



Graph Sparsification via Meta-Learning

Guihong Wan and Harsha Kokel

Department of Computer Science
The University of Texas at Dallas

Contributions

An edge sparsification algorithm for undirected graphs

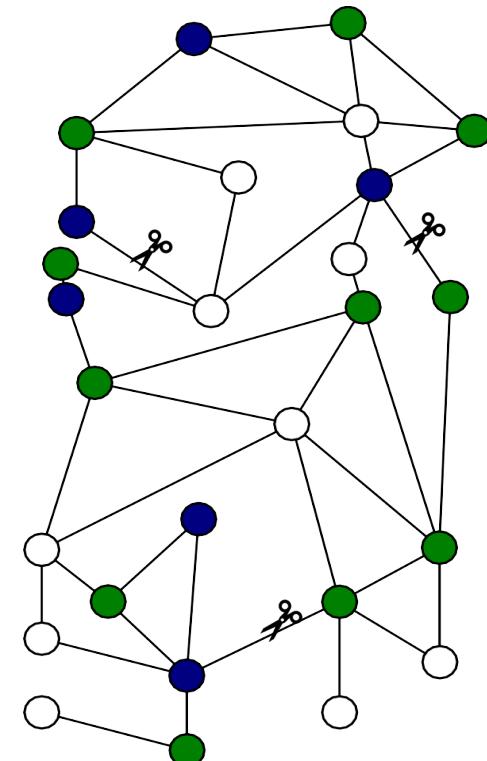
Delete edges \longleftrightarrow Preserve node classification accuracy

- Formulate as a bi-level optimization problem

$$\begin{aligned}\hat{G}^* &= \min_{\hat{G} \in \Phi(G)} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{G}), Y_U) \\ \text{s.t. } \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{G}), Y_L)\end{aligned}$$

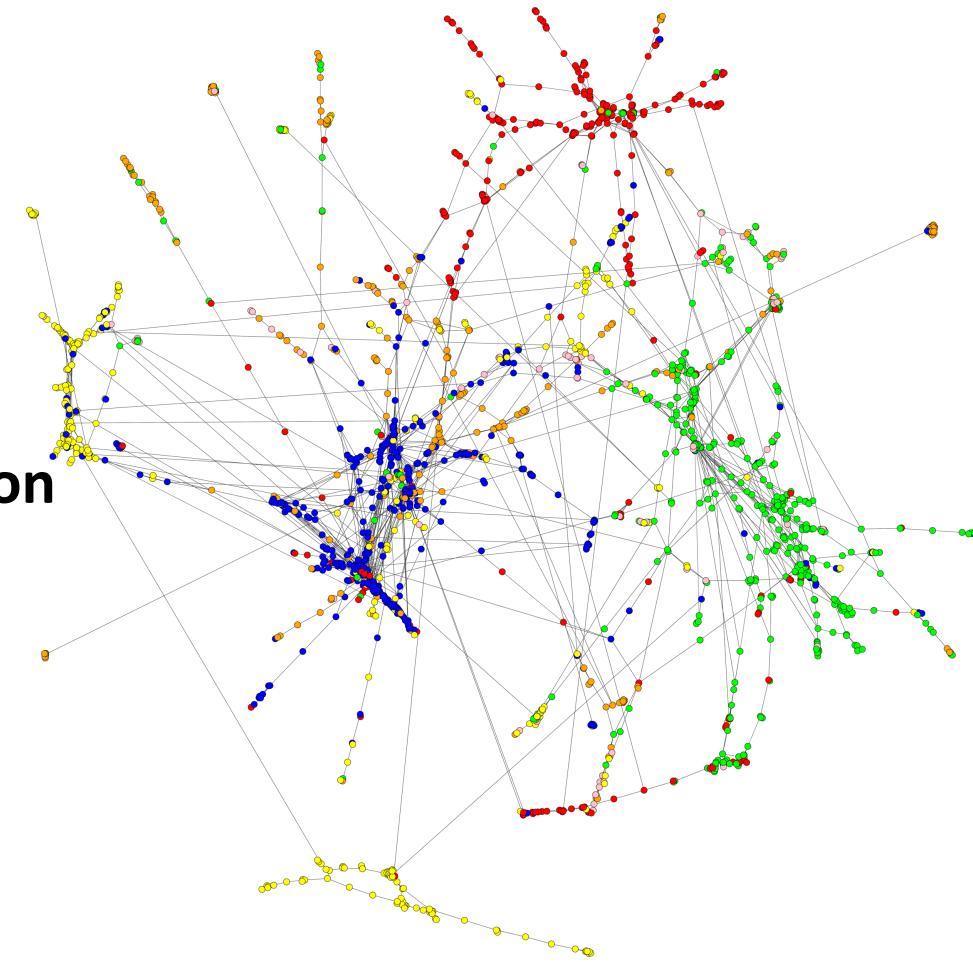
- Use meta-gradients to solve it.

$$\begin{aligned}\nabla_{\hat{A}}^{\text{meta}} &:= \nabla_{\hat{A}} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{A}, X), Y_U), \\ \text{s.t. } \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{A}, X), Y_L)\end{aligned}$$



Outline

- **Introduction**
 - Graph Sparsification
 - Semi-Supervised Node Classification
- **Our Approach**
 - Modeling the Problem
 - Meta-Gradients
 - Score Matrix
 - The Proposed Algorithm
- **Experimental Results**



Introduction

- Graph sparsification
- Semi-supervised node classification

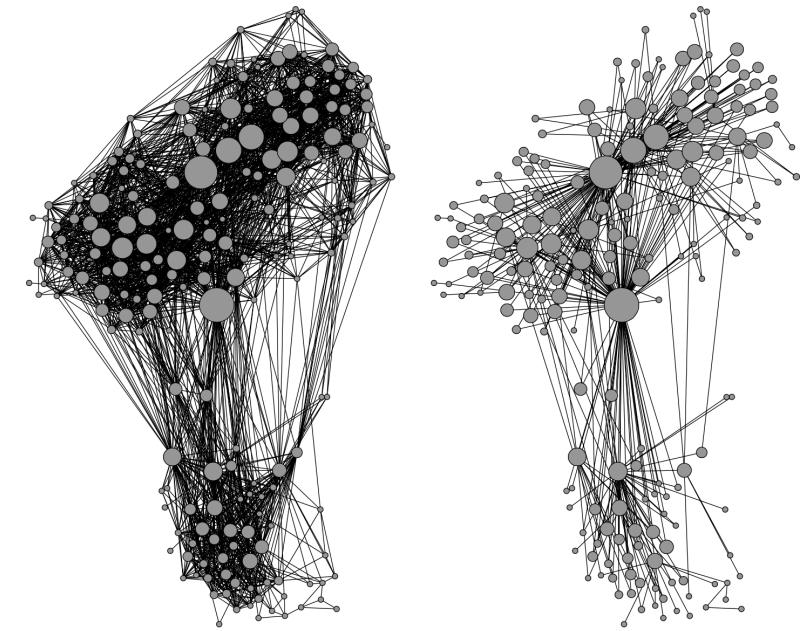
Graph Sparsification

- **Edge Sparsification:**

reduce the edges of a graph while preserving structural / statistical properties of interest.

- **Density** for undirected graphs:

$$\frac{2M}{N(N - 1)} \leftarrow \begin{array}{l} \text{number of edges} \\ \hline \text{number of nodes} \end{array}$$

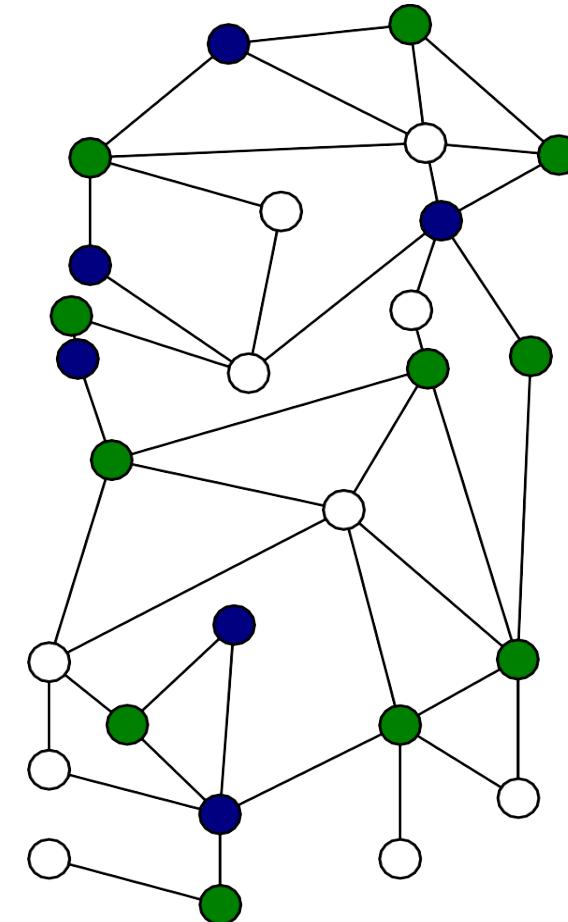


Hamann et al. 2016

We focus on **edge sparsification** while preserving **the node classification accuracy**.

Semi-Supervised Node Classification

- **Graph:** $G = (A, X), Y_L, Y_U \leftarrow$ to be predicted
 - adjacency matrix
 - attribute matrix
 - labels
- **Task:** Node classification
 - Input: (A, X, Y_L)
 - Output: predicted Y_U
- **Example:** Graph Convolution Network (GCN)



GCN

Two-layer GCN:

$$f_{\theta}(A, X) = \text{softmax}(A' \delta(A' X W_1) W_2)$$

The l^{th} GCN Layer: $H_{l+1} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H_l W_l)$

$H_1 = X$, otherwise the output of previous layer.

