



TEXAS A&M UNIVERSITY
Engineering



Fairness in Machine Learning: Metrics and Algorithms

Presenter: Zhimeng Jiang

Date: 10/19/2023

About Me

- Staff Research Scientist at Visa Research



- Just defended my thesis last week at Texas A&M University



- Research Interests

- Trustworthy ML
- Graph Neural Networks
- Generative AI

- Awards

- ICML 2022 Outstanding Paper Award
- CIKM 2022 Best Demo Paper Award



Algorithmic Bias in ML Systems

- Ethical challenges posed by ML systems
- Inherent bias presented in society
 - Reflected in training data
 - AI/ML models prone to amplifying such bias



COOKING

ROLE	VALUE
AGENT	▶ WOMAN
FOOD	▶ PASTA
HEAT	▶ STOVE
TOOL	▶ SPATULA
PLACE	▶ KITCHEN



COOKING

ROLE	VALUE
AGENT	▶ WOMAN
FOOD	▶ FRUIT
HEAT	▶ –
TOOL	▶ KNIFE
PLACE	▶ KITCHEN



COOKING

ROLE	VALUE
AGENT	▶ WOMAN
FOOD	▶ MEAT
HEAT	▶ GRILL
TOOL	▶ TONGS
PLACE	▶ OUTSIDE



COOKING

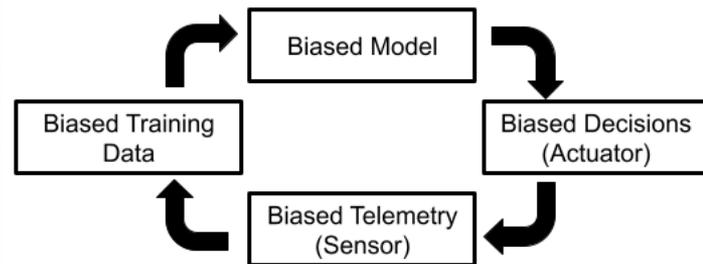
ROLE	VALUE
AGENT	▶ WOMAN
FOOD	▶ VEGETABLES
HEAT	▶ STOVE
TOOL	▶ TONGS
PLACE	▶ KITCHEN



COOKING

ROLE	VALUE
AGENT	▶ MAN
FOOD	▶ –
HEAT	▶ STOVE
TOOL	▶ SPATULA
PLACE	▶ KITCHEN

Image recognition systems



Feedback loop
Image from Medium: [link](#)

The Consequence of ML Systems' Bias

Impact of data bias on business



Source: DataRobot

- 62% lost revenue
- 61% lost customers
- 43% lost employees
- 35% incurred legal fees due to lawsuit
- 6% lost customer trust

Stacite

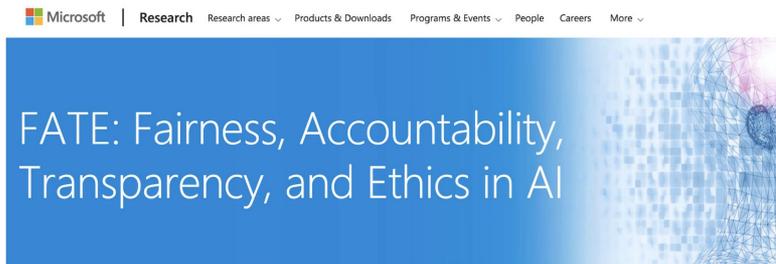


Business Impact

Microsoft's Tay AI chatbot

Bias increases the financial loss and production risk!

Ethical Regulation for ML Systems



Research ▾ Product ▾ Developers ▾ Safety Company ▾

We've created GPT-4, the latest milestone in OpenAI's effort in scaling up deep learning. GPT-4 is a large multimodal model (accepting image and text inputs, emitting text outputs) that, while less capable than humans in many real-world scenarios, exhibits human-level performance on various professional and academic benchmarks. For example, it passes a simulated bar exam with a score around the top 10% of test takers; in contrast, GPT-3.5's score was around the bottom 10%. **We've spent 6 months iteratively aligning GPT-4** using lessons from our adversarial testing program as well as ChatGPT, resulting in our best-ever results (though far from perfect) on factuality, steerability, and refusing to go outside of guardrails.

AI at Google: our principles

We will assess AI applications in view of the following objectives. We believe that AI should:

1. Be socially beneficial.

The expanded reach of new technologies increasingly touches society as a whole. AI has the potential for transformative impacts in a wide range of fields, including healthcare, security, energy, transportation, manufacturing, and entertainment. As we consider potential development and uses of AI technologies, we will take into account a broad range of social and economic factors, and will proceed where we believe that the overall likely benefits substantially exceed the foreseeable risks and downsides.

AI also enhances our ability to understand the meaning of content at scale. We will strive to make high-quality and accurate information readily available using AI, while continuing to respect cultural, social, and legal norms in the countries where we operate. And we will continue to thoughtfully evaluate when to make our technologies available on a non-commercial basis.

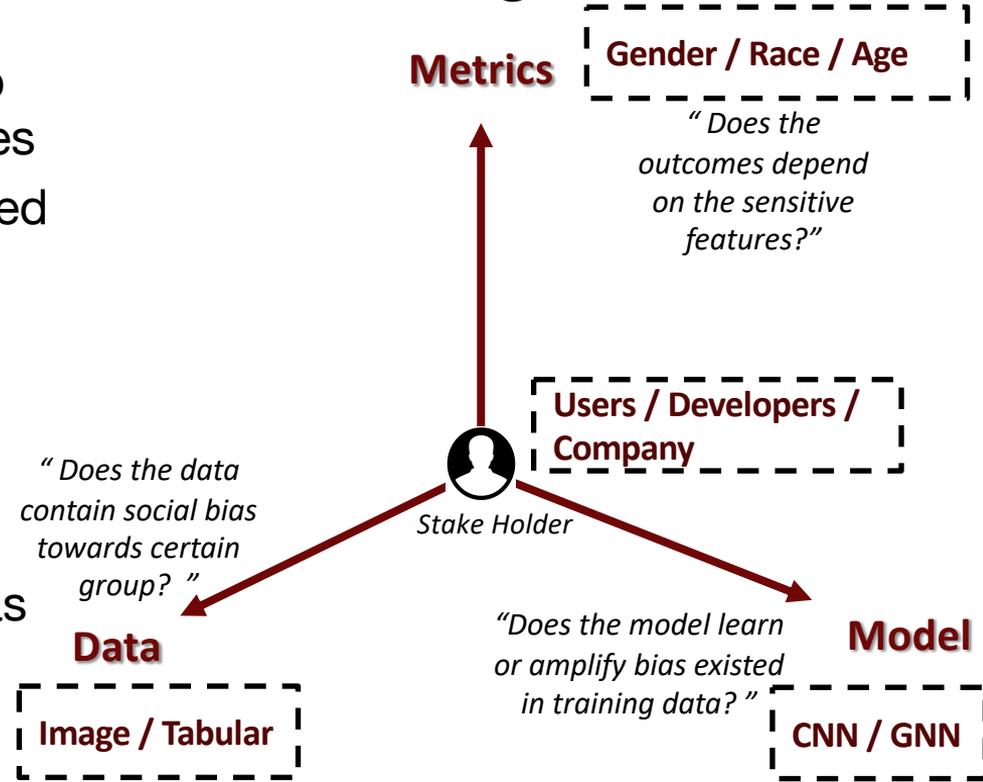
2. Avoid creating or reinforcing unfair bias.

AI algorithms and datasets can reflect, reinforce, or reduce unfair biases. We recognize that distinguishing fair from unfair biases is not always simple, and differs across cultures and societies. We will seek to avoid unjust impacts on people, particularly those related to sensitive characteristics such as race, ethnicity, gender, nationality, income, sexual orientation, ability, and political or religious belief.

Fairness in Machine Learning

Goal: Develop ML/AI systems that help make decisions leading to fair outcomes

- **Metrics:** Evaluate outcome bias based on protected attributes
- **Data:** human bias leading to biased training data
- **Model:** ML model even amplifies bias during training

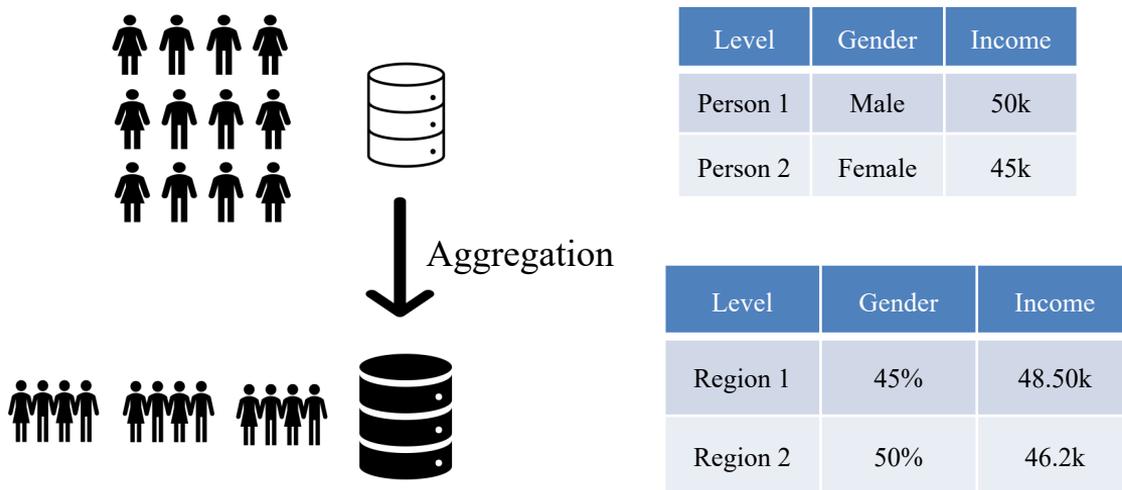


Overview

- [A] Generalized Demographic Parity for Group Fairness, ICLR'22
- [B] Learning fair graph representations via automated data augmentations, ICLR'23
- [C] Fair Graph Message Passing, under review
- [D] Topology Matters in Fair Graph Learning: a Theoretical Pilot Study, under review
- [E] Chasing Fairness under Distribution Shift: a Model Weight Perturbation Approach, NeurIPS'23

Generalized Demographic Parity [A]

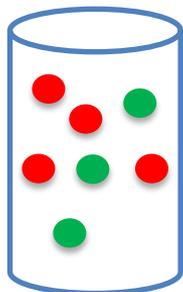
- Existing group fairness metrics are either inapplicable for continuous sensitive attribute or without tractable computation.



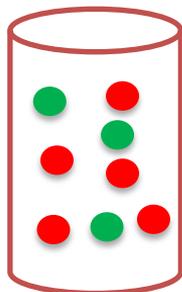
Observation: Data aggregation transforms binary sensitive attribute into **continuous** attributes

Why need a new fairness metric?

