LGD-GCN: Local and Global Disentangled Graph Convolutional Networks

Jingwei Guo^{1,2}

Kaizhu Huang¹

Xinping Yi²

Rui Zhang¹

¹ Xi'an Jiaotong-Liverpool University, Suzhou, China ² University of Liverpool, Liverpool, UK Jingwei.Guo@liverpool.ac.uk, Kaizhu.Huang@xjtlu.edu.cn Xinping.Yi@liverpool.ac.uk, Rui.Zhang02@xjtlu.edu.cn

Abstract

Disentangled Graph Convolutional Network (DisenGCN) is an encouraging framework to disentangle the latent factors arising in a real-world graph. However, it relies on disentangling information heavily from a local range (i.e., a node and its 1-hop neighbors), while the local information in many cases can be uneven and incomplete, hindering the interpretabiliy power and model performance of DisenGCN. In this paper^a, we introduce a novel Local and Global Disentangled Graph Convolutional Network (LGD-GCN) to capture both local and global information for graph disentanglement. LGD-GCN performs a statistical mixture modeling to derive a factor-aware latent continuous space, and then constructs different structures w.r.t. different factors from the revealed space. In this way, the global factor-specific information can be efficiently and selectively encoded via a message passing along these built structures, strengthening the intra-factor consistency. We also propose a novel diversity promoting regularizer employed with the latent space modeling, to encourage interfactor diversity. Evaluations of the proposed LGD-GCN on the synthetic and real-world datasets show a better interpretability and improved performance in node classification over the existing competitive models.

^aThis paper is a lighter version of "Learning Disentangled Graph Convolutional Networks Locally and Globally" where the results and analysis have been reworked substantially.

1 INTRODUCTION

Graphs are emerging as an insightful structured data modeling technique for capturing the similarity between data samples and the relationship between entities. To mine the domain-specific knowledge in graph structured data, Graph Convolutional Networks (GCNs) have been proposed to integrate the topological patterns and content features [Kipf and Welling, 2017], demonstrating excellent expressive power and growing popularity in various graph learning tasks, such as node classification, link prediction, and recommendation [Wu et al., 2020, Chen and Wong, 2020].

Most state-of-the-art methods, such as [Kipf and Welling, 2017, Hamilton et al., 2017, Velickovic et al., 2018], study node representations in a holistic approach, i.e., they interpret the node neighborhood as a whole without considering the within-distinctions. By contrast, a real-world graph typically contains multiple heterogeneous node relations which in many cases are implicitly determined by various latent factors shaping node aspects. For instance, a user in a social network, usually links with different persons for different reasons, such as family, work, and/or hobby, which potentially characterize the user from different perspectives. The existing holistic approaches usually fail to disentangle these latent factors, rendering the learned representations hardly explained and less informative.

Recently, Disentangled Graph Convolutional Network (DisenGCN) [Ma et al., 2019] offers a promising framework to disentangle the latent factors behind graph data via a neighborhood partition. Despite the novel design, Disen-GCN heavily relies on local node neighborhood, which may bring unexpected issues. First, the information from local ranges can be significantly varied across the entire graph. Solely depending on it, DisenGCN could easily produce latent representations losing consistent meaning of the associated factor. That may weaken the intra-factor correlation between disentangled features and leads to diminished interpretability. Second, the local neighborhood information can be scarce and limited especially considering sparse graphs, prohibiting DisenGCN from generating informative node aspects and yielding favourable performance boost. A detailed discussion can be seen later in Section 2.3.

To tackle this limitation, in this paper, we propose a novel framework, termed as Local and Global Disentangled Graph Convolutional Network (LGD-GCN), to learn disentangled node representations capturing both local and global graph information. The central idea is that disentangling the latent factors inherent in a graph can benefit from a latent continuous space which uncovers the underlying factor-aware node relations. Specifically, we first leverage the neighborhood routing mechanism to locally disentangle node representations into multiple latent units pertinent to different factors. Then, we propose to guide the disentanglement from a global perspective.

To this end, our approach performs a mixture statistical modeling over the locally disentangled latent units, to derive a factor-aware latent continuous space. This enables a different component or mode, specific to a latent factor, in a different region of the latent space [Ghahramani and Hinton, 1996]. After that, we manage to build a different structure by connecting near neighbors in a different region of the revealed latent space. These latent structures disclose the underlying factor-aware relations between nodes. Employing message passing along them can efficiently and selectively encode the global factor-specific information, which enhances intra-factor consistency, i.e., the consistent meaning of disentangled latent units w.r.t. the associated factor. Furthermore, we also design a novel diversity promoting regularizer to encourage inter-factor diversity. Practically, it enforces the disentangled latent units related to different factors to fall into separate clusters in the latent space so as to enhance the disentangled informativeness. In sharp contrast to DisenGCN, Fig. 1 clearly visualizes the benefit of learning disentangled node representations both locally and globally. Our contributions are summarized as below:

- We argue that DisenGCN may bring unexpected issues by heavily relying on local graph information. Empirical analysis shows that DisenGCN may learn latent representations with weakly disentangled factors, and especially its boost performance becomes minor while facing sparse graphs.
- We propose a novel Local and Global Disentangled framework for Graph Convolutional Networks (LGD-GCN) to infer the latent factors underlying the graph data. Incorporating both local and global information, LGD-GCN can disentangle node representations with enhanced intra-factor consistency and promoted interfactor diversity.
- Extensive evaluations on synthetic and real-world datasets demonstrate that LGD-GCN provides a better interpretability and improved performance in node classification compared to other state-of-the-arts.

