

The background of the slide features a repeating pattern of stylized, light pink flowers and leaves. The flowers have five petals and are arranged in a circular, wreath-like fashion around the central text.

Combining Graph Convolutional Neural Networks and Label Propagation

TOIS2021

Hongwei Wang and Jure Leskovec
Reporter: Minjie Cheng

March 8, 2023

Overview

- ▶ Problem Formulation
- ▶ GCNs and LPA
- ▶ GCN-LPA(an end-to-end model that combines GCN and LPA)
- ▶ Experiments
- ▶ Conclusion

Problem Formulation

Consider a graph $G = (V, \mathbf{A}, \mathbf{X}, \mathbf{Y})$

- ▶ Node set $V = \{v_1, \dots, v_n\}$
- ▶ Adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$
- ▶ Node feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$
- ▶ Label set $L = \{1, \dots, c\}$
- ▶ The first m nodes have labels
- ▶ Node label matrix $\mathbf{Y} \in \mathbb{R}^{n \times c}$, $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$ are one-hot label indicator vectors.
- ▶ Learn a mapping $\mathcal{M} : V \rightarrow L$ and predict labels of unlabeled $\{\mathbf{y}_{m+1}, \dots, \mathbf{y}_n\}$.

GCNs and LPA

- ▶ GCN propagates and transforms **node features** across the graph.
- ▶ LPA propagates **labels** along the edges.

The feature propagation scheme of GCN in layer k is as follows:

$$\mathbf{X}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{X}^{(k)}\mathbf{W}^{(k)}) \quad (1)$$

- ▶ $\tilde{\mathbf{A}}$ is the normalized adjacency matrix.
- ▶ $\mathbf{X}^{(k)}$ is the k th layer node representations.
- ▶ $\mathbf{W}^{(K)}$ is trainable matrix in the layer k .
- ▶ $\sigma(\cdot)$ is an activation function.

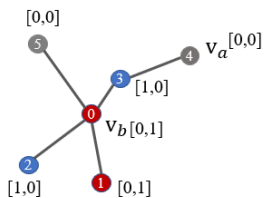
The feature propagation scheme of GCN in layer k is as follows:

$$\mathbf{Y}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{Y}^{(k)}) \quad (2)$$

$$\mathbf{y}^{(k+1)} = \mathbf{y}_i^{(0)}, \forall i \leq m \quad (3)$$

- Labels of all labels of nodes are reset to their initial values.

LPA-An Example



$$Y$$

0	1	0
1	1	0
2	0	1
3	0	1
4	0	0
5	0	0

A

	0	1	2	3	4	5
0	0	1	1	1	0	1
1	1	0	0	0	0	0
2	1	0	0	0	0	0
3	1	0	0	0	1	0
4	0	0	0	1	0	0
5	1	0	0	0	0	0

D

	0	1	2	3	4	5
0	4	0	0	0	0	0
1	0	1	0	0	0	0
2	0	0	1	0	0	0
3	0	0	0	2	0	0
4	0	0	0	0	1	0
5	0	0	0	0	0	1

$$\tilde{\mathbf{A}} = \mathbf{D}^{-1}\mathbf{A}$$

	0	1	2	3	4	5
0	0	1/4	1/4	1/4	0	1/4
1	1/4	0	0	0	0	0
2	1/4	0	0	0	0	0
3	1/4	0	0	0	1/4	0
4	0	0	0	1/4	0	0
5	1/4	0	0	0	0	0

$$Y^{(1)} = \tilde{A} Y^{(0)}$$

0	1/4	1/2
1	1/4	0
2	1/4	0
3	1/4	1/4
4	0	1/4
5	1/4	0



$$Y^{(1)} \quad y^{(k+1)} = y_i^{(0)}, \forall i \leq m$$

0	1	0
1	1	0
2	0	1
3	0	1
4	0	1/4
5	1/4	0

GCNs and LPA

The relationship between the two equations(GCNs and LPA)

► GCNs:

$$\mathbf{X}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{X}^{(k)}\mathbf{W}^{(k)}) \quad (4)$$

► LPA:

$$\mathbf{Y}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{Y}^{(k)}) \quad (5)$$

$$\mathbf{y}^{(k+1)} = \mathbf{y}_i^{(0)}, \forall i \leq m \quad (6)$$

GCNs and LPA

The relationship between the two equations(GCNs and LPA)

