
Fair & Disentangle Graph Mining Papers Organization

CAU
Junseo, Yu

DMAIS Lab Meeting
12.06.2024

Contents

01 Fairness

- Fairness notion
- Fairwalk
- Crosswalk
- InFoRM

02 Disentanglement

- Disentangle notion
- DisenGCN
- FactorGCN

03 Invariant Learning & Causal-based Learning

- Invariant Learning
- Causal-based Learning
- DIR
- DisC

04 Fairness & Disentanglement

- FairSAD
- FairINV

Weekly Meetings

1. Fairness

- Fairness notion
- Fairwalk
- Crosswalk
- InFoRM

Fairness

- ❑ The quality of treating people **equally** or in a way that is **right** or **reasonable**.
- ❑ How do we **define** the **equally, right, or reasonable**?
 - It might be the problem of the **philosophy, ethics**, or/and **sociology**
 - It continues to change over time. (e.g., **Golden Rule**)

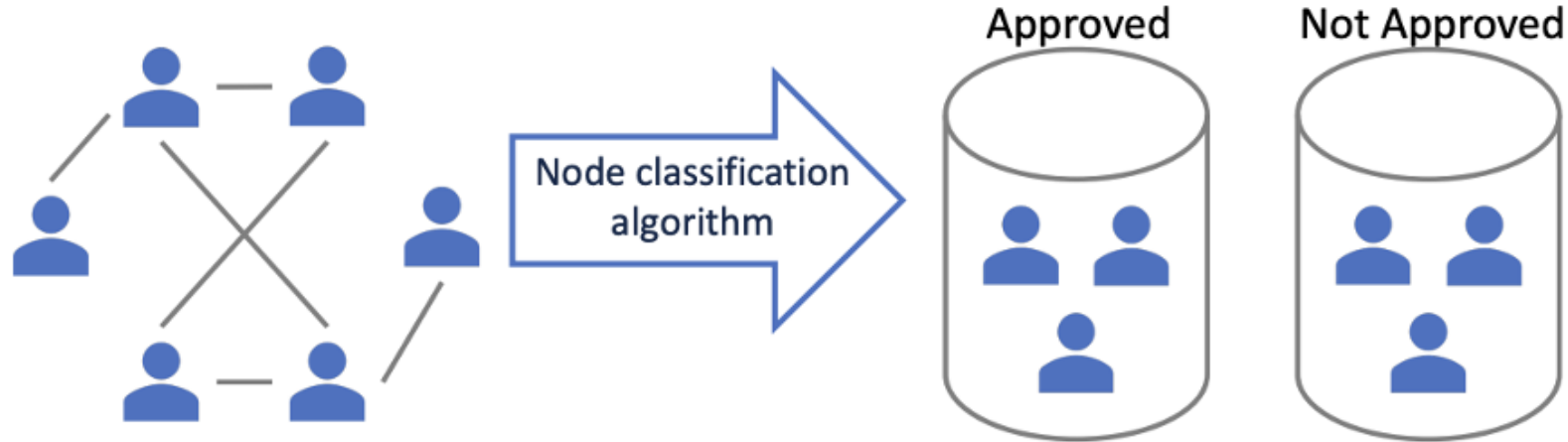


Fairness

Fairness notion

Fair & Disentangle graph mining

Fairness in Machine Learning(AI)



Loan Approval

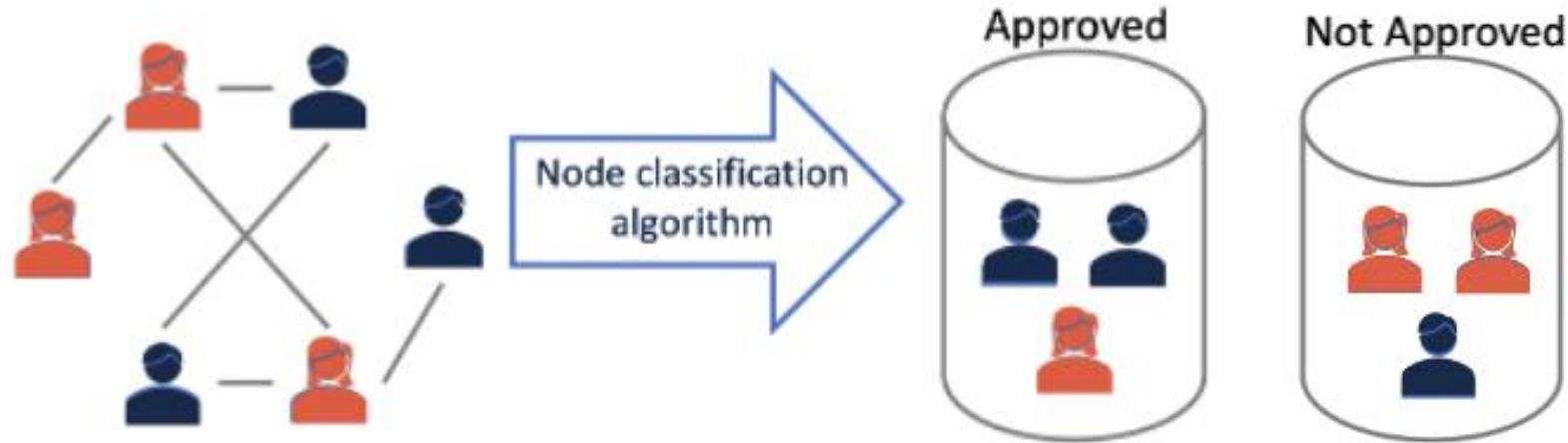
- ❑ The model try to classify the node(people) by using various techniques
- ❑ However, without considering fairness, the model can **learn in an unfair way**
- ❑ **Potential Cause:** Biased Data, Spurious Correlation, Biased learning strategy, and others

Fairness

Fairness notion

Fair & Disentangle graph mining

Fairness in Machine Learning(AI)



Loan Approval

- ❑ The upper Scenario **might** be **unfair**.
- ❑ Since Male has a higher approval rate than female.
- ❑ However, is it **absolute**? If not, how we can measure **the (un)fairness level**?

Fairness

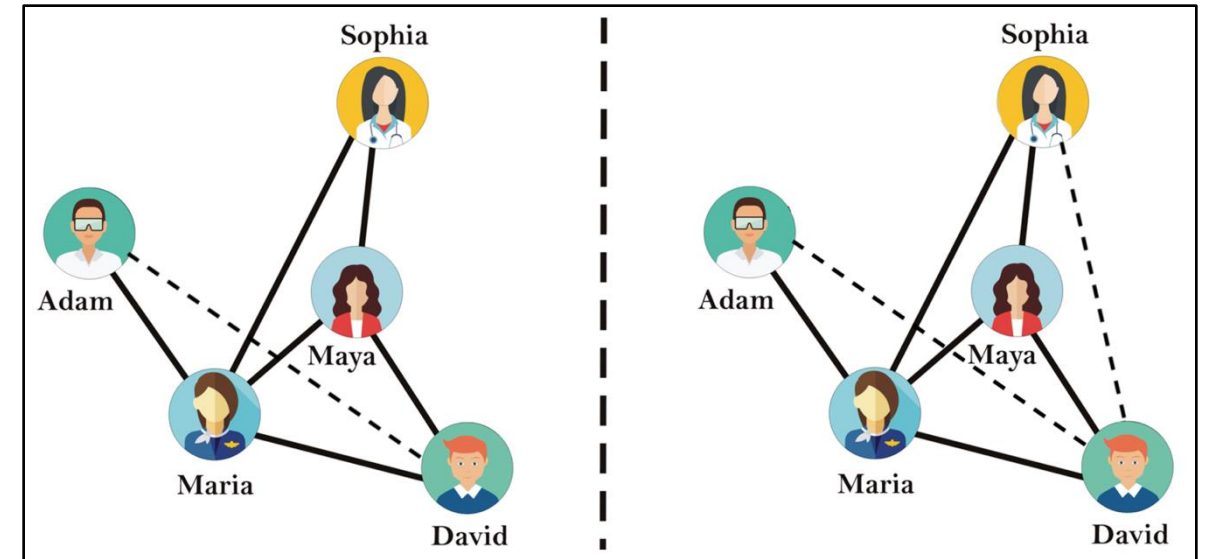
Fairness notion

Fairness in Machine Learning(AI)



Mosaic removal model

Fair & Disentangle graph mining



Link Prediction in SNS

Dashed line : Prediction

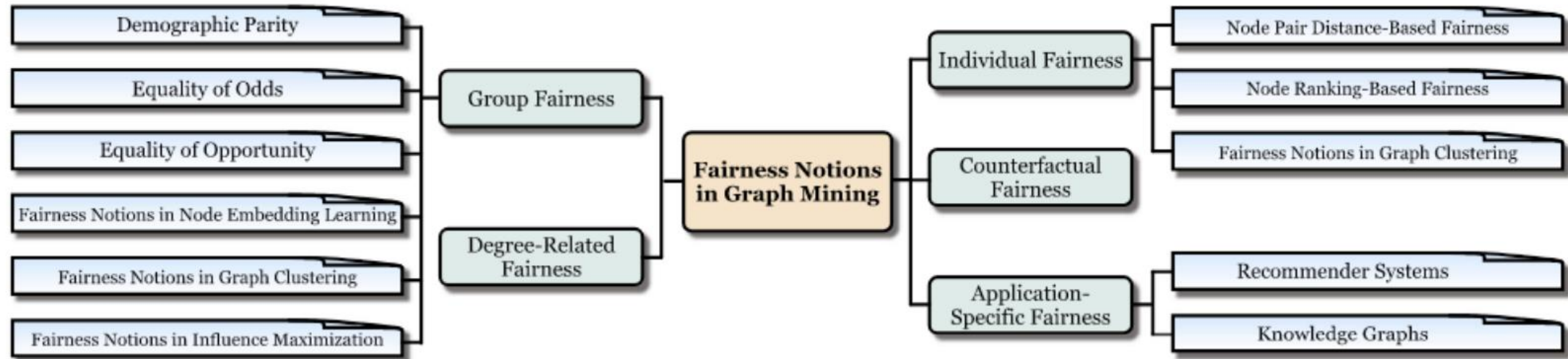
Bold line : Actual Line

Fairness

Fairness notion

Fair & Disentangle graph mining

Fairness Metrics



❑ Group Fairness

- DP, Demographic Parity
- EO, Equality of Opportunity

❑ Individual Fairness

Fairness

Fairness notion

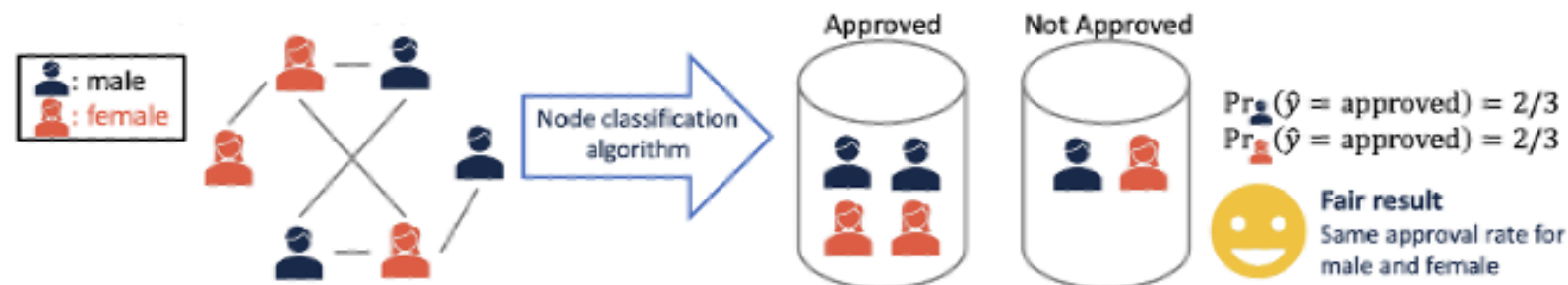
Fair & Disentangle graph mining

Fairness Metrics

Group Fairness – DP, Demographic Parity

$$\Pr_+(\hat{y} = c) = \Pr_-(\hat{y} = c)$$

- \Pr_+ : probability for the protected group
 - \Pr_- : probability for the unprotected group
- All groups should have an **equal** positive(acceptance) rate



