

Uncovering the Structural Fairness in Graph Contrastive Learning

Ruijia Wang¹, Xiao Wang¹, Chuan Shi¹, Le Song²

¹Beijing University of Posts and Telecommunications

²BioMap and MBZUAI



01 Background

02 Investigation and Analysis

03 **GRADE: The Proposed Model**

04 Experiments

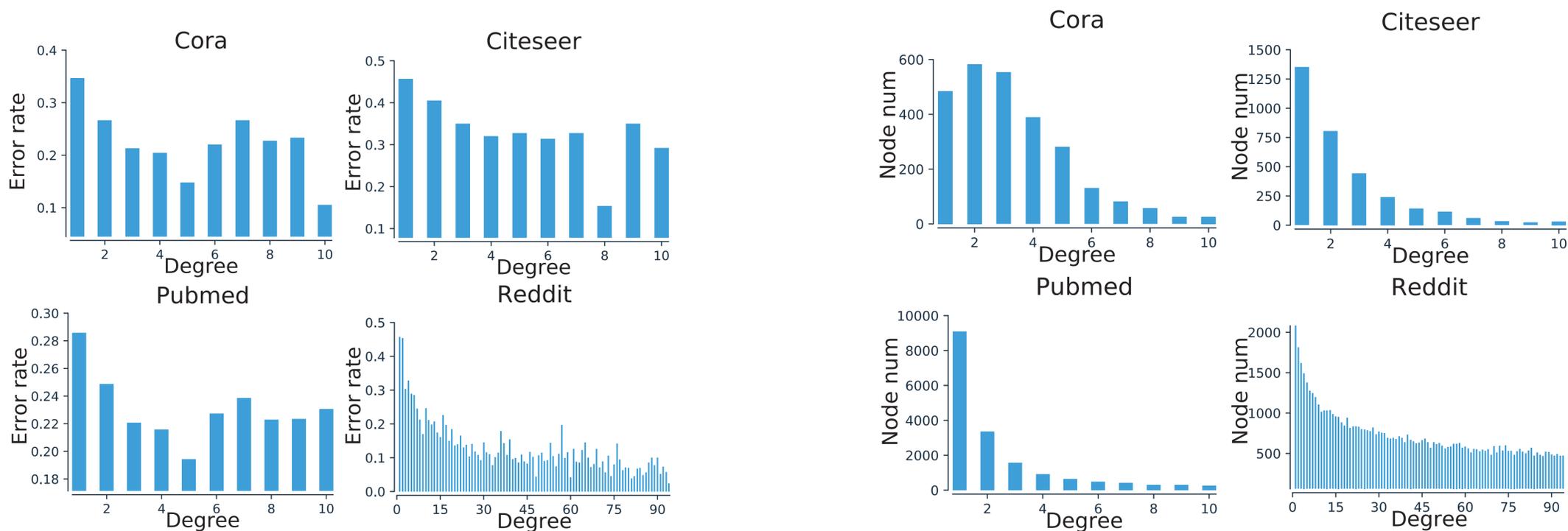
05 Conclusion



Background

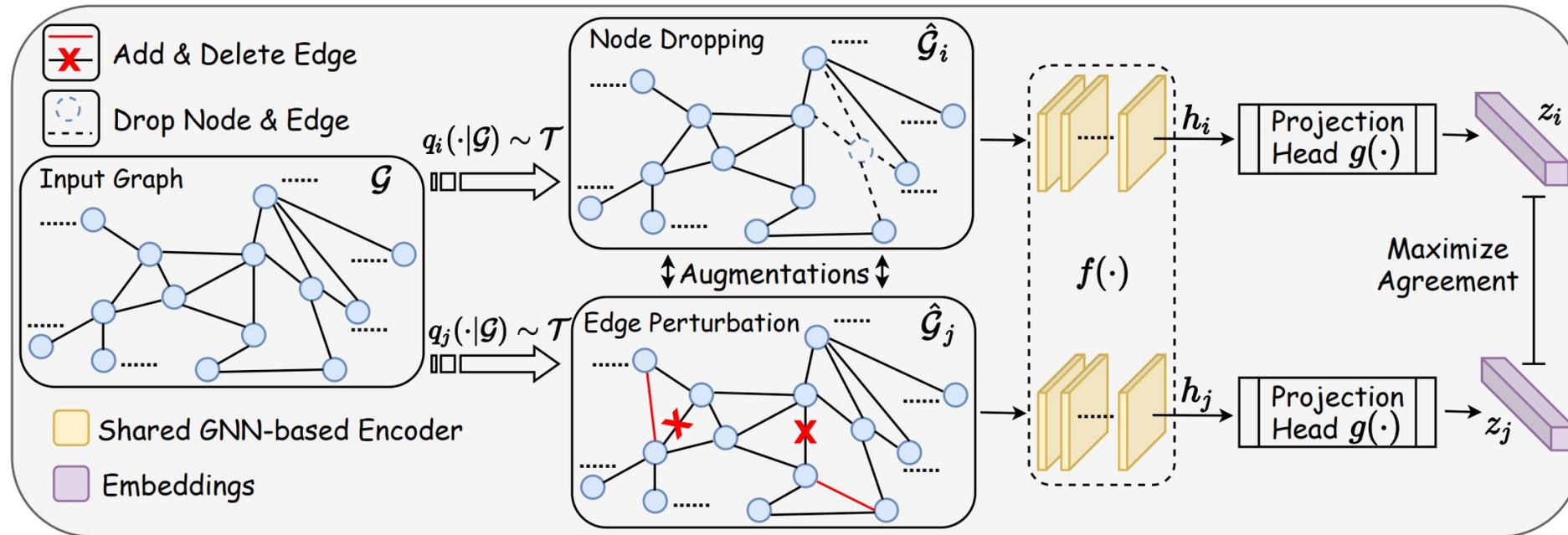
★ Graph Convolutional Network (GCN^[1])

- Node degrees of real-world graphs often follow a **long-tailed distribution**.



GCN exhibits a structural unfairness.