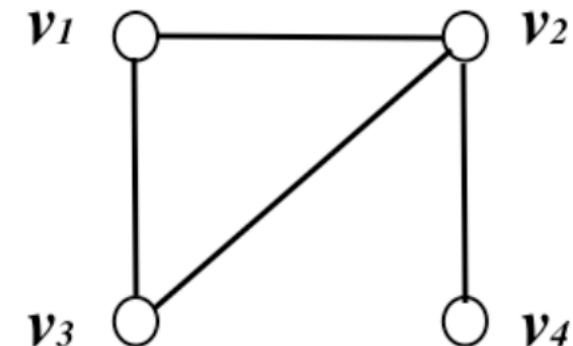
$$\tilde{A} = A + I_N, \quad \tilde{D} = \sum_j \tilde{A}_{ij}$$

W : what to learn.

GCN (cont.)

Understand intuitively:

$$\mathbf{H}_2 = \sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}_1) \text{ when no normalization.}$$



$$\begin{array}{c} \tilde{\mathbf{A}} \\ \mathbf{X} \\ \hline \end{array} \quad \begin{array}{c} \mathbf{x} \\ \mathbf{X} \\ \hline \end{array} \quad \begin{array}{c} \mathbf{AX} \\ \hline \end{array}$$

Matrix $\tilde{\mathbf{A}}$ (adjacency matrix with self-loops):

$$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

Matrix \mathbf{X} (feature matrix):

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}$$

Matrix \mathbf{AX} (product of $\tilde{\mathbf{A}}$ and \mathbf{X}):

$$\begin{bmatrix} x_{11} + x_{21} + x_{31} & x_{12} + x_{22} + x_{32} & \cdots \end{bmatrix}$$

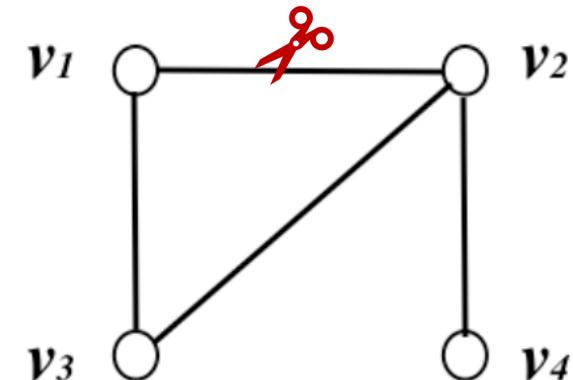
Arrows indicate the summation of features from node 1's neighbors (2, 3, 4) for each dimension of the feature vector x_1 , x_2 , x_3 , and x_4 .

Main idea: learn a node v 's representation by aggregating its own feature x_v and its neighbors' feature x_u , for all $u \in N(v)$.

GCN (cont.)

Understand intuitively:

$$\mathbf{H}_2 = \sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}_1) \text{ when no normalization.}$$



$$\begin{array}{c} \tilde{\mathbf{A}} \\ \mathbf{X} \\ \hline \end{array} \quad \begin{array}{c} \mathbf{x} \\ \mathbf{AX} \\ \hline \end{array}$$

$\tilde{\mathbf{A}}$	\mathbf{X}	\mathbf{ax}
$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ x_{41} & x_{42} & x_{43} & x_{44} \end{bmatrix}$	$\begin{bmatrix} x_{11} + x_{21} + x_{31} & x_{12} + x_{22} + x_{32} & \dots \end{bmatrix}$

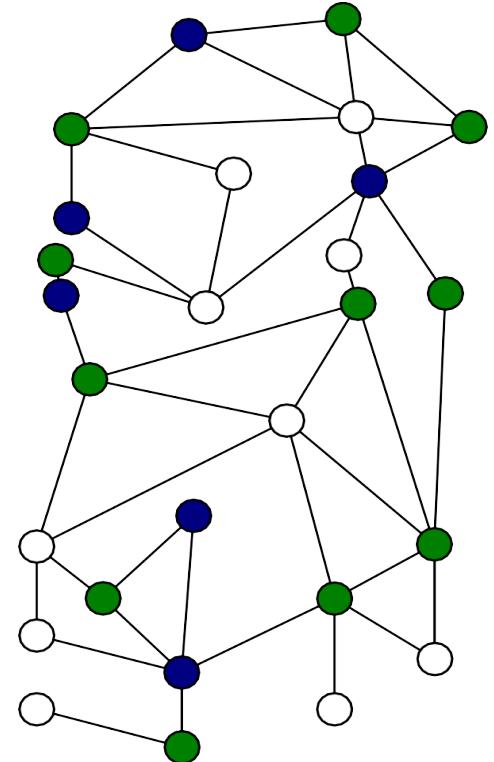
Main idea: learn a node v 's representation
by aggregating its own feature x_v and its neighbors' feature x_u ,
for all $u \in N(v)$.

Our Approach

- Modeling the problem
- Meta-gradients
- Score matrix

Modeling the Problem

- **Given:** $G = (A, X)$, labeled nodes: Y_L .
- **Goal of semi-supervised node classification:**
learn a function f_θ to map each node to a class.



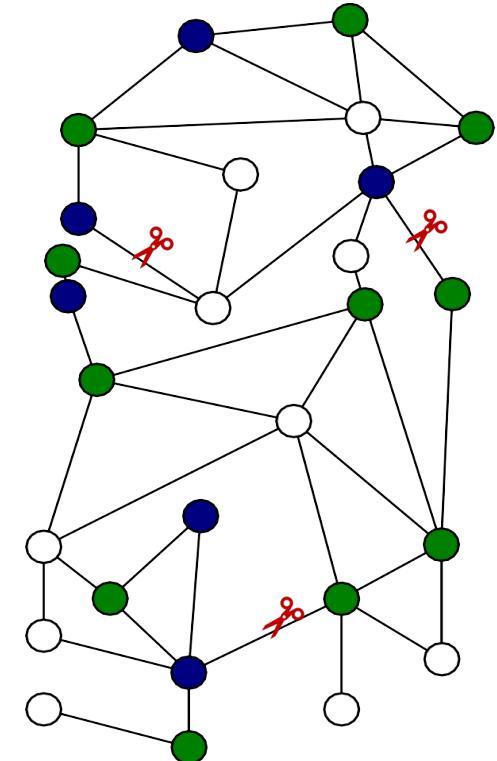
$$\theta^* = \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(G), Y_L)$$

Modeling the problem (cont.)

- **Given:** $G = (A, X)$,
labeled nodes: Y_L ,
number of edges to be deleted: ζ .
- **Goal:** delete edges
but reduce the loss of node classification
accuracy on unlabeled nodes:

$$L_{sps}(\tilde{Y}_U, Y_U)$$

Predicted labels True labels of unlabeled nodes



Modeling the problem (cont.)

- Formulate as a bi-level optimization problem

$$\begin{aligned}\hat{G}^* &= \min_{\hat{G} \in \Phi(G)} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{G}), Y_U) && \text{Outer} \\ s.t. \quad \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{G}), Y_L) && \text{Inner}\end{aligned}$$

- **Inner optimization**: train the model over labeled nodes for predicing labels of unlabeled nodes.
- **Outer objective**: for the sparsifier which aims to minimize the loss of classification accuracy.

Y_U is unknown to the sparsifier.

Modeling the problem

$$\begin{aligned}\hat{G}^* &= \min_{\hat{G} \in \Phi(G)} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{G}), Y_U) \quad \text{unknown to the sparsifier.} \\ s.t. \quad \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{G}), Y_L)\end{aligned}$$

Three options to approximate:

- $L_{\text{sps}} \approx L_{\text{train}}$: compute from Y_L .
- $L_{\text{sps}} \approx L_{\text{self}}$: the sparsifier can train a classifier on labeled data to estimate the labels of unlabeled nodes \hat{Y}_U .
- $L_{\text{sps}} \approx L_{\text{both}}$: combine Y_L and \hat{Y}_U .

Meta-Gradients

- Adjacency matrix → hyperparameters
- Compute the gradients of the sparsifier's loss w.r.t the adjacency matrix

$$\nabla_{\hat{A}}^{\text{meta}} := \nabla_{\hat{A}} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{A}, X), Y_U)$$

$$s.t. \quad \theta^* = \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{A}, X), Y_L)$$

- Indicate how the sparsifier's loss L_{sps} will change after training on the simplified graph.

Meta-Gradients (cont.)

- Adjacency matrix \rightarrow hyperparameters
- Compute the gradients of the sparsifier's loss w.r.t the adjacency matrix

$$\nabla_{\hat{A}}^{\text{meta}} := \nabla_{\hat{A}} \mathcal{L}_{\text{sps}}(f_{\theta^*}(\hat{A}, X), Y_U)$$
$$s.t. \quad \theta^* = \arg \min_{\theta} \mathcal{L}_{\text{train}}(f_{\theta}(\hat{A}, X), Y_L)$$

- Inner update: $\theta_{t+1} = \theta_t - \alpha \nabla_{\theta_t} \mathcal{L}_{\text{train}}(f_{\theta_t}(\hat{G}), Y_L)$
- Outer update: $\hat{A}^{k+1} = \hat{A}^k - \beta \nabla_{\hat{A}^k}^{\text{meta}}$, with $\hat{A}^0 = A$

Score Matrix

$$\hat{A}^{k+1} = \hat{A}^k - \beta \nabla_{\hat{A}^k}^{\text{meta}}, \text{ with } \hat{A}^0 = A$$

- 0/1 problem: an edge is either deleted or kept (A is discrete)
- Score matrix:

$$S = \nabla_{\hat{A}}^{\text{meta}} \odot \hat{A}$$

$$e^* = \arg \max_{\substack{e(i,j) \in \hat{A} \\ e(i,j) \in \Phi(G)}} S(i,j)$$

Score Matrix (cont.)