► GCNs:

$$\mathbf{X}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{X}^{(k)}\mathbf{W}^{(k)}) \quad (7)$$

► LPA:

$$\mathbf{Y}^{(k+1)} = \sigma(\tilde{\mathbf{A}}\mathbf{Y}^{(k)}) \quad (8)$$

$$\mathbf{y}^{(k+1)} = \mathbf{y}_i^{(0)}, \forall i \leq m \quad (9)$$

Two definitions

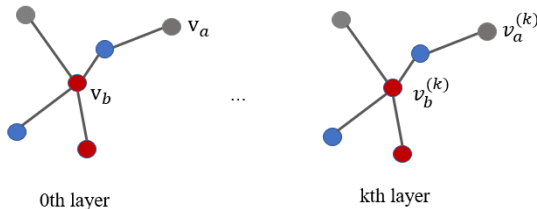
► Feature influence:

- The feature influence of node v_b on node v_a

$$I_f(v_a, v_b; k) = \|\mathbb{E}_{\mathbf{w}(\cdot)} \frac{\partial \mathbf{x}_a^{(k)}}{\partial \mathbf{x}_b}\|_1 \quad (10)$$

- The normalized feature influence is defined as

$$\tilde{I}_f(v_a, v_b; k) = \frac{I_f(v_a, v_b; k)}{\sum_{v_i \in V} I_f(v_a, v_i; k)} \quad (11)$$



Two definitions

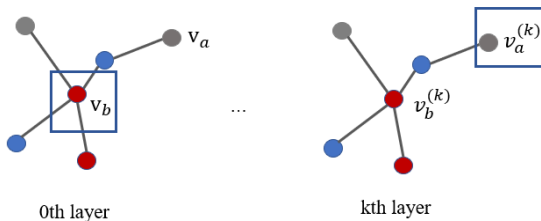
► Feature influence:

- The feature influence of node v_b on node v_a

$$I_f(v_a, v_b; k) = \|\mathbb{E}_{\mathbf{w}(\cdot)} \frac{\partial \mathbf{x}_a^{(k)}}{\partial \mathbf{x}_b}\|_1 \quad (12)$$

- The normalized feature influence is defined as

$$\tilde{I}_f(v_a, v_b; k) = \frac{I_f(v_a, v_b; k)}{\sum_{v_i \in V} \tilde{I}_f(v_a, v_i; k)} \quad (13)$$

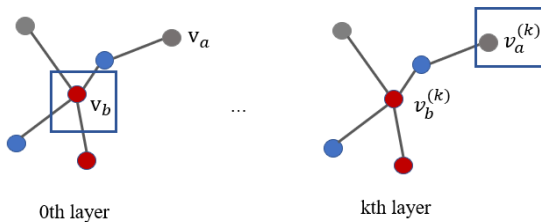


Two definitions

► Label influence:

- The label influence of labeled node v_b on **unlabeled node** v_a

$$I_l(v_a, v_b; k) = \frac{\partial \mathbf{y}_a^{(k)}}{\partial \mathbf{y}_b} \quad (14)$$



Theorem1

$$\mathbb{E}_{\mathbf{w}(\cdot)} [I_l(v_a, v_b; k)] = \sum_{j=1}^k \beta^j I_f(v_a, v_b; j)$$

- ▶ β is the fraction of unlabeled nodes.
- ▶ The labels influence of v_b after k iterations of LPA equals to the cumulative normalized feature influence of v_b on v_a after k layers of GCN.

- We defined **intra-class** feature influence:

$$\sum_{i \in L} \sum_{v_a, v_b: y_a=i, y_b=i} I_f(v_a, v_b) \quad (15)$$

- We can optimize the **intra-class** label influence instead:

$$\sum_{i \in L} \sum_{v_a, v_b: y_a=i, y_b=i} I_l(v_a, v_b) \quad (16)$$

$$\sum_{i \in L} \sum_{v_a, v_b: y_a=i, y_b=i} I_l(v_a, v_b) = \sum_{v_a} \sum_{v_b: y_b=y_a} I_l(v_a, v_b) \quad (17)$$

Theorem2

$$\sum_{v_b: y_b=y_a} I_l(v_a, v_b; k) \propto Pr(\hat{y}_a^{lpc} = y_a)$$

- \hat{y}_a^{LPA} is the predicted label of v_a using a k-iteration LPA.

