Appendix for "FairGLite: Fair Graph Representation Learning with Limited Demographics"

1 Theoretical analysis

This section provides additional theoretical analysis demonstrating the performance guarantees of FairGLite: Specifically, for a binary classification task, statistical parity is defined as: $\Delta_{DP} = |P\left(\hat{y} = 1|s = 0\right) - P\left(\hat{y} = 1|s = 1\right)|$. We consider the binary classification task and examine the properties of the Softmax function in this context. Let P_1 and P_2 represent the probabilities of class 1 (c_1) and class 2 (c_2) , respectively. The function $Softmax(\cdot)$ is Lipschitz continuous with a Lipschitz constant L. Due to this Lipschitz continuity, the difference in output probabilities can be bounded by the difference in input vectors:

$$||f(\mathbf{h_i}) - f(\mathbf{h_j})|| = |P_1 - P_2| + |(1 - P_1) - (1 - P_2)|$$

= $2|P_1 - P_2| \le L||\mathbf{h_i} - \mathbf{h_i}||$ (1)

where $\mathbf{h_i}$ is the node representation for $\forall v_i \in S_d$ and $\mathbf{h_j}$ is the node representation for $\forall v_j \in S_f$.

Building on this, we can rewrite the statistical parity as follows:

$$\Delta_{DP} = \left| \frac{1}{N_d} \sum_{i \in \mathcal{S}_d} f(\mathbf{z}_i)_1 - \frac{1}{N_f} \sum_{j \in \mathcal{S}_f} f(\mathbf{z}_j)_1 \right|$$
(2)

where $\mathbf{z}_i = W^l \mathbf{h}_i^{(l)}$, and $W^{(l)}$ is the weight matrix at layer l.

Hence, the upper bound of disparity of DP depends on the node representation disparity. Given that the Graph Attention Networks (GAT) adopt the message passing by assigning different weights to neighbor nodes. Hence, we can write the node aggregation process as follows:

$$\mathbf{h}_{u}^{(l)} = \xi \, \mathbf{h}_{u}^{(l-1)} + \sum_{k \in \mathcal{N}(u)} \alpha_{u,k}^{(l)} \operatorname{ReLU}(\mathbf{W}^{(l)} \, \mathbf{h}_{k}^{(l-1)}), \quad \alpha_{u,k}^{(l)} = \frac{\exp(e_{u,k}^{(l)})}{\sum_{k \in \mathcal{N}(u)} \exp(e_{u,k}^{(l)})}$$
(3)

Building on this, the node representation disparity arising from the node aggregation process can be measured. Specifically, for a node \$v_i\$, the neighbor information during the aggregation process is:

$$\begin{split} \mathbf{h}_{i}^{(l)} &= \xi \mathbf{h}_{i}^{(l-1)} + \sum_{u \in N(i) \cap \mathcal{S}_{d}} \alpha_{i,u}^{(l)} \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) + \sum_{u \in N(i) \cap \mathcal{S}_{f}} \alpha_{i,u}^{(l)} \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \ \, (4) \\ &- \frac{1}{N_{d}^{2}} \sum_{u \in \mathcal{S}_{d}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \\ &+ \frac{1}{N_{d} N_{f}} \sum_{u \in \mathcal{S}_{f}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \\ &+ \left[\frac{1}{N_{d}^{2}} \sum_{u \in \mathcal{S}_{d}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \right] \\ &- \frac{1}{N_{d} N_{f}} \sum_{u \in \mathcal{S}_{f}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \right] \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \end{split}$$

Let each node v_i have a node representation $h_i^{(l)}$ subject to $\mu^{(d)} - \Delta^l \leq h_i^{(l)} \leq \mu^{(d)} + \Delta^l$, where the parameter Δ^l serves as a tolerance per layer indicating the allowed deviation of the representation from group mean $(\mu^{(d)})$ along each coordinate. Hence, we can re-write the Equation 4 as follow:

$$\begin{split} \mathbf{h}_{i}^{(l)} &= \xi \mathbf{h}_{i}^{(l-1)} + \sum_{u \in N(i) \cap \mathcal{S}_{d}} \alpha_{i,u}^{(l)} \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) + \sum_{u \in N(i) \cap \mathcal{S}_{f}} \alpha_{i,u}^{(l)} \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \ \, (5) \\ &- \frac{1}{N_{d}^{2}} \sum_{u \in \mathcal{S}_{d}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \\ &+ \frac{1}{N_{d} N_{f}} \sum_{u \in \mathcal{S}_{f}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}) \\ &+ \left[\frac{1}{N_{d}^{2}} \sum_{u \in \mathcal{S}_{d}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \right] \\ &- \frac{1}{N_{d} N_{f}} \sum_{u \in \mathcal{S}_{f}} k(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)})) \right] \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \\ &\pm C^{(l)} \cdot \frac{1}{4} \left\| W^{(l)} \right\|_{2} \left(\left(1 + \frac{2}{N_{d}} \right) \sqrt{d_{h}} \Delta^{l-1} + \Delta_{\text{base}}^{(l)} \right) \end{split}$$

