

Unsupervised Deep Manifold Attributed Graph Embedding

Zelin Zang, Siyuan Li, Di Wu, Jianzhu Guo, Yongjie Xu, Stan Z. Li
Zhejiang University, Hangzhou, China
AI Lab, School of Engineering
Westlake University, Hangzhou, China

April 28, 2021

Abstract

Unsupervised attributed graph representation learning is challenging since both structural and feature information are required to be represented in the latent space. Existing methods concentrate on learning latent representation via reconstruction tasks, but cannot directly optimize representation and are prone to oversmoothing, thus limiting the applications on downstream tasks. To alleviate these issues, we propose a novel graph embedding framework named Deep Manifold Attributed Graph Embedding (DMAGE). A node-to-node geodesic similarity is proposed to compute the inter-node similarity between the data space and the latent space and then use Bergman divergence as loss function to minimize the difference between them. We then design a new network structure with fewer aggregation to alleviate the oversmoothing problem and incorporate graph structure augmentation to improve the representation’s stability. Our proposed DMAGE surpasses state-of-the-art methods by a significant margin on three downstream tasks: unsupervised visualization, node clustering, and link prediction across four popular datasets.

Attributed graphs are graphs with node attributes/features and are widely applied to represent network-structured data, e.g., social networks [10], citation networks [13], recommendation systems [36]. Tasks such as analyzing attributed graphs include node clustering, link prediction, and visualization are challenging in high dimensional non-Euclidean space. Graph embedding methods aim to solve these tasks in a low-dimensional latent space that preserves graph information.

Early approaches such as DeepWalk [24], and GraphSAGE [9] focus on local information but lack a global view. Autoencoder-based methods such as GAE/VGAE [12], AGC [38], and DAEGC [31] obtain embedding by enforcing the reconstruction constraint. However, the reconstruction task only requires preserving information beneficial of reconstruction and fails to optimize embedding directly (shown in Fig. 1) yet does not guarantee a good representation

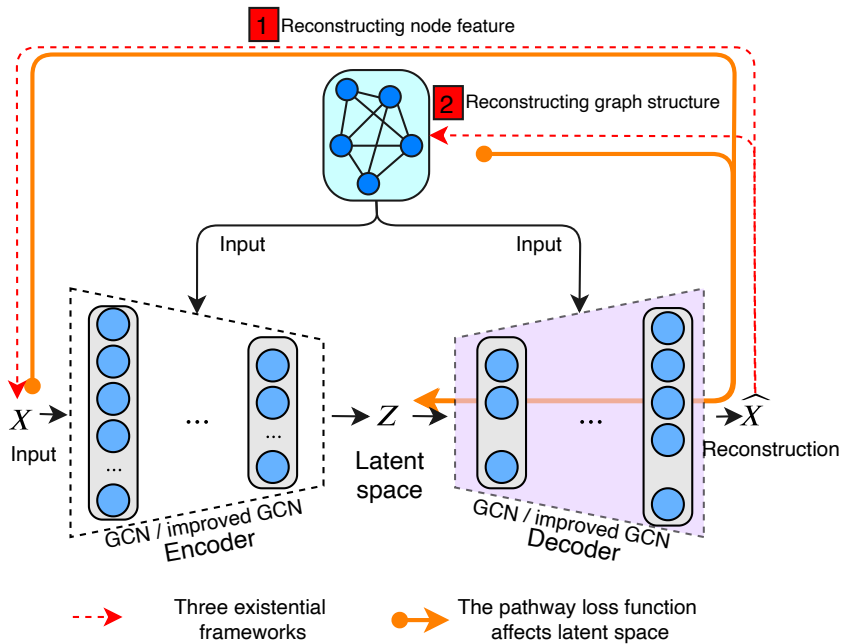


Figure 1: Frameworks of existing neural network-based methods: (1) reconstructing node feature, and (2) reconstructing graph structure. The latent space is optimized through decoders during backpropagation, thus the learned latent representation of all two frameworks are task-dependent and obscured, which leads to the lack of interpretability and subsequent performance guarantees (e.g., generalizability, transferability, and robustness, etc.)

generalizable to other tasks [6]. Graph Convolutional Network (GCN) based methods, e.g. ARGA [22], AGE [6] are usually used to combine graph structure information with node feature information. Such methods suffer from over-smoothing caused by the stacked aggregation layers in the graph network, which loses high-frequency components [34].

To solve the problems mentioned above, we propose Deep Manifold Attributed Graph Embedding (DMAGE), which is a similarity-based framework for attributed graph embedding, as shown in Fig. 2. Unlike reconstruction based methods, we aim to preserve inter-node similarity between non-Euclidean high dimensional space and Euclidean latent space. Firstly, we propose graph geodesic similarity to collect both graph structure information and node feature information, which measures a node-to-node relationship via the graph’s shortest paths. t -distribution is used as a kernel function to fit the neighborhoods between nodes instead of Gaussian distribution to balance intra-cluster and inter-cluster relationship. Secondly, to mitigate the over-smoothing phenomenon in the deeply stacked GCN networks, we use a fully connected neural

network structure with a less aggregated layer called Fully Connected Aggregation Layer (FCA). The proposed network layout allows a deeper network while taking into account the appropriate aggregation operation. Finally, we design a graph structure augmentation scheme during the training process to improve network embedding stability. We randomly add edges with the hop-1 neighbor and drop edges with the hop-2 neighbors to enhance the graph structure without overly changing the graph semantics, forcing the network to obtain a stable embedding mapping. Our contributions can be summarized as:

- Propose a node-to-node geodesic similarity-based graph embedding framework that preserves both graph structure and node feature information.
- Design a network structure to mitigate the oversmoothing problem and a graph structure augmentation method to improve the embedding’s stability.
- Achieve new state-of-the-art performance on visualization, node clustering, and link prediction tasks on benchmark through extensive experiments.

1 Related Works

Attributed graph embedding, also known as network embedding [3] or network representation learning [37], transfers graph nodes into vectors. From the perspective of information exploration, algorithms can be divided into model-free and model-based methods.

We summarize **model-free** methods into four types: (1) Laplacian eigenmaps-based models [20], (2) Local similarity-based models, (3) Matrix factorization-based models, and (4) Nonparametric Bayesian modeling-based models. The local similarity-based approach is based on local similarity with SkipGram model [15] for embedding, which can handle very large datasets, mainly including [24, 26, 8]. [16, 33, 35] is matrix factorization extension that add feature-related regularization terms. [1, 5] model features as latent variables in Nonparametric Bayesian networks.

And the **model-based** methods can be divided into two parts according to whether they use graph convolutional networks (GCN) framework [13] or not.