★ Graph Contrastive Learning (GCL)



- GCL integrates the power of GCN and contrastive learning.
- GCL relieves GCN from annotations, and displays SOTA performance in many tasks.

Will GCL present the same structure unfairness as GCN?

Background

Investigation and Analysis

GRADE

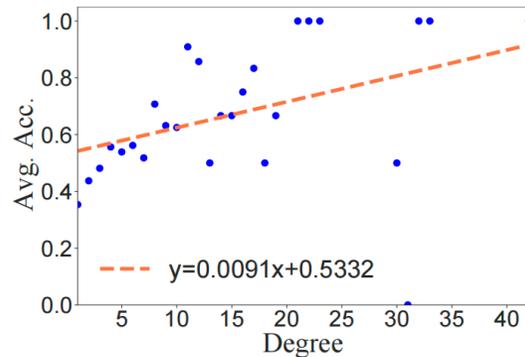
Experiments

Conclusion

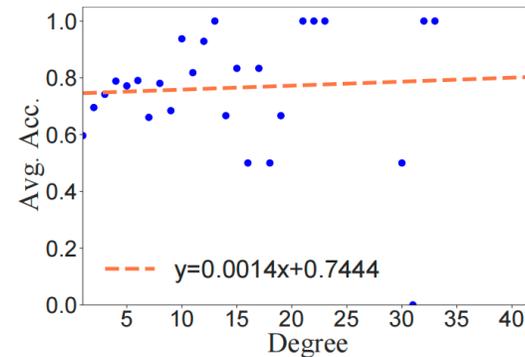


Investigation and Analysis

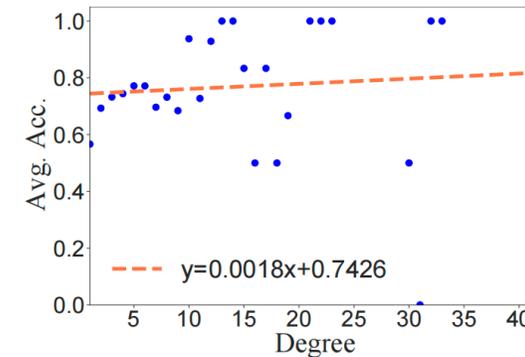
★ Exploring the Behavior of Graph Contrastive Learning on Degree Bias



