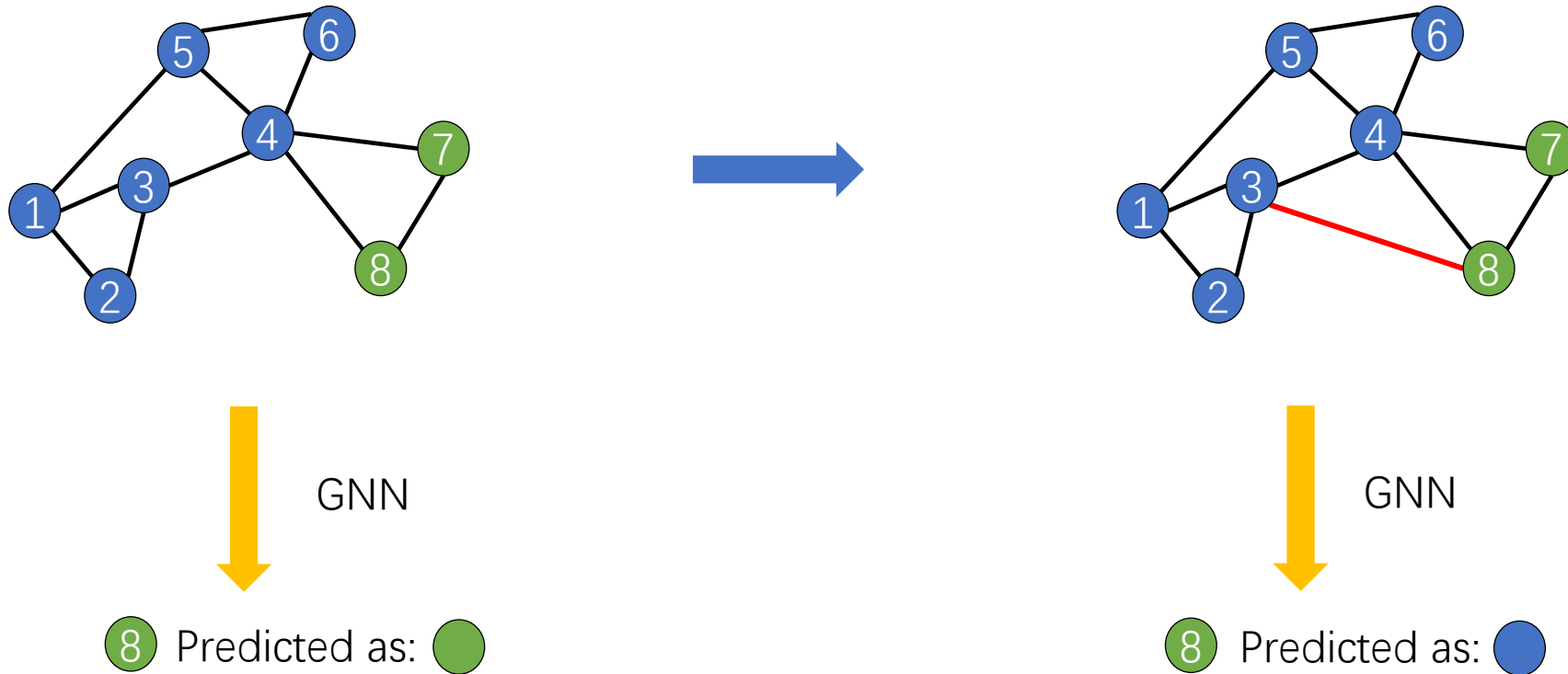


# Graph Structure Learning for Robust Graph Neural Networks

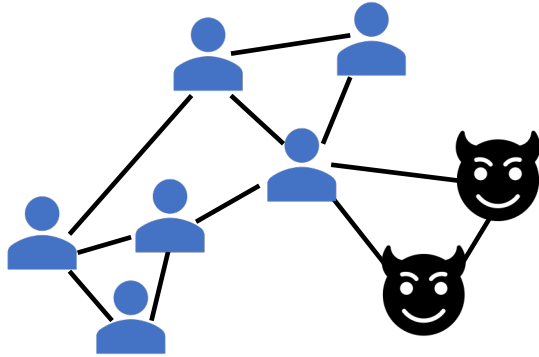
Wei Jin, Yao Ma, Xiaorui Liu, Xianfeng Tang, Suhang Wang, Jiliang Tang

KDD 2020

# Adversarial Attacks on GNN



# Consequences



- Financial Systems
  - Credit Card Fraud Detection
- Recommender Systems
  - Social Recommendation
  - Product Recommendation
- ...

