., 2022a) for node classification or link prediction tasks (Laclau et al., 2021; Li et al., 2021; Dong et al., 2021; 2023; Jiang et al., 2022). Interpretation techniques are also introduced in bias mitigation (Dong et al., 2022b). However, the inherent reason behind the observation that GNNs show higher prediction bias than MLP is still missing. In this work, we theoretically and experimentally reveal that many GNNs aggregation schemes boost node representation bias under topology bias. Furthermore, we develop a simple yet effective graph rewiring method, named FairGR, to achieve fair prediction.

## I DATASET STATISTICS AND VISUALIZATION

The data statistical information on three real-world datasets, including Pokec-n, Pokec-z, and NBA, are provided in Table 2. It is seen that the sensitive homophily is even higher than label homophily coefficient among three real-world datasets. Additionally, We visualize the topology for three real-world datasets (Pokec-n, Pokec-z, and NBA) in Figure 7, where different edge types are highlighted with different colors for the top-3 largest connected components in original graphs.

## J TRAINING ALGORITHMS

We summarize the training algorithm for FairGR and provide the pseudo codes in Algorithm 1.

Table 2: Statistical Information on Datasets

Dataset	# Nodes	# Node Features	# Edges	# Training Labels	# Training Sens	Label Homop	Sens Homop
Pokec-n	66569	265	1034094	4398	500	73.23%	95.30%
Pokec-z	67796	276	1235916	5131	500	71.16%	95.06%
NBA	403	95	21242	156	246	39.22%	72.37%

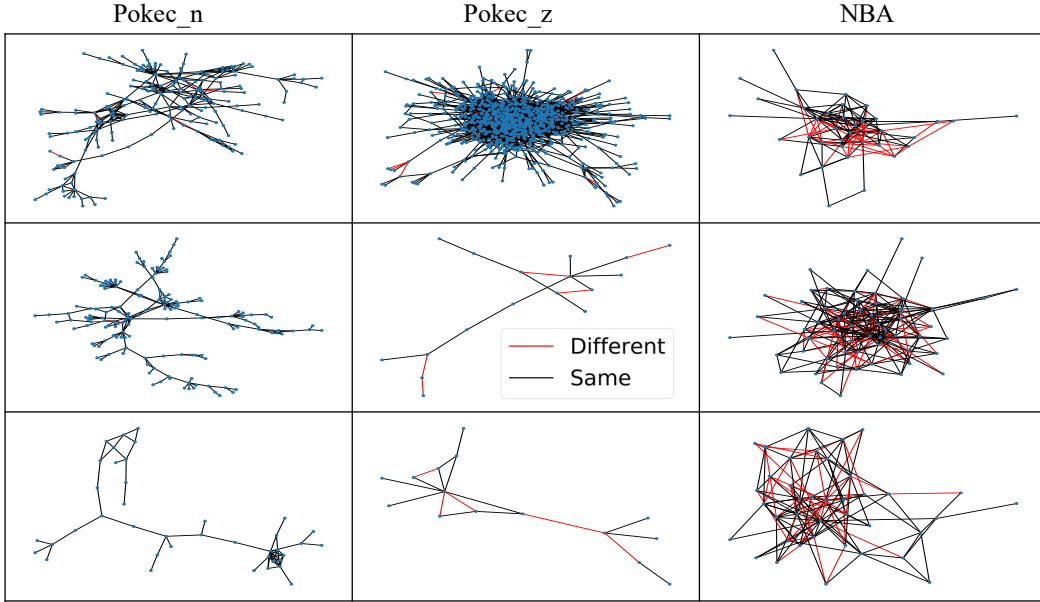


Figure 7: Visualization of topology bias in three real-world datasets, where black edges  $\bullet\text{---}\bullet$  and red edges  $\bullet\text{---}\bullet$  represent the edge with the same or different sensitive attributes for the connected node pair, respectively. We visualize the largest three connected components for each dataset. It is obvious that the sensitive homophily coefficients (the ratio of homo edges) are high in practice, i.e., 95.30%, 95.06%, and 72.37% for Pokec-n, Pokec-z, and NBA dataset, respectively.