According to **Theorem2**, We can therefore first learn the optimal edge weights \mathbf{A}^* by minimizing the loss of predicted labels by LPA.

$$\mathbf{A}^* = \arg \min_{\mathbf{A}} L_{lpa}(\mathbf{A}) = \arg \min_{\mathbf{A}} \frac{1}{m} \sum_{v_a: a \leq m} J(\hat{\mathbf{y}}_a^{lpa}, \mathbf{y}_a) \quad (18)$$

- ▶ J is the cross-entry loss.
- ▶ The optimal \mathbf{A}^* maximizes the probability that each node is correctly labeled by LPA, also maximizes the **intra-class label** influence (according to Theorem2) and **intra-class feature** influence (according to Theorem1).

GCNs-LPA

- We can plug \mathbf{A}^* into GCN to predict labels:

$$\mathbf{X}^{(k+1)} = \sigma(\mathbf{A}^* \mathbf{X}^{(k)} \mathbf{W}^{(k)}), k = 0, 1, \dots, K-1 \quad (19)$$

- The optimal transformation matrices can be learned by minimizing the loss of predicted labels by GCN:

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} L_{gc n}(\mathbf{W}, \mathbf{A}^*) = \arg \min_{\mathbf{W}} \frac{1}{m} \sum_{v_a: a \leq m} J(\hat{\mathbf{y}}_a^{gc n}, \mathbf{y}_a) \quad (20)$$

- Combine the above LPA and GCN and train the whole model in an end-to-end fashion:

$$\mathbf{W}^*, \mathbf{A}^* = \arg \min_{\mathbf{W}, \mathbf{A}} L_{gc n}(\mathbf{W}, \mathbf{A}) + \lambda L_{lp a}(\mathbf{A}) \quad (21)$$

GCNs-LPA

$$\mathbf{A}^* = \mathbf{M} \circ \mathbf{A}$$

- ▶ Each element $M_{i,j}$ can be a free variable during training.
 - ▶ It can only work in **transductive** setting while the trained model cannot be used for new graphs.
 - ▶ For large graphs, the number of **model parameters** increases linearly with the number of edges.
- ▶ Each element M_{ij} can be set as a function of features of two endpoints e.g.,
 $M_{ij} = k(\mathbf{x}_i^T \mathbf{H} \mathbf{x}_j)$
 - ▶ The learned \mathbf{H} can be applied to new graphs in inductive setting.
 - ▶ The model size is also independent with the graph size.

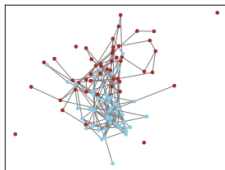
GCNs-LPA

$$\mathbf{A}^* = \mathbf{M} \circ \mathbf{A}$$

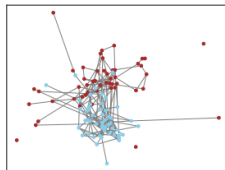
- ▶ Each element M_{ij} can be a free variable during training. **Semi-Supervised Node Classification**
 - ▶ It can only work in **transductive** setting while the trained model cannot be used for new graphs.
 - ▶ For large graphs, the number of **model parameters** increases linearly with the number of edges.
- ▶ Each element M_{ij} can be set as a function of features of two endpoints e.g., $M_{ij} = k(\mathbf{x}_i^T \mathbf{H} \mathbf{x}_j)$. **Knowledge-Graph-Aware Recommendation**
 - ▶ The learned \mathbf{H} can be applied to new graphs in inductive setting.
 - ▶ The model size is also independent with the graph size.

Experiments

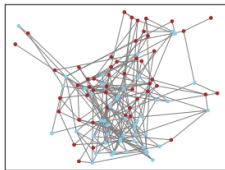
GCN-LPA is more robust to inter-class edges than GCNs



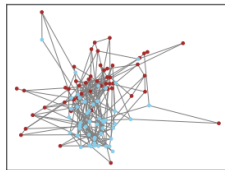
(a) GCN on the original network



(b) GCN-LPA on the original network



(c) GCN on the noisy network



(d) GCN-LPA on the noisy network

Fig. 2. Node embeddings of a subgraph of Cora trained on a node classification task (red vs. blue). Node coordinates in Figures 2(a)–2(d) are the embedding coordinates. Notice that GCN does not produce linearly separable embeddings (Figure 2(a) vs. Figure 2(b)), while GCN-LPA performs much better even in the presence of noisy edges (Figure 2(c) vs. Figure 2(d)). Additional visualizations are included in Appendix D.

