

Learning to Pre-train Graph Neural Networks

Background

Code and datasets can be found in <https://yuanfulu.github.io>

GNNs

- node-level representation

$$\mathbf{h}_v^l = \Psi(\psi; \mathcal{A}, \mathcal{X}, \mathcal{Z})^l$$

$$= \text{UPDATE}(\mathbf{h}_v^{l-1}, \text{AGGREGATE}(\{\mathbf{h}_v^{l-1}, \mathbf{h}_u^{l-1}, \mathbf{z}_{uv}\} : u \in \mathcal{N}_v\}))$$

- graph-level representation

$$\mathbf{h}_G = \Omega(\omega; \mathbf{H}^l) = \text{READOUT}(\{\mathbf{h}_v^l | v \in \mathcal{V}\})$$

Pre-train GNNs

- Pre-training

on a large graph-structured dataset (e.g., multiple small graphs or a large-scale graph)

$$\theta_0 = \arg \min_{\theta} \mathcal{L}^{pre}(f_{\theta}; \mathcal{D}^{pre})$$

- Fine-tuning

on downstream tasks

$$\theta_1 = \theta_0 - \eta \nabla_{\theta_0} \mathcal{L}^{fine}(f_{\theta_0}; \mathcal{D}^{tr})$$

θ_0 is pre-trained without accommodating the adaptation in fine-tuning

A gap between pre-training and fine-tuning!

L2PGNN: Learn to Pre-train GNNs

C1: How to narrow the gap caused by different optimization objectives? **C2: How to simultaneously preserve node- and graph-level information?**

- Existing methods fall into a two-step paradigm with a gap
- Solution: learn to pre-train (meta learning)

- SOTAs either only consider the node level or require supervision for graph-level pre-training
- Solution: intrinsic self-supervision