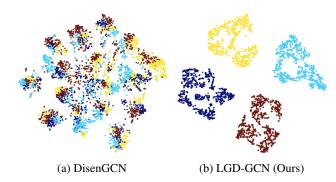


Figure 1: Visualization of the disentangled latent units w.r.t. four latent factors on a synthetic graph. Points with a different color mean the disentangled latent units (for all nodes) of a different latent factor. In sharp contrast to DisenGCN, our LGD-GCN displays a highly disentangled pattern with strong intra-factor consistency and inter-factor diversity; it indicates high (low) intra-factor (inter-factor) correlations between features

2 BACKGROUND AND MOTIVATION

2.1 CONVENTIONAL GNNS

Graph neural networks (GNNs) are powerful machine learning models in dealing with graph-structured data, where the input data are modeled as graphs. A graph is denoted as G=(V,E) with V being the set of nodes and E being the set of edges. Given two distinct nodes $u,v\in V$ with $u\neq v$, we define $(u,v)\in E$ if u and v are connected with an edge, and the neighborhood of node u as $N_u=\{v|(u,v)\in E\}$. For attributed graphs, each node u usually has an initial representation $\mathbf{h}_u^{(0)}\in \mathbb{R}^{d_0}$. GNNs are used to mine the underlying relationship between nodes according to their attributes $\{\mathbf{h}_u^{(0)}|\forall u\in V\}$ for (semi-)supervised learning tasks such as node classification and link prediction.

In the past years, an increasing number of GNN models have been proposed [Wu et al., 2020]. Most of them can be generalized by a message passing mechanism [Gilmer et al., 2017], where the node attributes are exchanged through the graph edges following a neighborhood aggregation strategy descibed below.

$$\mathbf{h}_u \leftarrow \mathbf{UPDATE}(\mathbf{AGGREGATE}(\{\mathbf{h}_v | \forall v \in N_u\}), \mathbf{h}_u),$$

where the **AGGREGATE** operation is to aggregate information from a node neighborhood, and the **UPDATE** operation is to combine these information to update the node's attributes. Such a strategy works in an iterative way to learn node representations. For graph-level representations, a readout operation, such as a simple mean or sum, can be applied to summarize the overall information.

2.2 DISENTANGLED NODE REPRESENTATION

Albeit promising several learning tasks, most GNNs treat the neighborhood as a whole and ignore the inner-differences, learning noninterpretable representations. To address this issue, DisenGCN [Ma et al., 2019] was proposed by hypothesizing nodes are connected due mainly to different kinds of relationship; and there are M inherent factors determining edge connections and potentially shaping nodes from M aspects.

DisenGCN aims to disentangle each node representation into multiple latent units w.r.t. different latent factors. In each layer, given a node u and its neighborhood N_u , the node representations, $\{\mathbf{h}_i | \forall i \in \{u\} \cup N_u\}$, will be first projected onto M subspaces using different channels. In each channel m (m = 1, 2, ..., M), the projected representation of node i is given by

$$\mathbf{z}_{i,m} = \frac{\sigma(\mathbf{W}_m^T \mathbf{h}_i + \mathbf{b}_m)}{\|\sigma(\mathbf{W}_m^T \mathbf{h}_i + \mathbf{b}_m)\|_2}$$
(1)

where $\mathbf{W}_m \in \mathbb{R}^{d_{im} \times \frac{d_{out}}{M}}$ and $\mathbf{b}_m \in \mathbb{R}^{\frac{d_{out}}{M}}$ are learnable parameters, and σ is an activation function. A neighborhood routing mechanism, detailed in [Ma et al., 2019, Algorithm-1], then iteratively partitions all the neighbors into different clusters. After that, independent information aggregations are applied over them in different channels, to attain the disentangled latent units for node u, $\{\hat{\mathbf{z}}_{u,m} \in \mathbb{R}^{\frac{d_{out}}{M}} | \forall m=1,2,...,M \}$. Finally, the disentangled node representation is obtained by concatenation, $\hat{\mathbf{h}}_u = \hat{\mathbf{z}}_{u,1} \oplus \hat{\mathbf{z}}_{u,2} \oplus ... \oplus \hat{\mathbf{z}}_{u,M}$ and $\hat{\mathbf{h}}_u \in \mathbb{R}^{d_{out}}$.

2.3 LIMITATIONS OF DISENGEN

While DisenGCN reveals certain latent factors, we argue that it tends to produce weakly disentangled representations and yield limited performance boost because of its heavy reliance on local graph information. To validate our argument, we further provide two experimental investigations over the graph synthesized with four latent factors (see details in Section 4.1).

First, we visualize the disentangled latent units of Disen-GCN using t-SNE in Fig. 1a. At the micro-level, we can observe the separability between points with different colors in some regions. But, when it comes to the macro-level, all points unexpectedly fall into discrete clusters and mixed together, indicating a weak disentanglement. This is because the disentangled latent units by DisenGCN may preserve some specific micro-meanings of the factor, but losing the consistent macro-meaning (*intra-factor consistency*). Additionally, DisenGCN only considers disentangling representations in different channels without ensuring the diversity between those w.r.t. different factors (*inter-factor diversity*). The learned representations thereby are prone to preserve the redundant information, partially explaining Fig. 1a.

Second, we further augment the synthetic graph by tuning

Table 1: Micro (Top) and Macro (Bottom) F1 scores (%) on graphs synthesized with four latent factors but different average neighborhood sizes

Methods	Average Neighborhood Sizes							
Methods	40	30	20	10	6			
GCN	79.5±0.8	75.5±0.7	66.1±0.9	47.2±0.6	37.2±0.9			
DisenGCN	84.1 ± 1.0	79.5 ± 0.7	69.0 ± 1.0	48.8 ± 0.9	$38.4 {\pm} 0.8$			
Improvements	+4.6%	+4.0%	+2.9%	+1.6%	+1.2%			
GCN	78.3±0.9	75.0±0.8	65.8±1.0	45.8±0.6	36.7±0.8			
DisenGCN	82.9 ± 1.1	78.9 ± 0.7	68.3 ± 1.0	47.4 ± 1.0	37.7 ± 0.8			
Improvements	+4.6%	+3.9%	+2.5%	+1.6%	+1.0%			

the *p* value (it controls the density of the synthetic graph as described in Section 4.1) to generate graphs with multiple average neighborhood sizes. We then apply GCN [Kipf and Welling, 2017] and DisenGCN to train for multi-label classification, and report the F1 scores in Table 1. From the table, the relative improvements reduce from approximately 5% to 1% as the average neighborhood size decreasing from 40 to 6. The result meets the expectation. DisenGCN may perform well on a dense graph by learning disentangled representations. However, sparsing the input graph (limiting the accessible local information) can negatively affect the performance boost, which reflects the heavy local reliance of DisenGCN.