# Pro-GNN: Defend Against Adversarial Attacks

## Attack Setting

- Untargeted structure attack
- Poisoning attack
- Node classification
  - Graph dataset  $G = (A, X)$
  - Graph neural network  $f: f(x_i) \rightarrow \hat{y}_i$

## Defense Goal

- Improve the overall performance of GNN on the perturbed graph

# Pro-GNN: Defend Against Adversarial Attacks

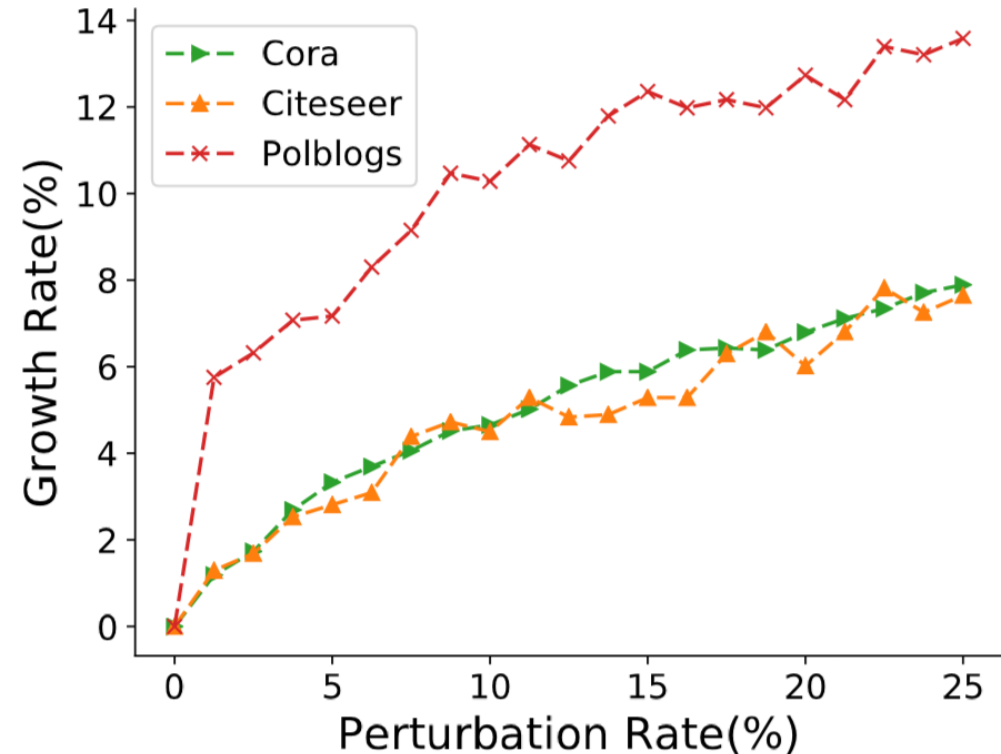
## Graph Properties

- Low-rank
- Sparsity
- Feature smoothness

# Pro-GNN: Defend Against Adversarial Attacks

## Graph Properties

- Low-rank
- Sparsity
- Feature smoothness

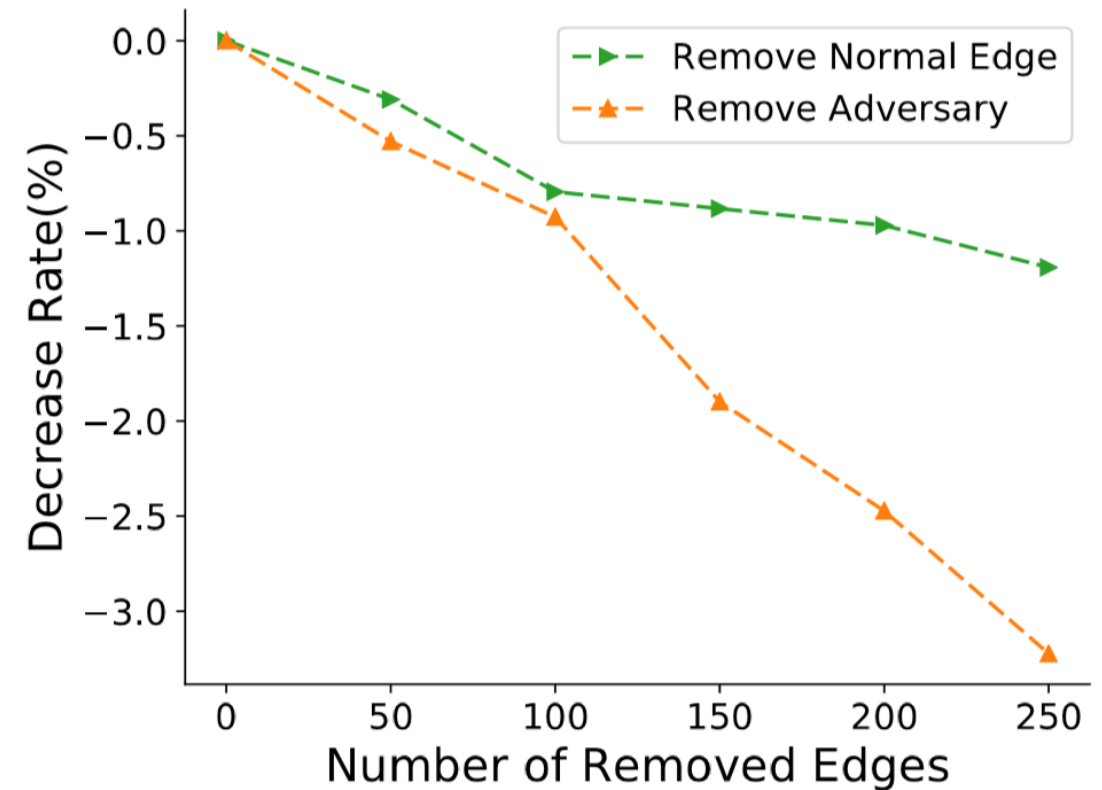


(b) Rank Growth

# Pro-GNN: Defend Against Adversarial Attacks

## Graph Properties

- Low-rank
- Sparsity
- Feature smoothness



# Pro-GNN: Defend Against Adversarial Attacks

## Graph Properties

- Low-rank