Task Construction

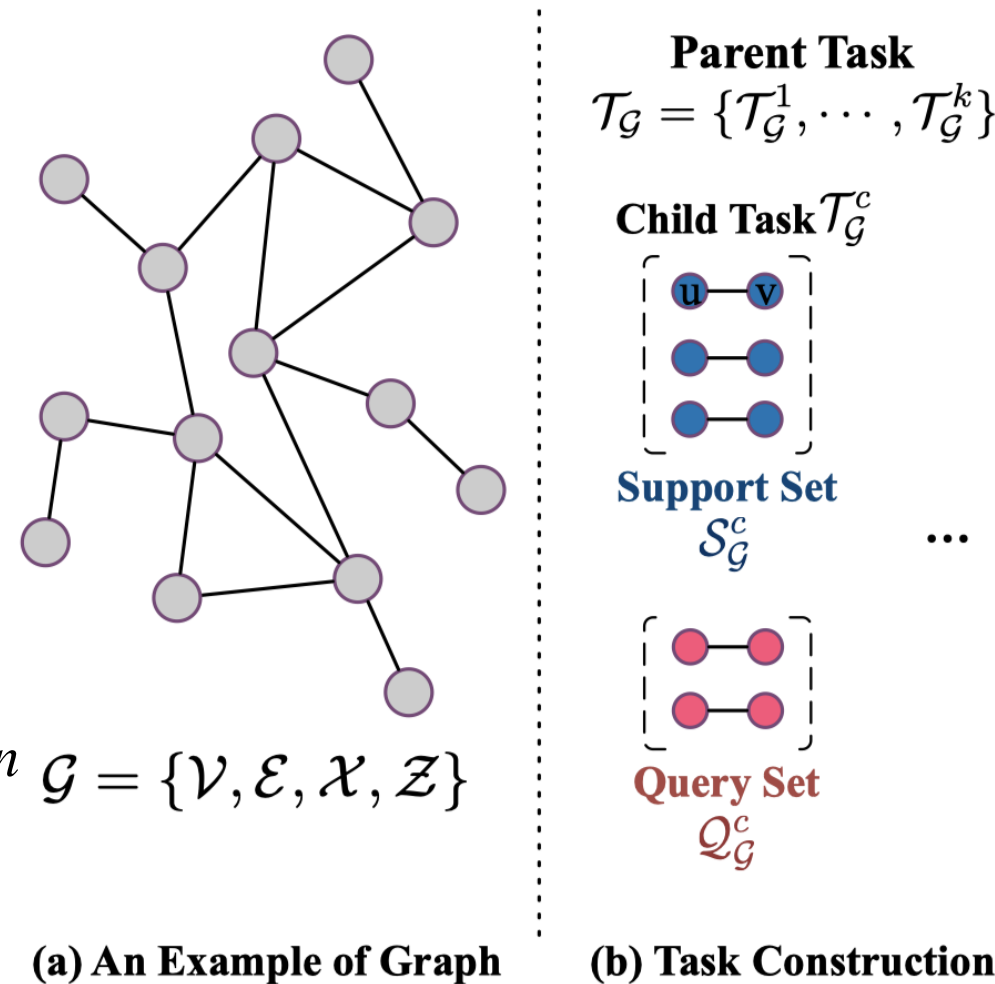
- the pre-training data

$$\mathcal{D}^{pre} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_N\}$$

- A task involving a graph

$$\mathcal{T}_G = (\mathcal{S}_G, \mathcal{Q}_G)$$

- gradient descent w.r.t. the loss on \mathcal{S}_G
- optimize the performance on \mathcal{Q}_G
- simulating the training and testing in $\mathcal{g} = \{\mathcal{V}, \mathcal{E}, \mathcal{X}, \mathcal{Z}\}$ the fine-tuning step



Self-supervised Base Model

- node-level aggregation

$$\mathcal{L}^{node}(\psi; \mathcal{S}_G^c) = \sum_{(u,v) \in \mathcal{S}_G^c} -\ln(\sigma(\mathbf{h}_u^T \mathbf{h}_v)) - \ln(\sigma(-\mathbf{h}_u^T \mathbf{h}_{v'}))$$

- graph-level pooling

$$\mathcal{L}^{graph}(\omega; \mathcal{S}_G) = \sum_{c=1}^k -\log(\sigma(\mathbf{h}_{\mathcal{S}_G^c}^T \mathbf{h}_G)) - \log(\sigma(-\mathbf{h}_{\mathcal{S}_G^c}^T \mathbf{h}_{G'}))$$

$$\mathcal{L}_{\mathcal{T}_G}(\theta; \mathcal{S}_G) = \mathcal{L}^{graph}(\omega; \mathcal{S}_G) + \frac{1}{k} \sum_{c=1}^k \mathcal{L}^{node}(\psi; \mathcal{S}_G^c)$$

Dual Adaptation

$$\psi' = \psi - \alpha \frac{\partial \sum_{c=1}^k \mathcal{L}^{node}(\psi; \mathcal{S}_G^c)}{\partial \psi}$$

$$\omega' = \omega - \beta \frac{\partial \mathcal{L}^{graph}(\omega; \mathcal{S}_G)}{\partial \omega}$$

$$\theta \leftarrow \theta - \gamma \frac{\partial \sum_{G \in \mathcal{D}^{pre}} \mathcal{L}_{\mathcal{T}_G}(\theta'; \mathcal{Q}_G)}{\partial \theta}$$

Experiments & Analysis

Datasets

A new dataset for pre-training GNNs

Dataset	Biology	PreDBLP
#subgraphs	394,925	1,054,309
#labels	40	6
#subgraphs for pre-training	306,925	794,862
#subgraphs for fine-tuning	88,000	299,447

Baselines

- EdgePred
- DGI
- ContextPred
- AttrMasking

GNN Architectures

- GCN, GraphSAGE, GAT, GIN

Table 2: Experimental results (mean \pm std in percent) of different pre-training strategies w.r.t. various GNN architectures. The improvements are relative to the respective GNN without pre-training.

Model	Biology				PreDBLP			
	GCN	GraphSAGE	GAT	GIN	GCN	GraphSAGE	GAT	GIN
No pre-train	63.22 \pm 1.06	65.72 \pm 1.23	68.21 \pm 1.26	64.82 \pm 1.21	62.18 \pm 0.43	61.03 \pm 0.65	59.63 \pm 2.32	69.01 \pm 0.23
EdgePred	64.72 \pm 1.06	67.39 \pm 1.54	67.37 \pm 1.31	65.93 \pm 1.65	65.44 \pm 0.42	63.60 \pm 0.21	55.56 \pm 1.67	69.43 \pm 0.07
DGI	64.33 \pm 1.14	66.69 \pm 0.88	68.37 \pm 0.54	65.16 \pm 1.24	65.57 \pm 0.36	63.34 \pm 0.73	61.30 \pm 2.17	69.34 \pm 0.09
ContextPred	64.56 \pm 1.36	66.31 \pm 0.94	66.89 \pm 1.98	65.99 \pm 1.22	66.11 \pm 0.16	62.55 \pm 0.11	58.44 \pm 1.18	69.37 \pm 0.21
AttrMasking	64.35 \pm 1.23	64.32 \pm 0.78	67.72 \pm 1.16	65.72 \pm 1.31	65.49 \pm 0.52	62.35 \pm 0.58	53.34 \pm 4.77	68.61 \pm 0.16
L2P-GNN (Improv.)	66.48\pm1.59 (5.16%)	69.89\pm1.63 (6.35%)	69.15\pm1.86 (1.38%)	70.13\pm0.95 (8.19%)	66.58\pm0.28 (7.08%)	65.84\pm0.37 (7.88%)	62.24\pm1.89 (4.38%)	70.79\pm0.17 (2.58%)