Semi- Supervised Node Classification

Three variants of GCN-TPA

- ▶ **GCN-LPA(T)** learns the optimal adjacency matrix \mathbf{A} by minimizing L_{lpa} first, then freezes \mathbf{A} and optimizes \mathbf{W} by minimizing L_{gcn} .
- ▶ **GCN-LPA(S)** optimizes \mathbf{A} and \mathbf{W} , but the gradient of \mathbf{A} only propagates back from L_{lpa} .
- ▶ **GCN+LPA** simply adds predictions of GCN and LPA together.

Semi- Supervised Node Classification

Table 3. Mean and the 95% Confidence Intervals of Test Set Accuracy for All Methods and Datasets in Node Classification Task

Method	Cora	Citeseer	Pubmed	Coauthor-CS	Coauthor-Phy
Logistic Regression	57.0 \pm 1.6	61.2 \pm 1.5	64.0 \pm 2.3	86.1 \pm 0.8	86.7 \pm 1.3
LPA	74.4 \pm 2.0	67.5 \pm 1.3	70.6 \pm 3.4	73.8 \pm 1.4	86.2 \pm 1.5
GCN	81.4 \pm 0.8	<u>71.8</u> \pm 1.3	<u>77.6</u> \pm 2.0	90.9 \pm 0.5	92.3 \pm 0.9
GAT	80.8 \pm 1.2	71.3 \pm 1.4	76.8 \pm 1.6	90.4 \pm 0.5	92.2 \pm 0.7
JK-Net	81.1 \pm 1.2	70.3 \pm 1.0	<u>77.4</u> \pm 0.5	90.3 \pm 0.3	90.8 \pm 0.6
GIN	74.2 \pm 0.9	60.5 \pm 1.1	73.3 \pm 1.0	84.2 \pm 1.1	87.1 \pm 0.9
GDC	83.2 \pm 0.8	<u>72.0</u> \pm 0.9	<u>77.8</u> \pm 0.7	91.1 \pm 0.6	92.1 \pm 0.4
GCN-LPA(T)	<u>82.6</u> \pm 0.7	<u>72.4</u> \pm 0.7	<u>78.2</u> \pm 1.0	<u>91.5</u> \pm 0.5	<u>93.1</u> \pm 0.7
GCN-LPA(S)	<u>82.9</u> \pm 0.8	<u>72.4</u> \pm 0.9	<u>78.4</u> \pm 0.9	91.8 \pm 0.4	<u>93.4</u> \pm 0.8
GCN+LPA	78.3 \pm 0.5	69.9 \pm 1.1	74.0 \pm 0.6	84.4 \pm 0.8	89.9 \pm 0.7
GCN-LPA	<u>83.1</u> \pm 0.7	72.6 \pm 0.8	78.6 \pm 1.3	91.8 \pm 0.4	93.6 \pm 1.0

The best result is highlighted in bold, while the result falling within the confidence interval of the highest one is highlighted with underline.

Efficacy of LPA regularization

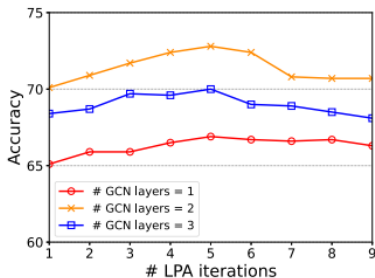


Fig. 3. Sensitivity to the number of LPA iterations on Citeseer dataset.

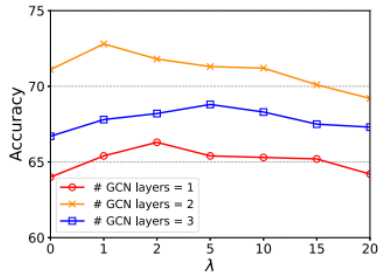


Fig. 4. Sensitivity to λ (weight of LPA loss) on Citeseer dataset.