- Existing group fairness can not be applicable to continuous sensitive attribute
 - Demographic parity: the same average positive prediction rate among demographic group
 - Continuous sensitive attribute: infinite demographic groups
- Existing fairness metrics for continuous attributes are computation-intractable
 - Mutual information
 - Metric estimation methods either rely on tractable bounds or neural network approximation
 - Therefore insufficiently trustful for algorithm ranking



Male



Female

MINE[1]: NN replaces optimization

$$I(X; Z) \geq I_{\Theta}(X, Z),$$

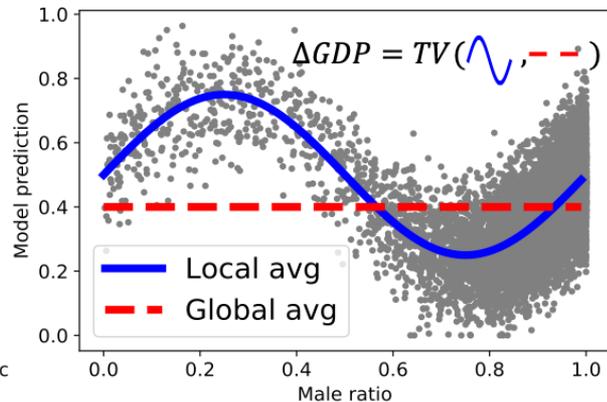
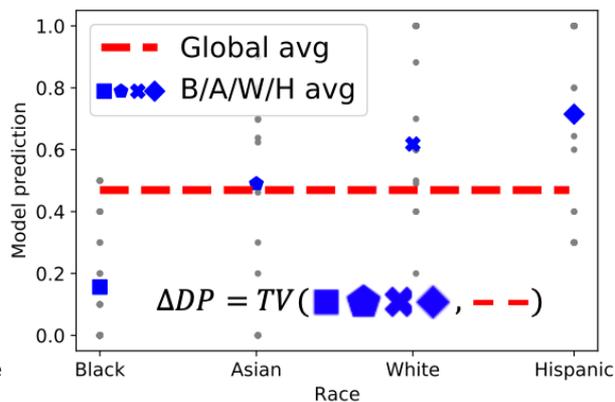
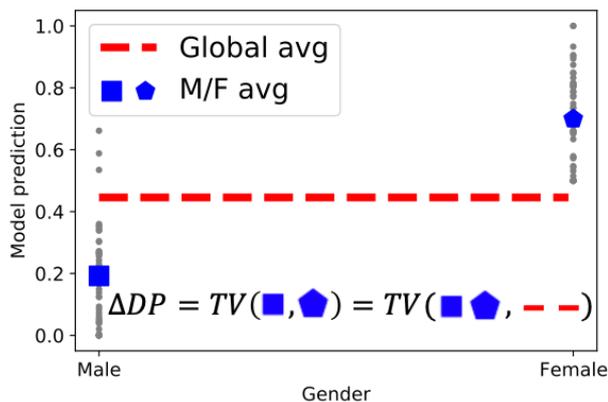
$$I_{\Theta}(X, Z) = \sup_{\theta \in \Theta} \mathbb{E}_{\mathbb{P}_{XZ}} [T_{\theta}] - \log(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Z} [e^{T_{\theta}}]).$$

Find by NN

[1] Mohamed Ishmael Belghazi, et al. “Mutual information neural estimation Mutual information neural estimation” ICML, 2018.

GDP Overview

- Demographic parity (DP) [2]: binary sensitive attribute
- Difference w.r.t. DP (DDP) [3]: categorical sensitive attribute
- Generalized DP (**GDP**): general version for binary/categorical/continuous sensitive attribute
 - local/global difference
 - Local average: average prediction given specific sensitive attribute



[2] Feldman, Michael, et al. "Certifying and removing disparate impact." proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining. 2015.

[3] Cho, Jaewoong, et al. "A fair classifier using kernel density estimation." Advances in Neural Information Processing Systems 33 (2020): 15088-15099.

GDP Justifications

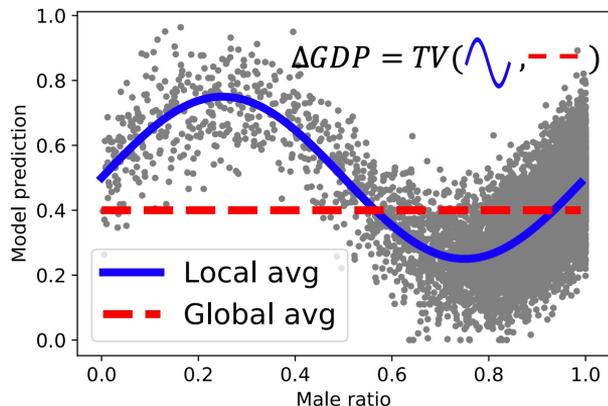
- GDP is a natural extension of DP/DDP for continuous attribute
 - GDP and DP are equivalent except for the dataset-dependent coefficient for binary attribute.
 - GDP is weighted DDP for categorical attributes.
- GDP understanding from a probabilistic view
 - Idea case: prediction \perp sensitive attribute
 - Joint distribution = Product marginal distribution
 - GDP is a necessary condition for independence
 - $GDP \leq TV \text{ distance}(\text{joint}, \text{product margin})$
- GDP regularizer v.s. adversarial debiasing
 - Adversarial debiasing leads to lower GDP

$$\mathcal{L}_{adv} \left(\boxed{g^*}(f(X)), S \right) \geq \Delta GDP.$$

Adversary: Predict sensitive attribute based on NN outputs

GDP Estimation

- Histogram estimation
 - Hard group: consecutive, non-overlapping intervals
 - Internal group average as local average
 - Estimation error v.s #samples: $Err_{hist} = O(N^{-\frac{2}{3}})$
- Kernel estimation
 - Soft group: closer attribute pair, higher weight
 - Normalized weighted average (Nadaraya–Watson kernel estimator)
 - Estimation error v.s #samples: $Err_{kernel} = O(N^{-\frac{4}{5}})$



$$\tilde{m}^h(s) = \frac{\sum_{n=1}^N \hat{y}_n K\left(\frac{s_n - s}{h}\right)}{\sum_{n=1}^N K\left(\frac{s_n - s}{h}\right)},$$

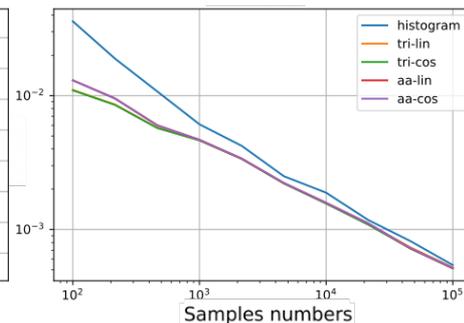
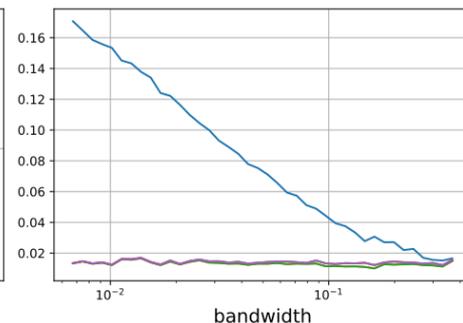
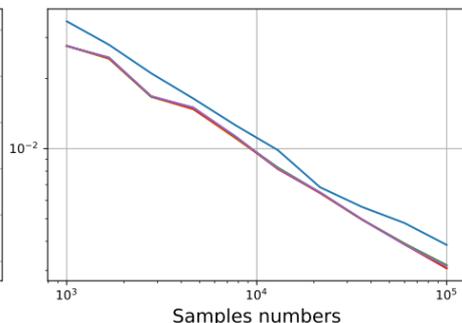
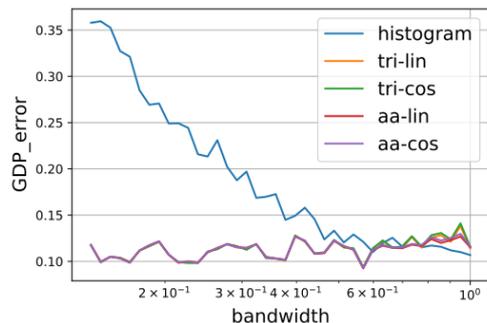
$$\tilde{m}_{avg}^h = \frac{\sum_{n=1}^N \hat{y}_n}{N}.$$

$$\tilde{\Delta GDP}(h) = \int_0^1 \left| \tilde{m}^h(s) - \tilde{m}_{avg}^h \right| \tilde{p}_S^h(s) ds.$$

Synthetic Experiments

GDP estimation error w.r.t. bandwidth/kernel/#samples

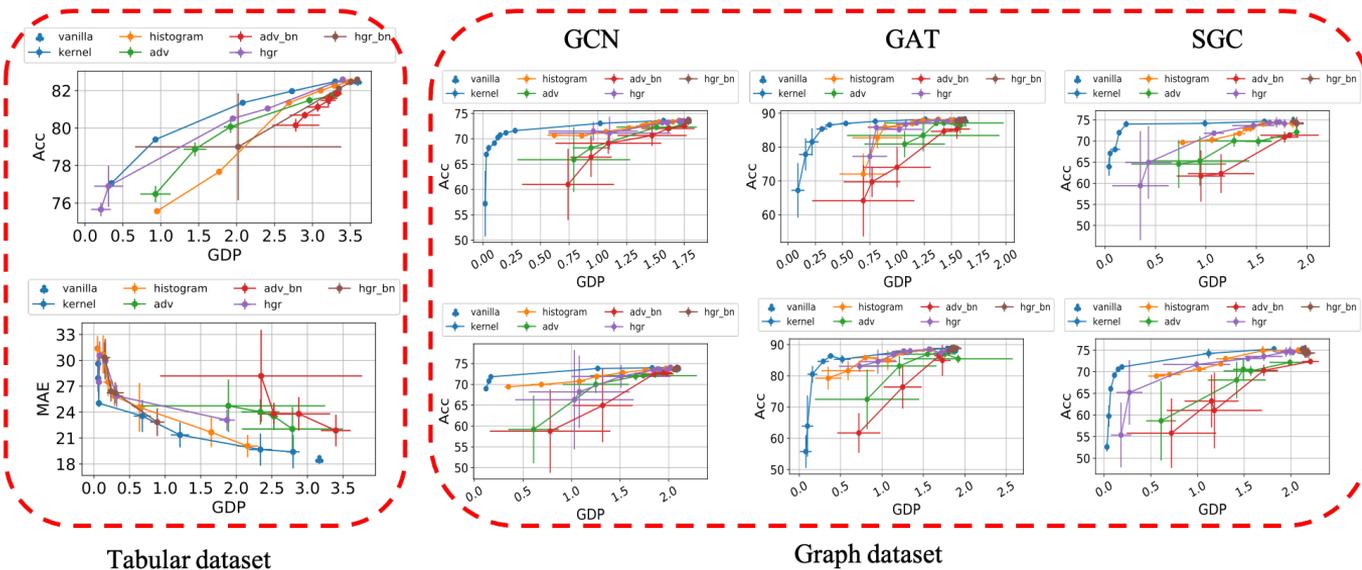
- Kernel estimation is robust over kernel/bandwidth
- Estimation error convergence rate
 - Kernel estimation > histogram estimation



Experiments on Real data

GDP regularizer achieves the best fairness-prediction trade-off performance

- Tabular/graph/dynamic graph data
- Classification/regression tasks
- Single/compositional sensitive attributes



Learning fair graph representations via automated data augmentations [B]

- Data Augmentation to mitigate data bias

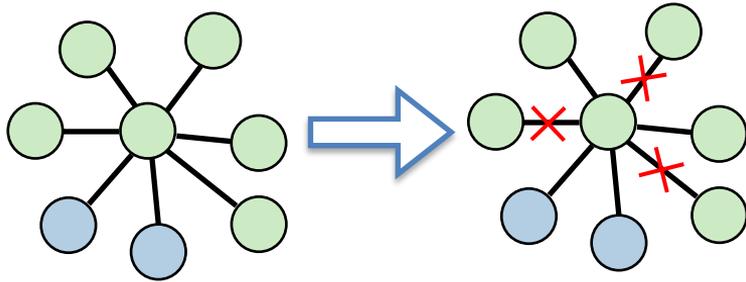


- Existing fairness-aware graph data augmentations
 - Heuristic graph properties that are beneficial to fair representation learning
 - Not optimal from data augmentation perspective.