Fairness

Fairness notion

Fair & Disentangle graph mining

Fairness Metrics

Group Fairness – **EO, Equality of Opportunity**

$$\Pr_+(\hat{y} = c | y = c) = \Pr_-(\hat{y} = c | y = c)$$

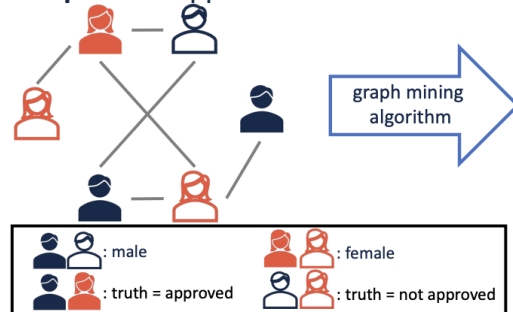
- \Pr_+ : probability for the protected group
- \Pr_- : probability for the unprotected group

→ All groups should have an **equal TPR**(true positive rate) regardless of their protected attributes

Group Fairness – **DP, Demographic Parity**

$$\Pr_+(\hat{y} = c) = \Pr_-(\hat{y} = c)$$

Example: loan approval



Approved



Not Approved



$$\Pr_{\text{male}}(\hat{y} = \text{approved} | \text{truth} = \text{approved}) = 1$$
$$\Pr_{\text{female}}(\hat{y} = \text{approved} | \text{truth} = \text{approved}) = 1$$



Fair result
Same true positive rate
for male and female

Fairness

Fairness notion

Fairness Metrics

Individual Fairness

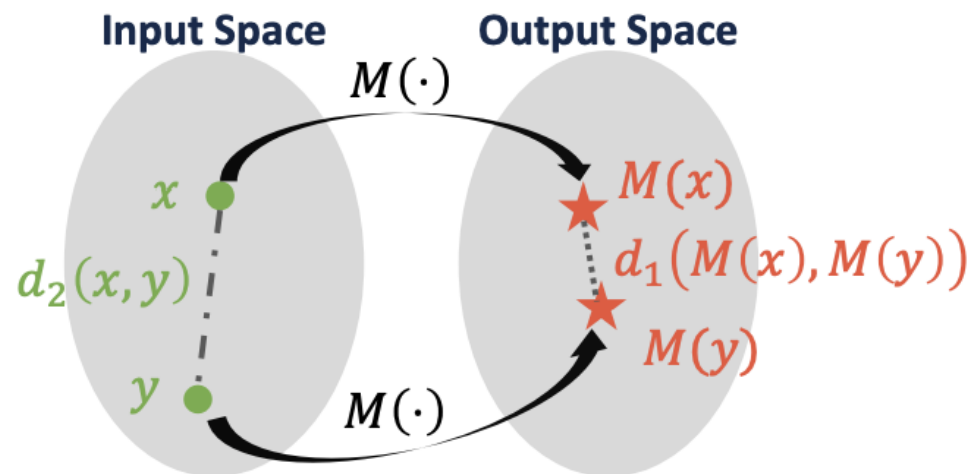
$$d_1(M(x), M(y)) \leq L d_2(x, y)$$

Lipschitz inequality

- M : a mapping from input to output
- d_1 : distance metric for output
- d_2 : distance metric for input
- L : a constant scala

→ **Similar individuals** should have **similar outcomes**

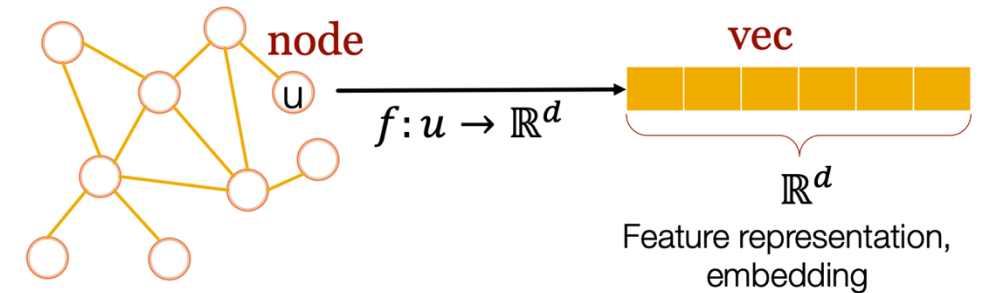
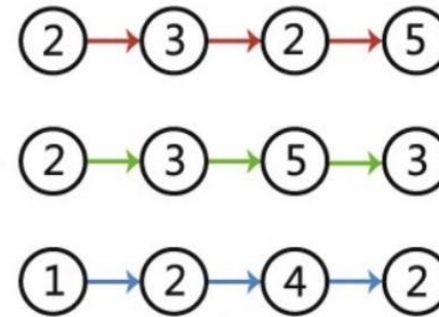
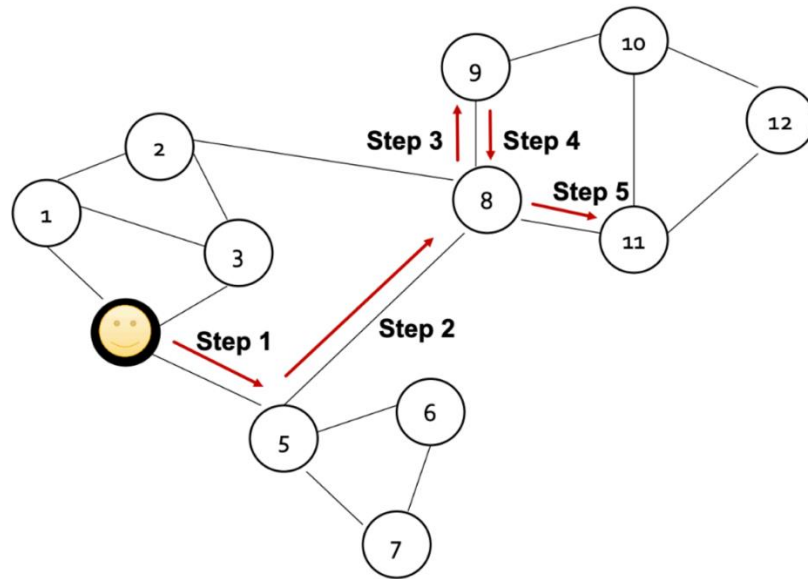
Fair & Disentangle graph mining



Fairness

Fairwalk (IJCAI' 19)

Fair & Disentangle graph mining



Random Walk

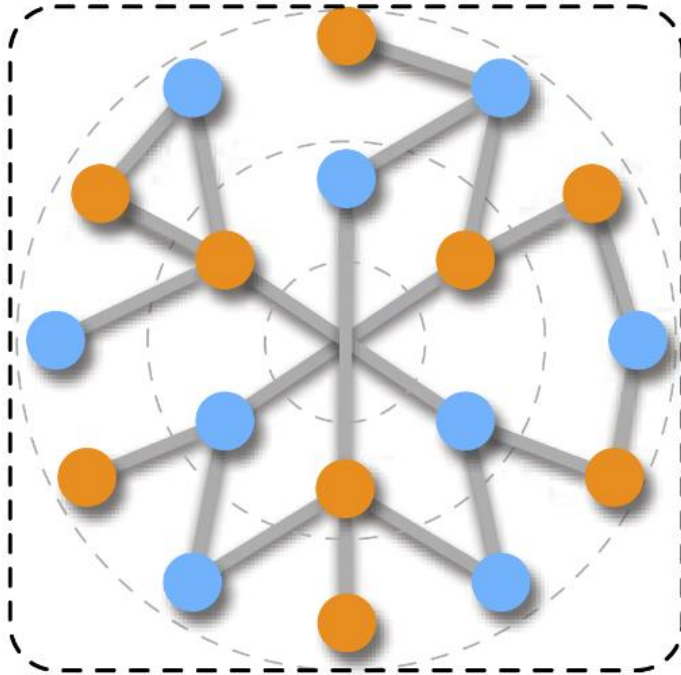
Goal : Capture **graph structure information** to generate **node embeddings**

Method: Visit neighbor node **randomly** and record the visiting order

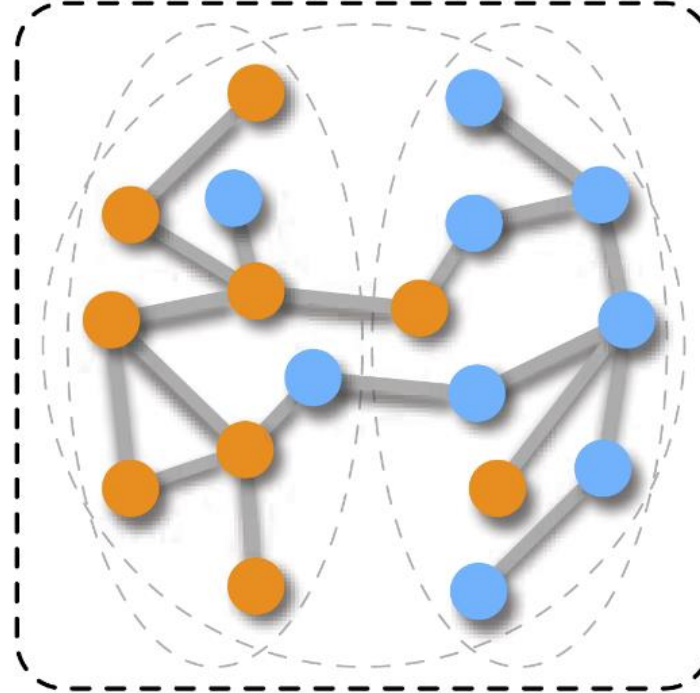
Fairness

Fairwalk (IJCAI' 19)

Fair & Disentangle graph mining



(a) Unbiased graph topology



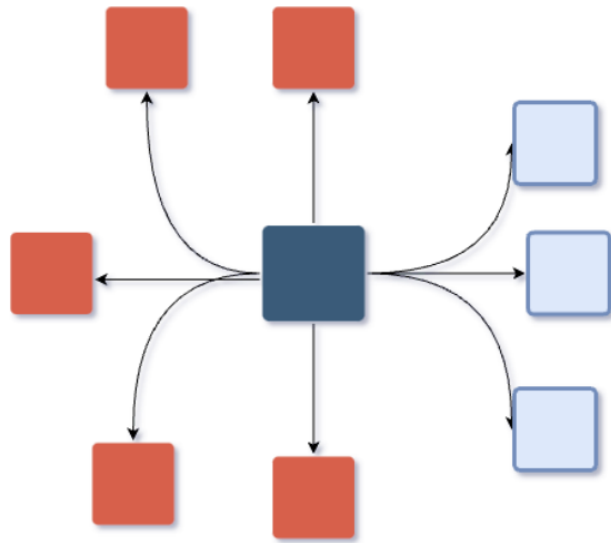
(b) Biased graph topology

Motivation:

In biased graph, random walk has a difficulty to obtain information of specific groups

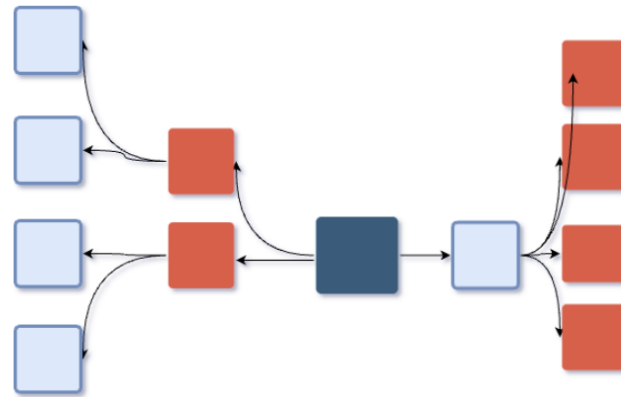
Fairness

Fairwalk (IJCAI' 19)



Probability of next walk

- Red : $5/8 \rightarrow 1/2$
- Blue: $3/8 \rightarrow 1/2$



Fair & Disentangle graph mining

Methodology:

Modify random walk process

1. Partition neighbors into groups
2. Give each group the **same probability** of being chosen regardless of their size

Limitations:

- Can not capture the information beyond the one-hop

Fairness

CrossWalk (WWW' 23)

$$m(v) = \frac{\sum_{j \in [r]} \sum_{u \in \mathcal{W}_v^j} \mathbb{I}[l_v \neq l_u]}{r \times d}.$$

$$w'_{vu} = \begin{cases} w_{vu}(1 - \alpha) \times \frac{m(u)^p}{\sum_{z \in N_v} w_{vz} m(z)^p} & \text{if } l_v = l_u \\ w_{vu} \alpha \times \frac{m(u)^p}{|R_v| \sum_{z \in N_v^c} w_{vz} m(z)^p} & \text{if } l_v \neq l_u = c. \end{cases}$$

Fair & Disentangle graph mining

Motivation:

- control more **elaborately** than fairwalk by hyper parameters
- consider fairness beyond one-hop nodes

Methodology :

Modify random walk process

- **More tendency** to the groups' peripheries and different groups

CrossWalk (WWW' 23)

$$m(v) = \frac{\sum_{j \in [r]} \sum_{u \in \mathcal{W}_v^j} \mathbb{I}[l_v \neq l_u]}{r \times d}.$$

m(v): The fraction of other groups in random walk process

➔ **High** m(v) means that the node v is **close to other groups**

Example

Did 10 iterations with random walk of 10 lengths. (Total Visit: 100)

Assume, the node v meet 30 nodes of different group with v.

