

# Dimensionwise Separable 2-D Graph Convolution for Unsupervised and Semi-Supervised Learning on Graphs

Qimai Li\*  
csqqli@comp.polyu.edu.hk  
The Hong Kong Polytechnic  
University

Xiaotong Zhang\*  
zhangxt@dlut.edu.cn  
School of Software  
Dalian University of Technology

Han Liu\*  
hanliu@dlut.edu.cn  
School of Software  
Dalian University of Technology

Quanyu Dai  
daiquanyu@huawei.com  
Huawei Noah's Ark Lab

Xiao-Ming Wu†  
xiao-ming.wu@polyu.edu.hk  
The Hong Kong Polytechnic  
University

## ABSTRACT

Graph convolutional neural networks (GCN) have been the model of choice for graph representation learning, which is mainly due to the effective design of graph convolution that computes the representation of a node by aggregating those of its neighbors. However, existing GCN variants commonly use 1-D graph convolution that solely operates on the object link graph without exploring informative relational information among object attributes. This significantly limits their modeling capability and may lead to inferior performance on noisy and sparse real-world networks. In this paper, we explore 2-D graph convolution to jointly model object links and attribute relations for graph representation learning. Specifically, we propose a computationally efficient dimensionwise separable 2-D graph convolution (DSGC) for filtering node features. Theoretically, we show that DSGC can reduce intra-class variance of node features on both the object dimension and the attribute dimension to learn more effective representations. Empirically, we demonstrate that by modeling attribute relations, DSGC achieves significant performance gain over state-of-the-art methods for node classification and clustering on a variety of real-world networks. The source code for reproducing the experimental results is available at <https://github.com/liqimai/DSGC>.

## CCS CONCEPTS

• Computing methodologies → Artificial intelligence.

## KEYWORDS

2-D graph convolution, node classification, node clustering, variance reduction

\*These authors contributed equally.

†Corresponding Author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

KDD '21, August 14–18, 2021, Virtual Event, Singapore

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-8332-5/21/08...\$15.00  
<https://doi.org/10.1145/3447548.3467413>

## ACM Reference Format:

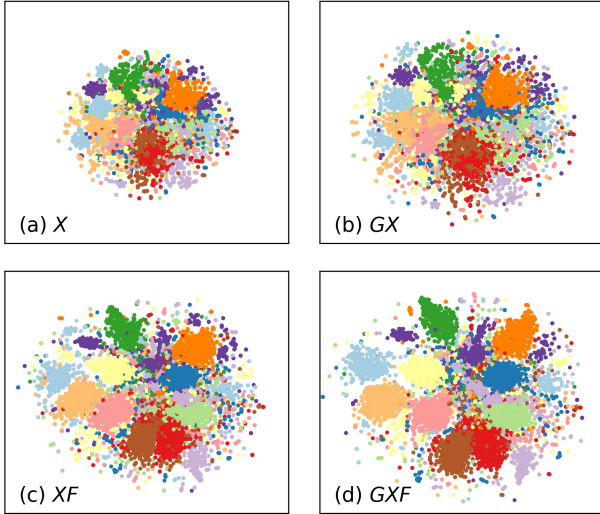
Qimai Li, Xiaotong Zhang, Han Liu, Quanyu Dai, and Xiao-Ming Wu. 2021. Dimensionwise Separable 2-D Graph Convolution for Unsupervised and Semi-Supervised Learning on Graphs. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3447548.3467413>

## 1 INTRODUCTION

Graph-structured data occurs in citation networks [58], social networks [60], traffic networks[55], protein networks[20], knowledge graphs [37, 51] and many other application fields. Analyzing graph data often leads to new insights and interesting discoveries. For example, unsupervised learning on graphs enables to detect user communities in social networks, while semi-supervised learning on graphs helps to predict protein functions in protein networks. Learning effective node representations by utilizing graph structures and other aspects of information such as node content has proved very useful for unsupervised and semi-supervised learning on graphs.

Previous approaches for unsupervised learning have explored applying non-negative matrix factorization, random walk statistics, and Laplacian eigenmaps on both graph structures and node attributes [23, 36] to learn node representations. For semi-supervised learning, a common way is to regularize a supervised classifier trained with node features by a Laplacian regularizer or an embedding-based regularizer to take into account graph structures, e.g., manifold regularization [5], manifold denoising [21], deep semi-supervised embedding [54], and Planetoid [59]. However, these methods model node connectivity and node content separately and hence may not be able to fully utilize these information.

In the past few years, graph convolutional neural networks (GCN) and variants have dominated the research of graph representation learning and achieved new state-of-the-art results in various learning tasks on graphs, especially in unsupervised learning [49, 63] and semi-supervised learning [20, 25, 29, 46]. The success is mainly attributed to graph convolution, a function that naturally combines node connectivity and node content for feature propagation and smoothing, which computes the representation of a node by aggregating the features of its neighbors. The effective utilization of both modalities of data gives the unique advantage to GCN-based



**Figure 1: t-SNE visualization of the “20 Newsgroups” dataset. (a) Raw features; (b) Filtered by the regular object graph convolution; (c) Filtered by our proposed attribute graph convolution; (d) Filtered by our proposed dimensionwise separable graph convolution (DSGC).**

methods over previous approaches, including topology-only models [19, 41] and graph regularization based methods [5, 59].

In spite of their empirical success, one major limitation of most existing GCN-based methods is that they commonly adopt one-dimensional (1-D) graph convolution that operates on the *object link graph* to model node (object) relations and features, whose performance critically relies on the quality of the graph. However, real-life networks are often noisy and sparse. For example, in a web graph such as Wikipedia, a hyperlink between two webpages does not necessarily indicate that they belong to the same category, and mixing their features could be harmful for webpage classification or clustering. Moreover, it has been shown that many real-world networks are scale-free and there exist many low-degree nodes [3]. Since these nodes may have very few or even no links to other nodes, it is difficult, and even impossible, to do feature propagation to endow them with similar features as other same-class nodes to facilitate downstream learning tasks.

To address the above limitation, we propose to explore data relations in another dimension, by constructing an *attribute affinity graph*, which encodes relations between object attributes. The underlying assumption is attributes that indicate the same category should have strong relations. For example, in a citation network, object attributes are words, and documents of AI category usually contain words such as “learning”, “robotics”, “machine”, “neural”, etc. These indicative words for AI category should be more closely related than other non-indicative words. These informative relations can then be utilized for feature smoothing, similar to the use of object links. Importantly, the attribute affinity graph can be a useful complement to the object link graph in learning node representations. For instance, consider a document that has no links to others (e.g., node e or d in Figure 2), and hence it is impossible to do feature propagation with object links. Yet, with meaningful

attribute relations, feature propagation can still be performed along the attribute dimension, and it is possible to obtain similar features for same-class documents.

To formalize the above insight, we propose to perform graph convolution on the attribute affinity graph (Figure 1 (c)). Further, we develop an efficient two-dimensional (2-D) graph convolution to perform feature smoothing on both the object link graph and the attribute affinity graph to learn more effective node representations (Figure 1 (d)). Our main contributions are described as follows.

- **Methodology:** We propose to use 2-D graph convolution to jointly model object links, attribute relations, and object features for node representation learning. Furthermore, we develop a computationally efficient dimensionwise separable 2-D graph convolutional filter (DSGC), which is equivalent to performing 1-D graph convolution alternately on the object dimension and the attribute dimension, as illustrated in Figure 2. Finally, we propose two learning frameworks based on DSGC for unsupervised node clustering and semi-supervised node classification.
- **Theoretical insight:** We show that the regular 1-D graph convolution on the object link graph can reduce intra-class variance of node features, which helps to explain the success of many existing methods. Further, we show that the same can be proved for graph convolution on a properly constructed attribute affinity graph. Jointly, they provide a theoretical justification of DSGC.
- **Empirical study:** We implement DSGC for semi-supervised node classification and unsupervised node clustering, and conduct extensive experiment on a variety of real-world networks including email networks, citation networks, and web graphs. The comparison with state-of-the-art methods demonstrate the advantages of DSGC over the regular 1-D graph convolution. Moreover, we show that DSGC can be easily plugged into some strong GCN-based methods to further improve their performance substantially.

## 2 RELATED WORKS

In this work, we focus on the studies of graph convolutional neural networks for unsupervised and semi-supervised learning. There are a large number of prior works for unsupervised and semi-supervised learning on graph-structured data. We can only briefly review a small portion of them due to space limitation. We first review non-GCN-based methods, then review GCN-based methods.

**Network Embedding** emerges as an effective way for learning node representations in recent years by leveraging different techniques, such as Markov random walks [41], matrix factorization [58], autoencoders [7], and generative adversarial nets [12]. The learned node representations can boost various downstream learning tasks, such as node clustering [50] and node classification [41].

**Graph-based Semi-supervised Learning.** Early works in this line explicitly or implicitly adopt the assumption that nearby vertices are likely to have the same label. The most popular framework is label propagation [9, 56, 64, 65], which formulates a quadratic regularization framework to implement the cluster assumption and enforce the consistency with labeled data. Other ideas include using graph partition techniques to place the cuts in low density regions [6], using spectral kernels to learn smooth low-dimensional

embeddings [8, 62], and modified adsorption [45] and iterative classification algorithm [44]. It was shown in [14] and [17] that the graph regularization in many of these methods can be interpreted as low-pass graph filters. However, these methods are limited in their ability to incorporate vertex features for prediction. To jointly model graph structures and vertex features, many methods adopt an idea to regularize a supervised learner (e.g., support vector machines, neural networks) with some regularizer such as a Laplacian regularizer or an embedding-based regularizer. Examples include manifold regularization (LapSVM) [5], deep semi-supervised embedding [54], and Planetoid [59].

**Graph Convolutional Neural Networks.** In the past few years, a series of works based on graph convolutional neural networks [30, 42] have demonstrated promising performance in both semi-supervised and unsupervised learning.