Building on this, for nodes $v_i \in S_d$, we have:

$$\frac{1}{N_{d}} \sum_{i \in \mathcal{S}_{d}} \mathbf{h}_{i}^{(l)} \in \left[\mu^{(d)} + \left(\mu_{l-1}^{(d)} + \frac{1}{N_{d}} \sum_{i \in \mathcal{S}_{d}} \left(\sum_{u \in \mathcal{N}(i) \cap \mathcal{S}_{f}} \alpha_{i,u}^{(l)} \right) \left(\mu_{l-1}^{(f)} - \mu_{l-1}^{(d)} \right) \right) \\
+ \frac{1}{N_{d}^{2} N_{f}} \sum_{i \in \mathcal{S}_{d}} \sum_{u \in \mathcal{S}_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \left(\mu_{l-1}^{(f)} - \mu_{l-1}^{(d)} \right) \right] \\
\pm \left[C^{(l)} \cdot \frac{1}{4} \left\| W^{(l)} \right\|_{2} \left(\left(1 + \frac{2}{N_{d}} \right) \sqrt{d_{h}} \Delta^{l-1} + \Delta_{\text{base}}^{(l)} \right) \right]$$

A similar way can be applied for S_f as follows:

$$\frac{1}{N_{f}} \sum_{i \in \mathcal{S}_{f}} \mathbf{h}_{i}^{(l)} \in \left[\mu^{(f)} + \left(\mu_{l-1}^{(f)} + \frac{1}{N_{f}} \sum_{i \in \mathcal{S}_{f}} \left(\sum_{u \in \mathcal{N}(i) \cap \mathcal{S}_{d}} \alpha_{i,u}^{(l)} \right) \left(\mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \right) \right) \\
+ \frac{1}{N_{f}^{2} N_{d}} \sum_{i \in \mathcal{S}_{f}} \sum_{u \in \mathcal{S}_{d}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \left(\mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \right) \right] \\
\pm \left[C^{(l)} \cdot \frac{1}{4} \left\| W^{(l)} \right\|_{2} \left(\left(1 + \frac{2}{N_{f}} \right) \sqrt{d_{h}} \Delta^{l-1} + \Delta_{\text{base}}^{(l)} \right) \right]$$

Building on this, the upper bound of the consequent representation discrepancy on node representations between two demographic groups is defined as follows:

$$\begin{aligned} \mathbf{h}_{D}^{(l)} &= \left\| \frac{1}{N_{d}} \sum_{i \in S_{d}} \mathbf{h}_{i}^{(l)} - \frac{1}{N_{f}} \sum_{j \in S_{f}} \mathbf{h}_{j}^{(l)} \right\| \\ &\leq \left| 1 - \left(\frac{1}{N_{d}} \sum_{i \in S_{d}} \sum_{u \in S_{f}} \alpha_{i,u}^{(l)} + \frac{1}{N_{f}} \sum_{j \in S_{f}} \sum_{u \in \mathcal{N}(j) \cap S_{d}} \alpha_{j,u}^{(l)} \right) \\ &- \frac{1}{N_{d}^{2} N_{f}} \sum_{i \in S_{d}} \sum_{u \in S_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \\ &- \frac{1}{N_{f}^{2} N_{d}} \sum_{j \in S_{f}} \sum_{u \in S_{d}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{j}^{(l-1)}) \right) \right| \left\| \mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \right\| \\ &+ \left\| \mu^{(d)} - \mu^{(f)} \right\| + \varepsilon_{d}^{(l)} + \varepsilon_{f}^{(l)} \\ &\leq \left(\frac{1}{N_{d}} \sum_{i \in S_{d}} \sum_{u \in \mathcal{N}_{f}} \alpha_{i,u}^{(l)} + \frac{1}{N_{f}} \sum_{j \in S_{f}} \sum_{u \in \mathcal{N}(j) \cap S_{d}} \alpha_{j,u}^{(l)} + 1 \\ &- \frac{1}{N_{d}^{2} N_{f}} \sum_{i \in S_{d}} \sum_{u \in S_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \right) \right\| \mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \| \\ &+ \left\| \mu^{(d)} - \mu^{(f)} \right\| + \varepsilon_{d}^{(l)} + \varepsilon_{f}^{(l)} \\ &\leq \left(3 - \frac{1}{N_{d}^{2} N_{f}} \sum_{i \in S_{d}} \sum_{u \in S_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \right\| \mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \| \\ &+ \left\| \mu^{(d)} - \mu^{(f)} \right\| \pm \left[C^{(l)} \cdot \frac{1}{4} \left\| W^{(l)} \right\|_{2} \left(\left(1 + \frac{2}{N_{f}} \right) \sqrt{d_{h}} \Delta^{l-1} + \Delta_{\text{base}}^{(l)} \right) \right] \end{aligned}$$

