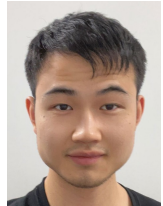


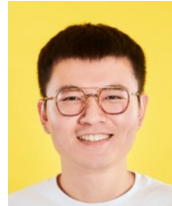
Data Quality-Aware Graph Machine Learning



Yu Wang^{1,6}



Yijun Tian²



Tong Zhao³



Xiaorui Liu⁴



Jian Kang⁵

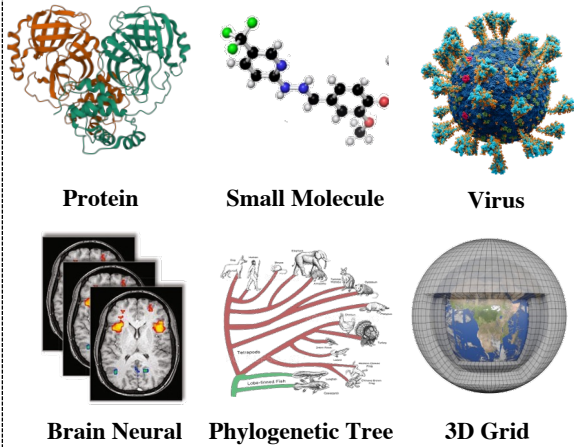


Tyler Derr¹

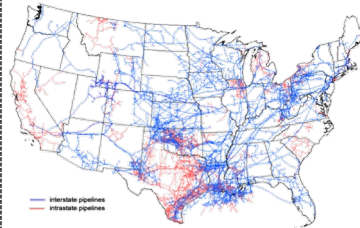
Vanderbilt University¹
University of Notre Dame²
Snap Research³
North Carolina State University⁴
University of Rochester⁵
University of Oregon⁶

Introduction and Background - Graph-Structured Data is Everywhere

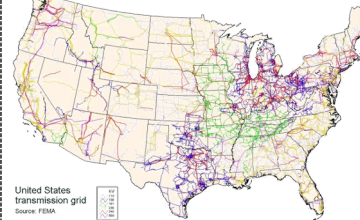
Scientific Graph



Gas Network



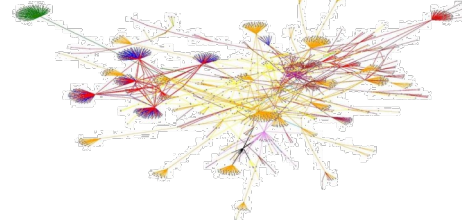
Power Network



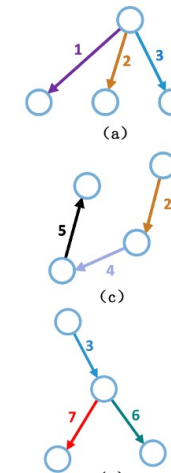
Comcast Nationwide Fiber Optic Network



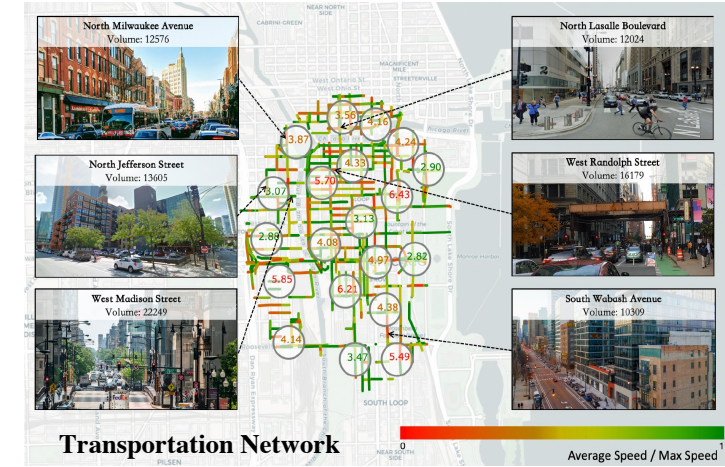
Communication Network



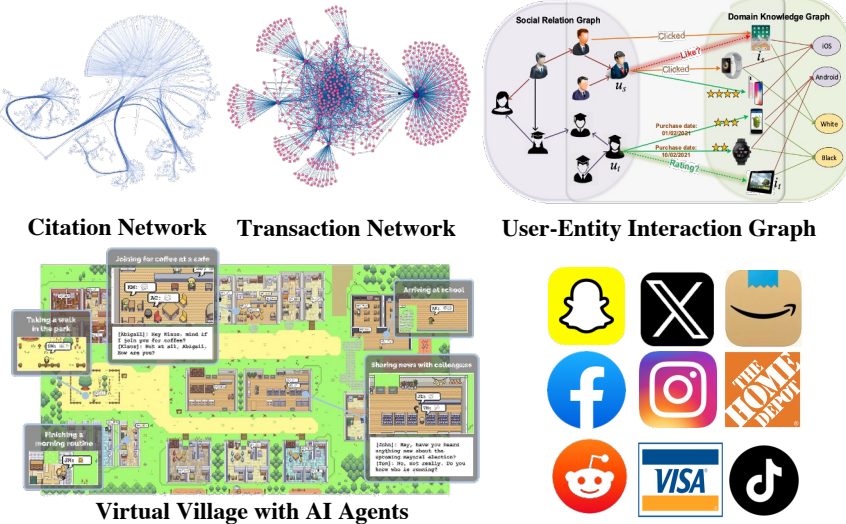
Traffic Trace



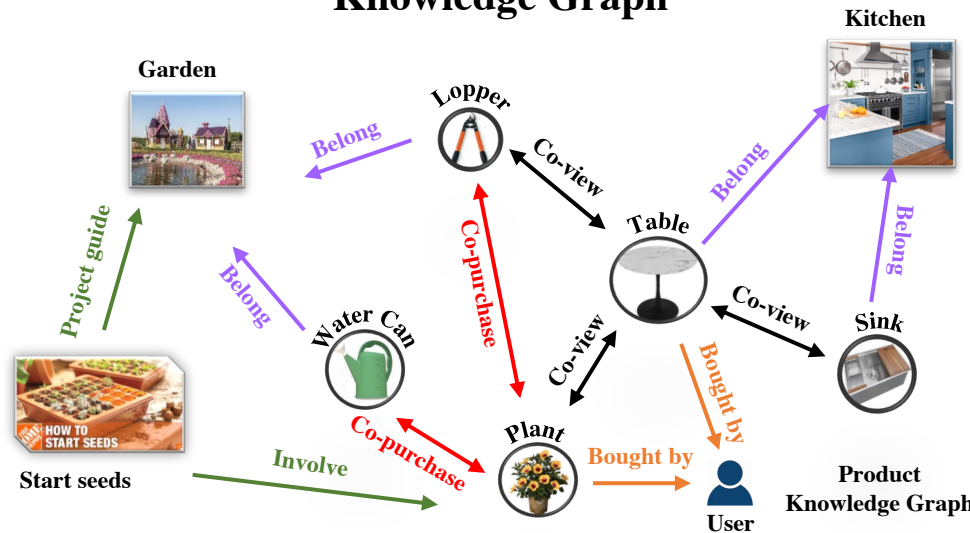
Infrastructure Graph



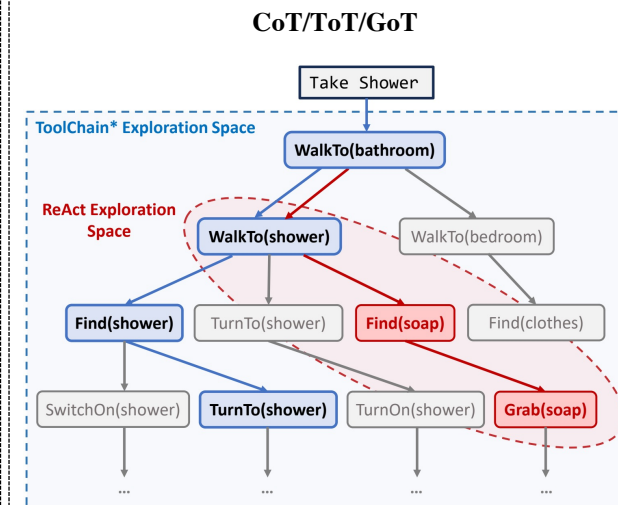
Social Interaction Graph



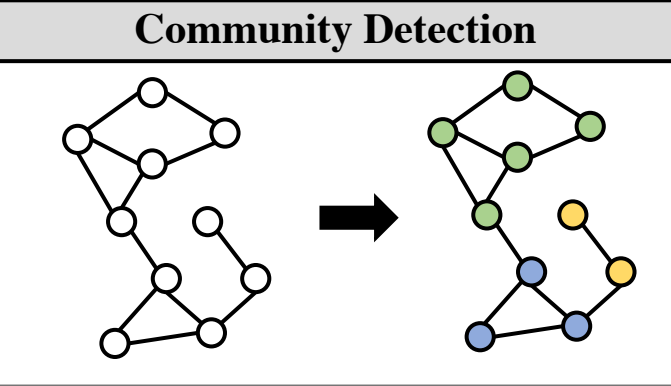
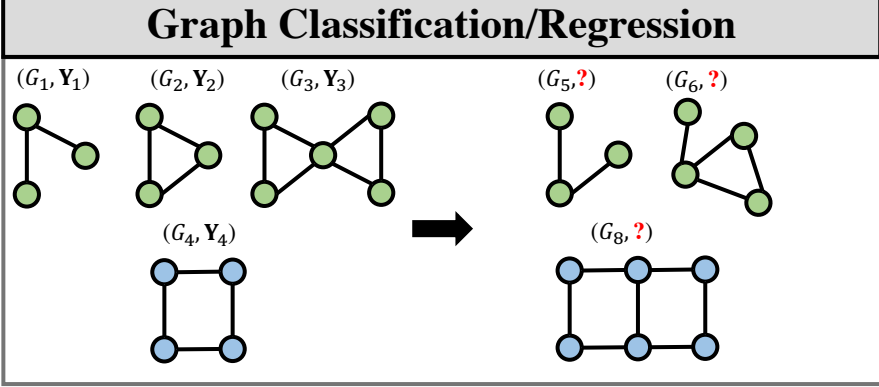
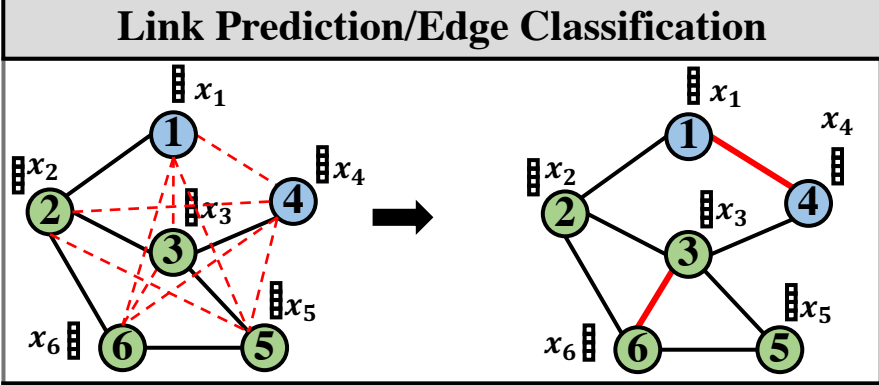
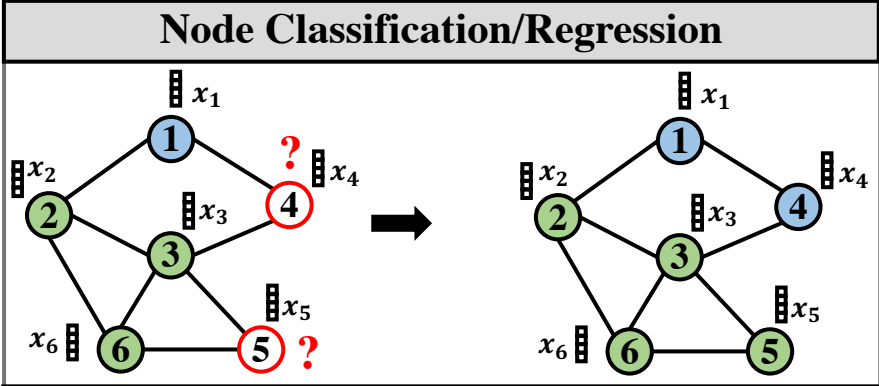
Knowledge Graph



Reasoning / Planning Graph



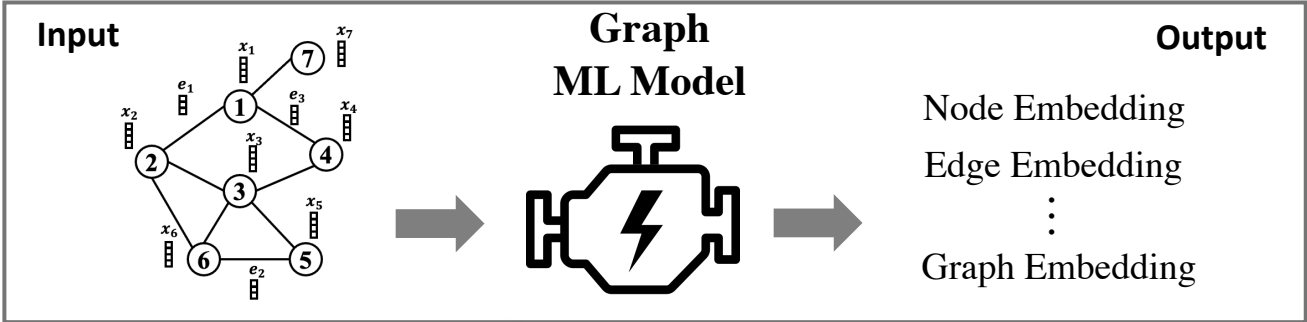
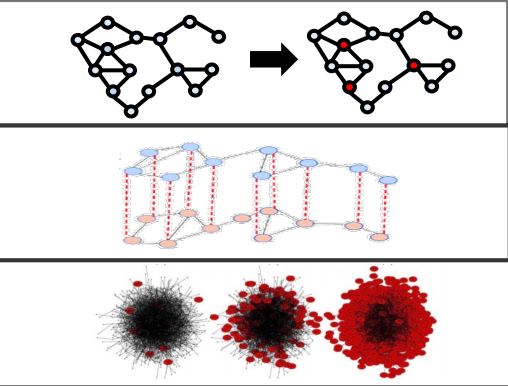
Introduction and Background - Graph-based Tasks and Graph Machine Learning



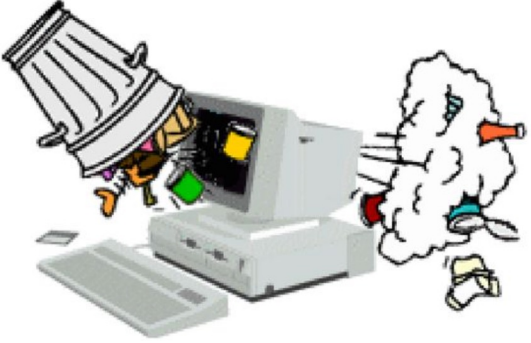
Network Dismantling

Network Alignment

Influence Maximization



Real-world graph data can have data quality challenges...



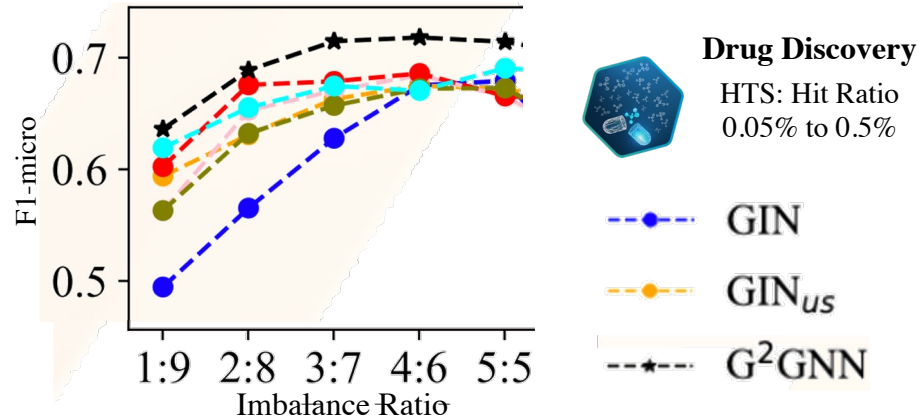
Garbage in, garbage out

Introduction and Background – Real-world Graphs have Data Quality Issues

Topological Issues e.g., Homophily vs Heterophily

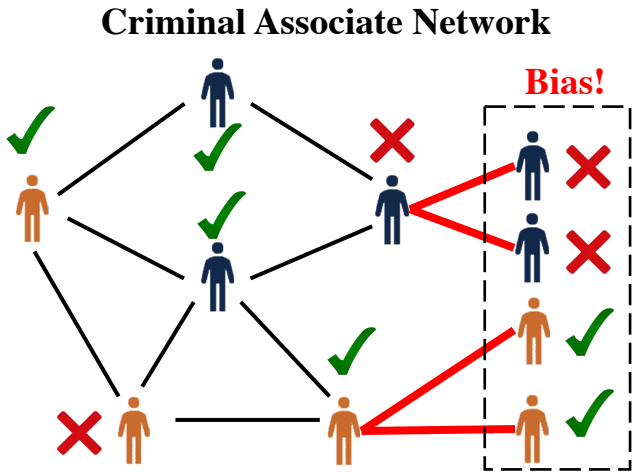


Imbalance Issues e.g., labeled data in chemistry

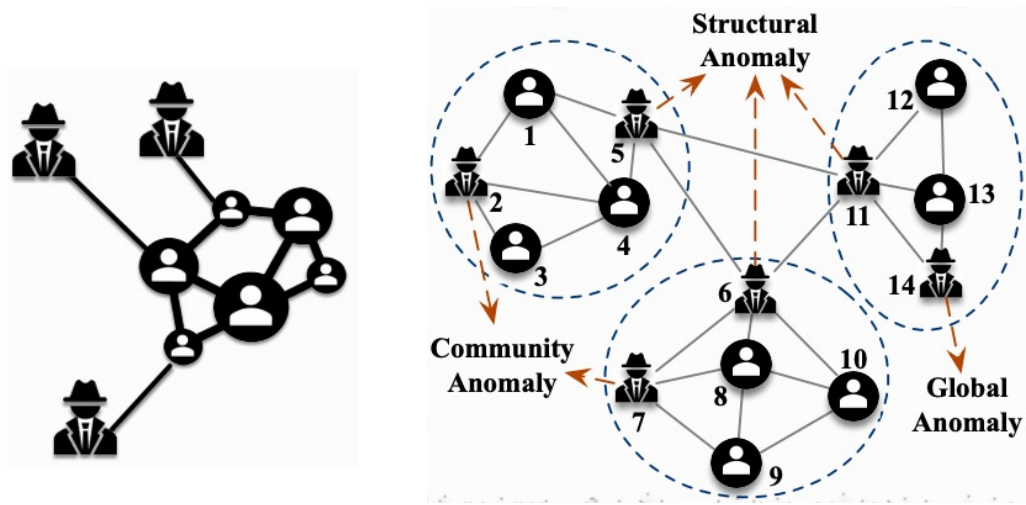


Bias Issues e.g., bail decision making

✓ Bail ✗ No-Bail 👤 Group 1 👤 Group 2

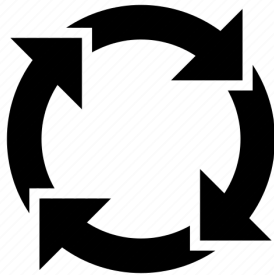
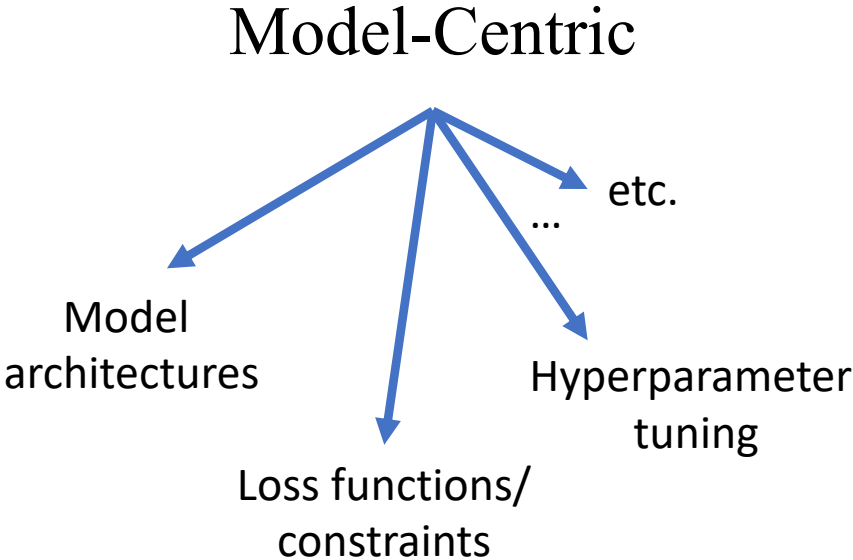


Abnormal Graph Data

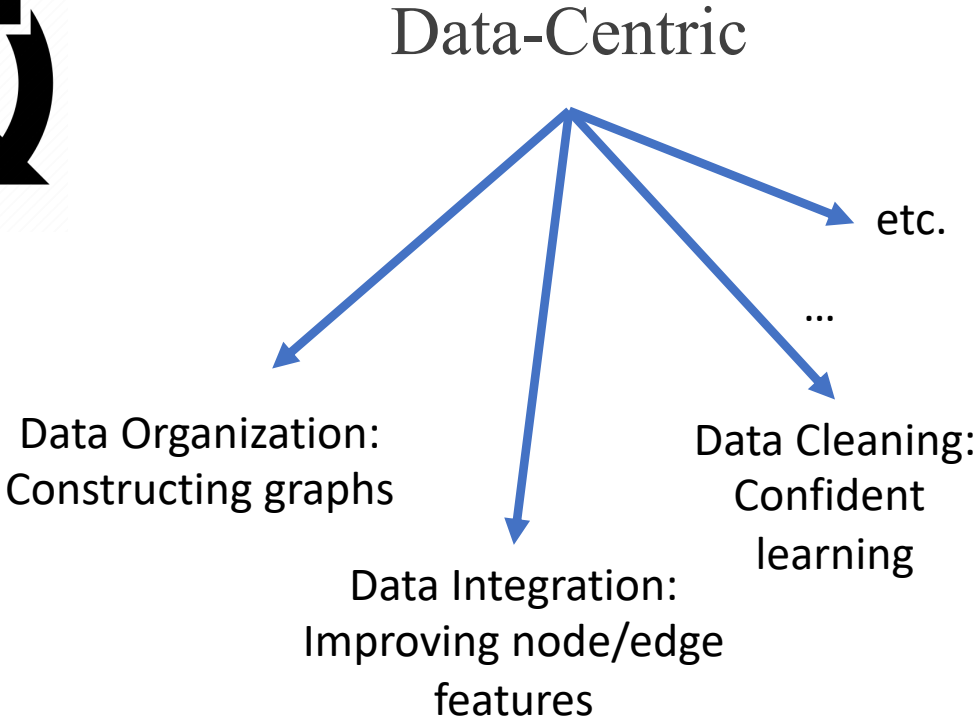


Introduction and Background – Model- vs. Data-Centric Methods

Find the *best model* for the given *fixed dataset*



Realize the *best dataset* for the given *prediction task*



Introduction and Background – Model- vs. Data-Centric Methods



Credit: MIT Introduction to Data-Centric AI course & Inspired by XKCD 2494 "Flawed Data"

Data Quality-Aware Graph Machine Learning

- Introduction and Background
- Topology Issues
- Imbalance Issues
- Short Break
- Bias and Fairness Issues
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- Summary

Outline

- **Introduction and Background**
- Topology Issues
- Imbalance Issues
- Short Break
- Bias and Fairness Issues
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- Summary

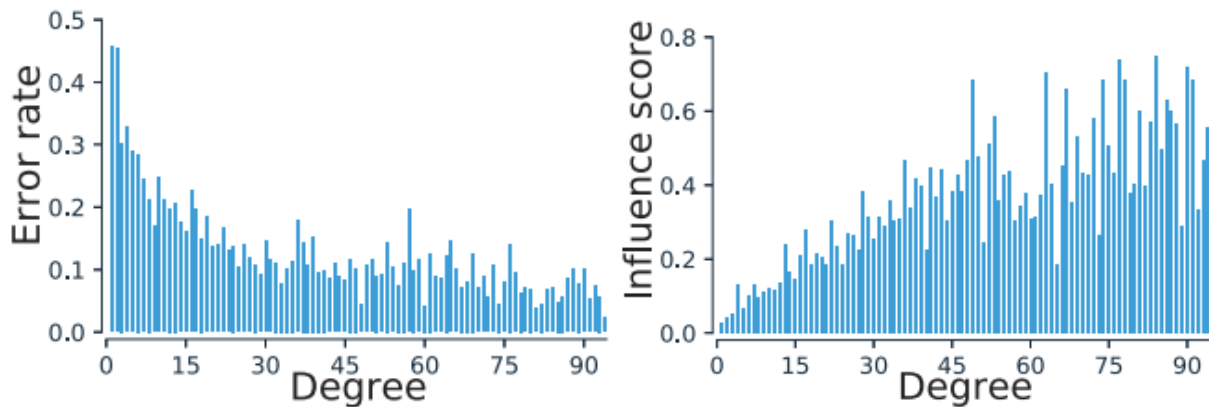
Outline

- Introduction and Background
- **Topology Issues**
- Imbalance Issues
- Short Break
- Bias and Fairness Issues
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- Summary

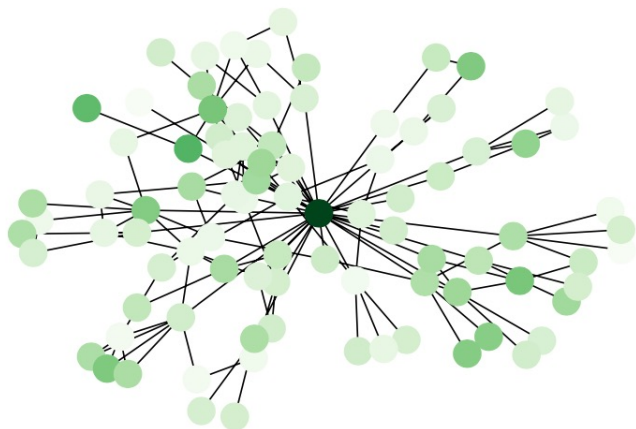
- Global Positional Issues
- Local Topology Issues
- Missing Graph Issues
- Future Directions and Q&A

Topology Issues – Global Topology Issues – Labeled Node Influence

Degree -> Influence



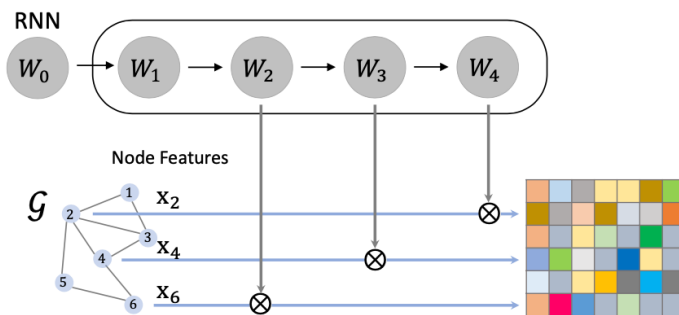
If $d_i > d_j$, v_i has higher influence than v_j on training GNNs



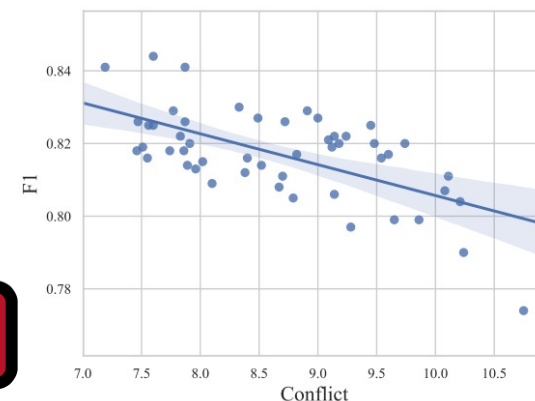
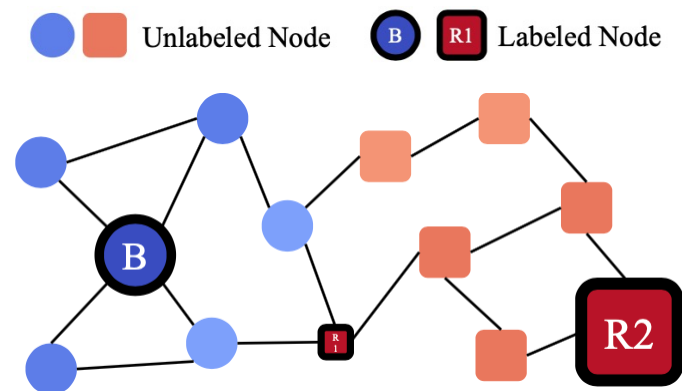
Darker colors - Higher influences.

$$\mathbf{x}_i^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} a_{ij} (\mathbf{W}^l + \mathbf{W}_{d_j}^l) \mathbf{x}_j^l \right)$$

Degree-dependent!