**Semi-supervised learning.** The pioneering work ChebyNet [13] exploits a  $k$ -th order polynomial filter via Chebyshev expansion to avoid the expensive eigen-decomposition in spectral graph convolution. GCN [25] simplifies ChebyNet by designing an efficient layer-wise propagation rule with a first-order approximation, which inspires many follow-up works. Specifically, graph attention networks [46] and gated attention networks [61] introduce attention mechanisms to assign different weights to different nodes in a neighborhood. Some researches [28] provide deeper insights into GCN by showing the Laplacian smoothing nature and low-pass filtering nature of GCN-like models from spatial and spectral views, respectively. Instead of using two hops of neighbors, some later works such as JKNet [57], LanczosNet [31], MixHop [1], and GCNII [10] propose to utilize multiple hops of neighbors.

**Unsupervised Learning.** One kind of methods stack several graph convolutional layers and learn node representations based on feature/structure preserving, e.g., graph autoencoder and graph variational autoencoder [24], GraphSAGE [20], marginalized graph autoencoder [49]. Since these methods lack additional constraints to enhance the robustness of representations, adversarially regularized (variational) graph autoencoder [35] further adopts GAN to match the learned embeddings with a Gaussian prior. Another type of methods exploit contrastive learning to train an encoder to be contrastive between similar nodes and dissimilar ones for node representation learning, e.g., deep graph infomax [47] and graphical mutual information [38].

However, existing GCN-based methods mainly concentrate on 1-D graph convolution. To our knowledge, the only exception is [33], which explores 2-D graph convolution for matrix completion for recommendation. In this work, we investigate 2-D graph convolution for unsupervised and semi-supervised learning.

### 3 2-D GRAPH CONVOLUTION

Graph signal processing is an active research field in recent years, which generalizes basic concepts in harmonic analysis, including signals, filters, Fourier transform, and convolution, to the graph domain. Given a graph  $\mathcal{G}$  with a vertex set  $\mathcal{V}$  and an adjacency matrix  $A$ , if we associate each vertex  $v_i$  with a real value  $x_i$ , then a signal on  $\mathcal{G}$  can be defined as a vector  $\mathbf{x} = [x_i]$ . Graph filters are defined as mappings between input and output signals. Let  $G$  be a polynomial of  $A$  (usually normalized), i.e.  $G = p(A)$ , then  $G$  is a

legit convolutional filter on graph  $\mathcal{G}$ , and the corresponding 1-D graph convolution is defined as

$$\mathbf{z} = G\mathbf{x}, \quad (1)$$

where  $\mathbf{z}$  is the output signal. Existing graph convolutional filters are all defined in this way and have achieved considerable success in various learning tasks on graphs.

However, previous research mainly focuses on the design and application of 1-D graph convolution. In this section, we present 2-D graph convolution for learning node representations. A comprehensive introduction to multi-dimensional graph convolution is provided by [26]. Based on the theory developed by [26], we propose a localized 2-D graph convolution to circumvent the computationally intensive graph Fourier transform. Furthermore, we propose an even simpler dimensionwise separable 2-D graph convolution to efficiently model both object links and attribute relations.

#### 3.1 2-D Graph Signal and Spectral Convolution

**2-D Graph Signal.** A 2-D graph signal is a function defined on the Cartesian product of the vertex sets of two graphs. Formally, given two graphs  $\mathcal{G}^{(1)}$  and  $\mathcal{G}^{(2)}$  with  $n$  and  $m$  nodes respectively, and denote the vertex sets by  $\mathcal{V}^{(1)}$  and  $\mathcal{V}^{(2)}$ . A real-valued signal defined on them is a function  $x : \mathcal{V}^{(1)} \times \mathcal{V}^{(2)} \rightarrow \mathbb{R}$ . For convenience, we simply denote  $x(v_i^{(1)}, v_j^{(2)})$  by  $x_{ij}$  and organize then as a matrix:

$$\mathbf{X} = (x_{ij}) \in \mathbb{R}^{n \times m}, \quad x_{ij} = x(v_i^{(1)}, v_j^{(2)}), \quad (2)$$

Signal  $\mathbf{X}$  associates each node pairs  $(v_i, v_j) \in \mathcal{V}^{(1)} \times \mathcal{V}^{(2)}$  with a real number  $x_{ij}$ . For example, the feature matrix given by usual node classification tasks is a 2-D signal defined on object link graph and attribute affinity graph.

**2-D Graph Fourier Transform** Define the graph Laplacian of  $\mathcal{G}^{(1)}$  and  $\mathcal{G}^{(2)}$  as  $L_{(1)} = D^{(1)} - A^{(1)}$  and  $L_{(2)} = D^{(2)} - A^{(2)}$ , where  $A^{(1)}, A^{(2)}$  are adjacency matrices and  $D^{(1)}, D^{(2)}$  are the corresponding degree matrices.<sup>1</sup> Assuming two Laplacian matrices have following eigen-decomposition

$$L_{(1)} = U\Lambda U^{-1}, \quad L_{(2)} = V\mathbf{M}V^{-1}, \quad (3)$$

where  $\Lambda, \mathbf{M}$  stores the eigenvalues  $\lambda_i, \mu_j$  of two Laplacian matrices in their main diagonals,  $U = [\mathbf{u}_1, \dots, \mathbf{u}_n]$  and  $V = [\mathbf{v}_1, \dots, \mathbf{v}_m]$  store the corresponding unit-length eigen-vectors in their columns. This eigen-decomposition is always attainable for undirected graphs and nearly always attainable for directed graphs.

All  $n \times m$  outer products  $\mathbf{u}_i \mathbf{v}_j^\top$  together form a basis for the linear space  $\mathbb{R}^{n \times m}$ . It is known as 2-D graph Fourier basis – an analogy of the Fourier basis in classical harmonic analysis in graph domain. A 2-D graph signal  $\mathbf{X}$  can be decomposed in this basis with coefficients  $s_{ij}$ :

$$\mathbf{X} = \sum_{ij} s_{ij} (\mathbf{u}_i \mathbf{v}_j^\top), \quad (4)$$

or in matrix form:

$$\mathbf{X} = U\mathbf{S}V^\top, \quad \text{where } \mathbf{S} = (s_{ij}) \in \mathbb{R}^{n \times m}. \quad (5)$$

<sup>1</sup>Discussion below also applies to row-normalized Laplacian  $L_r = I - D^{-1}A$ , column-normalized Laplacian,  $L_c = I - AD^{-1}$ , and symmetric normalized Laplacian  $L_s = I - D^{-1/2}AD^{-1/2}$ .

$S$  is called the spectrum or Fourier coefficients of signal  $X$  and can be obtained by formula

$$S = U^{-1}X(V^{-1})^\top. \quad (6)$$

Eq. (6) is so-called 2-D graph Fourier transform; Eq. (5) is the corresponding inverse transform.

**2-D Spectral Graph Convolution** Given above decomposition, now we can manipulate the spectrum of 2-D signals and define 2-D spectral graph convolution. Convolution is an operation that takes a signal as input and outputs another signal. By convolution theorem, convolution is equivalent to scaling entries of the spectrum. Thus, given a signal  $X$  with spectrum  $S$ , a 2-D spectral graph convolution with  $X$  as input is defined as:

$$Z = U(S \circ P)V^\top, \quad (7)$$

where  $P$  is the spectral kernel (parameters in spectral domain) of this convolution, and ‘ $\circ$ ’ is Hadamard (entry-wise) product.

### 3.2 Fast Localized 2-D Graph Convolution

Although Eq. (7) well defines 2-D graph convolution, it is often impractical to perform convolution in the spectral domain, due to the high cost of computing the eigenbasis  $U, V$  needed for Fourier transform. Similar to what [13] did to 1-D graph convolution, we propose 2-D spatial graph convolution here to avoid intensive computation. To achieve this goal, we need to parameterize the entries of spectral kernel  $P$  as a two-variable polynomial  $p(\cdot, \cdot)$  of eigenvalues  $\lambda, \mu$  such that  $P_{ij} = p(\lambda_i, \mu_j)$ . Denote coefficients of the polynomial by  $\Theta = (\theta_{k_1 k_2}) \in \mathbb{R}^{n \times m}$ , then  $P$  is parameterized as

$$P_{ij} = p(\lambda_i, \mu_j) = \sum_{k_1=0}^{n-1} \sum_{k_2=0}^{m-1} \theta_{k_1 k_2} \lambda_i^{k_1} \mu_j^{k_2}. \quad (8)$$

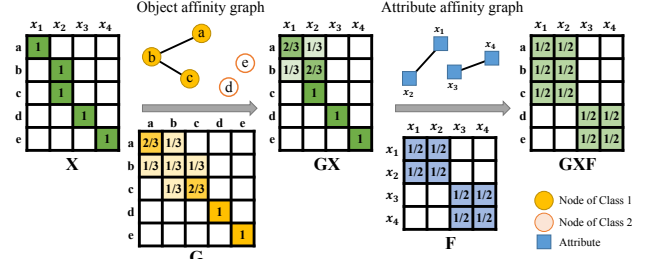
or in matrix form:

$$P = \sum_{k_1=0}^{n-1} \sum_{k_2=0}^{m-1} \theta_{k_1 k_2} \Lambda^{k_1} \mathbb{1} M^{k_2}, \quad (9)$$

where  $\mathbb{1}$  denotes all-one matrix of shape  $n \times m$ . Substitute Eq. (9) into Eq. (7) and rearrange it, 2-D graph convolution will become

$$\begin{aligned} Z &= \sum_{k_1 k_2} \theta_{k_1 k_2} U \left( S \circ (\Lambda^{k_1} \mathbb{1} M^{k_2}) \right) V^\top \\ &= \sum_{k_1 k_2} \theta_{k_1 k_2} U \Lambda^{k_1} (S \circ \mathbb{1}) M^{k_2} V^\top \\ &= \sum_{k_1 k_2} \theta_{k_1 k_2} (U \Lambda^{k_1} U^{-1}) X (V M^{k_2} V^{-1})^\top \\ &= \sum_{k_1 k_2} \theta_{k_1 k_2} L_{(1)}^{k_1} X (L_{(2)}^{k_2})^\top \end{aligned} \quad (10)$$