- Sparsity
- Feature smoothness

Dataset	$r(\%)$	edge+	edge-	edges	ranks	clustering coefficients
Cora	0	0	0	5069	2192	0.2376
	5	226	27	5268	2263	0.2228
	10	408	98	5380	2278	0.2132
	15	604	156	5518	2300	0.2071
	20	788	245	5633	2305	0.1983
	25	981	287	5763	2321	0.1943
Citeseer	0	0	0	3668	1778	0.1711
	5	181	2	3847	1850	0.1616
	1	341	25	3985	1874	0.1565
	15	485	65	4089	1890	0.1523
	20	614	119	4164	1902	0.1483
	25	743	174	4236	1888	0.1467
Polblogs	0	0	0	16714	1060	0.3203
	5	732	103	17343	1133	0.2719
	10	1347	324	17737	1170	0.2825
	15	1915	592	18038	1193	0.2851
	20	2304	1038	17980	1193	0.2877
	25	2500	1678	17536	1197	0.2723

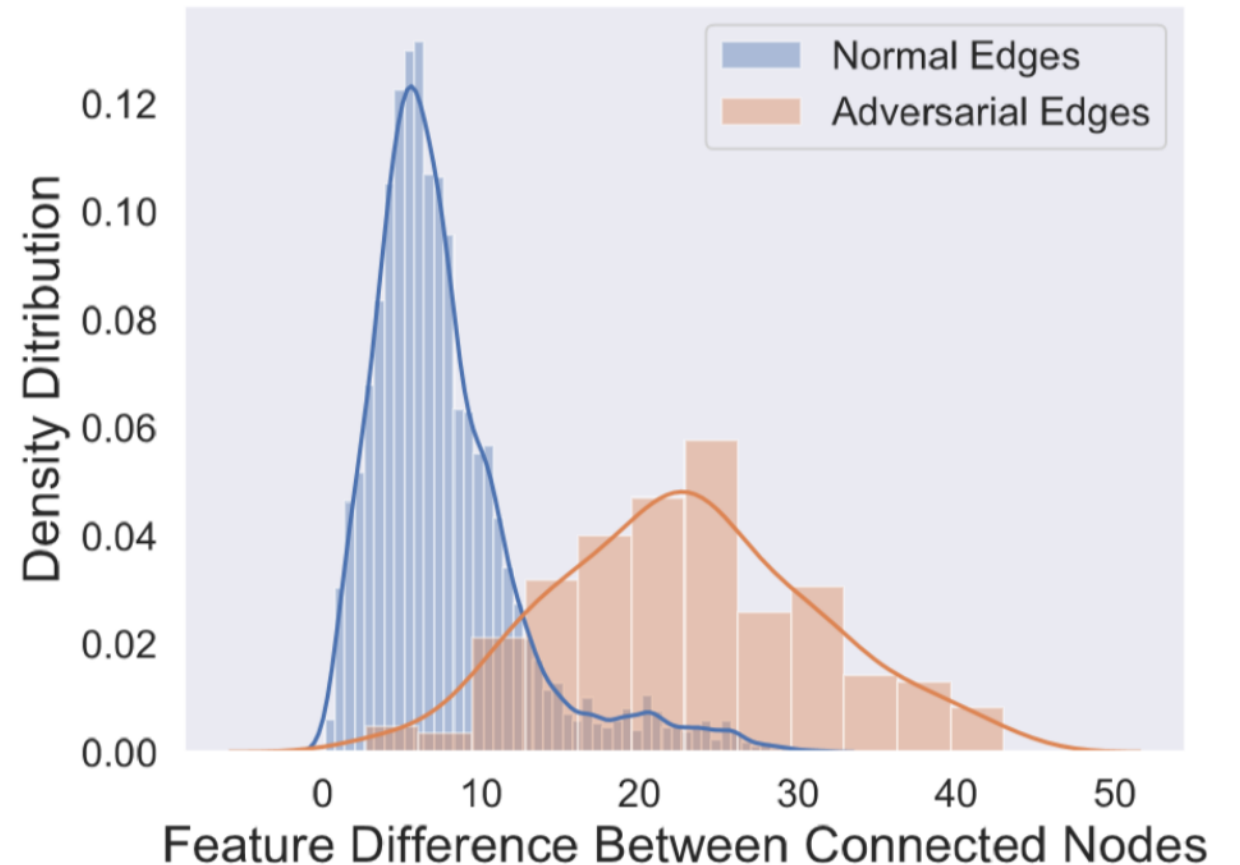
Table Credit: Adversarial Attacks and Defenses on Graphs: A Review and Empirical Study