Then, $m(v) = 80 / 100 = 0.8$

Fairness

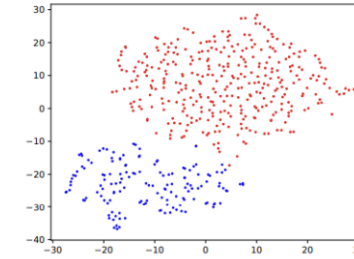
CrossWalk (WWW' 23)

Fair & Disentangle graph mining

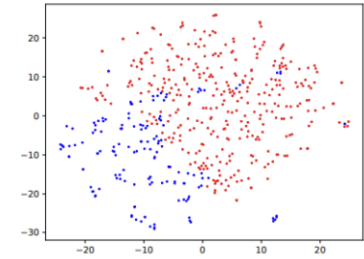
$$w'_{vu} = \begin{cases} w_{vu}(1 - \alpha) \times \frac{m(u)^p}{\sum_{z \in N_v} w_{vz} m(z)^p} & \text{if } l_v = l_u \\ w_{vu} \alpha \times \frac{m(u)^p}{|R_v| \sum_{z \in N_v^c} w_{vz} m(z)^p} & \text{if } l_v \neq l_u = c. \end{cases}$$

- ➔ Based on $m(v)$, **reweight** the edges to w_{uv}
- ➔ Then, transfer probabilities in random walks based on reweighting strategy

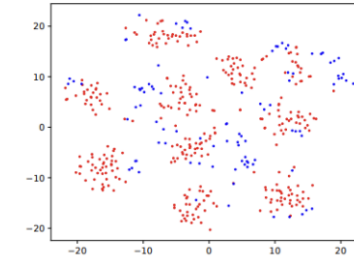
- α : manipulate strengths to the other group (The higher, the powerful)
- p : manipulate strengths to the group boundaries (The higher, the powerful)



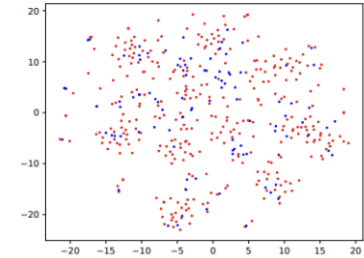
(a) Synthetic Layered Dataset
DeepWalk



(b) Synthetic Layered Dataset
CrossWalk ($\alpha = 0.5, p = 2$)



(c) Rice-Facebook Dataset
DeepWalk

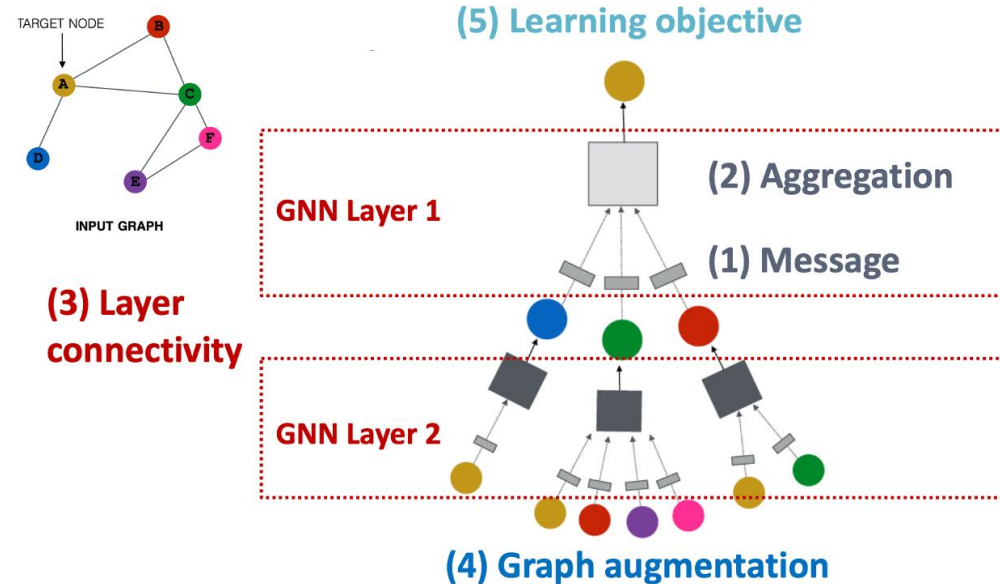


(d) Rice-Facebook Dataset
CrossWalk ($\alpha = 0.5, p = 4$)

Fairness

InFoRM (KDD '20)

Fair & Disentangle graph mining



GNN, Graph Neural Networks

Goal : To learn meaningful representations of nodes, edges, or entire graphs for tasks

Method:

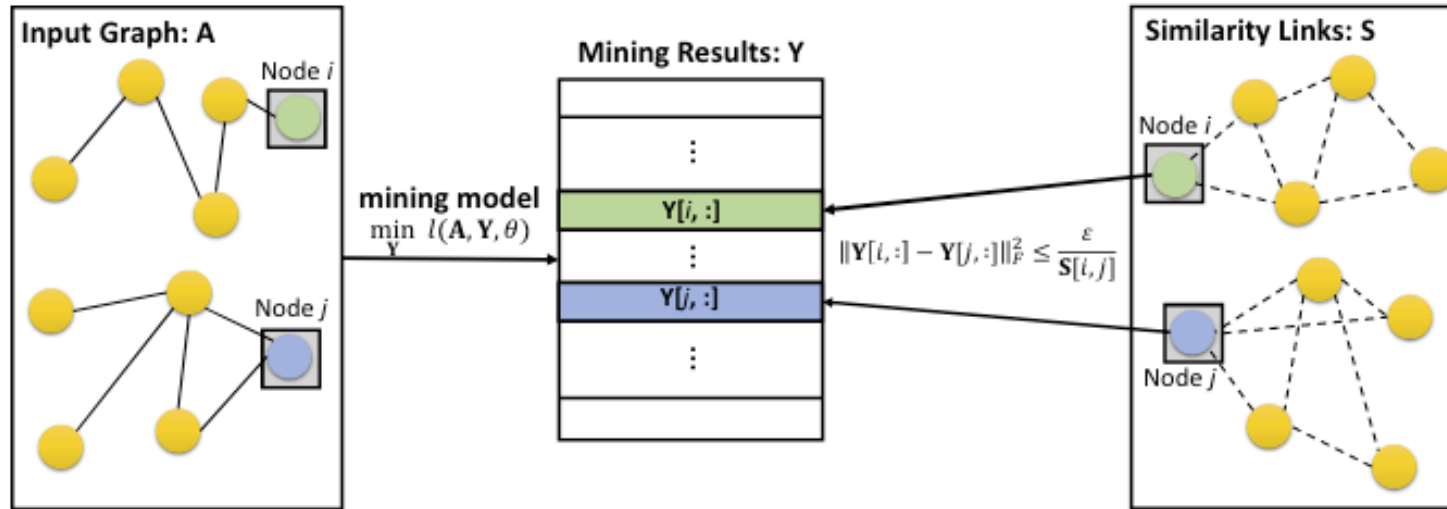
Message Passing : Each node gathers information from its neighbors

Aggregation : The gathered messages are combined using a function

Fairness

InFoRM (KDD '20)

Fair & Disentangle graph mining



Goal : Increase the **individual fairness** of graph mining tasks

Method: Add the term related with individual fairness. Then, solve the optimization problem

- Debiasing the input graph
- Debiasing mining model
- Debiasing mining results

Individual Fairness Term

$$\|Y[i, :] - Y[j, :]\|_F^2 \leq \frac{\epsilon}{S[i, j]} \quad \forall i, j = 1, \dots, n$$

- ❑ The higher the similarity, the smaller the difference
- ❑ Need to calculate about **all pairs** of i and j nodes

$$\sum_{i=1}^n \sum_{j=1}^n \|Y[i, :] - Y[j, :]\|_F^2 S[i, j] = 2\text{Tr}(Y' L_S Y) \leq m\epsilon = \delta$$

- ❑ Convert the equation into Trace function format (Can calculate **in single time**)
- ❑ The smaller the Tr term, the higher the individual fairness

InFoRM (KDD '20)

Debiasing the input graph

$$\min_{\tilde{A}} \|\tilde{A} - A\|_F^2 + \alpha \text{Tr}(Y' L_S Y) \quad \text{s.t.} \quad \partial_Y l(\tilde{A}, Y, \theta) = 0$$

Debiasing the mining model

$$Y^* = \underset{Y}{\operatorname{argmin}} \quad J = l(A, Y, \theta) + \alpha \text{Tr}(Y' L_S Y)$$

Debiasing the mining results

$$Y^* = \underset{Y}{\operatorname{argmin}} \quad J = \|Y - \bar{Y}\|_F^2 + \alpha \text{Tr}(Y' L_S Y)$$

Weekly Meetings

2. Disentanglement

- Disentangle notion
- DisenGCN
- FactorGCN

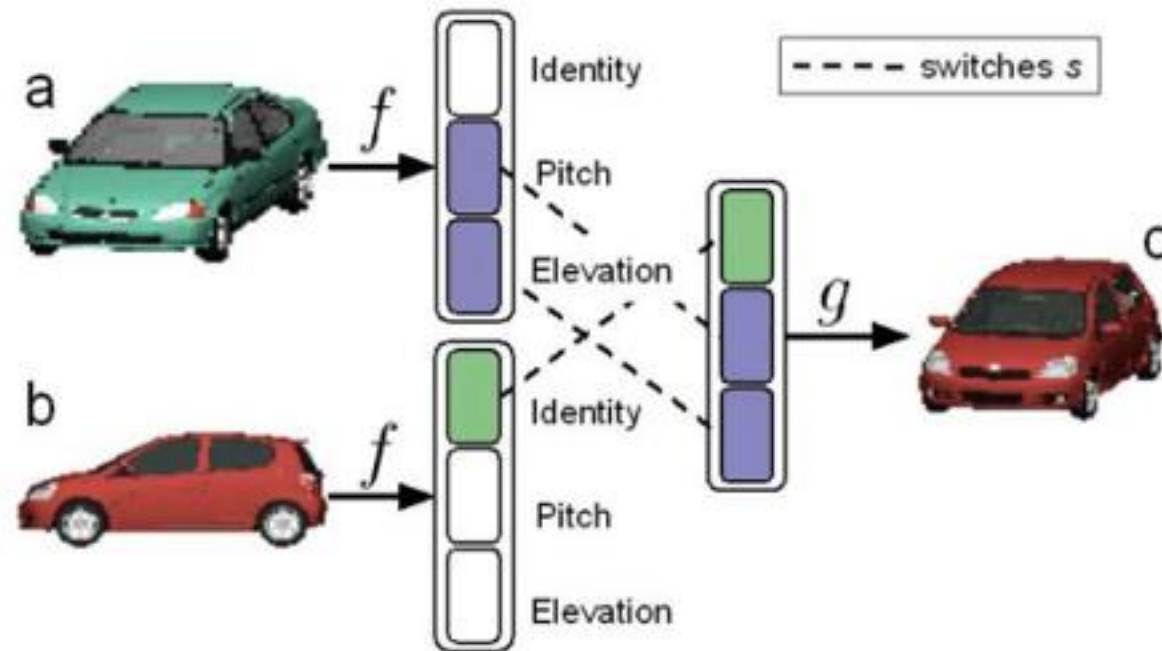
Disentanglement

Disentangle notion

Fair & Disentangle graph mining

Disentangled representation learning

→ Aims to encode **independent factors** of variation **into different dimensions** of the learned representation.