Eq. (10) is called 2-D spatial graph convolution, as it manipulates the signal  $X$  in the spatial domain. Eigen-decomposition of Laplacian and Fourier transformation of  $X$  is no longer required. Parameters  $\Theta = (\theta_{k_1 k_2})$  in Eq. (10) are called the (spatial) kernel of this convolution. Similar to [13], this spatial convolutional filter is localized. Denote by  $K_1$  and  $K_2$  the largest exponent of  $L_{(1)}$  and  $L_{(2)}$  in the polynomial, i.e.,  $k_1 > K_1$  or  $k_2 > K_2$  implies  $\theta_{k_1 k_2} = 0$ . The convoluted signal  $z_{ij}$  of vertex pair  $(v_i^{(1)}, v_j^{(2)})$  only depends on the neighbourhood of  $v_i^{(1)}$  within  $K_1$  hops and the neighbourhood of  $v_j^{(2)}$  within  $K_2$  hops, so the filter is said to be  $K_1$ -localized on



**Figure 2: Conceptual illustration of DSGC by a toy example. The node representations (row vectors) obtained by DSGC (GXF) is better than those in the originals feature matrix ( $X$ ) or filtered by 1-D graph convolution ( $GX$ ), in the sense that nodes of the same class have more similar features.**

$\mathcal{G}^{(1)}$  and  $K_2$ -localized on  $\mathcal{G}^{(2)}$ , and the kernel  $\Theta$  is said to be of size  $(K_1 + 1) \times (K_2 + 1)$ .

### 3.3 Dimensionwise Separable 2-D Graph Convolution (DSGC)

Although the above spatial graph convolution avoids the computationally expensive Fourier transform, its general form with kernel size  $K_1 \times K_2$  still involves at least  $K_1 \times K_2$  matrix multiplications. Inspired by the depth-wise separable convolution proposed in [22], we streamline spatial graph convolution by restricting the rank of  $\Theta$  to be one. Consequently,  $\Theta$  is able to be decomposed as an outer product of two vectors  $\theta^{(1)} \in \mathbb{R}^n$  and  $\theta^{(2)} \in \mathbb{R}^m$ , i.e.,  $\Theta = \theta^{(1)} \theta^{(2)\top}$  and  $\theta_{k_1 k_2} = \theta_{k_1}^{(1)} \theta_{k_2}^{(2)}$ . Finally, the 2-D spatial graph convolution in Eq. (10) becomes

$$Z = \sum_{k_1=0}^{K_1} \sum_{k_2=0}^{K_2} \theta_{k_1}^{(1)} \theta_{k_2}^{(2)} L_{(1)}^{k_1} X (L_{(2)}^{k_2})^\top \quad (11)$$

$$= \left( \sum_{k_1=0}^{K_1} \theta_{k_1}^{(1)} L_{(1)}^{k_1} \right) X \left( \sum_{k_2=0}^{K_2} \theta_{k_2}^{(2)} L_{(2)}^{k_2} \right)^\top = GXF, \quad (12)$$

$$\text{where } G = \sum_{k_1=0}^{K_1} \theta_{k_1}^{(1)} L_{(1)}^{k_1} \text{ and } F = \sum_{k_2=0}^{K_2} \theta_{k_2}^{(2)} (L_{(2)}^{k_2})^\top. \quad (13)$$

We refer to Eq. (12) as dimensionwise separable graph convolution (DSGC),  $G$  as the *object graph convolutional filter*, and  $F$  as the *attribute graph convolutional filter*. The fastest way to compute it only requires  $K_1 + K_2$  matrix multiplications, much less than the  $K_1 \times K_2$  matrix multiplications needed by a general 2-D spatial graph convolution. Fig. 2 illustrates how  $G$  and  $F$  can work together to learn better node representations with DSGC.

## 4 VARIANCE REDUCTION BY DSGC

Given a data distribution, the lowest possible error rate an classifier can achieve is the Bayes error rate [16], which is caused by the intrinsic overlap between different classes and cannot be avoided. In this section, we show that DSGC with proper filters can reduce intra-class variance of the data distribution while keeping class centers roughly unchanged, hence reducing the overlap between classes and improving learning performance.

**Intra-class Variance and Inter-class Variance** Suppose samples  $\mathbf{x}_i$  and their labels  $y_i$  are observations of a random vector  $\mathbb{X} = [\mathbb{X}_1, \dots, \mathbb{X}_m]^\top$  and a random variable  $\mathbb{Y}$  respectively. We define the variance of random vector  $\mathbb{X}$  to be the sum of the variance of each dimension  $\mathbb{X}_j$ , i.e., the trace of the covariance matrix of  $\mathbb{X}$ . According to law of total variance [18], the variance of  $\mathbb{X}$  can be divided into intra-class variance and inter-class variance:

$$\text{Var}(\mathbb{X}) = \underbrace{\text{E}[\text{Var}(\mathbb{X}|\mathbb{Y})]}_{\text{Intra-class Variance}} + \underbrace{\text{Var}(\text{E}[\mathbb{X}|\mathbb{Y}])}_{\text{Inter-class Variance}}, \quad (14)$$

where the conditional variance  $\text{Var}(\mathbb{X}|\mathbb{Y} = k)$  is the variance of class  $k$  and the conditional expectation  $\text{E}[\mathbb{X}|\mathbb{Y} = k]$  is the  $k$ -th class center. Intra-class variance (IntraVar) measures the average divergence within each class, while inter-class variance (InterVar) measures the divergence among class centers. We are interested in the IntraVar/InterVar ratio. Desired node representations should have low intra-class variance (i.e., compact and dense for each class), and high inter-classes variance (large margin between classes).

#### 4.1 Intra-class Variance Reduction by Object Graph Convolution

Commonly used object graph convolution reduces variance by averaging over neighborhood. For any node  $v_i$ , object graph convolution  $G\mathbf{X}$  produces a new feature vector  $\mathbf{z}_i = \sum_j G_{ij} \mathbf{x}_j$ . When  $G$  is a stochastic matrix, the output feature vector  $\mathbf{z}_i$  is a weighted average of the neighbours of  $\mathbf{x}_i$ . Denote by  $\mathbb{Z}$  a random vector of  $\mathbf{z}_i$ . Intuitively, as long as each node  $i$  has enough same-class neighbours,  $\mathbb{Z}$  will have a smaller IntraVar/InterVar ratio than  $\mathbb{X}$ .

Formally, assume that objects from different classes are connected with probability  $q$ , and classes are balanced, i.e.,  $\text{Pr}(\mathbb{Y} = k) = 1/K$  for each class  $k$ . Then, with the stochastic graph filter  $G = D^{-1}A^{(1)}$ , we have the following theorem.

**THEOREM 1.** *When  $q$  is sufficiently small, the IntraVar/InterVar ratio of  $\mathbb{Z}$  is less than or equal to that of  $\mathbb{X}$ , i.e.,*

$$\frac{\text{E}[\text{Var}(\mathbb{Z}|\mathbb{Y})]}{\text{Var}(\text{E}[\mathbb{Z}|\mathbb{Y}])} \leq \frac{\text{E}[\text{Var}(\mathbb{X}|\mathbb{Y})]}{\text{Var}(\text{E}[\mathbb{X}|\mathbb{Y}])}. \quad (15)$$

The proof is given in the Appendices B. This theorem tells that under the assumption that connected nodes are most likely to be of the same class, object graph convolution  $G\mathbf{X}$  can efficiently reduce the IntraVar/InterVar ratio.

#### 4.2 Intra-class Variance Reduction by Attribute Graph Convolution

A proper attribute graph convolutional filter  $F$  can also reduce the IntraVar/InterVar ratio. We use the convention that the random vector  $\mathbb{X}$  is a column vector, and hence the attribute graph convolution  $\mathbf{X}F$  results in a new random vector  $F^\top \mathbb{X}$ . We also assume that the node features are mean-centered, i.e.  $\text{E}[\mathbb{X}] = \mathbf{0}$ .

**THEOREM 2.** *If the attribute graph convolutional filter  $F$  is a doubly stochastic matrix, then the output of attribute graph convolution has an intra-class variance less than or equal to that of  $\mathbb{X}$ , i.e.,*

$$\begin{aligned} \sum_i F_{ij} &= \sum_j F_{ij} = 1 \text{ and } F_{ij} \geq 0, \forall i, j \\ \Rightarrow \text{E}[\text{Var}(F^\top \mathbb{X}|\mathbb{Y})] &\leq \text{E}[\text{Var}(\mathbb{X}|\mathbb{Y})]. \end{aligned}$$

The proof for Theorem 2 is given in the Appendices C.

To achieve a low IntraVar/InterVar ratio, in addition to reducing intra-class variance, we also need to keep the class centers apart after convolution, which then depends on the quality of the attribute affinity graph. A good attribute affinity graph should connect attributes that share similar expectations conditioned on  $\mathbb{Y}$ . Formally, each attribute  $\mathbb{X}_j$  has  $K$  conditional expectations w.r.t.  $\mathbb{Y}$ , which are denoted as a vector  $\mathbf{e}_j = (\text{E}[\mathbb{X}_j|\mathbb{Y} = 1], \dots, \text{E}[\mathbb{X}_j|\mathbb{Y} = K]) \in \mathbb{R}^K$ . We have the following.

**THEOREM 3.** *If  $\forall F_{ij} \neq 0, \|\mathbf{e}_i - \mathbf{e}_j\|_2 \leq \varepsilon$ , then the distance between  $\mathbf{e}_j$  and  $\hat{\mathbf{e}}_j = \sum_i F_{ij} \mathbf{e}_i$  is also less than or equal to  $\varepsilon$ , i.e.,*

$$\|\mathbf{e}_i - \mathbf{e}_j\|_2 \leq \varepsilon, \forall F_{ij} \neq 0 \Rightarrow \|\mathbf{e}_j - \hat{\mathbf{e}}_j\|_2 \leq \varepsilon,$$

and  $\varepsilon$  can be arbitrarily small with a proper  $F$ .