As a generalization of convolutional operations on graphs, there has been a surge of research interests in graph neural networks in recent years.

(1) Some of these approaches use GCNs as network components and rely on graph autoencoders to fuse graph structure information with feature information. For examples, graph autoencoder (GAE) and variational graph autoencoder (VGAE) [12] learn the node embeddings by using GCN as the encoder, then decode by inner product with cross-entropy loss. As variants of GAE (VGAE), [21] exploits adversarially regularized method to learn more robust node embeddings. [31] further employs graph attention networks [30] to differentiate the importance of the neighboring nodes to a target node.

(2) Some other methods use adjacency matrix as a filter to fuse graph structure information and feature information during forward-propagation of the network. [32] leverages marginalized denoising autoencoder to disturb the structure information. To build a symmetric graph autoencoder, [23] proposes Laplacian sharpening as the counterpart of Laplacian smoothing in the encoder. The authors claim that Laplacian sharpening is a process that makes the reconstructed feature of each node away from the centroid of its neighbors to avoid over-smoothing. [38] propose an adaptive graph convolution method for attributed graph clustering that exploits high-order graph convolution to capture global cluster structure and adaptively selects the appropriate order for different graphs. A better Adaptive Graph Encoder is designed in [6] to smooth and alleviate the high-frequency noises and iteratively strengthen the filtered features. [39] refine the graph topology by strengthening intra-class edges and reducing node connections between different classes based on cluster labels, which better preserves cluster structures in the embedding space.

Also, many schemes for graph node embedding without using GCN networks have emerged due to exploring the GCN over-smoothing problem. [7] uses two networks to learn structural information and feature information separately and then fuses the two kinds of information in latent space. [28] uses a neural network architecture to fuse first-order proximity and second-order proximity of hypergraph networks to complete the node embedding of hypergraph networks. [22] introduced the idea of generative adversarial networks to the graph embedding problem and designed an adversarially regularized graph autoencoder to improve the robustness of embedding. In addition, [27, 11] use the structure of RNN for node embedding related studies.

2 PROPOSED METHOD

In this section, we first formalize the embedding task on attributed graphs. Then we propose our DMAGE algorithm. Specifically, we propose graph geodesic similarity. Then we describe the neural network structure and graph structure augmentation methods. Finally, we define the loss function and algorithmic framework of DMAGE. The overall framework is shown in Fig. 2.

2.1 Problem Formalization

Given an attributed graph $\mathbf{G} = (\mathbf{V}, \mathbf{E}, \mathbf{X})$, where $\mathbf{V} = \{v_1, \dots, v_n\}$ is the vertex set with n nodes in total, \mathbf{E} is the edge set, and $\mathbf{X} = [x_1, \dots, x_n]^T$ is the feature matrix. The topology structure of graph G can be denoted by adjacency matrix \mathbf{A} :

$$\mathbf{A} = \{a_{ij}\} \in \mathcal{R}^{n \times n}, a_{ij} = 1 \text{ if } (v_i, v_j) \in \mathbf{E} \text{ else } 0. \quad (1)$$

We seek to find a set of low-dimensional embeddings $\mathbf{Z} = [z_1, \dots, z_n]^T$ which preserves both structure and feature information of \mathbf{G} and use \mathbf{Z} to accomplish a variety of downstream tasks like node clustering, link prediction, and visualization.

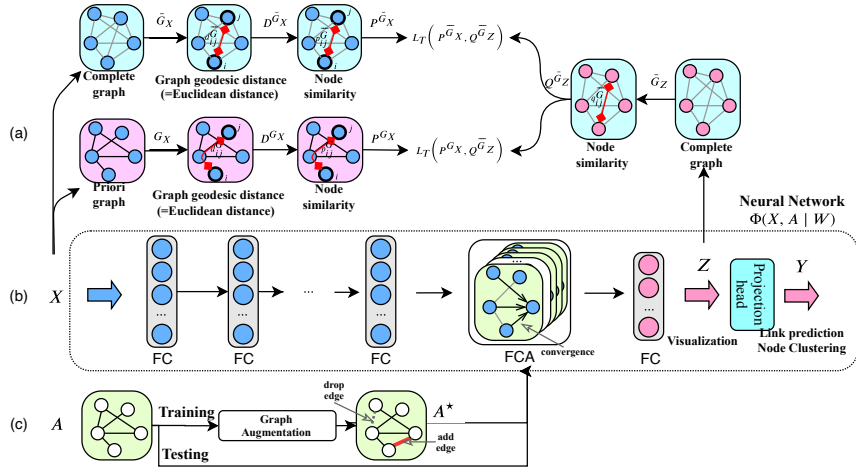


Figure 2: Illustration of Deep Manifold Attributed Graph Embedding (DMAGE) framework: Part (a) describes the fusion process of the structure and feature information in the loss calculation. Graph geodesic distance between nodes is first calculated using input data X and then transformed to similarity on two graphs (Priori graph G_X and complete graph \bar{G}_X) correspondingly. We measure the difference of inter-node similarity between the latent space Z and input space. Part (b) describes our proposed network structure and demonstrates the aggregation based on the adjacency matrix in the Fully Connected Aggregation (FCA) layer. Part (c) illustrates the graph augmentation Process. During the training phase, the augmented graph adjacency matrix is passed to FCA, while during the testing process, we directly use the prior graph adjacency matrix.

2.2 Graph Geodesic Similarity

Graph geodesic distance $\mathbf{D}^G = \{d_{ij}^G \mid i, j = 1, \dots, n\}$ describes the shortest distance between nodes on the graph $\mathbf{G}[\mathbf{V}, \mathbf{E}, \mathbf{X}]$, defined as

$$d_{ij}^G = \begin{cases} \pi(v_i, v_j) & \text{if } v_i, v_j \text{ are connected} \\ \Lambda \cdot \max(\mathbf{D}^G) & \text{otherwise} \end{cases} \quad (2)$$

where $\pi(v_i, v_j)$ is the shortest path based on any distance metric between node v_i and v_j . The distance metric can be chosen from various schemes, such as Euclidean distance, Manhattan distance, cosine distance, etc.

We assume that the distance of unconnected nodes is much greater than the maximum distance between other connected nodes. Thus, we denote the distance between unconnected nodes with the maximum distance of connected nodes multiplied by a large positive constant Λ .