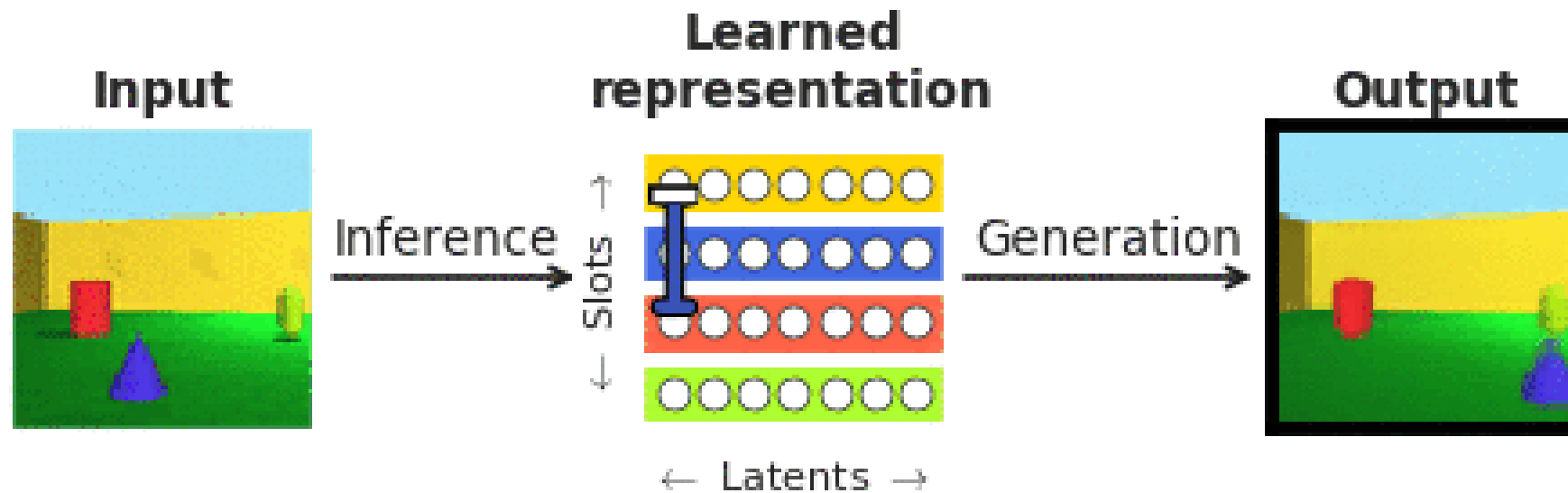
Disentanglement

Disentangle notion

Fair & Disentangle graph mining

Disentangled representation learning

→ Aims to encode **independent factors** of variation **into different dimensions** of the learned representation.



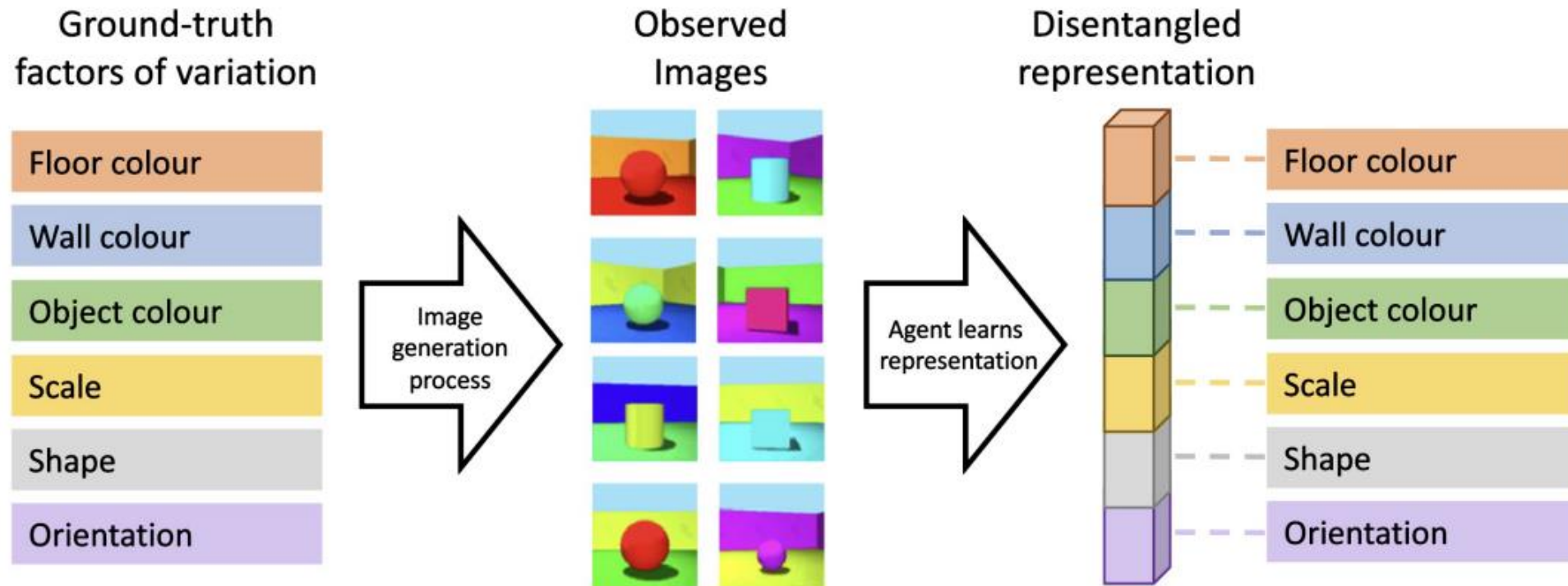
Disentanglement

Disentangle notion

Fair & Disentangle graph mining

Why do we need and What are benefits?

❑ Enhance **generalization ability** and **robustness**

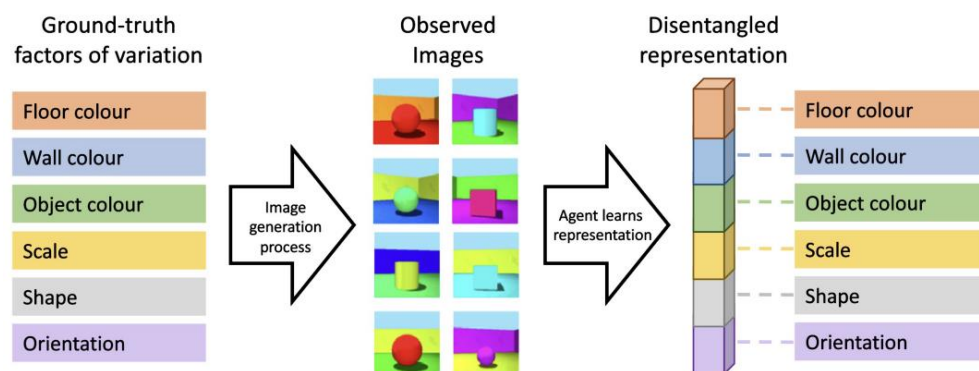


Disentanglement

Disentangle notion

Why do we need and What are benefits?

- ❑ Enhance **generalization ability** and **robustness**



- ❑ Enhances **Interpretability**

- Since we can distinguish the input

- ❑ Advances **Fairness**

- We may separate the sensitive data from input

Disentanglement

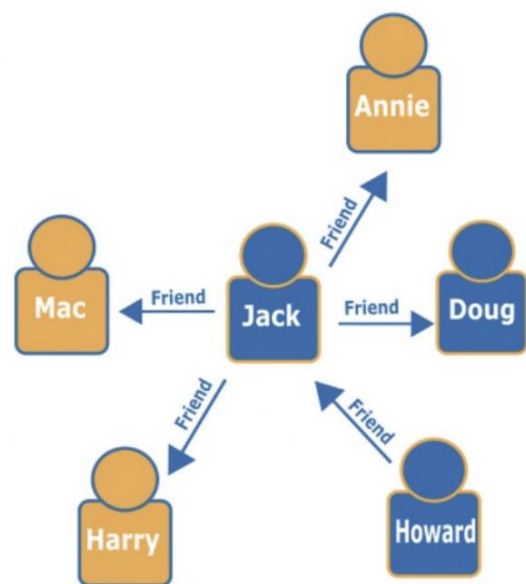
DisenGCN (ICML '19)

Fair & Disentangle graph mining

Motivation

To identify the subset of neighbors that are connected due to factor k .

e.g., Friends, co-worker, or subscribing



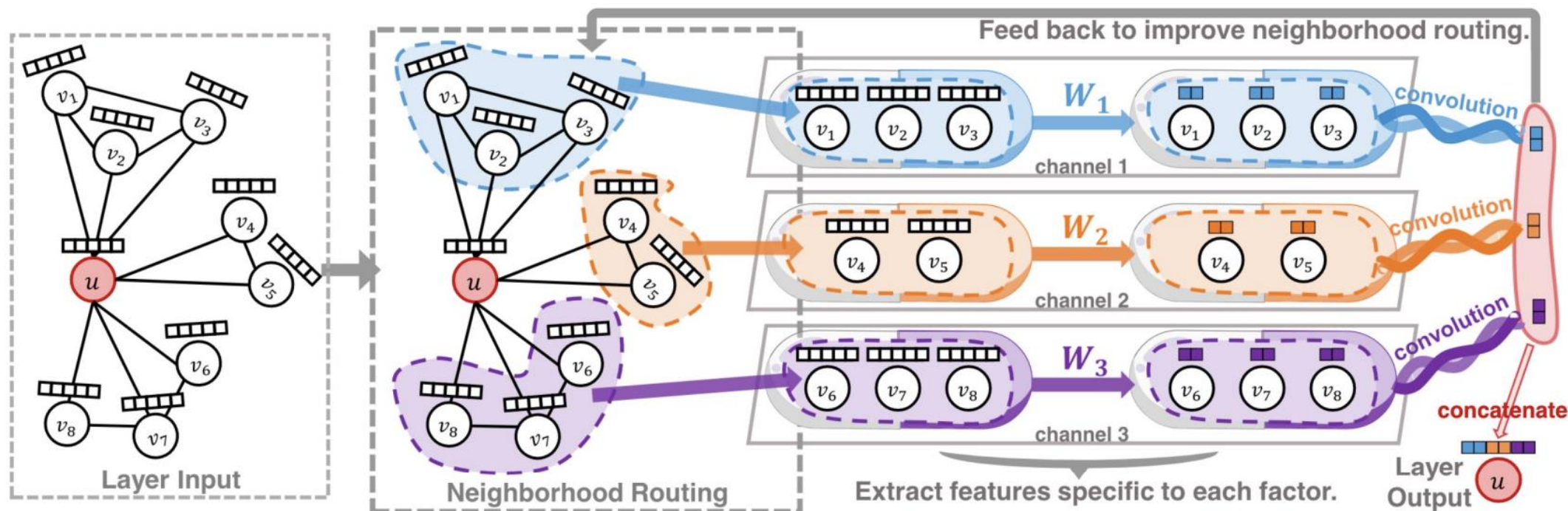
❖ Scenario

- ☐ Jack in a social network
- ☐ Mac, Harry, and Annie are high school friends
- ☐ Doug and Howard are co-workers
- ☐ Connects with others for various reasons

Disentanglement

DisenGCN (ICML '19)

Fair & Disentangle graph mining



Methodology

Neighborhood Routing : **Segments** the neighborhood according to the factors

Extract features specific : **Extract** the features from the input node

Disentanglement

DisenGCN (ICML '19)

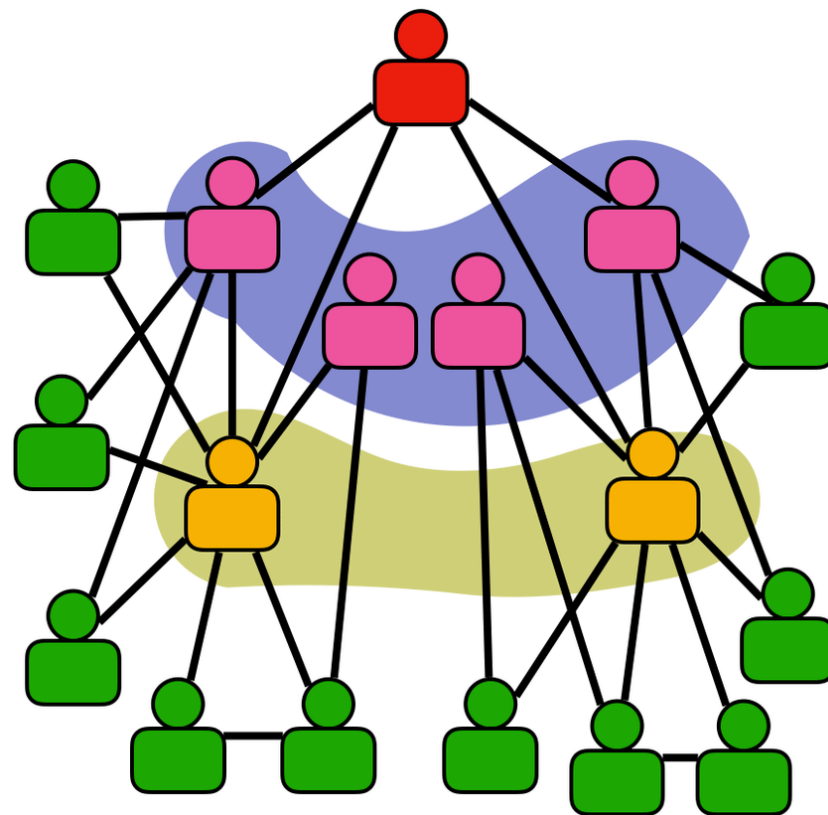
Fair & Disentangle graph mining

Output: a **disentangled** representation into **K** independent components

$$\mathbf{y}_u = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_K], \text{ where } \mathbf{c}_k \in \mathbb{R}^{\frac{d_{out}}{K}} \ (1 \leq k \leq K),$$