# Pro-GNN: Defend Against Adversarial Attacks

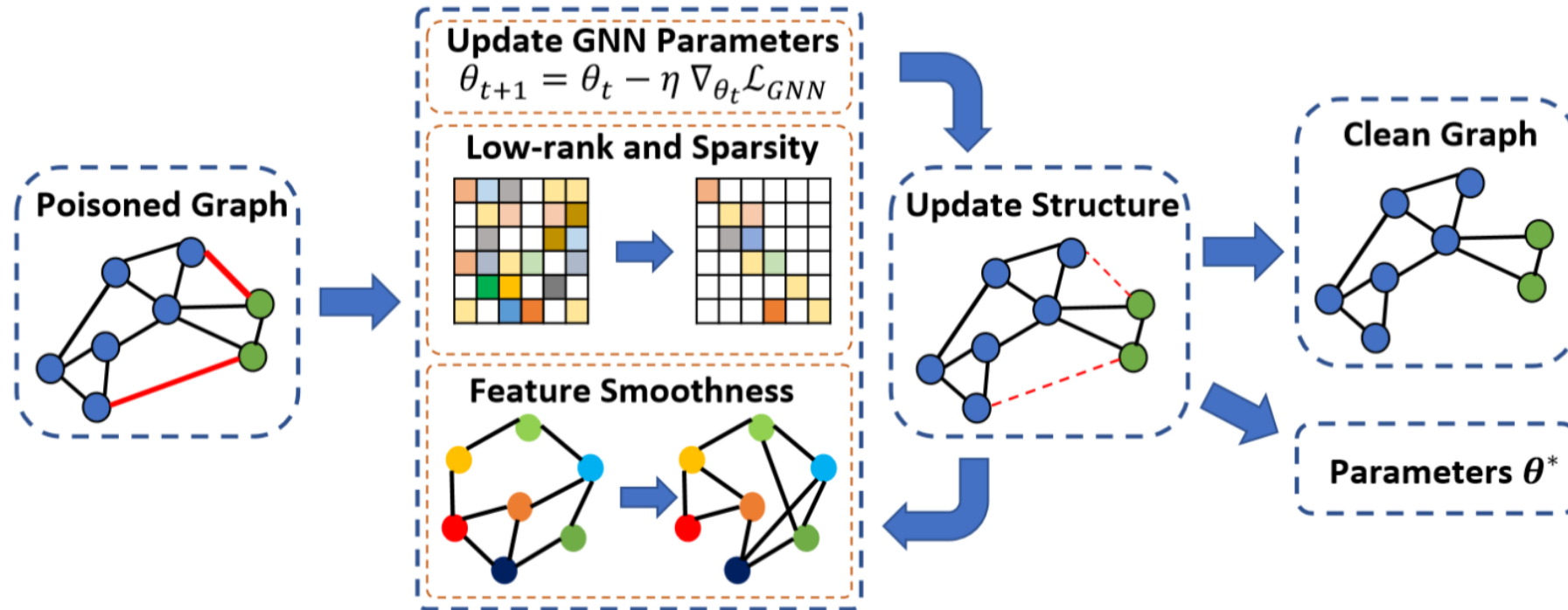
## Graph Properties

- Low-rank
- Sparsity
- Feature smoothness



**(d) Feature Smoothness**

# Pro-GNN: Framework



**Figure 2: Overall framework of Pro-GNN. Dash lines indicate smaller weights.**

# Pro-GNN: Modelling

- Low rank and sparsity

$$\arg \min_{\mathbf{S} \in \mathcal{S}} \mathcal{L}_0 = \|\mathbf{A} - \mathbf{S}\|_F^2 + \alpha \|\mathbf{S}\|_1 + \beta \|\mathbf{S}\|_*, \text{ s.t.}, \mathbf{S} = \mathbf{S}^\top$$