3 LOCAL AND GLOBAL DISENTANGLED GCN

We present an novel method for Graph Convolutional Networks (LGD-GCN) to disentangle node representations both locally and globally, as presented in Fig. 2. By hiring the neighborhood routing mechanism [Ma et al., 2019], we first attain disentangled latent units preserving local graph information w.r.t. different latent factors. However, these disentangled units are prone to be weakly disentangled without incorporating global information and being properly regularized. In the following, we show how to enhance the disentanglement from a global perspective, via disclosing the underling factor-aware relations between nodes, to learn a better disentangled representations with strengthened intrafactor consistency and promoted inter-factor diversity.

3.1 MODELING LATENT CONTINUOUS SPACE

Extending the hypothesis in DisenGCN from local neighborhood to global graph, we assume that the locally disentangled units for all nodes, $\{\hat{\mathbf{z}}_{i,m} \in \mathbb{R}^{\frac{d_{out}}{M}} | \forall i \in V, \forall m = 1,2,...,M\}$, are generated from a gaussian mixture distribution with equal mixture weights, s.t.

$$p(\hat{\mathbf{z}}_{i,m}) = \frac{1}{M} \sum_{e=1}^{M} \mathcal{N}(\hat{\mathbf{z}}_{i,m}; \boldsymbol{\mu}_e, \boldsymbol{\Sigma}_e), \tag{2}$$

where $\mu_e \in \mathbb{R}^{\frac{d_{out}}{M}}$ and $\Sigma_e \in \mathbb{R}^{\frac{d_{out}}{M} \times \frac{d_{out}}{M}}$ are the mean and the covariance associated with latent factor e. Then, we

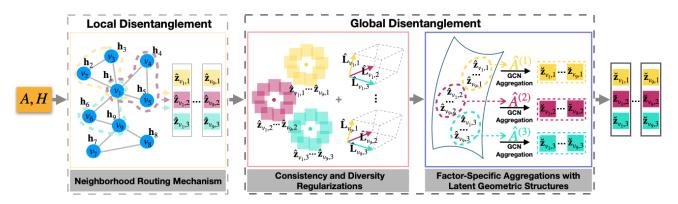


Figure 2: Illustrative example of the LGD-GCN layer with M=3 latent factors. First, the node representations are locally disentangled by leveraging the neighborhood routing mechanism. These disentangled representations are then modeled in a latent continuous space, and promoted with consistent and diverse latent factors globally, from which geometric structures are constructed for further aggregation.

employ this assumption to learn factor-specific means and covariances to regularize the disentangling of the latent units. Specifically, we maximize the conditional likelihood of the latent units $\hat{\mathbf{z}}_{i,m}$ (for each node i and each factor m) w.r.t. the associated factor m. It is equivalent to minimizing the negative log term expressed in Eq. (3) after removing constants.

$$\mathcal{L}_{i,m} = (\hat{\mathbf{z}}_{i,m} - \boldsymbol{\mu}_m)^T \boldsymbol{\Sigma}_m^{-1} (\hat{\mathbf{z}}_{i,m} - \boldsymbol{\mu}_m)$$
(3)

Minimizing the term $\mathcal{L}_{i,m}$ is equivalent to minimizing the Mahalanobis Distance [De Maesschalck et al., 2000] between the disentangled latent unit and its globally inferred center. It derives a latent continuous space where the disentangled latent units are encouraged to be more discriminative with respect to their centers, and to carry the type of global factor-specific information shared by all nodes.

3.2 CONSTRUCTING LATENT STRUCTURES

Although node relations are naturally presented in graph data, we believe that they are not always optimal for disentangled graph learning. Taking a huge and sparse graph as an example. It is difficult for most nodes to absorb sufficient information, coming from a small number of their neighbors (one or two in most cases), to learn disentangled representations w.r.t. latent factors in a larger number. On the other hand, the raw graph may not contain the desired topologies after projecting node features in different channels. As such, disclosing the underlying factor-aware relations between nodes from the disentangled latent space becomes a promising alternative.

The previously modeled latent space enables a different component or mode, specific to a different latent factor, in a different region. It would be reasonable to apply a proper graph construction algorithm over different regions to obtain latent structures specific to different factors. We expect these built structures uncovering the factor-aware relations between

nodes, and selecting a sufficient number of latent neighbors (from the entire graph) for each node in shaping node aspects. Accordingly, the global factor-specific information can be efficiently and selectively encoded, by employing a simple message passing scheme independently along these different structures, to strengthen the *intra-factor consistency*.