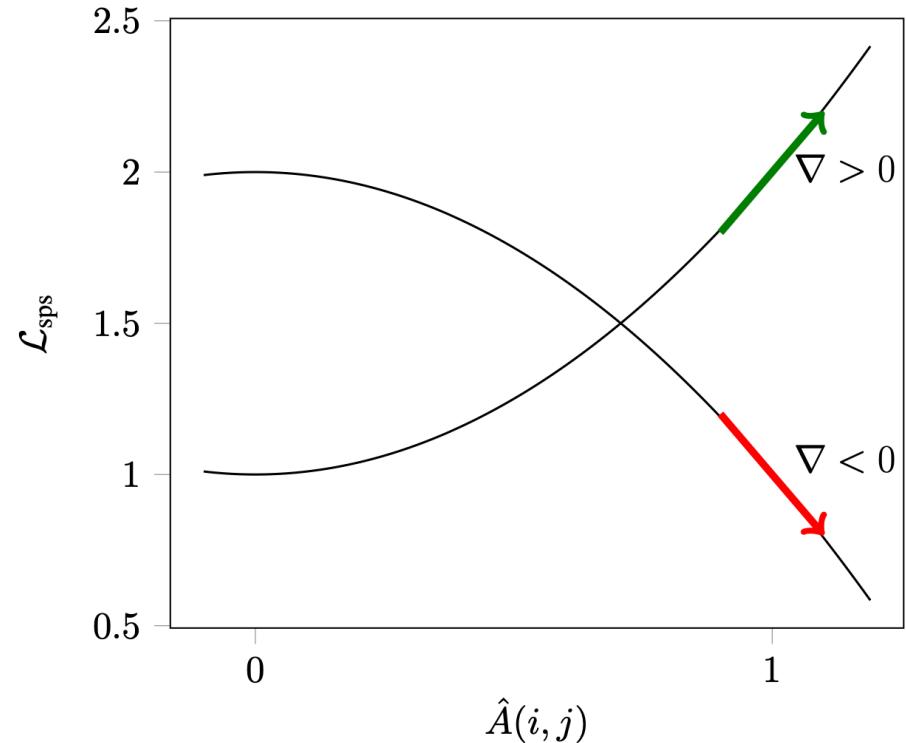
$$\hat{A}^{k+1} = \hat{A}^k - \beta \nabla_{\hat{A}^k}^{\text{meta}}, \text{ with } \hat{A}^0 = A$$

$$S = \nabla_{\hat{A}}^{\text{meta}} \odot \hat{A}$$

- Deletion: From 1 to 0;
- Positive gradients are preferred.

For weighted graphs, we can learn an indicator matrix initialized as 1 if there is an edge.

Figure 1. Illustration of the score matrix.



Algorithm

Algorithm 1 Graph sparsification via meta-gradients

Input: Graph $G = (A, X)$; labels Y_L ; number of edges to delete ζ ; number of training steps T ; learning rate α .

Output: $\hat{G}^* = (\hat{A}^*, X)$

```
1:  $\hat{Y}_U \leftarrow$  estimated labels of unlabeled nodes using self-training;  
2:  $\hat{A} \leftarrow A$ ;  
3: while  $\zeta > 0$  do  
4:    $\theta_0 \leftarrow$  initialize randomly;  
5:   for  $t$  in  $0 \dots T - 1$  do  
6:      $\theta_{t+1} = \theta_t - \alpha \nabla_{\theta_t} \mathcal{L}_{\text{train}}(f_{\theta_t}(\hat{A}, X), Y_L)$ ;  
7:   end for  
8:    $\nabla_{\hat{A}}^{\text{meta}} \leftarrow \nabla_{\hat{A}} \mathcal{L}_{\text{self}}(f_{\theta_T}(\hat{A}, X), \hat{Y}_U)$ ;  
9:    $S = \nabla_{\hat{A}}^{\text{meta}} \odot \hat{A}$ ;  
10:   $e^* \leftarrow$  the maximum entry  $(i, j)$  in  $S(i, j)$  that satisfies the constraints  $\Phi(G)$ ;  
11:   $\hat{A} \leftarrow$  remove edge  $e^*$ ;  
12:   $\zeta -= 1$ ;  
13: end while  
14:  $\hat{G}^* \leftarrow (\hat{A}, X)$ ;  
15: return  $\hat{G}^*$ .
```

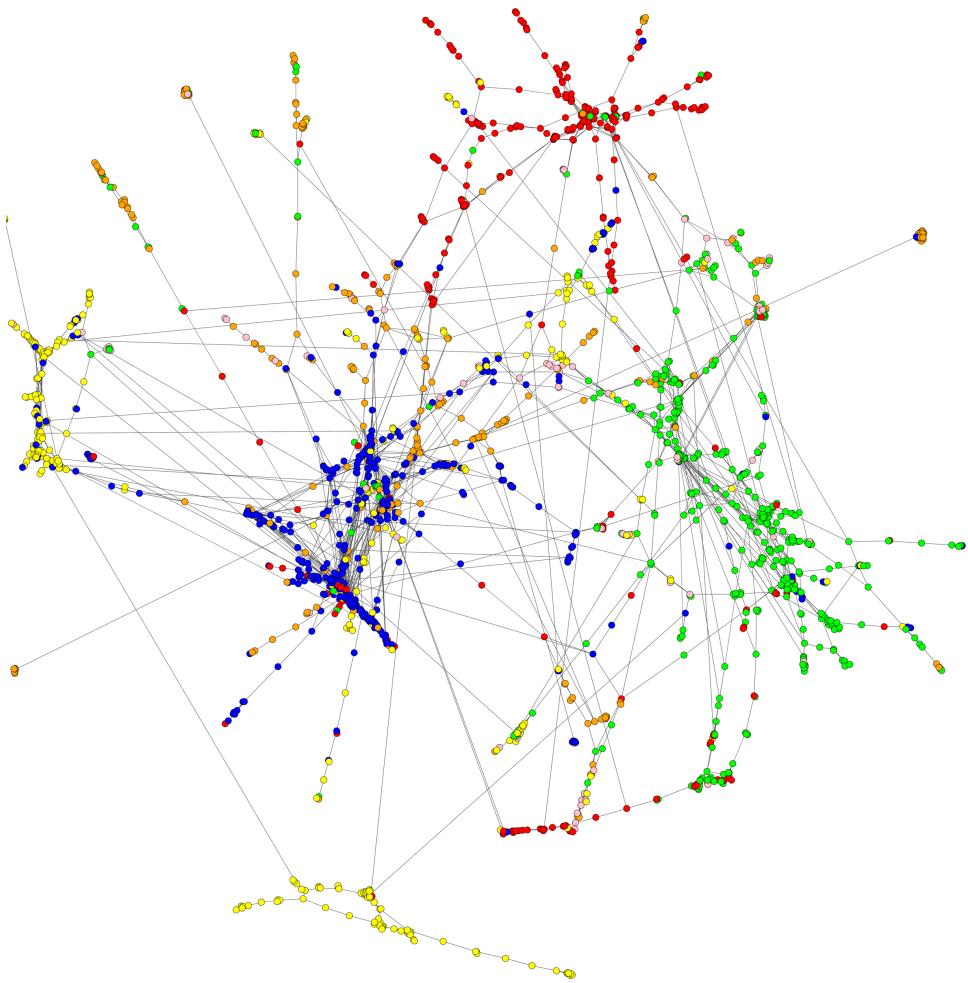
Experimental Results

- **Results on CiteSeer dataset**
- **Results on Cora-ML dataset**

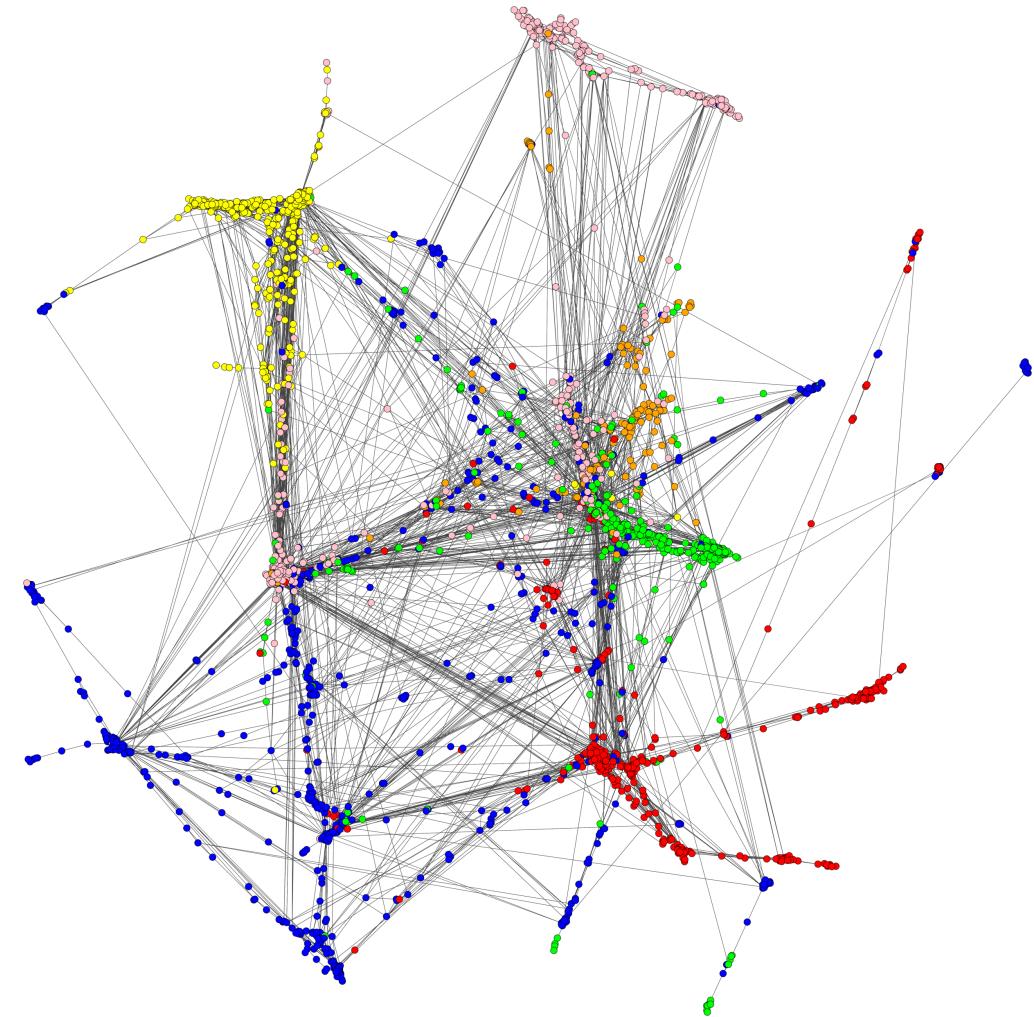
Datasets

	#Nodes	#Edges	Density	Avg degree	Max degree	Test Acc
CiteSeer	2,110	3,668	0.2	5.22	198	0.71
Cora-ML	2,810	7,981	0.4	11.36	492	0.85