The proof for Theorem 3 is also given in the Appendices C. By Theorem 3, the conditional expectations of each attribute (i.e., class means) may change little after attribute graph convolution, and so does the inter-class variance. Combining Theorems 2 & 3, it suggests that a proper attribute affinity graph should connect attributes that have similar class means, so as to achieve a low IntraVar/InterVar ratio and improve performance.

## 5 UNSUPERVISED AND SEMI-SUPERVISED LEARNING WITH DSGC

### 5.1 Learning Frameworks

Given an attributed graph with node feature matrix  $\mathbf{X}$ , we can learn node representations  $\mathbf{Z}$  in an unsupervised manner by applying DSGC on  $\mathbf{X}$ , i.e.,

$$\mathbf{Z} = G\mathbf{X}F, \quad (16)$$

and then perform various downstream learning tasks with the node representations  $\mathbf{Z}$ .

**Unsupervised Node Clustering** Any standard clustering algorithm can be applied on  $\mathbf{Z}$  for clustering, as long as it is suitable for present data. In experiments, we use the popular spectral clustering method [40, 48] along with linear kernel  $\mathbf{K} = \mathbf{Z}\mathbf{Z}^\top$ .

**Semi-supervised Node Classification** After obtaining the unsupervised node representations  $\mathbf{Z}$ , we may adopt any proper supervised classifier and train it with  $\mathbf{Z}$  and a small portion of labels for semi-supervised classification. This two-step framework is semi-supervised in nature. In experiments, we choose a multi-layer perceptron with a single hidden layer as our classifier. In addition to the two-step framework, we can also plug DSGC into existing end-to-end graph-convolution-based methods. In experiments, we improve several popular methods, including GCN, GAT and GraphSAGE, by replacing their 1-D graph convolution with DSGC. For example, to incorporate DSGC into the vanilla GCN, we can modify the first layer propagation of GCN as:

$$\mathbf{H}^{(1)} = \sigma(G\mathbf{X}F\mathbf{W}^{(1)}), \quad (17)$$

where  $\mathbf{H}^{(1)}$  is the hidden units in the first layer,  $\mathbf{W}^{(1)} \in \mathbb{R}^{m \times l}$  is the trainable parameters of GCN, and  $\sigma$  is an activation function such as ReLU.

Importantly, Eq. (17) can be considered as feeding a filtered feature matrix  $\mathbf{X}F$  instead of the raw feature matrix  $\mathbf{X}$  to GCN. By our above analysis, a proper attribute graph convolutional filter  $F$



**Table 1: Dataset statistics.**

Dataset	#vertices	#edges	#cls	#features	ratio of intra-class edges
20 NG	18,846	147,034	20	11,697	96.8%
Wiki	3,767	129,597	9	18,316	38.0%
L-Cora	11,881	64,898	10	3,780	76.5%
Corn.	247	384	5	3371	23.7%
Texa.	255	205	5	3371	19.8%
Wisc.	320	721	5	3371	26.3%
Wash.	265	417	5	3371	40.3%

can reduce intra-class variance, which makes  $XF$  much easier to classify and guarantees to help train a better model. Furthermore, parameters of GCN is freely chosen from the parameter space  $\mathbf{W}^{(1)}$ , while the model trained by Eq. (17) is restricted in a subspace  $\mathbf{FW}^{(1)}$ . Since the chosen filter  $F$  is low-pass (Section 5.2) and thus is nearly singular,  $\mathbf{FW}^{(1)}$  is a subspace of  $\mathbb{R}^{m \times l}$  projected by  $F$ . Model parameters in this subspace are generally better in terms of the generalization performance, due to the variance reduction property of  $F$ . However, the model learned by Eq. (17) can hardly be learned by GCN, since the subspace  $\mathbf{FW}^{(1)}$  has measure zero, which is a tiny subset of  $\mathbb{R}^{m \times l}$ .

## 5.2 Implementation of Filters

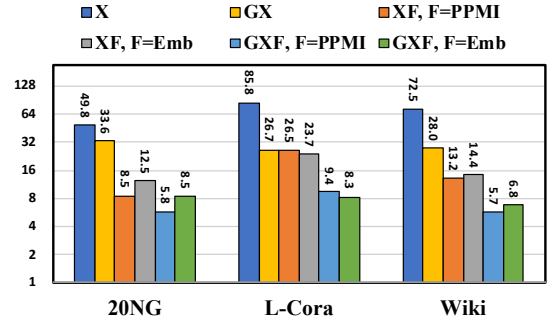
**Object Graph Convolutional Filter** In most cases, the object graph  $\mathbf{A}^{(1)}$  is given as part of the dataset, all we need is to design the filter. There are various graph convolutional filters available [13, 25, 29, 63], but the key principle of filter design for semi-supervised and unsupervised learning is low-pass [29, 34]. Following this principle, we use the 2-order row-normalized affinity matrix as the object graph filter, i.e.,

$$\mathbf{G} = (\mathbf{I} - \mathbf{L}_r^{(1)})^2 = (\mathbf{D}_1^{-1} \mathbf{A}^{(1)})^2. \quad (18)$$

**Attribute Graph Convolutional Filter** A key issue in implementing DSGC is to construct a suitable attribute affinity graph  $\mathbf{A}^{(2)}$ . Possible ways to construct attribute affinity graphs include extracting entity relation information from existing knowledge bases, building a similarity graph from features, or identifying correlations by domain knowledge. In experiments, we evaluate our methods on text dataset as in [20, 25, 46], and leverage two suitable attribute affinity graphs for text data described below.

*Positive point-wise mutual information (PPMI)* is a common tool for measuring the association between two words in computational linguistics [11]. PPMI between words  $w_i$  and  $w_j$  is defined by  $\text{PPMI}(w_i, w_j) = \left[ \log \frac{\Pr(w_i, w_j)}{\Pr(w_i) \Pr(w_j)} \right]_+$ , where  $\Pr(w_i)$  is the probability of occurrence of word  $w_i$ , and  $\Pr(w_i, w_j)$  is the probability of two words occurring together. If there is a semantic relation between two words, they usually tend to co-occur more frequently, and thus share a high PPMI value. Here, we use PPMI between words as the corresponding weights in the attribute affinity matrix  $\mathbf{A}^{(2)}$ , and symmetrically normalize it as [25].

*Word embedding based  $k$ -NN graphs.* Word embedding is a collection of techniques that map vocabularies to vectors in an Euclidean space. Embeddings of words are pre-trained vectors learned from corpus with algorithms such as GloVe [39]. Since word embeddings


**Figure 3: IntraVar/InterVar ratios. The lower, the better.**

capture semantic relations between words [4], they can be used for constructing an attribute affinity graph. With the embedding vectors, we can construct a  $k$ -NN graph with some proximity metric such as the Euclidean distance.

With the constructed attribute affinity graph  $\mathbf{A}^{(2)}$ , we use one-step lazy random walk filter [29] in our experiments, i.e.,

$$\mathbf{F} = (\mathbf{I} - \frac{1}{2} \mathbf{L}_s) = \frac{1}{2} (\mathbf{I} + \mathbf{D}_2^{-1/2} \mathbf{A}^{(2)} \mathbf{D}_2^{-1/2}). \quad (19)$$

## 6 EMPIRICAL STUDY

We conduct extensive experiments in semi-supervised node classification and node clustering on seven real-world networks including 20 Newsgroups (20 NG) [27], Large Cora (L-Cora) [29, 32], Wikispeedia (Wiki) [52, 53], and four subsets of WebKB (Cornell, Texas, Wisconsin, and Washington).<sup>2</sup> Due to space limitation, on the four subsets of WebKB, we only report the results of semi-supervised node classification with some representative baselines, but the results are indicative enough.

### 6.1 Datasets

The statistics of all datasets are summarized in Table 1, where the last row shows the intra-class edge ratio of the object link graph of each dataset, which can reflect the quality of the graph.

**20 Newsgroups** (20 NG) [27] is an email discussion group, where each object is an email and there are 18846 emails in total. Each email is represented by an 11697-dimension tf-idf feature vector. Two emails are connected by an edge if they replies the same one.

**Wikispeedia** (Wiki) [52, 53] is a webpage network in which the objects are 3767 Wikipedia webpages, and the edges are web hyperlinks. Each webpage is described by a 18316-dimension tf-idf vector. We removed several tiny classes, so the webpages distribute more evenly across the remaining 9 categories.

**Large Cora** (L-Cora) [32] is a citation network in which the objects are computer science research papers represented by 3780 dimension of tf-idf values. Two papers are connected by an undirected edge if and only if one cites the other. These citation links form a object graph. After removing the papers that belong to no topic and the ones that have no authors or title, a subset of 11881 papers is obtained [43]. We name this dataset ‘‘Large Cora’’ to distinguish it from the ‘‘Cora’’ dataset with 2708 papers used in [25, 44, 59].

<sup>2</sup>Note that we did not use the ‘‘Cora’’, ‘‘Citeseer’’ and ‘‘PubMed’’ datasets as in [25, 44, 59], since the attribute (word) lists are not provided.

Table 2: Classification accuracy on 20 NG, L-Cora, and Wiki.