Hypo1: First-order proximity

→ If u and v are similar in terms of factor k ,
then the factor k is likely to be the reason why
they are connected



Hypo1: First-order proximity

→ If u and v are similar in terms of factor k,

then the factor k is likely to be the reason why they are connected

$$p_{v,k}^{(t)} = \frac{\exp(\mathbf{z}_{v,k}^\top \mathbf{c}_k^{(t)} / \tau)}{\sum_{k'=1}^K \exp(\mathbf{z}_{v,k'}^\top \mathbf{c}_k^{(t)} / \tau)}, \quad \mathbf{z}_{i,k} = \frac{\sigma(\mathbf{W}_k^\top \mathbf{x}_i + \mathbf{b}_k)}{\|\sigma(\mathbf{W}_k^\top \mathbf{x}_i + \mathbf{b}_k)\|_2},$$

$\mathbf{z}_{i,k}$: Describes the aspect of node i that are related with the k-th factor

→ Kind of a embeddings in terms of k-th factor

$\mathbf{p}_{v,k}$: Probability that factor k is the reason why u and v are connected

→ Measure how similar with u and v by multiplying two embeddings

\mathbf{c}_k : The subset of the final output about node u

Hypo2: Second-order proximity

→ if the neighbors form a large cluster in the k-th subspace, reflecting similarity with respect to factor k.

$$\mathbf{c}_k^{(t)} = \frac{\mathbf{z}_{u,k} + \sum_{v:(u,v) \in G} p_{v,k}^{(t-1)} \mathbf{z}_{v,k}}{\|\mathbf{z}_{u,k} + \sum_{v:(u,v) \in G} p_{v,k}^{(t-1)} \mathbf{z}_{v,k}\|_2},$$

$\mathbf{z}_{i,k}$: Describes the aspect of node i that are related with the kth factor

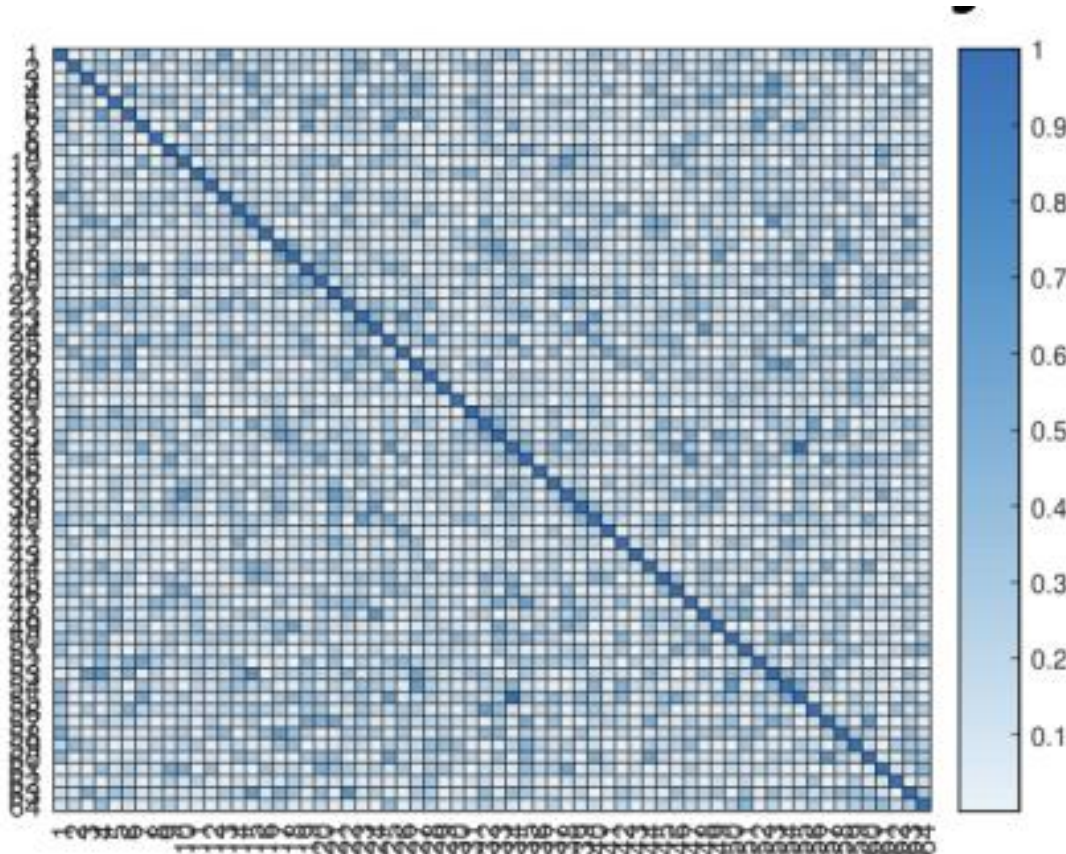
$p_{v,k}$: Probability that factor k is the reason why u and v are connected

$\mathbf{c}_{v,k}$: The final output of the node in terms of the factor k

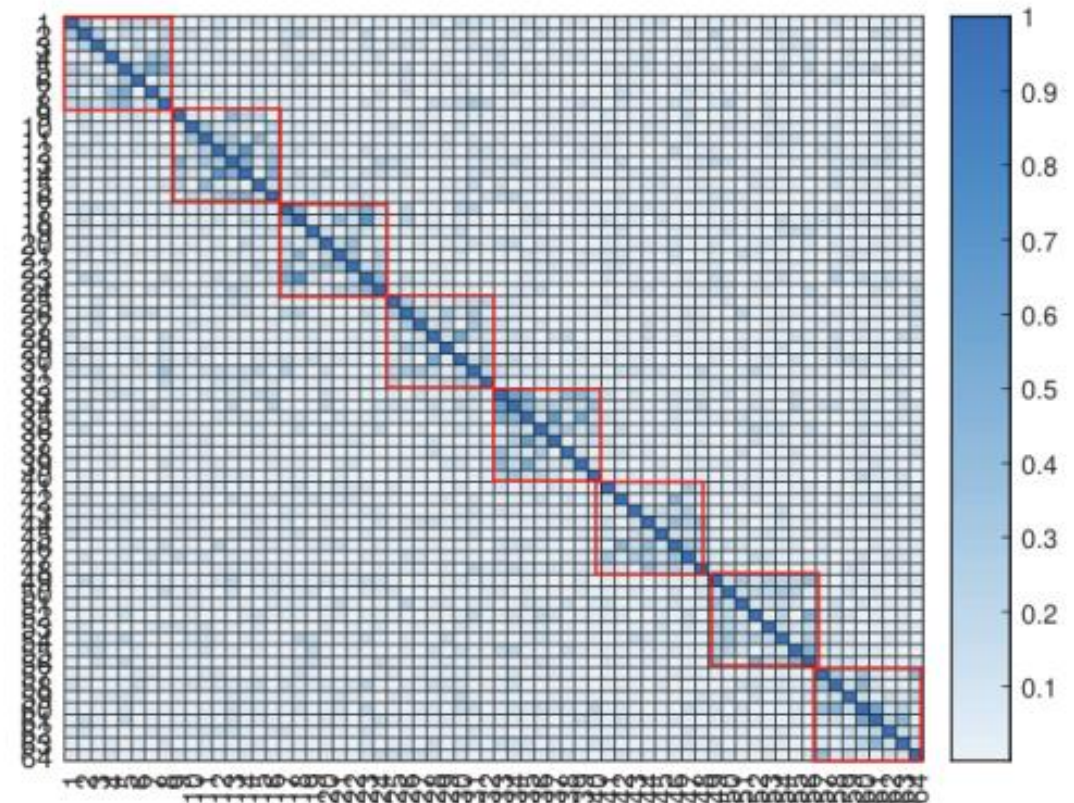
Disentanglement

DisenGCN (ICML '19)

Fair & Disentangle graph mining



(a) GCN.



(b) DisenGCN (this work).

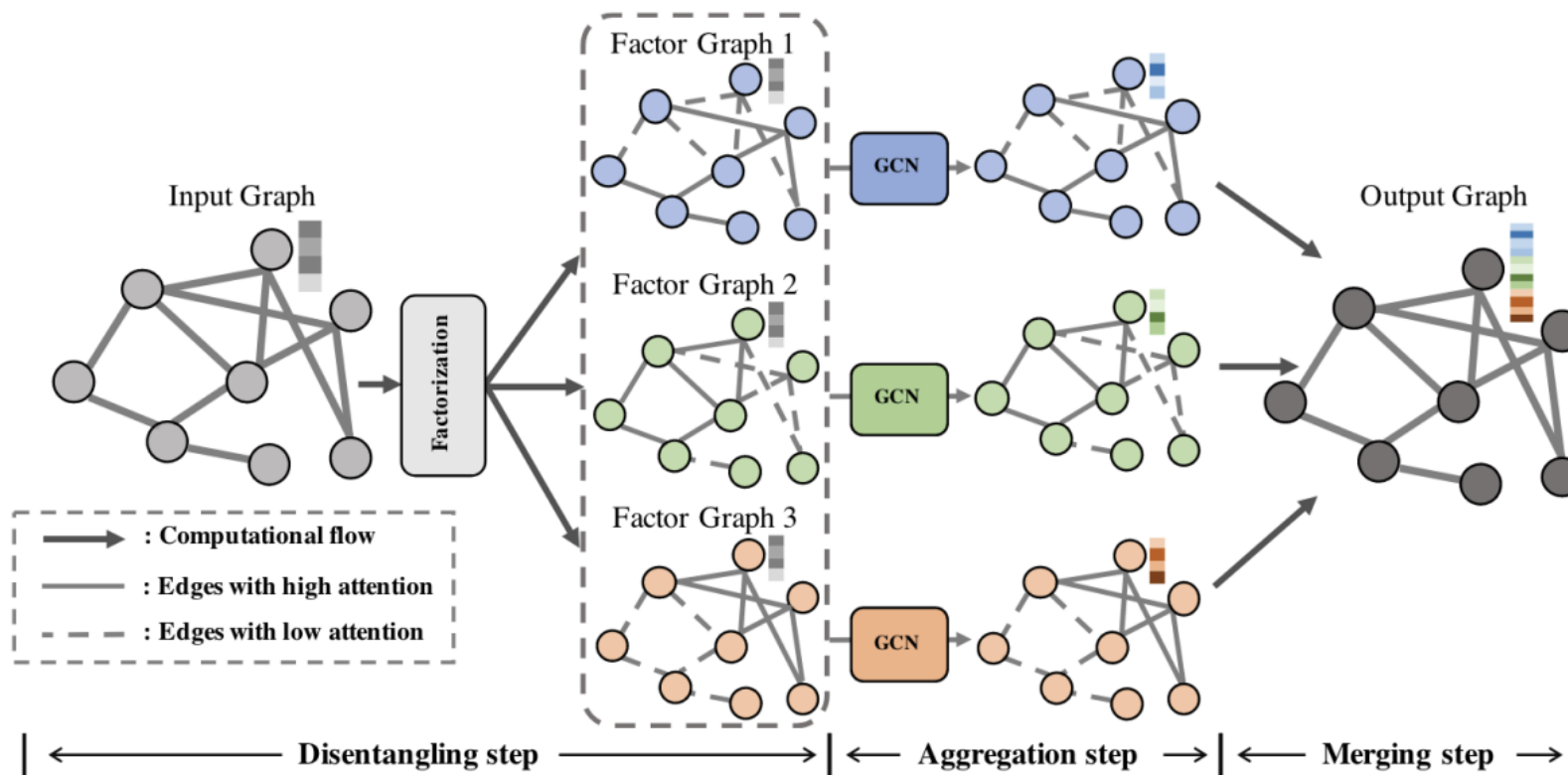
Disentanglement

FactorGCN (NerulIPS '20)

Fair & Disentangle graph mining

Motivation

- ❑ To focus **on the graph-level** partition not a node-level neighbor partition.
- ❑ To allow for **overlapping edges** where needed



Methodology – Attention score

$$E_{ije} = 1 / \left(1 + e^{-\Psi_e(h'_i, h'_j)} \right) ; h' = \mathbf{W}h,$$

- ❑ \mathbf{h} : the set of nodes with feature of F dimension
- ❑ \mathbf{W} : a linear transformation matrix
- ❑ Ψ : computes **the attention score** for factor graph e (**one-layer MLP**)
 - As Ψ increases, E_{ije} decreases, and vice versa
- ❑ E_{ije} : the coefficient of edge
- ❑ Notice there are no **softmax** → The sum of the E **does not need to be 1**.