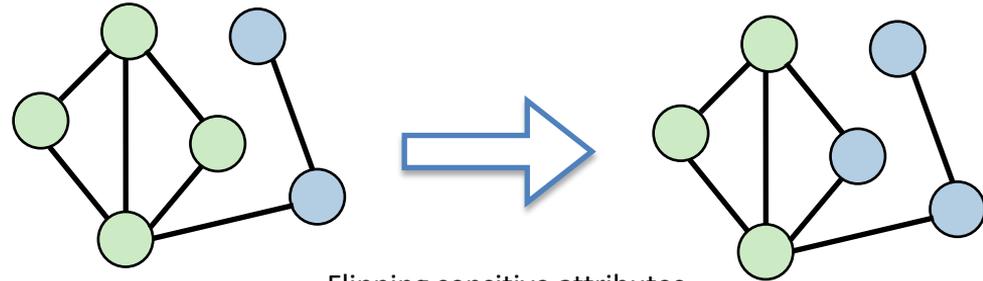
Automated Data Augmentation!

Heuristic Examples

Balanced inter/intra edges



NIFTY [4]



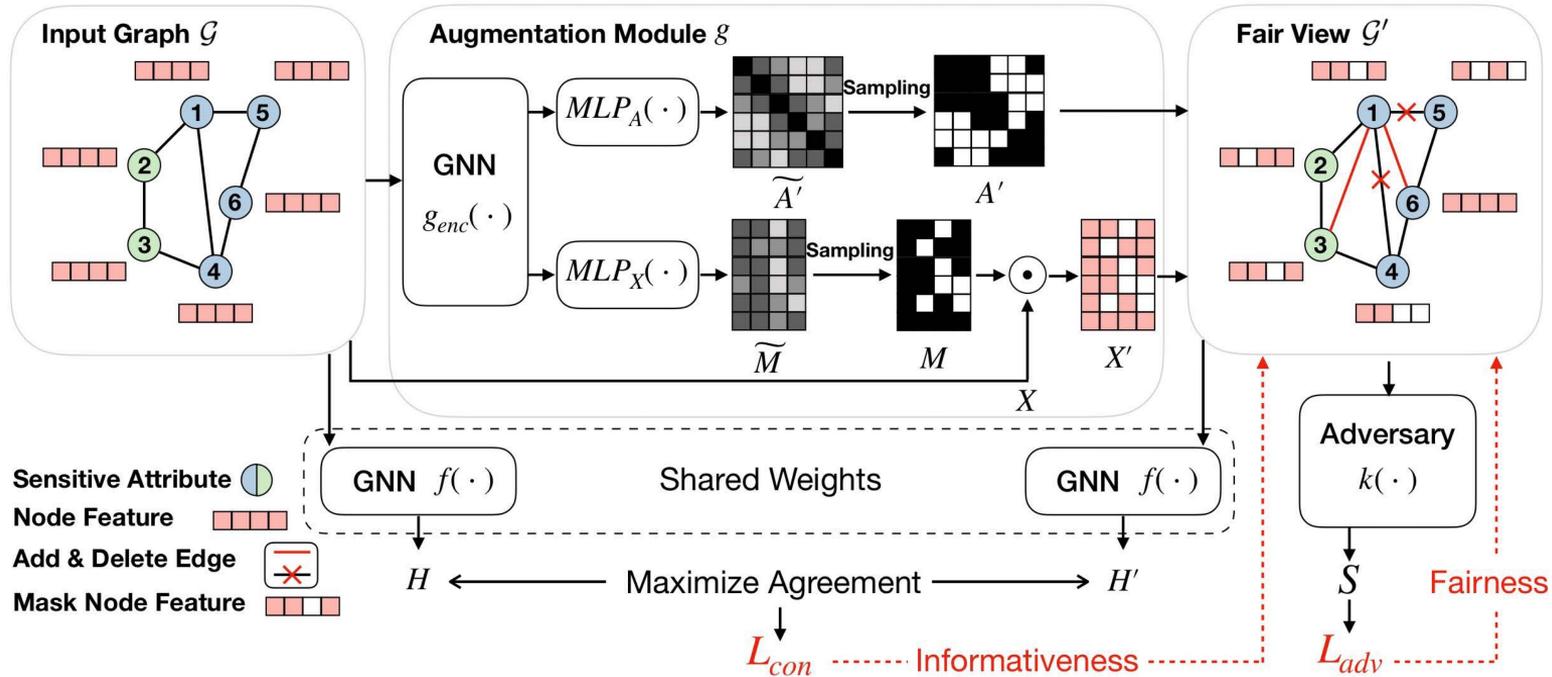
Flipping sensitive attributes

[4] Agarwal, Chirag, Himabindu Lakkaraju, and Marinka Zitnik. "Towards a unified framework for fair and stable graph representation learning." UAI 2021.

Graphair

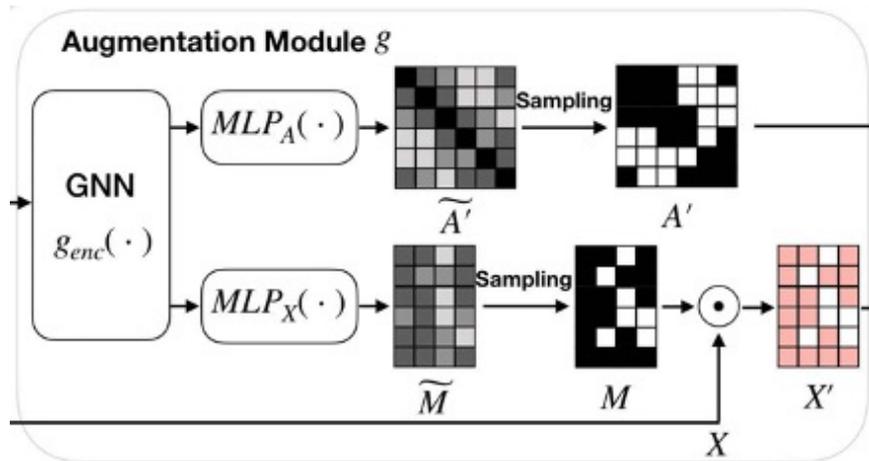
- Graphair is a novel automated graph augmentation method for fair graph representation learning
- Use two types of graph transformations:
 - Edge perturbation: removing existing edges and adding new edges
 - Node feature masking: setting some node features to zeros
- End-to-end training with multiple optimization objectives
 - Fairness: reducing the bias in the generated graphs
 - Informativeness: preserving the most informative components of the input graph in the generated graphs

Overview



Augmentation Module

- Graphair generate the fair view by a learnable augmentation module
 - Extract node embeddings by a GNN model
 - Predict probabilities for each graph element (i.e., node features and edge) by MLP models
 - Sample a new graph based on the predicted probabilities



Fairness with Adversarial Training

- Supervised training is impossible as there are no ground truths indicating which graph elements lead to prediction bias and should be modified
- Achieve fairness via adversarial learning
 - Use an adversary model to predict the sensitive attribute from graph generated by augmentation module
 - The adversary is optimized to maximize the prediction accuracy
 - The augmentation module is optimized to mitigate bias, so that it is difficult for the adversary model to identify sensitive attribute information

$$\min_g \max_k L_{\text{adv}} = \min_g \max_k \frac{1}{n} \sum_{i=1}^n \left[S_i \log \hat{S}_i + (1 - S_i) \log (1 - \hat{S}_i) \right]$$

Informativeness via Contrastive Training

- Only using the adversarial training may cause the augmentation module to collapse into trivial solutions
 - Ex. A complete graph with all zero node features
- Achieve informativeness by using contrastive learning
 - Maximizing the similarity between the representations of the same node

$$l(h_i, h'_i) = -\log \frac{\exp(\text{sim}(h_i, h'_i)/\tau)}{\sum_{j=1}^n \exp(\text{sim}(h_i, h'_j)/\tau) + \sum_{j=1}^n \mathbb{1}_{[j \neq i]} \exp(\text{sim}(h_i, h_j)/\tau)},$$

$$L_{\text{con}} = \frac{1}{2n} \sum_{i=1}^n [l(h_i, h'_i) + l(h'_i, h_i)].$$

Experiments

- RQ1: Does Graphair outperform state-of-the-art methods in terms of fairness while maintaining comparable accuracy?
- RQ2: Does Graphair has a better trade-off between accuracy and fairness than state-of-the-art methods?
- Three real-world datasets
 - NBA, Pokec-z and Pokec-n

Dataset	NBA	Pokec-z	Pokec-n
# Nodes	403	67,797	66,569
# Node features	39	59	59
# Edges	16,570	882,765	729,129
# Inter-group edges	4,401	39,804	31,515
# Intra-group edges	12,169	842,961	697,614

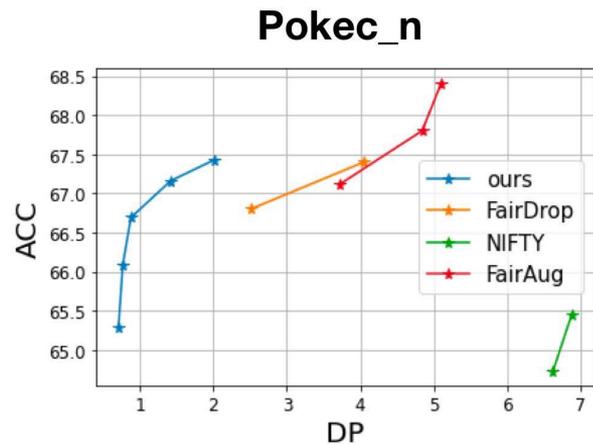
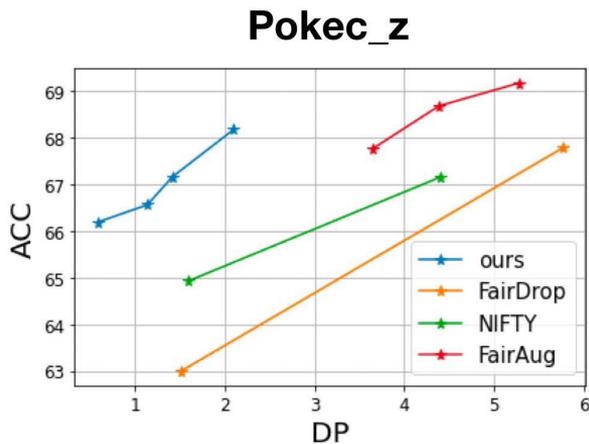
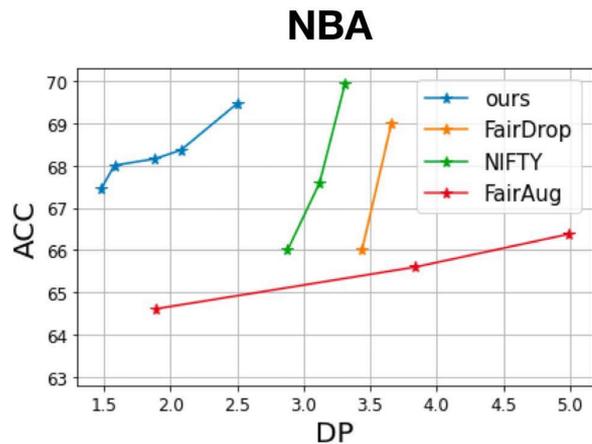
RQ1: Experiments on Fairness Performance

- Our proposed Graphair consistently achieves the best fairness performance in terms of demographic parity and equal opportunity