Position -> Influence



$$\mathbf{P} = \alpha (\mathbf{I} - (1 - \alpha) \mathbf{A}')^{-1}$$

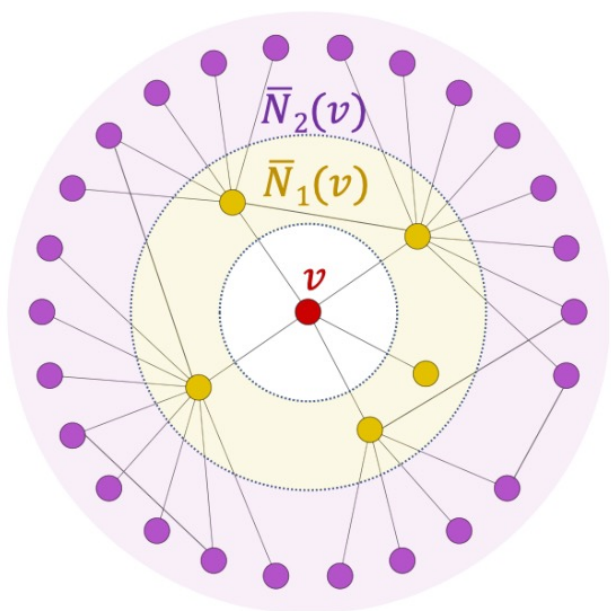
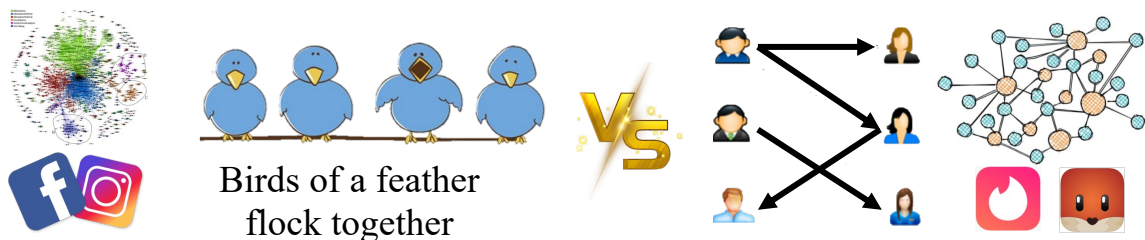
$$\mathbf{T}_v = \mathbb{E}_{\mathbf{x} \sim \mathbf{P}_v} \left(\sum_{j \in [1, k], j \neq y_v} |c_j|^{-1} \sum_{i \in \mathcal{C}_j} \mathbf{P}_{i, \mathbf{x}} \right)$$

$$L = -|\mathcal{L}|^{-1} \sum_{v \in \mathcal{L}} \mathbf{W}_v \sum_{c=1}^k y_v^c \log p_v^c$$

High T
High Conflicts,
low weight

Topology Issue – Local Topology Issues – Heterophily/Homophily

Homophily vs Heterophily



- Ego
- 1-order
- 2-order

Graph-level Homophily

$$h(\mathcal{G}, \{y_i; i \in \mathcal{V}\}) = \frac{1}{|\mathcal{E}|} \sum_{(j,k) \in \mathcal{E}} \mathbb{1}(y_j = y_k)$$

Ego-Neighbor Separation

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^{k-1}, \text{AGGR}(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v\}))$$

Higher-order Neighbor

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^{k-1}, \text{AGGR}_1(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v^1\}), \text{AGGR}_2(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v^2\} \dots))$$

Combination of Intermediate Representation

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^1, \mathbf{r}_v^2, \dots, \mathbf{r}_v^K)$$

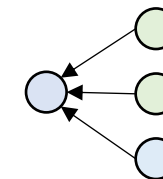
Class belief propagation

$$\mathbf{B}^k = \mathbf{B}^0 + \underbrace{\mathbf{A}\mathbf{B}^{k-1}}_{\text{Transition among Graph}} \mathbf{H}$$

Transition among Graph

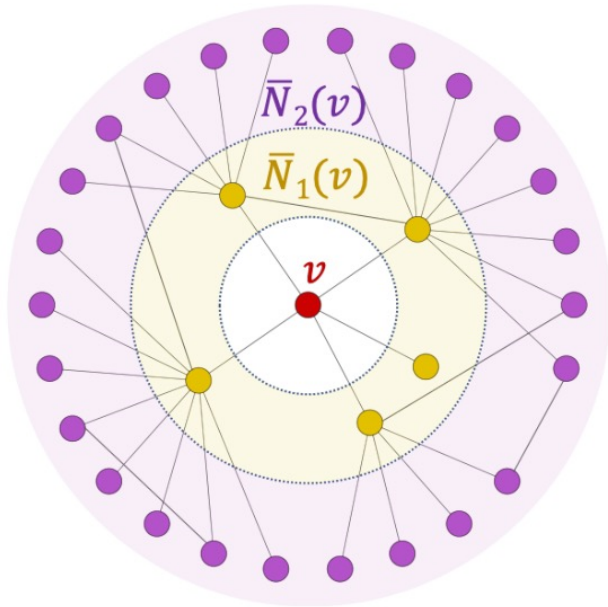
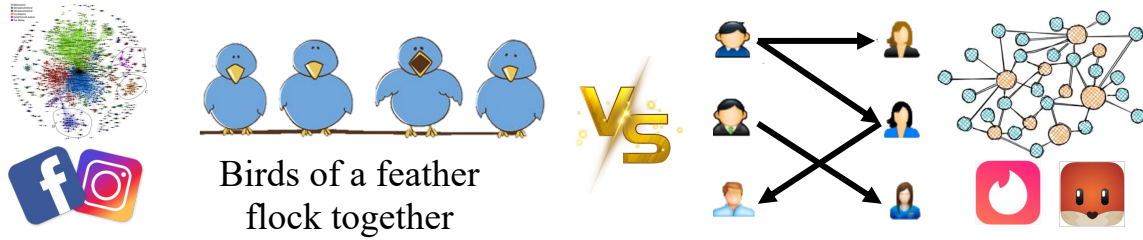
$$\mathbf{H} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$$

$$\mathbf{B} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{Y}|}$$



Topology Issue – Local Topology Issues – Heterophily/Homophily

Homophily vs Heterophily



- Ego
- 1-order
- 2-order

Graph-level Homophily

$$h(\mathcal{G}, \{y_i; i \in \mathcal{V}\}) = \frac{1}{|\mathcal{E}|} \sum_{(j,k) \in \mathcal{E}} \mathbb{1}(y_j = y_k)$$

Ego-Neighbor Separation

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^{k-1}, \text{AGGR}(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v\}))$$

Higher-order Neighbor

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^{k-1}, \text{AGGR}_1(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v^1\}), \text{AGGR}_2(\{\mathbf{r}_u^{k-1}; u \in \mathcal{N}_v^2\} \dots))$$

Combination of Intermediate Representation

$$\mathbf{r}_v^k = \text{COMBINE}(\mathbf{r}_v^1, \mathbf{r}_v^2, \dots, \mathbf{r}_v^K)$$

Class belief propagation

$$\mathbf{B}^k = \mathbf{B}^0 + \underbrace{\mathbf{A}\mathbf{B}^{k-1}}_{\text{Graph Transition}} \mathbf{H}$$

$$\mathbf{H} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{Y}|}$$

$$\mathbf{B} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{Y}|}$$

Graph Transition

Class Transition

$$\mathbf{A}\mathbf{B}^{k-1} \quad [0.1|0.1|0.8]$$

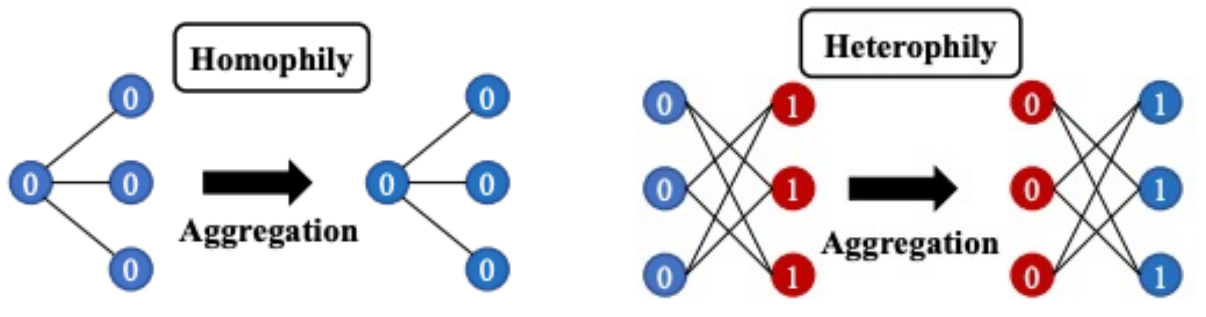


$$\mathbf{H} \quad \begin{bmatrix} 0.1 & 0.3 & 0.6 \\ 0.8 & 0.1 & 0.1 \\ 0.7 & 0.3 & 0 \end{bmatrix}$$

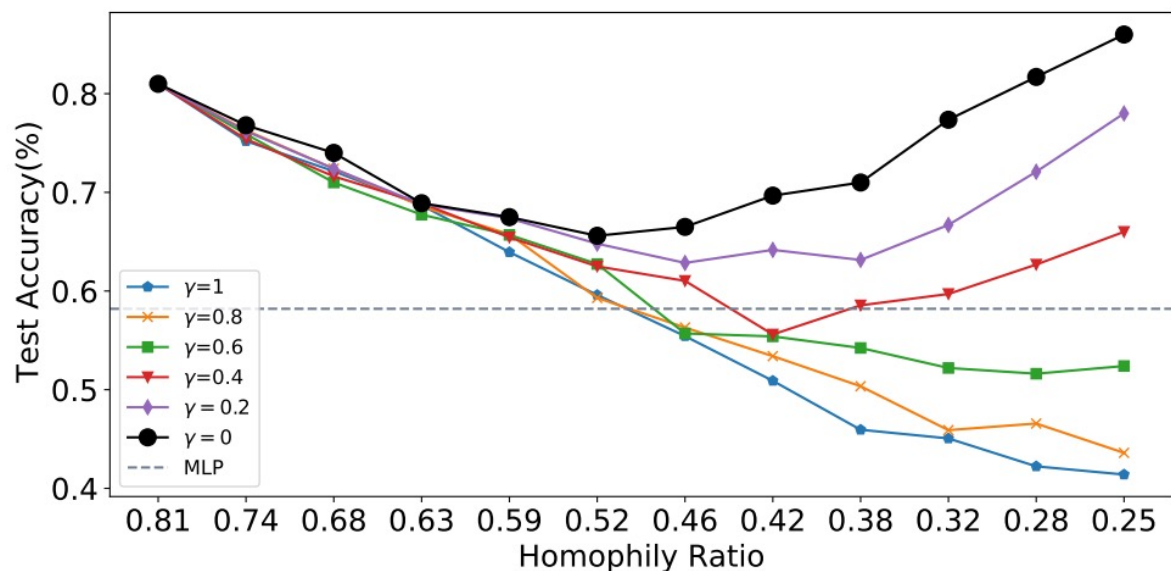
H

Topology Issue – Local Topology Issues – Heterophily/Homophily

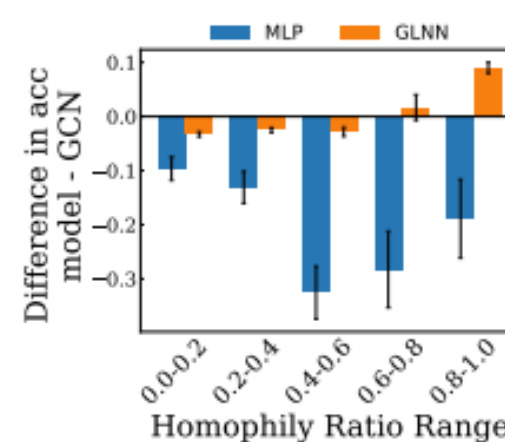
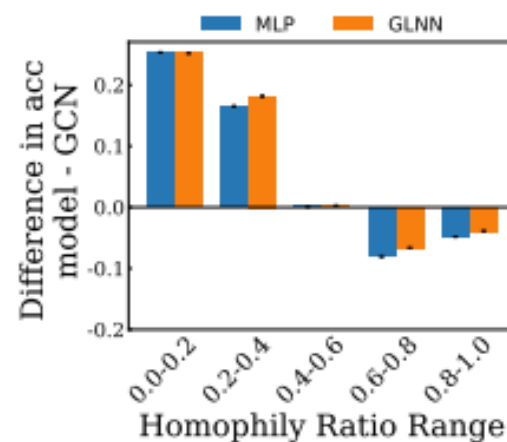
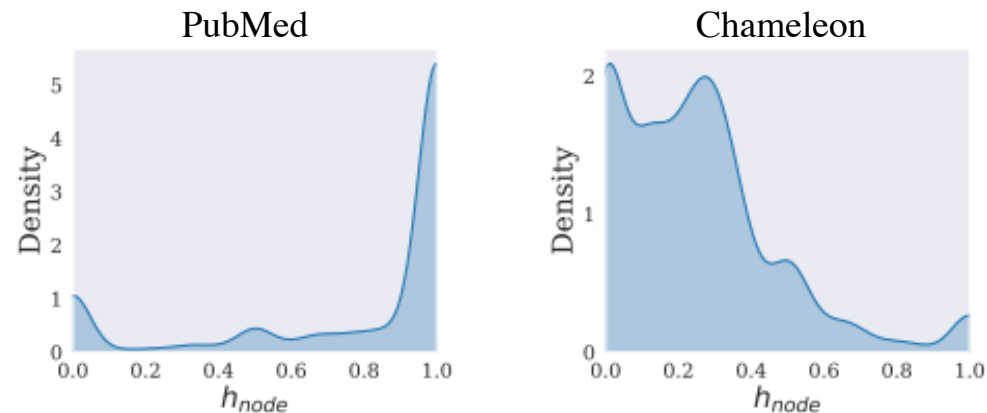
Across Different Graphs



Class Neighborhood Distribution Distinguishability



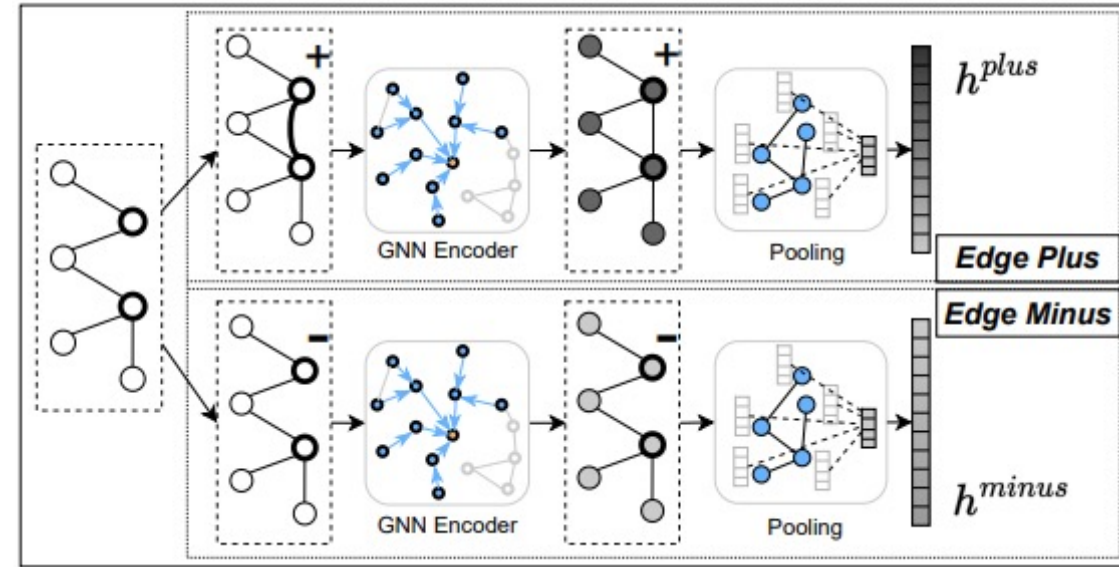
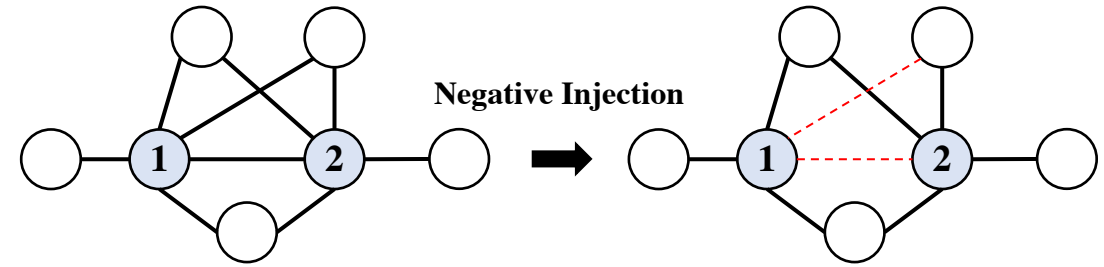
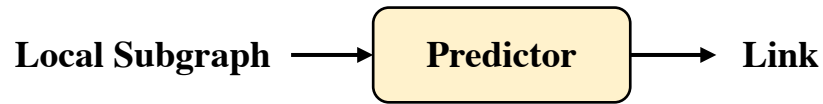
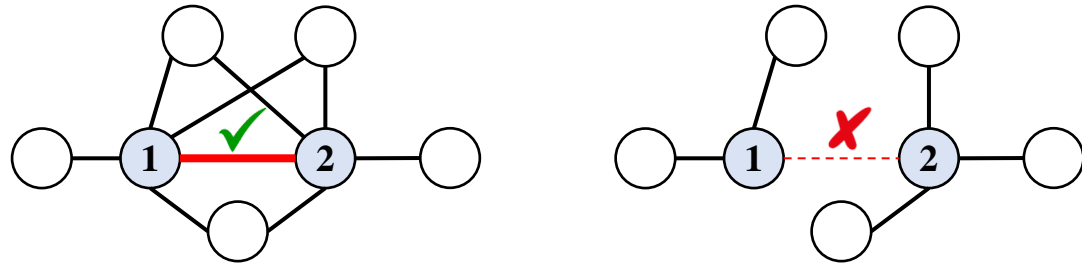
Within the Same Graph



In **homophily** graph, GNNs > MLP on **homophily** nodes
 In **heterophily** graph, GNNs > MLP on **heterophily** nodes

Topology Issue – Local Topology Issues – Training-to-Testing Topology Shift

Link Prediction



Edge Mean

Edge Attention

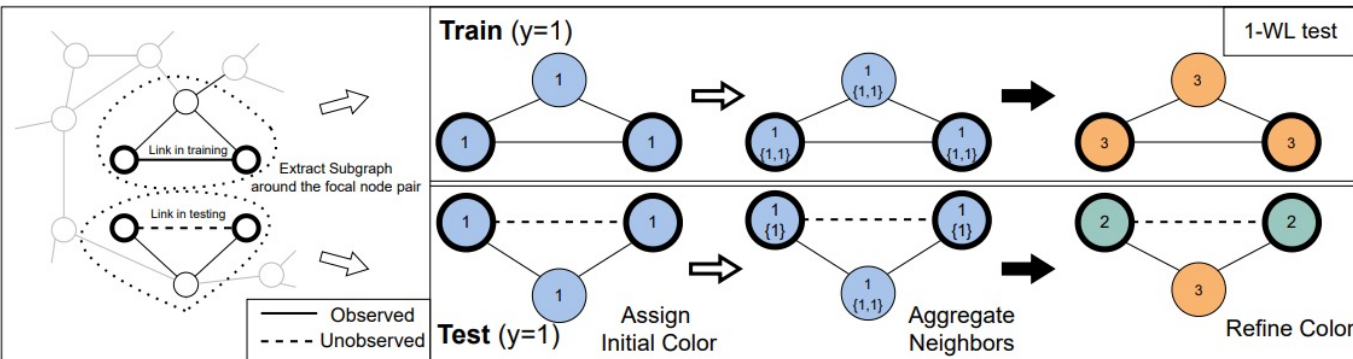
$$\mathbf{h}^{\text{mean}} = \frac{(\mathbf{h}^+ + \mathbf{h}^-)}{2}$$

$$\mathbf{h}^{\text{mean}} = w^+ \mathbf{h}^+ + w^- \mathbf{h}^-$$

$$w^+ = \sigma(\mathbf{q}^T \tanh(\mathbf{W}\mathbf{h}^+ + \mathbf{b}))$$

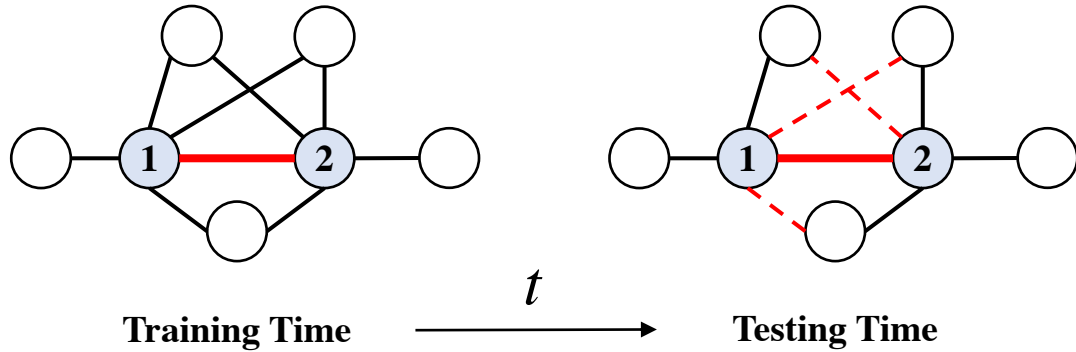
Focal Link is missing from training subgraph to testing subgraph

Distribution Shift

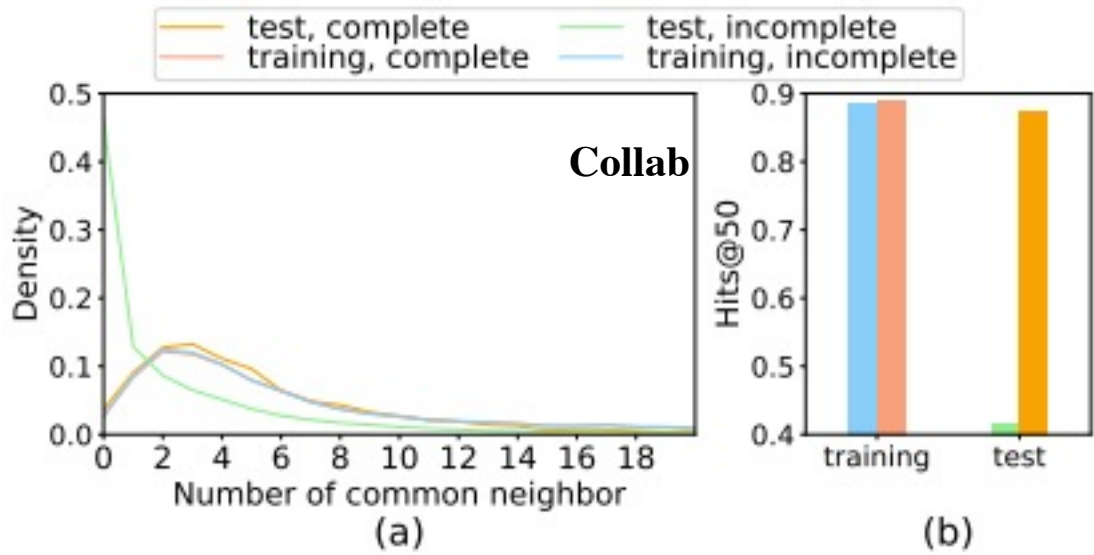
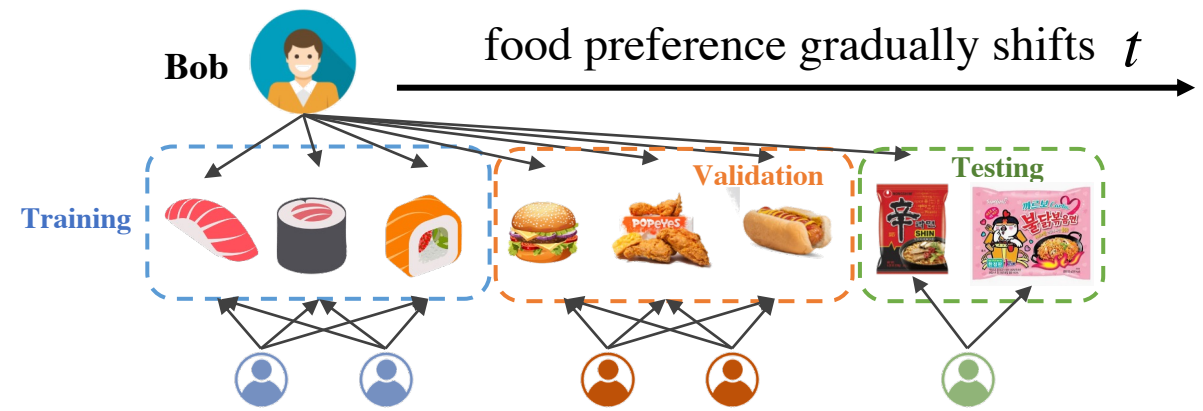


Topology Issue – Local Topology Issues – # of Common Neighbor Shift

Link-centric Perspective



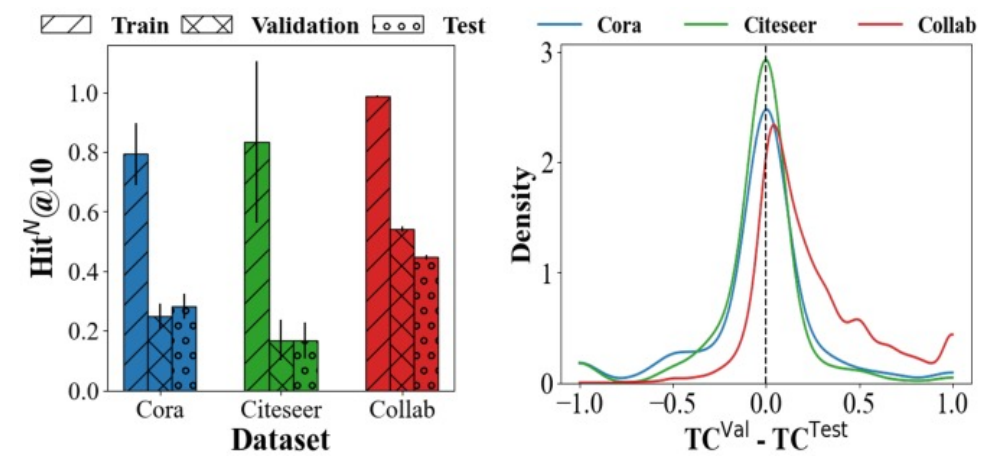
Node-centric Perspective



Time-based Split

Testing edges have more testing edges around

TC^{Val} : Common interaction between training and validation
 TC^{Test} : Common interaction between training and validation

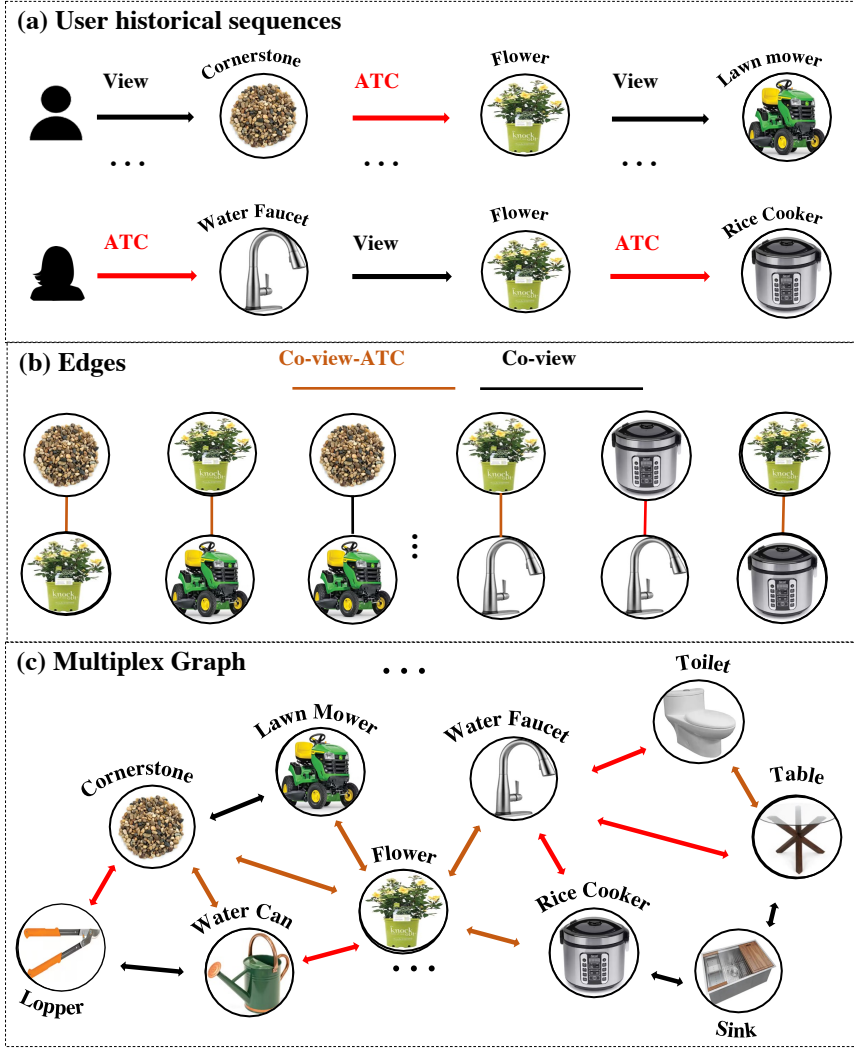


Topology Issue – Missing Topology Issues

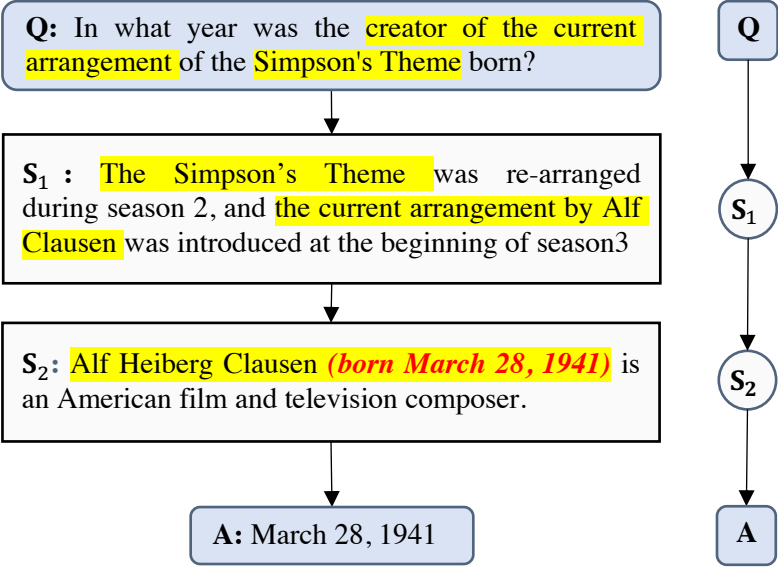
Sometimes Real-world Applications do not have Graphs!