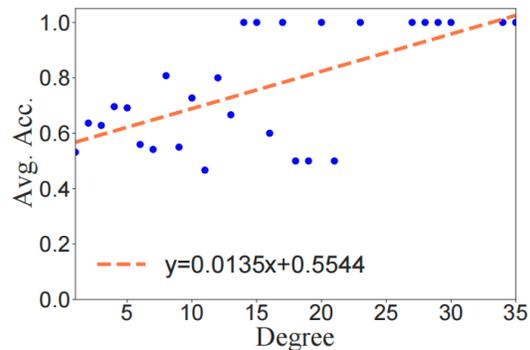
(a) GCN on Cora



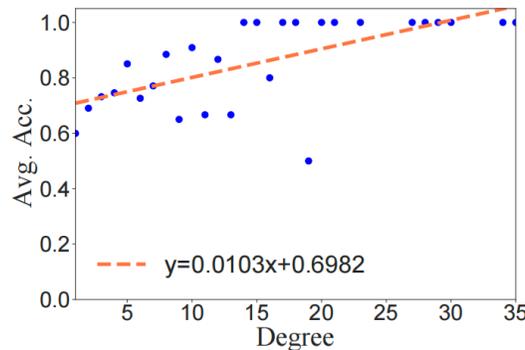
(b) DGI on Cora



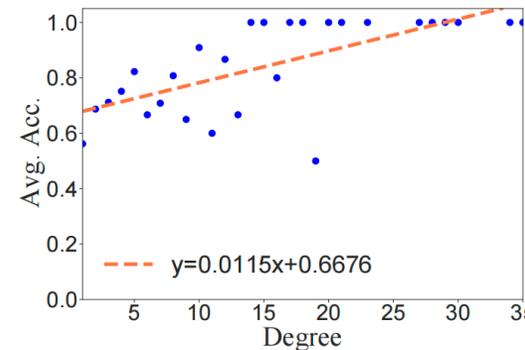
(c) GraphCL on Cora



(d) GCN on Citeseer



(e) DGI on Citeseer



(f) GraphCL on Citeseer

- A **smaller performance gap** exists in GCL methods than that of GCN.

Why is graph contrastive learning fairer to degree bias?

★ Analysis on the Structural Fairness of Graph Contrastive Learning

□ Preliminary Notations

- Let $G = (\mathcal{V}, \mathcal{E}, X)$ be a graph, $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{N \times B}$ is node feature matrix.
- The edges can be represented by an adjacency matrix $A \in \{0, 1\}^{N \times N}$.
- Assume **the augmentation set** \mathcal{T} consisting of all transformations on topology.
- **Positive samples** generated from ego network \mathcal{G}_i of node v_i denoted as $\mathcal{T}(\mathcal{G}_i)$.

- Here we focus on **topological augmentation and single-layer GCN**

$$f(\mathcal{G}_i) = \text{ReLU}(\tilde{L}_i X W) \quad \tilde{L} = \tilde{D}^{-1} \tilde{A}, \quad \tilde{A} = A + I, \quad \tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$$

- We consider a **community indicator** F_f

$$F_f(\mathcal{G}_i) = \arg \min_{k \in [K]} \|f(\mathcal{G}_i) - \mu_k\| \quad \mu_k = \mathbb{E}_{v_i \in C_k} \mathbb{E}_{\hat{\mathcal{G}}_i \in \mathcal{T}(\mathcal{G}_i)} [f(\hat{\mathcal{G}}_i)]$$

- The **error** of community indicator can be formulated as

$$\text{Err}(F_f) = \sum_{k=1}^K \mathbb{P}[F_f(\mathcal{G}_i) \neq k, \forall v_i \in C_k]$$

- we denote $S_\varepsilon = \{v_i \in \cup_{k=1}^K C_k : \forall \hat{\mathcal{G}}_i^1, \hat{\mathcal{G}}_i^2 \in \mathcal{T}(\mathcal{G}_i), \|f(\hat{\mathcal{G}}_i^1) - f(\hat{\mathcal{G}}_i^2)\| \leq \varepsilon\}$ as **nodes with ε -close representations among graph augmentations.**

★ Analysis on the Structural Fairness of Graph Contrastive Learning

□ Theoretical Analysis

- Assume the **nonlinear transformation has M-Lipschitz continuity**

$$\|f(\mathcal{G}_i) - f(\mathcal{G}_j)\| = \|\text{ReLU}(\tilde{L}_i X W) - \text{ReLU}(\tilde{L}_j X W)\| \leq M \|\tilde{L}_i X - \tilde{L}_j X\|$$

- Graph augmentations are **uniformly sampled** with m augmented edges

$$\mathbb{P}[\hat{\mathcal{G}}_i = \mathcal{T}(\mathcal{G}_i)] = 1/C(N - 1, m)$$

- Let there be a ball of **radius** βm such that for any augmentation

$$\|\tilde{L}_i X - \hat{L}_i X\|^2 \leq \beta m$$

Theorem 1 Intra-community Concentration. Let pre-transformation representations $\tilde{L}X$ be sub-Gaussian random variable with variance σ^2 . For all nodes $v_i \in S_\varepsilon$, if $\varepsilon^2 \leq \frac{\beta m}{6M^2\kappa}$, their representations $f(\mathcal{G}_i)$ fit sub-Gaussian distribution with variance $\sigma_{f,\varepsilon}^2 \leq \frac{1}{\kappa}\sigma^2$ with $\kappa \geq 1$ where κ is a coefficient that reflects the degree of concentration.