## K MORE EXPERIMENTAL RESULTS

### K.1 MORE SYNTHETIC EXPERIMENTAL RESULTS

In this subsection, we provide more experimental results on **different** covariance matrix. Although our theory is only derived for the same covariance matrix, we still observe similar results for the case of **different** covariance matrix. For **node attribute generation**, we generate node attribute with Gaussian distribution  $\mathcal{N}(\mu_1, \Sigma)$  and  $\mathcal{N}(\mu_2, \Sigma)$  for node with binary sensitive attribute, respectively,

---

#### Algorithm 1 FairGR Algorithm

---

**Input:** Graph topology  $\mathbf{A}$ , label  $\mathbf{y}$ , and sensitive attribute  $\mathbf{s}$ ; The total epochs  $T$ ; Hyperparameters  $\alpha$  and  $\beta$ .  
**for** epoch from 1 to  $T$  **do**  
  Calculate the gradient  $\frac{\partial \mathcal{L}(\hat{\mathbf{A}}|\mathbf{s}, \mathbf{y}, \mathbf{A})}{\partial \hat{\mathbf{A}}}$  ;  
  Update graph topology  $\tilde{\mathbf{A}}$  using gradient descent;  
  Project graph topology  $\tilde{\mathbf{A}}$  into a feasible region.  
**end for**

---

where  $\mu_1 = [0, 1]$ ,  $\mu_2 = [1, 0]$  and  $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ . We adopt the same evaluation metric and adjacency matrix generation scheme in Section 5.1

Figure 3 shows the DP difference during message passing with respect to sensitive homophily coefficient for different initial covariance matrices. We observe that a higher sensitive homophily coefficient generally leads to larger bias enhancement. Additionally, higher graph connection density  $\rho_d$ , larger node number  $n$ , and balanced sensitive attribute ratio  $c$  correspond to higher bias enhancement, which is consistent with our theoretical analysis in Theorem 3.5.

## K.2 ABLATION STUDY FOR FAIRGR

For investigating the effect of hyperparameters  $\alpha$  and  $\beta$ , we conduct experiments with different hyperparameters chosen from  $\alpha = \{0.0, 0.1, 0.5, 1.0, 5.0, 10.0\}$  and  $\beta = \{0.0, 0.1, 0.5, 1.0, 5.0, 10.0\}$  while the other one is selected as default. The default value for hyperparameters  $\alpha$  and  $\beta$  are 0.1 and 0.5, respectively. The results of hyperparameter study with respect to  $\alpha$  and  $\beta$  are shown in Figures 4 and 5, respectively. From these results, we can obtain the following observations:

- Hyperparameter  $\alpha$  and  $\beta$  demonstrate different influences on GNN backbone. For example, for Pokec-n dataset, hyperparameter  $\alpha$  only shows a negligible influence on the accuracy of GCN and GAT, while significant for GAT. As for DP, the bias metric is more sensitive to  $\alpha$  compared with accuracy.
- Hyperparameter  $\alpha$  and  $\beta$  demonstrate different influences on graph dataset. For example, GAT achieves the best accuracy and lowest DP with  $\alpha = 0.1$  in NBA dataset, while achieving the lowest accuracy and highest DP in Pokec-n dataset. Such observation indicates the importance of hyperparameter tuning for different datasets.

### K.2.1 EXPERIMENTAL SETTINGS ON REAL-WORLD DATASETS

**Datasets.** The experiments are conducted on three real-world datasets, including Pokec-z, Pokec-n, and NBA (Dai and Wang, 2021). Pokec-z and Pokec-n are sampled from a larger social network Pokec (Takac and Zabovsky, 2012) based on the province in Slovakia. We choose region information and the working field of the users as the sensitive attribute and the predicted label, respectively. NBA dataset includes around 400 NBA basketball players and is collected from a Kaggle dataset<sup>10</sup> and Twitter. The information of players includes age, nationality, and salary in the 2016-2017 season. We choose nationality (U.S. and overseas player) as the binary sensitive attribute, and the prediction label is whether the salary is higher than the median.