But Graph can actually encode some useful information

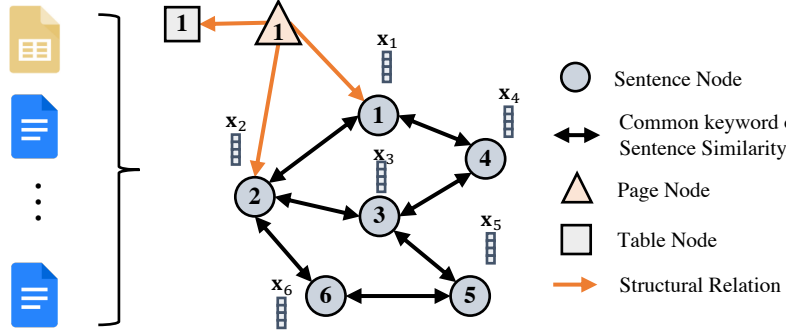
User/Item Interaction



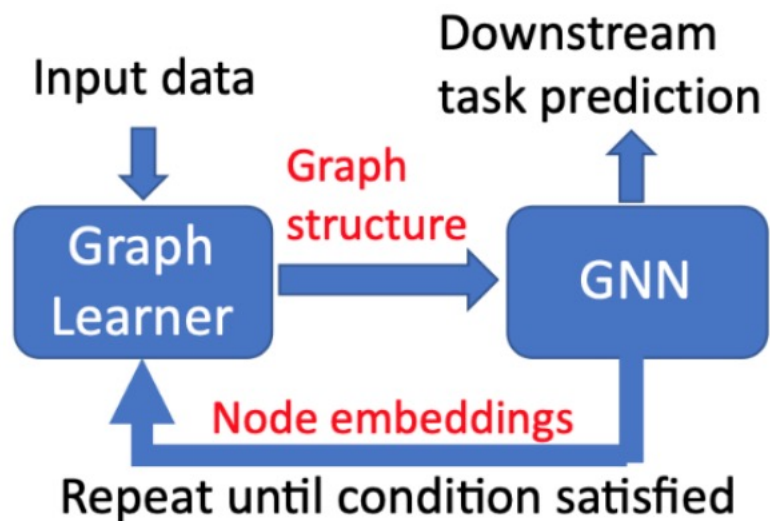
Multi-hop Reasoning



Document Graph Construction



Topology Issue – Missing Topology Issues



Better Graph Structure



Better Node Embeddings

GNN embeddings

$$A_{ij}^p = \cos(\mathbf{w}_p \odot \mathbf{v}_i, \mathbf{w}_p \odot \mathbf{v}_j), A_{ij} = m^{-1} \sum_{p=1}^m a_{ij}^p$$



Real-world Graph is sparse!

$$\mathbf{a}_{ij} = \begin{cases} A_{ij}, & \mathbf{a}_{ij} < \epsilon \\ 0, & \mathbf{a}_{ij} > \epsilon \end{cases} \quad \mathbf{A}^t = \lambda \mathbf{L}^0 + (1 - \lambda)(\eta f(\mathbf{A}^t) + (1 - \eta)f(\mathbf{A}^1))$$



Quadratic Computation $\mathcal{O}(n^2)$

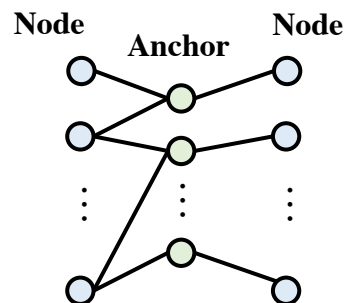
$$\mathbf{L}^0 = (\mathbf{D}^0)^{-0.5} \mathbf{A}^0 (\mathbf{D}^0)^{-0.5}$$

$$f(\mathbf{A})_{ij} = \mathbf{A}_{ij} / \sum_{j=1}^m \mathbf{A}_{ij}$$



$$\mathbf{r}_{ik}^p = \cos(\mathbf{w}_p \odot \mathbf{v}_i, \mathbf{w}_p \odot \mathbf{u}_k), \mathbf{r}_{ik} = m^{-1} \sum_{p=1}^m \mathbf{r}_{ik}^p$$

Anchor Selection $\mathcal{O}(nK), K \ll n$



$$\mathbf{F}^{0'} = \mathbf{\Lambda}^{-1} \mathbf{R}^T \mathbf{F}^0$$

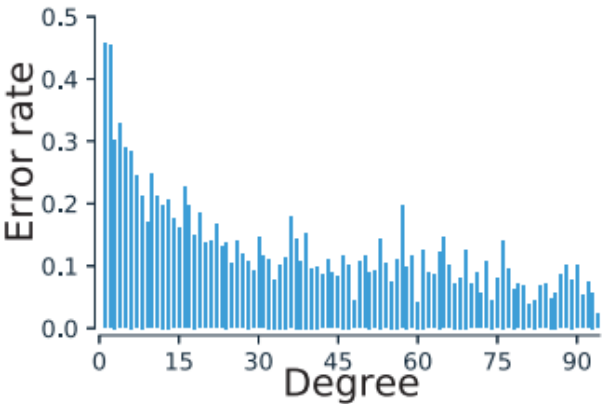
Node -> Anchor

$$\mathbf{F}^1 = \mathbf{\Lambda}^{-1} \mathbf{R}^T \mathbf{F}^{0'}$$

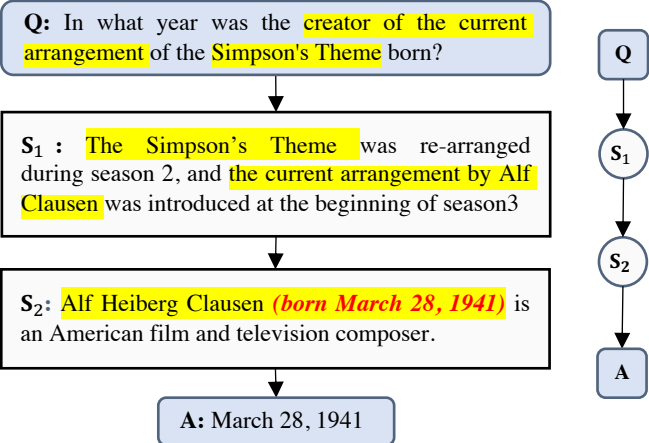
Anchor -> Node

Q&A and Future Work – Topology Issue

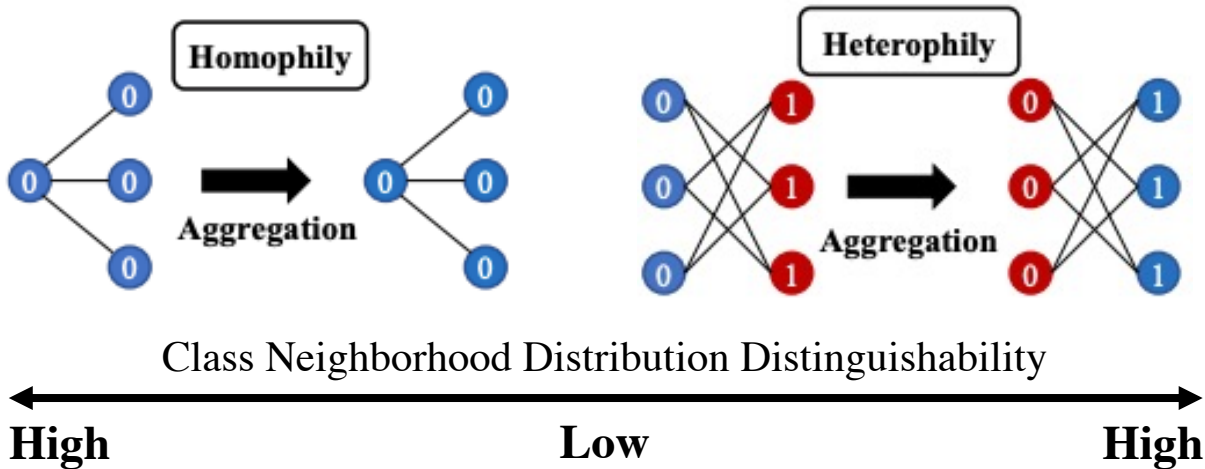
Global Topology Issue



Missing Topology Issue

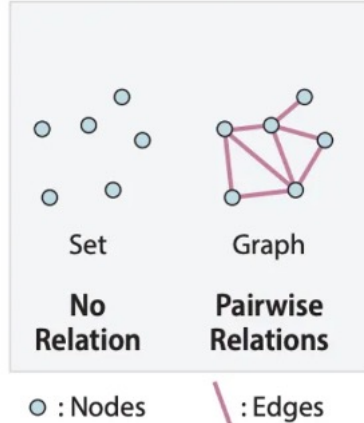


Local Topology Issue

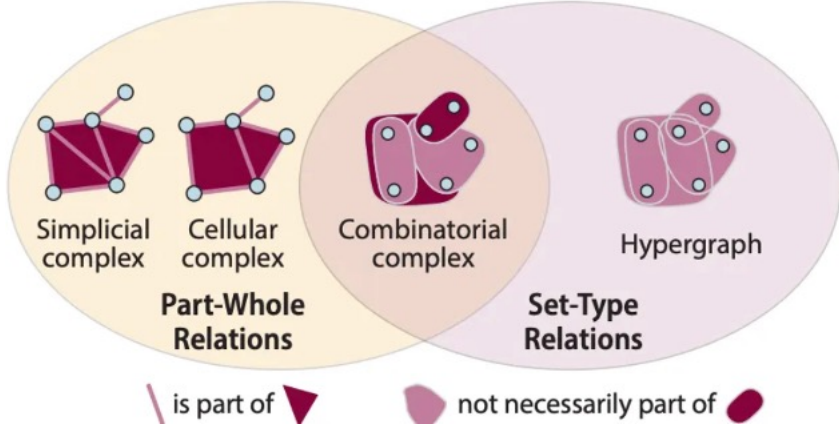


Topology Issue of Complex Graphs

Traditional Discrete Domains



Domains of Topological Deep Learning



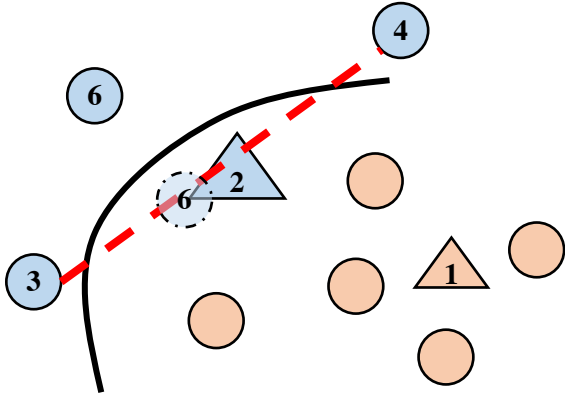
Outline

- Introduction and Background
- Topology Issues
- **Imbalance Issues**
- Short Break
- Bias and Fairness Issues
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- Summary

- Node-level Imbalance
- Graph-level Imbalance
- Edge-level Imbalance
- Future Directions and Q&A

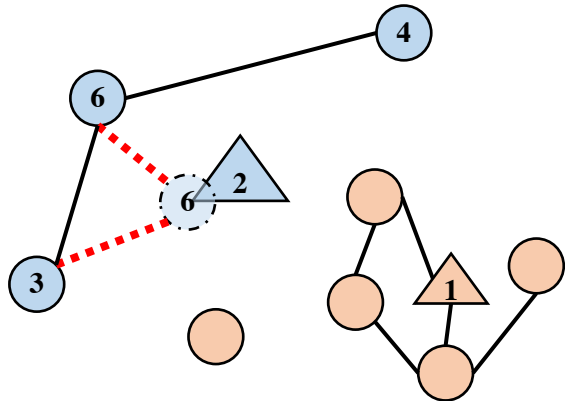
Imbalance Issues – Node-level imbalance

SMOTE



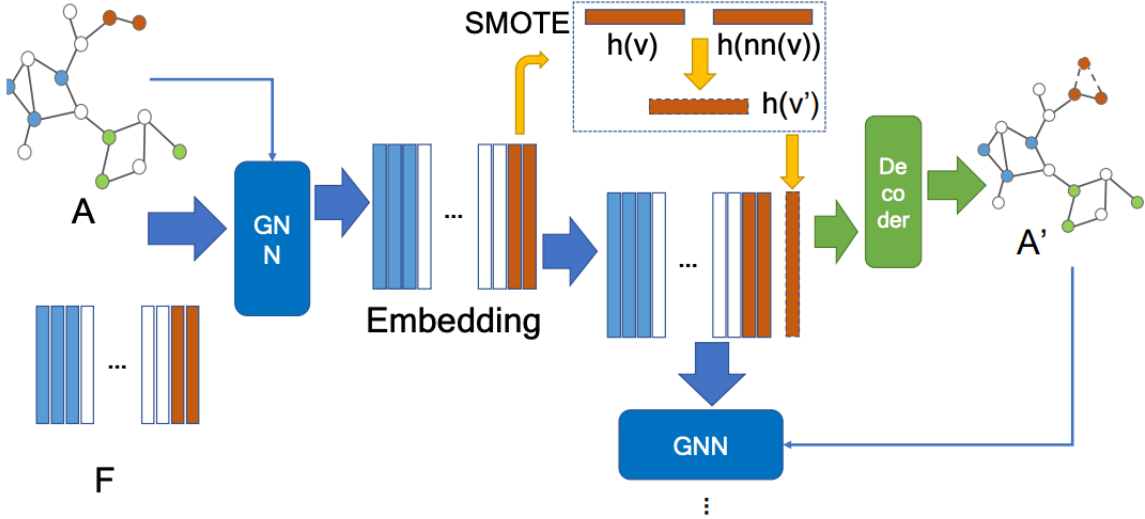
- Feature Interpolation
- Train – Major
- Train – Minor
- ▲ Test – Major
- ▲ Test – Minor

Graph-structured data has both feature and edge



- Feature Interpolation
- ⋯ Edge Generation

GraphSMOTE



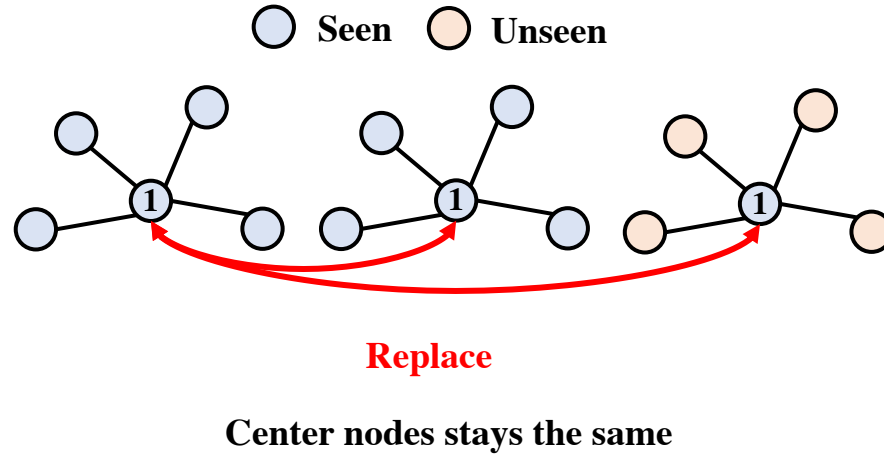
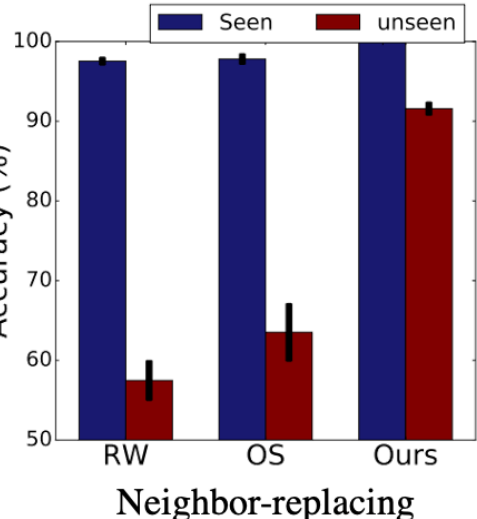
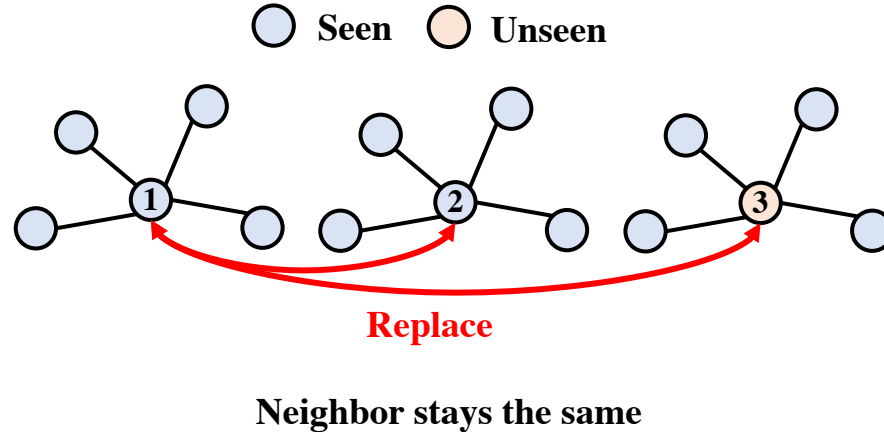
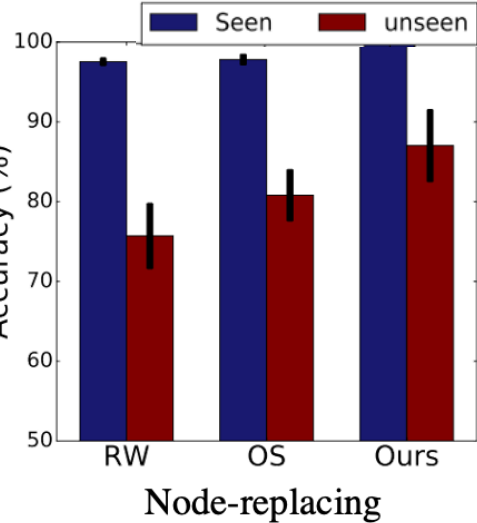
$$nn(v) = \operatorname{argmin}_u ||\mathbf{h}_u^1 - \mathbf{h}_v^1||, \text{ s. t. } \mathbf{Y}_u = \mathbf{Y}_v$$

$$\mathbf{h}_{v'}^1 = (1 - \delta)\mathbf{h}_v^1 + \delta\mathbf{h}_{nn(v)}^1$$

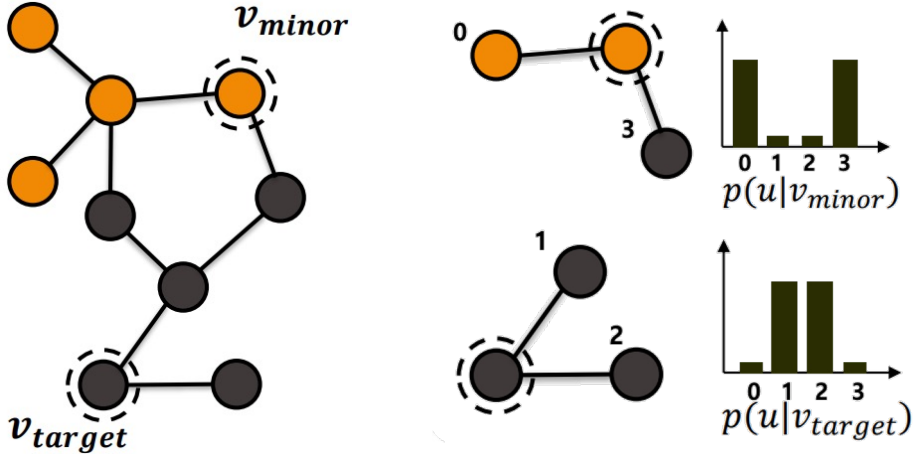
$$\mathbf{A}_{v'u} = \begin{cases} 1, & \text{if } \mathbf{E}_{v'u} \geq \eta \\ 0, & \text{otherwise} \end{cases} \quad \mathcal{L}_{edge} = ||\mathbf{E} - \mathbf{A}'||_F^2$$

$$\mathbf{E}_{vu} = \operatorname{softmax}(\sigma(\mathbf{h}_v^1 \mathbf{S} \mathbf{h}_u^1))$$

Imbalance Issues – Node-level imbalance



Neighborhood Memorization



$$p(u|v_{mixed}) = \hat{\phi}p(u|v_{minor}) + (1 - \hat{\phi})p(u|v_{target})$$

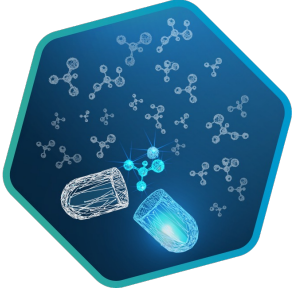
$$0.5 < \hat{\phi} = \frac{1}{1 + e^{-\phi}} < 1 \quad \phi = KL(\sigma(\mathbf{o}_{minor}) || \sigma(\mathbf{o}_{target}))$$

$$\mathbf{o}_{minor} = |\mathcal{N}_v|^{-1} \sum_{u \in \mathcal{N}_v} \mathbf{o}_{minor}$$



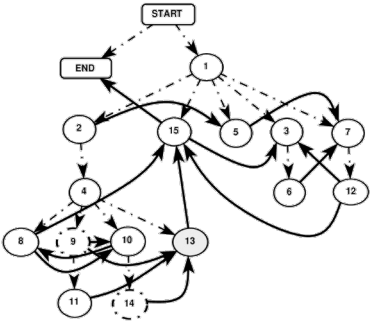
Imbalance Issues – Graph-level imbalance

Drug Discovery



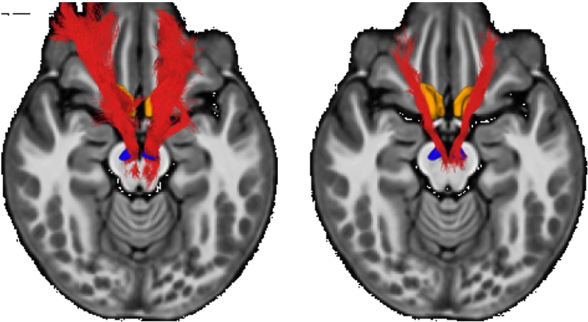
HTS Hit Ratio
0.05% to 0.5%

Malware Detection



0.01% Google, 2% Android,

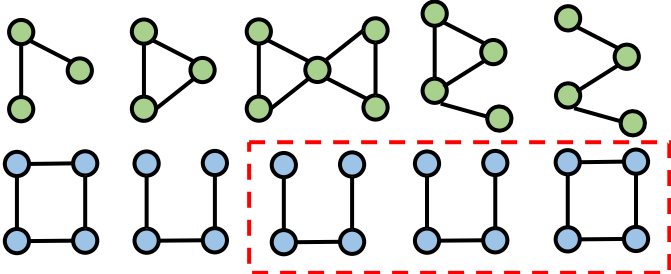
ASD Brain Classification



Normal : Autism
36 : 1

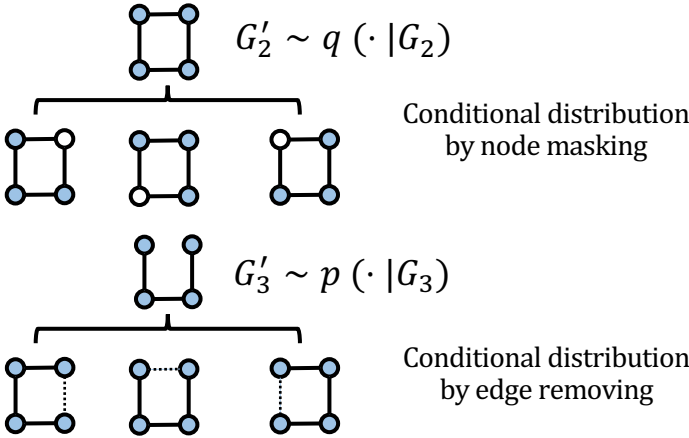
Autism Statistics. 2023

Quantity Augmentation

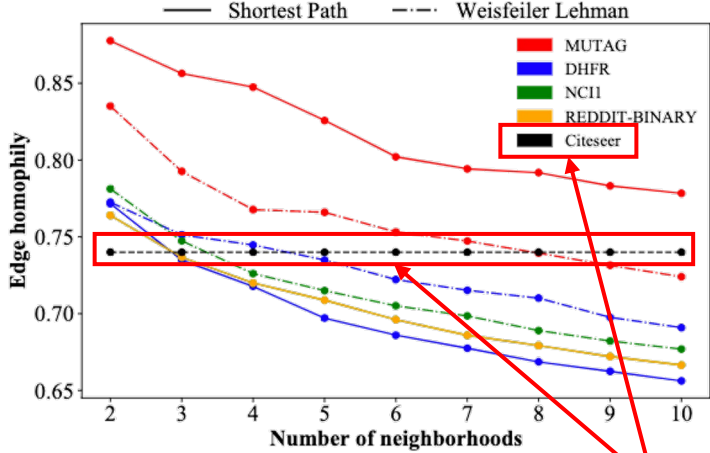
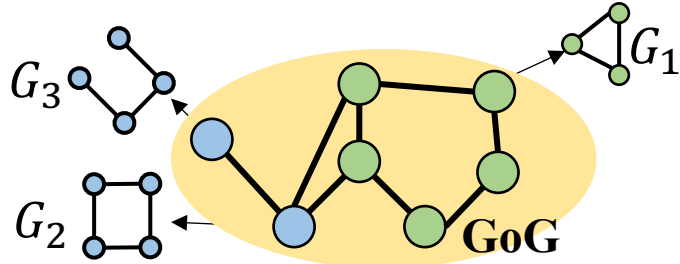


Structure Augmentation

SPP - Structurally Similar Molecules tend to have similar properties

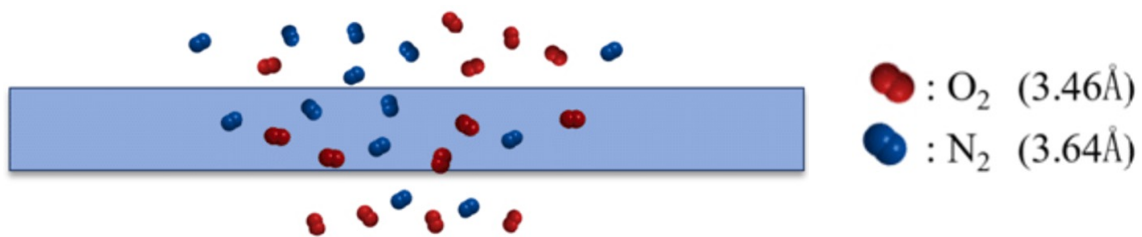


Graph-of-Graphs (GoG)



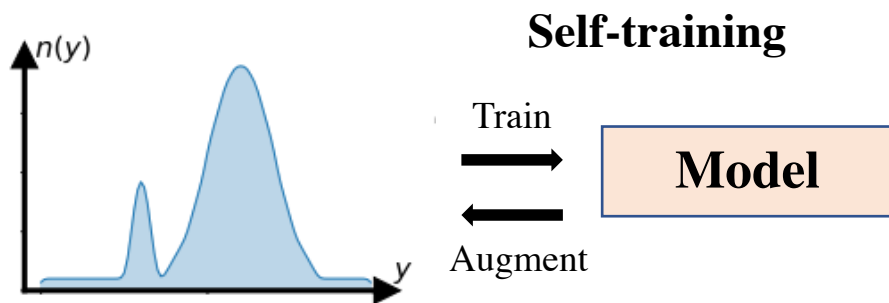
Constructed GoG demonstrates high homophily!