FactorGCN (NerulPS '20)

Methodology – Aggregation & Merge & Loss function

Aggregation

$$h_i^{(l+1)e} = \sigma\left(\sum_{j \in \mathcal{N}_i} E_{ije}/c_{ij} h_j^{(l)} \mathbf{W}^{(l)}\right), c_{ij} = (|\mathcal{N}_i| |\mathcal{N}_j|)^{1/2},$$

Merge

$$h_i^{(l+1)} = \big\|_{e=1}^{N_e} h_i^{(l+1)e},$$

Loss Function

$$\mathcal{L} = \mathcal{L}_t + \lambda * \mathcal{L}_d$$

Methodology – Separating

$$G_e = \text{Softmax}\left(f\left(\text{Readout}(\mathcal{A}(\mathbf{E}_e, \mathbf{h}'))\right)\right).$$
$$\mathcal{L}_d = -\frac{1}{N} \sum_i \left(\sum_{c=1}^{N_e} \mathbb{1}_{e=c} \log(G_i^e[c]) \right),$$

N_e : the number of factor graphs

- ❑ Without **any other constraints**, some factor graphs **may become similar**.
- ❑ Need to be distinguished from the rest.
- ❑ By assigning unique labels to the factor graphs
and optimizing them as a graph classification problem **like a discriminator**
- ❑ **A classifier f**: consist of one fully connected layer
- ❑ **$G_i[c]$** : represents the probability that the generated factor graph has label c

Weekly Meetings

3. Invariant Learning^[SEP] & Causal-based Learning

- Invariant Learning^[SEP]
- Causal-based Learning^[SEP]
- DIR^[SEP]
- DisC

Invariant Learning

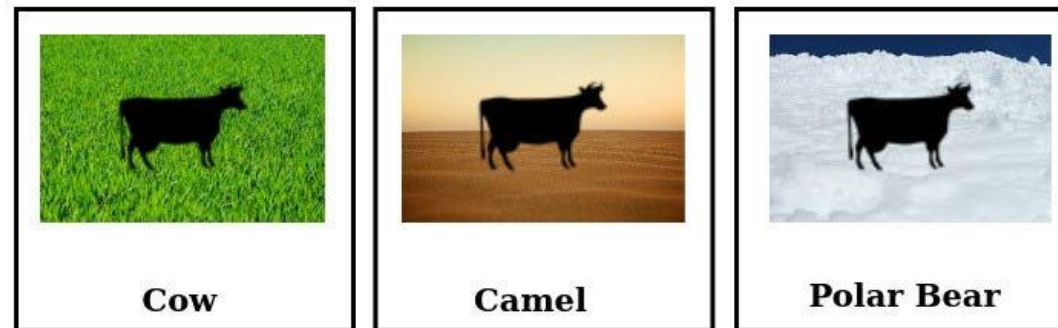
Invariant Learning

What: Aims to identify patterns that **remain consistent**
across different environments

How: By optimizing for predictive **performance that remains stable**
across multiple environments or domains.

Benefits: Enhances **Robustness** to Distribution Shifts (Better performance in OOD)

Reduces **Overfitting** to Spurious Correlations



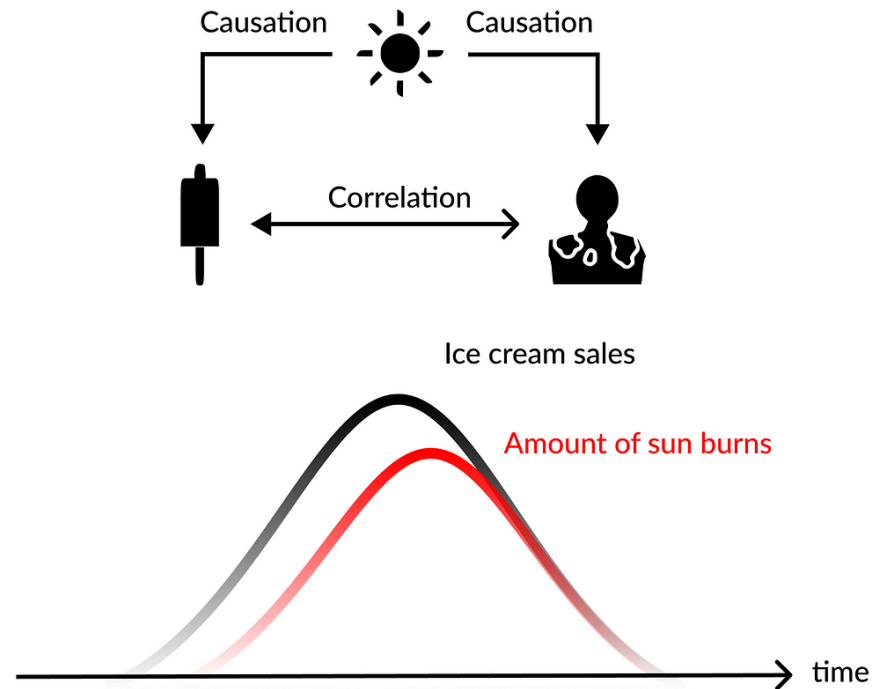
Neural Network Predictions

Causal-based learning

Causal-based learning

What:

focuses on uncovering **cause-effect** relationships



Causal-based learning

Causal-based learning

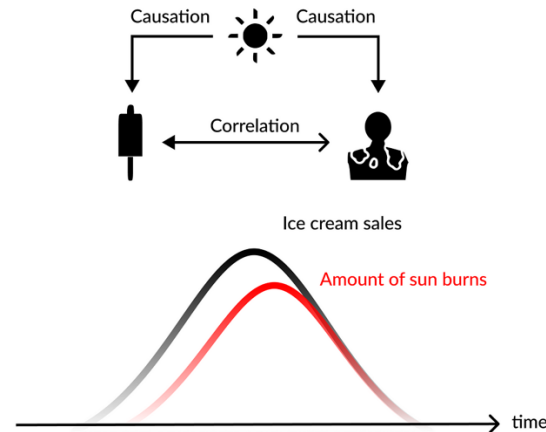
How:

identifies and leverages **causal structures** in the data

Benefits:

Enhances **Generalization**

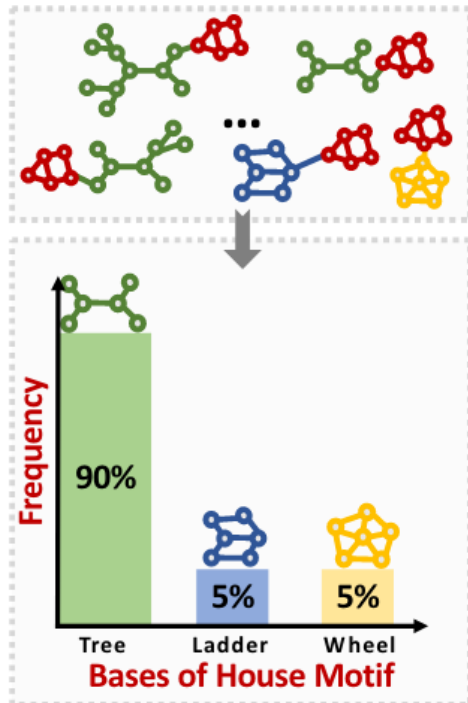
Improve **Fairness** & Facilitates **Interpretability**



DIR (ICLR '22)

Motivation

- ❑ Risk of learning from the **statistical shortcuts** can lead **poor generalization**
- ❑ Aim to identify rationales that capture **the environment-invariant causal patterns**



Example

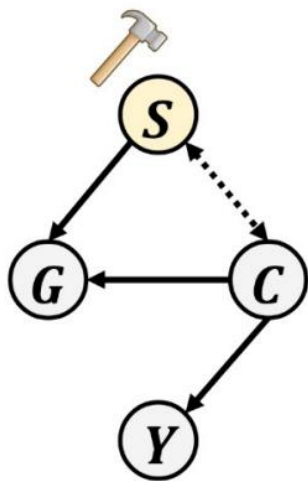
Want to classify **the house motif (the red ones)**

If the model capture **the tree motif** as a sign of house motif, the model might have **poor generalization**

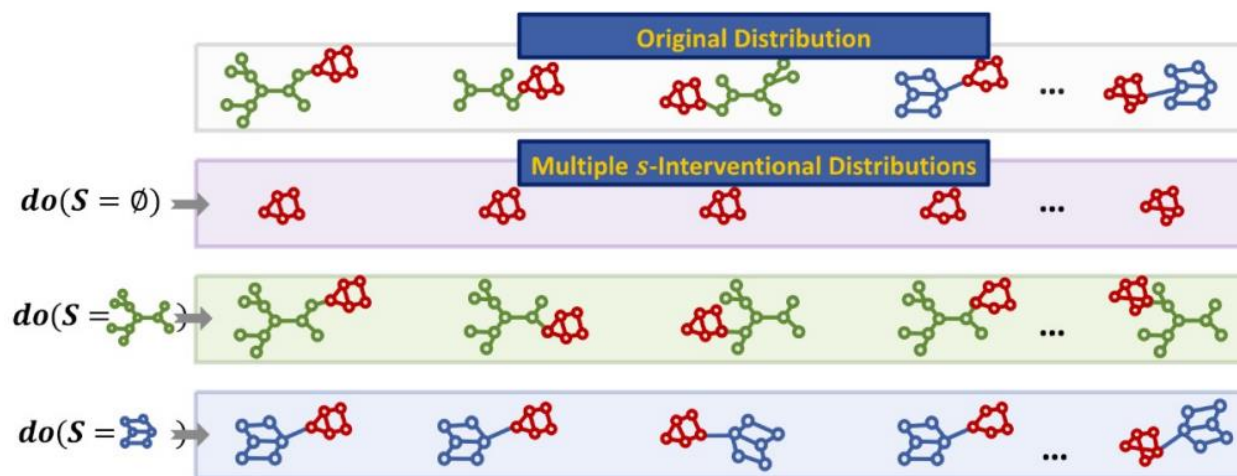
How

□ By Do-Calculus developed by Pearl

- **S: Shortcut (\approx suspicious) Part**
- Introduce Interventional Distribution (Iteratively replace the **S**)



(a) SCM

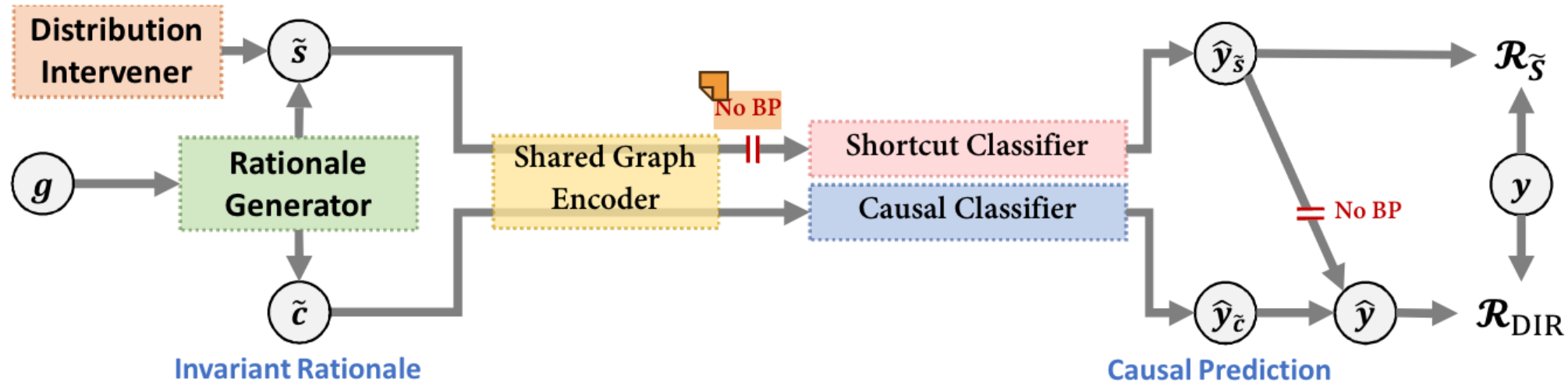


(b) Interventional Distributions.