- We only consider the largest connected component.
- 10% labeled nodes for training;
90% unlabeled nodes for testing.



CiteSeer
Density: 0.2



Cora-ML
Density: 0.4

Main Observations

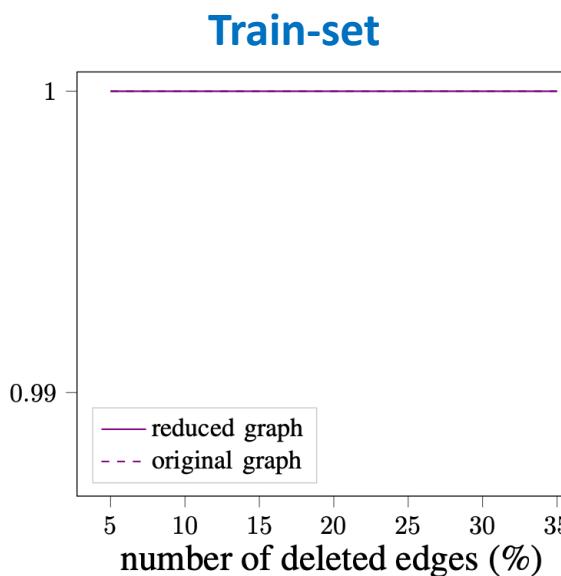
- Our algorithm works better than the conventional methods.
- L_{train} works better when overfitting;
- L_{self} works better when underfitting.

CiteSeer

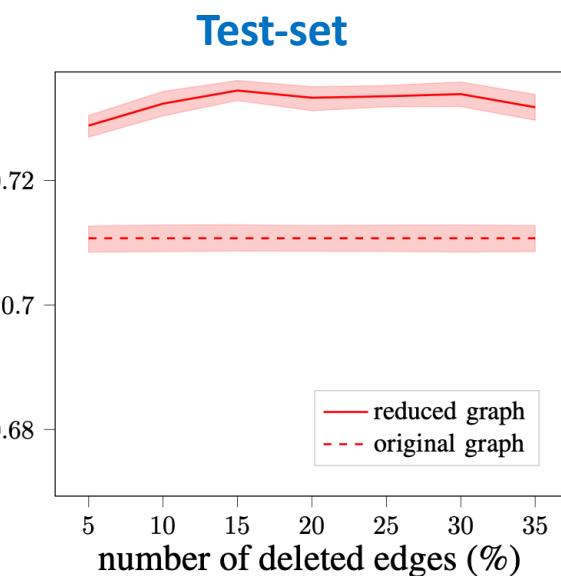
L_{train}



accuracy



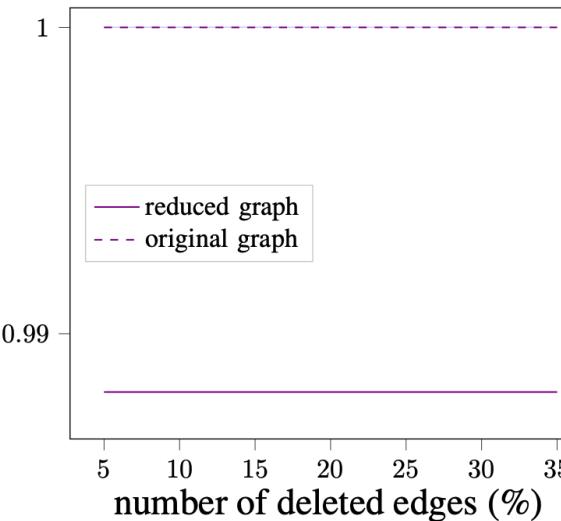
accuracy



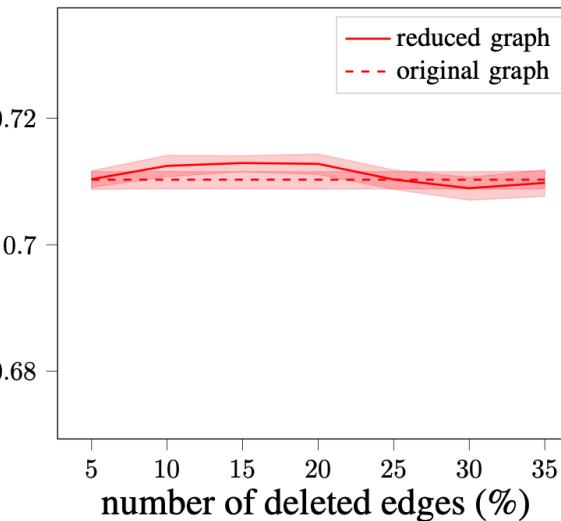
L_{self}



accuracy

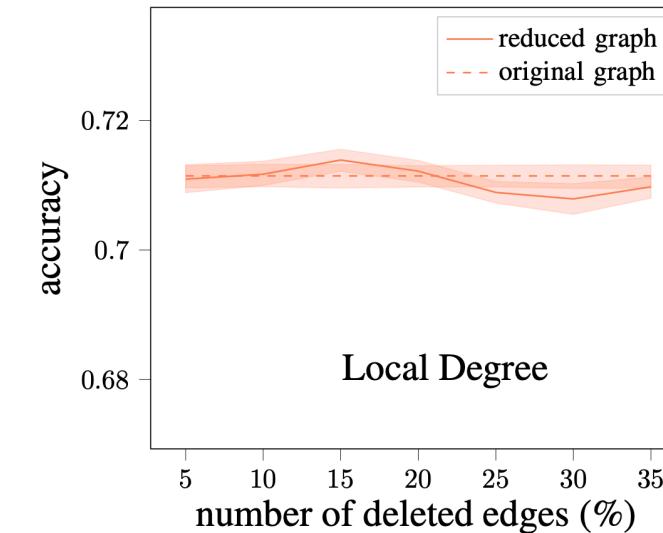
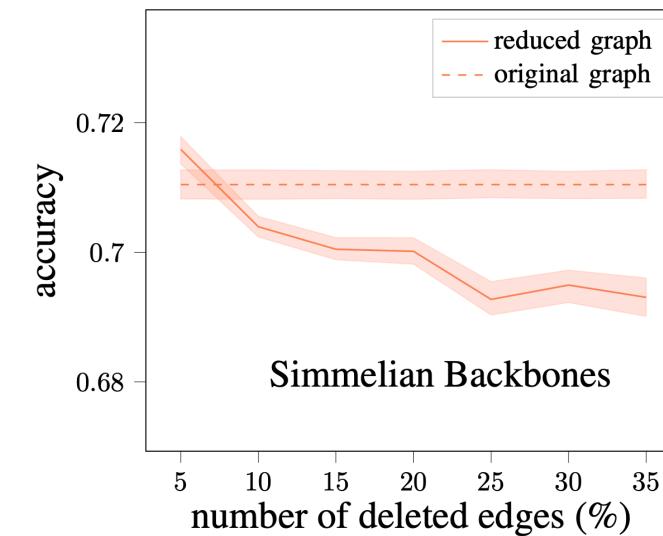
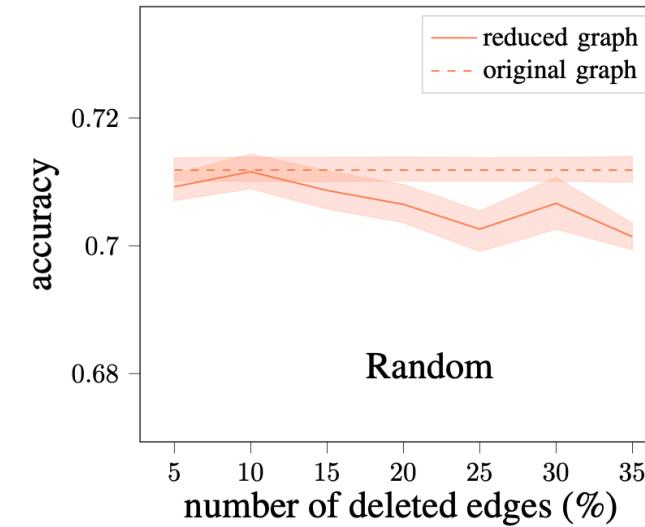
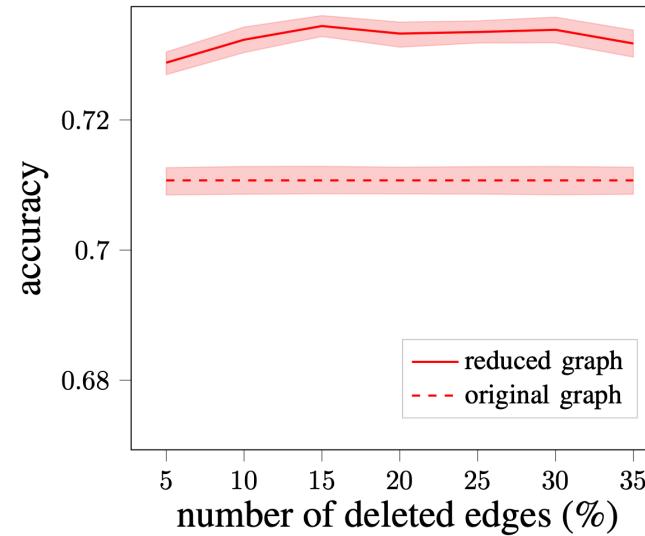