To avoid the adverse effects of outliers and neighborhood inhomogeneity in real data on the characterization of the manifold, we transform the distance d_{ij}^G

to $d_{i|j}^{\mathbf{G}}$ using ρ_i and σ_i^* .

$$d_{i|j}^{\mathbf{G}} = \frac{d_{ij}^{\mathbf{G}} - \rho_i}{\sigma_i^*}, \quad (3)$$

where $\rho_i = \min([d_{i0}^{\mathbf{G}}, \dots, d_{in}^{\mathbf{G}}])$ is deducted from distances of all others nodes to node v_i for the purpose of alleviating possible skewed embedding caused by outliers. $d_{i|j}^{\mathbf{G}}$ and $d_{j|i}^{\mathbf{G}}$ have different ρ_i and σ_i , so $d_{i|j}^{\mathbf{G}} \neq d_{j|i}^{\mathbf{G}}$. In addition, the transformation in distance will make the distance between each node and its nearest node to 0, which means that the similarity is normalized to 1. The optimal σ_i^* is determined by a binary search method, under the objective function.

$$\sigma_i^* = \arg \min_{\sigma_i} : |Q_p - 2 \sum_j \kappa((d_{ij}^{\mathbf{G}} - \rho_i) / \sigma_i, \nu)^2|, \quad (4)$$

where Q_p and ν are hyperparameters that controls the compactness of neighbors, and $\kappa(\cdot)$ is a t -distribution kernel function which map the distance d to similarity

$$\kappa(d, \nu) = \sqrt{2\pi} \cdot \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{d^2}{\nu}\right)^{-\frac{(\nu+1)}{2}}, \quad (5)$$

where ν is degrees of freedom in t -distribution and related discussion is in the ablation study.

We further transform the graph geodesic distance $d_{i|j}^{\mathbf{G}}$ to graph geodesic similarity $p_{i|j}^{\mathbf{G}}$ by

$$p_{i|j}^{\mathbf{G}} = \kappa(d_{i|j}^{\mathbf{G}}, \nu). \quad (6)$$

Then, **Graph geodesic similarity** $p_{i|j}^{\mathbf{G}}$ is symmetrized in the form of joint probability using the following equation:

$$p_{ij}^{\mathbf{G}} = p_{i|j}^{\mathbf{G}} + p_{j|i}^{\mathbf{G}} - 2p_{i|j}^{\mathbf{G}}p_{j|i}^{\mathbf{G}}. \quad (7)$$

In summary, we write the calculation of graph geodesic similarity in matrix form.

$$\mathbf{P}^{\mathbf{G}} = S(\mathbf{G}[\mathbf{V}, \mathbf{E}, \mathbf{X}], \nu) = \{p_{ij}^{\mathbf{G}} | i, j = 1, \dots, n\}. \quad (8)$$

2.3 Network Structure and Graph Augmentation

Excessive aggregation operations and undesired filtering in traditional GNNs bring about the problem of oversmoothing. We use the conventional fully connected layer (FC) and fully connected aggregation layer (FCA) to assemble our neural network to combat the oversmoothing problem. Also, the structure we apply allows for a deeper network and, therefore, better mapping capabilities.

Fully Connected Aggregation layer (FCA):

$$\begin{aligned} \mathbf{Z}^{l+1} &= \text{FCA}^l(\mathbf{Z}^l, \mathbf{A} | \mathbf{W}^l, \mathbf{B}^l) \\ &= \sum_{(i,j) \in \mathbf{E}} \sqrt{|\mathcal{N}(i)|} \sqrt{|\mathcal{N}(j)|} (w_{ij}^l z_j^l + b_j^l). \end{aligned} \quad (9)$$

where $|\mathcal{N}(i)|$ is number of neighbors of node v_i , \mathbf{Z}^l being the feature matrix of layer l . \mathbf{W}^l , \mathbf{B}^l is the weight and bias of layer l respectively. FCA can be viewed as a GCN [13] layer without activation function.

A deeper network structure with L layers $\Phi(\mathbf{X}, \mathbf{A}|\mathbf{W}, \mathbf{B})$ with fewer aggregation operations can be defined as

$$\begin{aligned} \mathbf{Z}^1 &= \text{FC}^1(\mathbf{X}|\mathbf{W}^1, \mathbf{B}^1) \\ \mathbf{Z}^2 &= \text{FC}^2(\mathbf{Z}^1|\mathbf{W}^2, \mathbf{B}^2) \\ &\dots \\ \mathbf{Z}^{L-1} &= \text{FCA}^{L-1}(\mathbf{Z}^{L-2}, \mathbf{A}|\mathbf{W}^{L-1}, \mathbf{B}^{L-1}) \\ \mathbf{Z} &= \text{FC}^L(\mathbf{Z}^{L-1}|\mathbf{W}^L, \mathbf{B}^L). \end{aligned} \tag{10}$$

In general, the forward propagation of our model can be summarized as

$$\mathbf{Z} = \Phi(\mathbf{X}, \mathbf{A}|\mathbf{W}), \tag{11}$$

In complex networks, such as social networks, the relationship of nodes is often considered stable. A graph’s semantic information does not change dramatically due to changes in particular edges. A useful graph embedding should also behave durable against edge changes. Therefore, we introduce a graph augmentation method during the training process to improve the embedding’s stability and reduce the dependence of the model on the correctness of the graph structure. For descriptive convenience, the enhanced edge \mathbf{E}^* is defined as (can be easily transformed into an adjacency matrix \mathbf{A}^* by Eq. (1))

$$\mathbf{E}^* = \mathbf{E} - \mathbf{E}^- + \mathbf{E}^+ \tag{12}$$

where \mathbf{E}^- is a set of edges to be removed, \mathbf{E}^+ is a set of edges to be added.

$$\begin{aligned} \mathbf{E}^- &= \{\mathbf{1}_{r_{ij} > p^-}(i, j) \mid i \in V, j \in \mathbf{H}_1(i)\} \\ \mathbf{E}^+ &= \{\mathbf{1}_{r_{ij} > p^+}(i, j) \mid i \in V, j \in \mathbf{H}_2(i)\} \\ r_{ij} &\in \mathcal{U}(0, 1), \end{aligned} \tag{13}$$

where $\mathbf{H}_1(i)$ and $\mathbf{H}_2(i)$ are the hop-1 neighborhood and hop-2 neighborhood of node i . p^- and p^+ are the probability of augmentation. In this paper, the number of added and removed edges is kept equal in our experiments. During the training process, edge augmentation is performed randomly at the beginning of each epoch under uniform distribution $\mathcal{U}(0, 1)$.

2.4 Bregman divergence and Information Fusion

Based on the graph geodesic similarity, this paper uses Bregman divergence as loss function to minimize the similarity difference between two spaces. Bregman divergence is defined as follows:

$$L_T(x, y, \mid F(\cdot)) = \{F(x) - F(y) - \langle \nabla F(x), x - y \rangle\}. \tag{14}$$

where $F(\cdot)$ is any continuously differentiable, strictly convex function defined on a closed convex set. Our model minimizes the difference in node similarity between the input space and the latent space. We find that it is unnecessary to carefully choose the metric difference ground loss function $F(\cdot)$. We give the following examples of Bregman divergence with different choices.