**Evaluation Metrics.** For the node classification task, we adopt accuracy to evaluate the classification performance. As for fairness metric, we adopt two most common-used group fairness metrics, including *demographic parity* and *equal opportunity*, to measure the prediction bias (Louizos et al., 2015; Beutel et al., 2017). Specifically, *demographic parity* is defined as the average prediction difference over different sensitive attribute groups, i.e.,  $\Delta_{DP} = |\mathbb{P}(\hat{y} = 1|s = -1) - \mathbb{P}(\hat{y} = 1|s = 1)|$ . Similarly, *equal opportunity* is given by  $\Delta_{EO} = |\mathbb{P}(\hat{y} = 1|s = -1, y = 1) - \mathbb{P}(\hat{y} = 1|s = 1, y = 1)|$ , where  $y$  and  $\hat{y}$  represent the ground-truth label and predicted label, respectively.

### K.3 VANILLA TRAINING BEHAVIORS FOR GNNs AND MLP

For vanilla training across different GNN backbones, we plot the training curve with respect to epochs to investigate the training behaviors for GNNs and MLP. For the training behavior, GNNs, and MLP models both converge in terms of accuracy for the high label homophily dataset Pokec-n and Pokec-z dataset. For high sensitive homophily Pokec-n and Pokec-z dataset, GNNs demonstrate higher prediction than that MLP, while the prediction bias difference is relatively small for the low-sensitive-homophily dataset NBA.

<sup>10</sup><https://www.kaggle.com/noahgift/social-power-nba>

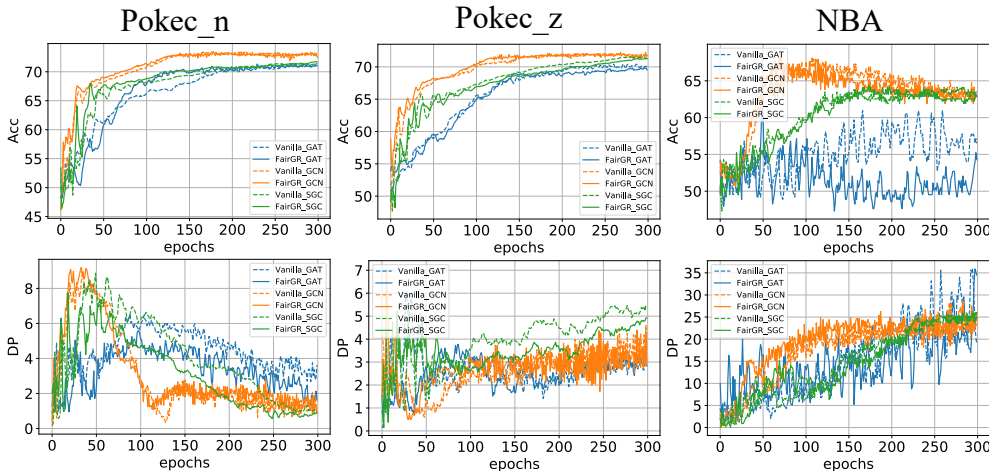


Figure 8: Accuracy (top) and DP (bottom) training curve w.r.t. epochs for different backbones, including GAT, GCN, SGC, and MLP, on three real-world datasets.

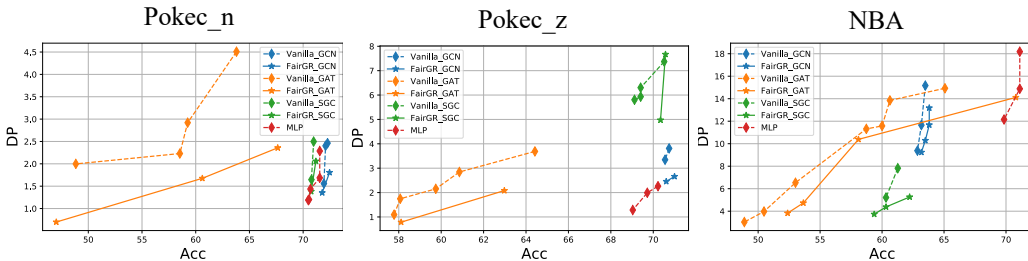


Figure 9: DP and Accuracy trade-off performance on three real-world datasets, including Poken-n, Pokec-z, and NBA, in (manifold) Fair Mixup setting.