- Figure 4 shows that training without the LPA loss term ($\lambda=0$) is more difficult than the case where $\lambda = 1 - 5$

Influence of labeled node rate

Table 4. Accuracy of LPA, GCN, and GCN-LPA on Citeseer Dataset with Different Labeled Node Rate

Labeled node rate	2%	5%	10%	20%	40%	60%
LPA	65.0 ± 1.1	67.7 ± 1.5	68.2 ± 1.0	70.5 ± 1.4	71.6 ± 1.3	73.9 ± 0.8
GCN	69.6 ± 0.9	72.1 ± 1.2	72.4 ± 1.0	74.4 ± 0.6	75.7 ± 0.5	79.1 ± 0.8
GCN-LPA	70.2 ± 0.7	72.7 ± 1.0	73.3 ± 0.8	75.3 ± 1.1	77.3 ± 0.8	80.9 ± 1.0

- GCN-LPA outperforms GCN and LPA consistently, the improvement achieved by GCN-LPA increases when labeled node rate is larger.

Visualization of learned edge weights

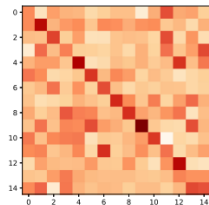


Fig. 5. Visualization of learned edge weights on Coauthor-CS dataset. The numbers along axes denote different node labels.

- ▶ It is clear that values along the diagonal(**intra-class** edges weights) are significantly larger than off-diagonal values(**inter-class** edge weights).
- ▶ GCN-LPA is able to identify the importance of edges and distinguish inter-class and intra-class edges.

Knowledge-Graph-Aware Recommendation

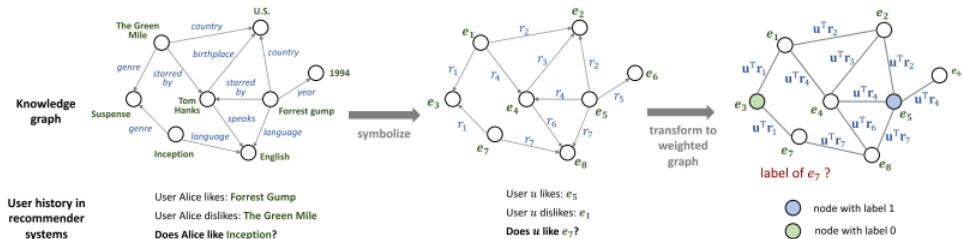


Fig. 7. Left: An instance of movie knowledge graph and the rating history of user Alice in recommender systems. Middle: The symbolized knowledge graph and user history. Right: Transforming the knowledge graph into a weighted graph, in which the weight of edge with relation r is $u^T r$, where u and r are the trainable embeddings for user u and relation r , respectively. Nodes that user u likes are assigned with label 1 and nodes that user u dislikes are assigned with label 0. In this way, the KG-aware recommendation problem is transformed to a binary node classification problem.

- We apply GCN-LPA to KG-aware recommendation problem.

Knowledge-Graph-Aware Recommendation

Table 8. The Results of *Recall@K* for Book-Crossing Dataset in Recommendation Task

Model	R@2	R@10	R@50	R@100
SVD	2.5 ± 0.3	4.6 ± 0.5	7.8 ± 0.4	10.9 ± 0.4
LibFM	3.1 ± 0.5	6.1 ± 0.5	9.2 ± 0.5	12.5 ± 0.6
LibFM + TransE	3.7 ± 1.3	6.5 ± 1.2	9.7 ± 1.2	13.0 ± 1.3
PER	2.1 ± 0.8	4.0 ± 1.0	6.4 ± 1.1	7.0 ± 1.0
CKE	2.6 ± 0.4	5.2 ± 0.4	7.8 ± 0.5	11.2 ± 0.6
RippleNet	3.5 ± 0.6	7.4 ± 0.6	10.7 ± 0.7	12.8 ± 0.7
KGIN	3.3 ± 0.4	7.0 ± 0.5	10.2 ± 0.7	12.6 ± 0.5
GCN-LPA	4.5 ± 0.4	8.3 ± 0.5	11.6 ± 0.4	14.9 ± 0.4

- GCN-LPA outperforms baselines by a significant margin.

Conclusion

- ▶ We study the theoretical relationship between GCNs and LPA.
- ▶ We propose a unified model, which learns transformation matrices and edge weights simultaneously on GCN with the assistance of **LPA regularizer**.

The connection to my work (OTGNN)

$$\min_{\mathbf{q}, \theta} \|\mathbf{W}_e \mathbf{q}\|_1 + \lambda \|\mathbf{Y}_L - (\text{GNN}_L(\mathbf{X}, \mathbf{q}, \theta) + \mathbf{J}_L \mathbf{q})\|_2^2 \quad (22)$$

- ▶ We design the **update strategy** for edge weights based on the labeled nodes from the OT on graph perspective.