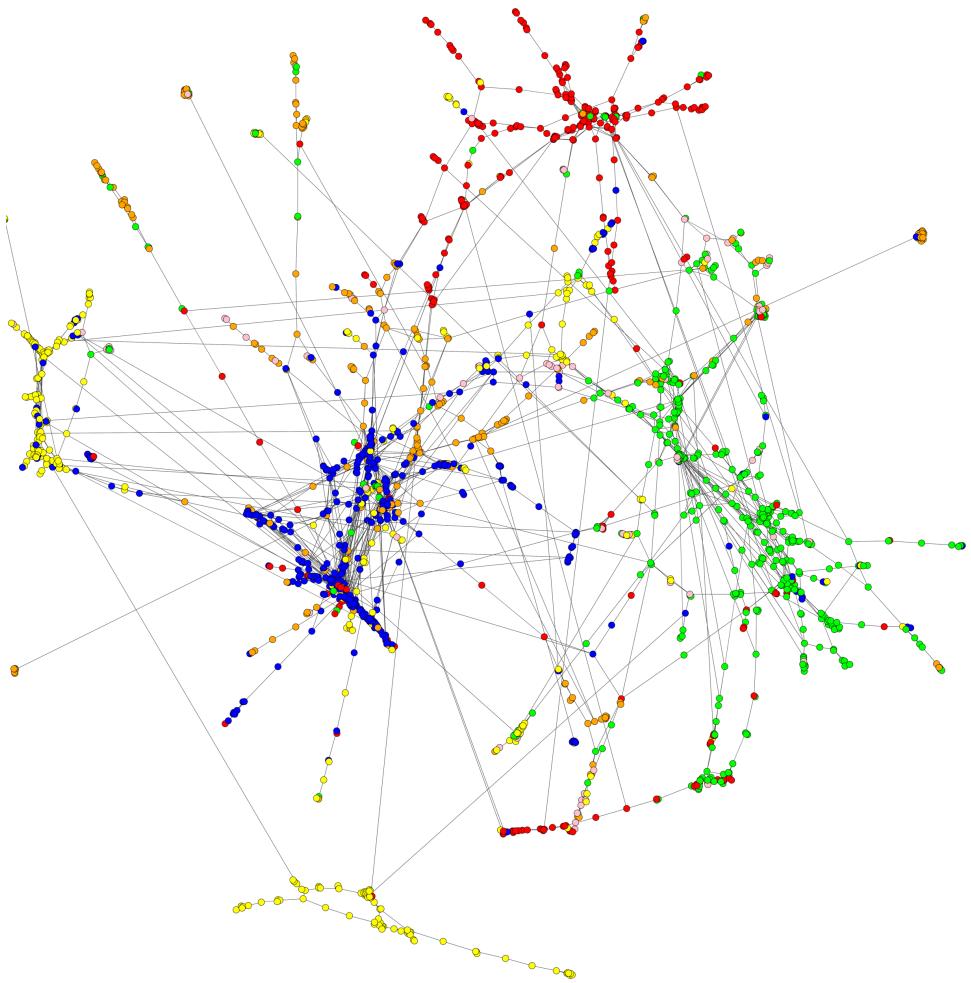
accuracy



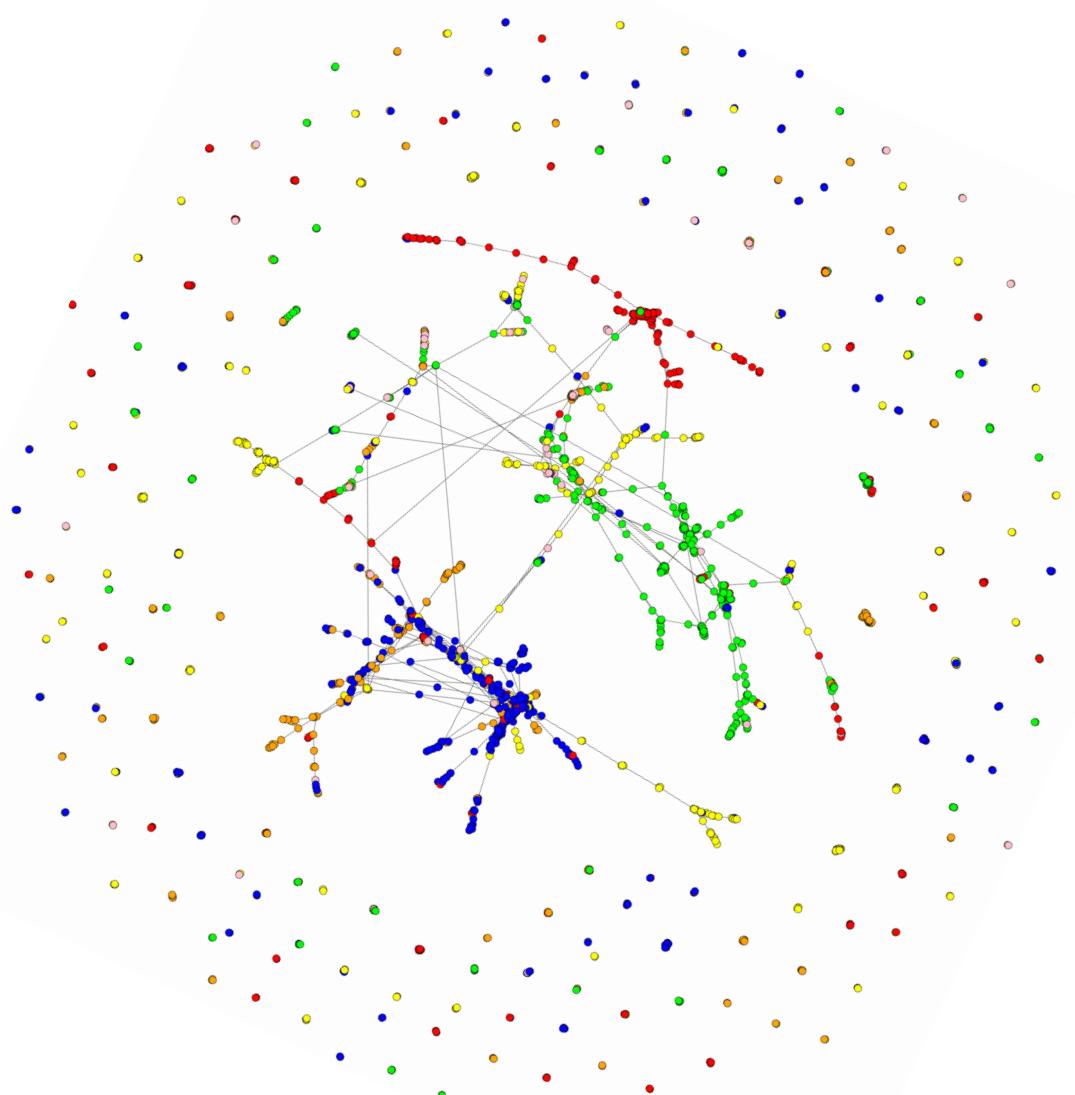
Comparison

Ours

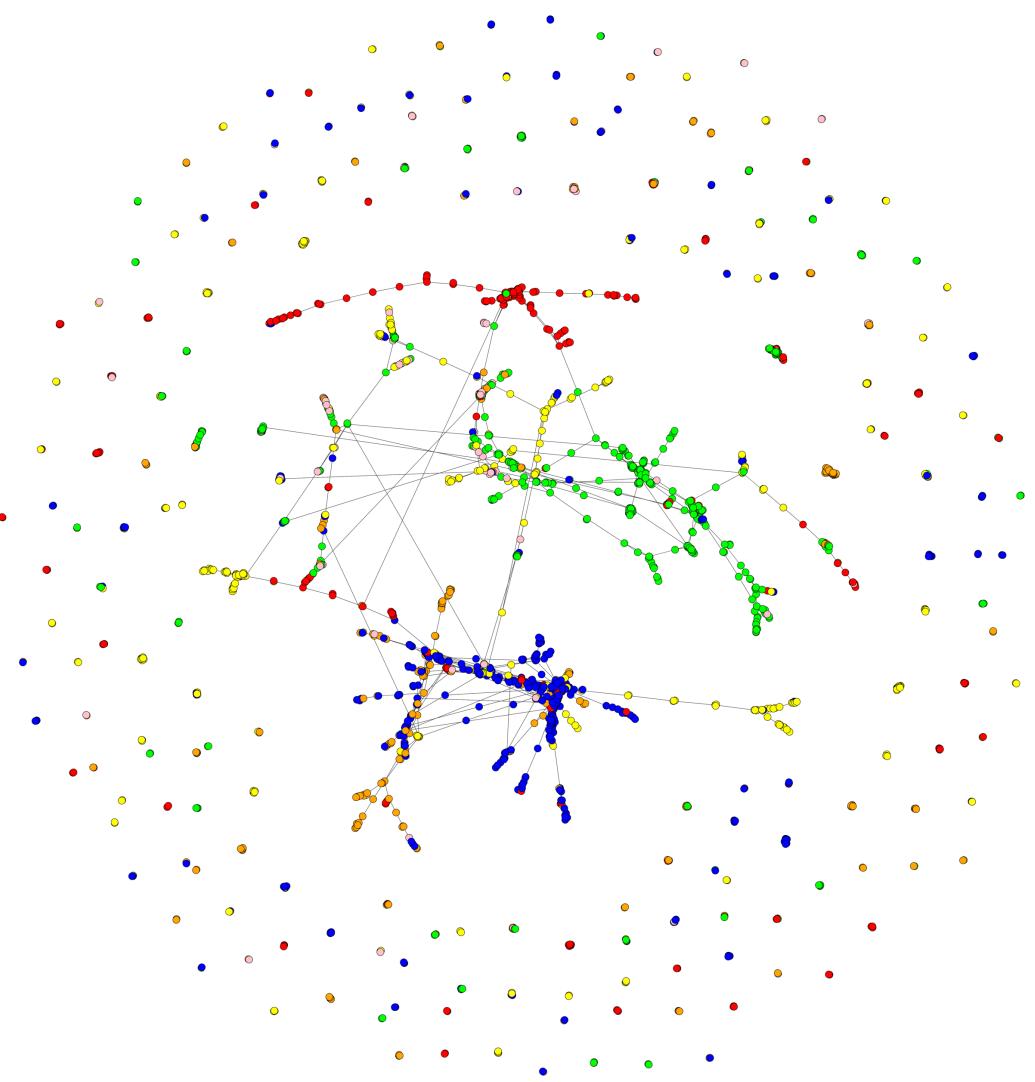




Original Graph
Test Acc : 0.71



30% edges deleted
Test Acc : 0.73

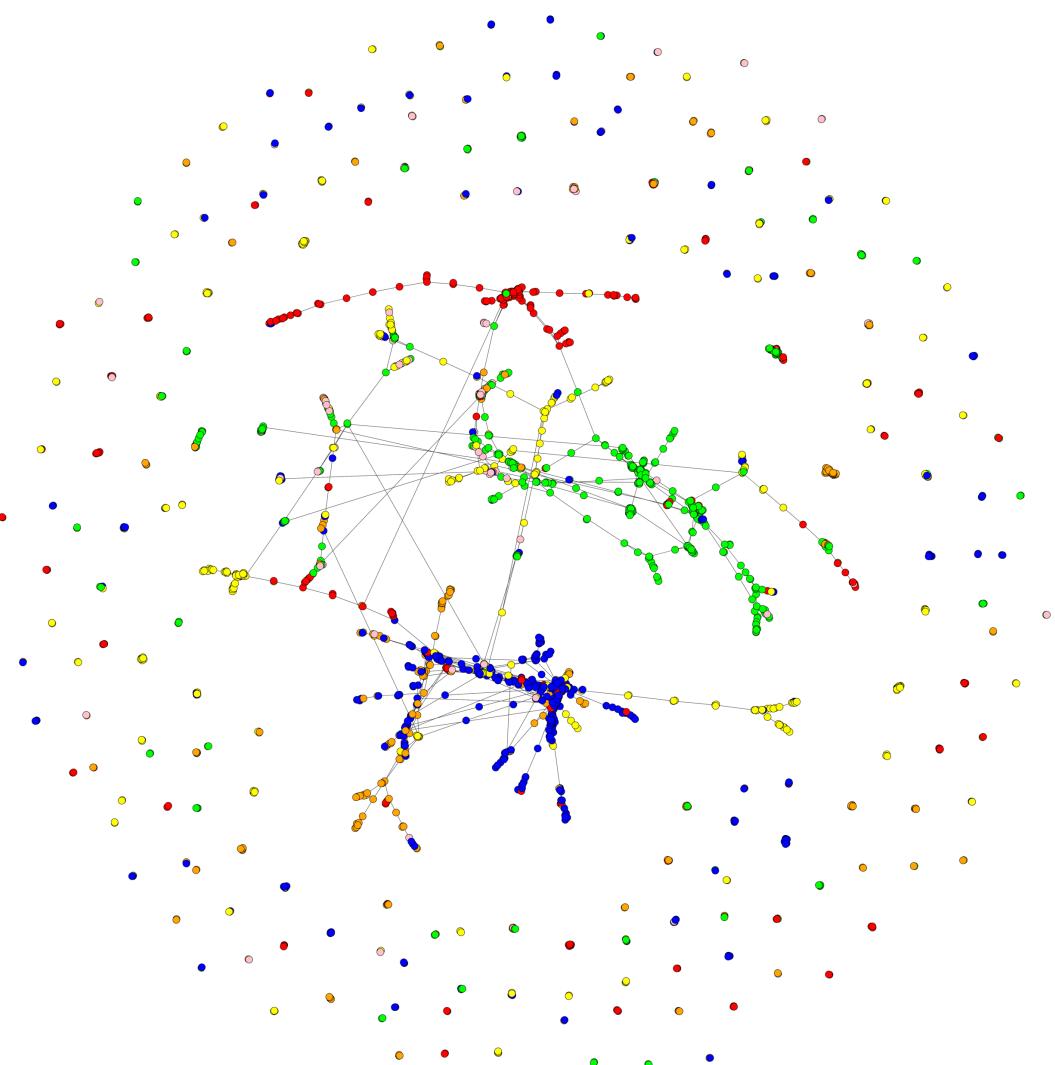


Ours

30% edges deleted