Method	NBA			Pokec-z			Pokec-n		
	ACC \uparrow	Δ_{DP} \downarrow	Δ_{EO} \downarrow	ACC \uparrow	Δ_{DP} \downarrow	Δ_{EO} \downarrow	ACC \uparrow	Δ_{DP} \downarrow	Δ_{EO} \downarrow
FairWalk	64.54 \pm 2.35	3.67 \pm 1.28	9.12 \pm 7.06	67.07 \pm 0.24	7.12 \pm 0.74	8.24 \pm 0.75	65.23 \pm 0.78	4.45 \pm 1.25	4.59 \pm 0.86
FairWalk+ X	69.74 \pm 1.71	14.61 \pm 4.98	12.01 \pm 5.38	69.01 \pm 0.38	7.59 \pm 0.96	9.69 \pm 0.09	67.65 \pm 0.60	4.46 \pm 0.38	6.11 \pm 0.54
GRACE	70.14 \pm 1.40	7.49 \pm 3.78	7.67 \pm 3.78	68.25 \pm 0.99	6.41 \pm 0.71	7.38 \pm 0.84	67.81 \pm 0.41	10.77 \pm 0.68	10.69 \pm 0.69
GCA	70.43 \pm 1.19	18.08 \pm 4.80	20.04 \pm 4.34	69.34 \pm 0.20	6.07 \pm 0.96	7.39 \pm 0.82	67.07 \pm 0.14	7.90 \pm 1.10	8.05 \pm 1.07
FairDrop	69.01 \pm 1.11	3.66 \pm 2.32	7.61 \pm 2.21	67.78 \pm 0.60	5.77 \pm 1.83	5.48 \pm 1.32	67.32 \pm 0.61	4.05 \pm 1.05	3.77 \pm 1.00
NIFTY	69.93 \pm 0.09	3.31 \pm 1.52	4.70 \pm 1.04	67.15 \pm 0.43	4.40 \pm 0.99	3.75 \pm 1.04	65.52 \pm 0.31	6.51 \pm 0.51	5.14 \pm 0.68
FairAug	66.38 \pm 0.85	4.99 \pm 1.02	6.21 \pm 1.95	69.17 \pm 0.18	5.28 \pm 0.49	6.77 \pm 0.45	68.61 \pm 0.19	5.10 \pm 0.69	5.22 \pm 0.84
Graphair	69.36 \pm 0.45	2.56 \pm 0.41	4.64 \pm 0.17	68.17 \pm 0.08	2.10 \pm 0.17	2.76 \pm 0.19	67.43 \pm 0.25	2.02 \pm 0.40	1.62 \pm 0.47

RQ2: Trade-off between Accuracy and Fairness

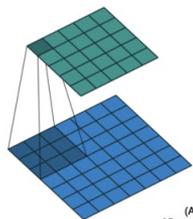
- Graphair achieves the best ACC-DP trade-off
 - The upper-left corner point represents the ideal performance, i.e., highest accuracy and lowest prediction bias.



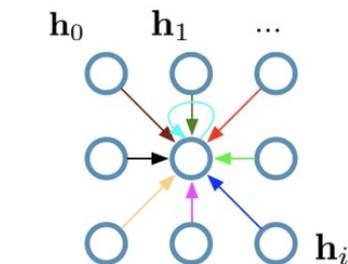
Fairness Graph Message Passing [C]

- What is Graph Message Passing?

Single CNN layer
with 3x3 filter:

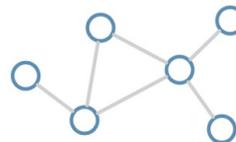


(Animation by
Vincent Dumoulin)

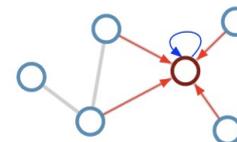


CNN

Consider this
undirected graph:



Calculate update
for node in red:



Update rule:
$$h_i^{(l+1)} = \sigma \left(h_i^{(l)} W_0^{(l)} + \sum_{j \in N_i} \frac{1}{c_{ij}} h_j^{(l)} W_1^{(l)} \right)$$

Scalability: subsample messages [Hamilton et al., NIPS 2017]

GNN

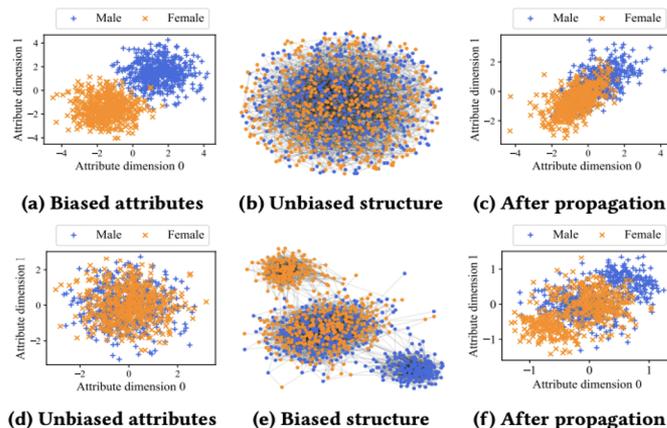
GNN is a general version of CNN!

Empirical Observations

- Aggregations in GNNs amplify bias compared with MLP.
 - GNNs > MLP in terms of prediction bias [5]
 - Representation bias after propagation even with unbiased input [6]

Table 2: Results of models w/ and w/o utilizing graph.

Dataset	Metrics	MLP	MLP-e	GCN	GAT
Pokey-z	ACC (%)	65.3 ±0.5	68.6 ±0.3	70.2 ±0.1	70.4 ±0.1
	AUC (%)	71.3 ±0.3	74.8 ±0.3	77.2 ±0.1	76.7 ±0.1
	Δ_{SP} (%)	3.8 ±1.3	6.9 ±1.0	9.9 ±1.1	9.1 ±0.9
	Δ_{EO} (%)	2.2 ±0.7	4.0 ±1.5	9.1 ±0.6	8.4 ±0.6
Pokey-n	ACC (%)	63.1 ±0.4	66.3 ±0.6	70.5 ±0.2	70.3 ±0.1
	AUC (%)	68.2 ±0.3	72.4 ±0.6	75.1 ±0.2	75.1 ±0.2
	Δ_{SP} (%)	3.3 ±0.6	8.7 ±1.0	9.6 ±0.9	9.4 ±0.7
	Δ_{EO} (%)	7.1 ±0.9	9.9 ±0.6	12.8 ±1.3	12.0 ±1.5



[5] Dai, Enyan, et al. "Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information." WSDM, 2021.

[6] Dong, Yushun, et al. "Edits: Modeling and mitigating data bias for graph neural networks." WWW, 2022.

Challenges

- Message passing is problematic from fairness aspect
 - Try fair message passing!

- Fair message passing scheme is challenging
 - Effective and fair with cross entropy loss and without data pre-processing
 - Compatible with backward propagation training
 - Grounded theoretical support (hand-craft architecture design is usually intuitive)

A Unified Optimization Framework

GNNs are graph signal denoising [7]

$$\arg \min_{\mathbf{F}} \mathcal{L}(\mathbf{F}) := \|\mathbf{F} - \mathbf{X}_{\text{in}}\|_F^2 + \mathcal{R}(\mathbf{F}, \tilde{\mathbf{L}})$$

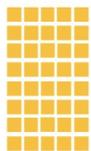
Close to the input

Smoothness prior

$$\mathcal{R}(\mathbf{F}, \tilde{\mathbf{L}}) = \lambda \text{tr}(\mathbf{F}^\top \tilde{\mathbf{L}} \mathbf{F}) = \lambda \sum_{(v_i, v_j) \in \mathcal{E}} \left\| \frac{\mathbf{F}_i}{\sqrt{d_i + 1}} - \frac{\mathbf{F}_j}{\sqrt{d_j + 1}} \right\|_2^2$$

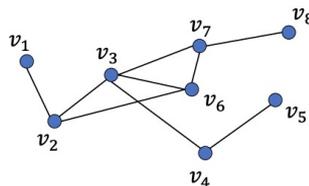
Define Prior \Rightarrow Optimization Solver \Rightarrow Message Passing

“Noisy Signal”



\mathbf{X}_{in}

Graph



“Nodes are similar to their neighbors”

“Clean Signal”



\mathbf{F}

- GCN
- PPNP
- APPNP/GCNI

$$\mathbf{X}_{\text{out}} = \tilde{\mathbf{A}} \mathbf{X}_{\text{in}}$$

$$\mathbf{X}_{\text{out}} = \alpha (\mathbf{I} - (1 - \alpha) \tilde{\mathbf{A}})^{-1} \mathbf{X}_{\text{in}}$$

$$\mathbf{X}^{(k+1)} = (1 - \alpha) \tilde{\mathbf{A}} \mathbf{X}^{(k)} + \alpha \mathbf{X}_{\text{in}}$$

[7] Ma, Yao, et al. “A unified view on graph neural networks as graph signal denoising.” CIKM 2021

Fair Message Passing

Define Prior \rightarrow Optimization Solver \rightarrow Message Passing

- Objective design

$$\min_{\mathbf{F}} \underbrace{\frac{\lambda_s}{2} \text{tr}(\mathbf{F}^T \tilde{\mathbf{L}} \mathbf{F}) + \frac{1}{2} \|\mathbf{F} - \mathbf{X}_{trans}\|_F^2}_{h_s(\mathbf{F})} + \underbrace{\lambda_f \|\Delta_s S F(\mathbf{F})\|_1}_{h_f(\Delta_s S F(\mathbf{F}))} \rightarrow \text{Fairness prior}$$

- Optimization solver

- Avoid L1 norm objective via Fenchel conjugate $\min_{\mathbf{F}} \max_{\mathbf{u}} h_s(\mathbf{F}) + \langle \mathbf{p}, \mathbf{u} \rangle - h_f^*(\mathbf{u})$
- Proximal Alternating Predictor-Corrector Solver [8]

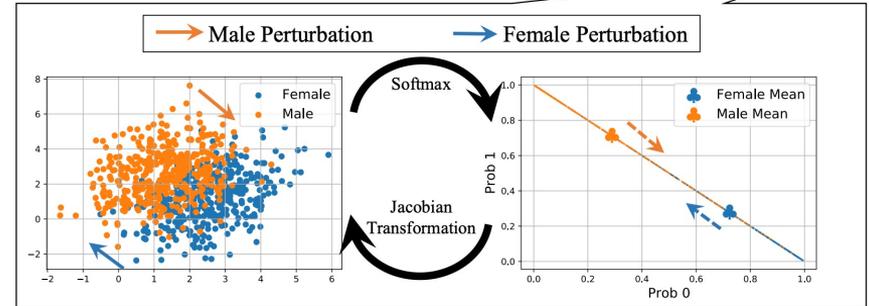
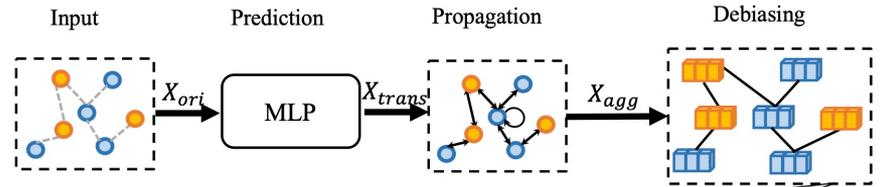
- Fair Message passing

$$\left\{ \begin{array}{l} \mathbf{X}_{agg}^{k+1} = \gamma \mathbf{X}_{trans} + (1 - \gamma) \tilde{\mathbf{A}} \mathbf{F}^k, \\ \bar{\mathbf{F}}^{k+1} = \mathbf{X}_{agg}^{k+1} - \gamma \frac{\partial \langle \mathbf{p}, \mathbf{u}^k \rangle}{\partial \mathbf{F}} \Big|_{\mathbf{F}^k}, \\ \bar{\mathbf{u}}^{k+1} = \mathbf{u}^k + \beta \Delta_s S F(\bar{\mathbf{F}}^{k+1}), \\ \mathbf{u}^{k+1} = \min(|\bar{\mathbf{u}}^{k+1}|, \lambda_f) \cdot \text{sign}(\bar{\mathbf{u}}^{k+1}), \\ \mathbf{F}^{k+1} = \mathbf{X}_{agg}^{k+1} - \gamma \frac{\partial \langle \mathbf{p}, \mathbf{u}^{k+1} \rangle}{\partial \mathbf{F}} \Big|_{\mathbf{F}^k}. \end{array} \right. \begin{array}{l} \text{Step ①} \rightarrow \text{Aggregation with skip connection} \\ \text{Step ②} \\ \text{Step ③} \rightarrow \text{Learn and reshape perturbation vector } \mathbf{u} \\ \text{Step ④} \\ \text{Step ⑤} \end{array}$$

[8] Ignace Loris, et al. On a generalization of the iterative soft-thresholding algorithm for the case of non-separable penalty. Inverse Problems, 27(12):125007, 2011.