Imbalance Issues – Graph-level imbalance



70 years, ~600 polymers, oxygen permeability ,
 Polymer Gas Separation Membrane Database

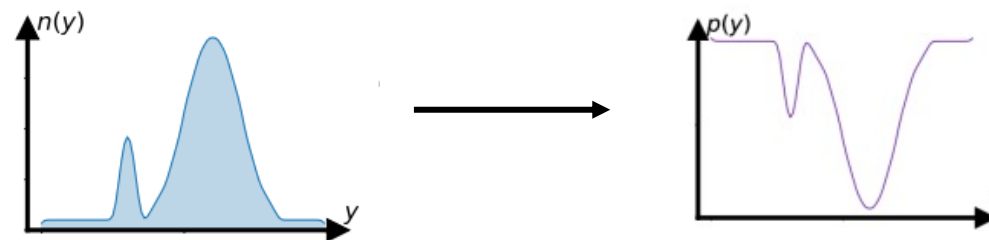
Imbalance Graph
 Regression!



(1) Use Model to predict on unlabeled graphs and
 select those high-quality-one

$$\sigma_i = \frac{1}{\text{Var} \left(\{f(g(G_{(i,j)}))\}_{j=1,2,\dots,B} \right)}$$

(2) Sample more for label interval with less training samples



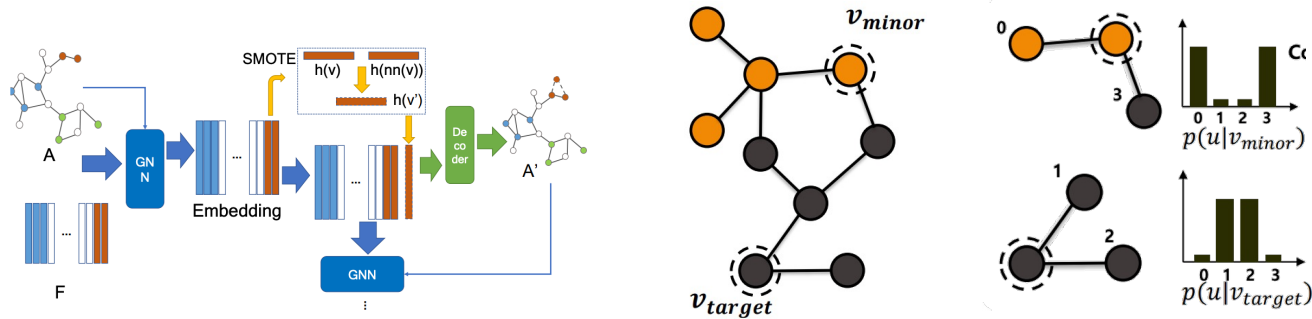
(3) Anchor-based Mix-up

a_i, \mathbf{z}_i : anchor-label
 and embedding

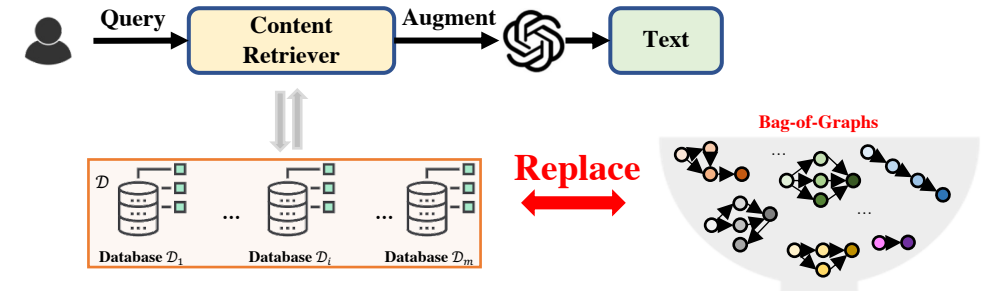
$$\begin{cases} \tilde{\mathbf{h}}_{(i,j)} &= \lambda \cdot \mathbf{z}_i + (1 - \lambda) \cdot \mathbf{h}_j, \\ \tilde{y}_{(i,j)} &= \lambda \cdot a_i + (1 - \lambda) \cdot y_j, \end{cases}$$

Q&A and Future Work – Imbalance Issues

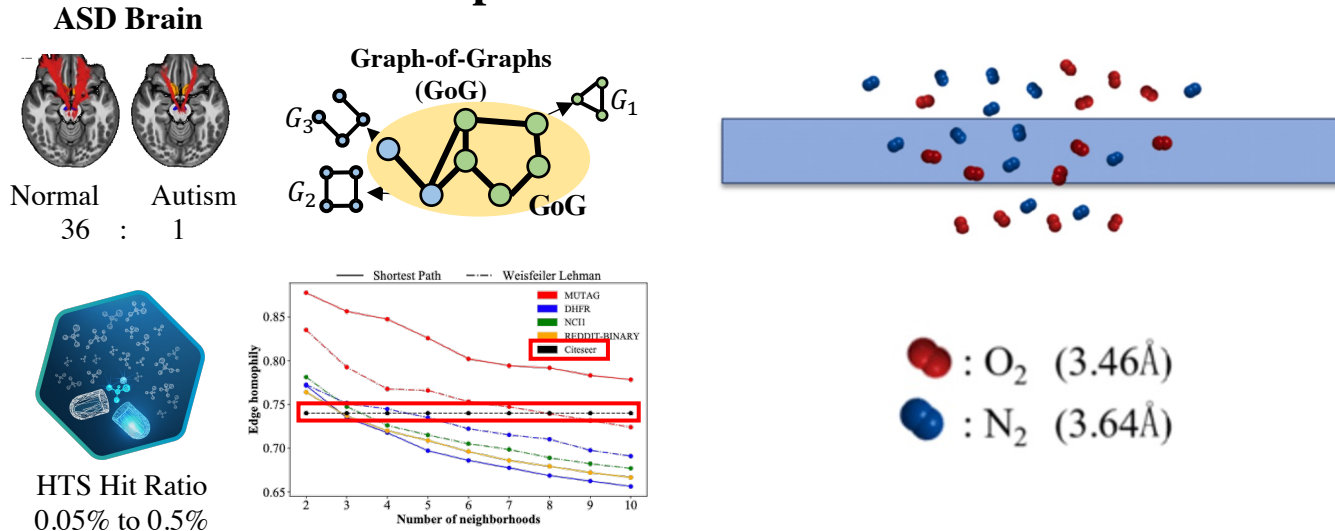
Node-level Imbalance



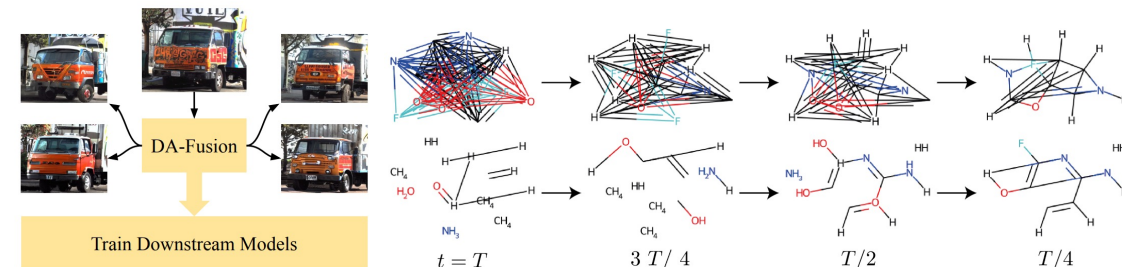
Retrieval Additional Supervision



Graph-level Imbalance



Generate Additional Supervision



Short Break (4 min)

Outline

- Introduction and Background
- Topology Issues
- Imbalance Issues
- Short Break
- **Bias and Fairness Issues**
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- Summary

Bias and Fairness Issues - Suicide Prevention

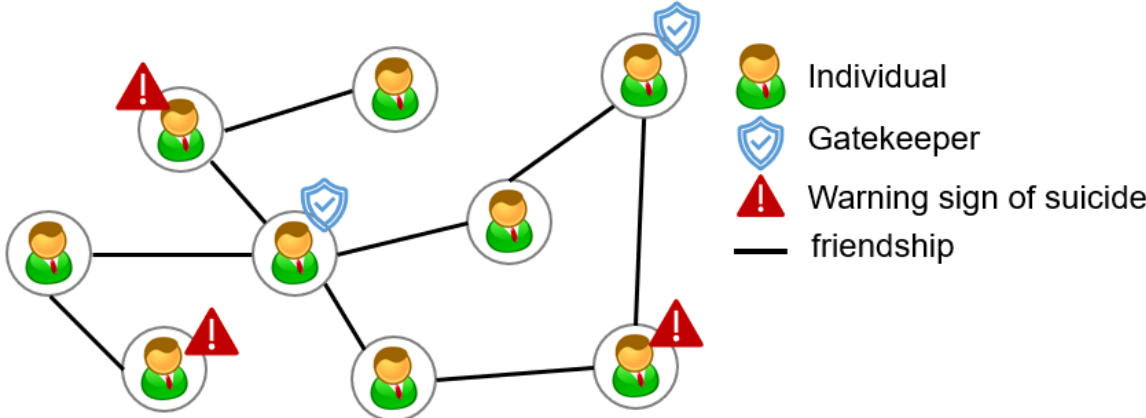
- **Why suicide prevention?**

- Suicide is one of the leading causes of death in United States

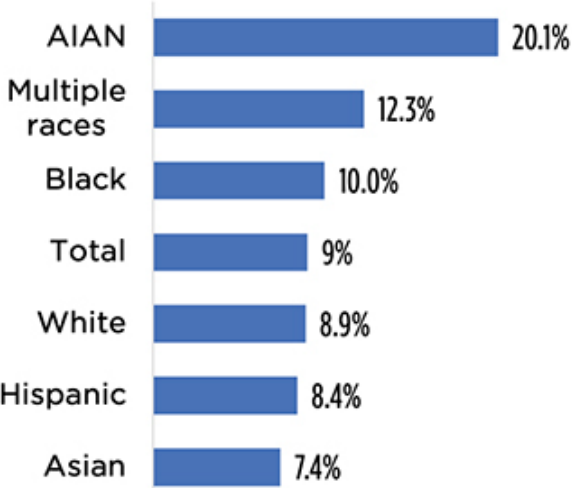
Percentage of high schoolers reporting a suicide attempt in the past 12 months, by race/ethnicity



Gatekeeper training programs



Toy example of a gatekeeper training program



Suicide attempts by race/ethnicity

- Existing prevention strategies **disproportionately** affect different groups

- **Key question**

- How to correct the bias and ensure fairness on graphs?

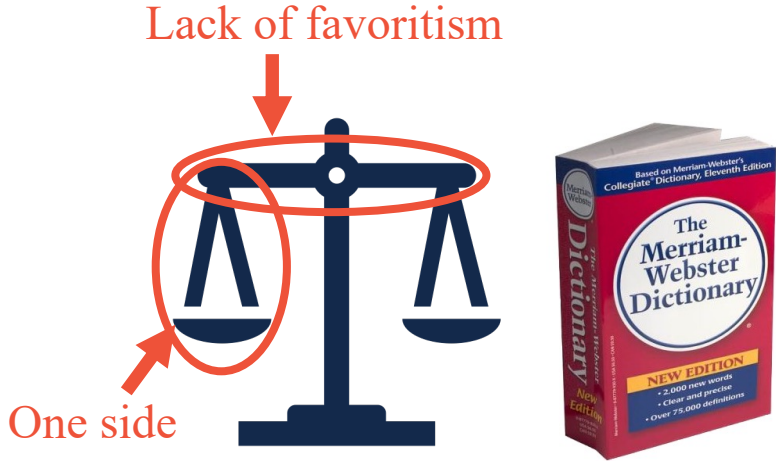
Bias and Fairness Issues - Fairness Definition

- **Principle**

- Lack of favoritism from one side or another

- **Rich fairness definitions**

- **Group fairness**
 - Statistical parity
 - Equal opportunity
 - Equalized odds
 - Accuracy parity
 - ...
- **Individual fairness**
- Counterfactual fairness
- **Degree fairness (on graphs)**



Fairness definition	Two sides
Group fairness	Two demographic groups
Individual fairness	Two data points
Counterfactual fairness	A data point and its counterfactual version
Degree fairness	Two group of nodes with same degree

- **Group Fairness on Graphs**
- Individual Fairness on Graphs
- Degree Fairness on Graphs
- Future Directions and Q&A

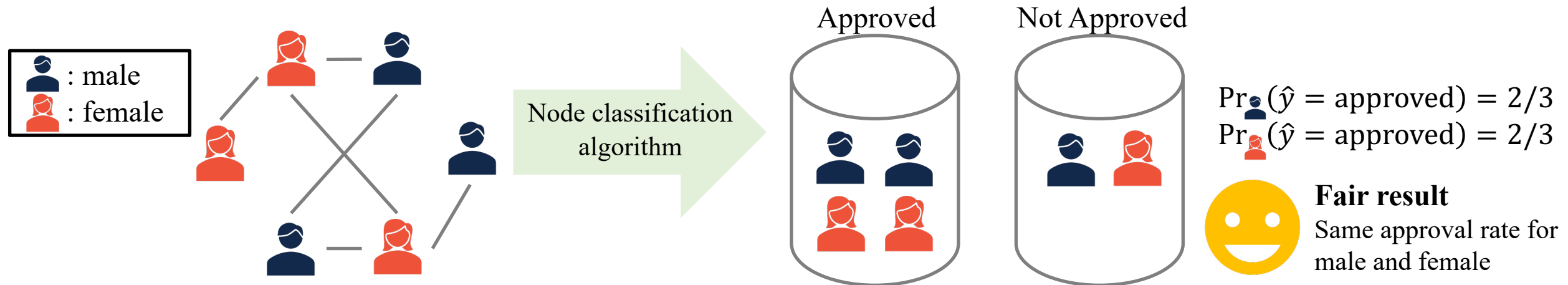
Group Fairness: Statistical Parity

- **Statistical parity = equal acceptance rate**

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

- \hat{y} : model prediction
- \Pr_+ : probability for the protected group
- \Pr_- : probability for the unprotected group
- Also known as demographic parity, disparate impact

- **Example: clinical trial participation**



Group Fairness: Equal Opportunity

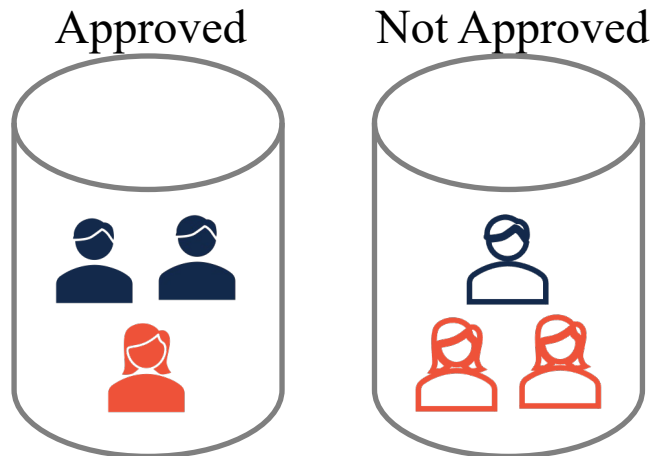
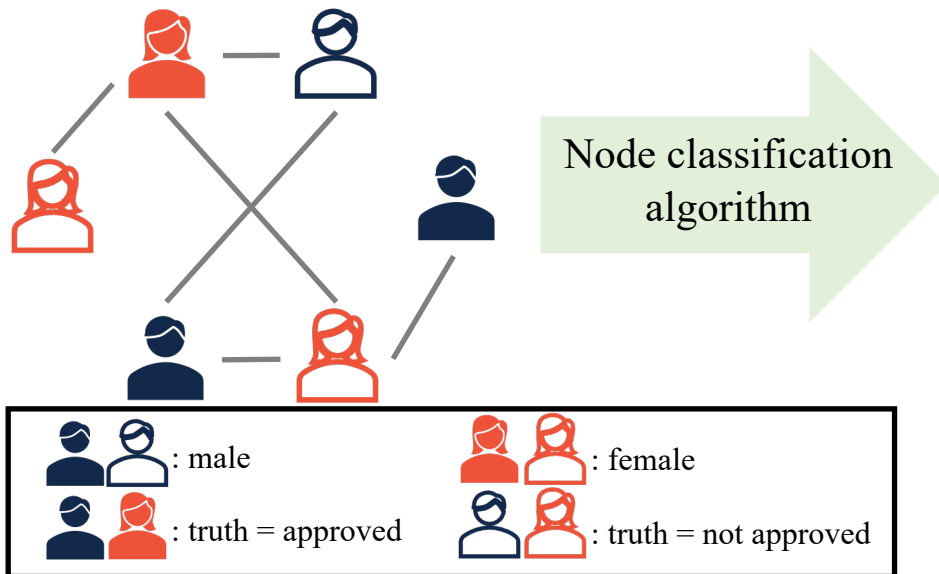
- **Equal opportunity = equal true positive rate**

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

- y : true label
- \hat{y} : model prediction
- \Pr_+ : probability for the protected group
- \Pr_- : probability for the unprotected group

If hold for **all** classes, it is called **equalized odds**

- **Example: clinical trial participation**



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

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



Fair result

Same true positive rate for male and female

Adversarial Learning for Fair Representation Learning

- **Statistical parity**

- Independence between the learned embedding \mathbf{z} and a sensitive attribute a

$$\mathbf{z}_u \perp a_u, \forall \text{ node } u$$

where a_u is the sensitive value of node u

- **Formulation**

- Mutual information minimization

$$I(\mathbf{z}_u, a_u) = 0, \forall \text{ node } u$$

- Analogous to statistical parity in classification task
- Fail to predict a_u using \mathbf{z}_u ← no information about a_u in \mathbf{z}_u

- **Solution**

- Adversarial learning
- Encoder: encode node into low-dimensional embedding space for downstream tasks
- Discriminator: fail to predict a_u using \mathbf{z}_u

Corresponding to
'adversarial'

Limitation #1: Full Access to Sensitive Attribute Information

- **Adversarial learning**

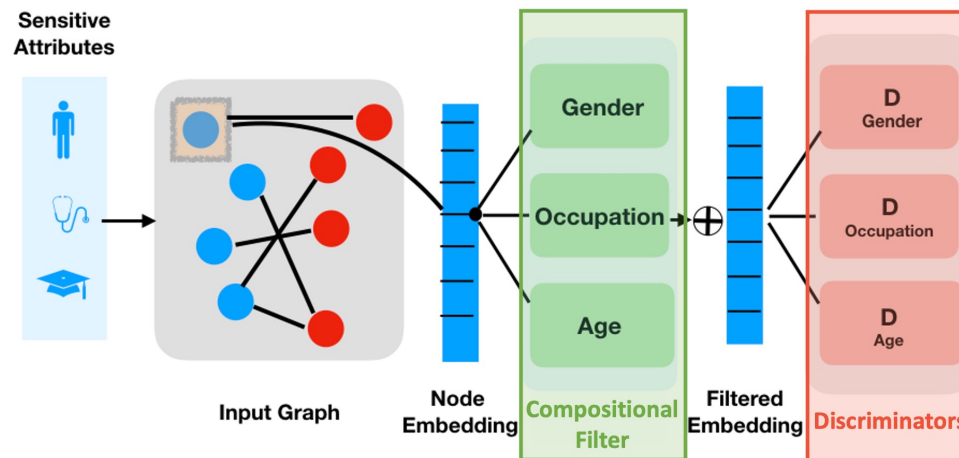
- Minimize a task-specific loss function to learn 'good' representations
- Maximize the error of predicting sensitive feature to learn 'fair' representations

- **Limitations**

- Require the sensitive attribute of all training nodes to train a good discriminator
- Ignore the fact that sensitive information is hard to obtain due to privacy

- **Question**

- What if we only have **limited** sensitive attribute information?



FairGNN: Additional Supervision Signal

- **Observation**

- Adversarial learning is unstable to train \leftarrow even worse with limited sensitive attribute
- Failure to converge may also cause discrimination

- **Key idea**

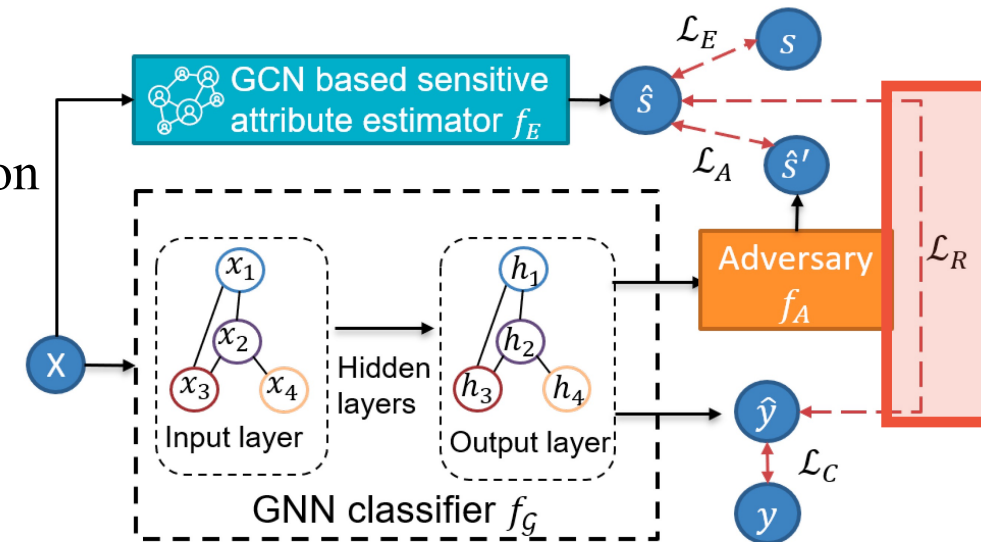
- Additional prerequisite of independence for additional supervision
- Independence \rightarrow zero covariance

- **Solution**

- Pseudo sensitive attribute from a sensitive attribute estimator
 - Not embedding from encoder
 - Offer pseudo-label for covariance minimization
- Absolute covariance minimizer to minimize absolute covariance between model prediction \hat{y} and pseudo sensitive attribute \hat{s}

$$\mathcal{L}_R = |\text{cov}(\hat{s}, \hat{y})| = |\mathbb{E}[(\hat{s} - \mathbb{E}[\hat{s}])(\hat{y} - \mathbb{E}[\hat{y}])]|$$

- Absolute covariance to avoid minimizing negative covariance



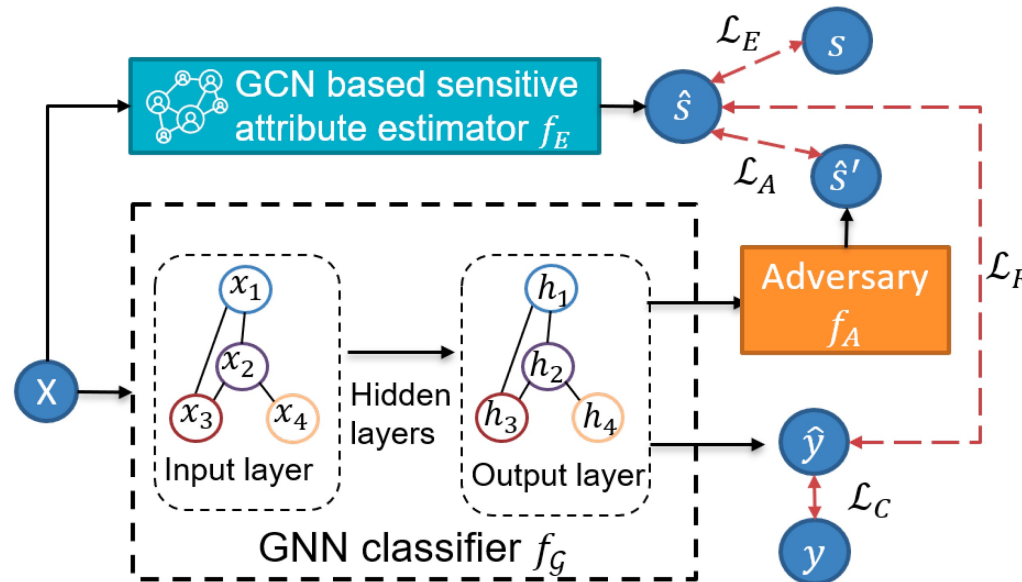
FairGNN: Overall Framework

- Overall loss function

$$\mathcal{L} = \mathcal{L}_C + \mathcal{L}_E - \alpha \mathcal{L}_A + \beta \mathcal{L}_R$$

- Intuition

- \mathcal{L}_C : classification loss (e.g., cross entropy) for learning representative node representation
- \mathcal{L}_E : sensitive attribute estimation loss for generating accurate pseudo sensitive attribute information
- \mathcal{L}_A : adversarial loss for debiasing the learned node representation
- \mathcal{L}_R : covariance minimizer to stabilize the adversary training