Test Acc : 0.73

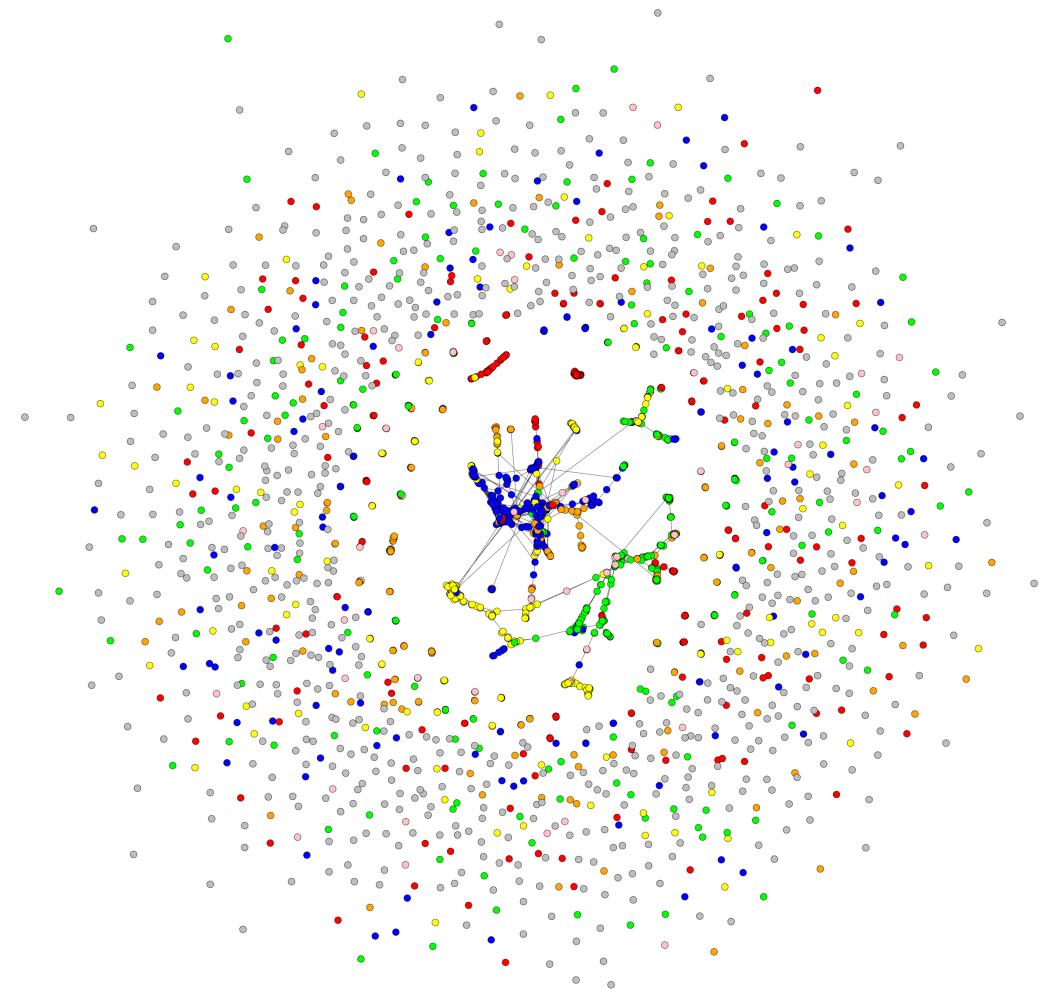
Local Degree

30% edges deleted
Test Acc : 0.71



Ours

30% edges deleted
Test Acc : 0.73



Simmelian

30% edges deleted
Test Acc : 0.69

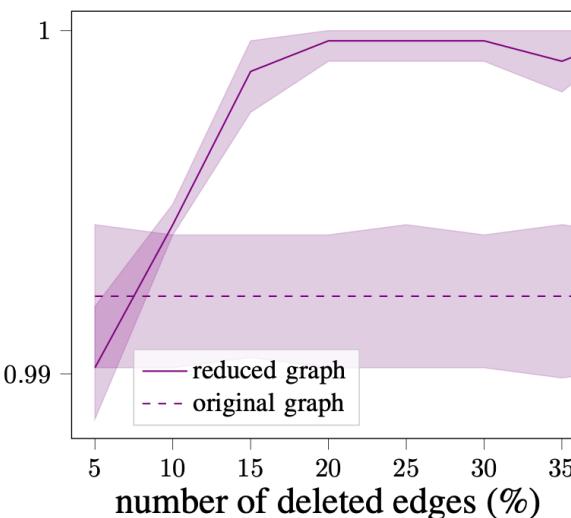
Cora-ML

L_{train}

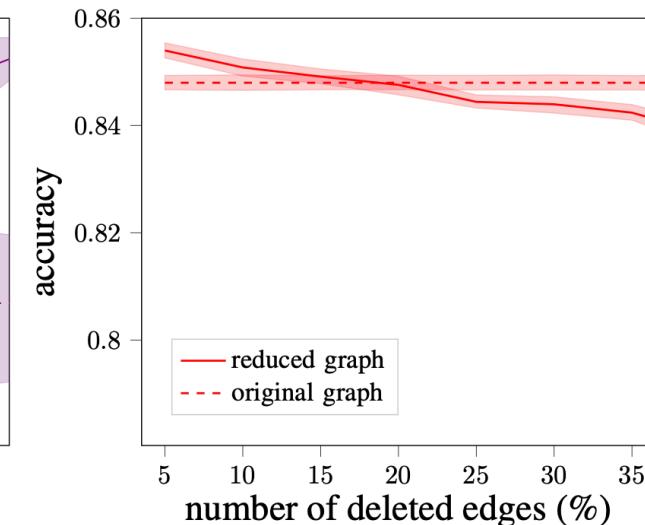