- Squared Euclidean Distance (DMAGE-SED for short)

$$L_T(\mathbf{x}, \mathbf{y} \mid \|\mathbf{x}\|^2) = \|\mathbf{x} - \mathbf{y}\|_2. \quad (15)$$

- Logistic Loss (DMAGE-LOGI for short)

$$L_T(\mathbf{x}, \mathbf{y} \mid \mathbf{x} \log(\mathbf{x}) + (1 - \mathbf{x}) \log(1 - \mathbf{x})) = \left\{ \mathbf{x} \log \frac{\mathbf{x}}{\mathbf{y}} + (1 - \mathbf{x}) \log \frac{1 - \mathbf{x}}{1 - \mathbf{y}} \right\}. \quad (16)$$

Several decisions of $F(\cdot)$ are compared in the experiment section.

Finally, we complete the definition of overall loss function defined with the fusion of graph structure information and feature information as follows:

$$L = L_T(\mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{x}}}, \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}}) + \alpha \cdot L_T(\mathbf{P}^{\mathbf{G}^{\mathbf{x}}}, \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}}), \quad (17)$$

where α is a balancing parameter between two terms. A bar on top of a graph $\mathbf{G}[\mathbf{V}, \mathbf{E}, \mathbf{X}]$ (short for \mathbf{G}) means the new graph $\overline{\mathbf{G}}[\mathbf{V}, \overline{\mathbf{E}}, \mathbf{X}]$ (short for $\overline{\mathbf{G}}$) is a complete graph with same nodes defined in \mathbf{G} . $\mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{x}}}$, $\mathbf{P}^{\mathbf{G}^{\mathbf{x}}}$ and $\mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}}$ is geodisic similarity defined on each corresponding graph using Eq. (8):

$$\begin{aligned} \mathbf{P}^{\mathbf{G}^{\mathbf{x}}} &= S(\mathbf{G}_X(\mathbf{V}, \mathbf{E}, \mathbf{X}), \nu_{input}), \\ \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{x}}} &= S(\overline{\mathbf{G}}_X(\mathbf{V}, \overline{\mathbf{E}}, \mathbf{X}), \nu_{input}), \\ \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}} &= S(\overline{\mathbf{G}}_Z(\mathbf{V}, \overline{\mathbf{E}}, \mathbf{Z}), \nu_{latent}). \end{aligned} \quad (18)$$

Term $L_T(\mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{x}}}, \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}})$ computes inter-node similarity on the complete graph ignoring any graph structure information. It ensures node feature information is preserved in the embedding. Term $L_T(\mathbf{P}^{\mathbf{G}^{\mathbf{x}}}, \mathbf{P}^{\overline{\mathbf{G}}^{\mathbf{z}}})$ takes graph structure into account while minimizing similarity difference. Combining two loss terms guarantees the latent embedding captures both feature and structure information. We use different t -distribution degrees of freedom coefficients ν_{input} and ν_{latent} to deal with the crowding problem[17] caused by the different dimensions of the two spaces, as we will illustrate in the ablation experiments. Also, this loss function design with multiple information fusion can be extended to multi-graph node embedding.

2.5 Algorithm and Complexity

Algorithm. 1 shows how to train a DMAGE model and obtain the attributed graph representation. The overall complexity is lower than $O(|V|^2)$, where $|V|$

Algorithm 1: The DMAGE algorithm

Input : Graph: $G[V, E, X]$,
Learning rate: l_r , Epochs: E_p , Batch size: B_s ,
Network structure: N_S , Weight hyperparameter: α ,
Degree of freedom: ν_{latent} , Graph augmentation rate: p^+, p^-
Neighbor prior distribution parameter: Q_p
Output : Graph Embedding: Z ,
Initialization
Calculate $d_{ij}^{G^x}$ and $d_{ij}^{\bar{G}^x}$ by Eq. (2)
Calculate σ^* by Eq. (4)
Calculate $p_{ij}^{G^x}$ and $p_{ij}^{\bar{G}^x}$ by Eq. (5) - Eq. (6)
Calculate P^{G^x} and $P^{\bar{G}^x}$ by Eq. (8)
Initialize the network $\phi(\cdot | W)$ with N_S
epoch loop
while $i = 0; i < E_p; i++$ **do**
 # batch loop
 while $b = 0; b < \lfloor |V|/B_s \rfloor; b++$ **do**
 Graph structure augmentation by Eq. (13), $A^* \leftarrow A_{ug}(A, p^+, p^-)$
 Network forward propagation by Eq. (10), $Z \leftarrow \phi(X, A^* | W)$
 Calculate $D^{\bar{G}^z}$ by Eq. (2).
 Calculate $Q^{\bar{G}^z}$ by Eq. (6) and (7)
 Calculate the loss by Eq. (17)
 $L \leftarrow L_T(P^{\bar{G}^x}, Q^{\bar{G}^z}) + \alpha L_T(P^{G^x}, Q^{\bar{G}^z})$
 Update parameters: $W \leftarrow W - l_r \frac{\partial L}{\partial W_{Enc}}$
 end
end
Get the final embedding result, $Z \leftarrow \phi(X, A | W)$

is the number of edges. We divide the algorithm into two parts: the initialization part and the network training part. The initialization part calculates the graph geodesic similarity. It costs $O(|V|^2)$ to compute the pair-wise distance matrix of a fully connected graph. To reduce the complexity, we use a K -nearest neighbor graph instead of the fully connected graph. Thus, the K -nearest neighbor graph's time complexity is determined as $O(K|V|^{1.14})$ via the Nearest-Neighbor-Descent algorithm [14]. The traditional Dijkstra algorithm to solve the shortest path, the time complexity is $O(|E| + |V| \log |V|)$, where $|E|$ is the number of nodes. The network training part uses a well-established mini-batch-based training method, and the complexity is $O(B_s|V|)$, where B_s is the batch size.

3 EVALUATION ON PUBLIC DATASETS

In this section, we evaluate the effectiveness of DMAGE with the state-of-the-art methods on different unsupervised attributed graph representation tasks. Then, we explore the effectiveness of the proposed modules and robustness of hyperparameters.

3.1 Datasets and Experimental Setup

Datasets. We conduct experiments on four widely used network datasets (CORA, CiteSeer, PubMed [13] and Wiki [35]). Features in CORA and CiteSeer are binary word vectors, while in Wiki and PubMed, nodes are associated with tf-idf weighted word vectors. The statistics of the four datasets are shown in table 1.