#### K.4 FAIRGR RESULTS ON VANILLA TRAINING

For vanilla training, Figure 8 shows the test accuracy and demographic parity curve during training for different GNNs backbones and whether FairGR is adopted for graph topology rewiring. From these results, we can obtain the following observations:

- Different GNNs demonstrate different accuracy and demographic parity performance. For example, for Pokec-n dataset, GCN has the highest accuracy performance and lowest demographic parity, which implies that message passing also matters even for the same graph topology.
- Our proposed FairGR consistently achieves lower demographic parity and comparable accuracy performance on all datasets and backbones.

#### K.5 TRADEOFF PERFORMANCE ON FAIR MIXUP

We also demonstrate that FairGR can achieve better tradeoff performance for different GNN backbones with Fair mixup (Chuang and Mroueh, 2021) in Figure 9. Specifically, since input fair mixup requires calculating model prediction for mixed input batch, it is non-trivial to adopt input fair mixup in our experiments on the node classification task. This is because, for GNN aggregations of neighborhoods’ information, the neighborhood information for the mixed input batch is missing. Instead, we adopt manifold fair mixup for the logit layer in our experiments. Experimental results show that FairGR can achieve better accuracy-fairness tradeoff performance for many GNNs backbones on three datasets.



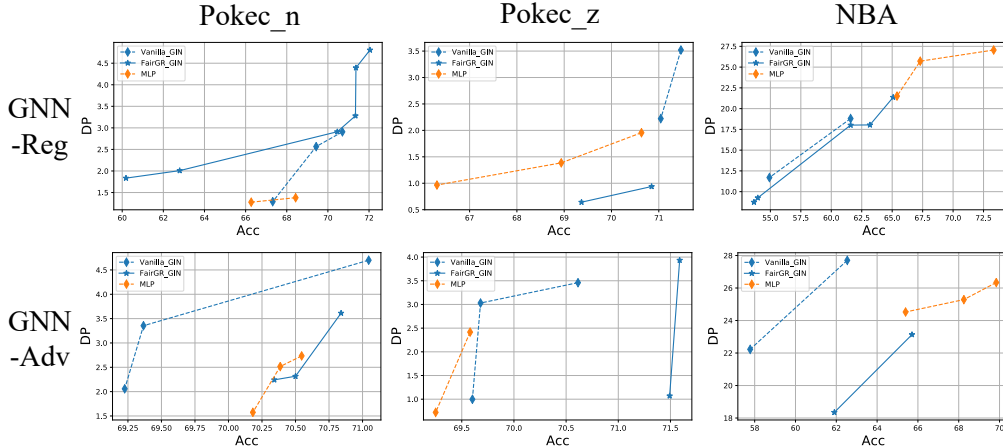


Figure 10: DP and Accuracy trade-off performance for GIN on three real-world datasets compared with adding regularization (Top) and adversarial debiasing (Bottom). The trade-off curve closer to the right bottom corner represents better trade-off performance.

### K.6 MORE EXPERIMENTAL RESULTS ON GIN AND H2GCN

We also demonstrate that FairGR can achieve better tradeoff performance for GIN backbones (Xu et al., 2018) by adding demographic parity regularization and adversarial debiasing in Figure 10. Specifically, GIN is originally designed for graph classification tasks, and the graph pooling layer is adopted in GIN for graph readout. In our paper, we investigate the fair node classification task, which is different from that in (Xu et al., 2018). Therefore, we remove the graph pooling layer and tailer GIN for the node classification task in our experiments. Experimental results show that FairGR can achieve better accuracy-fairness tradeoff performance for GIN on three datasets and two in-processing fairness methods, including adding demographic parity regularization and adversarial debiasing.

Additionally, we compare the performance of GIN and H2GCN (Zhu et al., 2020) over MLP on various real-world datasets and summarize the results in Table 3. It is seen that our proposed FairGR can mitigate bias for GIN on three real-world datasets on vanilla training. More importantly, GIN demonstrates a higher prediction bias compared with MLP model on all three datasets. For example, in terms of demographic parity, the prediction bias of MLP is lower than GIN by 50.15%, 51.74%, and 18.77% on Pokec-z, Pokec-n and NBA datasets. Although popular GNNs, such as GCN and SGC, achieve lower accuracy than MLP in NBA dataset due to the low label homophily coefficient (Ma et al., 2022), H2GCN and H2GCN-GR both achieve higher accuracy than that of MLP, which indicates GNNs backbone design is also important to achieve better tradeoff performance.