accuracy

Train-set



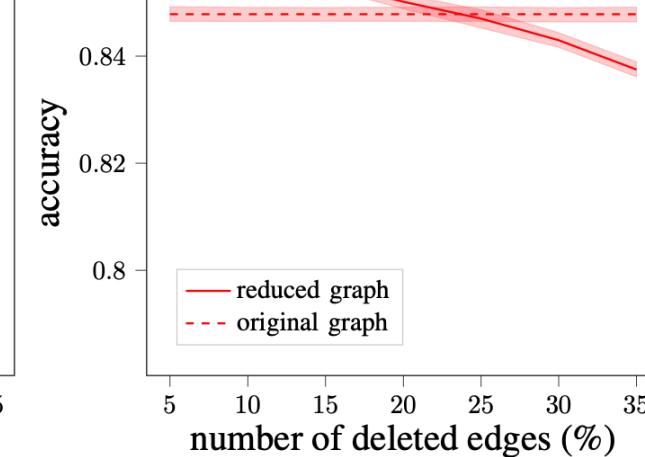
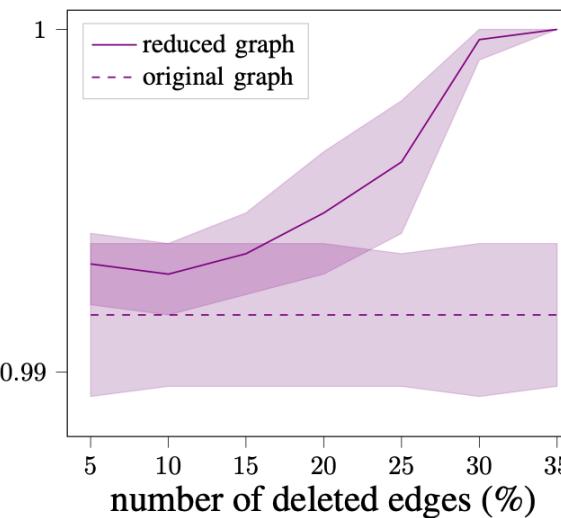
Test-set



L_{self}

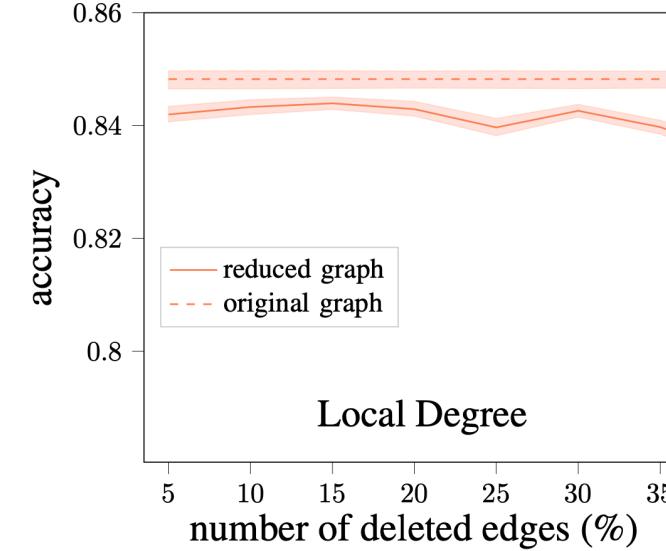
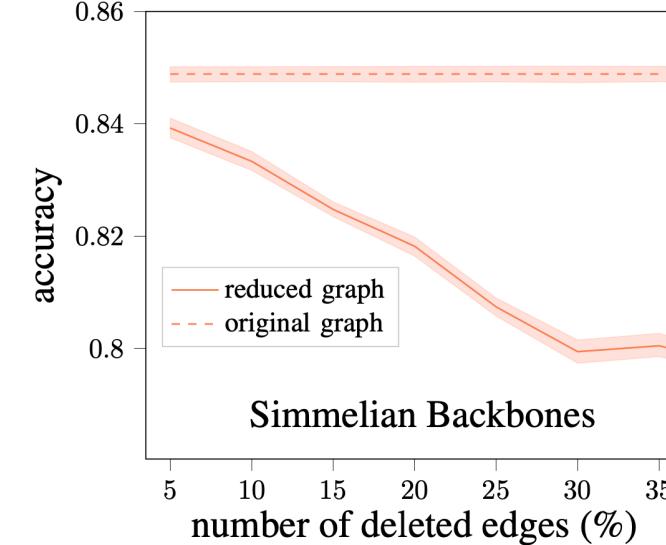
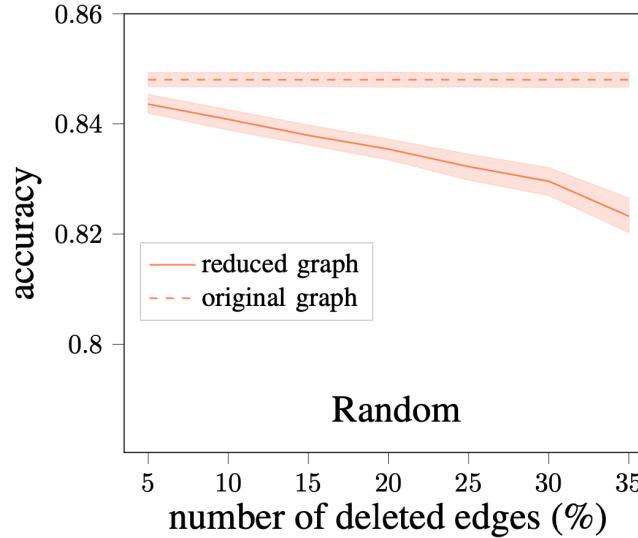
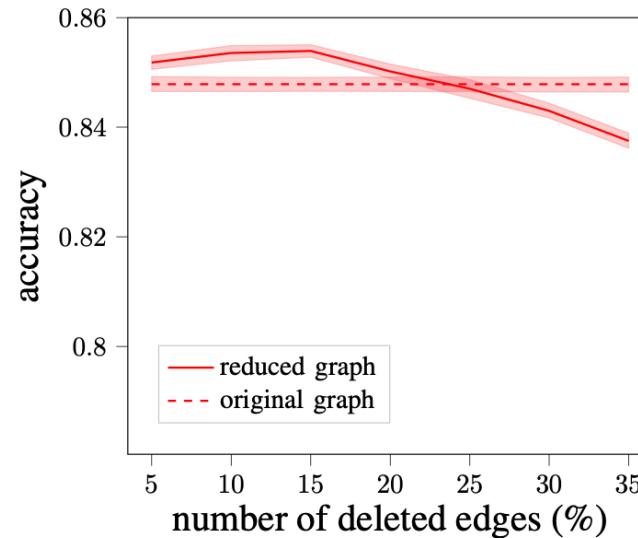