I.R. & Causal Based

DIR (ICLR '22)

Fair & Disentangle graph mining

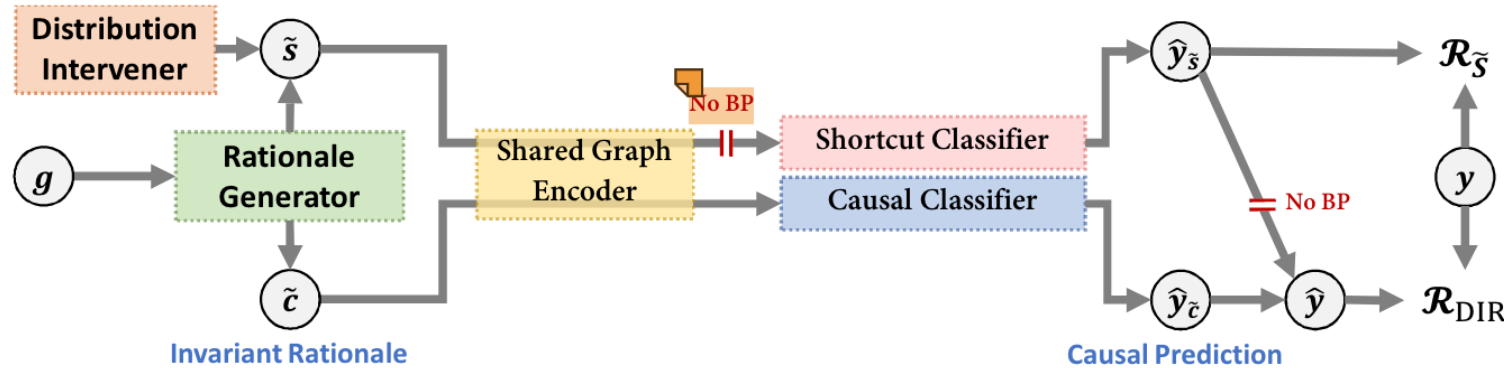


Rational Generator : Split the input graph into causal and non causal parts by **GNN**

Distribution Intervener : Creating interventional distributions. Memory bank of **S**

Shared Encoder : Generate node representations independently

Two Classifiers : Make a probability distribution over class labels by **GNN**



Loss function

$$\min \mathcal{R}_{\text{DIR}} = \mathbb{E}_s[\mathcal{R}(h(G), Y | do(S = s))] + \lambda \text{Var}_s(\{\mathcal{R}(h(G), Y | do(S = s))\}),$$

Try to **minimize** both the **risk itself** and the **variability of risk** in shortcut environments by **(do(S=s))**

$$\mathcal{R}(h(G), Y | do(S = \tilde{s})) = \mathbb{E}_{(g, y) \in \mathcal{O}, S = \tilde{s}, C = h_{\tilde{c}}(g)} l(\hat{y}, y), \quad \hat{y} = \hat{y}_{\tilde{c}} \odot \sigma(\hat{y}_{\tilde{s}}),$$

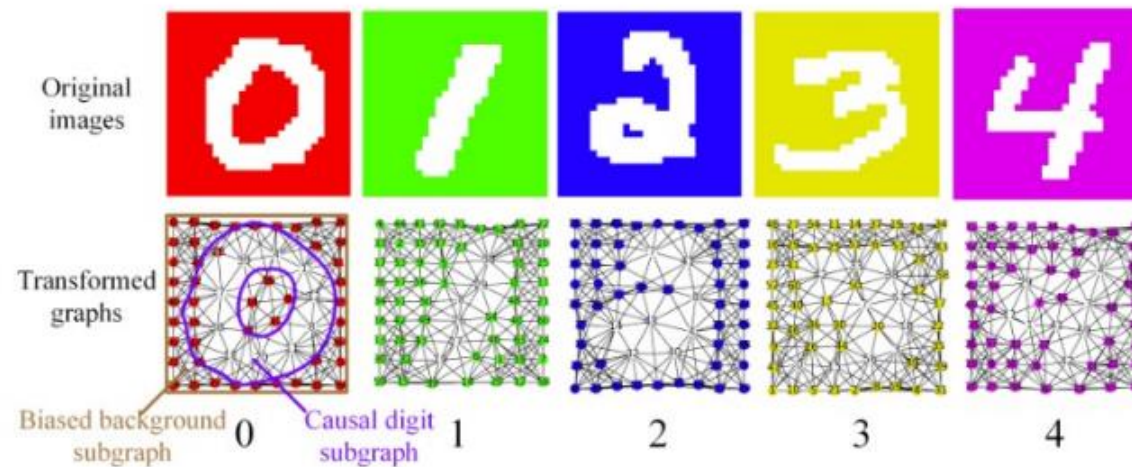
Risk: The **expectation(Mean)** of the loss between \hat{y} and y through shortcut environments

DisC (NeurIPS '22)

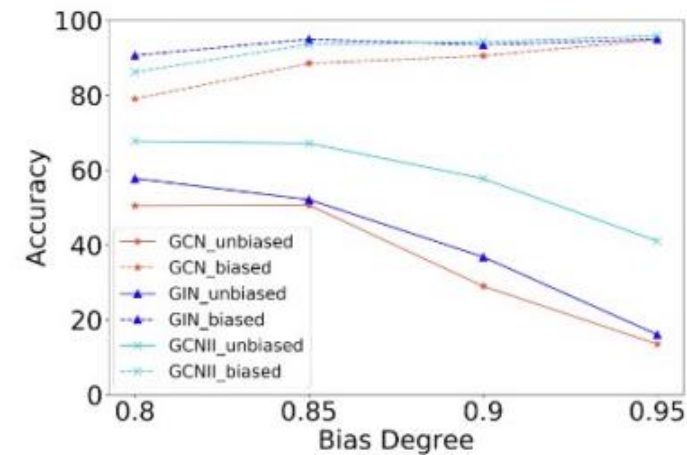
Motivation

- GNNs tend to learn bias information, especially **in severe bias situation**
- **Assume** bias part usually has **simpler structure** than **meaningful causal part**
 - ➔ Therefore, the **bias part** is **easy to learn** than causal part

Example

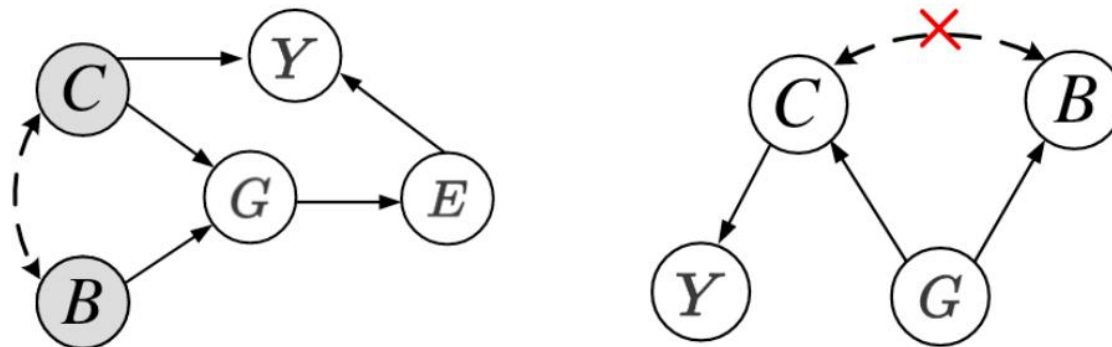


(a) Examples of graphs in CMNIST-75sp.



(b) Performance of GNNs.

How



SVM, Structural Causal Model on GNNs' prediction process

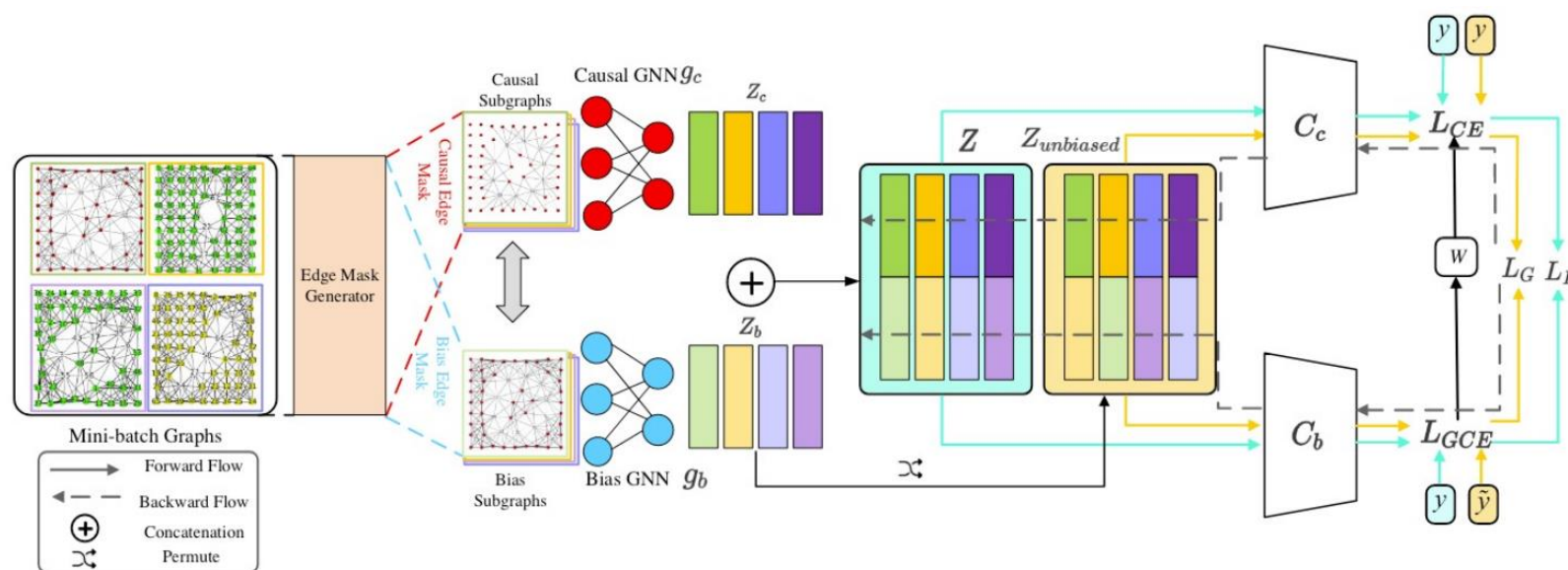
Problem & Solving

- ❑ We do not know the actual causal and bias part
 - ➔ Try to **disentangle** the latent variables C and B from the input G
- ❑ There are two paths that would induce the spurious correlation
 - ➔ **Cut off** the correlation between C and B

I.R. & Causal Based

DisC (NeurIPS '22)

Fair & Disentangle graph mining



Edge mask generator

: **Generate** Causal and Bias substructure

Two separate GNN modules

: **Learn** Disentangled Graph Representations

Unbiased Sample Generation

: **Permute** the bias representations (**Intervention**)

How to ensure they are **causal** and **bias** sub graph, respectively?