Table 3: The performance on Node Classification (GR represents graph topology rewire).

Models	Pokec-z			Pokec-n			NBA		
	Acc (%) ↑	$\Delta_{DP}$ (%) ↓	$\Delta_{EO}$ (%) ↓	Acc (%) ↑	$\Delta_{DP}$ (%) ↓	$\Delta_{EO}$ (%) ↓	Acc (%) ↑	$\Delta_{DP}$ (%) ↓	$\Delta_{EO}$ (%) ↓
MLP	70.48 ± 0.77	1.61 ± 1.29	2.22 ± 1.01	72.48 ± 0.26	1.53 ± 0.89	3.39 ± 2.37	65.56 ± 1.62	22.37 ± 1.87	18.00 ± 3.52
GIN	69.89 ± 2.13	3.23 ± 2.41	2.40 ± 1.16	69.59 ± 1.66	3.17 ± 1.51	3.08 ± 2.80	60.95 ± 3.56	27.54 ± 8.73	27.67 ± 8.36
GIN-GR	72.27 ± 0.65	<b>1.51</b> ± 0.58	2.10 ± 1.52	70.86 ± 0.80	<b>0.85</b> ± 0.43	2.36 ± 1.47	61.27 ± 2.94	25.86 ± 10.96	20.17 ± 10.51
H2GCN	71.57 ± 0.83	3.83 ± 2.02	2.98 ± 1.96	73.02 ± 1.16	1.66 ± 0.86	3.85 ± 0.46	67.30 ± 3.27	35.65 ± 2.91	25.67 ± 5.31
H2GCN-GR	71.14 ± 0.70	1.58 ± 0.54	<b>1.48</b> ± 1.50	72.00 ± 1.11	1.23 ± 0.57	<b>1.83</b> ± 0.76	66.98 ± 1.85	<b>21.89</b> ± 7.03	<b>15.33</b> ± 4.52

### K.7 MORE EXPERIMENTAL RESULTS ON PARTIAL NODES WITH ACCESSIBLE SENSITIVE ATTRIBUTES

We conduct experiments to compare the performance of many representative GNNs over MLP with partial labeling nodes in Table 4, where the sensitive attribute for 50% training nodes are available for

graph topology. Note that the label for training nodes are available for GNNs training. We observe that our proposed GR method can mitigate prediction bias on various GNNs backbones and datasets, and is even lower than that of MLP in most cases. Therefore, our proposed method can also work well in the scenario of partial nodes with accessible sensitive attributes.

Table 4: The performance on Node Classification (GR represents graph topology rewire) with partial labeling nodes in GR. The boldface indicates that the prediction bias of our GR methods is lower than that without GR and MLP.