Fair Message Passing

- FMP Interpretation
 - Three stages in FMP
 - Four steps in Debiasing
- Efficiency
 - Negligible additional computation
- White-box sensitive attribute usage
 - Explicit usage in FMP
 - Implicit encoding in parameters for fair training



Experiments

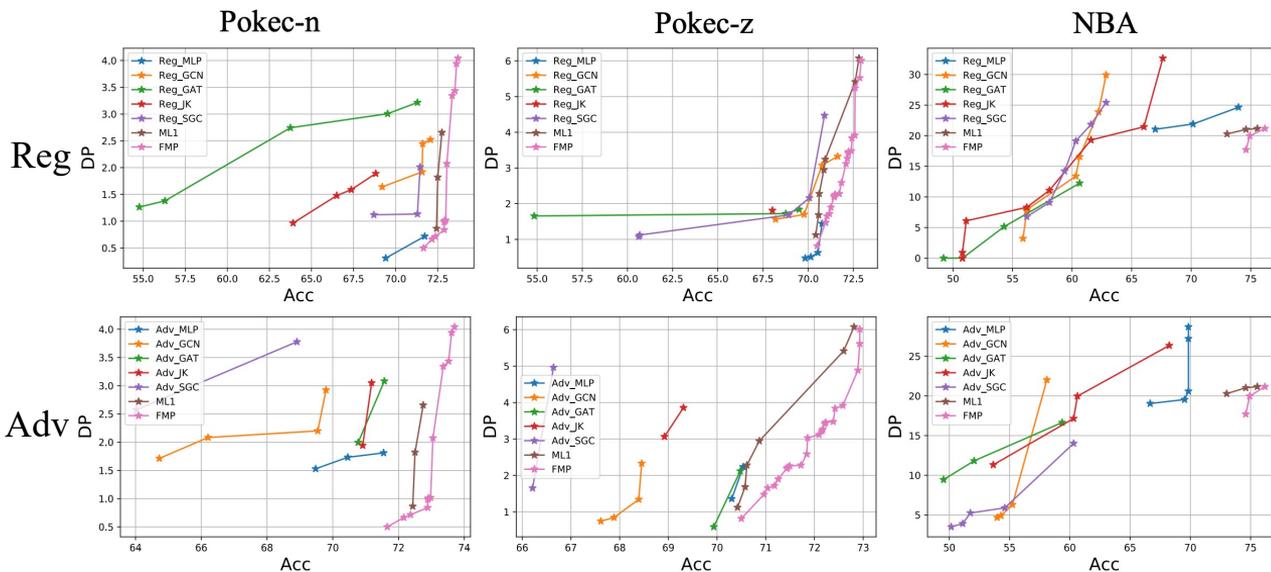
- Various GNNs architectures with cross-entropy loss
 - Lowest DP and EO with comparable accuracy

Table 1: Comparative Results with Baselines on Node Classification.

Models	Pokec-z			Pokec-n			NBA		
	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{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
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
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
APPNP	66.91 \pm 1.46	3.90 \pm 0.69	5.71 \pm 1.29	69.80 \pm 0.89	1.98 \pm 1.30	4.01 \pm 2.36	63.80 \pm 1.19	26.51 \pm 3.33	20.00 \pm 4.56
JKNet	66.89 \pm 3.79	1.28 \pm 0.96	1.79 \pm 0.82	63.59 \pm 6.36	1.91 \pm 2.14	0.70 \pm 0.92	67.94 \pm 2.73	27.80 \pm 8.41	20.33 \pm 7.52
ML1	70.42 \pm 0.40	2.35 \pm 0.83	2.00 \pm 0.50	72.36 \pm 0.26	1.47 \pm 1.12	3.03 \pm 1.77	72.70 \pm 1.19	26.46 \pm 4.93	25.50 \pm 8.38
FMP	70.50 \pm 0.50	0.81 \pm 0.40	1.73 \pm 1.03	72.16 \pm 0.33	0.66 \pm 0.40	1.47 \pm 0.87	73.33 \pm 1.85	18.92 \pm 2.28	13.33 \pm 5.89

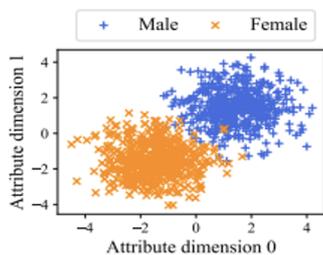
Experiments

- Compared with fair training
 - adding regularization and adversarial debiasing as baselines
 - FMP still achieves better tradeoff performance

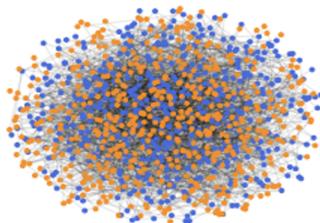


Topology Matters in Fair Graph Learning [D]

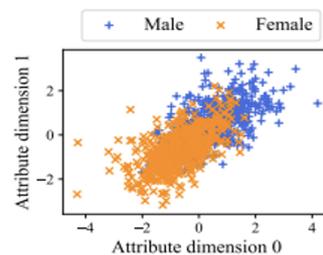
- Topology serves as an additional bias source in graph learning.



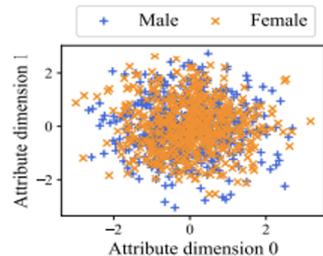
(a) Biased attributes



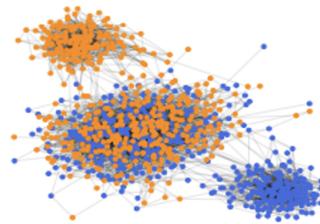
(b) Unbiased structure



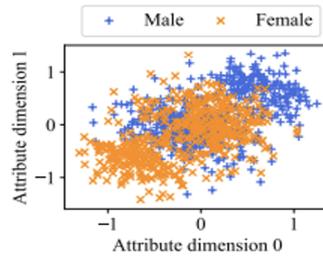
(c) After propagation



(d) Unbiased attributes



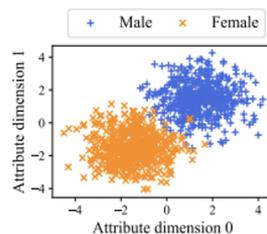
(e) Biased structure



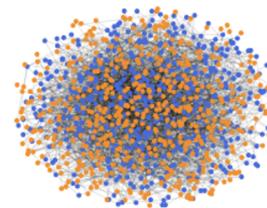
(f) After propagation

Understanding Graph Data Bias

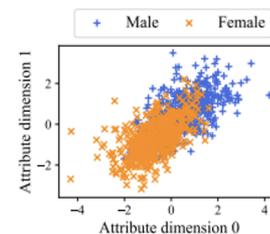
- Understanding the bias in graph neural networks (GNNs)
 - GNNs demonstrate empirically higher prediction bias than peer multilayer perception (MLP) [4] but without theoretical understanding.
 - Bias representation after propagation for bias structure even with unbiased attributes
 - When and Why does aggregation enhance the bias?



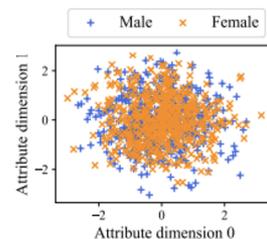
(a) Biased attributes



(b) Unbiased structure



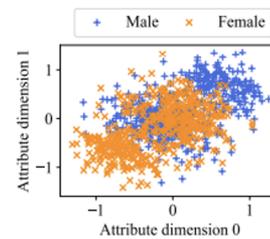
(c) After propagation



(d) Unbiased attributes



(e) Biased structure



(f) After propagation

[5] Dai, Enyan, et al. "Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information." WSDM, 2021.

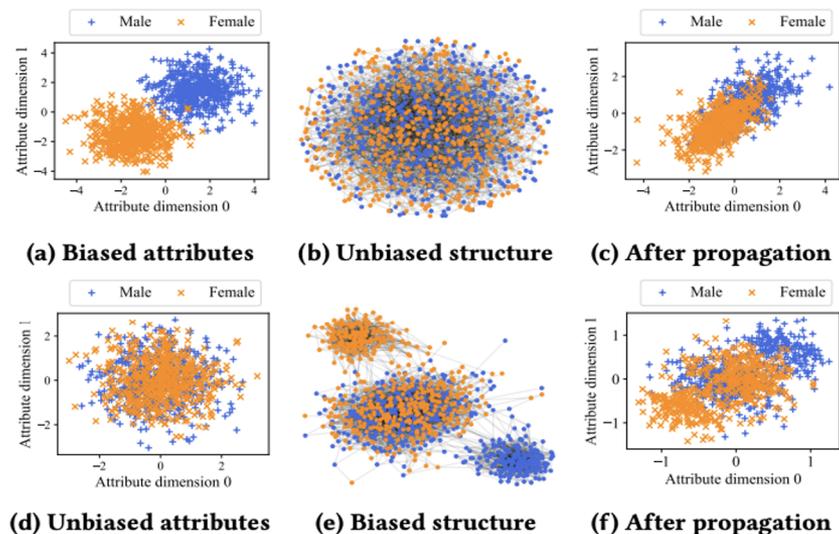
[6] Dong, Yushun, et al. "Edits: Modeling and mitigating data bias for graph neural networks." WWW, 2022.

Empirical Observations

- Aggregations in GNNs amplify bias compared with MLP.
 - GNNs > MLP in terms of prediction bias [5]
 - Representation bias after propagation even with unbiased input [6]

Table 2: Results of models w/ and w/o utilizing graph.

Dataset	Metrics	MLP	MLP-e	GCN	GAT
Pokey-z	ACC (%)	65.3 \pm 0.5	68.6 \pm 0.3	70.2 \pm 0.1	70.4 \pm 0.1
	AUC (%)	71.3 \pm 0.3	74.8 \pm 0.3	77.2 \pm 0.1	76.7 \pm 0.1
	Δ_{SP} (%)	3.8 \pm 1.3	6.9 \pm 1.0	9.9 \pm 1.1	9.1 \pm 0.9
	Δ_{EO} (%)	2.2 \pm 0.7	4.0 \pm 1.5	9.1 \pm 0.6	8.4 \pm 0.6
Pokey-n	ACC (%)	63.1 \pm 0.4	66.3 \pm 0.6	70.5 \pm 0.2	70.3 \pm 0.1
	AUC (%)	68.2 \pm 0.3	72.4 \pm 0.6	75.1 \pm 0.2	75.1 \pm 0.2
	Δ_{SP} (%)	3.3 \pm 0.6	8.7 \pm 1.0	9.6 \pm 0.9	9.4 \pm 0.7
	Δ_{EO} (%)	7.1 \pm 0.9	9.9 \pm 0.6	12.8 \pm 1.3	12.0 \pm 1.5



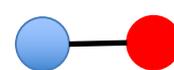
[5] Dai, Enyan, et al. "Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information." WSDM, 2021.

[6] Dong, Yushun, et al. "Edits: Modeling and mitigating data bias for graph neural networks." WWW, 2022.