Here, we list two popular methods for building graphs from data using local neighborhood in latent space:

1) **k-Nearest-Neighbor** (**kNN**): It connects every point to its k^{th} nearest neighbors, given a pairwise distance $d(\mathbf{z}_i, \mathbf{z}_j)$. Formally, the adjacency matrix $\mathbf{A}^{kNN} \in \{0,1\}^{N \times N}$ is defined as:

$$\mathbf{A}_{i,j}^{knn} = \begin{cases} 1 & d(\mathbf{z}_i, \mathbf{z}_j) \leq d(\mathbf{z}_i, \mathbf{z}_j^{(k)}) \text{ or } d(\mathbf{z}_i, \mathbf{z}_j) \leq d(\mathbf{z}_i^{(k)}, \mathbf{z}_j) \\ 0 & \text{otherwise} \end{cases}$$

where $\mathbf{z}_i^{(k)}$ and $\mathbf{z}_j^{(k)}$ denote the \mathbf{k}^{th} nearest neighbors of \mathbf{z}_i and \mathbf{z}_j , respectively.

2) **Continuous k-Nearest-Neighbor** (**CkNN**) [Berry and Sauer, 2016]: It is a discrete version of kNN for removing kNN's sensitivity to the density parameter k. Similarly, the adjacency matrix $\mathbf{A}^{CkNN} \in \{0,1\}^{N \times N}$ is defined as:

$$\mathbf{A}_{i,j}^{cknn} = \begin{cases} 1 & d(\mathbf{z}_i, \mathbf{z}_j) < \sqrt{d(\mathbf{z}_i, \mathbf{z}_i^{(k)}) d(\mathbf{z}_j, \mathbf{z}_j^{(k)})} \\ 0 & \text{otherwise} \end{cases}$$

In this paper, we apply the same message passing function in GCN [Kipf and Welling, 2017] to aggregate the factor-specific node information along these constructed structures, following Eq. (4).

$$\mathbf{\breve{Z}}^{(m)} = \mathbf{\breve{D}}^{(m)^{-\frac{1}{2}}} \mathbf{\breve{A}}^{(m)} \mathbf{\breve{D}}^{(m)^{-\frac{1}{2}}} \mathbf{\hat{Z}}^{(m)}$$
(4)

Here, $\hat{\mathbf{A}}^{(m)}$ refers to the built structures w.r.t. latent factor m from $\{\hat{\mathbf{z}}_{i,m}|\forall i\in V\}$, $\check{\mathbf{A}}^{(m)}=\hat{\mathbf{A}}^{(m)}+\mathbf{I}$, $\check{\mathbf{D}}_{i,i}^{(m)}=\sum_{j}\check{\mathbf{A}}_{i,j}^{(m)}$, $\check{\mathbf{D}}_{i,j}^{(m)}=0$ in case of $i\neq j$, and $\hat{\mathbf{Z}}^{(m)}$ is the feature matrix

with each column being $\hat{\mathbf{z}}_{i,m}$ for node i in V. Particularly, we adopt the Euclidean distance as the pairwise distance d(,), and denote this proposed module as \mathbb{LG}_{agg} .

3.3 PROMOTING INTER-FACTOR DIVERSITY

Diversity-promoting learning aims to encourage different components in latent space models to be mutually uncorrelated and different, and has been widely studied [Xie, 2018, Xie et al., 2016]. In the previous sections, we derived a factor-aware latent continuous space and built structures for encoding factor-specific node information from a global range, to strengthen the *intra-factor consistency*. However, without being regularized to be different with respect to different latent factors, the disentangled latent units may preserve redundant information of other irrelevant latent factors.

In this paper, we propose to promote the *inter-factor diver*sity to capture the unique information in disentangled latent units. Particularly, we define the diversity on the conditional likelihoods (given different factors) of each disentangled latent unit in latent space. Inspired by the Determinant Point Process [Kulesza and Taskar, 2012], we formulate the disentanglement diversity for each node *i* as

$$\mathbb{DD}_i = \det(\hat{\mathbf{F}}_i^T \hat{\mathbf{F}}_i), \tag{5}$$

where $\hat{\mathbf{F}}_i = \langle \hat{\mathbf{L}}_{i,1}, ..., \hat{\mathbf{L}}_{i,m}, ..., \hat{\mathbf{L}}_{i,M} \rangle$, $\hat{\mathbf{L}}_{i,m} = \|\mathbf{L}_{i,m}\|_2$, and $\mathbf{L}_{i,m} = \langle \mathcal{N}(\hat{\mathbf{z}}_{i,m}; \mu_1, \Sigma_1), ..., \mathcal{N}(\hat{\mathbf{z}}_{i,m}; \mu_M, \Sigma_M) \rangle$ contains the conditional likelihoods (given M factors) of the disentangled latent unit $\hat{\mathbf{z}}_{i,m}$.

By the property of Determinant [Bernstein, 2005], \mathbb{DD}_i is equal to the volume spanned by $\{\hat{\mathbf{L}}_{i,m}|\forall m=1,2,...,M\}$, elegantly providing an intuitive geometric interpretation as shown in Fig. 2 with M=3. Maximizing \mathbb{DD}_i encourages $\mathbf{L}_{i,1}, \mathbf{L}_{i,2},..., \mathbf{L}_{i,M}$ to be orthogonal to each other, i.e., enforcing the disentangled latent units to fall into separated regions of the statistical latent space; it essentially enhances the disentangled informativeness and promotes the *inter-factor diversity*.

3.4 NETWORK ARCHITECTURE

In this section, we detail the general network architecture of the proposed LGD-GCN for performing node-level tasks. The pseudocode of a LGD-GCN's layer is presented in Algorithm 1, and it is desirable to stack multiple LGD-GCN's layers to sufficiently exploit the graph data.