BeMap: Fair Topology View Generation

- Motivation**
 - Message passing could be unfair

Classification Task for Job Application

Notations

✓ : qualified

X : unqualified

Unfair Prediction

$$P(\checkmark | \text{black}) = \frac{2}{3}$$

$$P(\checkmark | \text{red}) = \frac{1}{3}$$

- Empirical evidence**
 - Predict node sensitive attribute using embeddings learned from GCN and MLP (no MP)

A Toy Graph

Demographic groups Majority: Minority:

How to calculate Majority Neighbor Ratio?

$$\text{Majority Neighbor Ratio}(\text{blue}) = \frac{(\# \text{ of } \text{blue})}{(\# \text{ of } \text{yellow}) + (\# \text{ of } \text{blue})}$$

$$= \frac{3}{4}$$

- Theoretical analysis**

node embedding = fair embedding + bias residual

Space of the Bias Residual $b^{(l)}$

Space of the Bias Residual $b^{(l+1)}$

Vanilla message passing

● : bias residual in majority group ● : bias residual in minority group

● : centroids of the majority or minority group ● : fair centroid in BeMap

- Method: BeMap**
 - (In every training epoch) neighbor sampling for balanced neighborhood and MP on it
 - Up to 80% bias reduction
 - Comparable or even better classification accuracy
 - More details in the paper

- Group Fairness on Graphs
- **Individual Fairness on Graphs**
- Degree Fairness on Graphs
- Future Directions and Q&A

Individual Fairness

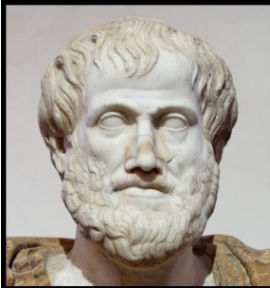
- **Definition**

- Similar individuals should have similar outcomes
- Rooted in Aristotle’s conception of justice as consistency

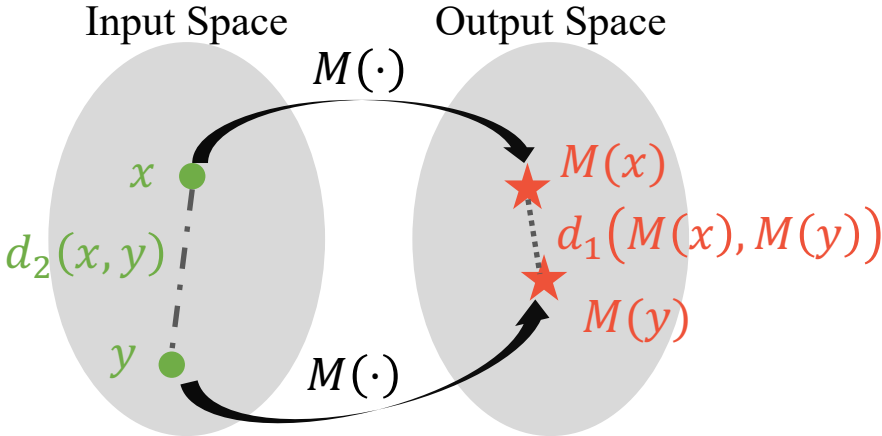
- **Formulation: Lipschitz inequality (most common)**

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

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



“ Equality consists in the same treatment of similar persons, and no government can stand which is not founded upon justice. ”



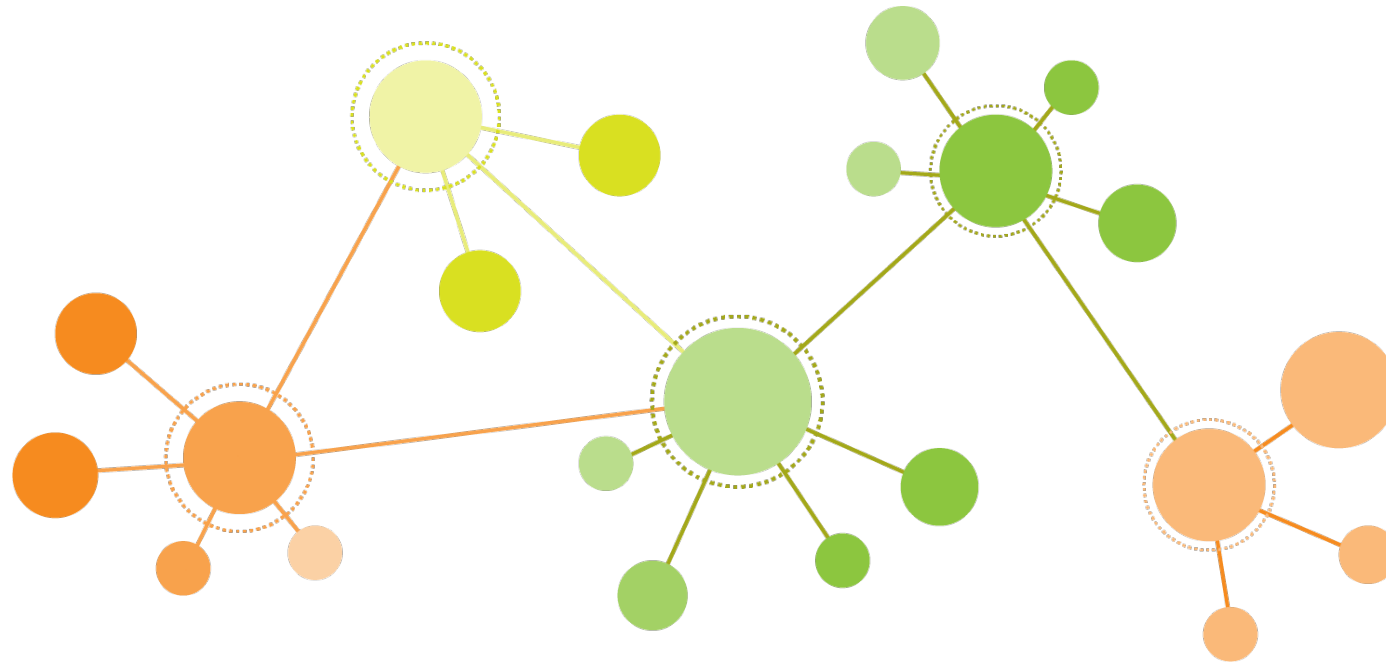
InFoRM: Individual Fairness on Graph Mining

- **Research questions**

RQ1. Measure: how to quantitatively measure individual bias?

RQ2. Algorithms: how to ensure individual fairness?

RQ3. Cost: what is the cost of individual fairness?



InFoRM Measure: Quantifying Individual Bias

- **Principle**

- Similar nodes \rightarrow similar mining results

- **Mathematical formulation**

(1) For any node pair (i, j)
 $\|Y[i, :] - Y[j, :]\|_F^2 S[i, j] \leq \epsilon$

(2) Sum it up for all node pairs

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

Similarity between node i and node j

- If $S[i, j]$ is high, $\frac{\epsilon}{S[i, j]}$ is small \rightarrow push $Y[i, :]$ and $Y[j, :]$ to be more similar
- Inequality should hold for **every** pairs of nodes i and j \rightarrow too restrictive

- **Relaxed criteria**

$$\sum_{i=1}^n \sum_{j=1}^n \|Y[i, :] - Y[j, :]\|_F^2 S[i, j] \leq m\epsilon$$

$$2\text{Tr}(Y^T L_S Y) \leq \delta$$

Overall individual bias of the graph

- m : number of edges in the graph
- $\delta = m\epsilon$



Alternative Measure: Ranking-Based Individual Fairness

- **Key challenge in InFoRM measure**

- Lipschitz condition (used in InFoRM)

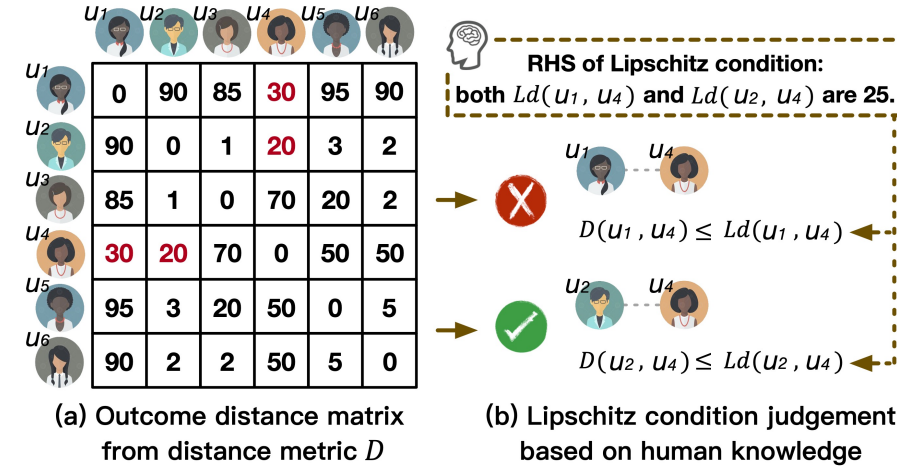
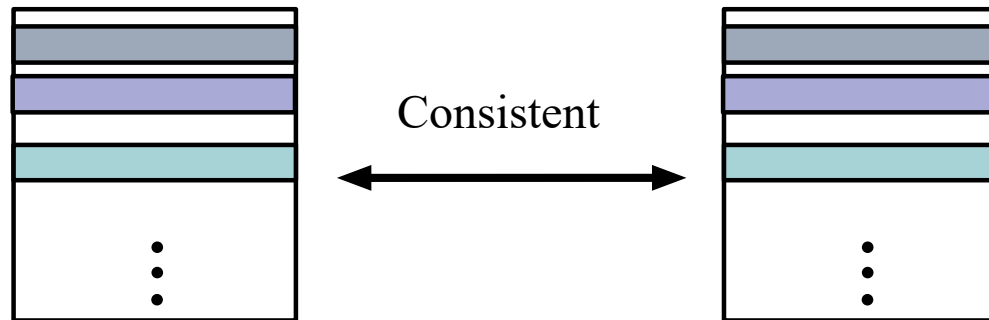
$$d_1(M(x), M(y)) \leq Ld_2(x, y)$$

- Distance comparison fails to calibrate between different individuals

- **Definition**

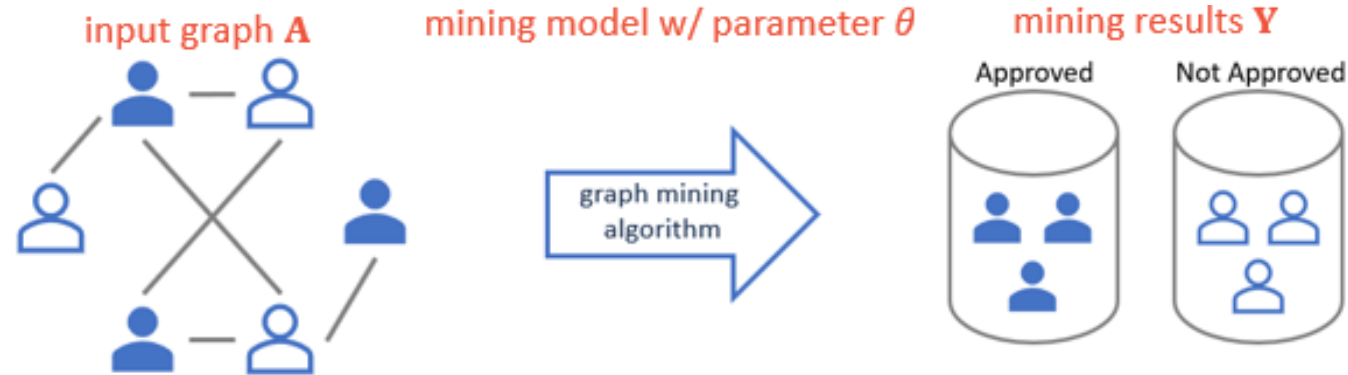
- Given
 - (1) the node similarity matrix \mathbf{S}_G of the input graph G
 - (2) the similarity matrix $\mathbf{S}_{\hat{Y}}$ of the GNN output \hat{Y}
- \hat{Y} is individually fair if, for each node i , it satisfies that

ranking list derived by $\mathbf{S}_G[i, :] =$ ranking list derived by $\mathbf{S}_{\hat{Y}}[i, :]$



InFoRM Measure: Mitigating Individual Bias

- Graph mining workflow



- Debiasing methods

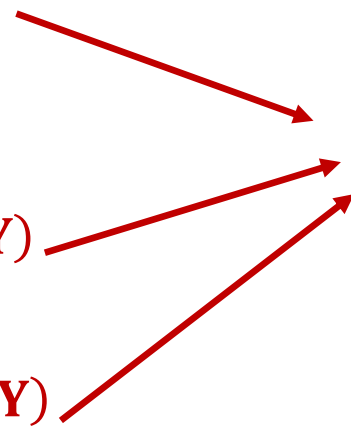
- Debiasing the input graph: $\min_{\mathbf{Y}} J = \underbrace{\|\tilde{\mathbf{A}} - \mathbf{A}\|_F^2}_{\text{topology consistency}} + \alpha \text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$

s. t. $\partial_{\mathbf{Y}} l(\tilde{\mathbf{A}}, \mathbf{Y}, \theta) = 0$

- Debiasing the mining model: $\min_{\mathbf{Y}} J = \underbrace{l(\mathbf{A}, \mathbf{Y}, \theta)}_{\text{task-specific loss function}} + \alpha \text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$

- Debiasing the mining results: $\min_{\mathbf{Y}} J = \underbrace{\|\mathbf{Y} - \bar{\mathbf{Y}}\|_F^2}_{\text{mining results consistency}} + \alpha \text{Tr}(\mathbf{Y}^T \mathbf{L}_S \mathbf{Y})$

Individual bias
(InFoRM measure)



InFoRM Cost: Characterizing Individual Bias

- **Main focus**

- Debiasing the mining results (model-agnostic)

- **Given**

- A graph with n nodes and adjacency matrix \mathbf{A}
- A node-node similarity matrix \mathbf{S}
- Vanilla mining results $\bar{\mathbf{Y}}$
- Debaised mining results $\mathbf{Y}^* = (\mathbf{I} + \alpha\mathbf{S})^{-1}\bar{\mathbf{Y}}$

- If $\|\mathbf{S} - \mathbf{A}\|_F = \Delta$, we have

$$\|\bar{\mathbf{Y}} - \mathbf{Y}^*\|_F \leq 2\alpha\sqrt{n} \left(\Delta + \sqrt{\text{rank}(\mathbf{A})\sigma_{\max}(\mathbf{A})} \right) \|\bar{\mathbf{Y}}\|_F$$

- **Key factors**

- The number of nodes n (i.e., size of the input graph)
- The difference Δ between \mathbf{A} and \mathbf{S}
- The rank of \mathbf{A} \rightarrow could be small due to (approximate) low-rank structures in real-world graphs
- The largest singular value of \mathbf{A} \rightarrow could be small if \mathbf{A} is normalized

- Group Fairness on Graphs
- Individual Fairness on Graphs
- **Degree Fairness on Graphs**
- Future Directions and Q&A

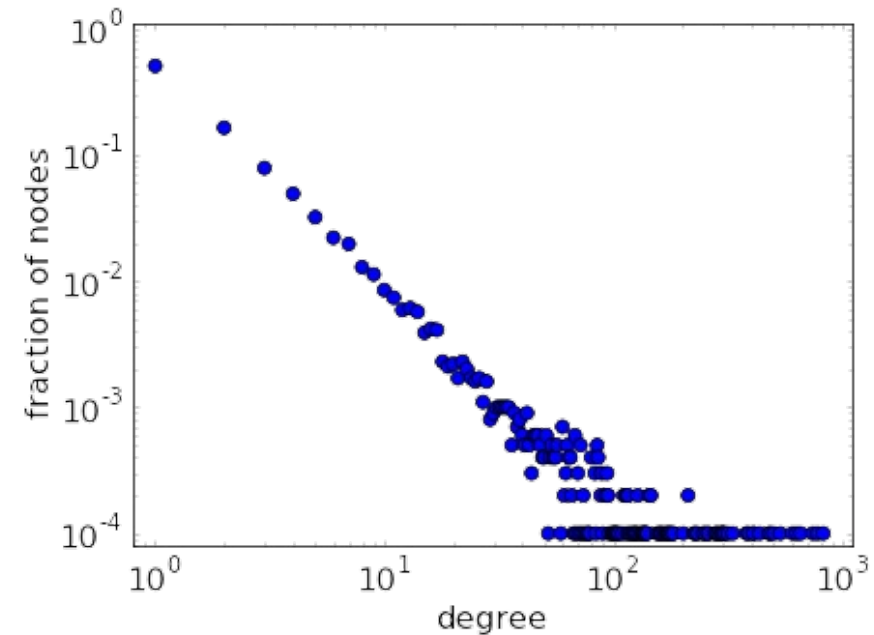
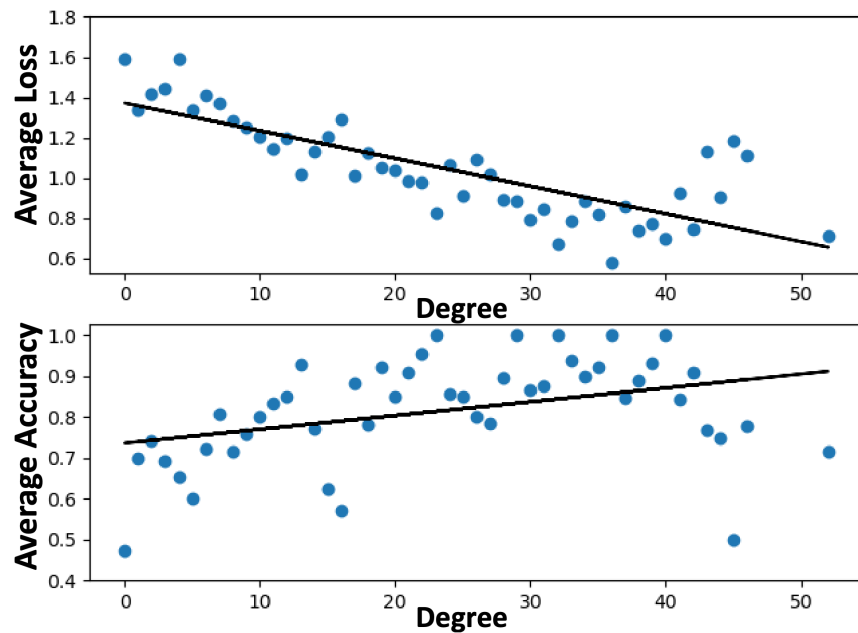
Degree Fairness: Definition and Motivation

- **Definition**

- Nodes of different degrees should have balanced utility on a graph mining task

- **Example: online advertising**

- (A small portion of) celebrities often enjoy high-quality model performance
- (A large portion of) grassroots users often suffer from bad model performance



Degree Unfairness: Pitfall of Graph Neural Networks

- **Given**

- (1) $\mathcal{G} = (\mathbf{A}, \mathbf{X})$
- (2) Any test node i in \mathcal{G} with label c
- (3) A graph learning model M which output (before softmax) \mathbf{Z}
- (4) Any wrong prediction $c' \neq c$

- **Our results**

- Misclassification rate

$$\Pr(\Pr(\hat{y} = c|i, M) > \Pr(\hat{y} = c'|M, i)) \leq \frac{1}{1 + R_{i,c'}}$$

where $R_{i,c'} = \frac{(\mathbb{E}[\mathbf{Z}[i,c'] - \mathbf{Z}[i,c]])^2}{\text{Var}[\mathbf{Z}[i,c'] - \mathbf{Z}[i,c]]}$ (reciprocal of measure of dispersion from economics)

- $R_{i,c'}$ is positively correlated with the degree of node i

- **Conclusion**

- high-degree nodes often have **lower misclassification rate!**

Causes #1: High-Degree Nodes with High Influence in Node Embeddings

- **Given**

- $\mathcal{V}_{\text{labeled}}$: a set of labeled nodes $\mathcal{V}_{\text{labeled}}$
- $\mathbf{W}^{(L)}$: the weight of L -th layer in an L -layer GCN
- d_i : degree of node i
- \mathbf{x}_i : input node feature of node i
- $\mathbf{h}_i^{(L)}$: output embeddings of node i learned by the L -layer GCN

- **Influence of node i on GCN training**

$$S(i) = \sum_{k \in \mathcal{V}_{\text{labeled}}} \left\| \mathbb{E} \left[\partial \mathbf{h}_i^{(L)} / \partial \mathbf{x}_k \right] \right\| \propto \sqrt{d_i} \|\mathbf{W}^{(L)}\| \sum_{k \in \mathcal{V}_{\text{labeled}}} \sqrt{d_k}$$

- **Remark**

- For two nodes i and j , if $d_i > d_j$, then $S(i) > S(j)$
→ Node with higher degree will have higher influence on GCN training

Solution #1: Degree-Specific Graph Convolution

- **Key idea**

- Degree-specific weights to encode degree information

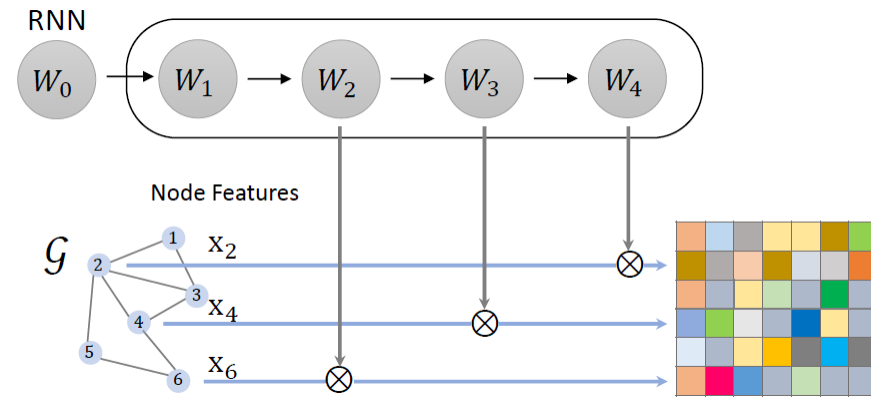
- **Given**

- d_i : the degree of node i
- $\mathbf{W}_{d_j}^{(l)}$: the degree-specific weight w.r.t. degree of node j

- **Degree-specific graph convolution**

$$\mathbf{h}_i^{(l+1)} = \sigma \left(\sum_{j \in \mathcal{N}_i \cup \{i\}} a_{ij} \left(\mathbf{W}^{(l)} + \mathbf{W}_{d_j}^{(l)} \right) \mathbf{h}_j^{(l)} \right)$$

- DEMO-Net $\rightarrow \mathbf{W}_{d_j}^{(l)}$ is generated randomly
- SL-DSGCN $\rightarrow \mathbf{W}_{d_j}^{(l)}$ is generated using a recurrent neural network

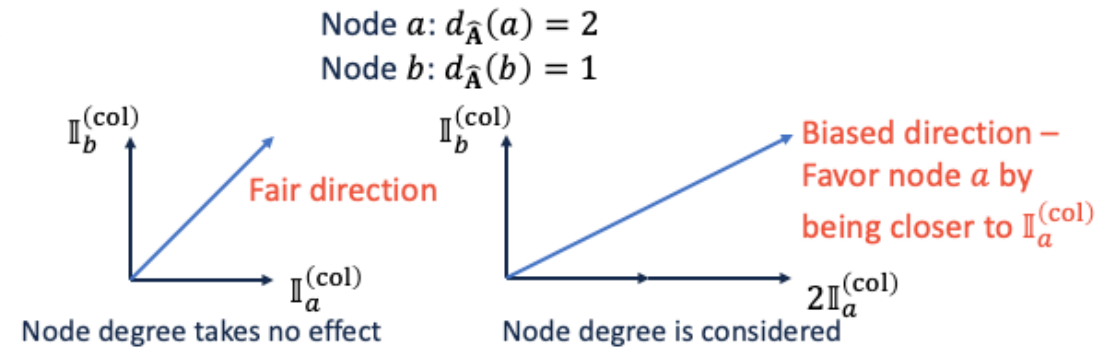


Causes #2: High-Degree Nodes with High Influence in Gradient

- Gradient of loss w.r.t. weight

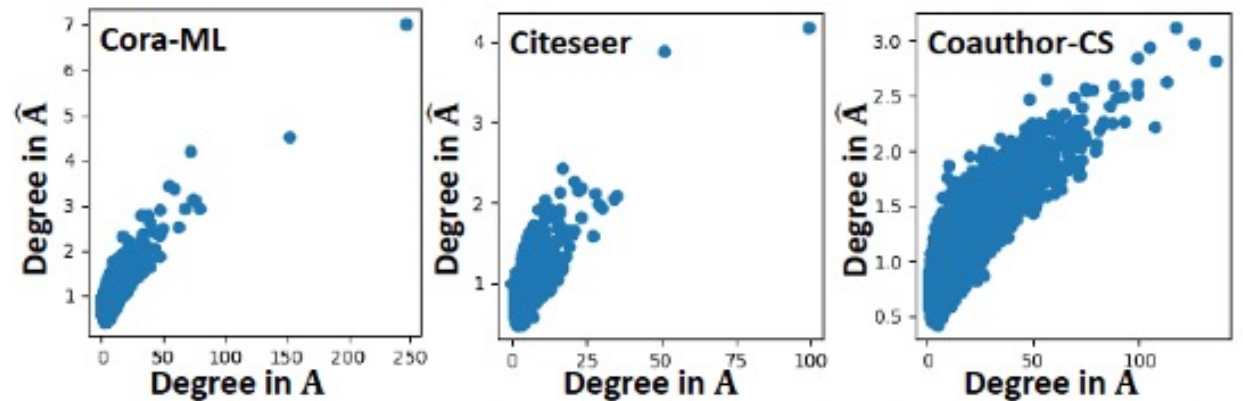
$$\frac{\partial J}{\partial \mathbf{W}^{(l)}} = \sum_{i=1}^n \overset{\text{Row sum in } \hat{\mathbf{A}}}{d_{\hat{\mathbf{A}}}(i)} \mathbb{I}_i^{(\text{col})} = \sum_{j=1}^n \overset{\text{Column sum in } \hat{\mathbf{A}}}{d_{\hat{\mathbf{A}}}(j)} \mathbb{I}_j^{(\text{row})}$$

- $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})\tilde{\mathbf{D}}^{-\frac{1}{2}} \rightarrow$ symmetric normalization kernel
- $\mathbb{I}_i^{(\text{col})}$ and $\mathbb{I}_j^{(\text{row})} \rightarrow$ the directions for gradient descent
- $d_{\hat{\mathbf{A}}}(i)$ and $d_{\hat{\mathbf{A}}}(j) \rightarrow$ the importance of the direction
- High degree \rightarrow more focus on that direction



Symmetric normalization

- Normalize the largest eigenvalue but not degree
- High degree in $\mathbf{A} \rightarrow$ high degree in $\hat{\mathbf{A}}$



Solution #2: Graph Normalization

- **Key idea**

- Mitigate impacts of node degree by normalizing it to constant (i.e., 1)
- Normalize the graph to a doubly stochastic graph

- **Sinkhorn-Knopp (SK) algorithm**

- Iteratively normalize row and columns
- **(Our result)** SK **always** finds the **unique** doubly stochastic form of symmetric normalization kernel

- **Fair gradient computation**

$$\left(\frac{\partial J}{\partial \mathbf{W}^{(l)}}\right)_{\text{fair}} = \left(\mathbf{H}^{(l-1)}\right)^T \hat{\mathbf{A}}_{\text{DS}}^T \frac{\partial J}{\partial \mathbf{E}^{(l)}}$$

- $\hat{\mathbf{A}}_{\text{DS}} \rightarrow$ doubly-stochastic normalization of $\hat{\mathbf{A}}$

- **RawlsGCN family**

- RawlsGCN-Graph: during **data pre-processing**, compute $\hat{\mathbf{A}}_{\text{DS}}$ and treat it as the input of GCN
- RawlsGCN-Grad: during **optimization (in-processing)**, treat $\hat{\mathbf{A}}_{\text{DS}}$ as a normalizer to equalize the importance of node influence

- Group Fairness on Graphs
- Individual Fairness on Graphs
- Degree Fairness on Graphs
- **Future Directions and Q&A**

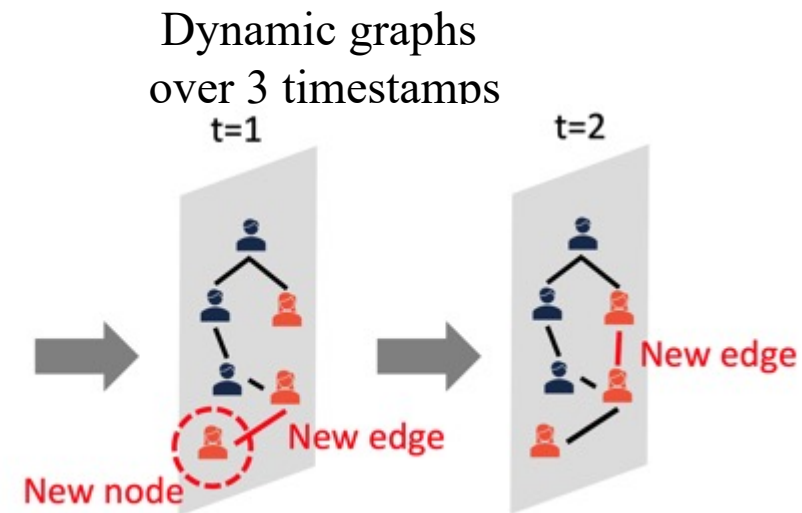
Future Direction #1: Fairness beyond Plain and Static Graphs

- **Observation**

- Real-world graphs are often dynamic and/or multi-sourced

- **Research questions**

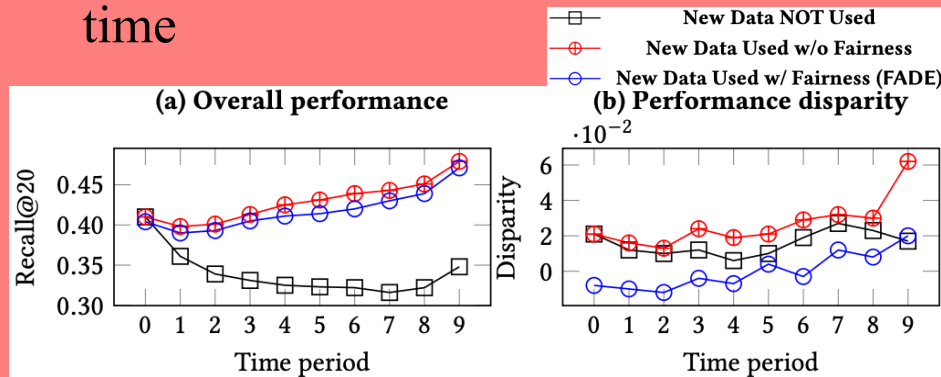
- How to ensure fairness for multiple type of nodes/edges or multi-graphs?
- How to efficiently update the fair mining results at each timestamp?
- How to characterize the impact of graph dynamics and multiple sources over the bias measure?