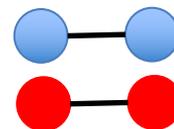
Why does Aggregation suffer?

Intuition

- Graph topology with high sensitive homophily coefficient
 - Definition: $\# \text{sensitive homo links} / \# \text{links}$
 - E.g., 95.30% for Pokec-n dataset
 - Higher than label homophily coefficient
- Graph concentration (over-smoothing)
 - More similar representation within the demographic group
 - Conditionally happens: no bias for fully over-smoothing

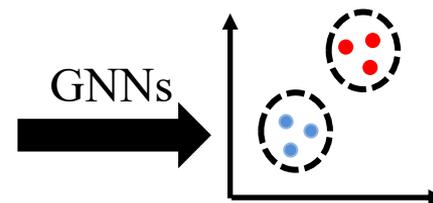
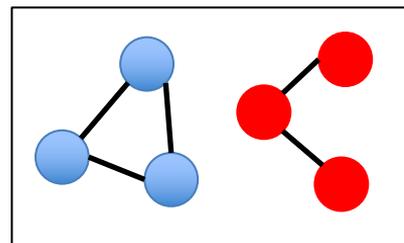


Inter link



Intra link

$$\text{Sensitive Homophily} = \frac{\# \text{ Intra links}}{\# \text{ all links}}$$



How can we theoretically understanding such GNNs behavior?

A Pilot Theoretical Study

Goal: find a **sufficient condition** of bias enhancement after aggregation

- Synthetic graph data: contexture stochastic block model
 - Topology with intra/inter-connect probability
 - Features with Gaussian Mixture Model
- GCN-like Aggregation
- Bias difference before/after aggregation

When does bias enhancement happen?

- large sensitive homophily coefficient & node number & connection density
- Balanced demographic size

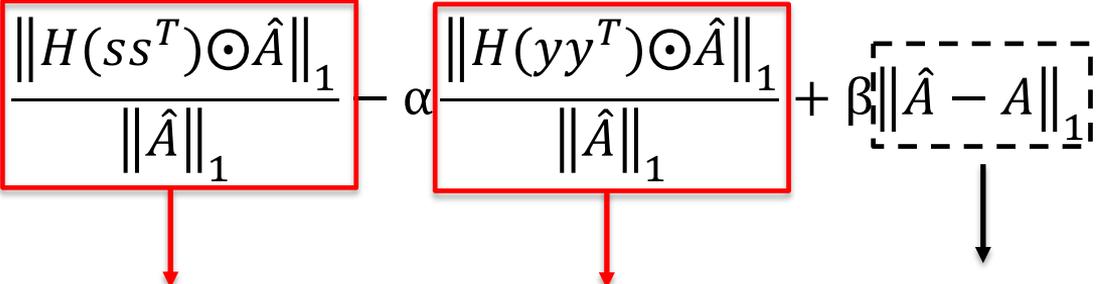
Topology matters in fair graph learning!

Fair Graph Rewiring

Preprocessing: rewire graph topology to achieve graph fairness

- Large label homophily coefficient
- Low sensitive homophily coefficient
- Low topology perturbation

$$L(\hat{A}|s, y, A) = \frac{\|H(ss^T) \odot \hat{A}\|_1}{\|\hat{A}\|_1} - \alpha \frac{\|H(yy^T) \odot \hat{A}\|_1}{\|\hat{A}\|_1} + \beta \| \hat{A} - A \|_1$$



Sensitive Homophily Label Homophily Topology Perturbation

Synthetic Experiments

Observations:

- 1) Bias enhancement happens conditionally
- 2) The bias enhancement tendency is consistent with our theory

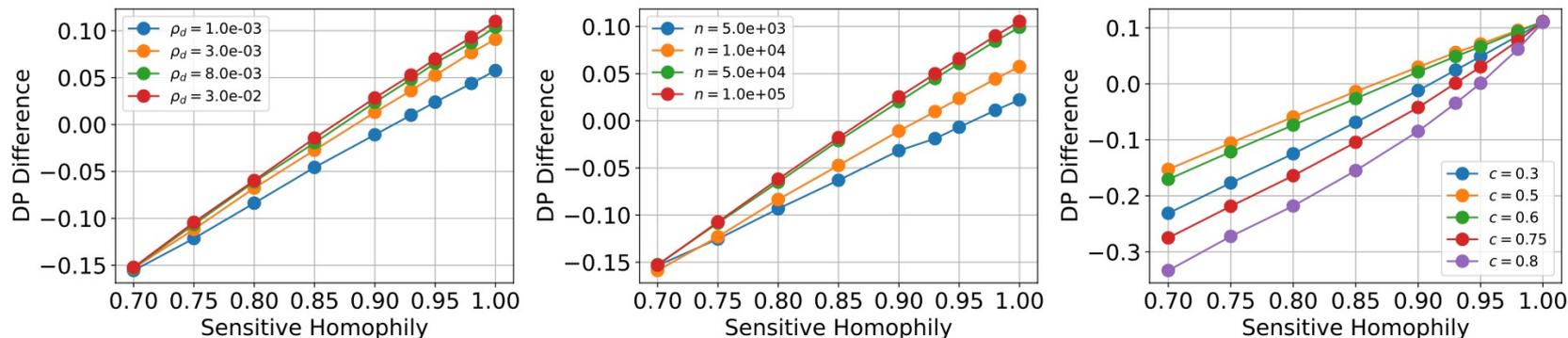


Figure 2: The difference of demographic parity for message passing. **Left:** DP difference for different graph connection density ρ_d with sensitive attribute ratio $c = 0.5$ and number of nodes $n = 10^4$; **Middle:** DP difference for different number of nodes n with sensitive attribute ratio $c = 0.5$ 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 $n = 10^4$;

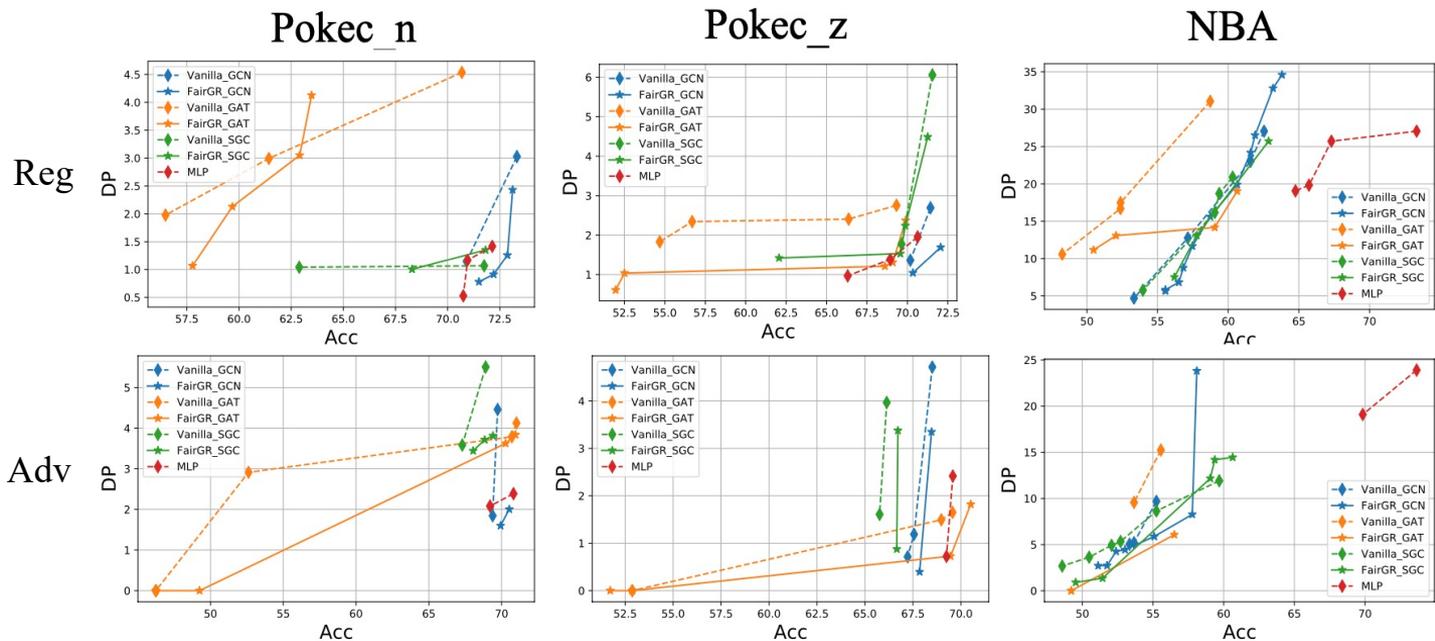
Experiments on Real Data

- Graph rewiring facilitates fairness in vanilla GNNs
 - Lower demographic parity and equal opportunities for various GNNs architectures

Models	Pokec-z			Pokec-n			NBA		
	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{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	1.04 \pm 0.80	1.14 \pm 1.02	61.27 \pm 9.34	0.54 \pm 0.51	2.27 \pm 1.55	53.65 \pm 10.31	4.16 \pm 5.13	3.67 \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	1.27 \pm 0.71	72.68 \pm 0.44	0.47 \pm 0.39	0.82 \pm 0.78	61.59 \pm 1.85	20.24 \pm 4.41	9.50 \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	0.71 \pm 0.65	2.39 \pm 0.69	62.54 \pm 1.62	18.56 \pm 2.81	2.50 \pm 1.66

Experiments on Real Data

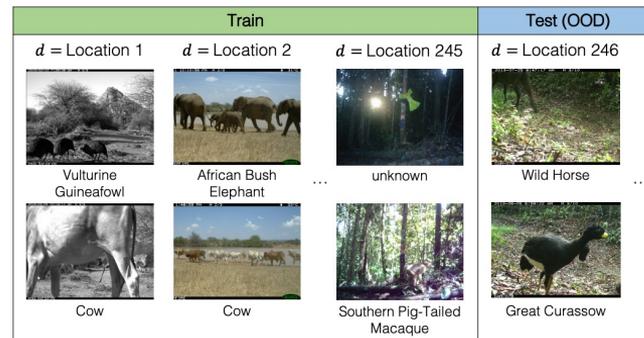
- Graph rewiring is orthogonal with other fair learning methods
 - E.g., adding regularization, adversarial debiasing, and Fair Mixup
 - Better fairness-prediction tradeoff performance with GR



Chasing Fairness under Distribution Shift [E]

- When do distribution shifts happen [9]?

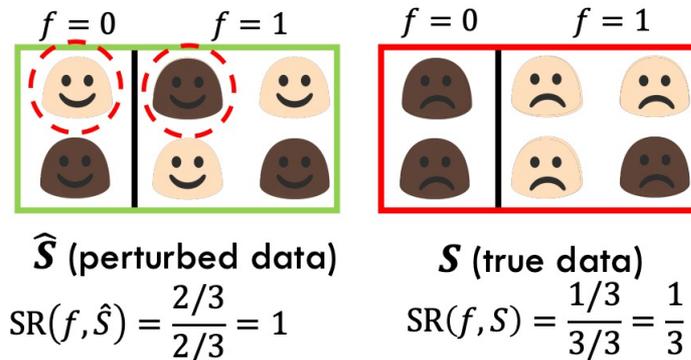
- Different locations/hospitals;
- Different experiments;
- Different time periods;
- Different devices.



- Robust Fairness under Distribution Shift

- Many metrics (acc, fairness) under distribution shifts
- Fairness metric is more vulnerable under distribution shifts [10]

- Statistical rate (SR) v.s. sensitive attribute perturbation



[9] Koh, Pang Wei, et al. "Wilds: A benchmark of in-the-wild distribution shifts." *International Conference on Machine Learning*. PMLR, 2021.