Specifically, we adopt ReLU activation function in Eq. (1) and append a dropout layer [Srivastava et al., 2014] in the end of each LGD-GCN's layer which is only enabled in training. We can then have the output of layer l as $\{\check{\mathbf{h}}_i^{(l)}|\forall i\in V\}$ = Dropout $(F^{(l)}(\{\check{\mathbf{h}}_i^{(l-1)}|\forall i\in V\}))$, where $1\leq l\leq L$,

In this work, we focus on the task of node classification. To incorporate Eq. (3) and Eq. (5) into the final optimization problem, we leverage them into two regularization terms for each node i, as expressed below.

$$\mathcal{L}_{space}^{i} = \frac{1}{M} \sum_{m=1}^{M} \mathcal{L}_{i,m}, \ \mathcal{L}_{div}^{i} = -\log(\mathbb{DD}_{i})$$
 (6)

For $\mathcal{L}_{i,m}$, we update μ_m and Σ_m (m=1,2,...,M) in each layer iteratively with the newly computed values after each training epoch by an update rate U_r . To adaptively modify the influential power of these two regularization terms in different layers, we apply a layer loss weight, $\lambda^{(l)} = 10^{l-L}$. It makes the influence of the regularization terms grows bigger as the layer goes deeper within a proper range. Then, we can formulate the final loss in Eq. (7) with coefficients λ_{space} and λ_{div} for trade-off.

$$\mathcal{L}_{total} = \mathcal{L}_{cls} + \sum_{l=1}^{L} \lambda^{(l)} (\lambda_{space} \mathcal{L}_{space}^{(l)} + \lambda_{div} \mathcal{L}_{div}^{(l)})$$
 (7)

Here, for single-label node classification, $\mathcal{L}_{cls} = -\frac{1}{|V|} \sum_{i}^{V} \mathbb{Y}_{i}^{T} \log(\operatorname{softmax}(\mathbf{Y}_{i}^{(L+1)}))$, and for multi-label classification, $\mathcal{L}_{cls} = -\frac{1}{|V|} \sum_{i}^{V} \mathbb{Y}_{i}^{T} \log(\operatorname{sigmoid}(\mathbf{Y}_{i}^{(L+1)})) + (1 - \mathbb{Y}_{i})^{T} \log(1 - \operatorname{sigmoid}(\mathbf{Y}_{i}^{(L+1)}))$, given $\mathbb{Y}_{i} \in \mathbb{R}^{C}$ being the ground truth label of node i in one hot encoding. The end-to-end optimization procedures are displayed in the supplemental material.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTING

Datasets Cora, Citeseer, and Pubmed are three citation benchmark networks widely used in [Kipf and Welling, 2017, Hamilton et al., 2017, Velickovic et al., 2018], where nodes and edges denote documents and undirected citations respectively; each node is assigned with one topic and associated with bags-of-words features. We synthesize graphs with latent factors following [Ma et al., 2019]. In detail, we first generate m Erdős-Rényi random graphs with 1,000 nodes and 16 classes, where nodes connect each other with probability p if they are in the same class, with probability q otherwise. Then, we merge these generated graphs by summing the adjacency matrix and turning the element-value bigger than zero to one, to obtain the final synthetic graphs

Algorithm 1: LGD-GCN's Layer

Input: $\{\mathbf{h}_i \in \mathbb{R}^{d_m} | \forall i \in V\}$; M: the number of latent factors; T: the routing iterations' number of the Neighborhood Routing Mechanism (NRM);

Parameter:

$$\mathbf{W}_{m} \in \mathbb{R}^{d_{im} \times \frac{d_{out}}{M}}, \mathbf{b}_{m} \in \mathbb{R}^{\frac{d_{out}}{M}}, \boldsymbol{\mu}_{m} \in \mathbb{R}^{\frac{d_{out}}{M}}, \boldsymbol{\Sigma}_{m} \in \mathbb{R}^{\frac{d_{out}}{M} \times \frac{d_{out}}{M}}, \\ \forall m = 1, 2, ..., M$$
 for $i \in V$ do
$$\mid \mathbf{z}_{i,1}, \mathbf{z}_{i,2}, ..., \mathbf{z}_{i,M} \leftarrow \mathbf{h}_{i} \text{ by Eq. (1)}.$$

end

for $i \in V$ do $\begin{vmatrix}
\hat{\mathbf{z}}_{i,m} \leftarrow \{\mathbf{z}_{i,m}\} \cup \{\mathbf{z}_{v,m} | \forall v \in N_i\}, \forall m = 1, 2, ..., M \text{ by NRM} \\
\text{with } T \text{ routing iterations.} \\
\text{Minimize } \mathcal{L}_{space}^i \text{ and } \mathcal{L}_{div}^i \text{ by Eq. (6).}
\end{vmatrix}$

end

for m = 1, 2, ..., M do

Construct structure $\mathbf{G}^{(m)}$ with $\mathbf{A}^{(m)}$ from $\{\hat{\mathbf{z}}_{i,m}|\forall i\in V\}$. $\{\check{\mathbf{z}}_{i,m}|\forall i\in V\}\leftarrow \{\hat{\mathbf{z}}_{i,m}|\forall i\in V\}$ by Eq. (4).

end

Output: $\{ \breve{\mathbf{z}}_{i,1} \oplus \breve{\mathbf{z}}_{i,2} \oplus \cdots \oplus \breve{\mathbf{z}}_{i,M} | \forall i \in V \}$

Table 2: Semi-supervised classification accuracies (%)

Method	Splits	Datasets			
Methou	Spirts	Cora	Citeseer	Pubmed	
MLP		51.5±1.0	46.5	71.4	
MoNet		82.2 ± 0.7	70.0 ± 0.6	77.7 ± 0.6	
GCN	Standard	81.9 ± 0.8	69.5 ± 0.9	79.0 ± 0.5	
GAT		82.5 ± 0.5	71.0 ± 0.6	77.0 ± 1.3	
DisenGCN		83.7	73.4	80.5	
LGD-GCN (ours)		84.9 ± 0.4	74.5 ± 0.8	81.3 ± 0.6	
MLP		58.2 ± 2.1	59.1±2.3	70.0±2.1	
MoNet		81.3 ± 1.3	71.2 ± 2.0	78.6 ± 2.3	
GCN	Random	81.5 ± 1.3	71.9 ± 1.9	77.8 ± 2.9	
GAT	Kandom	81.8 ± 1.3	71.4 ± 1.9	78.7 ± 2.3	
DisenGCN		81.4 ± 1.6	69.5 ± 1.4	79.1 ± 2.3	
LGD-GCN (ours)		84.0 ±1.3	72.0 ±1.3	79.8 ±2.3	

with m latent factors. We set q to $3e^{-5}$ following [Ma et al., 2019], and tune p value such that the average neighborhood size is between 39.5 and 40.5. Each node is initialized with the row of the adjacency matrix as the features and has m labels. The data statistics are listed in the supplementary material.