Preliminary Work: Dynamic Group Fairness in Recommender Systems

• Observation

- performance disparity is getting larger over time



• Method: FADE

- Model-agnostic
- Fine-tuning with newly observed data
- Periodically re-training to keep historical information
- Linear complexity w.r.t. # new data

• Theory

- Fine-tuning is better than re-training for fairness over time

Re-training

- (1) $L_{t_{\text{test}}}^{\text{rt}}$ = real loss of re-training at test time; (2) $L_{t_{\text{test}}}^*$ = optimal loss at time test; (3) m_0 = #. edges at time 0; (4) m_t = #. edge changes at time t ; (5) $0 < \gamma < 1$

$$L_{t_{\text{test}}}^{\text{rt}} \leq L_{t_{\text{test}}}^* + 2 \frac{m_0 d_{0,t_{\text{test}}} + \sum_{t=1}^{t_{\text{test}}-1} m_1 d_{1,t_{\text{test}}}}{m_0 + (t_{\text{test}} - 1)m_1} + 4 \sqrt{\frac{1}{m_0 + (t_{\text{test}} - 1)m_1} \log \frac{2}{\delta}}$$

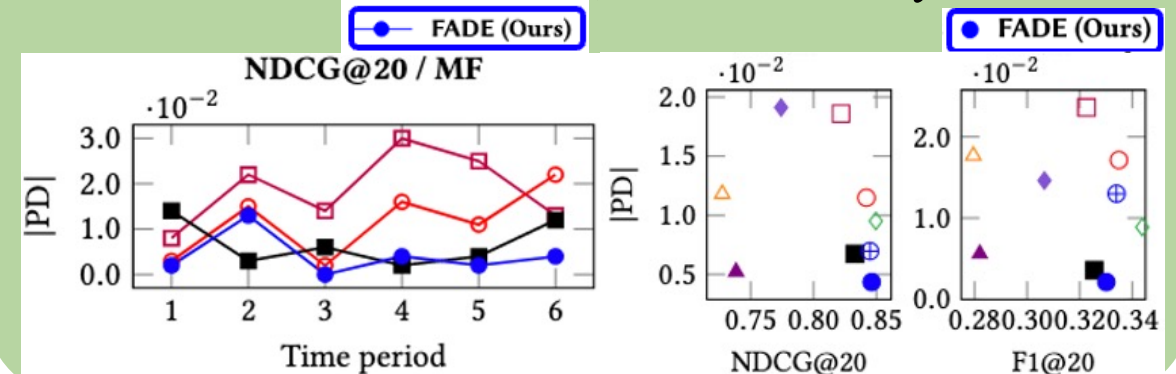
Fine-tuning

- Similar settings as re-training but $L_{t_{\text{test}}}^{\text{ft}}$ = real loss of fine-tuning at test time

$$L_{t_{\text{test}}}^{\text{ft}} \leq L_{t_{\text{test}}}^* + 2 \frac{(1-\gamma) \left(2 \sum_{t=0}^{t_{\text{test}}-1} \gamma^{t_{\text{test}}-t-1} d_{t,t_{\text{test}}} + 4 \sqrt{\frac{\gamma^{2t_{\text{test}}-2} + \frac{1-\gamma^{2t_{\text{test}}-2}}{(1-\gamma^2)m_1}}{\log \frac{2}{\delta}}} \right)}{1-\gamma^{t_{\text{test}}}}$$

• Results

- Fairness over time, small accuracy decrease



Future Direction #2: Fairness on Graphs → Fairness with Graphs

- **Fairness on graphs**

- Graph as data
- Nodes = entities
- Social networks → nodes = users
- Citation networks → nodes = papers
- Web graph → nodes = webpages

- **Fairness with graphs**

- Graph as context
- Nodes = models/datasets/modalities

- **Example: supply chain**

1. Demand + supply for medical resources
2. Models to allocate medical resources



- How can we leverage demand + supply + model collectively for fair supply chain?

Future Direction #3: Benchmark and Evaluation Metrics

- **Observation**

- No consensus on the experimental settings for fair graph learning
- Which data to compare? What sensitive attribute to consider?
- Which evaluation metrics for each type of fairness?

- **Consequences**

- Different settings for different research works
- Hardly fair comparison among fair graph learning methods
- Hardly deployable methods in real-world scenarios

- **Call for community effort**

- Evaluation benchmark for consistent experimental settings and fair comparison
- Collection of large-scale, realistic, but challenging dataset for evaluation

Outline

- Introduction and Background
- Topology Issues
- Imbalance Issues
- Short Break
- Bias and Fairness Issues
- **Limited Labeled Data Issues**
- Abnormal Graph Data Issues
- Summary

- Graph Data Augmentations
- Self-supervised Learning on Graphs

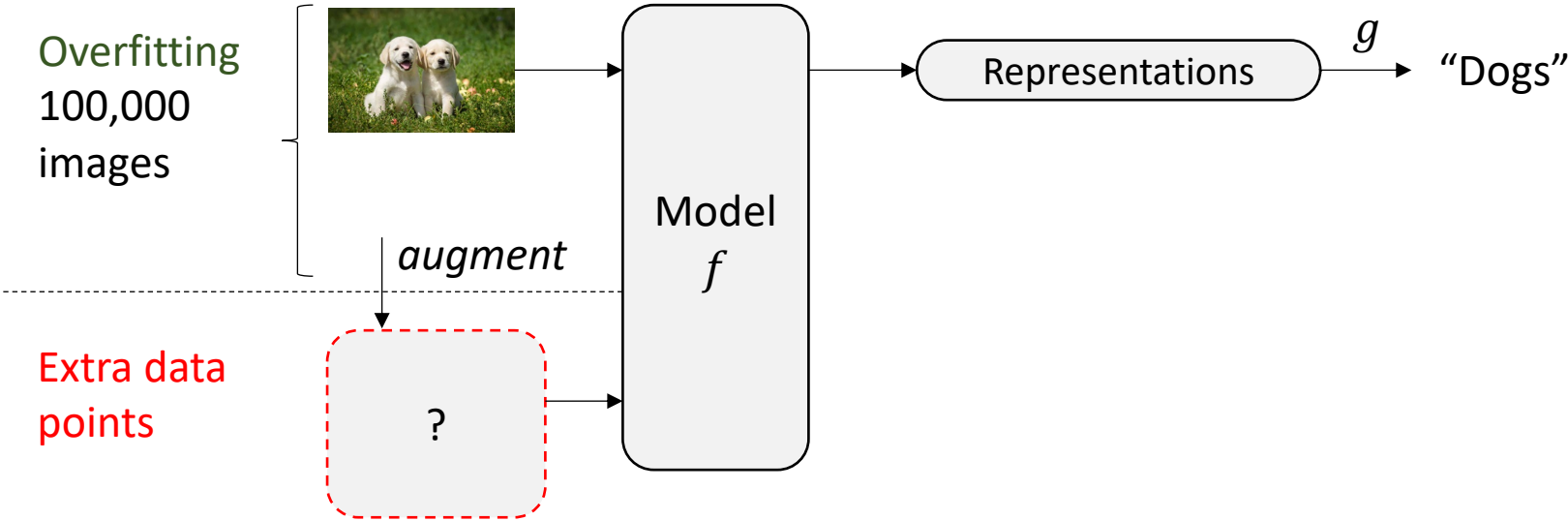
Data Augmentation

Wikipedia: Techniques used to **increase the amount of data** by adding *slightly modified* copies of already existing data or *newly created* synthetic data **from existing data**.

- Why data augmentation?
 - It helps reduce overfitting when training a machine learning model.
 - The acquisition of labeled graph data can be expensive.

Data Augmentation

Wikipedia: Techniques used to **increase the amount of data** by adding *slightly modified* copies of already existing data or *newly created* synthetic data **from existing data**.



Data Augmentation

Wikipedia: Techniques used to **increase the amount of data** by adding *slightly modified* copies of already existing data or *newly created* synthetic data **from existing data**.

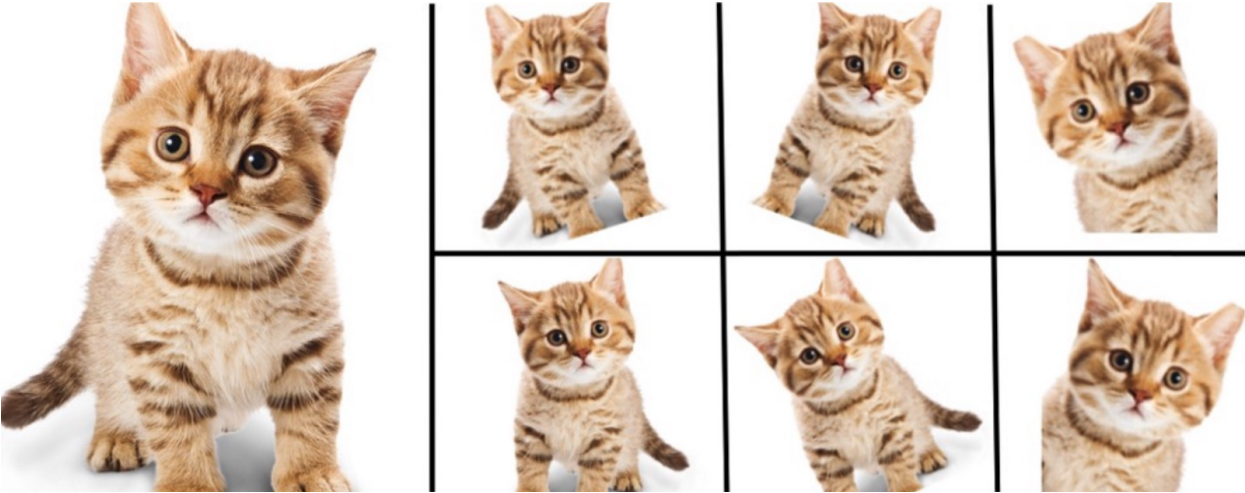
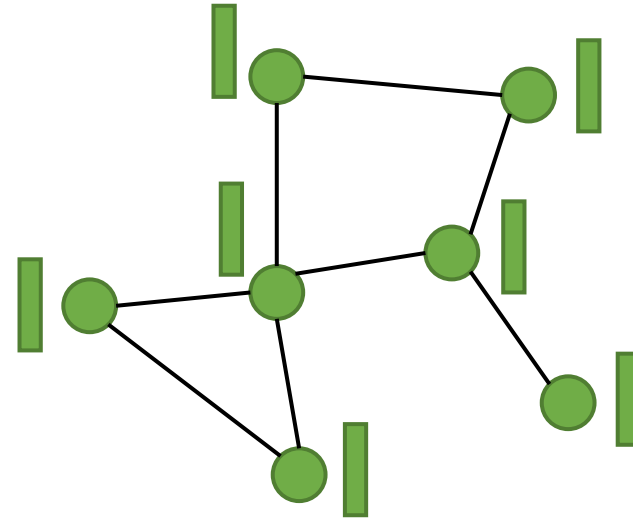


Image sources:
<https://www.kdnuggets.com/2018/05/data-augmentation-deep-learning-limited-data.html>
<https://amitnss.com/2020/05/data-augmentation-for-nlp/>

Graph Data Augmentation

- Structure Augmentation
 - Drop/add nodes/edges, etc.
- Feature Augmentation
 - Mask off features, etc.
- Label Augmentation
 - Label propagation, etc.



Graph Data Augmentation

- Rule-based augmentations
 - Designed based on heuristic rules
 - Usually efficient and scalable
 - Simple and easy to implement
 - Commonly used in self-supervised learning
- Learned augmentations
 - Involve learning during augmentation
 - Augmented data better fits GML models
 - Better performances in supervised learning



Rule-based Graph Data Augmentation Approaches

Methodology	Representative Works	Task Level			Augmented Data		
		Node	Graph	Edge	Structure	Feature	Label
Rule-based GDA	Stochastic Dropping/Masking	DropEdge [87]	✓			✓	
		DropNode [27]		✓			✓
		NodeDropping [127]		✓		✓	
		Feature Masking [100]	✓				✓
		Feature Shuffling [106]	✓				✓
		DropMessage [23]	✓		✓		✓
	Subgraph Cropping/Substituting	Subgraph Masking [127]		✓	✓	✓	✓
		GraphCrop [111]		✓		✓	
		M-Evolve [145]		✓		✓	
	Virtual Node	MoCL [97]		✓		✓	✓
Graphormer [125]			✓		✓		
GNN-CM ⁺ /CM [45]			✓	✓			
Mixup	Graph Mixup [115]	✓	✓				✓
	ifMixup [37]		✓		✓	✓	✓
	Graph Transparent [85]		✓		✓	✓	✓
	G-Mixup [39]		✓		✓	✓	✓
SMOTE	GraphSMOTE [140]	✓				✓	
	GATSMOTE [75]	✓			✓		
	GNN-CL [70]	✓			✓	✓	
Diffusion	GDA [60]	✓			✓		
Counterfactual Augmentation	CFLP [141]			✓	✓		✓
Attribute Augmentation	LA-GNN [74]	✓				✓	
	SR+DR [93]	✓				✓	
Pseudo-labeling	Label Propagation [147]	✓					✓
	PTA [21]	✓					✓

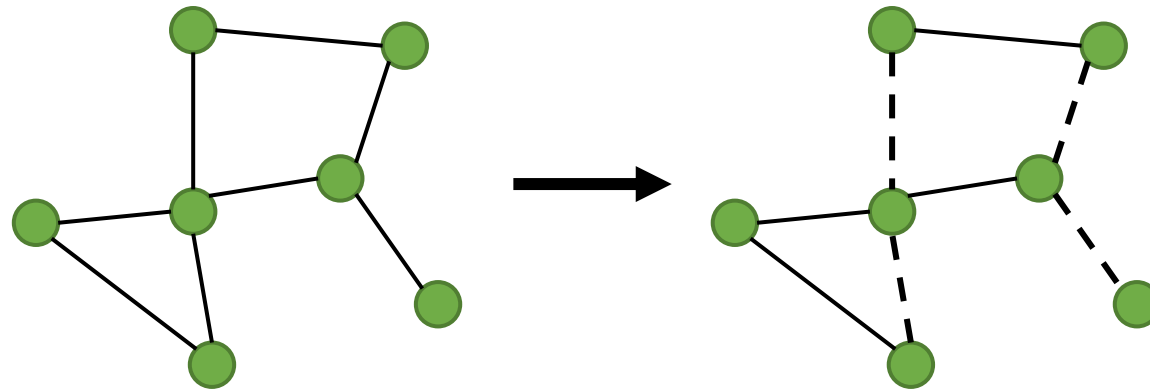
DropEdge

- Dropout on edges: randomly remove some edges at the beginning of every training epoch.

$$\tilde{\mathbf{A}} = \mathbf{M} \odot \mathbf{A}$$

$$\mathbf{M} \in \{0, 1\}^{N \times N} \text{ s.t. } M_{i,j} = \text{Bernoulli}(\varepsilon)$$

- Prevents overfitting and over-smoothing.



Other Stochastic Masking/Dropping Methods

- Node Dropping
 - Randomly removing part of the nodes.
- Feature Masking
 - Randomly mask off node features.
 - Random row-shuffling on node feature matrix \mathbf{X} .
- Subgraph Masking
 - Randomly mask off a connected subgraph.

Feng, et al. Graph Random Neural Networks for Semi-supervised Learning on Graphs. NeurIPS 2020.

You, et al. Graph Contrastive Learning with Augmentations. NeurIPS 2020.

Thakoor, et al. Large-scale Representation Learning on Graphs via Bootstrapping. ICLR 2022.

Velickovic, et al. Deep Graph Infomax. ICLR 2019.

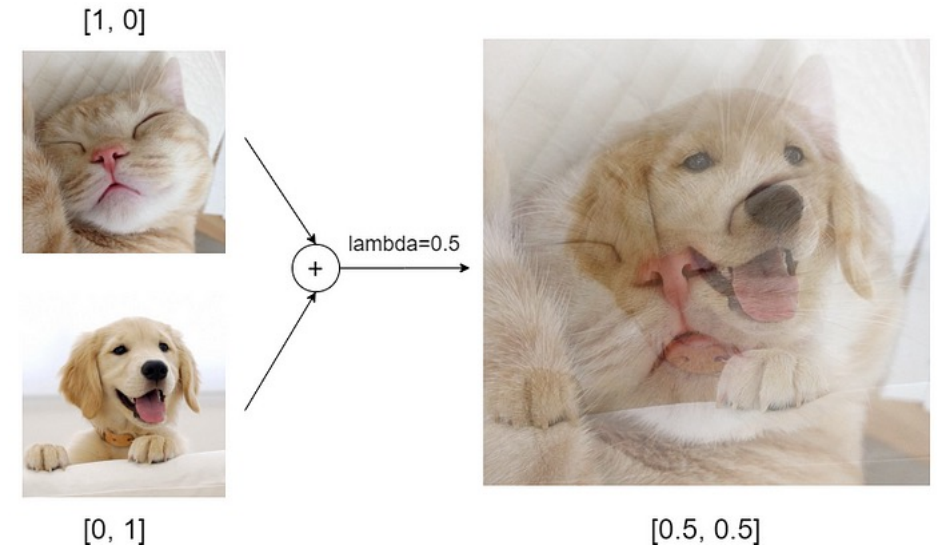


Mixup

- Mixup: generates a weighted combination of random pairs from the training data.

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j,$$
$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j.$$

- Manifold Mixup: interpolating hidden states.

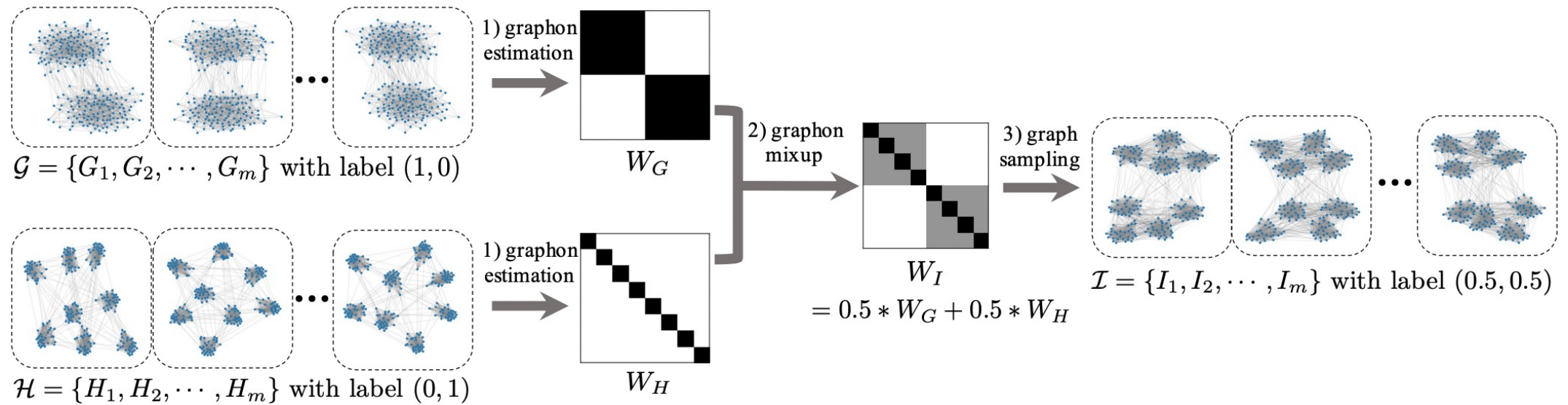


Zhang, et al. Mixup: Beyond Empirical Risk Minimization. ICLR 2018.

Verma, et al. Manifold Mixup: Better Representations by Interpolating Hidden States. ICML 2019.

Image source: <https://medium.com/@wolframalphav1.0/easy-way-to-improve-image-classifier-performance-part-1-mixup-augmentation-with-codes-33288db92de5>

G-Mixup



1. Graphon estimation:
2. Graphon Mixup:
3. Graph Generation:
4. Label Mixup:

$$\mathcal{G} \rightarrow W_G, \mathcal{H} \rightarrow W_H$$

$$W_I = \lambda W_G + (1 - \lambda) W_H$$

$$\{I_1, I_2, \dots, I_m\} \stackrel{\text{i.i.d}}{\sim} \mathbb{G}(K, W_I)$$

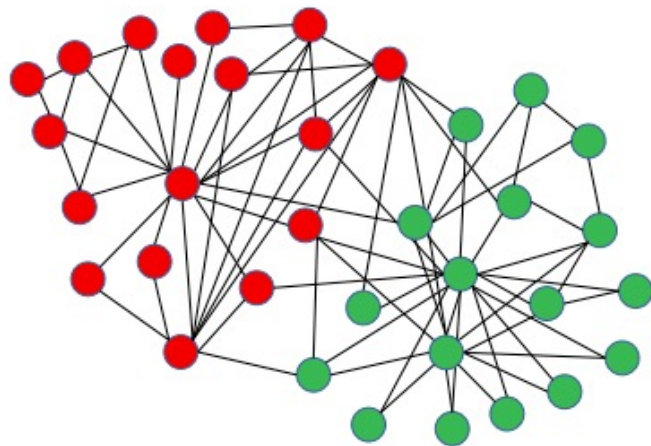
$$\mathbf{y}_I = \lambda \mathbf{y}_G + (1 - \lambda) \mathbf{y}_H$$

Learned Graph Data Augmentation Approaches

Learned GDA	Graph Structure Learning	GAug [140] GLCN [47] LDS [28] ProGNN [50] Eland [141]	✓ ✓ ✓ ✓ ✓		✓ ✓ ✓ ✓ ✓
	Graph Adversarial Training	RobustTraining [125] AdvT [18] FLAG [63] GraphVAT [25]	✓ ✓ ✓ ✓	✓ ✓	✓ ✓
	Graph Rationalization	GREA [71] AdvCA [97]		✓ ✓	✓ ✓
	Automated Augmentation	AutoGDA [144] GraphAug [79] JOAO [130] MolCLE [116]	✓		✓ ✓ ✓ ✓

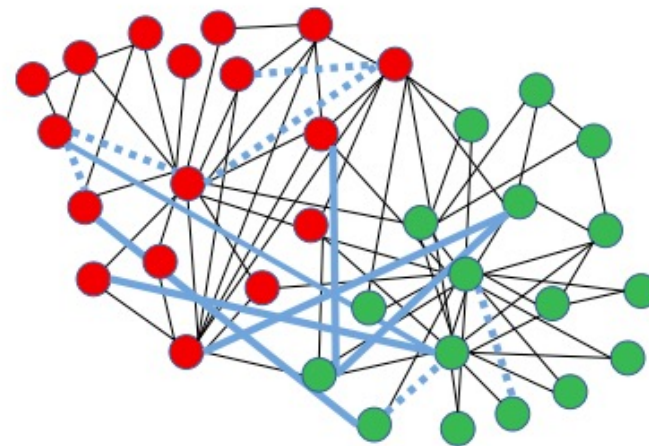
Limitations of Rule-based Approaches

Do not leverage task information and could hurt the downstream performance



(a) Original graph.

F1 Score: 92.4



(b) Random mod.

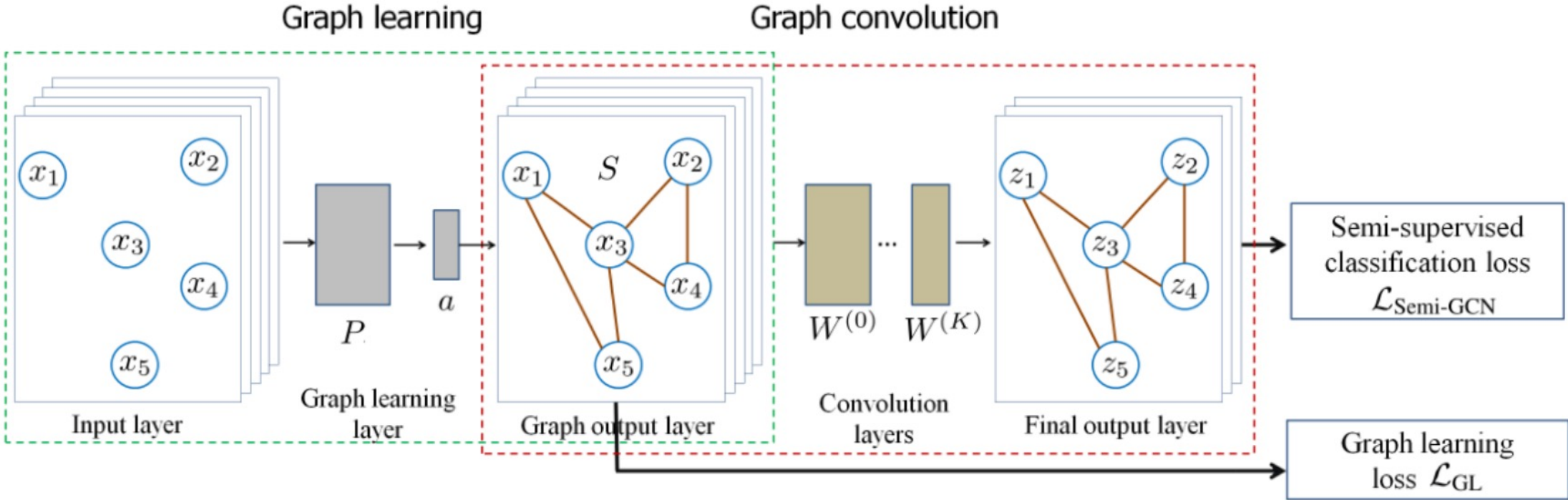
F1 Score: 91.0

Learned Graph Data Augmentation Approaches

- Graph Structure Learning
 - Augment data with good graph structures
- Adversarial Training
 - Augment data with adversarial examples
- Rationalization
 - Augment data by changing graph environment
- Automated Augmentation
 - Automatically combine different augmentations

Graph Structure Learning

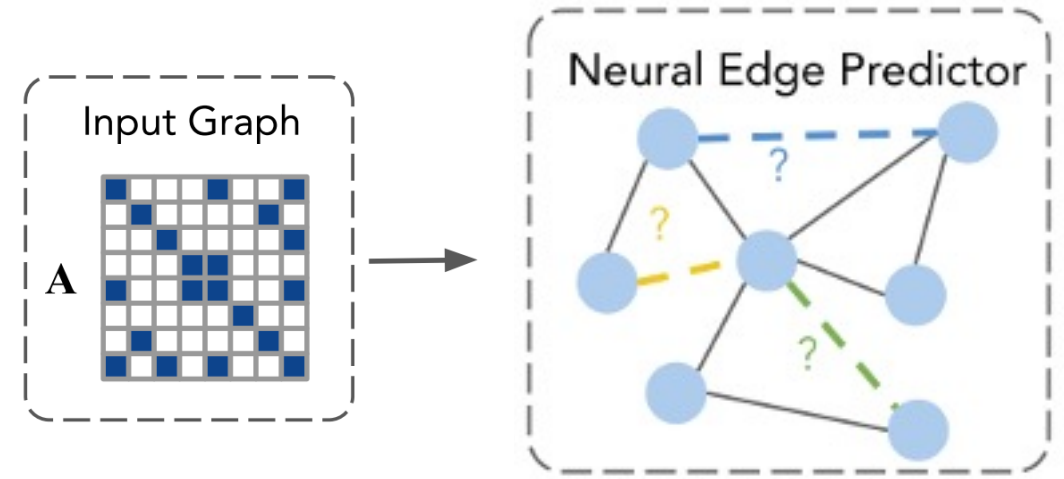
Graph Learning + Graph Convolution



GAug: Neural Edge Predictor

What are better graph structures?