Dataset			20 NG		L-Cora		Wiki	
Method	<i>G</i>	<i>F</i>	20 labels/cls.	5 labels/cls.	20 labels/cls.	5 labels/cls.	20 labels/cls.	5 labels/cls.
MLP	✗	✗	63.76 ± 0.17	38.67 ± 0.38	52.97 ± 0.41	39.56 ± 0.85	67.23 ± 0.25	54.41 ± 0.66
LP [65]	✓	✗	16.39 ± 0.20	8.62 ± 0.20	55.77 ± 0.97	38.97 ± 3.15	9.53 ± 0.05	10.54 ± 0.19
GLP [29]	✓	✗	74.99 ± 0.11	52.62 ± 0.45	68.95 ± 0.29	56.42 ± 0.85	60.05 ± 0.11	48.45 ± 0.54
GCN [25]	✓	✗	76.25 ± 0.11	53.78 ± 0.49	67.75 ± 0.33	54.27 ± 0.82	59.81 ± 0.30	47.93 ± 0.56
GAT [46]	✓	✗	76.33 ± 0.16	56.02 ± 0.57	68.88 ± 0.78	56.89 ± 1.53	50.97 ± 0.54	46.99 ± 0.83
DGI [47]	✓	✗	73.34 ± 0.27	<b>66.57</b> ± 0.63	61.39 ± 0.50	54.77 ± 1.24	49.70 ± 1.63	43.64 ± 1.89
GraphSAGE [20]	✓	✗	65.73 ± 0.17	42.48 ± 0.77	57.28 ± 0.71	46.79 ± 1.91	65.52 ± 0.62	48.81 ± 0.76
GCNII [10]	✓	✗	<b>77.41</b> ± 0.12	58.10 ± 0.85	68.18 ± 0.19	57.02 ± 0.55	43.65 ± 1.03	35.98 ± 4.45
JK-MaxPool [57]	✓	✗	71.00 ± 0.19	49.09 ± 0.27	67.44 ± 0.18	51.63 ± 0.52	45.26 ± 0.37	44.13 ± 0.38
JK-Concat [57]	✓	✗	72.24 ± 0.13	49.76 ± 0.27	67.47 ± 0.19	51.96 ± 0.46	47.24 ± 0.30	45.21 ± 0.29
GRAND [15]	✓	✗	74.45 ± 0.72	57.97 ± 2.79	69.30 ± 0.59	52.12 ± 0.73	62.25 ± 0.93	47.17 ± 3.38
DSGC ( <i>GX</i> )	✓	✗	75.60 ± 0.13	53.84 ± 0.46	67.74 ± 0.30	55.67 ± 0.72	58.73 ± 0.34	47.34 ± 0.54
DSGC ( <i>XF</i> )	✗	Emb	66.27 ± 0.13	48.04 ± 0.38	58.70 ± 0.30	46.41 ± 0.55	<b>69.76</b> ± 0.20	<b>59.76</b> ± 0.58
	✗	PPMI	75.36 ± 0.11	59.61 ± 0.34	61.01 ± 0.23	48.31 ± 0.62	<b>69.91</b> ± 0.21	<b>60.13</b> ± 0.61
DSGC ( <i>GXF</i> )	✓	Emb	76.53 ± 0.15	59.91 ± 0.31	<b>69.81</b> ± 0.26	<b>58.63</b> ± 0.75	60.50 ± 0.26	49.69 ± 0.56
	✓	PPMI	<b>81.69</b> ± 0.12	<b>68.94</b> ± 0.32	<b>70.20</b> ± 0.24	<b>59.43</b> ± 0.68	58.84 ± 0.26	48.51 ± 0.54

\* ✓ and ✗ indicate using/not using *G* or *F*.

Table 3: Classification accuracy on four subsets of WebKB.

Method	<i>F</i>	Corn.	Texa.	Wisc.	Wash.
GCN [25]	✗	50.25 ± 0.66	60.31 ± 1.06	52.74 ± 0.92	53.83 ± 1.36
GLP [29]	✗	50.86 ± 0.75	59.40 ± 1.27	55.28 ± 0.97	55.87 ± 1.32
GAT [46]	✗	51.21 ± 1.40	60.06 ± 1.08	52.92 ± 1.18	56.85 ± 2.00
GRAND [15]	✗	49.01 ± 1.04	57.38 ± 0.26	49.30 ± 0.22	42.09 ± 1.25
DSGC ( <i>GX</i> )	✗	51.22 ± 0.77	59.35 ± 1.26	56.26 ± 0.93	55.77 ± 1.30
DSGC ( <i>XF</i> )	Emb	<b>62.14</b> ± 1.02	<b>68.00</b> ± 0.75	<b>73.50</b> ± 0.81	<b>65.78</b> ± 1.27
	PPMI	<b>60.10</b> ± 0.91	<b>67.34</b> ± 0.94	<b>72.88</b> ± 0.78	<b>65.40</b> ± 1.35
DSGC ( <i>GXF</i> )	Emb	53.02 ± 0.67	61.89 ± 0.94	58.64 ± 1.12	59.05 ± 1.21
	PMI	52.35 ± 0.52	61.83 ± 0.91	56.40 ± 1.04	57.01 ± 1.20

WebKB<sup>3</sup> is a webpage dataset collected by Carnegie Mellon University. We use its four subdatasets: Cornell, Texas, Wisconsin, and Washington. The objects are webpages, and the edges are hyperlinks between them. The webpages are represented by tf-idf feature vectors, and manually classified into five categories: student, project, course, staff, and faculty.

## 6.2 Variance Reduction and Visualization

First of all, to verify our analysis in section 4, we illustrate the variance reduction effect of both object graph convolution and attribute graph convolution. As shown in Figure 3, 1-D graph convolution (*GX* or *XF*) already greatly reduces the IntraVar/InterVar ratio, and 2-D graph convolution (*GXF*) reduces it even further.

In Figure 1, we visualize the results of performing graph convolution on the object features of 20 NG by t-SNE. It can be seen that both object graph convolution and attribute graph convolution can successfully reduce the overlap among classes, and 2-D graph convolution (DSGC) is more effective than 1-D.

<sup>3</sup><http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/>

Table 4: Baselines improved by DSGC.

Methods	<i>F</i>	20 NG	L-Cora	Wiki
GAT [46]	✗	76.33 ± 0.16	68.88 ± 0.78	50.97 ± 0.54
	PPMI	78.01 ± 0.30	67.38 ± 0.65	55.43 ± 0.51
GCN [25]	✗	76.25 ± 0.11	67.75 ± 0.33	59.81 ± 0.30
	PPMI	81.60 ± 0.10	67.87 ± 0.25	61.33 ± 0.28
GraphSAGE [20]	✗	65.73 ± 0.17	57.28 ± 0.71	65.52 ± 0.62
	PPMI	76.27 ± 0.33	60.23 ± 1.81	67.26 ± 0.52

\* ✗ indicates not using *F*.

## 6.3 Semi-supervised Node Classification

**Baselines** We compare DSGC with the following baselines: label propagation (LP) [56], multi-layer perceptron (MLP), graph convolutional networks (GCN) [25], generalized label propagation (GLP) [29], GraphSAGE [20], graph attention networks (GAT) [46], deep graph infomax (DGI) [47], GCNII [10], jumping knowledge networks (JK-MaxPool & JK-Concat) [57], and GRAND[15]. We also try to improve GCN, GAT, and GraphSAGE by DSGC as described in Section 5.1. Our methods with mere object graph filter (*G*) or mere attribute graph filter (*F*) are also tested, for the purpose of ablation study. PPMI and Emb denote attribute affinity graphs constructed by positive point-wise mutual information and word embedding respectively as described in Section 5.2).

**Settings** For 20 NG, L-Cora and Wiki, we test two scenarios – 20 labels/class and 5 labels/class. We follow GCN [25] and many others to set aside a validation set containing 500 samples for hyperparameter tuning. For the four small subdatasets of WebKB, we randomly split them into 5/30/65% as train/valid/test set, and ensure each class has at least 1 label. Hyper-parameters of all methods, including ours and baselines, are tuned by grid search according to validation. The reported results of all methods are averaged over 50 runs. More experimental details are provided in Appendix A.

**Table 5: Clustering performance.**

Dataset			20 NG		L-Cora		Wiki	
Method	<i>G</i>	<i>F</i>	Acc(%)	NMI(%)	Acc(%)	NMI(%)	Acc(%)	NMI(%)
Spectral	$\times$	$\times$	25.29 $\pm$ 1.01	28.18 $\pm$ 0.74	28.22 $\pm$ 1.01	11.61 $\pm$ 0.04	29.25 $\pm$ 0.00	21.83 $\pm$ 0.00
GAE [24]	$\checkmark$	$\times$	38.92 $\pm$ 1.39	44.58 $\pm$ 0.40	34.45 $\pm$ 0.76	22.38 $\pm$ 0.18	33.78 $\pm$ 0.32	22.88 $\pm$ 0.20
VGAE [24]	$\checkmark$	$\times$	25.04 $\pm$ 0.81	25.72 $\pm$ 0.77	29.45 $\pm$ 1.25	17.53 $\pm$ 0.15	33.83 $\pm$ 0.45	21.46 $\pm$ 0.19
MGAE [49]	$\checkmark$	$\times$	47.83 $\pm$ 2.33	<b>56.14</b> $\pm$ 1.00	35.87 $\pm$ 0.97	30.57 $\pm$ 0.98	32.73 $\pm$ 1.16	27.95 $\pm$ 2.29
ARGE [35]	$\checkmark$	$\times$	42.04 $\pm$ 0.50	44.13 $\pm$ 0.91	36.07 $\pm$ 0.05	27.74 $\pm$ 0.01	26.49 $\pm$ 0.10	17.17 $\pm$ 0.05
ARVGE [35]	$\checkmark$	$\times$	21.10 $\pm$ 0.61	21.79 $\pm$ 0.49	26.45 $\pm$ 0.03	12.94 $\pm$ 0.01	33.82 $\pm$ 0.13	21.42 $\pm$ 0.11
AGC [63]	$\checkmark$	$\times$	38.83 $\pm$ 0.84	47.08 $\pm$ 1.57	<b>41.76</b> $\pm$ 0.01	<b>33.65</b> $\pm$ 0.01	32.74 $\pm$ 0.01	24.90 $\pm$ 0.01
DSGC ( <i>GX</i> )	$\checkmark$	$\times$	38.42 $\pm$ 0.66	46.28 $\pm$ 0.93	38.26 $\pm$ 0.02	30.66 $\pm$ 0.02	31.43 $\pm$ 0.09	24.16 $\pm$ 0.18
DSGC ( <i>XF</i> )	$\times$	Emb	28.99 $\pm$ 0.06	33.22 $\pm$ 0.10	30.80 $\pm$ 0.56	17.46 $\pm$ 0.21	<b>35.45</b> $\pm$ 0.91	<b>33.44</b> $\pm$ 0.66
	$\times$	PPMI	<b>48.36</b> $\pm$ 2.40	53.27 $\pm$ 2.17	36.46 $\pm$ 0.06	22.53 $\pm$ 0.03	<b>38.10</b> $\pm$ 0.01	<b>36.07</b> $\pm$ 0.02
DSGC ( <i>GXF</i> )	$\checkmark$	Emb	43.40 $\pm$ 0.66	50.97 $\pm$ 0.58	40.75 $\pm$ 0.02	<b>33.05</b> $\pm$ 0.04	30.50 $\pm$ 0.01	25.48 $\pm$ 0.03
	$\checkmark$	PPMI	<b>52.25</b> $\pm$ 1.97	<b>61.34</b> $\pm$ 1.07	<b>41.24</b> $\pm$ 0.04	30.92 $\pm$ 0.01	31.37 $\pm$ 0.08	26.06 $\pm$ 0.20

\*  $\checkmark$  and  $\times$  indicate using/not using *G* or *F*.