↑ Perturbed Graph  
↓ Graph Variable to be Recovered

$$\|\mathbf{S}\|_1 = \sum_{ij} |\mathbf{S}_{ij}| \quad \|\mathbf{S}\|_* = \sum_{i=1}^{\text{rank}(\mathbf{S})} \sigma_i$$

# Pro-GNN: Modelling

- Feature smoothness

$$\mathcal{L}_s = \text{tr}(\mathbf{X}^T \hat{\mathbf{L}} \mathbf{X}) = \frac{1}{2} \sum_{i,j=1}^N \mathbf{S}_{ij} (\mathbf{x}_i - \mathbf{x}_j)^2.$$

# Pro-GNN: Modelling

- Overall objective

$$\arg \min_{S \in \mathcal{S}, \theta} \mathcal{L} = \mathcal{L}_0 + \lambda \mathcal{L}_s + \gamma \mathcal{L}_{GNN}$$

# Pro-GNN: Optimization

- Overall objective

$$\begin{aligned} \arg \min_{\mathbf{S} \in \mathcal{S}, \theta} \mathcal{L} &= \mathcal{L}_0 + \lambda \mathcal{L}_s + \gamma \mathcal{L}_{GNN} & (9) \\ &= \|\mathbf{A} - \mathbf{S}\|_F^2 + \alpha \|\mathbf{S}\|_1 + \beta \|\mathbf{S}\|_* + \gamma \mathcal{L}_{GNN}(\theta, \mathbf{S}, \mathbf{X}, \mathcal{Y}_L) + \lambda \text{tr}(\mathbf{X}^T \hat{\mathbf{L}} \mathbf{X}) \\ \text{s.t.} \quad &\mathbf{S} = \mathbf{S}^T, \end{aligned}$$

# Pro-GNN: Optimization

- Overall objective

$$\arg \min_{\mathbf{S} \in \mathcal{S}, \theta} \mathcal{L} = \mathcal{L}_0 + \lambda \mathcal{L}_s + \gamma \mathcal{L}_{GNN} \quad (9)$$

$$= \underbrace{\|\mathbf{A} - \mathbf{S}\|_F^2}_{\text{data fit}} + \alpha \|\mathbf{S}\|_1 + \beta \|\mathbf{S}\|_* + \underbrace{\gamma \mathcal{L}_{GNN}(\theta, \mathbf{S}, \mathbf{X}, \mathcal{Y}_L) + \lambda \text{tr}(\mathbf{X}^T \hat{\mathbf{L}} \mathbf{X})}_{\text{regularization}}$$

*s.t.*      $\mathbf{S} = \mathbf{S}^T,$

# Pro-GNN: Optimization

- Alternating Optimization

**Update  $\theta$ :** 
$$\min_{\theta} \mathcal{L}_{GNN}(\theta, \mathbf{S}, \mathbf{X}, \mathcal{Y}_L) = \sum_{u \in \mathcal{V}_L} \ell(f_{\theta}(\mathbf{X}, \mathbf{S})_u, y_u)$$

**Update  $\mathbf{S}$ :** 
$$\min_{\mathbf{S}} \mathcal{L}(\mathbf{S}, \mathbf{A}) + \alpha \|\mathbf{S}\|_1 + \beta \|\mathbf{S}\|_* \quad s.t., \quad \mathbf{S} = \mathbf{S}^T, \mathbf{S} \in \mathcal{S},$$

where 
$$\mathcal{L}(\mathbf{S}, \mathbf{A}) = \|\mathbf{A} - \mathbf{S}\|_F^2 + \mathcal{L}_{GNN}(\theta, \mathbf{S}, \mathbf{X}, \mathbf{Y}) + \lambda tr(\mathbf{X}^T \hat{\mathbf{L}} \mathbf{X}).$$



# Pro-GNN: Optimization

- Incremental Proximal Descent method

$$\min_{\mathbf{S}} \mathcal{L}(\mathbf{S}, \mathbf{A}) + \alpha \|\mathbf{S}\|_1 + \beta \|\mathbf{S}\|_*$$

For each iteration, do

$$\begin{cases} \mathbf{S}^{(k)} = \mathbf{S}^{(k-1)} - \eta \cdot \nabla_{\mathbf{S}} (\mathcal{L}(\mathbf{S}, \mathbf{A})), \\ \mathbf{S}^{(k)} = \text{prox}_{\eta\beta \|\cdot\|_*} (\mathbf{S}^{(k)}), \\ \mathbf{S}^{(k)} = \text{prox}_{\eta\alpha \|\cdot\|_1} (\mathbf{S}^{(k)}). \end{cases}$$

# Pro-GNN: Algorithm

---

**Algorithm 1:** Pro-GNN

---

**Data:** Adjacency matrix  $\mathbf{A}$ , Attribute matrix  $\mathbf{X}$ , Labels  $\mathcal{Y}_L$ ,  
Hyper-parameters  $\alpha, \beta, \gamma, \lambda, \tau$ , Learning rate  $\eta, \eta'$