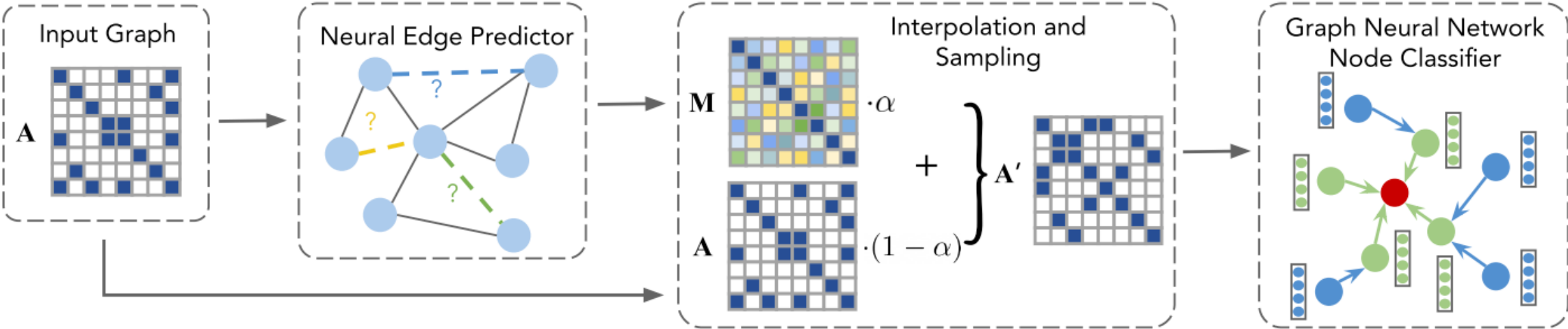
- “Noisy” edges should be removed
Inter-class edges
- “Missing” edges should be added
Intra-class edges



$$\underline{\mathbf{M}} = \sigma(\mathbf{Z}\mathbf{Z}^T), \text{ where } \mathbf{Z} = f_{GCL}^{(1)}\left(\mathbf{A}, f_{GCL}^{(0)}(\mathbf{A}, \mathbf{X})\right)$$

M models node similarities

GAug: Interpolation and Sampling



$$P_{ij} = \alpha M_{ij} + (1 - \alpha) A_{ij}$$

↓ Bernoulli sampling

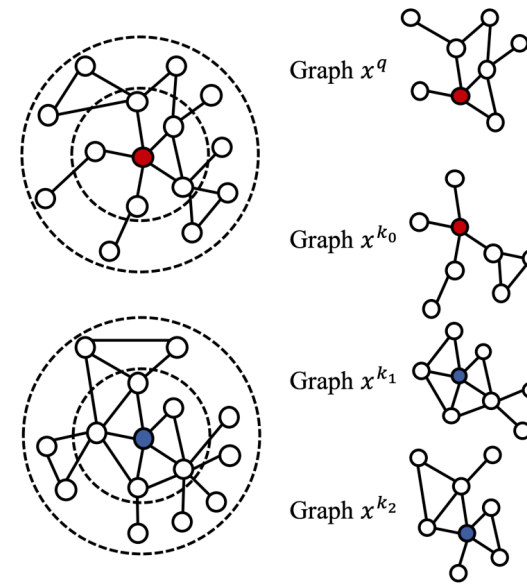
$$A'_{ij}$$

Graph Self-supervised Learning

- **Graph Self-Supervised Learning** aims to learn generalizable **node/edge/graph** representations **without** using any human-annotated labels
 - **Graph Generative Modeling**
 - Learn generalizable representations by **reconstructing** the node features or/and graph structure

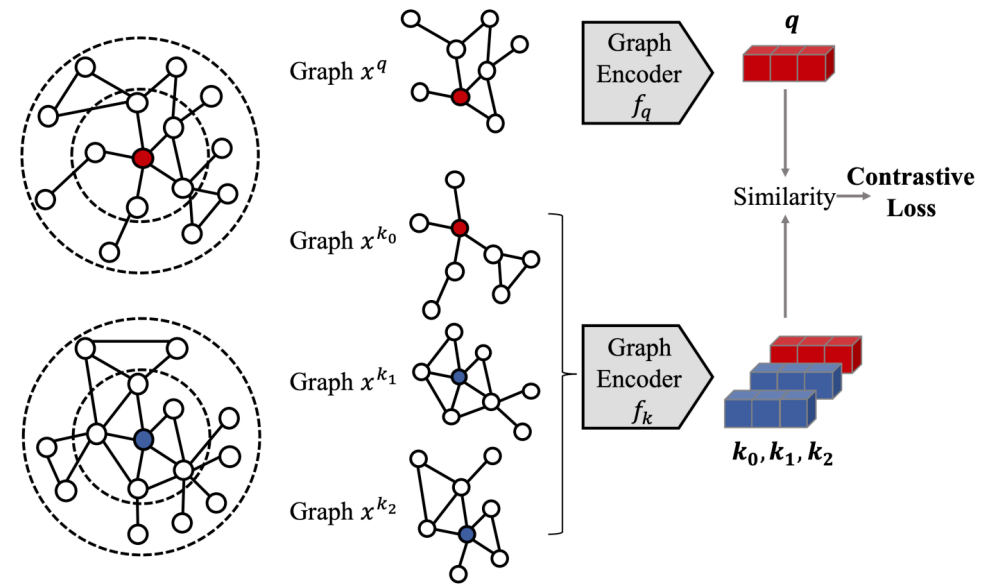
Graph Self-supervised Learning

- **Graph Self-Supervised Learning** aims to learn generalizable **node/edge/graph** representations **without** using any human-annotated labels
 - **Graph Generative Modeling**
 - Learn generalizable representations by **reconstructing** the node features or/and graph structure
 - **Graph Contrastive Learning (GCL)**
 - Create different views from the unlabeled input graph via data augmentation



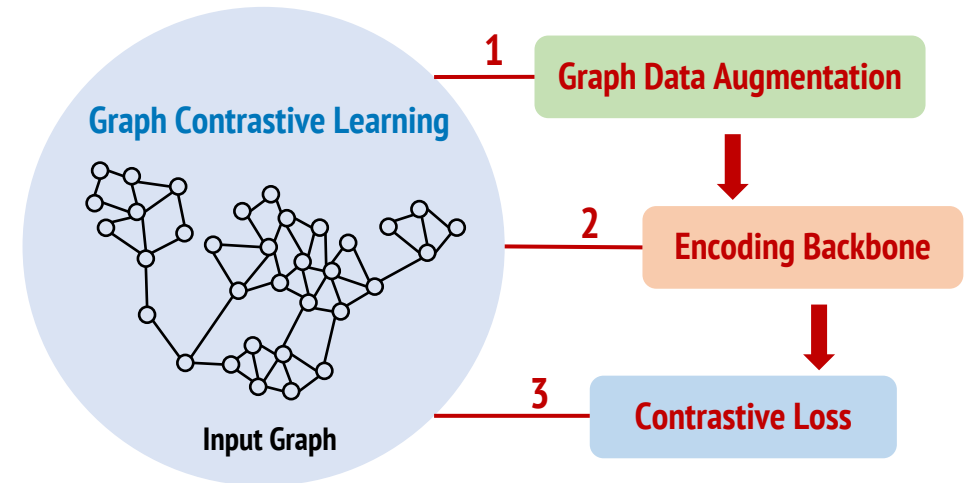
Graph Self-supervised Learning

- **Graph Self-Supervised Learning** aims to learn generalizable **node/edge/graph** representations **without** using any human-annotated labels
 - **Graph Generative Modeling**
 - Learn generalizable representations by **reconstructing** the node features or/and graph structure
 - **Graph Contrastive Learning (GCL)**
 - Create different views from the unlabeled input graph via data augmentation
 - Maximize the agreement between representations of different augmented views of the same instance



Typical Unsupervised Graph Contrastive Learning

- **Graph Data Augmentation**
 - Create different views of each instance (e.g., node, subgraph)
 - Arbitrary graph data augmentation (e.g., edge dropping, feature masking)
- **Encoding Backbone**
 - Encode different augmented views
 - Shallow GNNs (e.g., 2-layer GCN)
- **Contrastive Loss**
 - Maximize the agreement between representations learned from different augmented views
 - Instance-level contrastive learning



Outline

- Introduction and Background
- Topology Issues
- Imbalance Issues
- Short Break
- Bias Issue
- Limited Labeled Data Issues
- **Abnormal Graph Data Issues**
- Summary

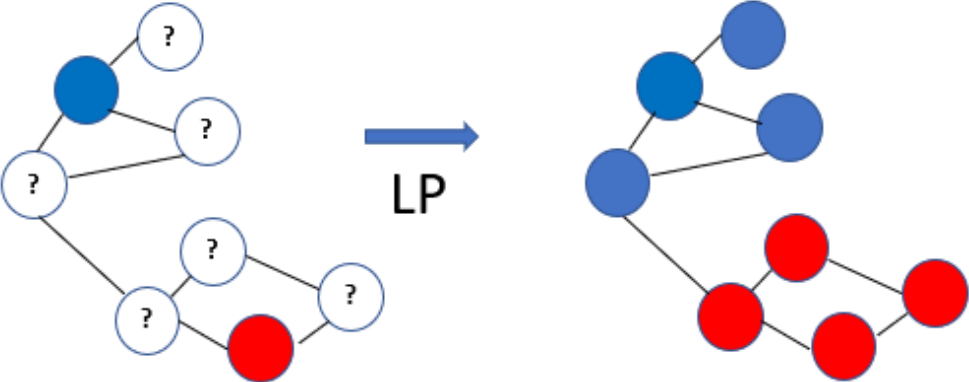
Abnormal Graph Data Issues

- Missing Data
- Adversarial Attacked Data

Missing Data

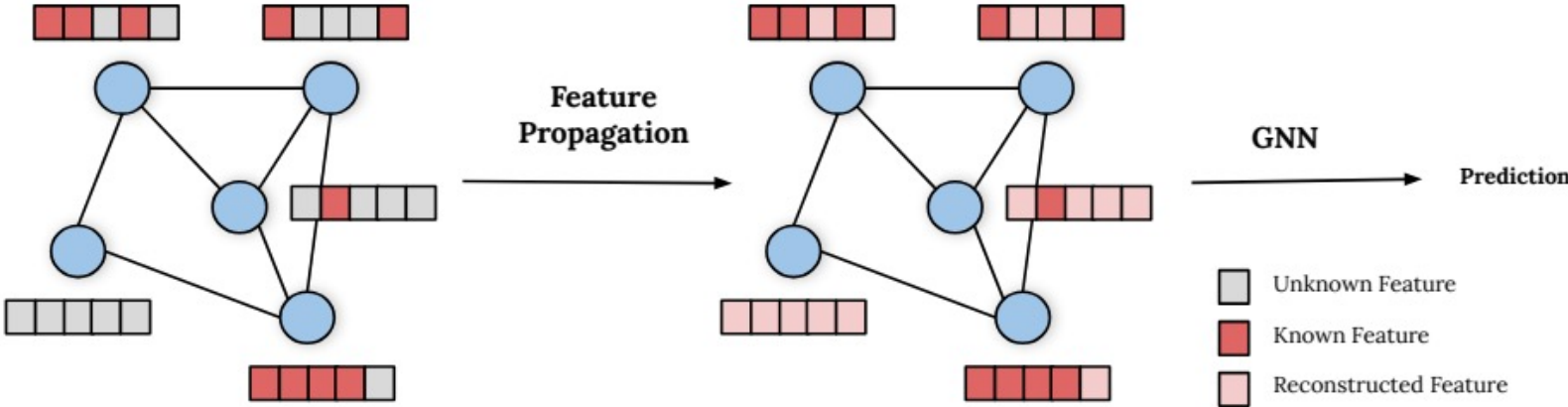
There are various solutions to deal with **missing labels**:

- Label propagation (LP)
- Self-supervised learning
- Unsupervised learning
- ...



What if we have **missing features**?

- Feature propagation



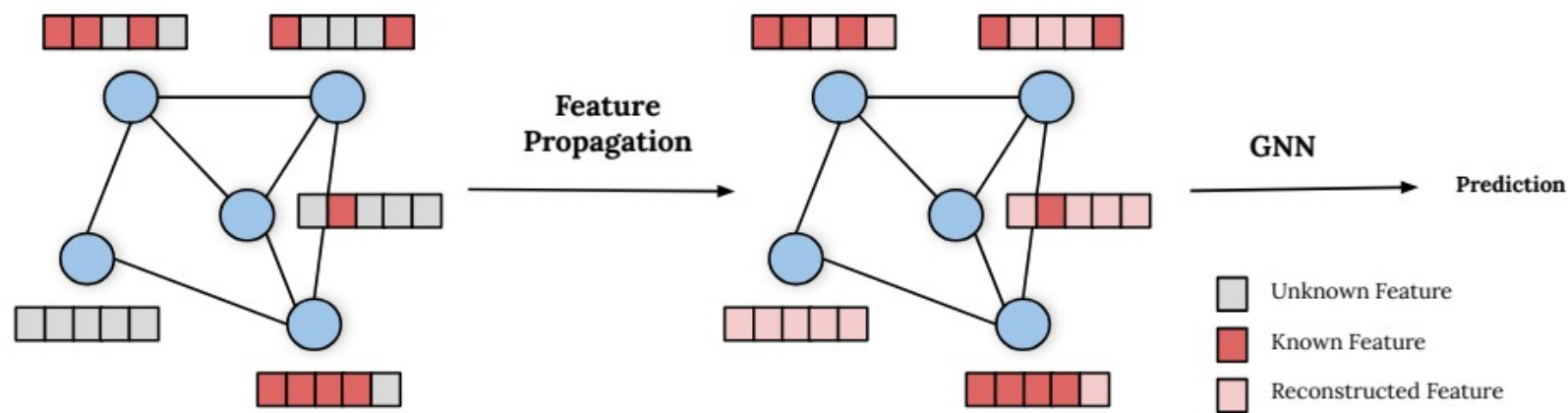
Missing Data

What if we have **missing features**?

- Feature propagation

Algorithm 1 Feature Propagation

- 1: **Input:** feature vector \mathbf{x} , diffusion matrix $\tilde{\mathbf{A}}$
 - 2: $\mathbf{y} \leftarrow \mathbf{x}$
 - 3: **while** \mathbf{x} has not converged **do**
 - 4: $\mathbf{x} \leftarrow \tilde{\mathbf{A}}\mathbf{x}$ ▷ Propagate features
 - 5: $\mathbf{x}_k \leftarrow \mathbf{y}_k$ ▷ Reset known features
 - 6: **end while**
-



Missing Data

Comparison of Feature Propagation to Label Propagation

Feature Propagation:

- Propagates features (continuous)
- Prediction is made by a GNN on top of the propagated features
- Uses features, and a low % of them being present is enough for good performance

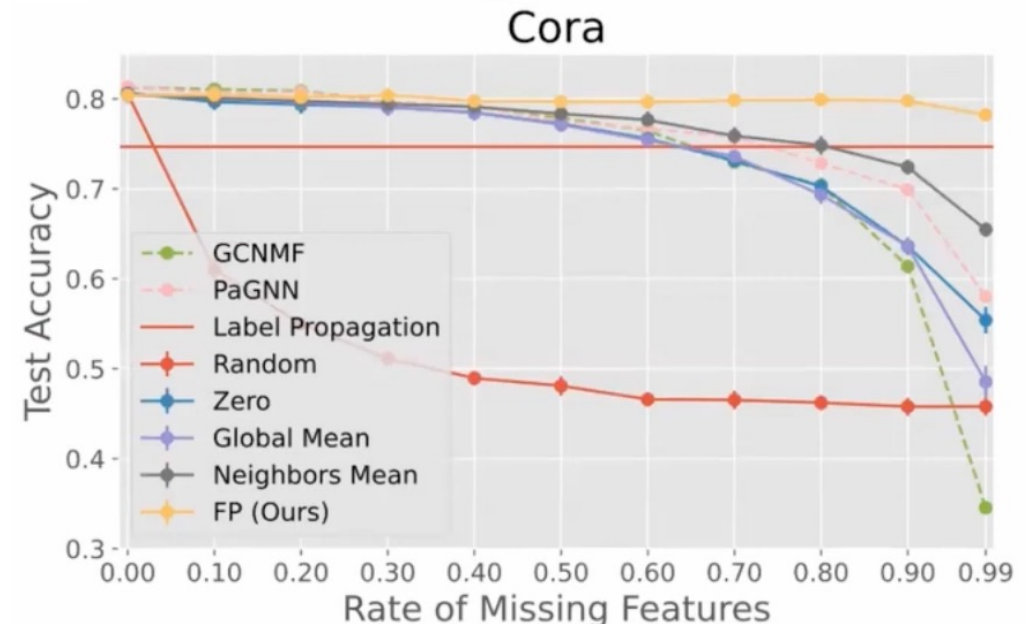
Experiment Results

Across different levels of missing features, Feature Propagation achieves the best performance



Label Propagation:

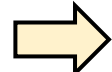
- Propagates class labels (discrete)
- Prediction is obtained directly from propagating class labels
- Feature-agnostic



Missing Data

Beyond missing features on graphs, **can we solve the general missing data problem?**

Data Matrix with Missing Values					Labels
	F_1	F_2	F_3	F_4	Y
O_1	0.3	0.5	NA	0.1	y_1
O_2	NA	NA	0.6	0.2	y_2
O_3	0.3	NA	NA	0.5	?

- Two ways of approaching missing data problems:
- 
- **Feature imputation:** missing feature values are estimated based on observed values
 - **Label prediction:** downstream labels are learned directly from incomplete data

Issues:

- Existing methods fail to make full use of feature values from other observations
- Existing methods tend to make biased assumptions about the missing values by initializing them with special default values

Missing Data






GRAPE: reformulate the tasks as graph tasks

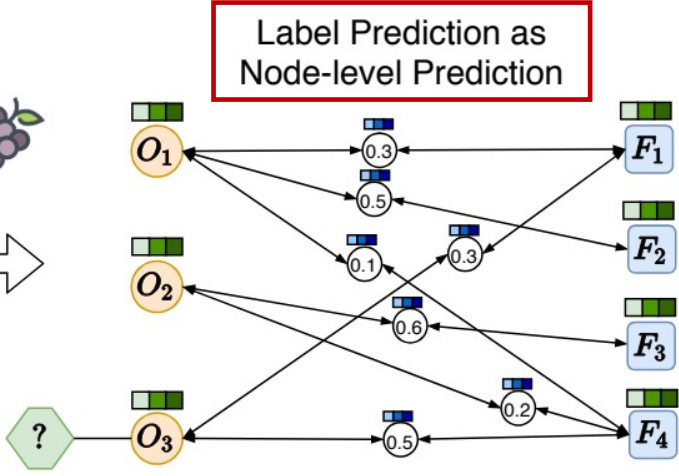
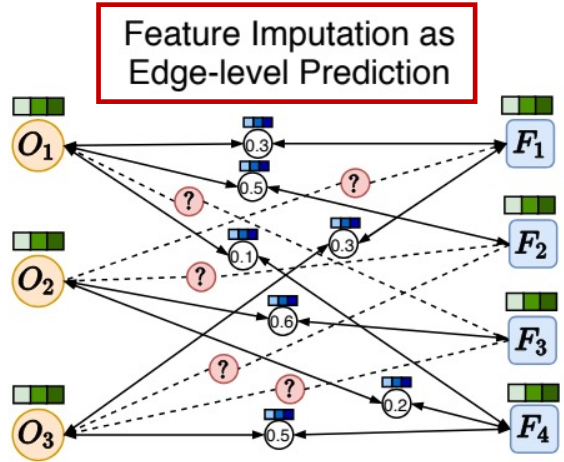
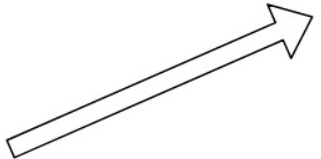
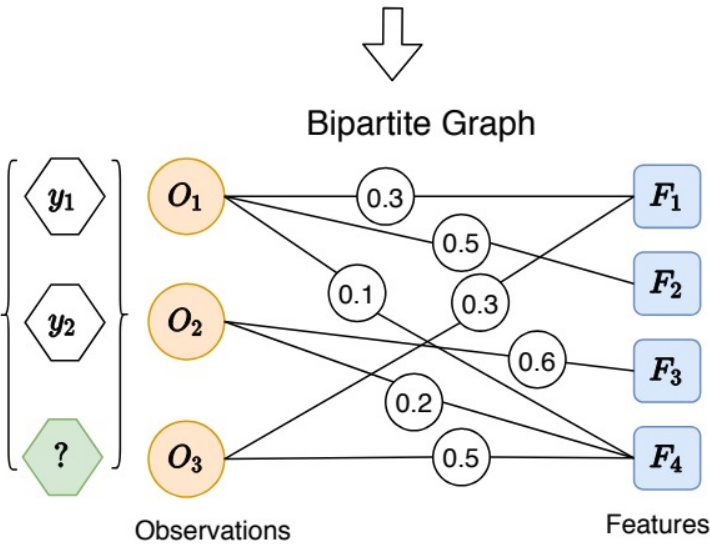
Data Matrix with Missing Values

	F_1	F_2	F_3	F_4
O_1	0.3	0.5	NA	0.1
O_2	NA	NA	0.6	0.2
O_3	0.3	NA	NA	0.5

Labels

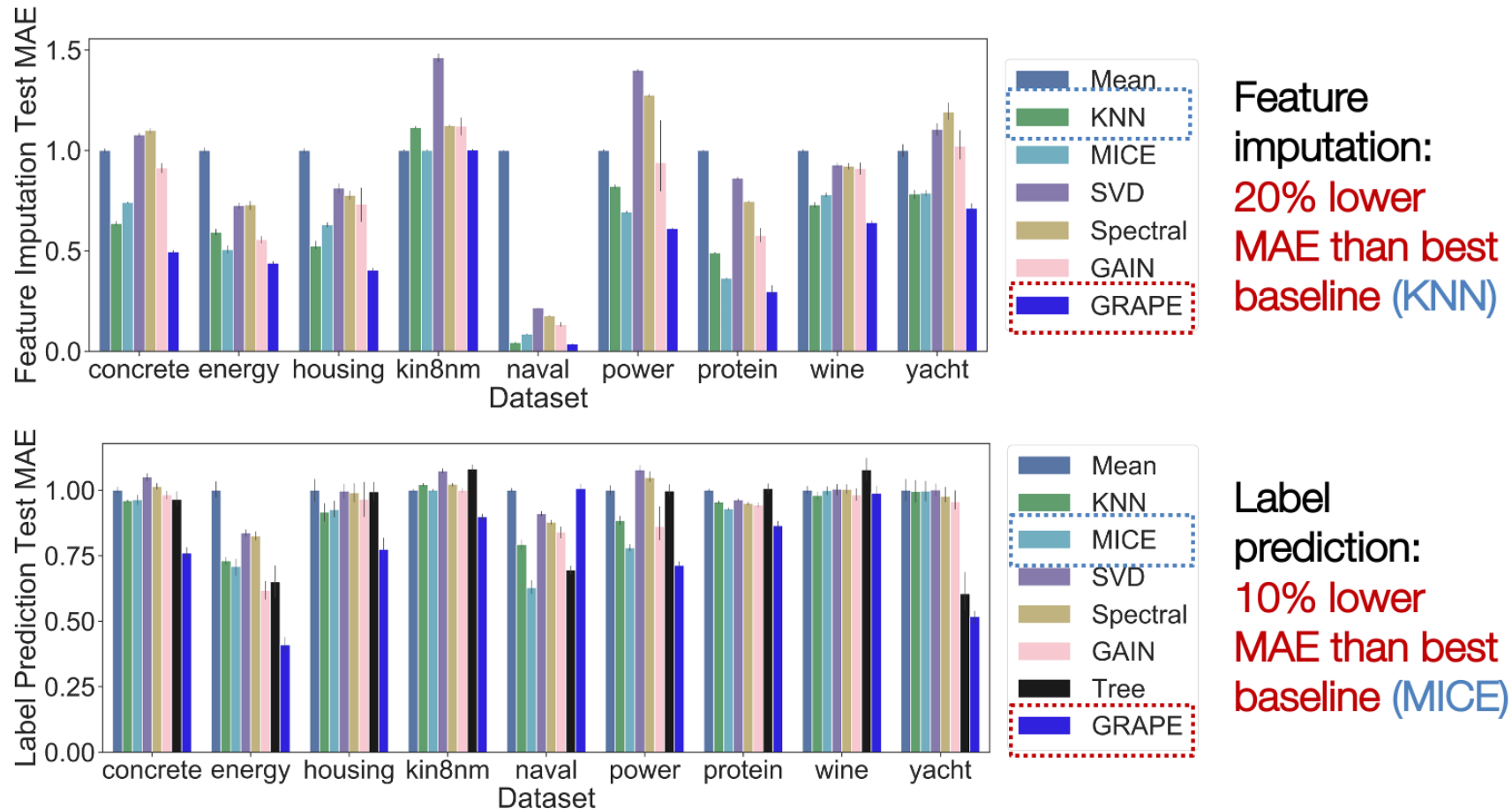
Y
y_1
y_2
?

