## Adversarial Graph Contrastive Learning with Information Regularization

(WWW, 2022)

paper

#### Graph Representation Learning

Given an attributed graph:  $G = \{\mathcal{V}, \mathcal{E}, X\}$ , Where  $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$ ,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ ,

 $X \in \mathbb{R}^{n \times d}$ , Adjacency matrix can be defined as :  $A \in \{0, 1\}^{n \times n}$ 

Objective is to learn an encoder  $f : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d'}$ ,

 $\mathbf{H} = f(\mathbf{A}, \mathbf{X})$ , where  $\mathbf{H}[i, :] \in \mathbb{R}^{d'}$  is the embedding of node  $v_i$ .



## Contrastive Learning for GRL



## Motivation

#### Limitation of Existing Augment methods for GCL:

- Uncontrollable: The data augmentation on the graph, could be either too similar to or totally different from the original graph.
- Vulnerable : GNNs are known to be vulnerable to adversarial attacks

How to generate a new graph that is hard enough for the model to discriminate from the original one, and in the meanwhile also maintains the desired properties?

Introduce adversarial training into GCL.

On the one hand, since the perturbation is under the constraint, the adversarial sample still stays close enough to the original one.

On the other hand, the adversarial attack makes sure the adversarial sample is hard to discriminate from the other view by increasing the contrastive loss.





Projected Gradient Descent Attack

$$\Delta_t = \Pi_{\|\Delta\|_{\infty} \le \delta} (\Delta_{t-1} + \eta \cdot \operatorname{sgn}(\nabla_{\Delta_{t-1}} L(\mathbf{Z} + \Delta_{t-1})), \qquad (7)$$

where L() is loss function of input matrix Z,  $\Delta_t$  is the perturbation matrix. **sgn**() is sign function,  $\eta$  is step size.

Adversarial Objective: 
$$G_{adv} = \arg \max_{G'} L_{con}(G_1, G'),$$
 (8)

where G' = (A', X') is generated from the original graph G and satisfy:

$$\sum_{i,j} |\mathbf{A}'[i,j] - \mathbf{A}[i,j]| \le \Delta_{\mathbf{A}},\tag{9}$$

$$\sum_{i,j} |\mathbf{X}'[i,j] - \mathbf{X}[i,j]| \le \Delta_{\mathbf{X}}.$$
(10)

Attack on Structure:

 $A_{adv} = A + C \circ L_A, \qquad (11)$  $C = \overline{A} - A, \qquad (12)$ 

Where  $\overline{A} = \mathbf{1}_{n \times n} - \mathbf{I}_n - \mathbf{A}$ ,  $\mathbf{L}_{\mathbf{A}} \in \{0, 1\}^{n \times n}$ 

 $L_A$  is relaxed to its convex hull  $\tilde{L}_A \in [0,1]^{n \times n}$ 

Attack on Feature matrix:

$$\mathbf{X}_{\mathrm{adv}} = \mathbf{X} + \mathbf{L}_{\mathbf{X}},\tag{13}$$

Where  $L_X \in \mathbb{R}^{n \times d}$  is the perturbation on the feature matrix.

# Adversarial Training

Update for adversarial perturbation:

$$\begin{split} \tilde{\mathbf{L}}_{\mathbf{A}}^{(t)} &= \Pi_{\mathcal{S}_{\mathbf{A}}} [\tilde{\mathbf{L}}_{\mathbf{A}}^{(t-1)} + \alpha \cdot \mathbf{G}_{\mathbf{A}}^{(t)}], \qquad (14) \\ \mathbf{L}_{\mathbf{X}}^{(t)} &= \Pi_{\mathcal{S}_{\mathbf{X}}} [\mathbf{L}_{\mathbf{X}}^{(t-1)} + \beta \cdot \operatorname{sgn}(\mathbf{G}_{\mathbf{X}}^{(t)})], \qquad (15) \\ \mathbf{G}_{\mathbf{A}}^{(t)} &= \nabla_{\tilde{\mathbf{L}}_{\mathbf{A}}^{(t-1)}} L_{\operatorname{con}}(G_{1}, G_{\operatorname{adv}}^{(t-1)}), \qquad (16) \\ \mathbf{G}_{\mathbf{X}}^{(t)} &= \nabla_{\mathbf{L}_{\mathbf{X}}^{(t-1)}} L_{\operatorname{con}}(G_{1}, G_{\operatorname{adv}}^{(t-1)}), \qquad (17) \\ G_{\operatorname{adv}}^{(t-1)} &= \{\mathbf{A} + \mathbf{C} \circ \tilde{\mathbf{L}}_{\mathbf{A}}^{(t-1)}, \mathbf{X} + \mathbf{L}_{\mathbf{X}}^{(t-1)}\} \\ \mathcal{S}_{\mathbf{A}} &= \{\tilde{\mathbf{L}}_{\mathbf{A}} |\sum_{i,j} \tilde{\mathbf{L}}_{\mathbf{A}} \leq \Delta_{\mathbf{A}}, \tilde{\mathbf{L}}_{\mathbf{A}} \in [0, 1]^{n \times n}\} \\ \mathcal{S}_{\mathbf{X}} &= \{\mathbf{L}_{\mathbf{X}} |\|\mathbf{L}_{\mathbf{X}}\|_{\infty} \leq \delta_{\mathbf{X}}, \mathbf{L}_{\mathbf{X}} \in \mathbb{R}^{n \times d}\} \end{split}$$

Adversarial Graph Contrastive Learning

$$L(G_1, G_2, G_{adv}) = L_{con}(G_1, G_2) + \epsilon_1 L_{con}(G_1, G_{adv}), \quad (19)$$

The projection operation  $\Pi_{S_X}$  simply clips  $L_X$  into the range  $[-\delta_X, \delta_X]$ 

The projection operation  $\Pi_{\mathcal{S}_A}$  is calculated as

$$\Pi_{\mathcal{S}_{A}}(\mathbf{Z}) = \begin{cases} P_{[0,1]}[\mathbf{Z} - \mu \mathbf{1}_{n \times n}], & \text{if } \mu > 0, \text{ and} \\ \sum_{i,j} P_{[0,1]}[\mathbf{Z} - \mu \mathbf{1}_{n \times n}] = \Delta_{A}, \\ P_{[0,1]}[\mathbf{Z}], & \text{if } \sum_{i,j} P_{[0,1]}[\mathbf{Z}] \le \Delta_{A}, \end{cases}$$
(18)

## Information Regularization

**THEOREM3.1**. For two graph views  $G_1$  and  $G_2$  independently transformed from the graph G, the density ratio of their node embeddings  $H_1$  and  $H_2$  should satisfy  $g(H_2[i, :], H_1[i, :]) \leq g(H_2[i, :], H[i, :])$  and  $g(H_1[i, :], H_2[i, :]) \leq g(H_1[i, :], H[i, :])$ , where H is the node embeddings of the original graph.

The node embeddings of two views  $H_1$ ,  $H_2$  and of original graph H satisfy the Markov relation as follow:

$$H_1 \to H \to H_2$$
$$H_1[i,:] \to H[i,:] \to H_2[i,:]$$

Then, we get:

$$p(\mathbf{H}_{2}[i,:]|\mathbf{H}_{1}[i,:]) = p(\mathbf{H}_{2}[i,:]|\mathbf{H}[i,:])p(\mathbf{H}[i,:]|\mathbf{H}_{1}[i,:]) \qquad (24)$$
$$\leq p(\mathbf{H}_{2}[i,:]|\mathbf{H}[i,:]), \qquad (25)$$

$$\frac{p(\mathbf{H}_{2}[i,:]|\mathbf{H}_{1}[i,:])}{p(\mathbf{H}_{2}[i,:])} \le \frac{p(\mathbf{H}_{2}[i,:]|\mathbf{H}[i,:])}{p(\mathbf{H}_{2}[i,:])}.$$
(26)

Since  $g(\mathbf{a}, \mathbf{b}) \propto \frac{p(\mathbf{a}|\mathbf{b})}{p(\mathbf{a})}$ , then  $g(\mathbf{H}_2[i, :], \mathbf{H}_1[i, :]) \leq g(\mathbf{H}_2[i, :], \mathbf{H}[i, :])$ .