★ Analysis on the Structural Fairness of Graph Contrastive Learning

□ Theoretical Analysis

- Define the **augmentation distance** between nodes as the minimum distance between their pre-transformation representations

$$d_{\mathcal{T}}(v_i, v_j) = \min_{\hat{\mathcal{G}}_i \in \mathcal{T}(\mathcal{G}_i), \hat{\mathcal{G}}_j \in \mathcal{T}(\mathcal{G}_j)} \|\hat{L}_i X - \hat{L}_j X\| = \min_{\hat{\mathcal{G}}_i \in \mathcal{T}(\mathcal{G}_i), \hat{\mathcal{G}}_j \in \mathcal{T}(\mathcal{G}_j)} \left\| \left(\frac{\hat{A}_i}{\hat{d}_i} - \frac{\hat{A}_j}{\hat{d}_j} \right) X \right\|$$

- Introduce **the definition of $(\alpha, \gamma, \hat{d})$ -augmentation** to measure the concentration of pre-transformation representations

Definition 1 $(\alpha, \gamma, \hat{d})$ -Augmentation. The augmentation set \mathcal{T} is a $(\alpha, \gamma, \hat{d})$ -augmentation, if for each community C_k , there exists a subset $C_k^0 \subset C_k$ such that the following two conditions hold

- $\mathbb{P}[v_i \in C_k^0] \geq \underline{\alpha} \mathbb{P}[v_i \in C_k]$ where $\alpha \in (0, 1]$,
- $\sup_{v_i, v_j \in C_k^0} d_{\mathcal{T}}(v_i, v_j) \leq \gamma \left(\frac{B}{\underline{\hat{d}}_{\min}^k} \right)^{\frac{1}{2}}$ where $\gamma \in (0, 1]$,

where $\hat{d}_{\min}^k = \min_{v_i \in C_k^0, \hat{\mathcal{G}}_i \in \mathcal{T}(\mathcal{G}_i)} \hat{d}_i$, and B is the feature dimension.

★ Analysis on the Structural Fairness of Graph Contrastive Learning

□ Theoretical Analysis

- Assume that the representation is **normalized** by $\|f(\mathcal{G}_i)\| = r$ and let $p_k = \mathbb{P}[v_i \in C_k]$
- **Bound the inter-community distance and the error of the community indicator**

Theorem 2 *Inter-community Scatter.* For a $(\alpha, \gamma, \hat{d})$ -augmentation, if

$$\mu_\ell^\top \mu_k < r^2(1 - \rho_{\max}(\alpha, \gamma, \hat{d}, \varepsilon) - \sqrt{2\rho_{\max}(\alpha, \gamma, \hat{d}, \varepsilon)} - \frac{\Delta_\mu}{2}) \quad (5)$$

holds for any pair of (ℓ, k) with $\ell \neq k$, then the error of the community indicator F_f can be

bounded by $(1 - \alpha) + R_\varepsilon$, where $\rho_{\max}(\alpha, \gamma, \hat{d}, \varepsilon) = 2(1 - \alpha) + \max_\ell \left(\frac{2R_\varepsilon}{p_\ell} + \frac{M\alpha\gamma\sqrt{B}}{r\sqrt{\hat{d}_{\min}^\ell}} \right) + \frac{2\alpha\varepsilon}{r}$

and $\Delta_\mu = 1 - \min_{k \in [K]} \|\mu_k\|^2 / r^2$. **② concentration of representations**

Theorem 3 *The term R_ε is upper bounded by*

$$R_\varepsilon \leq \frac{[C(N-1, m)]^2}{\varepsilon} \mathbb{E}_{v_i} \mathbb{E}_{\hat{\mathcal{G}}_i^1, \hat{\mathcal{G}}_i^2 \in \mathcal{T}(\mathcal{G}_i)} \|f(\hat{\mathcal{G}}_i^1) - f(\hat{\mathcal{G}}_i^2)\|. \quad (6)$$