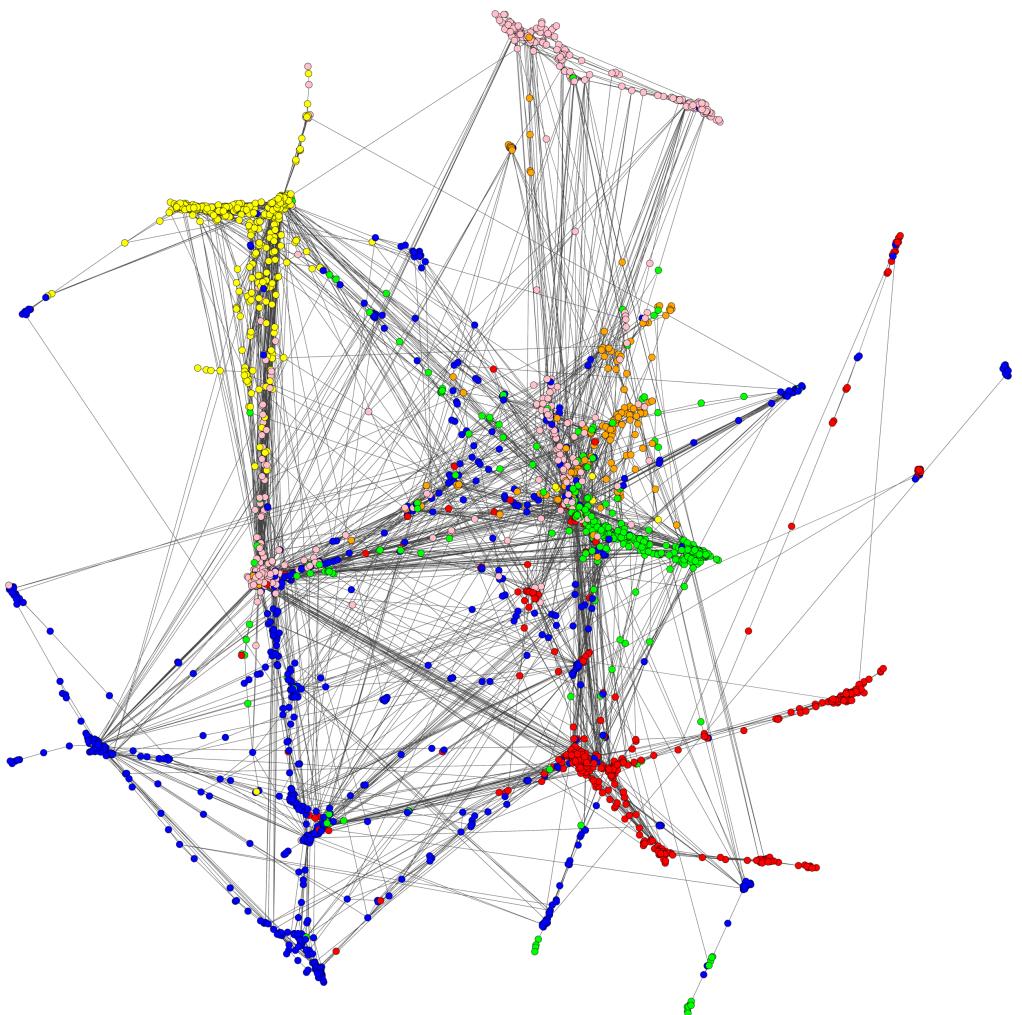
accuracy



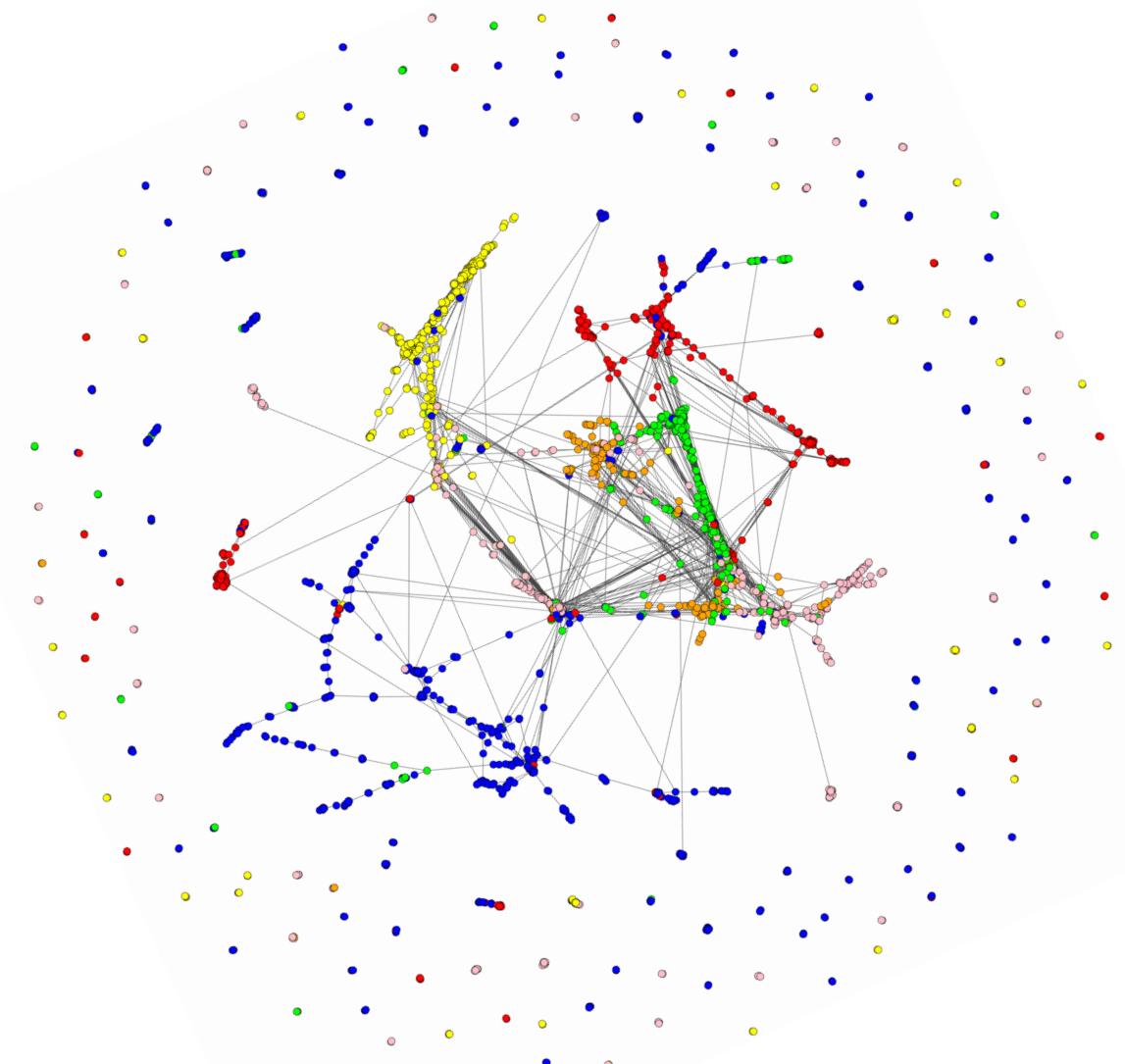
Comparison

Ours →

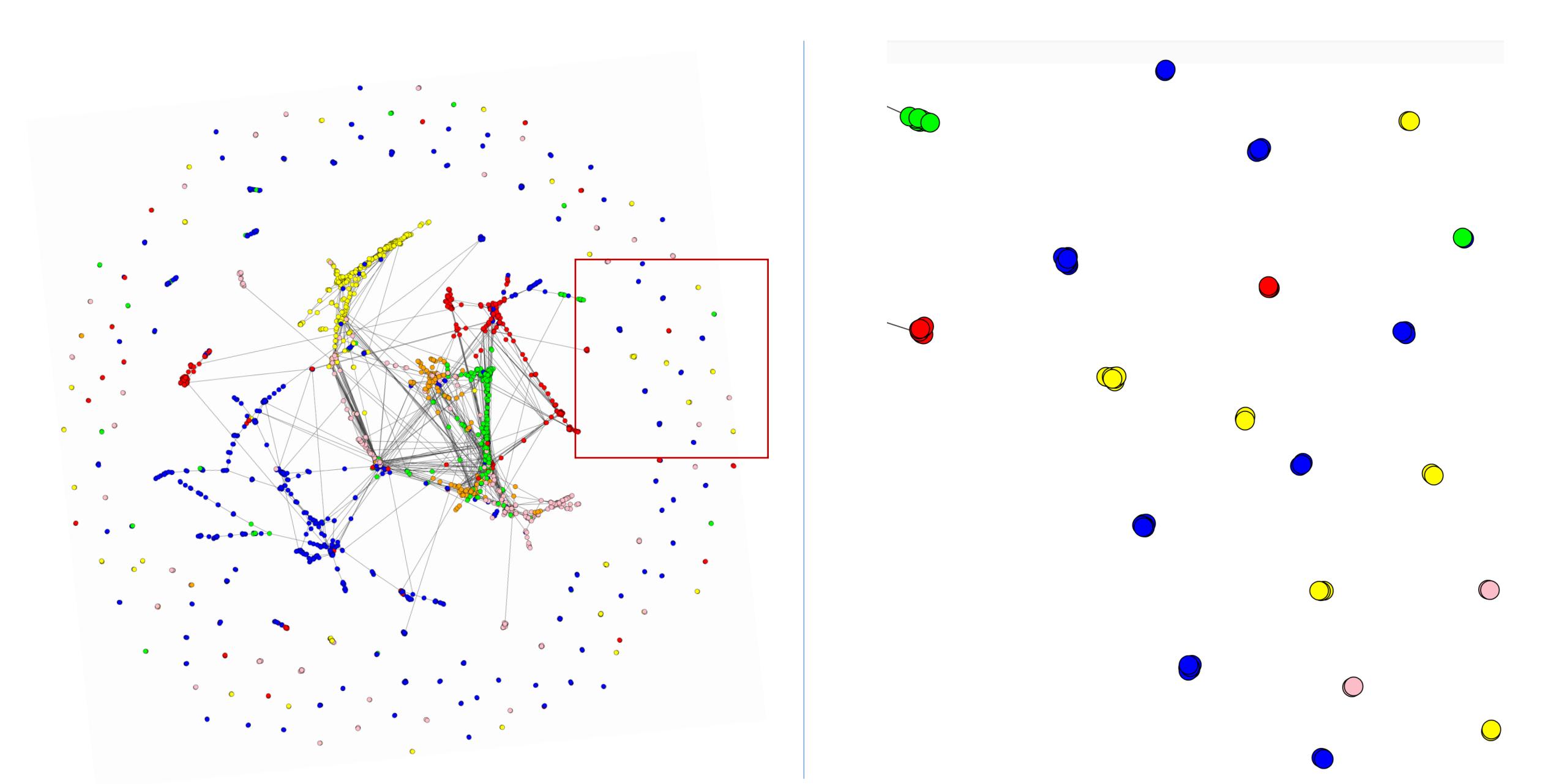




Original Graph
Test Acc : 0.848



30% edges deleted
Test Acc : 0.843



Thank You!