 $g(H_2[i, :], H_1[i, :]) \leq g(H_2[i, :], H[i, :])$ 

```
g(H_1[i,:], H_2[i,:]) \leq g(H_1[i,:], H[i,:])
```

According to the previous definition,  $g(a, b) = e^{\theta(a,b)/\tau}$ , simply replace  $g(\cdot, \cdot)$  with  $\theta(\cdot, \cdot)$ , then combine two upper bounds into one:

 $2 \cdot \theta(\mathbf{H}_1[i,:], \mathbf{H}_2[i,:]) \le \theta(\mathbf{H}_2[i,:], \mathbf{H}[i,:]) + \theta(\mathbf{H}_1[i,:], \mathbf{H}[i,:]).$ 

 $d_i = 2 \cdot \theta(\mathbf{H}_1[i,:], \mathbf{H}_2[i,:]) - (\theta(\mathbf{H}_2[i,:], \mathbf{H}[i,:]) + \theta(\mathbf{H}_1[i,:], \mathbf{H}[i,:]))$ 

$$L_I(G_1, G_2, G) = \frac{1}{n} \sum_{i=1}^n \max\{d_i, 0\}.$$

 $L(G_1, G_2, G_{adv}) = L_{con}(G_1, G_2) + \epsilon_1 L_{con}(G_1, G_{adv}) + \epsilon_2 L_I(G_1, G_2, G), \quad (30)$ 

#### Training pseudocode

Algorithm 1 Algorithm of ARIEL

```
Input data: Graph G = (\mathbf{A}, \mathbf{X})

Input parameters: \alpha, \beta, \Delta_{\mathbf{A}}, \delta_{\mathbf{X}}, \epsilon_1, \epsilon_2, \gamma and T

Randomly initialize the graph encoder f

for iteration k = 0, 1, \cdots do

Sample a subgraph G_s from G

Generate two views G_1 and G_2 from G_s

Generate the adversarial view G_{adv} according to Equations

(15), (14)

Update model f to minimize L(G_1, G_2, G_{adv}) in Equation (30)

if (k + 1) \mod T = 0 then

Update \epsilon_1 \leftarrow \gamma * \epsilon_1

end if

end for

return: Node embedding matrix \mathbf{H} = f(\mathbf{A}, \mathbf{X})
```

#### Experiments

*RQ1.* How effective is the proposed ArieL in comparison with previous graph contrastive learning methods on the node classification task?

Method	Cora	CiteSeer	Amazon-Computers	Amazon-Photo	Coauthor-CS	<b>Coauthor-Physics</b>
GCN	$84.14\pm0.68$	$69.02 \pm 0.94$	$88.03 \pm 1.41$	$92.65 \pm 0.71$	$92.77 \pm 0.19$	$95.76 \pm 0.11$
GAT	$83.18 \pm 1.17$	$69.48 \pm 1.04$	$85.52\pm2.05$	$91.35 \pm 1.70$	$90.47 \pm 0.35$	$94.82 \pm 0.21$
DeepWalk+features	$79.82 \pm 0.85$	$67.14 \pm 0.81$	$86.23 \pm 0.37$	$90.45 \pm 0.45$	$85.02\pm0.44$	$94.57 \pm 0.20$
DGI	$84.24 \pm 0.75$	$69.12 \pm 1.29$	$88.78 \pm 0.43$	$92.57 \pm 0.23$	$92.26 \pm 0.12$	$95.38 \pm 0.07$
GMI	$82.43 \pm 0.90$	$70.14 \pm 1.00$	$83.57 \pm 0.40$	$88.04 \pm 0.59$	OOM	OOM
MVGRL	$\textbf{84.39} \pm \textbf{0.77}$	$69.85 \pm 1.54$	$89.02\pm0.21$	$92.92 \pm 0.25$	$92.22\pm0.22$	$95.49 \pm 0.17$
GRACE	$83.40 \pm 1.08$	$69.47 \pm 1.36$	$87.77 \pm 0.34$	$92.62 \pm 0.25$	$93.06\pm0.08$	$95.64 \pm 0.08$
GCA-DE	$82.57 \pm 0.87$	$72.11 \pm 0.98$	$88.10\pm0.33$	$92.87 \pm 0.27$	$93.08 \pm 0.18$	$95.62 \pm 0.13$
GCA-PR	$82.54 \pm 0.87$	$72.16 \pm 0.73$	$88.18 \pm 0.39$	$92.85 \pm 0.34$	$93.09\pm0.15$	$95.58 \pm 0.12$
GCA-EV	$81.80 \pm 0.92$	$67.07 \pm 0.79$	$87.95 \pm 0.43$	$92.63 \pm 0.33$	$93.06 \pm 0.14$	$95.64 \pm 0.08$
ArieL	$84.28 \pm 0.96$	$\textbf{72.74} \pm \textbf{1.10}$	$91.13 \pm 0.34$	$94.01 \pm 0.23$	$\textbf{93.83} \pm \textbf{0.14}$	$95.98 \pm 0.05$