**Result:** Learned adjacency  $\mathbf{S}$ , GNN parameters  $\theta$

- 1 Initialize  $\mathbf{S} \leftarrow \mathbf{A}$
- 2 Randomly initialize  $\theta$
- 3 **while** *Stopping condition is not met* **do**
- 4      $\mathbf{S} \leftarrow \mathbf{S} - \eta \nabla_{\mathbf{S}} (\|\mathbf{S} - \mathbf{A}\|_F^2 + \gamma \mathcal{L}_{GNN} + \lambda \mathcal{L}_s)$
- 5      $\mathbf{S} \leftarrow \text{prox}_{\eta\beta\|\cdot\|_*}(\mathbf{S})$
- 6      $\mathbf{S} \leftarrow \text{prox}_{\eta\alpha\|\cdot\|_1}(\mathbf{S})$
- 7      $\mathbf{S} \leftarrow P_{\mathcal{S}}(\mathbf{S})$
- 8     **for**  $i=1$  to  $\tau$  **do**
- 9          $g \leftarrow \frac{\partial \mathcal{L}_{GNN}(\theta, \mathbf{S}, \mathbf{X}, \mathcal{Y}_L)}{\partial \theta}$
- 10          $\theta \leftarrow \theta - \eta' g$
- 11 **Return**  $\mathbf{S}, \theta$