Baseline Models. We compare our model with several methods, including the state-of-the-art, as the baselines: MLP is a multi-layer perception; MoNet [Monti et al., 2017] is a mixture model CNN generalizing convolutional neural network to non-Euclidean graph data structure; GCN [Kipf and Welling, 2017] approximates graph Laplacian with Chebyshev expainsion; GAT [Velickovic et al., 2018] combines the attention mechanism with graph neural networks to aggregate information with selective neighbors; DisenGCN attempts to learn disentangled node representations via a neighborhood routing mechanism.

Hyper-parameters. We set $d_{out} = 64$ as the output dimension of each LGD-GCN's hidden layer and T = 7 as the number of routing iterations, to follow GAT [Velick-ovic et al., 2018] and DisenGCN respectively. For semi-

Table 3: Micro-F1 (Top) and Macro-F1 (Bottom) scores (%) on synthetic graphs with different number of latent factors

	Number of Latent Factors					
Method	4	6	8	10	12	
MLP	79.3 ± 0.5	55.5±0.4	37.0 ± 0.8	25.9 ± 0.6	21.2 ± 0.8	
GCN	74.5 ± 0.8	56.3 ± 0.7	38.2 ± 0.9	28.0 ± 0.7	23.1 ± 0.8	
DisenGCN	84.1 ± 1.0	60.4 ± 0.9	41.4 ± 1.3	29.4 ± 0.7	24.2 ± 0.8	
LGD-GCN (ours)	87.2 ± 0.5	$65.0 {\pm} 0.5$	43.6 ± 0.7	$30.2 {\pm} 0.5$	26.1 ± 0.5	
MLP	77.9 ± 0.7	54.8±0.6	36.0 ± 0.8	24.5 ± 0.7	20.1±0.9	
GCN	78.3 ± 0.9	55.6 ± 0.9	37.2 ± 1.0	26.9 ± 0.5	22.2 ± 0.9	
DisenGCN	82.9 ± 1.1	59.9 ± 1.0	40.2 ± 1.2	28.1 ± 0.7	23.4 ± 0.7	
LGD-GCN (ours)	86.1 ± 0.5	64.2 ± 0.6	$42.5 {\pm} 0.6$	$28.8{\pm}0.5$	$25.1 {\pm} 0.5$	

supervised node classification on real-world datasets, we fix the number of channels M as 4 for simplification. We use dropout $\sim [0,1]$, learning rate $\sim [3e-3,1]$, weight decay $\sim [5e-5,0.2]$, update rate $\sim [0.1,0.9]$ for μ_k and Σ_k , and the number of layers $\sim \{1,2,...,10\}$. For multi-label classification on the synthetic datasets, with a slight difference, we fix dropout as 0.5, learning rate $\sim [5e-4,5e-3]$, weight decay $\sim [1e-3,1e-2]$, and $M \sim \{2,4,...,16\}$.

Additionally, the regularization coefficients λ_{space} and λ_{div} as well as the density parameter k are empirically searched from different ranges for different datasets as provided in the supplementary material. Then, we carefully tune the hyper-parameters defined above on the validation set using optuna [Akiba et al., 2019]. With the best hyper-parameters, we train the model in 1,000 epochs using the early-stopping strategy with a patience of 100 epochs, and report the average performance in 10 runs on the test split.

4.2 QUANTITATIVE EVALUATION

In this section, we evaluate our model quantitatively in tasks of semi-supervised node classification and multi-label node classification.

Semi-supervised Node Classification. In this task, we follow the experimental protocal established by Kipf and Welling [2017], Velickovic et al. [2018], and consider both standard split [Yang et al., 2016] and random split. For random split, we uniformly sample the same number of instances as in the standard split in 10 times.

The results are listed in Table 2 measured in classification accuracy. Since Shchur et al. [2018] have conducted extensive evaluations in their work, we will quote their reported results for baseline methods. For DisenGCN, we not only collect their results, but also optimize and evaluate the model on the random splits using their source codes. For our model, considering the non-linear complexity of the real-world datasets, we adopt CkNN [Berry and Sauer, 2016] in the module \mathbb{LG}_{agg} .

From the results, the proposed LGD-GCN consistently outperforms other baselines. Especially, our model is able to improve upon DisenGCN by a margin of 1.2% and 2.6% on Cora in standard and random splits, respectively. This

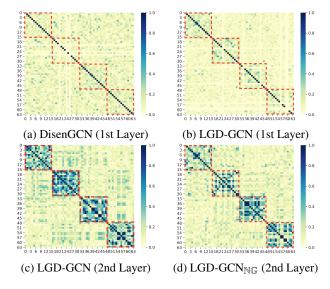


Figure 3: Features correlation analysis

demonstrates the benefits brought by absorbing rich and diverse global information. More importantly, real graphs are typically highly sparse as observed in Cora, Citeseer and Pubmed whose graphs contain an average neighbor number of 3.9, 2.8 and 4.5 for each node. In this case, our model is more effective in capturing long-range dependencies via the created shortcuts in the built geometric structures, which further explains the performance improvement.