Table 2: Node classification accuracy in percentage on six real-world datasets. We bold the results with the best mean accuracy. The methods above the line are the supervised ones, and the ones below the line are unsupervised. OOM stands for Out-of-Memory on our 32G GPUs.

Method	Cora	CiteSeer	Amazon-Computers	Amazon-Photos	<b>Coauthor-CS</b>	<b>Coauthor-Physics</b>
GCN	$80.03\pm0.91$	$62.98 \pm 1.20$	$84.10 \pm 1.05$	$91.72 \pm 0.94$	$80.32 \pm 0.59$	$87.47 \pm 0.38$
GAT	$79.49 \pm 1.29$	$63.30 \pm 1.11$	$81.60 \pm 1.59$	$90.66 \pm 1.62$	$77.75\pm0.80$	$86.65 \pm 0.41$
DeepWalk+features	$74.12 \pm 1.02$	$63.20\pm0.80$	$79.08 \pm 0.67$	$88.06 \pm 0.41$	$49.30 \pm 1.23$	$79.26 \pm 1.38$
DGI	$80.84 \pm 0.82$	$64.25\pm0.96$	$83.36\pm0.55$	$91.27 \pm 0.29$	$78.73 \pm 0.50$	$85.88 \pm 0.37$
GMI	$79.17 \pm 0.76$	$65.37 \pm 1.03$	$77.42 \pm 0.59$	$89.44 \pm 0.47$	$80.92 \pm 0.64$	$87.72 \pm 0.45$
MVGRL	$\textbf{80.90} \pm \textbf{0.75}$	$64.81 \pm 1.53$	$83.76 \pm 0.69$	$91.76\pm0.44$	$79.49 \pm 0.75$	$86.98 \pm 0.61$
GRACE	$78.55 \pm 0.81$	$63.17 \pm 1.81$	$84.74 \pm 1.13$	$91.26 \pm 0.37$	$80.61 \pm 0.63$	$85.71 \pm 0.38$
GCA	$76.79\pm0.97$	$64.89 \pm 1.33$	$85.05\pm0.51$	$91.71 \pm 0.34$	$82.72\pm0.58$	$89.00\pm0.31$
ArieL	$80.33 \pm 1.25$	$69.13 \pm 0.94$	$88.61 \pm 0.46$	$\textbf{92.99} \pm \textbf{0.21}$	$\textbf{84.43} \pm \textbf{0.59}$	$89.09 \pm 0.31$

*RQ2.* To what extent does ArieL gain robustness over the attacked graph?

Table 3: Node classification accuracy in percentage on the graphs under Metattack, where subgraphs of Amazon-Computers, Amazon-Photo, Coauthor-CS and Coauthor-Physics are used for attack and their results are not directly comparable to those in Table 2. We bold the results with the best mean accuracy. GCA is evaluated on its best variant on each clean graph.

#### Experiments

*RQ3.* How does each part of ArieL contribute to its performance?



#### Thanks