① alignment of positive pairs

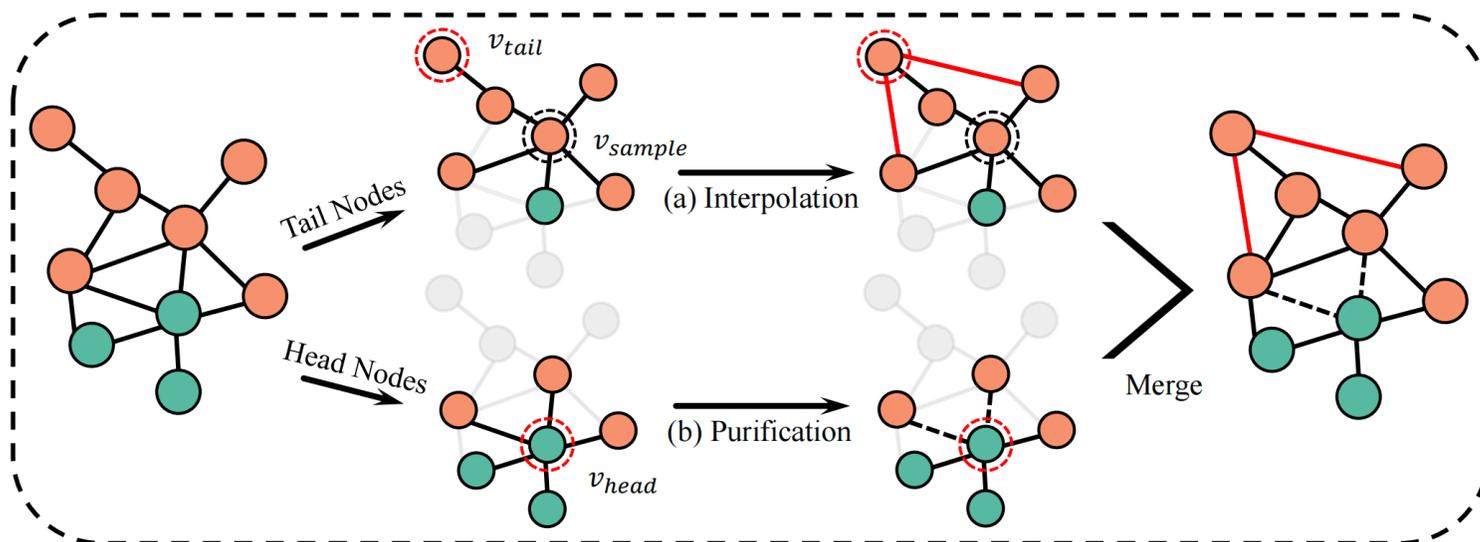
GCL conform to a clearer community structure



The Proposed Model

★ Overview

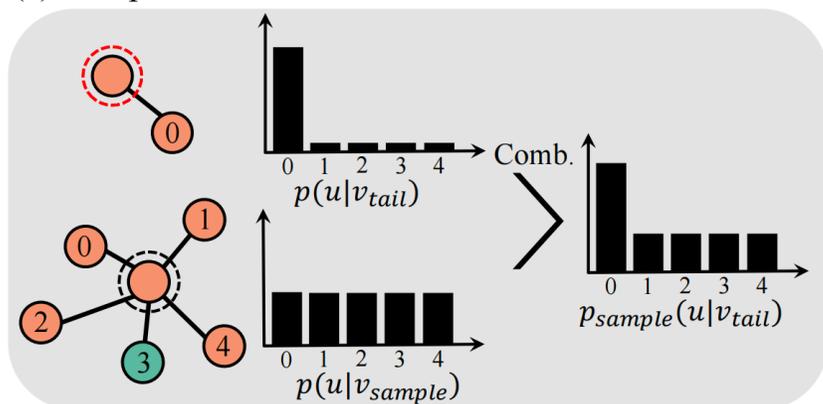
Aim to increase intra-community edges while decreasing inter-community edges