Performance Comparison

- 6.27% and 3.52% improvements compared to the best baseline
- 8.19% and 7.88% gains relative to non-pretrained models
- negative transfer harms the generalization of the pre-trained GNNs

Comparative Analysis

- Centered Kernel Alignment (CKA) similarity between the parameters
- Smaller similarity, larger changes of model parameters
- changes in loss and performance (delta loss and RUC-AUC/Micro-F1)

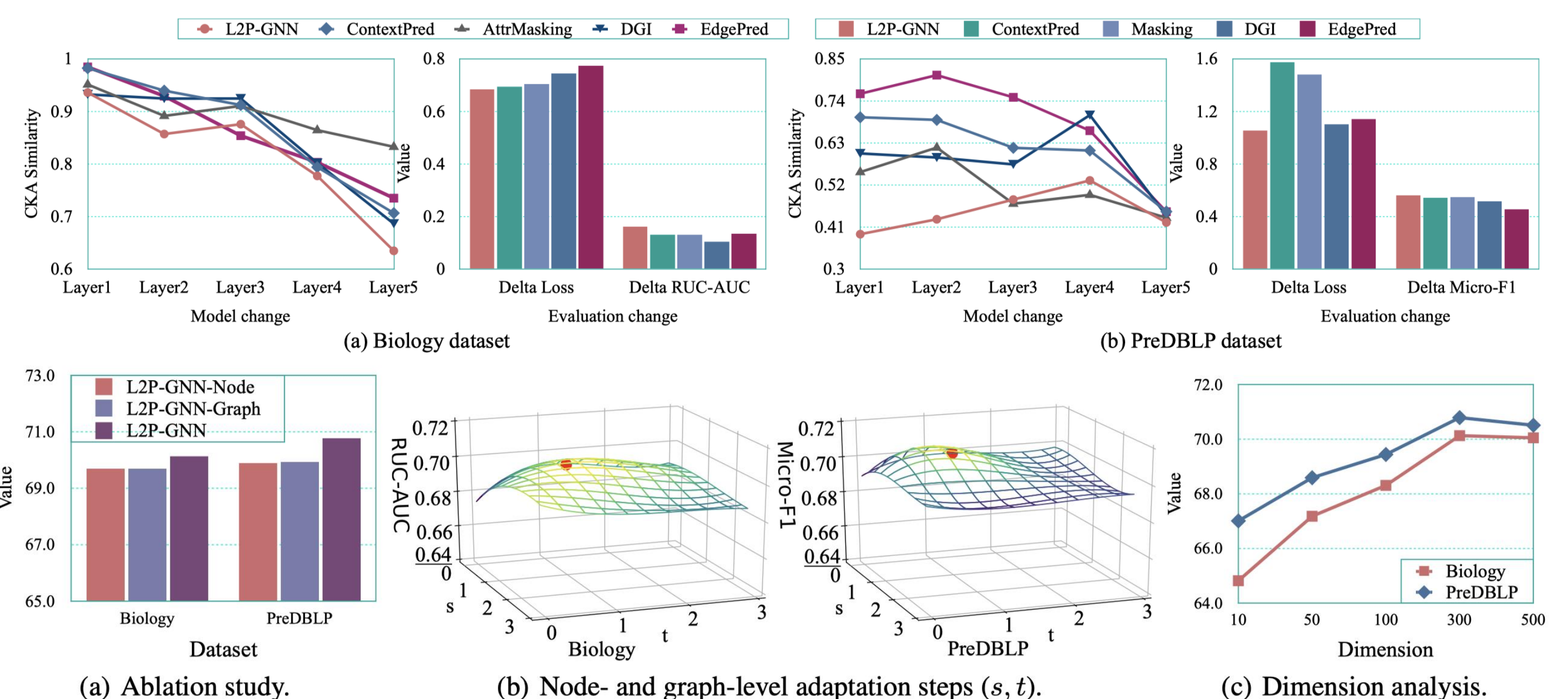
Smaller change, more easily achieve the optimal point

Ablation Study

- L2P-GNN-Node with only node-level adaptation
- L2P-GNN-Graph with only graph-level adaptation

Parameter Analysis

- the number of node- and graph-level adaptation steps (s, t)
- the dimension of node representations



Conclusions

- Problem:** There exists a divergence between the pre-training and fine-tuning objectives, resulting in suboptimal pre-trained GNN models
- Solution:** A self-supervised pretraining strategy for GNNs, L2P-GNN, which attempts to learn how to fine-tune in the pre-training process in the form of transferable prior knowledge
- Dataset:** A new large-scale graph structured data for pre-training GNNs