Armed with this node representation disparity, the effect on statistical parity can then be measured. Given that:

$$\Delta_{\text{DP}} = \left| \frac{1}{N_d} \sum_{i \in \mathcal{S}_d} f(\mathbf{z}_i)_1 - \frac{1}{N_f} \sum_{j \in \mathcal{S}_f} f(\mathbf{z}_j)_1 \right|$$
 (10)

As discussed previously, given that $||f(\mathbf{h_i}) - f(\mathbf{h_j})|| = 2|P_1 - P_2| \le L||\mathbf{h_i} - \mathbf{h_j}||$. Hence, we can rewrite it as:

$$2 | f(\mathbf{z}_i)_1 - f(\mathbf{z}_{\mu^{(d)}})_1 | \le L | |\mathbf{z}_i - \mathbf{z}_{\mu^{(d)}} | |$$
 (11)

Therefore, the following inequality holds:

$$f(\mathbf{z}_{\mu^{(d)}})_1 - \frac{L}{2} \|\mathbf{z}_i - \mathbf{z}_{\mu^{(d)}}\| \le f(\mathbf{z}_i)_1 \le f(\mathbf{z}_{\mu^{(d)}})_1 + \frac{L}{2} \|\mathbf{z}_i - \mathbf{z}_{\mu^{(d)}}\|$$
(12)

Let $\mathbf{z}_i = \mathbf{W}^{(l)} \mathbf{h}_i^{(l)}$ for node i, and $\mathbf{z}_{\mu^{(d)}} = \mathbf{W}^{(l)} \mu_l^{(d)}$, $\mathbf{z}_{\mu^{(f)}} = \mathbf{W}^{(l)} \mu_l^{(f)}$ be the group means in logits space.

$$f(\mathbf{z}_{\mu^{(d)}})_{1} - f(\mathbf{z}_{\mu^{(f)}})_{1} - \frac{1}{N_{d}} \sum_{i=1}^{N_{d}} \frac{L}{2} \|\mathbf{z}_{i} - \mathbf{z}_{\mu^{(d)}}\| - \frac{1}{N_{f}} \sum_{j=1}^{N_{f}} \frac{L}{2} \|\mathbf{z}_{j} - \mathbf{z}_{\mu^{(f)}}\|$$

$$\leq \frac{1}{N_{d}} \sum_{i \in \mathcal{S}_{d}} f(\mathbf{z}_{i})_{1} - \frac{1}{N_{f}} \sum_{j \in \mathcal{S}_{f}} f(\mathbf{z}_{j})_{1} \leq f(\mathbf{z}_{\mu^{(d)}})_{1} - f(\mathbf{z}_{\mu^{(f)}})_{1} + \frac{1}{N_{d}} \sum_{i=1}^{N_{d}} \frac{L}{2} \|\mathbf{z}_{i} - \mathbf{z}_{\mu^{(d)}}\| + \frac{1}{N_{f}} \sum_{j=1}^{N_{f}} \frac{L}{2} \|\mathbf{z}_{j} - \mathbf{z}_{\mu^{(f)}}\|$$

Considering the obtained h_D , the following inequality holds:

$$\begin{aligned} \left\| \mathbf{z}_{i} - \mathbf{z}_{\mu^{(d)}} \right\| &= \left\| \mathbf{W}^{(l)} \left(\mathbf{h}_{i}^{(l)} - \mu_{l}^{(d)} \right) \right\| \leq \left\| \mathbf{W}^{(l)} \right\|_{2} \left\| \mathbf{h}_{i}^{(l)} - \mu_{l}^{(d)} \right\| \end{aligned} \tag{14}$$

$$\leq \left\| \mathbf{W}^{(l)} \right\|_{2} \left[\sqrt{d_{h}} \Delta^{l} + C^{(l)} \left(\left(3 - \frac{1}{N_{d}^{2} N_{f}} \sum_{p \in \mathcal{S}_{d}} \sum_{q \in \mathcal{S}_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{q}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{p}^{(l-1)}) \right) \right) - \frac{1}{N_{f}^{2} N_{d}} \sum_{j \in \mathcal{S}_{f}} \sum_{u \in \mathcal{S}_{d}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{j}^{(l-1)}) \right) \right) \left\| \mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \right\| + \left\| \mu^{(d)} - \mu^{(f)} \right\| + \left[C^{(l)} \cdot \frac{1}{4} \left\| W^{(l)} \right\|_{2} \left(\left(1 + \frac{2}{N_{f}} \right) \sqrt{d_{h}} \Delta^{l-1} + \Delta_{\text{base}}^{(l)} \right) \right] \right) \end{aligned}$$