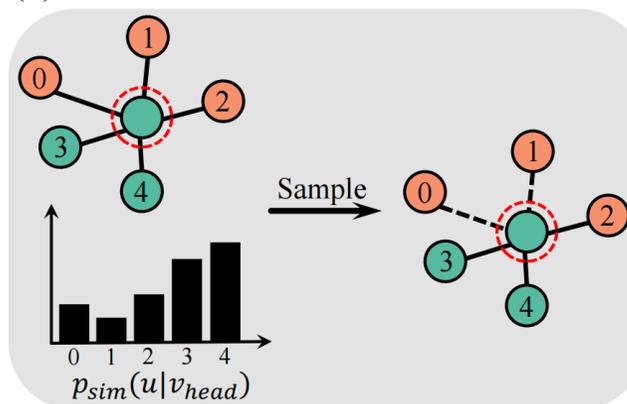
□ Tail nodes

- **Interpolate the ego network** of the anchor tail node with that of a similar node.

(a) Interpolation



(b) Purification



□ Head nodes

- **Purify the neighborhood** by similarity-based sampling.

★ Graph Augmentation

□ Topology Augmentation

- We build the similarity matrix S based on **cosine similarity** of representations
$$S_{ij} = \text{sim}(\mathbf{h}_i, \mathbf{h}_j) \text{ for } i \neq j \text{ and } S_{ii} = 0 \text{ otherwise}$$
- For any tail node v_{tail} , we sample a node v_{sample} from the distribution $\text{Multi}(\mathbf{s}_{tail})$.
- The similarity $\text{sim}(\mathbf{h}_{tail}, \mathbf{h}_{sample})$ is used as the **interpolation** ratio ϕ
$$p_{sample}(u|v_{tail}) = \phi p(u|v_{tail}) + (1 - \phi)p(u|v_{sample})$$
- For each head node v_{head} , we define the **similarity distribution for purification**
$$p_{sim}(u|v) = \text{sim}(\mathbf{h}_u, \mathbf{h}_v) \text{ if } u \in \mathcal{N}(v) \text{ and } p(u|v) = 0$$
- We sample $d_{head}(1 - p_{edr})$ neighbors without replacement.

□ Feature Augmentation

- We randomly sample a **mask** $\mathbf{m} \in \{0, 1\}^B$ from a Bernoulli distribution $\text{Ber}(1 - p_{fdr})$
$$\hat{X} = [\mathbf{x}_1 \circ \mathbf{m}, \mathbf{x}_2 \circ \mathbf{m}, \dots, \mathbf{x}_N \circ \mathbf{m}]$$

★ Optimization Objective

- Node representations \mathbf{h}_i and \mathbf{o}_i from different graph augmentations form the positive pair.
- Node representations of other nodes in graph augmentations are regarded as negative pairs.
- We define the **pairwise objective for each positive pair^[1]** $(\mathbf{h}_i, \mathbf{o}_i)$ as

$$\ell(\mathbf{h}_i, \mathbf{o}_i) = \log \frac{e^{\theta(\mathbf{h}_i, \mathbf{o}_i)/\tau}}{e^{\theta(\mathbf{h}_i, \mathbf{o}_i)/\tau} + \sum_{k \neq i} e^{\theta(\mathbf{h}_i, \mathbf{o}_k)/\tau} + \sum_{k \neq i} e^{\theta(\mathbf{h}_i, \mathbf{h}_k)/\tau}}$$

where τ is a temperature parameter, and the critic $\theta(\mathbf{h}, \mathbf{o})$ is defined as $\text{sim}(g(\mathbf{h}), g(\mathbf{o}))$.

- The overall objective to be maximized is **the average of all positive pairs**

$$\mathcal{J} = \frac{1}{2N} \sum_{i=1}^N [\ell(\mathbf{h}_i, \mathbf{o}_i) + \ell(\mathbf{o}_i, \mathbf{h}_i)]$$