-  Node Embeddings
-  Edge Embeddings
-  Message Passing
-  Missing Feature Values
-  Downstream Labels



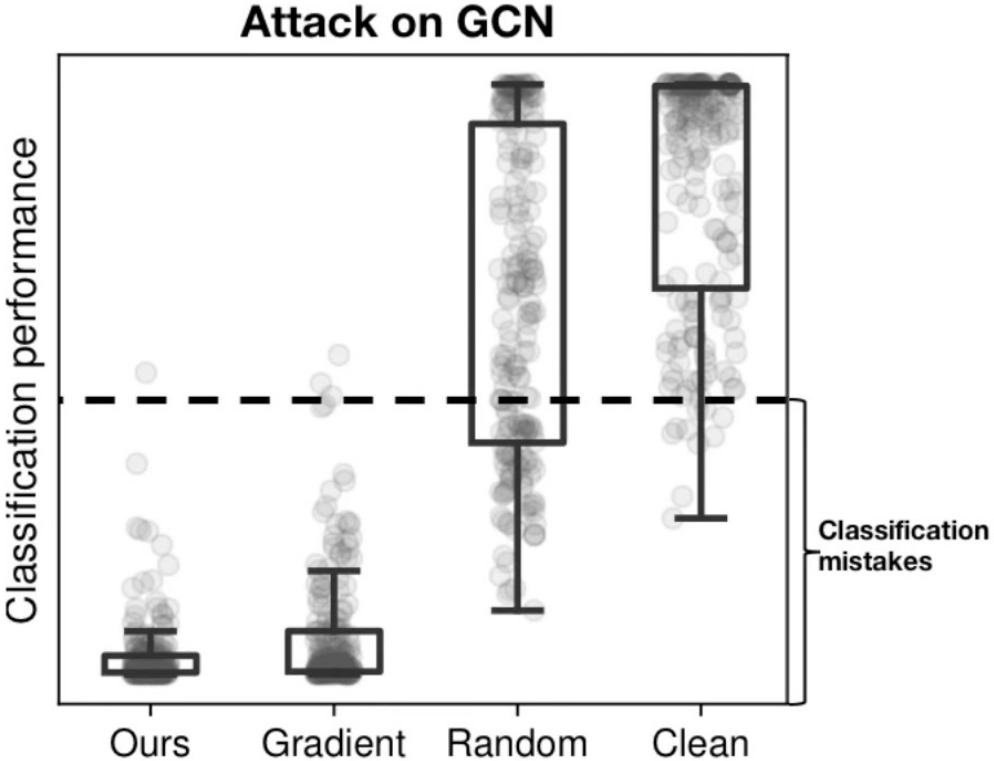
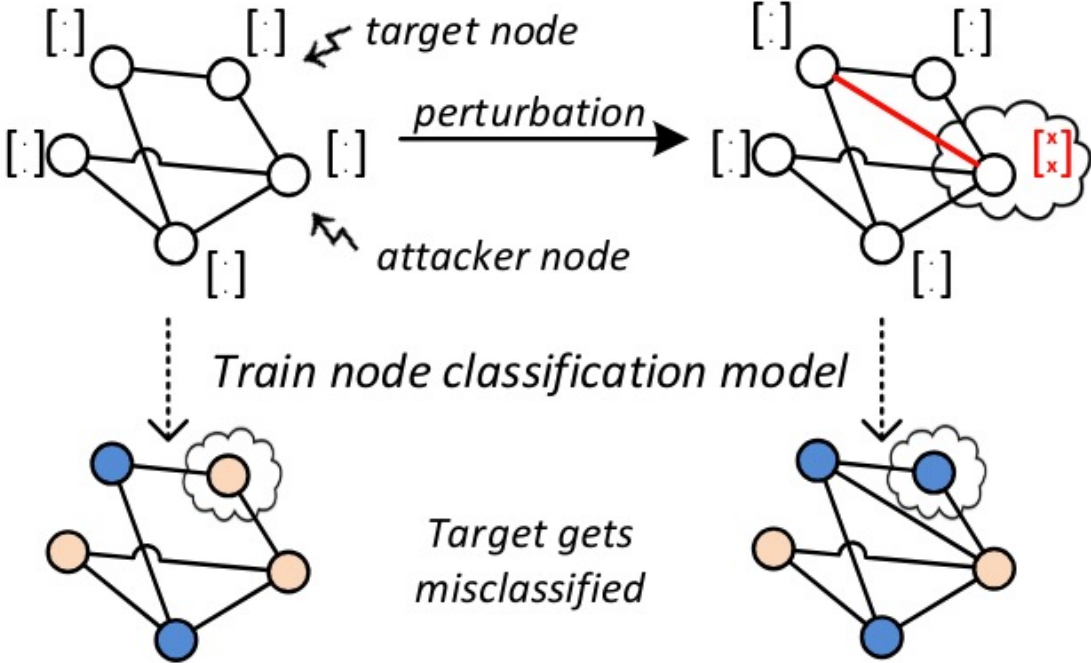
Missing Data

GRAPE yields 20% lower mean absolute error for feature imputation, and 10% lower MAE for label prediction



Adversarial Attacked Data

Observation: Small perturbations of the graph structure and node features lead to misclassification of the target



Adversarial Attacked Data

Can we leverage small data perturbations to **improve performance**?

Yes, adversarial training

Adversarial training is the process of crafting adversarial data points, and then injecting them into training data

$$\min_{\theta} E_{(x,y) \sim \mathcal{D}} \left[\max_{\|\delta\|_p \leq \epsilon} L(f_{\theta}(x + \delta), y) \right]$$



Find the optimal perturbation sample to achieve maximum loss

Find the optimal model parameters to resist the attack of perturbation sample

D: distribution

$\|\cdot\|_p$: l_p -norm distance metric

ϵ : perturbation budget

Adversarial Attacked Data

Can we leverage small data perturbations to **improve performance?**

Yes, adversarial training

Node Classification

Backbone	ogbn-products Test Acc	ogbn-proteins Test ROC-AUC	ogbn-arxiv Test Acc
GCN	-	72.51±0.35	71.74±0.29
+FLAG	-	71.71±0.50	72.04±0.20
GraphSAGE	78.70±0.36	77.68 ±0.20	71.49±0.27
+FLAG	79.36±0.57	76.57±0.75	72.19±0.21
GAT	79.45±0.59	-	73.65±0.11
+FLAG	81.76±0.45	-	73.71±0.13
DeeperGCN	80.98±0.20	85.80±0.17	71.92±0.16
+FLAG	81.93±0.31	85.96±0.27	72.14±0.19

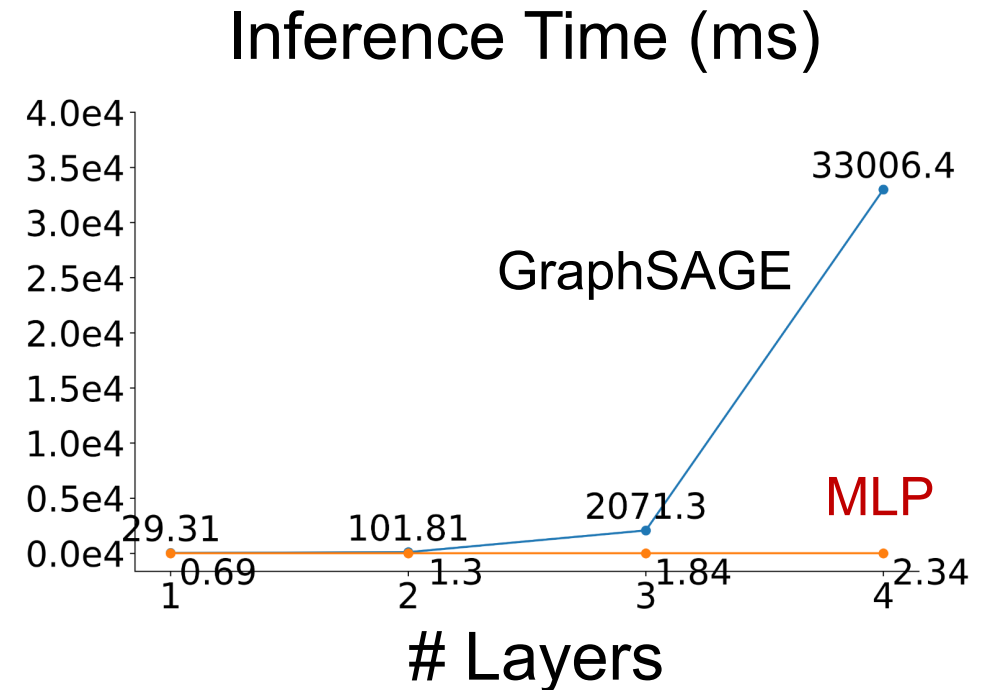
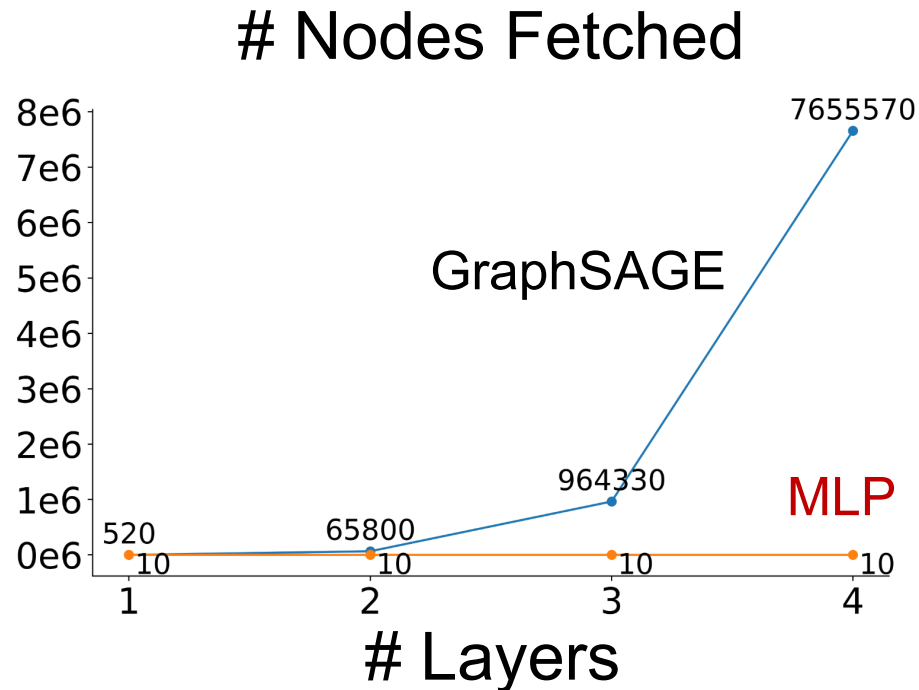
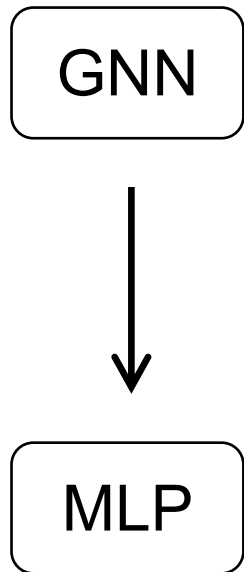
Adversarial Attacked Data

Can we leverage small data perturbations to **improve robustness**?

Yes, adversarial training

A use case: training an MLP on graphs

Reason: to avoid the computation-intensive message passing mechanism

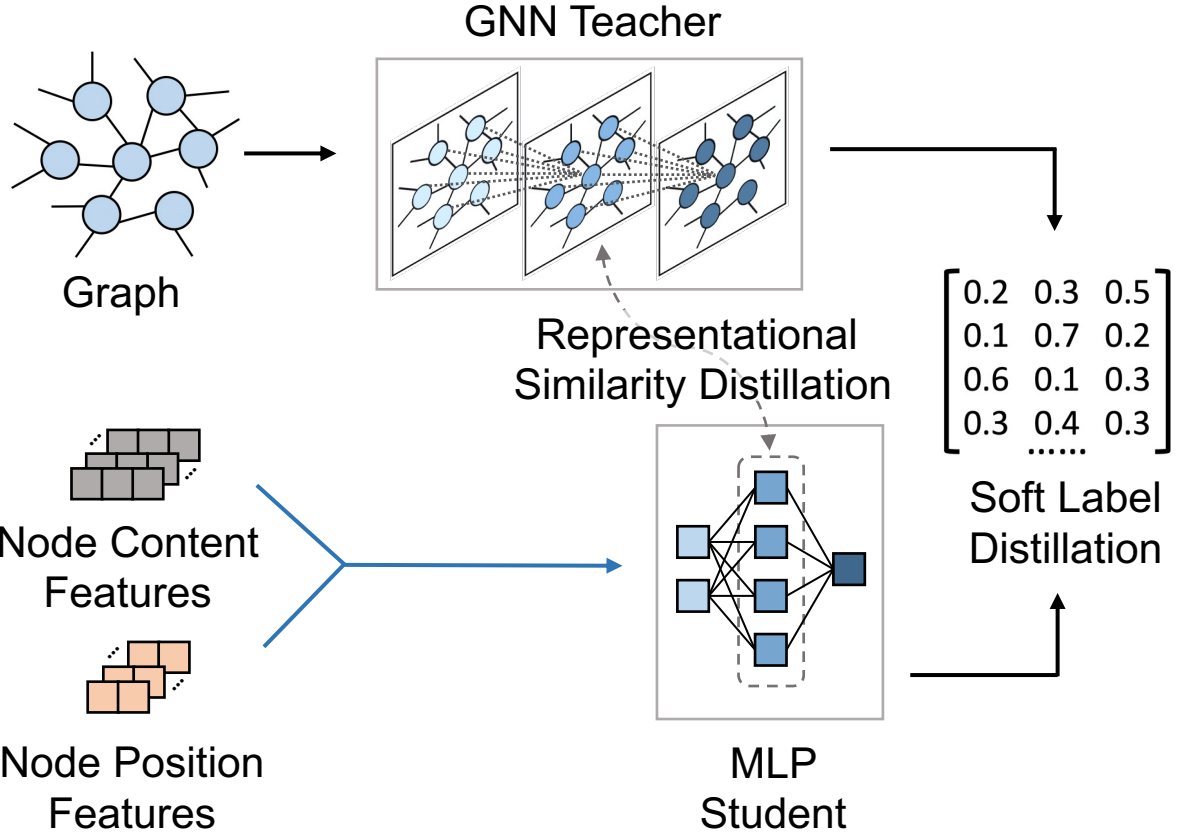


Adversarial Attacked Data

Can we leverage small data perturbations to **improve robustness**?

Yes, adversarial training

A use case: training an MLP on graphs



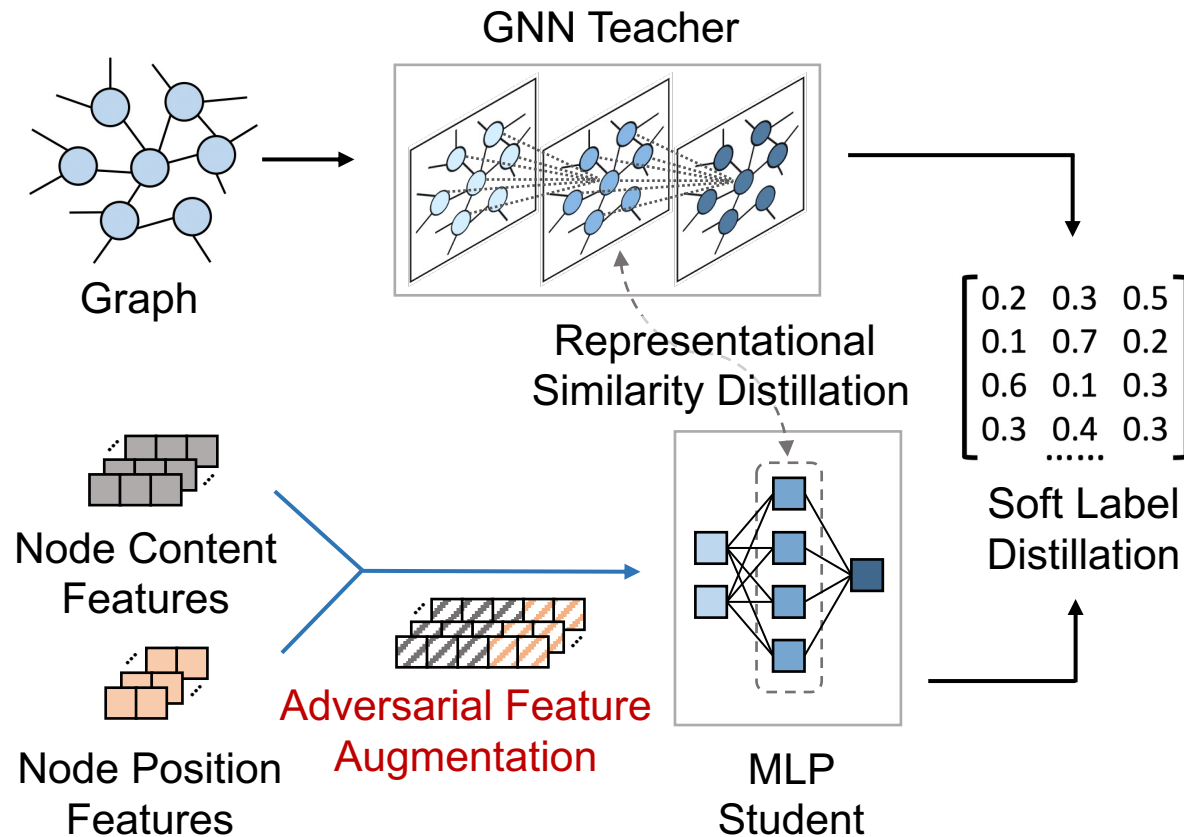
The problem of training an MLP on graphs:
sensitive to features

Adversarial Attacked Data

Can we leverage small data perturbations to **improve robustness**?

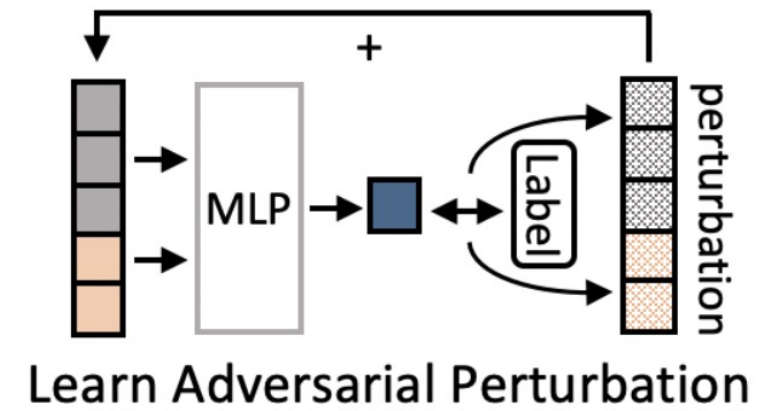
Yes, adversarial training

A use case: training an MLP on graphs



The problem of training an MLP on graphs: sensitive to features

Overcome this problem with adversarial training

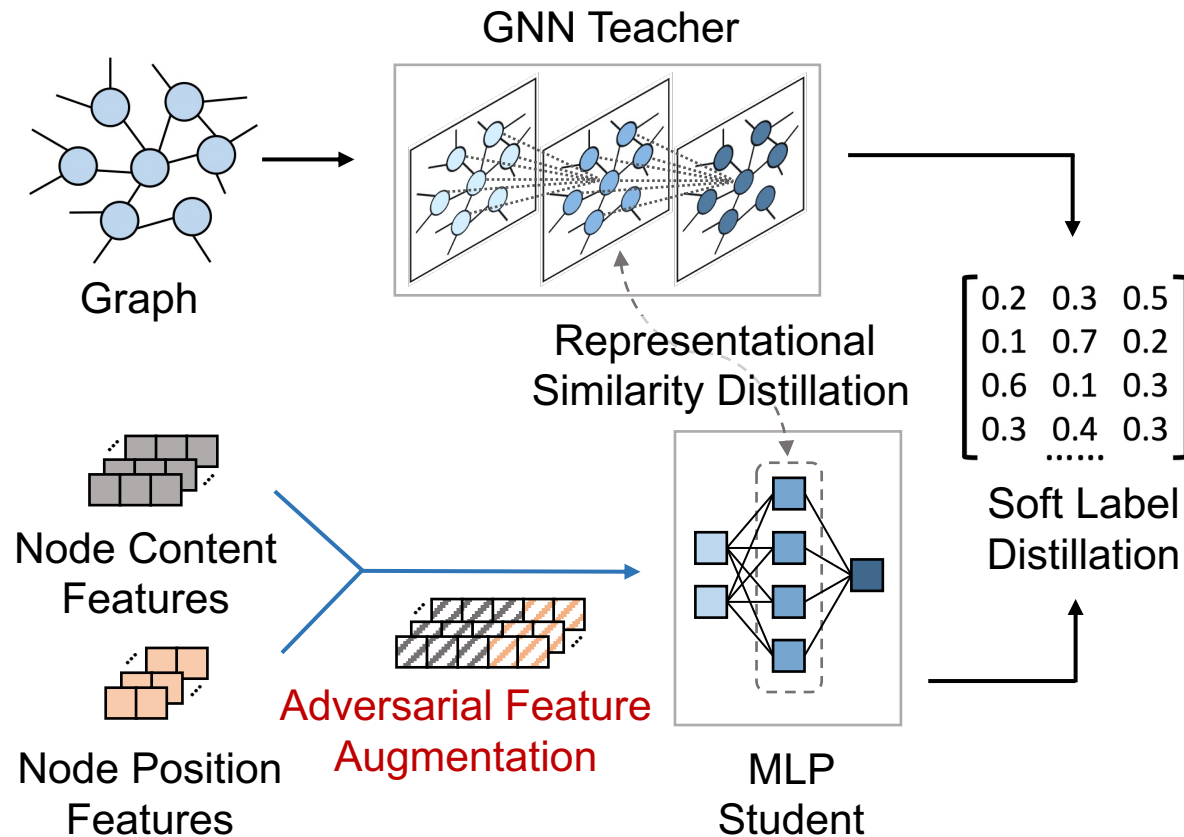


Adversarial Attacked Data

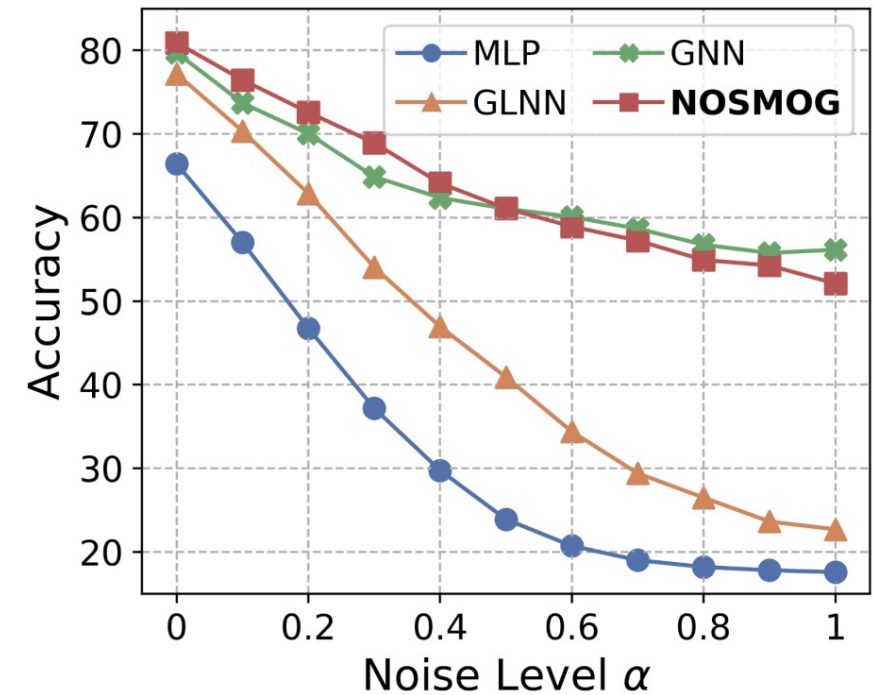
Can we leverage small data perturbations to **improve robustness**?

Yes, adversarial training

A use case: training an MLP on graphs



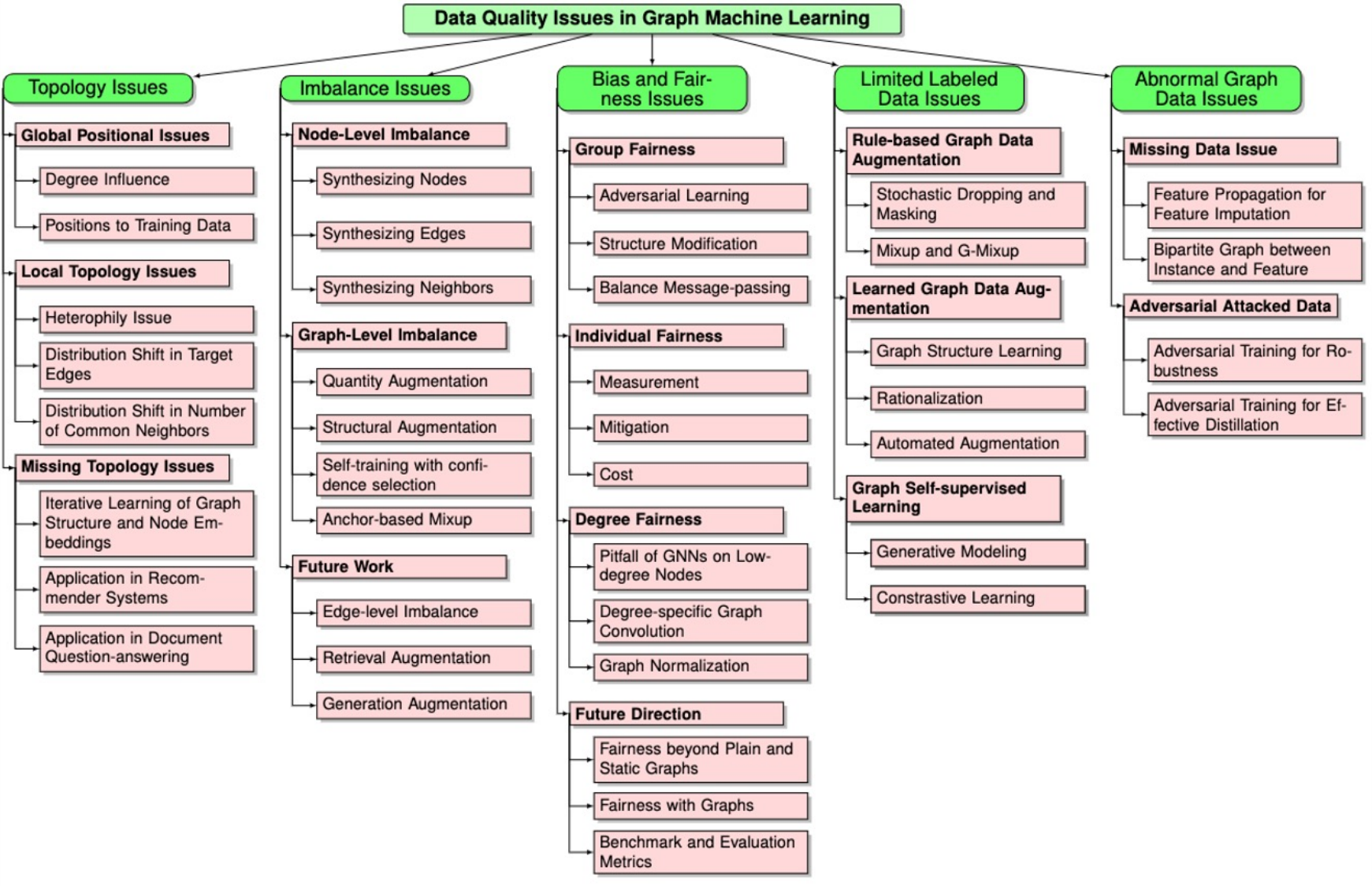
NOSMOG is as robust as GNNs



Summary

- Introduction and Background
- Topology Issues
- Imbalance Issues
- Short Break
- Bias Issue
- Limited Labeled Data Issues
- Abnormal Graph Data Issues
- **Summary**

Summary

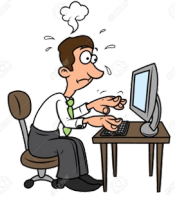


Future Directions

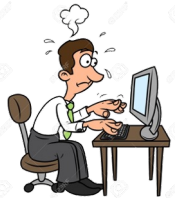
Topology Issue



Imbalance Data Issue



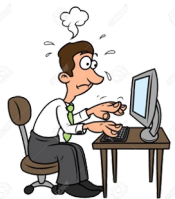
Bias and Fairness Issue



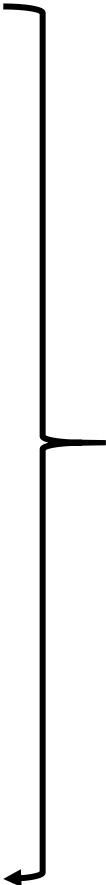
Limited Labeled Data Issue



Abnormal Graph Data Issue



Existing Data Processing is very time-consuming and labor extensive!



Summary

Intelligent Data Processing Tool