Multi-label Node Classification. To further demonstrate our model's disentangling ability quantitatively, we apply MLP, GCN, DisenGCN, and our model to train graphs synthesized with various number of latent factors for multi-label node classification. Specifically, we randomly split each synthetic dataset into train/validation/test as 0.6/0.2/0.2, adopt kNN in the module \mathbb{LG}_{agg} , measure model performance in Micro-F1 and Macro-F1 scores, and report the results in Table 3. It can be observed that our model consistently outperforms others while varying the number of latent factors, and especially achieves significant performance gains by (micro-f1) 4.6% and (macro-f1) 4.3% upon DisenGCN on the graph synthesized with six latent factors.

4.3 QUALITATIVE EVALUATION

The qualitative evaluation focuses on disentanglement performance and informativeness of learned embeddings.

Visualization of disentangled representations. We give in Fig. 1b a 2D visualization of the learned representations w.r.t. four latent factors on the synthetic graph. Compared to that of DisenGCN in Fig. 1a, our model displays a highly disentangled pattern with consistent and diverse latent factors, evidenced by the intra-factor compactness and inter-factor separability; it also indicates the nodes carry the common type of factor-specific global information.

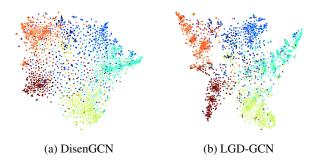


Figure 4: Visualization of node embedding on Citeseer

Correlation of disentangled features. The correlation analysis of the latent features, learned by DisenGCN and our model on test split of the graph synthesized with four latent factors, is presented in Fig. 3. As observed, our model showcases a more block-wise correlation pattern, which becomes denser in the second layer, indicating the enhanced interpretability. We also analyze the feature correlation of our model while ablating the module \mathbb{LG}_{agg} , denoted as $\mathrm{LGD}\text{-}\mathrm{GCN}_{\mathbb{NG}}$ in Fig. 3d. Though the block-wise pattern in Fig. 3d can still be observed, it is obviously weaker than that of $\mathrm{LGD}\text{-}\mathrm{GCN}$ in Fig. 3c. This verifies the significance of \mathbb{LG}_{agg} ; the captured factor-specific global information strengthens the factor-aware feature correlation, and enhances the interpretability of the learned representations.

Visualization of node embeddings. Fig. 4 provides a intuitive comparison between the learned node embeddings of DisenGCN and our model on Citeseer dataset. It can be observed that the proposed LGD-GCN learns better node embeddings and shows a high inter-class similarity and intraclass difference. This is because our model learns more informative node aspects by absorbing rich factor-specific global information, leading to increasing discriminative power.

4.4 PARAMETER AND ABLATION ANALYSIS

We investigate the sensitivity of hyper-parameters, and perform ablation analysis over the proposed modules on realworld and synthetic datasets.

Analysis of consistency coefficient λ_{space} . We plot the learning performance of our model w/o \mathcal{L}_{div} while varying λ_{space} in Eq. (7) e.g. from 0 to 5 on Cora in Fig. 5a. The accuracy goes up first and drops slowly. Practically, promising performance can be attained on Cora by choosing λ_{space} from [0.1, 1].

Analysis of diversity coefficient λ_{div} . We then test the effect of λ_{div} in Eq. (7), and vary it from e.g. 0 to 0.5 on Citeseer. λ_{div} is relatively robust within a certain range e.g. [0, 0.1] for Citeseer in Fig. 6b. Once out of that range, the results drops to a low point, suggesting overly focusing on diversity is harmful to model performance.

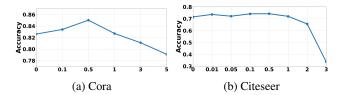


Figure 5: Analysis of parameter λ_{space}

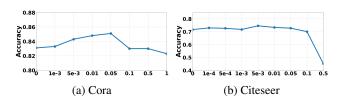


Figure 6: Analysis of parameter λ_{div}

Analysis of density parameter k. Fig. 7 displays the impact of k from 1 to 12 on Cora and Citeseer. The results are relatively stable while selecting k from a wide range, e.g. 1 to 8 on Cora and 1 to 10 on Citeseer. However, as k getting larger, the accuracy performance deteriorates obviously. It probably because larger k may introduce noisy edges, leading to inappropriate information sharing.

Analysis of the number of channels M. We study the influence of the number of channels M on the synthetic graphs generated with eight latent factors. From Fig. 8, our model performs the best when the number of channels is around eight, the true number of the latent factors.

Ablation analysis. We validate the contributions of the proposed modules denoted by \mathcal{L}_{space} , \mathcal{L}_{div} , and \mathbb{LG}_{agg} in node classification. From Table 4, we can see that both modules can independently and jointly improve the accuracy.

Table 4: Ablation analysis in node classification accuracies.

Components		Cora	Citeseer	Pubmed
-		81.9 ± 1.1	70.4 ± 1.6	78.9 ± 0.7
\mathscr{L}_{space}		83.0 ± 0.5	72.4 ± 1.3	79.1 ± 0.8
\mathcal{L}_{space} + \mathcal{L}_{div}	%	83.6 ± 0.6	72.4 ± 1.1	79.1 ± 0.4
\mathscr{L}_{space} + $\mathbb{L}\mathbb{G}_{agg}$	%	84.4 ± 0.3	74.0 ± 0.7	81.3 ±0.6
\mathcal{L}_{space} + \mathcal{L}_{div} + $\mathbb{L}\mathbb{G}_{agg}$		84.9 ± 0.4	74.5 \pm 0.8	81.2 ± 0.7

5 CONCLUSION

In this paper, we propose a novel framework, termed Local and Global Disentangled Graph Convolutional Network (LGD-GCN), to disentangle node representations with strengthened intra-factor consistency and promoted interfactor diversity. Extensive experiments demonstrate the improved performance in node classification and enhanced interpretability of the proposed LGD-GCN over existing state-of-the-art approaches.