- Use the **assumption** that is “The bias part is easy to learn”
- Train the **bias** classifier by using **GCE**, **generalized cross entropy**

$$GCE(C_b(z; \alpha_b), y) = \frac{1 - C_b^y(z; \alpha_b)^q}{q},$$

- **C(•)** are softmax **output** of the bias classifier
- **High C(•)** can be interpreted by **high confidence** about its prediction
- **High C(•)** leads a **low GCE**. That is high confidence leads a low GCE.
- This loss function encourages the bias classifier to focus on data that is high confidence and easy to learn. (it might be bias ones)

DisC (NeurIPS '22)

Train the **causal** classifier by using **Unbias score, $W(z)$**

$$W(z) = \frac{CE(C_b(z), y)}{CE(C_c(z), y) + CE(C_b(z), y)}.$$

- **High CE** from bias classifies means that the **C_b fail to predict** on data z
- **Relatively** high CE loss from C_b than C_c can be regarded as the unbiased
 - **Relatively**, C_b fail to predict and C_c success to predict

$$L_D = W(z)CE(C_c(z), y) + GCE(C_b(z), y).$$

- Combine with unbiased score and CE from C_c for **causal** classifier

Problem & Solving

- ❑ We do not know the actual causal and bias part
 - ➔ Try to **disentangle** the latent variables C and B from the input G
- ❑ There are two paths that would induce the spurious correlation
 - ➔ **Cut off** the correlation between C and B

Until now, we did a first part: disentangle the latent variables

➔ We need to **cut off** the correlation between **casual** and **bias** variables ($C \leftrightarrow \text{X} \rightarrow B$)

➔ Use **Do-calculus**; Randomly **permute(change)** bias part

$$L_G = W(z)CE(C_c(z_{unbiased}), y) + GCE(C_b(z_{unbiased}), \hat{y}),$$

$$L = L_D + \lambda_G L_G,$$

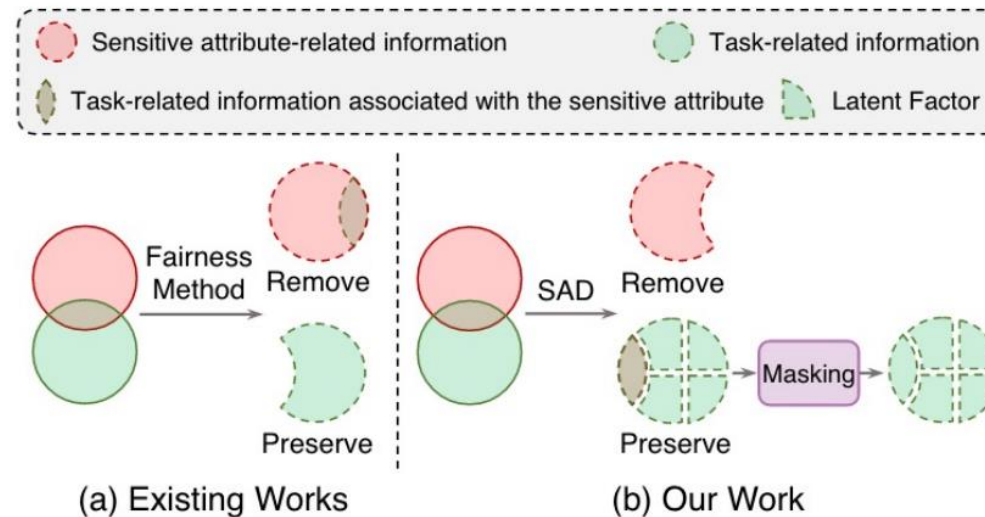
Weekly Meetings

4. Fairness & Disentanglement

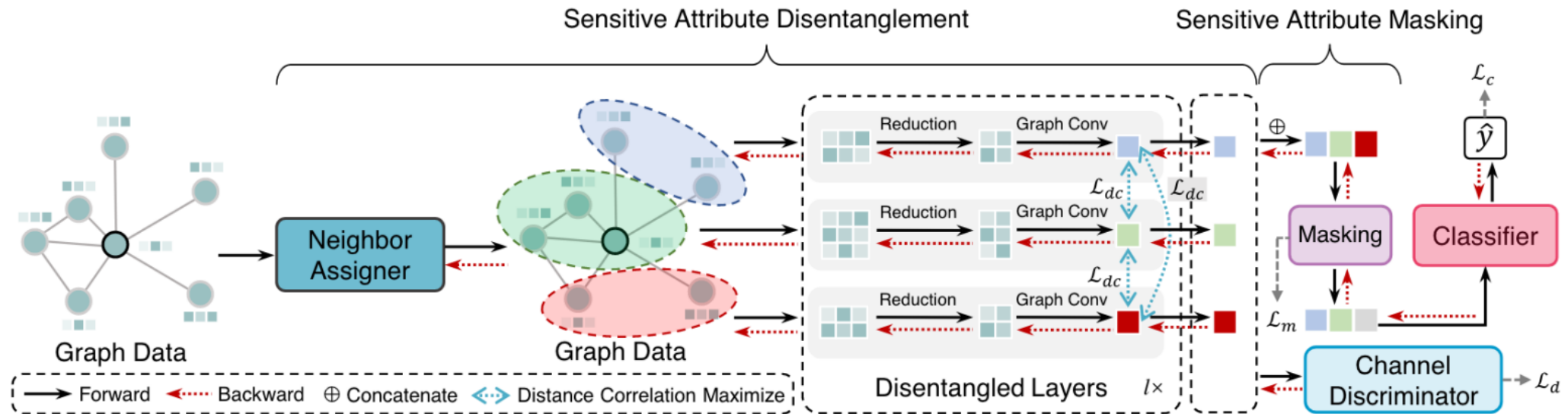
- FairSAD
- FairINV

Motivation

- ❑ Focus on **addressing performance degradation** when improving fairness.
- ❑ Employ the **Disentangle Learning** for two potential advantages
 - **Reduces correlations** between the sensitive attribute and others
 - Simplifies downstream tasks and leads to **better utility performance**



How



SAD, Sensitive Attribute Disentanglement

Disentangles the sensitive attribute into independent components

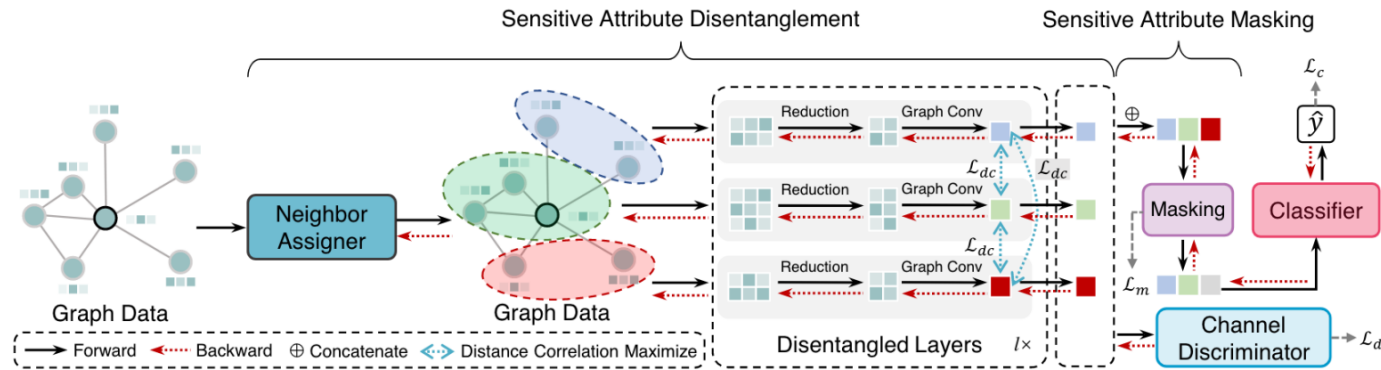
Sensitive attribute masking

Employs a **channel masking** to identify the sensitive attribute

Fair & Disentangle

FairSAD (WWW '24)

Fair & Disentangle graph mining



SAD

Neighbor Assigner : **Separate** the sensitive attribute **by MLP**

Disentangled layers : Perform graph convolution in multi-channel

Output : $H = [\underbrace{c_1, c_2, c_3, \dots, c_{d_h-2}}_{Z^1}, \underbrace{c_{d_h-1}, c_{d_h}}_{Z^k}]$ (each channel consist of three columns)

Sensitive Attribute Masking

Try to **mask(remove)** the sensitive column

$$\begin{aligned}\tilde{H} &= H \odot m = [c_1 m_1, c_2 m_2, c_3 m_3, \dots, c_{d_h} m_{d_h}] \\ &= [\underbrace{\tilde{c}_1, \tilde{c}_2, \tilde{c}_3, \dots, \tilde{c}_{d_h-2}}_{\tilde{Z}^1}, \underbrace{\tilde{c}_{d_h-1}, \tilde{c}_{d_h}}_{\tilde{Z}^k}],\end{aligned}$$

$$\min_{\theta} \mathcal{L} = \mathcal{L}_c + \alpha(\mathcal{L}_{dc} + \mathcal{L}_d) + \beta \mathcal{L}_m,$$

Optimization Purposes

- **Downstream Tasks** : Ensure that the learned representations is informative
- **Disentanglement** : Ensure the independence between latent factors
- **Decorrelation** : Weaken the impact of the sensitive attribute-related component

Decorrelation

- **Goal:** Aims to assign the minimum masking value to the sensitive attribute-related component
- **How:**

$$\mathcal{L}_m = \sum_{i=1}^{d_h} |Cov(\mathbf{s}, \tilde{\mathbf{c}}_i)| = \sum_{i=1}^{d_h} |\mathbb{E}[(\mathbf{s} - \mathbb{E}(\mathbf{s}))(\tilde{\mathbf{c}}_i - \mathbb{E}(\tilde{\mathbf{c}}_i))]|,$$

Motivation

- Most works **necessitates prior knowledge** of considered sensitive attributes
- Need to **re-training from scratch** when faced **with fairness requirement alternations**
- ➔ Need to train fair GNNs across various sensitive attributes **in a single training session**

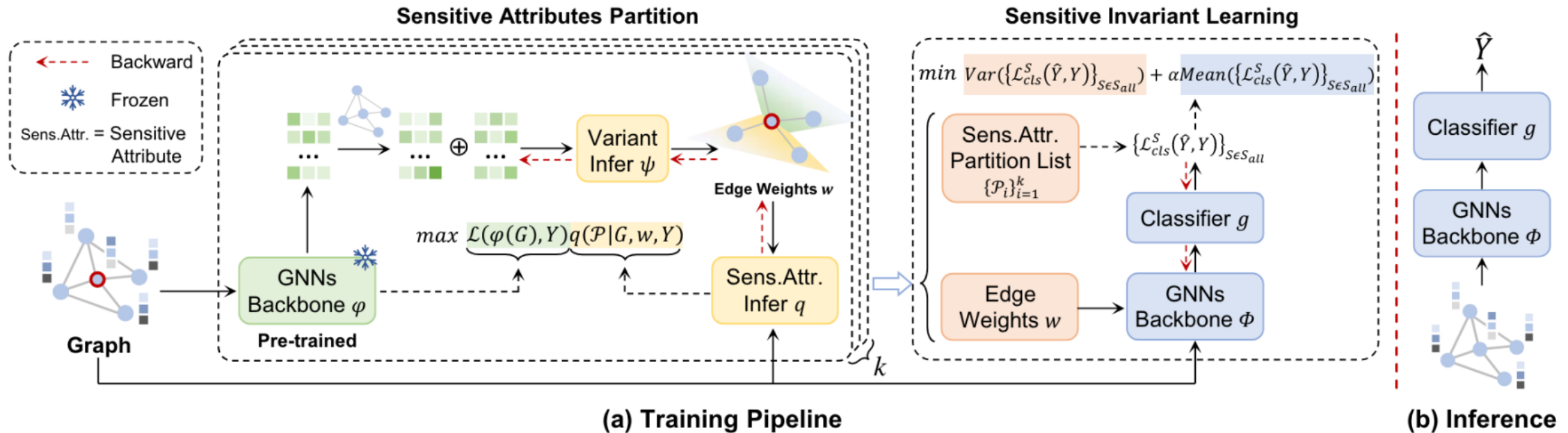
How

- The core idea behind **Invariant Learning** is similar with **Fairness Learning**
- The goal of **Invariant Learning** is to treat different **environments** equally
- The goal of **Fairness Learning** is to treat different **demographic groups** equally

Fair & Disentangle

FairINV (KDD '24)

Fair & Disentangle graph mining



SAP

Try to **automatically partitions** nodes into different subset

SIL

Learns a **GNN** invariant across different sensitive attribute partitions

SAP Goal:

- Try to capture **variant patterns** that result in significant performance differences across different environments
 - **Hope** that the environments is related with **sensitive attribute**

SAP How:

$$\max_{\theta_\psi, \theta_q} \|\nabla_{\overline{\mathbf{w}}} \mathcal{R}^S(\overline{\mathbf{w}} \circ \varphi, q)\|,$$

$$\mathcal{R}^S(\varphi, q) = \sum_{v \in \mathcal{V}} \mathbf{q}_v(S) \mathcal{L}(\varphi(\mathcal{G}_v), y_v),$$

$$\mathbf{q}_v(S) : q_v(S | \mathcal{G}_v, \mathbf{w}^i, y_v)$$

- $\mathbf{q}_v(\mathbf{S})$ means some soft partition
- By **calculating** and **maximizing the gradient of \mathcal{R}^S**
- $\mathbf{q}_v(\mathbf{S})$ is optimized in a way that partitions to maximize the performance differences across environments.

SIL Goal:

- Aims to **minimize the performance difference** between various sensitive attribute groups and ensures **the predicted accuracy** across all sensitive attribute groups

$$\min_{\theta_f} Var(\{\mathcal{L}_{cls}^S(\hat{Y}, Y)\}_{S \in \mathcal{S}_{all}}) + \alpha Mean(\{\mathcal{L}_{cls}^S(\hat{Y}, Y)\}_{S \in \mathcal{S}_{all}}),$$

- The sensitive attribute group S is derived from $\mathbf{q}(\mathbf{S})$
- \mathbf{L}_{cls}^S is the classification loss function under \mathbf{S}

- Until Now, I focus on **disentangled way to improve the fairness**
- However, there are many ways to improve **the fairness**
 - Could be extended to other ways like **counterfactual, adversarial, regularization, and others**