**Performance** Classification accuracies are summarized in Tables 2 and 3, and the top 2 accuracies are highlighted. Results of improved GAT, GCN, and GraphSAGE are shown in Table 4. The following observations can be made. Firstly, object graph convolution (*G*) does not always help. On 20 NG and L-Cora, methods based on it like DSGC (*GX*), GCN and GAT all outperform MLP significantly. However, on Wiki and the four subsets of WebKB, object graph convolution severely harms the performance. This is because the hyperlink graphs are highly noisy. Only a small portion of edges connect nodes of the same class, much lower than that of 20 NG (96.8%) and L-Cora (76.5%) (see Table 1). This shows the limitation of object graph convolution. Secondly, attribute graph convolution works. As shown in Table 2, DSGC with mere *F* already outperforms MLP significantly. Especially, on Wiki and WebKB, object graph convolution fails while attribute graph convolution is still effective. Thirdly, 2-D graph convolution is useful. On datasets with good object link graphs like 20 NG and L-Cora, DSGC with both *G* and *F* performs much better than with either one of them only. Especially, DSGC (*GXF*) with PPMI achieves the best performance among all methods. On datasets with bad object link graphs such as Wiki and the four subsets of WebKB, DSGC with both *G* and *F* improves upon DSGC with mere *G* and outperforms most baselines. Especially, DSGC with mere *F* achieves the best performance, which improves upon the best baseline by 3.63% and 5.72% in absolute accuracy in two scenarios.

Remarkably, it can be seen from Table 4 that by incorporating DSGC, the performance of baselines including GCN, GAT, GraphSAGE is improved substantially in most cases, which again confirms that attribute graph convolution is a useful complement to object graph convolution.

## 6.4 Node Clustering

**Baselines** We test the proposed node clustering method with DSGC (Section 5.1) with or without *G* and *F* in five cases, and compare them with existing strong baselines including GAE and VGAE [24], MGAE [49], ARGE and ARVGE [35] and AGC [63]. We also compare with the spectral clustering (Spectral) method

that operates on a similarity graph constructed by a linear kernel. Detailed experiment settings are included in Appendix A.

**Performance** We adopt two widely-used clustering measures [2]: clustering accuracy (Acc) and normalized mutual information (NMI), and the results are shown in Table 5 with the top 2 results highlighted. We can make the following observations. 1) Attribute graph convolution is highly effective. On 20 NG, DSGC (*XF*) with PPMI outperforms most baselines by a very large margin. On Wiki, DSGC (*XF*) with PPMI or Emb significantly outperforms all the baselines. 2) 2-D graph convolution is beneficial as validated in the classification experiments. On 20 NG, DSGC (*GXF*) with PPMI can further improve upon the already very strong performance of DSGC (*XF*) with PPMI and performs the best; On L-Cora, DSGC (*GXF*) with PPMI or Emb improves upon either DSGC (*GX*) or DSGC (*XF*) and outperforms most baselines significantly. On Wiki, DSGC (*XF*) performs better than DSGC (*GXF*), due to the low-quality object link graph as explained above.

## 7 CONCLUSION

We have proposed a simple and efficient dimensionwise separable 2-D graph convolution (DSGC) for unsupervised and semi-supervised learning on graphs. We have demonstrated theoretically and empirically that by exploiting attribute relations in addition to object relations, DSGC can learn better node representations than existing methods based on the regular 1-D graph convolution, leading to promising performance on node classification and clustering tasks. We believe DSGC can be applied to a wide variety of applications such as action recognition, malware detection, and recommender systems. In future work, we plan to apply DSGC to solve more practical problems.

## ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their insightful comments. This research was supported by the General Research Fund No.15222220 funded by the UGC of Hong Kong.



## REFERENCES

- [1] Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr Harutyunyan, Greg Ver Steeg, and Aram Galstyan. 2019. MixHop: Higher-Order Graph Convolutional Architectures via Sparsified Neighborhood Mixing. In *ICML*. 21–29.
- [2] Charu C Aggarwal and Chandan K Reddy. 2014. *Data Clustering: Algorithms and Applications*. CRC Press, Boca Raton.
- [3] Réka Albert and Albert-László Barabási. 2002. Statistical mechanics of complex networks. *Reviews of modern physics* 74, 1 (2002), 47.
- [4] Amir Bakarov. 2018. A Survey of Word Embeddings Evaluation Methods. *CoRR* abs/1801.09536 (2018).
- [5] M. Belkin, P. Niyogi, and V Sindhwani. 2006. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research* 7, 1 (2006), 2399–2434.
- [6] A. Blum, J. Lafferty, M.R. Rwebangira, and R Reddy. 2004. Semi-supervised learning using randomized mincuts. In *ICML*. 13.
- [7] Shaosheng Cao, Wei Lu, and Qiongkai Xu. 2016. Deep Neural Networks for Learning Graph Representations. In *AAAI*. 1145–1152.
- [8] O. Chapelle, J. Weston, and B Scholkopf. 2003. Cluster kernels for semi-supervised learning. In *NeurIPS*. 601–608.
- [9] O. Chapelle and A Zien. 2005. Semi-supervised classification by low density separation. In *International Workshop on Artificial Intelligence and Statistics*. 57–64.
- [10] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and Deep Graph Convolutional Networks. In *ICML*.
- [11] Kenneth Ward Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics* 16, 1 (1990), 22–29.
- [12] Quanyu Dai, Qiang Li, Jian Tang, and Dan Wang. 2018. Adversarial network embedding. In *AAAI*.
- [13] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *NeurIPS*. 3844–3852.
- [14] Venkatesan N Ekambaram, Giulia Fanti, Babak Ayazifar, and Kannan Ramchandran. 2013. Wavelet-regularized graph semi-supervised learning. In *Global Conference on Signal and Information Processing*. 423–426.
- [15] Wenzheng Feng, Jie Zhang, Yuxiao Dong, Yu Han, Huanbo Luan, Qian Xu, Qiang Yang, Evgeny Kharlamov, and Jie Tang. 2020. Graph Random Neural Networks for Semi-Supervised Learning on Graphs. *Advances in Neural Information Processing Systems* 33 (2020).
- [16] Keinosuke Fukunaga. 2013. *Introduction to statistical pattern recognition*. Elsevier.
- [17] Benjamin Girault, Paulo Gonçalves, Eric Fleury, and Arashpreet Singh Mor. 2014. Semi-supervised learning for graph to signal mapping: A graph signal wiener filter interpretation. In *Conference on Acoustics, Speech and Signal Processing*. 1115–1119.
- [18] Charles M Grinstead and James Laurie Snell. 2012. *Introduction to probability*. American Mathematical Soc.
- [19] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *ACM SIGKDD*. ACM, 855–864.
- [20] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*. 1024–1034.
- [21] Matthias Hein and Markus Maier. 2007. Manifold Denoising. In *NeurIPS*. 561–568.
- [22] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).
- [23] Xiao Huang, Jundong Li, and Xia Hu. 2017. Accelerated attributed network embedding. In *SIAM ICDM*. 633–641.
- [24] Thomas N Kipf and Max Welling. 2016. Variational Graph Auto-Encoders. *NIPS Workshop on Bayesian Deep Learning* (2016).
- [25] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *ICLR*.
- [26] Takashi Kurokawa, Taihei Oki, and Hiromichi Nagao. 2017. Multi-dimensional Graph Fourier Transform. *arXiv preprint arXiv:1712.07811* (2017).
- [27] Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *ICML*. 331–339.
- [28] Qimai Li, Zhichao Han, and Xiao-Ming Wu. 2018. Deeper Insights into Graph Convolutional Networks for Semi-Supervised Learning. In *AAAI*. 3538–3545.
- [29] Qimai Li, Xiao-Ming Wu, Han Liu, Xiaotong Zhang, and Zhichao Guan. 2019. Label Efficient Semi-Supervised Learning via Graph Filtering. In *CVPR*. 9582–9591.
- [30] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. 2016. Gated Graph Sequence Neural Networks. In *ICLR*.
- [31] Renjie Liao, Zhizhen Zhao, Raquel Urtasun, and Richard S Zemel. 2019. LanczosNet: Multi-Scale Deep Graph Convolutional Networks. *CoRR* abs/1901.01484 (2019).
- [32] Andrew McCallumzy, Kamal Nigamy, Jason Rennie, and Kristie Seymorey. 1999. Building domain-specific search engines with machine learning techniques. In *Proceedings of the AAAI Spring Symposium on Intelligent Agents in Cyberspace*. Citeseer, 28–39.
- [33] Federico Monti, Michael M. Bronstein, and Xavier Bresson. [n.d.]. Geometric Matrix Completion with Recurrent Multi-Graph Neural Networks. In *NIPS'17*.
- [34] Hoang NT and Takanori Maehara. 2021. Revisiting Graph Neural Networks: All We Have is Low-Pass Filters. In *25th International Conference on Pattern Recognition (ICPR'20)*.
- [35] Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang. 2018. Adversarially Regularized Graph Autoencoder for Graph Embedding. In *IJCAI*. 2609–2615.
- [36] Shirui Pan, Jia Wu, Xingquan Zhu, Chengqi Zhang, and Yang Wang. 2016. Tri-Party Deep Network Representation. In *IJCAI*. 1895–1901.
- [37] Shichao Pei, Lu Yu, Guoxian Yu, and Xiangliang Zhang. [n.d.]. REA: Robust Cross-lingual Entity Alignment Between Knowledge Graphs. In *KDD*.
- [38] Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. 2020. Graph Representation Learning via Graphical Mutual Information Maximization. In *WWW*. 259–270.
- [39] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *EMNLP*. 1532–1543.
- [40] Pietro Perona and William Freeman. 1998. A factorization approach to grouping. In *ECCV*. 655–670.
- [41] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *ACM SIGKDD*. 701–710.
- [42] Meng Qu, Yoshua Bengio, and Jian Tang. 2019. GMNN: Graph Markov Neural Networks. In *ICML*. 5241–5250.
- [43] Claudio Saccà, Michelangelo Diligenti, and Marco Gori. 2013. Collective Classification Using Semantic Based Regularization. In *IEEE ICMLA*. Vol. 1. 283–286.
- [44] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. 2008. Collective classification in network data. *AI Magazine* 29, 3 (2008), 93–106.
- [45] Partha Pratim Talukdar and Koby Crammer. 2009. New regularized algorithms for transductive learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 442–457.
- [46] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR*.
- [47] Petar Velickovic, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. 2019. Deep Graph Infomax. In *ICLR*.
- [48] Ulrike Von Luxburg. 2007. A tutorial on spectral clustering. *Statistics and computing* 17, 4 (2007), 395–416.
- [49] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jing Jiang. 2017. Mgae: Marginalized graph autoencoder for graph clustering. In *CIKM*. 889–898.
- [50] Daixin Wang, Peng Cui, and Wenwu Zhu. 2016. Structural deep network embedding. In *ACM SIGKDD*. 1225–1234.
- [51] Pengyang Wang, Kunpeng Liu, Lu Jiang, Xiaolin Li, and Yanjie Fu. 2020. Incremental Mobile User Profiling: Reinforcement Learning with Spatial Knowledge Graph for Modeling Event Streams. In *KDD*. 853–861.
- [52] Robert West and Jure Leskovec. 2012. Human wayfinding in information networks. In *WWW*. 619–628.
- [53] Robert West, Joelle Pineau, and Doina Precup. 2009. Wikispeedia: An Online Game for Inferring Semantic Distances between Concepts. In *IJCAI*. 1598–1603.
- [54] Jason Weston, Frédéric Ratle, Hossein Mobahi, and Ronan Collobert. 2008. Deep learning via semi-supervised embedding. In *ICML*. 1168–1175.
- [55] Ning Wu, Wayne Xin Zhao, Jingyuan Wang, and Dayan Pan. [n.d.].
- [56] Xiao-Ming Wu, Zhenguo Li, Anthony M. So, John Wright, and Shih-fu Chang. 2012. Learning with Partially Absorbing Random Walks. In *NeurIPS*. 3077–3085.
- [57] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation Learning on Graphs with Jumping Knowledge Networks. In *Proceedings of the 35th ICML (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 5453–5462.
- [58] Cheng Yang, Zhiyuan Liu, Deli Zhao, Maosong Sun, and Edward Y Chang. 2015. Network representation learning with rich text information. In *IJCAI*. 2111–2117.
- [59] Zhilin Yang, William W Cohen, and Ruslan Salakhutdinov. 2016. Revisiting semi-supervised learning with graph embeddings. In *ICML*. 40–48.
- [60] Fanghua Ye, Chuan Chen, and Zibin Zheng. 2018. Deep Autoencoder-like Non-negative Matrix Factorization for Community Detection. In *CIKM*. 1393–1402.
- [61] Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, and Dit-Yan Yeung. 2018. GaAN: Gated Attention Networks for Learning on Large and Spatiotemporal Graphs. In *UAI*. 339–349.
- [62] T. Zhang and R.K Ando. 2006. Analysis of spectral kernel design based semi-supervised learning. In *NeurIPS*. 1601–1608.
- [63] Xiaotong Zhang, Han Liu, Qimai Li, and Xiao-Ming Wu. 2019. Attributed Graph Clustering via Adaptive Graph Convolution. In *IJCAI*. 4327–4333.
- [64] Denny Zhou, Olivier Bousquet, Thomas N Lal, Jason Weston, and Bernhard Schölkopf. 2004. Learning with local and global consistency. In *NeurIPS*. 321–328.
- [65] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. 2003. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*. 912–919.