Background

Investigation and Analysis

GRADE

Experiments

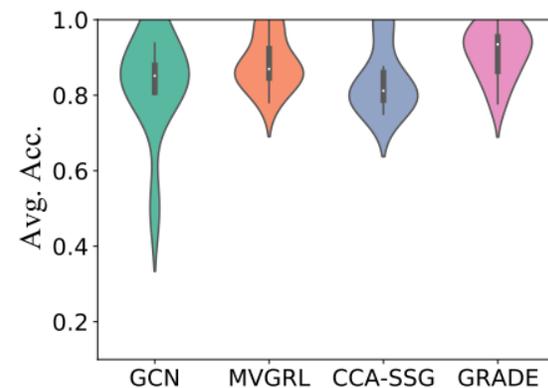
Conclusion



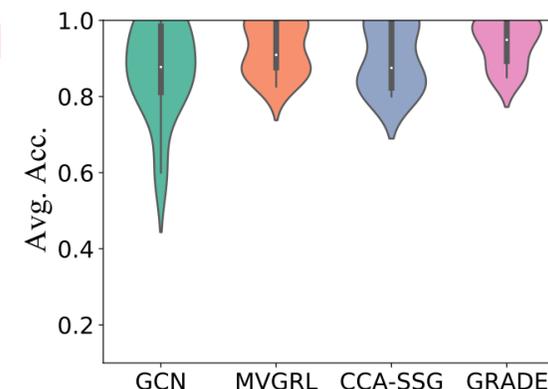
Experiments

★ Node Classification

	Cora		Citeseer		Photo		Computer	
	<i>Micro-F1</i>	<i>Macro-F1</i>	<i>Micro-F1</i>	<i>Macro-F1</i>	<i>Micro-F1</i>	<i>Macro-F1</i>	<i>Micro-F1</i>	<i>Macro-F1</i>
	Supervised Split							
GCN	82.30 \pm 0.49	76.87 \pm 0.34	65.84 \pm 0.55	59.62 \pm 0.64	93.52 \pm 0.82	78.88 \pm 2.01	89.14 \pm 0.75	72.61 \pm 3.05
DGI	82.28 \pm 0.84	77.23 \pm 0.90	65.64 \pm 0.63	59.47 \pm 1.24	92.98 \pm 1.12	78.83 \pm 1.66	88.96 \pm 0.96	72.30 \pm 1.80
GraphCL	81.78 \pm 0.67	76.01 \pm 1.07	65.16 \pm 1.02	58.72 \pm 1.37	—	—	—	—
GRACE	82.32 \pm 0.45	76.78 \pm 0.87	64.16 \pm 2.07	59.73 \pm 1.94	93.12 \pm 0.40	78.60 \pm 3.12	88.22 \pm 1.04	71.74 \pm 3.05
MVGRL	83.22 \pm 1.02	77.84 \pm 1.35	66.26 \pm 0.72	60.30 \pm 0.95	94.10 \pm 0.31	78.36 \pm 2.22	—	—
CCA-SSG	82.70 \pm 0.86	77.35 \pm 1.06	65.96 \pm 1.36	58.81 \pm 1.67	94.36 \pm 0.25	79.34 \pm 3.42	89.22 \pm 0.95	73.82 \pm 1.80
GRADE	83.40 \pm 0.80	78.54 \pm 1.15	67.14 \pm 1.07	61.04 \pm 2.07	94.72 \pm 0.30	78.86 \pm 2.77	89.42 \pm 0.53	74.71 \pm 1.30
Semi-supervised Split								
GCN	74.18 \pm 0.40	69.84 \pm 0.56	53.80 \pm 0.94	50.15 \pm 0.69	91.04 \pm 0.65	65.47 \pm 1.20	78.58 \pm 0.93	61.80 \pm 1.43
DGI	75.92 \pm 0.86	70.04 \pm 0.53	54.52 \pm 1.44	51.92 \pm 1.23	90.78 \pm 0.78	66.27 \pm 0.76	79.00 \pm 0.80	62.00 \pm 1.70
GraphCL	75.68 \pm 2.84	69.86 \pm 2.41	54.06 \pm 1.93	51.75 \pm 1.78	—	—	—	—
GRACE	75.12 \pm 1.41	69.66 \pm 1.29	53.56 \pm 3.42	49.83 \pm 1.74	91.12 \pm 0.31	65.07 \pm 1.28	79.10 \pm 1.79	61.76 \pm 1.97
MVGRL	76.44 \pm 1.17	70.52 \pm 1.63	56.84 \pm 1.26	53.79 \pm 1.25	92.01 \pm 0.87	66.16 \pm 2.13	—	—
CCA-SSG	75.74 \pm 1.96	71.70 \pm 1.59	57.90 \pm 1.82	54.70 \pm 1.54	91.68 \pm 0.50	67.08 \pm 1.08	82.20 \pm 0.47	65.04 \pm 1.16
GRADE	77.20 \pm 0.94	73.37 \pm 1.27	59.44 \pm 0.78	56.47 \pm 0.64	92.04 \pm 0.30	66.62 \pm 2.27	82.50 \pm 1.04	67.50 \pm 1.80



(a) Tail nodes



(b) Head nodes

GRADE outperforms all baselines in most cases regardless of tail nodes or head nodes.