Table 1: Basis statistics of datasets

Dataset	#Nodes	#Edges	#Features	#Classes
CORA	2,708	5,429	1,433	7
CiteSeer	3,327	4,732	3,703	6
PubMed	19,717	44,338	500	3
Wiki	2,405	17,981	4,973	17

Experimental Setup. For all datasets, the neural network structure is set to be [FC(input dimension) \rightarrow FC(500) \rightarrow FC(250) \rightarrow FCA(250) \rightarrow FC(200)]. The latent space’s dimension is 200. The t -distribution degrees of freedom of the input data are set to $v_{input} = 100$. The dropping edge rate $p^- = 0.01$, the adding edge rate $p^+ = |E^-|/|H_2|$, where $|H_2|$ is the number of the edges that have the potential to be added. $\Lambda = 10$. Cosine distance is used in Eq. (2). We directly set $\sigma_i^* = 1$ and $\rho_i = 0$ in the latent space as a trade-off between the speed and performance. Only the latent space degrees of freedom v_{latent} , weights α and Q_p are selected accordingly for each dataset, details are shown in Appendix. All codes are implemented using PyTorch library and run on NVIDIA v100 GPU.

3.2 Baseline Methods

Table 2: Node clustering performance on PubMed, Cora, CiteSeer and Wiki. Our DMAGE has the highest ACC (accuracy) scores in 11 of 12 metrixs. Concurrent Work are marked with [olive](#).

Method	Input	CORA			CiteSeer			PubMed			Wiki		
		ACC	NMI	F1	ACC	NMI	F1	ACC	NMI	F1	ACC	NMI	F1
k-means	Feature	34.7	16.7	25.4	38.5	17.0	30.5	57.3	29.1	57.4	33.4	30.2	24.5
Spectral-f	Feature	36.3	15.1	25.6	46.2	21.2	33.7	59.9	32.6	58.6	41.3	44.0	25.2
Spectral-g	Graph	34.2	19.5	30.2	25.9	11.8	29.5	39.7	3.5	52.0	23.6	19.3	17.2
DeepWalk(2014) [24]	Graph	46.7	31.8	38.1	36.2	9.7	26.7	61.9	16.7	47.1	38.5	32.4	25.7
DNGR(2016) [4]	Graph	49.2	37.3	37.3	32.6	18.0	44.2	45.4	15.4	17.9	37.6	35.9	25.4
GAE(2016) [12]	Both	53.3	40.7	42.0	41.3	18.3	29.1	64.1	23.0	49.3	17.3	11.9	15.4
VGAE(2016) [12]	Both	56.0	38.5	41.5	44.4	22.7	31.9	65.5	25.1	51.0	28.7	30.3	20.5
MGAE(2017) [32]	Both	63.4	45.6	38.0	63.6	39.8	39.5	43.9	8.2	42.0	50.1	48.0	39.2
ARGE(2018) [21]	Both	64.0	44.9	61.9	57.3	35.0	54.6	59.1	23.2	58.4	41.4	39.5	38.3
DANE(2018) [7]	Both	70.3	54.3	66.8	48.0	25.6	34.8	69.4	31.2	67.58	47.3	46.2	35.9
DGI (2018) [29]	Both	71.3	56.4	68.2	68.8	44.4	65.7	53.3	18.1	18.6	-	-	-
AGC (2019) [38]	Both	68.9	53.7	65.6	67.0	41.1	62.5	69.8	31.6	68.7	47.7	45.2	35.7
DAEGC(2019) [31]	Both	70.4	52.8	68.2	67.2	39.7	63.6	67.1	26.6	65.9	-	-	-
ARGA(2020) [22]	Both	71.1	52.6	69.3	58.1	33.8	52.5	69.0	29.0	67.8	-	-	-
AGE (2020) [6]	Both	71.2	55.9	68.2	56.9	34.8	54.4	-	-	-	51.9	49.4	40.8
GIC (2020) [18]	Both	72.5	53.7	69.4	69.6	45.3	65.4	67.3	31.9	70.4	50.5	48.6	43.8
Ours(DMAGE-SED)	Both	74.2	58.0	69.8	69.6	44.1	66.3	73.3	35.8	73.2	52.3	49.0	46.8
Ours(DMAGE-LOGI)	Both	74.3	57.7	69.7	69.7	43.8	66.3	76.0	36.8	75.8	53.8	50.4	47.6
Ours(DMAGE-SED+LOGI)	Both	74.3	58.2	70.0	69.9	43.5	66.5	74.1	36.0	73.1	52.1	50.0	44.6

We compare our model with eight more algorithms. The baselines can be categorized into three groups:

(1) **Methods using features only.** k-means [20] and Spectral Clustering (Spectral-f) are two traditional clustering algorithms. Spectral-f takes the cosine similarity of node features as input.

(2) **Methods using graph structure only.** Spectral-g is Spectral Clustering with the adjacency matrix as the input similarity matrix. DeepWalk [24] learns node embeddings using SkipGram on generated random walk paths on graphs.

(3) **Methods using both features and graph.** We compared some classical algorithms such as GAE, VGAE [12], MGAE [32], ARGE [21], DANE [7], DGI [29], AGC [38] and DAEGC [31], ARGAs [22]. And some emerging methods: AGE [6] design a non-parametric Laplacian smoothing filter which preserves optimal denoising properties to filter out high-frequency noises. GIC [18] (Concurrent Work) learn a graph representation by maximizes the mutual information between similar nodes.

3.3 Task of Node Clustering

In the node clustering task, the computed embeddings are clustered into $K = \#classes$ clusters with K-means in an unsupervised manner. Then we evaluate clustering performance using external labels. Also, we use three specific choices of $F(\cdot)$ in Bregman divergence in DMAGE: DMAGE-SED, DMAGE-SED, and DMAGE-SED+LOGI (DMAGE-SED+LOGI is the combination of DMAGE-SED and DMAGE-LOGI). We report classification accuracy (ACC), normalized mutual information (NMI), and balanced F-score (F1) in Table 2. The best results for each indicator are shown in bold. For the method with the same testing protocol, we directly use the results reported in their paper. For a fair comparison, we used k-means instead of AGE’s spectral clustering in the official code to re-run the test. Table 2 reports the results obtained when the random seed is set to 1. The average results for the 20 runs are shown in Appendix.

DMAGE outperformed state-of-the-art methods in 11 of all 12 tests on PubMed, CORA, and Kiwi. In particular, performance is improved by over 6% on average on the relatively large PubMed dataset. On Citeseer, DMAGE yields better performance in Acc and F1 scores with a slightly lower NMI result. It’s suggested that by preserving similarity information in the embedding space, DMAGE has the advantage in clustering nodes over other reconstruction and adversarial based approaches. The performance of our proposed three loss functions is similar.