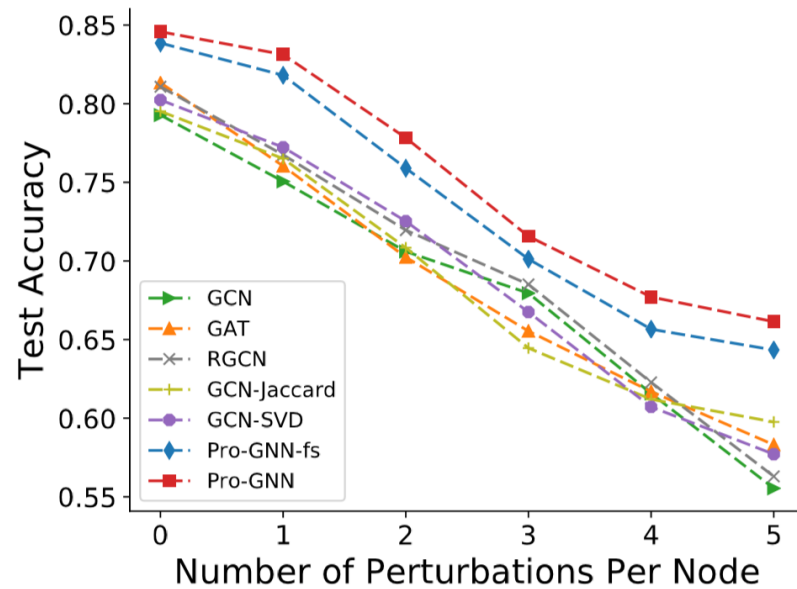
---

# Pro-GNN: Experiments

**Table 2: Node classification performance (Accuracy $\pm$ Std) under non-targeted attack (*metattack*).**

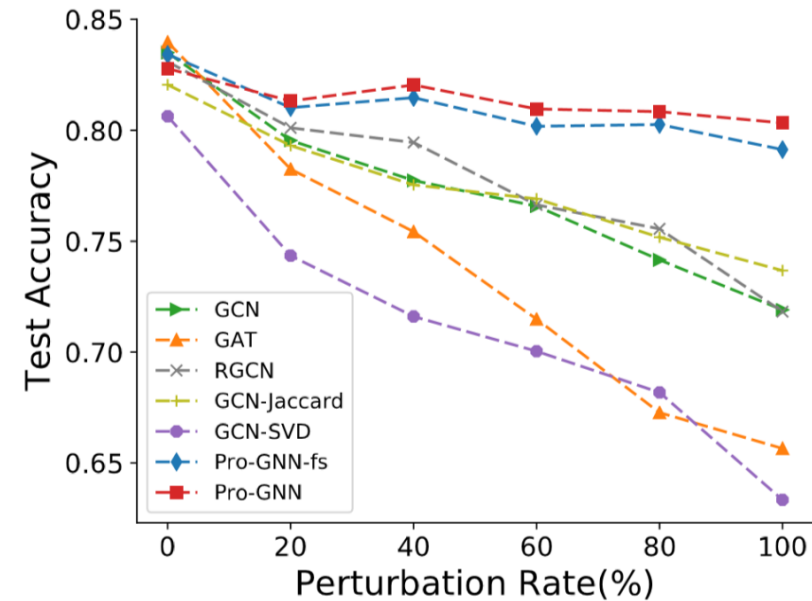
Dataset	Ptb Rate (%)	GCN	GAT	RGCN	GCN-Jaccard <sup>2</sup>	GCN-SVD	Pro-GNN-fs	Pro-GNN <sup>3</sup>
Cora	0	83.50 $\pm$ 0.44	<b>83.97<math>\pm</math>0.65</b>	83.09 $\pm$ 0.44	82.05 $\pm$ 0.51	80.63 $\pm$ 0.45	83.42 $\pm$ 0.52	82.98 $\pm$ 0.23
	5	76.55 $\pm$ 0.79	80.44 $\pm$ 0.74	77.42 $\pm$ 0.39	79.13 $\pm$ 0.59	78.39 $\pm$ 0.54	<b>82.78<math>\pm</math>0.39</b>	82.27 $\pm$ 0.45
	10	70.39 $\pm$ 1.28	75.61 $\pm$ 0.59	72.22 $\pm$ 0.38	75.16 $\pm$ 0.76	71.47 $\pm$ 0.83	77.91 $\pm$ 0.86	<b>79.03<math>\pm</math>0.59</b>
	15	65.10 $\pm$ 0.71	69.78 $\pm$ 1.28	66.82 $\pm$ 0.39	71.03 $\pm$ 0.64	66.69 $\pm$ 1.18	76.01 $\pm$ 1.12	<b>76.40<math>\pm</math>1.27</b>
	20	59.56 $\pm$ 2.72	59.94 $\pm$ 0.92	59.27 $\pm$ 0.37	65.71 $\pm$ 0.89	58.94 $\pm$ 1.13	68.78 $\pm$ 5.84	<b>73.32<math>\pm</math>1.56</b>
	25	47.53 $\pm$ 1.96	54.78 $\pm$ 0.74	50.51 $\pm$ 0.78	60.82 $\pm$ 1.08	52.06 $\pm$ 1.19	56.54 $\pm$ 2.58	<b>69.72<math>\pm</math>1.69</b>