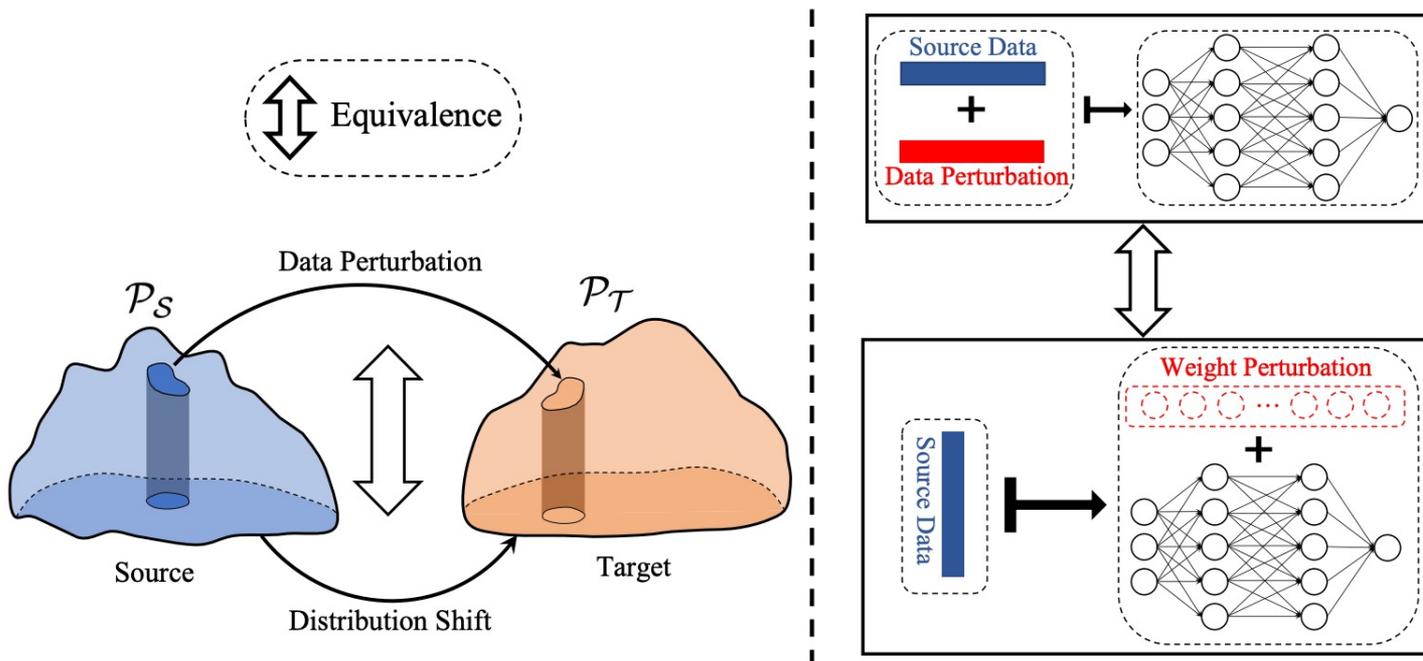
[10] Celis, L. Elisa, et al. "Fair classification with adversarial perturbations." *Advances in Neural Information Processing Systems*. (2021).

Challenges

- The distribution shifts are unknown during training
 - Can not access feature distributions in the test dataset
- Fairness is vulnerable under distribution shifts
 - Many methods are originally designed for performance
 - What can we do from model perspective to tackle data distribution shifts problem?

Rethinking Distribution Shifts

Distribution shifts \longleftrightarrow Data Perturbation \longleftrightarrow Model Weight Perturbation

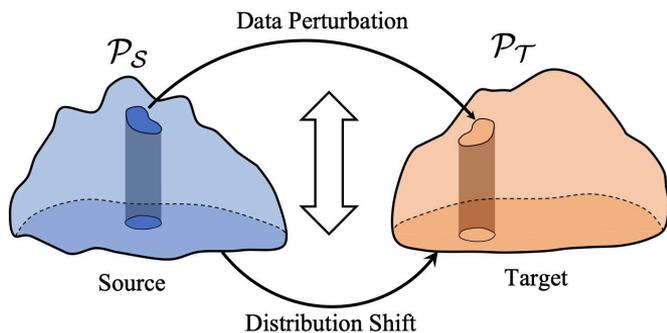


Fairness under Distribution Shifts

- Distribution shifts are data perturbation
 - Holds for any loss function and model architectures

Corollary 2.2. *Given source and target datasets with probability distribution \mathcal{P}_S and \mathcal{P}_T , there exists data perturbation δ so that the training loss of any neural network $f_\theta(\cdot)$ for target distribution equals that for source distribution with data perturbation δ , i.e.,*

$$\mathbb{E}_{(X,Y)\sim\mathcal{P}_T}[l(f_\theta(X), Y)] = \mathbb{E}_{\delta_X, \delta_Y} \mathbb{E}_{(X,Y)\sim\mathcal{P}_S}[l(f_\theta(X + \delta_X), Y + \delta_Y)]. \quad (3)$$



Data Perturbation

Fairness under Distribution Shifts

- Data perturbation equals model weight perturbation
 - Holds for any loss and model architectures

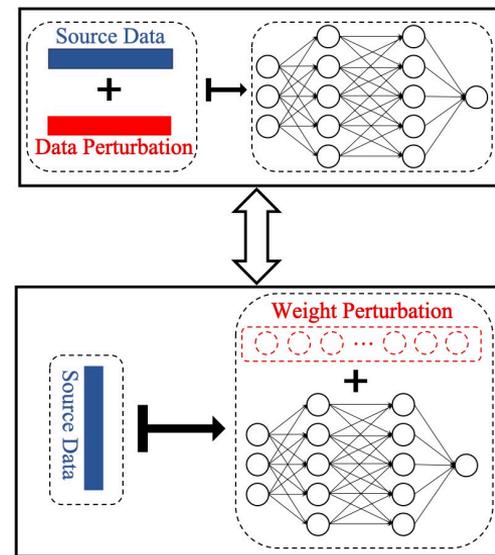
Theorem 2.3. *Considering the source dataset with distribution \mathcal{P}_S , suppose the source dataset is perturbed with data perturbation δ , and the neural network is given by $f_\theta(\cdot)$, there exists model weight perturbation $\Delta\theta$ so that the training loss on perturbed source dataset is the same with that for model weight perturbation $\Delta\theta$ on source distribution:*

$$\mathbb{E}_{\delta_X, \delta_Y} \mathbb{E}_{(X, Y) \sim \mathcal{P}_S} [l(f_\theta(X + \delta_X), Y + \delta_Y)] = \mathbb{E}_{(X, Y) \sim \mathcal{P}_S} [l(f_{\theta + \Delta\theta}(X), Y)]. \quad (4)$$

Data Perturbation



Model Weight Perturbation



Fairness under Distribution Shifts

- Fairness under distribution shifts
 - What are the conditions guarantee such robust fairness?
- Take demographic parity (DP) as an example
 - Low DP at source dataset
 - **Low average prediction gap** between source/target dataset at the same sensitive group

$$\begin{aligned} DP_{\mathcal{T}} &\stackrel{(a)}{\leq} DP_S + \left| |\mathbb{E}_{\mathcal{T}_0}[f_{\theta}(\mathbf{x})] - \mathbb{E}_{\mathcal{T}_1}[f_{\theta}(\mathbf{x})]| - |\mathbb{E}_{S_0}[f_{\theta}(\mathbf{x})] - \mathbb{E}_{S_1}[f_{\theta}(\mathbf{x})]| \right| \\ &\stackrel{(b)}{\leq} \boxed{DP_S} + \boxed{|\mathbb{E}_{S_0}[f_{\theta}(\mathbf{x})] - \mathbb{E}_{\mathcal{T}_0}[f_{\theta}(\mathbf{x})|} + \boxed{|\mathbb{E}_{S_1}[f_{\theta}(\mathbf{x})] - \mathbb{E}_{\mathcal{T}_1}[f_{\theta}(\mathbf{x})|} \end{aligned}$$

Δ_0 Δ_1

Loss function-agnostic

Fairness under Distribution Shifts

- How can we achieve low prediction gap for each demographic group?
 - Model weight perturbation: bi-level optimization problem
 - The worst case within the model weight perturbation ball for each sensitive attribute group
 - Can be accelerated with two forward-backward propagation

$$\begin{aligned}\Delta_0 &\leq \max_{\|\epsilon_0\|_p \leq \rho} |\mathbb{E}_{\mathcal{S}_0}[f_{\theta+\epsilon_0}(\mathbf{x})] - \mathbb{E}_{\mathcal{S}_0}[f_{\theta}(\mathbf{x})]| \\ &\approx \max_{\|\epsilon_0\|_p \leq \rho} \mathbb{E}_{\mathcal{S}_0}[f_{\theta+\epsilon_0}(\mathbf{x})] - \mathbb{E}_{\mathcal{S}_0}[f_{\theta}(\mathbf{x})] \quad \text{Model Weight Perturbation} \\ &\triangleq \mathcal{L}_{RFR, \mathcal{S}_0},\end{aligned}$$

Fairness under Distribution Shifts

- Robust Fairness Regularization (RFR)

$$\mathcal{L}_{all} = \mathcal{L}_{CLF} + \lambda \cdot (\mathcal{L}_{DP} + \mathcal{L}_{RFR}),$$

Classification loss

Low DP on source

Low prediction gap

Experiments

Distribution shifts type

- Synthetic distribution shifts
 - Use different sample selection prior distribution across training and test dataset [11]
 - Dataset: Adult; ACS-I; ACS-E
- Real-world distribution shifts
 - New Adult dataset [12]
 - Distribution shift across states
 - Distribution shift across times

Evaluation metrics: Accuracy & DP tradeoff

Baselines:

- Regularization & Adversarial debiasing
- Fair consistency regularization (FCR)

[11] Ashkan Rezaei, et al. Robust fairness under 375 covariate shift. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 9419–9427, 2021.

[12] Frances Ding, et al. Retiring adult: New datasets for fair machine learning. NeurIPS, 2021.

Experiments

Synthetic distribution shifts

- Low prediction bias at low-intensity distribution shifts
- Comparable at high-intensity distribution shifts

Table 1: Performance Comparison with Baselines on Synthetic Dataset. (α, β) control distribution shift intensity, and $(0, 1)$ represents no distribution shift. The best/second-best results are highlighted in **boldface**/underlined, respectively.