★ Fairness Analysis

- We define the **group mean** as the mean of degree-specific average accuracy
- The **bias** is the variance.

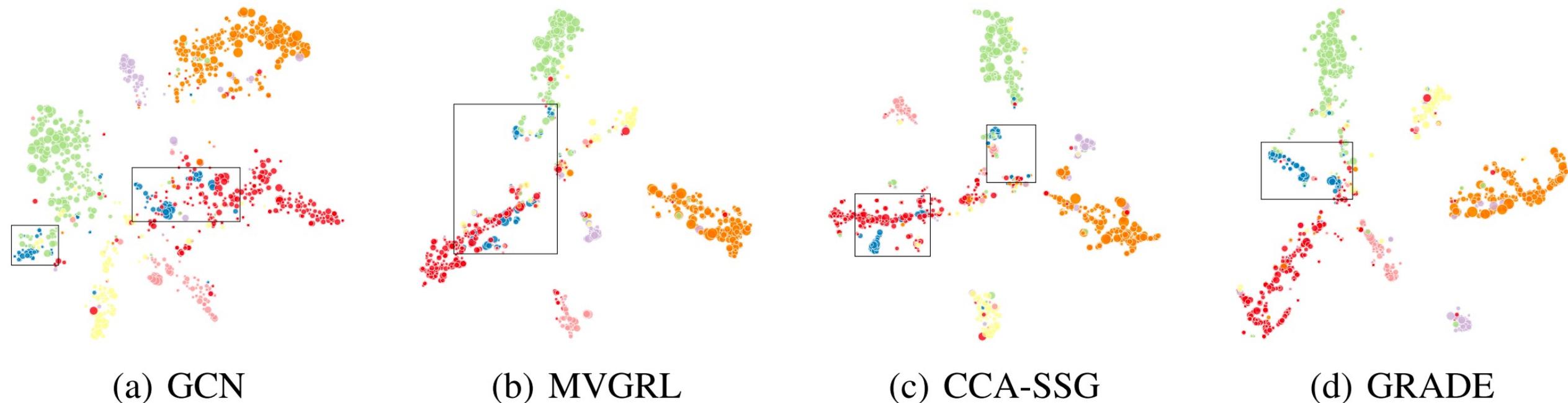
$$\text{Avg. Acc.}(k) = \mathbb{E}[\{\text{Acc}(v_i), \forall \text{ node } v_i \text{ such that } d_i = k\}],$$

$$G.Mean = \mathbb{E}[\{\text{Avg. Acc.}(k), \forall \text{ node degree } k\}], \text{Bias} = \text{Var}(\{\text{Avg. Acc.}(k), \forall \text{ node degree } k\})$$

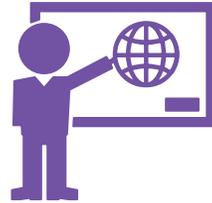
	Cora		Citeseer		Photo		Computer	
	<i>G. Mean</i> ↑	<i>Bias</i> ↓						
GCN	86.04	1.70	84.00	1.85	97.41	0.28	96.30	0.50
DGI	89.26	0.67	84.79	1.71	98.23	0.27	96.94	0.45
GraphCL	90.80	0.59	84.13	1.80	—	—	—	—
GRACE	89.91	0.70	85.44	1.67	98.28	0.23	96.92	0.47
MVGRL	91.01	0.54	83.86	1.83	98.39	0.27	—	—
CCA-SSG	90.86	0.63	84.35	1.73	98.44	0.24	97.17	0.39
GRADE	92.87	0.48	85.88	1.52	98.52	0.20	97.42	0.35

GRADE reduces the bias across all datasets and maintain the highest group mean.

★ Visualization



GRADE pulls same-community node representations more concentrated.



Conclusions

01 **New insights for structural fairness**

We are the first to discover that GCL methods exhibit more structural fairness than GCN. This discovery inspires a new path for alleviating structural unfairness based on contrastive learning.

02 **Deeper understanding for graph contrastive learning**

We theoretically validate the reason for structural fairness in GCL is that it stimulates intra-community concentration.

03 **A novel framework**

We propose a method GRADE to further improve the structural fairness by enriching the neighborhood of tail nodes while purifying neighbors of head nodes.



Thanks for listening!

E-mail: wangruijia@bupt.edu.cn

code & data: <http://shichuan.org/>