# Appendices

## Appendix A EXPERIMENTAL SETTINGS

### A.1 Semi-supervised Node Classification

For constructing the PPMI graph, we set the context window size to 20 and use the inverse distance as co-occurrence weight. For constructing the Emb graph, we use GloVe [39] to learn word embeddings and set the number of nearest neighbours  $k = 20$ . The classifier of DSGC we use is a multi-layer perceptron (MLP) with a 64-unit hidden layer. The MLP is trained for 200 epochs by Adam Optimizer.

Hyperparameters of all models, including our methods and baselines, are tuned by grid search based on validation. Dropout rate is selected from  $\{0, 0.2, 0.5\}$ , and 0.2 is chosen. Learning rate is selected from  $\{1, 0.1, 0.01, 0.001\}$ , and 0.1 is chosen. Weight decay is selected from  $\{5e-3, 5e-4, 5e-5, 5e-6, 5e-7\}$ , and  $5e-4$  is chosen for WebKB,  $5e-5$  for L-Cora,  $5e-6$  for 20 NG,  $5e-7$  for Wiki. Results of our methods and baselines are averaged over 50 runs.

### A.2 Node Clustering

For our method DSGC, the attribute affinity graphs PPMI and Emb are constructed in the same way as in object classification. For other baselines, we follow the parameter settings described in the original papers. In particular, for GAE and VGAE [24], we construct encoders with a 32-neuron hidden layer and a 16-neuron embedding layer, and train the encoders for 200 iterations using the Adam algorithm with a learning rate of 0.01.

For MGAE [49], the corruption level  $p$  is 0.4, the number of layers is 3, and the parameter  $\lambda$  is  $10^{-5}$ . For ARGE and ARVGE [35], we construct encoders with a 32-neuron hidden layer and a 16-neuron embedding layer. The discriminators are built by two hidden layers with 16 neurons and 64 neurons respectively. We train all the autoencoder-related models for 200 iterations and optimize them using the Adam algorithm. The learning rates of encoder and discriminator are both 0.001. For AGC [63], the maximum iteration number is 60. For fair comparison, these baselines also adopt spectral clustering as our DSGC to obtain the clustering results. We repeat each method for 10 times and report the average clustering results and standard deviations.

## Appendix B

Assume that objects from the same class are connected with probability  $r$ , and objects from different classes are connected with probability  $q$ , i.e., the adjacency matrix  $\mathbf{A}^{(1)}$  of object graph obeys following distribution:

	if $y_i = y_j$	if $y_i \neq y_j$
$\Pr(a_{ij} \neq 0)$	$r$	$q$
$\Pr(a_{ij} = 0)$	$1 - r$	$1 - q$

We also assume that classes are balanced, i.e.,  $\Pr(\mathbb{Y} = k) = 1/K$  for all  $k$ . Then, with the stochastic graph filter  $G = D^{-1}\mathbf{A}^{(1)}$ , we have the following theorem.

**THEOREM 1.** *When  $q$  is sufficiently small, the IntraVar/InterVar ratio of  $\mathbb{Z}$  is less than or equal to that of  $\mathbb{X}$ , i.e.,*

$$\frac{\mathbb{E}[\text{Var}(\mathbb{Z}|\mathbb{Y})]}{\text{Var}(\mathbb{E}[\mathbb{Z}|\mathbb{Y}])} \leq \frac{\mathbb{E}[\text{Var}(\mathbb{X}|\mathbb{Y})]}{\text{Var}(\mathbb{E}[\mathbb{X}|\mathbb{Y}])}. \quad (20)$$

**PROOF.** The proof consists of two parts. In the first part, we prove that inter-class variance is unchanged after object graph convolution, when  $q$  approximates 0, i.e.,

$$\lim_{q \rightarrow 0} \text{Var}(\mathbb{E}[\mathbb{Z}|\mathbb{Y}]) = \text{Var}(\mathbb{E}[\mathbb{X}|\mathbb{Y}]). \quad (21)$$

In the second part, we prove that intra-class variance becomes smaller after object graph convolution, i.e.,

$$\mathbb{E}[\text{Var}(\mathbb{Z}|\mathbb{Y})] \leq \mathbb{E}[\text{Var}(\mathbb{X}|\mathbb{Y})], \quad (22)$$

when  $G$  is a stochastic matrix.