A similar expression can be derived for $|\mathbf{z}i - \mathbf{z}\mu^{(f)}|$. Hence, we can write the upper bound of the DP as:

$$\Delta_{\mathrm{DP}} \leq \left| f(\mathbf{z}_{\mu^{(d)}})_{1} - f(\mathbf{z}_{\mu^{(f)}})_{1} \right| + \frac{L}{2} \left(\frac{1}{N_{d}} \sum_{i} \left\| \mathbf{z}_{i} - \mathbf{z}_{\mu^{(d)}} \right\| + \frac{1}{N_{f}} \sum_{j} \left\| \mathbf{z}_{j} - \mathbf{z}_{\mu^{(f)}} \right\| \right) \tag{15}$$

$$\leq \frac{L}{2} \left\| \mathbf{W}^{(l)} \right\|_{2} \left[\left(3 - \frac{1}{N_{d}^{2} N_{f}} \sum_{i \in \mathcal{S}_{d}} \sum_{u \in \mathcal{S}_{f}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{i}^{(l-1)}) \right) \right) - \frac{1}{N_{f}^{2} N_{d}} \sum_{j \in \mathcal{S}_{f}} \sum_{u \in \mathcal{S}_{d}} k \left(\sigma(\mathbf{W}^{(l)} \mathbf{h}_{u}^{(l-1)}), \sigma(\mathbf{W}^{(l)} \mathbf{h}_{j}^{(l-1)}) \right) \right) \left\| \mu_{l-1}^{(d)} - \mu_{l-1}^{(f)} \right\| + \varepsilon_{d}^{(l)} + \varepsilon_{f}^{(l)} + 2\sqrt{d_{h}} \Delta^{l} \right]$$

This completes the proof.

2 Additional Experimental Results

We further investigated the sensitivity of FairGLite to hyperparameters a and b on the NBA and Pokec-n datasets, as shown in Figure 1. Similar to the observations in Credit and Pokec-z datasets, increasing a consistently enhanced both fairness and prediction performance until reaching a threshold, beyond which improvements stabilized. For hyperparameter b, the results demonstrated three distinct phases: initially, at low values, fairness constraints had negligible effects, reflecting minimal regularization. As b increased to moderate levels, there was a clear improvement in fairness metrics, accompanied by a gradual decline in predictive accuracy, reflecting a stronger regularization impact. Eventually, beyond a certain threshold (approximately e^1 for NBA and e^3 for Pokec-n), further increasing b resulted in stabilized or slightly deteriorated fairness performance.

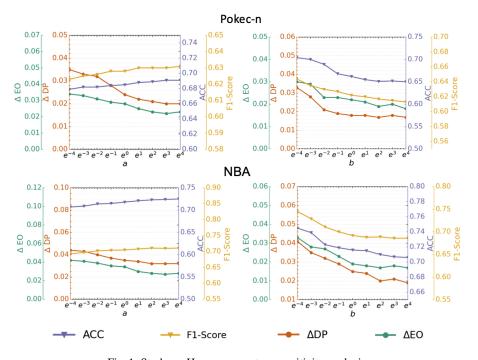


Fig. 1: Study on Hyper-parameters sensitivity analysis.

We further conducted ablation studies to assess the contributions of each FairGLite module on the NBA and Pokec-n datasets, as shown in Figure 2. Consistent with the Credit and Pokec-z datasets observations, removing the Fairness Constraint (FairGLite-NF) reduced model fairness due to biases propagating directly into predictions. The variant without the Graph Reconstruction Constraint (FairGLite-NG) exhibited a less severe drop in fairness but demonstrated reduced predictive accuracy, underscoring the

importance of structural information. Lastly, removing the Adaptivity Confidence Strategy Module (FairGLite-NA) led to decreased overall performance. This reduction illustrates the module's critical function in dynamically weighting fairness constraints based on prediction confidence, effectively balancing predictive accuracy and fairness enforcement.

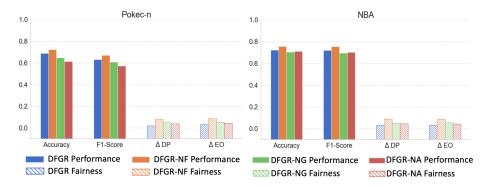


Fig. 2: Ablation study results for FairGLite, FairGLite-NF, FairGLite-NG, and FairGLite-NA.