# Pro-GNN: Experiments



(a) Cora

Under Nettack



(a) Cora

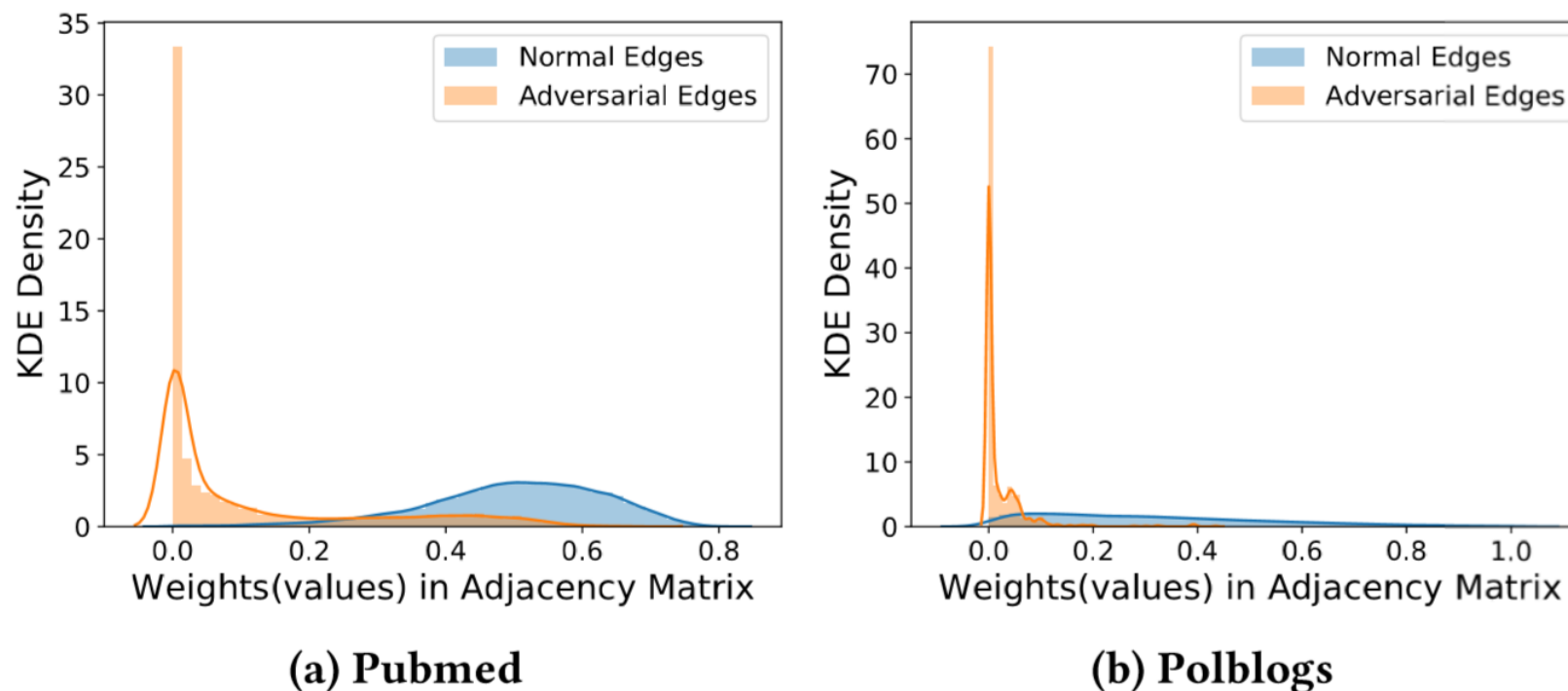
Under Random Noise

# Pro-GNN: Importance of Graph Structure Learning

**Table 3: Node classification accuracy given the graph under 25% perturbation by *metattack*.**

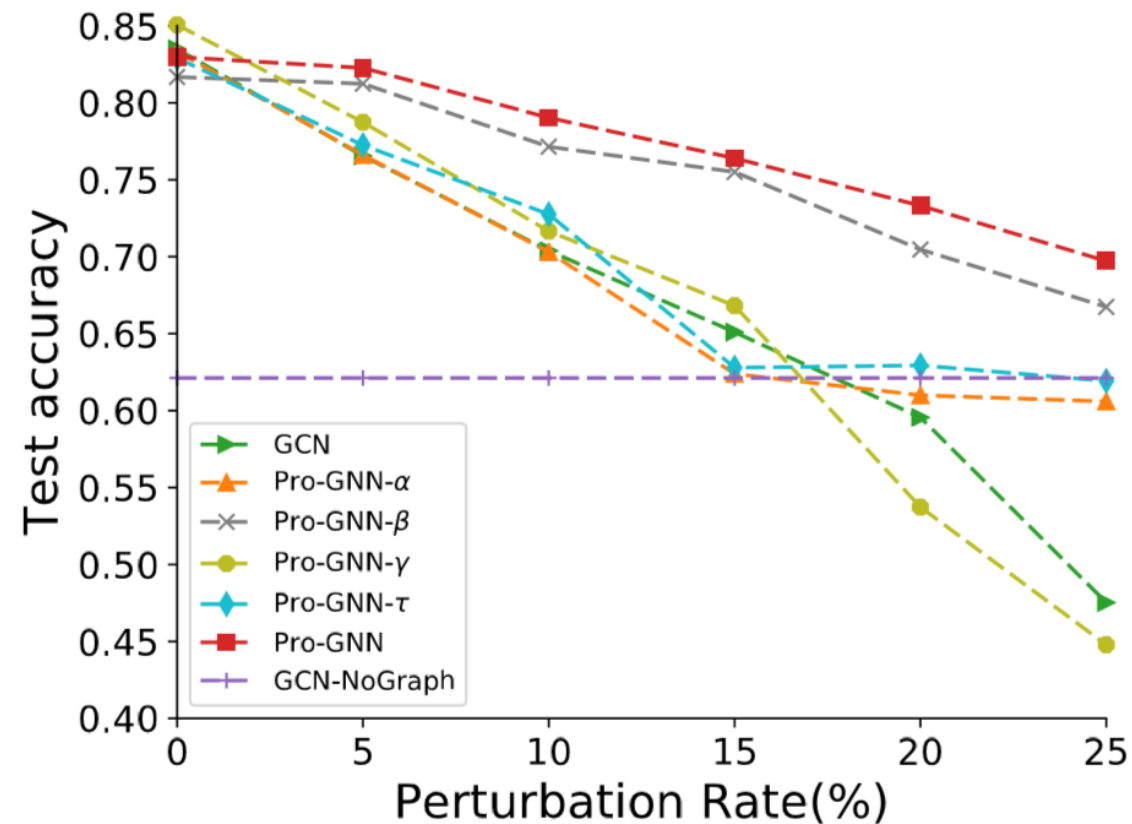
	GCN	GCN-NoGraph	Pro-GNN
Cora	47.53±1.96	62.12±1.55	<b>69.72±1.69</b>
Citeseer	56.94±2.09	63.75±3.23	<b>68.95±2.78</b>
Polblogs	49.23±1.36	51.79±0.62	<b>63.18±4.40</b>
Pubmed	75.50±0.17	84.14±0.11	<b>86.86±0.19</b>

# Pro-GNN: Importance of Graph Structure Learning



**Figure 5: Weight density distributions of normal and adversarial edges on the learned graph.**

# Pro-GNN: Ablation Study



(a) Cora

# Conclusion

- We found that graph adversarial attack can break important graph properties
- We introduced a novel defense approach Pro-GNN that learns the graph structure and GNN parameters simultaneously
- Our experiments show that our model consistently improves the overall robustness under various adversarial attacks.

Paper Link:

<https://arxiv.org/abs/2005.10203>

Code:

<https://github.com/ChandlerBang/Pro-GNN>





THANK YOU