(α, β)	Methods	Adult			ACS-I			ACS-E		
		Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow
(1.0, 2.0)	MLP	82.09 \pm 0.05	15.11 \pm 0.04	14.33 \pm 0.05	77.95 \pm 0.52	3.51 \pm 0.59	3.77 \pm 0.55	80.95 \pm 0.10	1.10 \pm 0.06	1.43 \pm 0.06
	REG	80.60 \pm 0.05	3.79 \pm 0.06	3.27 \pm 0.08	77.77 \pm 0.09	2.28 \pm 0.32	2.59 \pm 0.23	80.44 \pm 0.07	0.86 \pm 0.09	1.05 \pm 0.10
	ADV	78.80 \pm 0.68	0.83 \pm 0.26	0.79 \pm 0.14	75.72 \pm 0.63	1.96 \pm 0.38	2.00 \pm 0.35	79.39 \pm 0.15	1.09 \pm 0.26	0.95 \pm 0.26
	FCR	79.06 \pm 0.09	<u>9.98\pm0.06</u>	<u>9.47\pm0.07</u>	76.99 \pm 0.47	<u>2.94\pm0.34</u>	<u>2.95\pm0.28</u>	79.74 \pm 0.11	0.97 \pm 0.21	<u>1.00\pm0.22</u>
	RFR	78.84 \pm 0.09	0.44\pm0.05	0.12\pm0.06	74.15 \pm 0.81	1.84\pm0.27	1.60\pm0.33	80.08 \pm 0.08	0.71\pm0.10	0.06\pm0.11
(1.5, 3.0)	MLP	82.05 \pm 0.05	15.16 \pm 0.09	14.33 \pm 0.09	77.85 \pm 0.25	3.73 \pm 0.53	3.70 \pm 0.56	80.42 \pm 0.10	1.14 \pm 0.07	1.10 \pm 0.07
	REG	80.64 \pm 0.08	3.74 \pm 0.11	3.23 \pm 0.10	77.87 \pm 0.18	2.25 \pm 0.28	2.37 \pm 0.27	80.21 \pm 0.13	0.72\pm0.04	0.75 \pm 0.03
	ADV	78.71 \pm 0.41	1.07 \pm 0.87	<u>0.87\pm0.96</u>	75.79 \pm 0.68	2.22 \pm 0.53	2.44 \pm 0.48	79.58 \pm 0.13	1.07 \pm 0.19	<u>1.26\pm0.18</u>
	FCR	79.05 \pm 0.12	<u>10.01\pm0.07</u>	<u>9.51\pm0.06</u>	77.06 \pm 0.68	<u>3.39\pm0.33</u>	3.10 \pm 0.36	79.59 \pm 0.26	1.17 \pm 0.24	1.08 \pm 0.23
	RFR	78.91 \pm 0.03	0.46\pm0.10	0.16\pm0.09	74.19 \pm 0.58	1.82\pm0.29	2.17\pm0.32	80.47 \pm 0.03	0.72\pm0.04	0.71\pm0.05
(3.0, 6.0)	MLP	82.07 \pm 0.05	15.23 \pm 0.14	14.45 \pm 0.15	77.89 \pm 0.45	3.35 \pm 0.36	3.47 \pm 0.41	80.30 \pm 0.04	1.17 \pm 0.04	1.13 \pm 0.04
	REG	80.62 \pm 0.07	3.72 \pm 0.05	3.21 \pm 0.04	78.19 \pm 0.12	1.60\pm0.48	1.84\pm0.44	80.36 \pm 0.09	0.70\pm0.09	0.68 \pm 0.11
	ADV	78.97 \pm 0.49	1.28\pm0.74	1.09\pm0.50	75.71 \pm 0.68	2.28 \pm 0.39	2.24 \pm 0.41	79.66 \pm 0.16	1.34 \pm 0.14	<u>1.16\pm0.13</u>
	FCR	79.03 \pm 0.13	<u>10.00\pm0.05</u>	<u>9.50\pm0.05</u>	76.71 \pm 0.39	2.97 \pm 0.34	3.28 \pm 0.31	79.89 \pm 0.22	1.06 \pm 0.14	1.14 \pm 0.18
	RFR	80.15 \pm 0.07	<u>1.75\pm0.15</u>	<u>1.30\pm0.14</u>	74.22 \pm 0.56	<u>1.80\pm0.26</u>	<u>1.89\pm0.24</u>	80.28 \pm 0.12	<u>0.74\pm0.04</u>	0.51\pm0.04

Experiments

Real distribution shifts

- Low(comparable) prediction bias under temporal(spatial) distribution shift

Table 2: Performance comparison with baselines on real temporal (the year 2016 to the year 2018) and spatial (Michigan State to California State) distribution shift. The best and second-best results are highlighted with **bold** and underline, respectively.

Real	Methods	ACS-I			ACS-E		
		Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow	Acc (%) \uparrow	Δ_{DP} (%) \downarrow	Δ_{EO} (%) \downarrow
2016 \rightarrow 2018	MLP	77.75 \pm 0.44	3.26 \pm 0.38	3.48 \pm 0.41	80.46 \pm 0.05	1.07 \pm 0.10	1.02 \pm 0.10
	REG	77.74 \pm 0.62	2.09 \pm 0.21	2.27 \pm 0.24	80.37 \pm 0.12	<u>0.77\pm0.08</u>	<u>0.74\pm0.08</u>
	ADV	75.94 \pm 0.40	<u>2.41\pm0.49</u>	<u>2.53\pm0.55</u>	79.62 \pm 0.14	<u>1.17\pm0.14</u>	<u>1.10\pm0.14</u>
	FCR	76.40 \pm 0.45	2.81 \pm 0.30	2.96 \pm 0.30	79.59 \pm 0.38	0.95 \pm 0.42	0.91 \pm 0.34
	RFR	77.49 \pm 0.32	1.36\pm0.17	1.49\pm0.17	80.36 \pm 0.05	0.61\pm0.11	0.58\pm0.10
MI \rightarrow CA	MLP	75.62 \pm 0.80	5.22 \pm 0.86	3.60 \pm 0.34	79.02 \pm 0.20	0.73 \pm 0.07	0.94 \pm 0.05
	REG	75.52 \pm 0.78	2.88 \pm 0.44	2.17 \pm 0.22	75.34 \pm 1.11	0.42\pm0.09	0.61\pm0.11
	ADV	73.38 \pm 1.07	1.04\pm0.58	0.54\pm0.38	77.56 \pm 0.41	0.61 \pm 0.18	0.80 \pm 0.13
	FCR	74.28 \pm 0.35	5.06 \pm 0.62	3.67 \pm 0.51	77.96 \pm 0.22	0.44 \pm 0.14	0.67 \pm 0.38
	RFR	74.63 \pm 0.45	<u>1.35\pm0.39</u>	<u>1.30\pm0.24</u>	78.84 \pm 0.21	<u>0.44\pm0.09</u>	<u>0.65\pm0.07</u>

Experiments

Pareto tradeoff performance for synthetic distribution shifts

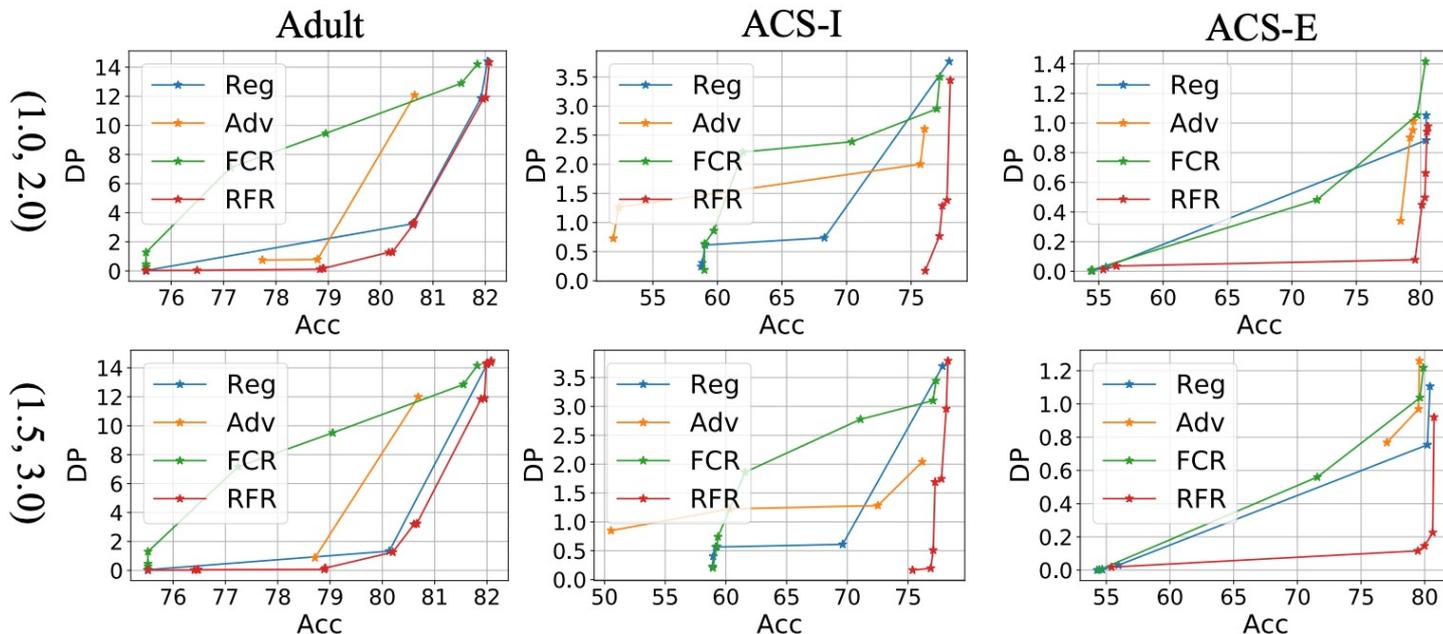
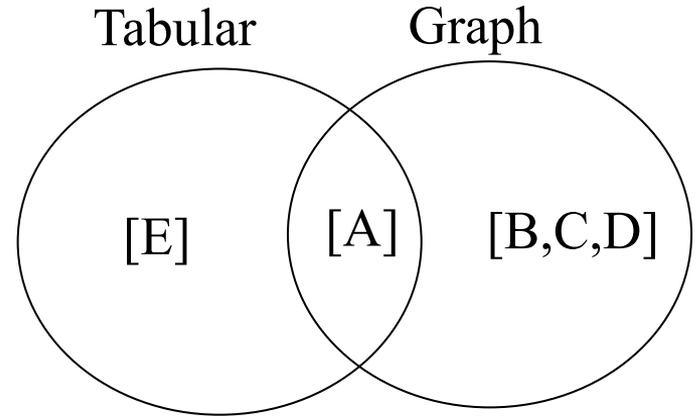
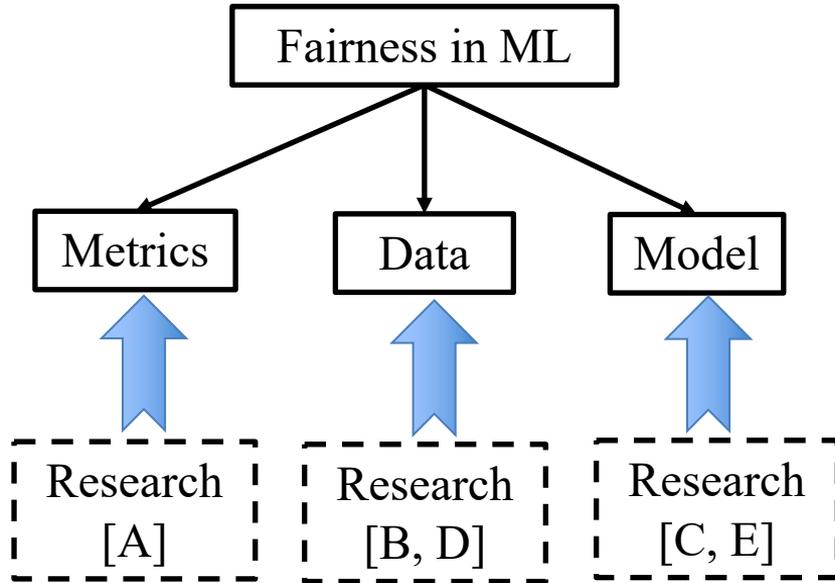


Figure 2: The fairness (DP) and prediction (Acc) trade-off performance on three datasets with different synthetic distribution shifts. The units for x- and y-axis are percentages (%).

Summary

- Algorithmic fairness is considered from **three** aspects: evaluation metrics, data, and model.
- Propose Generalized Demographic Parity (GDP) to broaden evaluation metrics for continuous sensitive attributes with tractable computation.
- Propose Graphair to conduct automatic graph data augmentation via learnable feature masking and adjacency matrix.
- Propose Fair Message Passing (FMP) to achieve graph fairness from model architecture perspective.
- Understand the role of topology in fair graph learning and propose Fair Graph Rewiring (FairGR) to mitigate bias from the data perspective.
- Rethink the distribution shifts problem and propose Robust Fairness Regularization (RFR) to achieve fairness under distribution shifts.

Summary



Q & A