3.4 Task of Link Prediction

Table 3: Link prediction performance on Cora and CiteSeer, and PubMed. Our DMAGE has the highest scores in Cora and CiteSeer. Concurrent Work are marked with [olive](#).

	CORA		CiteSeer		PubMed	
	AUC	AP	AUC	AP	AUC	AP
Spectral-g	84.6 (± 0.01)	88.5 (± 0.01)	80.5 (± 0.01)	85.0 (± 0.01)	84.2 (± 0.02)	87.8 (± 0.01)
DeepWalk(2014) [24]	83.1 (± 0.01)	85.0 (± 0.01)	80.5 (± 0.02)	83.6 (± 0.01)	84.4 (± 0.00)	84.1 (± 0.01)
VGAE(2016) [12]	91.4 (± 0.01)	92.6 (± 0.01)	90.8 (± 0.02)	92.0 (± 0.02)	96.4 (± 0.00)	96.5 (± 0.01)
DGI (2018) [29]	89.8 (± 0.8)	89.7 (± 1.0)	95.5 (± 1.0)	95.7 (± 1.0)	91.2 (± 0.6)	92.2 (± 0.5)
AGE (2020) [6]	92.4 (± 0.004)	93.2 (± 0.003)	92.4 (± 0.003)	93.0 (± 0.003)	96.8 (± 0.001)	97.1 (± 0.001)
GIC (2020) [18]	93.5 (± 0.6)	93.3 (± 0.7)	97.0 (± 0.5)	96.8 (± 0.5)	93.7 (± 0.3)	93.5 (± 0.3)
Ours(DMAGE-SED)	96.0 (± 0.3)	96.8 (± 0.2)	98.17 (± 0.15)	98.32 (± 0.11)	94.2 (± 0.2)	93.8 (± 0.4)
Ours(DMAGE-LOGI)	96.7 (± 0.2)	96.0 (± 0.2)	97.8 (± 0.09)	97.42 (± 0.1)	94.6 (± 0.3)	94.7 (± 0.2)
Ours(DMAGE-SED+LOGI)	95.4 (± 0.1)	96.0 (± 0.3)	97.8 (± 0.11)	97.92 (± 0.08)	94.1 (± 0.3)	95.2 (± 0.4)

In the link prediction task, some edges are hidden randomly in the input graph, and the goal is to predict the existence of hidden edges based using the computed embeddings. We follow the setup described in [12, 21]: We use 5% of edges and negative edges as validation set, 10% of edges and negative edges as the test set. Results are averaged over 20 runs and are shown in table 3. We report the area under the ROC curve (AUC) score [2], which equals the probability that a randomly chosen edge is ranked higher than a randomly chosen negative edge. The average precision (AP) score [25], which is the area under the precision-recall curve. We use the same hyperparameters selected for the node clustering as a demonstration of robustness towards hyperparameters.

As shown in table 3, our method achieves the highest average value in CORA and CiteSeer data with relatively high stability. VGAE and ARGAE perform slightly better than our DMAGE on the PubMed dataset. We discover that VGAE and ARGAE are both autoencoder based method. We suspect that embedding acquired with reconstruction supervision preserves more information of the input space, which helps predict hidden edges.

We also find that DMAGE yields similar results for node clustering and link prediction tasks using three specific Bregman divergence as loss functions (DMAGE-SED: squared Euclidean distance, DMAGE-LOGI: logistic loss, DMAGE-SED+LOGI: their combination). It shows that DMAGE is robust to different $F(\cdot)$ in Bregman divergence. Three different Bregman divergence loss functions have similar performance, so we choose only one loss function (DMAGE-LOGI) for visualization, ablation study, and parameter sensitivity analysis.

3.5 Task of Unsupervised Visualization

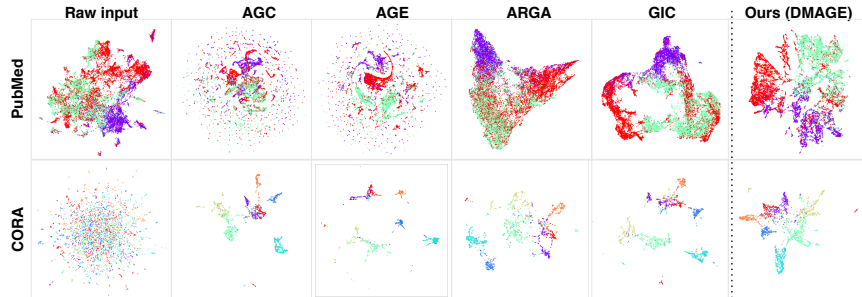


Figure 3: Comparison of visualization results generated by DMT and other methods. Colors represent label categories. Our embedding has clear class boundaries and still preserve the relationship between clusters, which may be a reason for the advantage in both node clustering and link prediction tasks.

To evaluate our proposed method, we visualize the distribution of learned latent representations compared to each node’s input features in two-dimensional space using UMAP [19]. We show our method’s visualization result compared with other methods on the PubMed and CORA in Fig. 3. We generated the

Table 4: Ablation study result on node clustering ACC

	CORA	CiteSeer	PubMed	Wiki
DMAGE	74.15	69.7	76.35	51.51
W/o augmentation	74.11	69.61	71.33	48.64
W/o FCA	73.89	68.37	69.59	50.14
W/o “soft” similarity	65.32	66.78	65.00	43.45

embedding results of the baseline method with the official code. The visualization results show that the DMAGE method can produce more clear category boundaries while preserving the interrelationships between clusters. Visualization results of the all datasets are shown in Appendix.

3.6 Ablation Study

To analyze the effectiveness of each module in the proposed methods, we remove them one by one on the node clustering task in Table 4.

(W/o augmentation): Removing the graph augmentation in the training process and directly using the graph structure given in the aggregation operation dataset. Performance degradation was observed for all datasets. Performance degradation is significant for PubMed and Wiki.

(W/o FCA): We observed a performance decrease on most datasets after removing the FCA layer and replacing it with a regular FC layer. In addition, if all five neural network layers are replaced with GCN layers, the embedding result collapses due to oversmoothing.

(W/o “soft” similarity): Removing graph geodesic similarity and using “hard” connection relations (1 if there is an edge between two nodes, 0 otherwise).