Models	Pokec-z			Pokec-n			NBA		
	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$	Acc (%) $\uparrow$	$\Delta_{DP}$ (%) $\downarrow$	$\Delta_{EO}$ (%) $\downarrow$
MLP	70.48 $\pm$ 0.77	1.61 $\pm$ 1.29	2.22 $\pm$ 1.01	72.48 $\pm$ 0.26	1.53 $\pm$ 0.89	3.39 $\pm$ 2.37	65.56 $\pm$ 1.62	22.37 $\pm$ 1.87	18.00 $\pm$ 3.52
GAT	69.76 $\pm$ 1.30	2.39 $\pm$ 0.62	2.91 $\pm$ 0.97	71.00 $\pm$ 0.48	3.71 $\pm$ 2.15	7.50 $\pm$ 2.88	57.78 $\pm$ 10.65	20.12 $\pm$ 16.18	13.00 $\pm$ 13.37
GAT-GR	57.10 $\pm$ 6.02	<b>1.56</b> $\pm$ 1.76	<b>2.15</b> $\pm$ 2.40	57.30 $\pm$ 4.96	<b>1.12</b> $\pm$ 1.01	<b>2.43</b> $\pm$ 2.47	58.62 $\pm$ 1.49	<b>0.48</b> $\pm$ 0.97	<b>0.84</b> $\pm$ 1.69
GCN	71.78 $\pm$ 0.37	3.25 $\pm$ 2.35	2.36 $\pm$ 2.09	73.09 $\pm$ 0.28	3.48 $\pm$ 0.47	5.16 $\pm$ 1.38	61.90 $\pm$ 1.00	23.70 $\pm$ 2.74	17.50 $\pm$ 2.63
GCN-GR	71.04 $\pm$ 0.26	<b>1.38</b> $\pm$ 0.68	2.33 $\pm$ 1.28	71.11 $\pm$ 0.38	<b>1.11</b> $\pm$ 0.53	<b>1.12</b> $\pm$ 0.88	70.00 $\pm$ 0.86	<b>4.80</b> $\pm$ 2.43	<b>4.22</b> $\pm$ 1.09
SGC	71.24 $\pm$ 0.46	4.81 $\pm$ 0.30	4.79 $\pm$ 2.27	71.46 $\pm$ 0.41	2.22 $\pm$ 0.29	3.85 $\pm$ 1.63	63.17 $\pm$ 0.63	22.56 $\pm$ 3.94	14.33 $\pm$ 2.16
SGC-GR	69.92 $\pm$ 0.21	3.72 $\pm$ 1.15	4.42 $\pm$ 0.67	70.92 $\pm$ 0.62	<b>0.91</b> $\pm$ 1.17	3.54 $\pm$ 1.21	69.47 $\pm$ 0.54	<b>2.35</b> $\pm$ 1.41	<b>10.55</b> $\pm$ 2.47
GIN	69.89 $\pm$ 2.13	3.23 $\pm$ 2.41	2.40 $\pm$ 1.16	69.59 $\pm$ 1.66	3.17 $\pm$ 1.51	3.08 $\pm$ 2.80	60.95 $\pm$ 3.56	27.54 $\pm$ 8.73	27.67 $\pm$ 8.36
GIN-GR	68.57 $\pm$ 0.98	<b>1.65</b> $\pm$ 1.23	<b>2.14</b> $\pm$ 1.45	69.90 $\pm$ 0.47	<b>1.00</b> $\pm$ 0.74	<b>2.08</b> $\pm$ 0.77	68.62 $\pm$ 4.04	<b>10.98</b> $\pm$ 5.81	<b>6.09</b> $\pm$ 8.27
H2GCN	71.57 $\pm$ 0.83	3.83 $\pm$ 2.02	2.98 $\pm$ 1.96	73.02 $\pm$ 1.16	1.66 $\pm$ 0.86	3.85 $\pm$ 0.46	67.30 $\pm$ 3.27	35.65 $\pm$ 2.91	25.67 $\pm$ 5.31
H2GCN-GR	70.50 $\pm$ 0.16	2.16 $\pm$ 0.76	<b>1.84</b> $\pm$ 1.06	70.61 $\pm$ 0.79	<b>1.52</b> $\pm$ 0.91	<b>1.96</b> $\pm$ 1.40	66.04 $\pm$ 2.09	<b>13.23</b> $\pm$ 3.67	<b>8.78</b> $\pm$ 5.24

## L FUTURE WORK

There are two lines of follow-up research directions. Firstly, The generalization of the theoretical analysis on why aggregation enhances bias in GNN can be further extended. As a pilot study, we theoretically investigate why this phenomenon happens for GCN-like aggregation under random graph topology generated by stochastic block model and Gaussian mixture feature distribution. The more general analysis of other aggregation operations, random graph models, and other feature distributions can be extended. The other line focuses on the graph topology rewire algorithmic perspective, including improving the efficiency and effectiveness of FairGR via designing different objectives for graph topology rewire. Additionally, neighborhood sampling before aggregation is also beneficial for fairness (Lin et al., 2023). For the partial labeling nodes case, the advanced GR method can be developed. For example, a sensitive attribute estimator can be trained using partial labeling nodes and provide surrogate sensitive attributes for unlabeled nodes. Subsequently, our proposed GR methods can be adopted using both sensitive attributes and surrogates for labeling and unlabeled nodes.