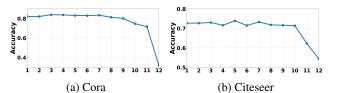


Figure 7: Analysis of parameter k

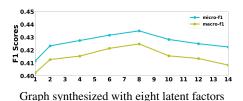


Figure 8: Analysis of parameter *M*

References

Takuya Akiba, Shotaro Sano, T. Yanase, Takeru Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019.

D. Bernstein. Matrix mathematics: Theory, facts, and formulas with application to linear systems theory. 2005.

Tyrus Berry and Timothy Sauer. Consistent manifold representation for topological data analysis. *arXiv* preprint *arXiv*:1606.02353, 2016.

Tianwen Chen and Raymond Chi-Wing Wong. Handling information loss of graph neural networks for session-based recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1172–1180, 2020.

Roy De Maesschalck, Delphine Jouan-Rimbaud, and Désiré L Massart. The mahalanobis distance. *Chemometrics and intelligent laboratory systems*, 50(1):1–18, 2000.

Zoubin Ghahramani and Geoffrey E. Hinton. The em algorithm for mixtures of factor analyzers. 1996.

Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272. PMLR, 2017.

William L. Hamilton, Zhitao Ying, and J. Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.

- Thomas Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *ArXiv*, abs/1609.02907, 2017.
- A. Kulesza and B. Taskar. Determinantal point processes for machine learning. *Found. Trends Mach. Learn.*, 5: 123–286, 2012.
- Jianxin Ma, P. Cui, Kun Kuang, X. Wang, and Wenwu Zhu. Disentangled graph convolutional networks. In *ICML*, 2019.
- Federico Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. 2017 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5425–5434, 2017.
- O. Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *ArXiv*, abs/1811.05868, 2018.
- N. Srivastava, Geoffrey E. Hinton, A. Krizhevsky, Ilya Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 15:1929–1958, 2014.
- Petar Velickovic, Guillem Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. *ArXiv*, abs/1710.10903, 2018.
- Zonghan Wu, Shirui Pan, Fengwen Chen, G. Long, C. Zhang, and P. Yu. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 2020.
- B. Xie, Yingyu Liang, and L. Song. Diversity leads to generalization in neural networks. *ArXiv*, abs/1611.03131, 2016.
- Pengtao Xie. Diversity-promoting and large-scale machine learning for healthcare. 2018.
- Z. Yang, W. Cohen, and R. Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. *ArXiv*, abs/1603.08861, 2016.

A SUPPLEMENTARY MATERIAL

In the supplementary material, for reproducibility, we provide the dataset information, algorithm and optimization procedure, and hyper-parameters' searching ranges. Finally, we show more visualization results to validate our conclusion in the manuscript.

A.1 DATASET STATISTICS

We list the information of datasets evaluated in the manuscript in Table 5 and Table 6. For real-world datasets, we only use 20 labeled nodes per class but with all the rest nodes unlabeled for training, another 500 nodes for validation and early-stopping, and 1,000 nodes from the rest for testing. For synthetic datasets, we specify the parameters including probability p and probability q for generating synthetic graphs with various number of latent factors.

A.2 ALGORITHM AND OPTIMIZATION

Algorithm 2 illustrates the optimization procedures in pseudo-codes.

A.3 ADDITIONAL RESULTS

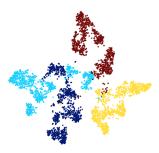
To further verify the importance of the module \mathbb{LG}_{agg} , we ablate it from our model and visualize the disentangled representations on the synthetic graph with four latent factors in Fig. 9. Comparing to that of the original model, we can witness an evident performance drop by the weakened intrafactor compactness. Even worse, the integrated blue set in Fig. 1b is broken into two disjoint clusters in Fig. 9 by turning off the module \mathbb{LG}_{agg} , indicating its effectiveness. We also show the visualization of node embedding learned on Cora dataset in Fig. 10. Similar to that on Citeseer dataset in Fig. 4, our model learns better embeddings, evidenced by intra-class compactness and inter-class separability.

Table 5: Real-world Dataset Statistics

Dataset	Cora	Citeseer	Pubmed
Nodes	2708	3327	19717
Avg-Neighbors	3.9	2.8	4.5
Features	1433	3703	500
Classes	7	6	3
Train	140	120	60
Validation	500	500	500
Test	1000	1000	1000

Table 6: Parameters for generating synthetic datasets

# Latent Factors	4	6	8	10	12
probablity p	0.164	0.110	0.082	0.065	0.055
probability q	3e-5	3e-5	3e-5	3e-5	3e-5



LGD-GCN w/o \mathbb{LG}_{agg}

Figure 9: Visualization of disentangled representations on a synthetic graph with four latent factors

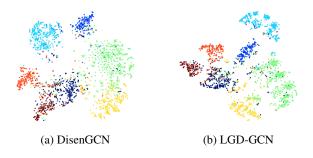


Figure 10: Visualization of node embedding on Cora

Algorithm 2: LGD-GCN's Optimization Procedure

Input: { $\mathbf{h}_i^{(0)} \in \mathbb{R}^{d_{in}} | \forall i \in V$ }; l_r be the learning rate; U_r be the update rate for $\mu_m^{(l)} \in \mathbb{R}^{\frac{d_{out}}{M}}$ and $\Sigma_m^{(l)} \in \mathbb{R}^{\frac{d_{out}}{M} \times \frac{d_{out}}{M}}$, $\forall m \in \{1, 2, ..., M\}, \forall l = 1, 2, ..., L$, \mathbf{F}_{Θ} refers to the proposed LGD-GCN with learnable weights Θ.

for number of training epochs do

end

Compute \mathcal{L}_{total} by Eq. (7); Update Θ using Adam optimizer [Kingma and Ba, 2015] with learning rate l_r .