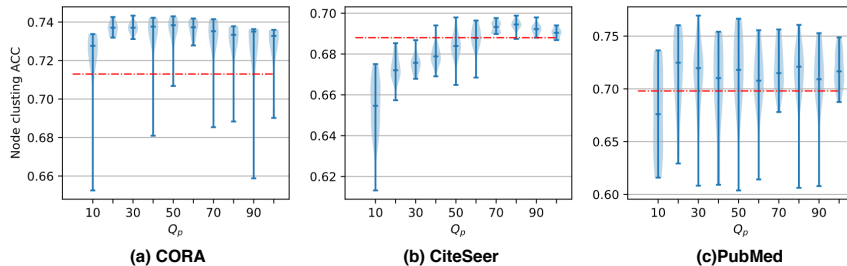


Figure 4: Violin diagram of the node clustering ACC (accuracy) by changing Q_p . The points inside each bar are the median of 20 replicate experiments. When submitting this paper, the highest ACC reported by other methods are marked with the red line (exclude concurrent work).

All three innovations result in performance improvements, among which graph geodesic similarity brings about most performance gain. The similarity

calculation method allows points that are not connected by edges to be related through the graph’s node similarity measurement. The experiments also show that the graph node embedding framework based on the inter-node similarity proposed in this paper can improve performance.

3.7 Parameter Sensitivity Analysis

We illustrate the mechanism of the method and the robustness of the method by varying the hyperparameters in DMAGE.

Effect of ν_{latent} . The degree of freedom ν of the t -distribution controls the sharpness of the kernel function. When high-dimensional data is mapped into a low-dimensional space, a crowding problem[17] may occur. It is essentially due to the lack of expressiveness of the low-dimensional space. The advantage of using the t -distribution as a kernel function is that when a crowding problem is encountered, we can change the shape of the mapping function by increasing the latent space degrees of freedom ν_{latent} (in latent space, this reduces the similarity of neighboring nodes and increases the similarity of non-neighboring nodes), eventually pushing the node’s neighbors farther away and alleviating the crowding problem. To illustrate the above idea, we show the visualization results for four datasets with different ν_{latent} in Fig. 5, and more comments are provided in the accompanying notes below the figure.

Effect of Q_p . The hyperparameter Q_p can help DMT to make a sufficient balance between local and global structure preservation. When Q_p is small, a relatively small σ is obtained by binary search, which results in a smaller graph geodesic similarity p_{ij}^G will show more local information. As Q_p increases, the σ obtained becomes larger, so p_{ij}^G becomes larger, and the embedding will present structural information more globally. Fig. 4 shows the effect of Q_p on the node clustering average ACC over 20 experiments on three datasets (CORA, CiteSeer, and PubMed). The Q_p can effectively affect the performance of the clustering task but is not very sensitive. Experiments show that we can find an interval of Q_p that allows the DMAGE’s average performance to exceed baseline method.

4 Conclusions and further work

We proposed an attributed graph node embedding framework based on the node-to-node similarity. We offer an inter-node similarity ground measure and use Bergman divergence as a loss function to optimize the latent space ground embedding. Also, we propose a depth model with less aggregation to avoid over-smoothing and use graph augmentation methods to improve the embedding’s stability. We studied our designed components and showed their effectiveness. The experiments proved the method’s effectiveness and stability and demonstrated that our innovative parts could substantially lead to enhancements.

In the future, we expect to extend the method to semi-supervised and supervised domains and to be able to solve some more challenging tasks, such as label propagation and label noise.

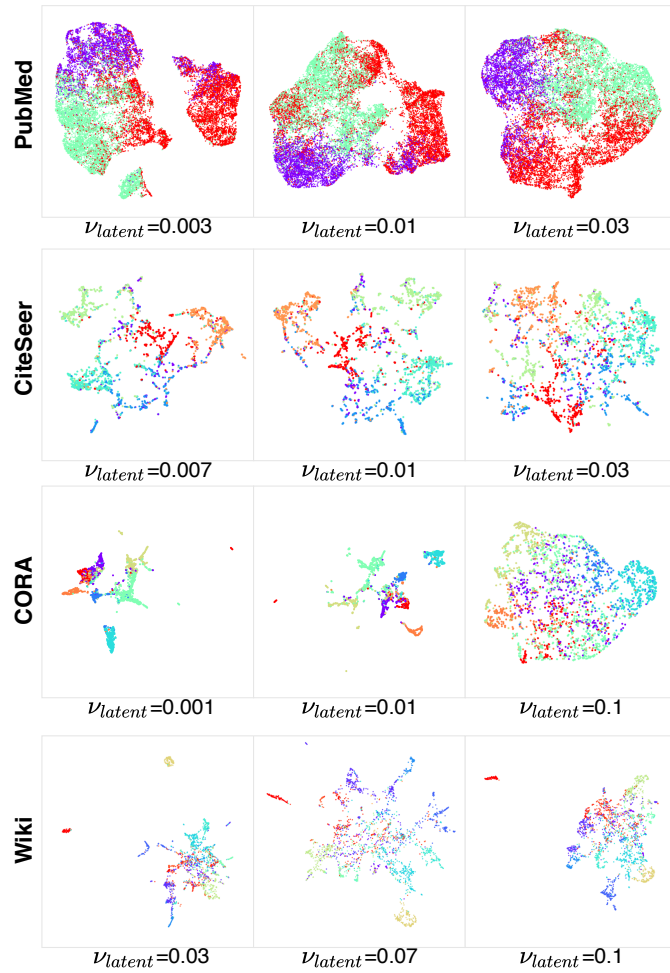


Figure 5: Visualization results of different ν_{latent} . When ν_{latent} is large, the nodes are close to each other, and the clusters will tend to be distributed together. When ν_{latent} is small, nodes that are not similar will be pushed further away, and clusters will tend to be separated from each other.

References

- [1] Aleksandar Bojchevski and Stephan Günnemann. Bayesian robust attributed graph clustering: Joint learning of partial anomalies and group structure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, pages 265–271, AAAI, 2018. AAAI.
- [2] A. P. Bradley. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern R.*, 30(7):1145–1159, July 1997.
- [3] H. Cai, V. W. Zheng, and K. C. Chang. A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications. *IEEE Transactions on Knowledge and Data Engineering*, 30(9):1616–1637, September 2018. Conference Name: IEEE Transactions on Knowledge and Data Engineering.
- [4] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Deep neural networks for learning graph representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, pages 1–9, Phoenix, Arizona USA, 2016. AAAI.
- [5] Jonathan Chang and David Blei. Relational topic models for document networks. In David van Dyk and Max Welling, editors, *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, volume 5 of *Proceedings of Machine Learning Research*, pages 81–88, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 16–18 Apr 2009. PMLR.
- [6] G. Cui and J. Zhou. Adaptive Graph Encoder for Attributed Graph Embedding. In *KDD*, pages 976–985, Virtual Event CA USA, 2020. Association for Computing Machinery.
- [7] Hongchang Gao and Heng Huang. Deep Attributed Network Embedding. In *IJCAI*, pages 3364–3370, Stockholm, Sweden, July 2018. IJCAI.
- [8] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *KDD*, pages 855–864, New York, NY, USA, 2016. Association for Computing Machinery.
- [9] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30, pages 1024–1034, Montreal CANADA, 2017. Curran Associates, Inc.
- [10] M Hastings. Community Detection as an Inference Problem. *Physical review. E, Statistical, nonlinear, and soft matter physics*, 74:035102, October 2006.