**Part 1. Inter-class variance is unchanged.** Since  $z_i = \sum_j G_{ij}x_j$ , we have

$$\begin{aligned} \mathbb{E}[z_i | y_i = k] &= \sum_j \mathbb{E}[G_{ij}] \mathbb{E}[x_j] \\ &= \sum_{j, y_j = k} \mathbb{E}[G_{ij}] \mathbb{E}[x_j] + \sum_{j, y_j \neq k} \mathbb{E}[G_{ij}] \mathbb{E}[x_j] \\ &= \frac{\sum_{j, y_j = k} \mathbb{E}[a_{ij}] \mathbb{E}[x_j] + \sum_{j, y_j \neq k} \mathbb{E}[a_{ij}] \mathbb{E}[x_j]}{\sum_j \mathbb{E}[a_{ij}]} \\ &= \frac{r \sum_{j, y_j = k} \mathbb{E}[\mathbb{X}|\mathbb{Y} = k] + q \sum_{j, y_j \neq k} \mathbb{E}[x_j]}{\frac{N}{K}(r - q) + Nq} \\ &= \frac{\frac{N}{K}r\mathbb{E}[\mathbb{X}|\mathbb{Y} = k] + q \sum_j \mathbb{E}[x_j] - q \sum_{j, y_j = k} \mathbb{E}[x_j]}{\frac{N}{K}(r - q) + Nq} \\ &= \frac{\frac{N}{K}(r - q)\mathbb{E}[\mathbb{X}|\mathbb{Y} = k] + Nq\mathbb{E}[\mathbb{X}]}{\frac{N}{K}(r - q) + Nq} \\ &= \frac{(r - q)\mathbb{E}[\mathbb{X}|\mathbb{Y} = k] + Kq\mathbb{E}[\mathbb{X}]}{(r - q) + Kq} \end{aligned} \quad (23)$$

When  $q$  approximates 0, Eq. (23) approximates  $\mathbb{E}[\mathbb{X}|\mathbb{Y} = k]$ , so

$$\begin{aligned} \mathbb{E}[\mathbb{Z}|\mathbb{Y} = k] &= \sum_{i, y_i = k} \Pr(\mathbb{Z} = z_i | y_i = k) \mathbb{E}[z_i | y_i = k] \\ &= \sum_{i, y_i = k} \Pr(\mathbb{Z} = z_i | y_i = k) \mathbb{E}[\mathbb{X}|\mathbb{Y} = k] = \mathbb{E}[\mathbb{X}|\mathbb{Y} = k]. \end{aligned}$$

Take variance of both side, we get

$$\text{Var}(\mathbb{E}[\mathbb{Z}|\mathbb{Y}]) = \text{Var}(\mathbb{E}[\mathbb{X}|\mathbb{Y}]). \quad (24)$$

**Part 2. Intra-class variance becomes smaller.** Denote by  $\text{Cov}(\cdot, \cdot)$  the covariance of two random variables. We have following inequality about variance.

$$\begin{aligned} \text{Var}\left(\sum_j G_{ij}x_j\right) &= \sum_j G_{ij}^2 \text{Var}(x_j) + \sum_{j, l} G_{ij}G_{il} \text{Cov}(x_j, x_l) \\ &\leq \sum_{j, l} G_{ij}G_{il} \sqrt{\text{Var}(x_j)} \sqrt{\text{Var}(x_l)} \\ &= \left(\sum_j G_{ij} \sqrt{\text{Var}(x_j)}\right)^2. \end{aligned} \quad (25)$$

Consider the variance of filtering result  $\mathbf{z}_i$  for each sample in class  $k$ , it is less than variance of  $\mathbb{X}$  of that class:

$$\begin{aligned}
& \text{Var}(\mathbf{z}_i | y_i = k) \\
&= \text{Var} \left( \sum_j G_{ij} \mathbf{x}_j \middle| y_j = k \right) \\
&\leq \left( \sum_j G_{ij} \sqrt{\text{Var}(\mathbf{x}_j | y_j = k)} \right)^2 \quad \# \text{ by inequality (25)} \\
&= \left( \sum_j G_{ij} \sqrt{\text{Var}(\mathbb{X} | \mathbb{Y} = k)} \right)^2 \\
&= \left( \sqrt{\text{Var}(\mathbb{X} | \mathbb{Y} = k)} \right)^2 \quad \# \text{ since } \sum_j G_{ij} = 1 \\
&= \text{Var}(\mathbb{X} | \mathbb{Y} = k),
\end{aligned}$$

Then variance of random vector  $\mathbb{Z}$  for each class is less than variance of  $\mathbb{X}$  of that class:

$$\begin{aligned}
\text{Var}(\mathbb{Z} | \mathbb{Y} = k) &= \sum_{i, y_i = k} \text{Pr}(\mathbb{Z} = \mathbf{z}_i | y_i = k) \text{Var}(\mathbf{z}_i | y_i = k) \\
&\leq \text{Var}(\mathbb{X} | \mathbb{Y} = k).
\end{aligned}$$

Sum them over all classes:

$$\begin{aligned}
\mathbb{E}[\text{Var}(\mathbb{Z} | \mathbb{Y})] &= \sum_k \text{Pr}(\mathbb{Y} = k) \text{Var}(\mathbb{Z} | \mathbb{Y} = k) \\
&\leq \sum_k \text{Pr}(\mathbb{Y} = k) \text{Var}(\mathbb{X} | \mathbb{Y} = k) \\
&= \mathbb{E}[\text{Var}(\mathbb{X} | \mathbb{Y})]. \tag{26}
\end{aligned}$$

Combining Eq. (24) and Eq. (26), we prove that when  $q$  is sufficiently small,

$$\frac{\mathbb{E}[\text{Var}(\mathbb{Z} | \mathbb{Y})]}{\text{Var}(\mathbb{E}[\mathbb{Z} | \mathbb{Y}])} \leq \frac{\mathbb{E}[\text{Var}(\mathbb{X} | \mathbb{Y})]}{\text{Var}(\mathbb{E}[\mathbb{X} | \mathbb{Y}])}. \tag{27}$$

□

## Appendix C

**THEOREM 2.** *If the attribute graph convolutional filter  $F$  is a doubly stochastic matrix, then the output of attribute graph convolution has an intra-class variance less than or equal to that of  $\mathbb{X}$ , i.e.,*

$$\begin{aligned}
\sum_i F_{ij} &= \sum_j F_{ij} = 1 \text{ and } F_{ij} \geq 0, \forall i, j \\
&\Rightarrow \mathbb{E}[\text{Var}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y})] \leq \mathbb{E}[\text{Var}(\mathbb{X} | \mathbb{Y})].
\end{aligned}$$

**PROOF.** We first prove a lemma that variance of each class will not increase after attribute graph convolution, i.e.,  $\text{Var}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y} = k) \leq \text{Var}(\mathbb{X} | \mathbb{Y} = k)$ . Denote by  $\text{Cov}(\cdot)$  the covariance matrix of a random vector. Based on our definition of variance

at the beginning of section 4, we have

$$\begin{aligned}
& \text{Var}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y} = k) \\
&= \text{Tr}(\text{Cov}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y} = k)) \\
&= \text{Tr}(\mathbf{F}^\top \text{Cov}(\mathbb{X} | \mathbb{Y} = k) \mathbf{F}) \quad \# \text{ property of covariance} \\
&= \text{Tr}(\text{Cov}(\mathbb{X} | \mathbb{Y} = k) \mathbf{F} \mathbf{F}^\top) \quad \# \text{ cyclic property of trace} \\
&= \sum_{ij} \text{Cov}(\mathbb{X}_i, \mathbb{X}_j | \mathbb{Y} = k) (\mathbf{F} \mathbf{F}^\top)_{ij} \quad \# \text{ property of trace} \\
&\leq \sum_{ij} \sqrt{\text{Var}(\mathbb{X}_i | \mathbb{Y} = k)} \sqrt{\text{Var}(\mathbb{X}_j | \mathbb{Y} = k)} (\mathbf{F} \mathbf{F}^\top)_{ij} \\
&= \sum_{ij} \sigma_i \sigma_j (\mathbf{F} \mathbf{F}^\top)_{ij} \quad \# \sigma \in \mathbb{R}^m, \sigma_i \triangleq \sqrt{\text{Var}(\mathbb{X}_i | \mathbb{Y} = k)} \\
&= \sigma^\top \mathbf{F} \mathbf{F}^\top \sigma \\
&\leq \|\sigma\|_2^2 \quad \# \text{ eigenvalues of } \mathbf{F} \text{ is no more than } 1 \\
&= \sum_i \text{Var}(\mathbb{X}_i | \mathbb{Y} = k) \\
&= \text{Var}(\mathbb{X} | \mathbb{Y} = k).
\end{aligned}$$

Next, we prove the theorem with the above lemma.

$$\begin{aligned}
& \mathbb{E}[\text{Var}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y})] \\
&= \sum_k \text{Pr}(\mathbb{Y} = k) \text{Var}(\mathbf{F}^\top \mathbb{X} | \mathbb{Y} = k) \\
&\leq \sum_k \text{Pr}(\mathbb{Y} = k) \text{Var}(\mathbb{X} | \mathbb{Y} = k) \\
&= \mathbb{E}[\text{Var}(\mathbb{X} | \mathbb{Y})]
\end{aligned}$$

□

## Appendix D

**THEOREM 3.** *If  $\forall F_{ij} \neq 0, \|\mathbf{e}_i - \mathbf{e}_j\|_2 \leq \varepsilon$ , then the distance between  $\mathbf{e}_j$  and  $\hat{\mathbf{e}}_j = \sum_i F_{ij} \mathbf{e}_i$  is also less than or equal to  $\varepsilon$ , i.e.,*

$$\|\mathbf{e}_i - \mathbf{e}_j\|_2 \leq \varepsilon, \forall F_{ij} \neq 0 \Rightarrow \|\mathbf{e}_j - \hat{\mathbf{e}}_j\|_2 \leq \varepsilon,$$

and  $\varepsilon$  can be arbitrarily small with a proper  $F$ .

**PROOF.**

$$\begin{aligned}
\|\mathbf{e}_j - \hat{\mathbf{e}}_j\|_2 &= \left\| \mathbf{e}_j - \sum_i F_{ij} \mathbf{e}_i \right\|_2 \\
&= \left\| \sum_i F_{ij} (\mathbf{e}_j - \mathbf{e}_i) \right\|_2 \quad \# \text{ since } \sum_i F_{ij} = 1 \\
&\leq \sum_i F_{ij} \|\mathbf{e}_j - \mathbf{e}_i\|_2 \quad \# \text{ Cauchy-Schwarz inequality} \\
&\leq \sum_i F_{ij} \varepsilon = \varepsilon
\end{aligned}$$

Next, we prove that there exists such an  $F$  that  $\varepsilon$  is 0. This is equivalent to finding a doubly stochastic  $F$  satisfying  $\sum_i F_{ij} \mathbf{e}_i = \mathbf{e}_j$  for all  $j$ . Given trivial solution  $F = I$ , this equation is solvable. In most real-world attributed networks, the number of attributes is far greater than the number of classes, so the number of variables in this linear system is greater than the number of equations. Given that it is solvable, it must have infinite number of solutions other than  $I$ . Thus,  $\varepsilon$  can be arbitrarily small with a proper  $F$ . □