- [11] Chaojie Ji, Hongwei Chen, Ruxin Wang, Yunpeng Cai, and Hongyan Wu. Smoothness sensor: Adaptive smoothness-transition graph convolutions for attributed graph clustering, 2020.
- [12] Thomas N. Kipf and Max Welling. Variational graph auto-encoders, 2016.
- [13] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017.
- [14] Dmitry Kobak and George C. Linderman. UMAP does not preserve global structure any better than t-SNE when using the same initialization. preprint, Bioinformatics, December 2019.
- [15] Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*, ICML'14, page II-1188-II-1196, Beijing, China, 2014. JMLR.org.
- [16] Y. Li and C. Sha. Community Detection in Attributed Graphs: An Embedding Approach. In *AAAI*, volume 0, pages 338-345, New Orleans, Louisiana, USA, 2018. AAAI.
- [17] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(Nov):2579-2605, 2008.
- [18] Costas Mavromatis and George Karypis. Graph infoclust: Leveraging cluster-level node information for unsupervised graph representation learning, 2020.
- [19] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction, 2020.
- [20] M. E. J. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical Review E*, 74(3):036104, September 2006. arXiv: physics/0605087.
- [21] S. Pan and R. Hu. Adversarially Regularized Graph Autoencoder for Graph Embedding. In *IJCAI*, pages 2609-2615, Stockholm, Sweden, July 2018. IJCAI.
- [22] S. Pan and R. Hu. Learning Graph Embedding with Adversarial Training Methods. *IEEE Transactions on Cybernetics*, 50(6):2475-2487, 2020.
- [23] Jiwoong Park, Minsik Lee, Hyung Jin Chang, Kyuewang Lee, and Jin Young Choi. Symmetric Graph Convolutional Autoencoder for Unsupervised Graph Representation Learning. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6518-6527, Seoul, Korea (South), October 2019. IEEE.

- [24] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk. *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 0:0–10, Aug 2014.
- [25] Wanhua Su, Yan Yuan, and Mu Zhu. A relationship between the average precision and the area under the roc curve. In *Proceedings of the 2015 International Conference on The Theory of Information Retrieval, ICTIR '15*, page 349–352, New York, NY, USA, 2015. Association for Computing Machinery.
- [26] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line. *Proceedings of the 24th International Conference on World Wide Web*, 0:0–10, May 2015.
- [27] K. Tu and P. Cui. Deep Recursive Network Embedding with Regular Equivalence. In *KDD*, pages 2357–2366, London United Kingdom, 2018. ACM.
- [28] Ke Tu, Peng Cui, Xiao Wang, Fei Wang, and Wenwu Zhu. Structural deep embedding for hyper-networks, 2018.
- [29] P. Veličković and W. Fedus. Deep Graph Infomax. *arXiv:1809.10341 [cs, math, stat]*, 0, December 2018. arXiv: 1809.10341.
- [30] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018.
- [31] C. Wang and S. Pan. Attributed Graph Clustering: A Deep Attentional Embedding. *arXiv*, 0, 2019.
- [32] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jing Jiang. Mgae: Marginalized graph autoencoder for graph clustering. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM '17*, page 889–898, New York, NY, USA, 2017. Association for Computing Machinery.
- [33] X. Wang and D. Jin. Semantic community identification in large attribute networks. In *AAAI*, pages 265–271, Phoenix, 2016. Citeseer.
- [34] F. Wu, T. Zhang, and Souza Jr. Simplifying Graph Convolutional Networks. *arXiv:1902.07153*, 0, 2019.
- [35] C. Yang, Z. Liu, and D. Zhao. Network representation learning with rich text information. In *IJCAI*, volume 2015, pages 2111–2117, Buenos Aires, 2015. IJCAI.
- [36] R. Ying, R. He, and K. Chen. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. *KDD*, 0:974–983, 2018.
- [37] Daokun Zhang, Jie Yin, Xingquan Zhu, and Chengqi Zhang. Network Representation Learning: A Survey. *IEEE Transactions on Big Data*, 6(1):3–28, 2020.

- [38] X. Zhang and H. Liu. Attributed Graph Clustering via Adaptive Graph Convolution. *arXiv:1906.01210*, 0, 2019.
- [39] Y. Zhu and Y. Xu. CAGNN: Cluster-Aware Graph Neural Networks for Unsupervised Graph Representation Learning. *arXiv:2009.01674*, 0, 2020.

A Statistical results of node clustering tasks

Since most of the compared methods only report the highest results, we report the best results under 20 different random seeds in the main paper. The results of the mean and standard deviation are shown in Table. 5.

Table 5: Statistical results of node clustering tasks

	Acc(%)	NMI	F1
CORA	73.56 (± 0.78)	56.86 (± 0.58)	69.38 (± 0.96)
CiteSeer	69.4 (± 0.25)	43.46 (± 0.27)	65.93 (± 0.513)
PubMed	70.13 (± 8.08)	31.94 (± 5.60)	68.01 (± 8.39)
Wiki	47.63 (± 2.13)	48.86 (± 0.99)	43.31 (± 1.56)

The values in parentheses are the standard deviations of 20 seeds. The experiments prove that our method has high stability in all three metrics. However, due to the excessive number of categories in the wiki dataset, the results have large fluctuations.

B Experimental Setup details

For all datasets, the neural network structure is set to be [FC(input dimension) \rightarrow FC(500) \rightarrow FC(250) \rightarrow FCA(250) \rightarrow FC(200)]. The latent space’s dimension is 200. The t -distribution degrees of freedom of the input data are set to $v_{input} = 100$. The dropping edge rate $p^- = 0.01$, the adding edge rate $p^+ = |E^-|/|H_2|$, where $|H_2|$ is the set of all edges that have the potential to be added. $\Lambda = 10$. Cosine distance is used in Eq. (2). To increase the computational speed, we directly set $\sigma_i^* = 1$ and $\rho_i = 0$ in the latent space. All codes are implemented using PyTorch library and run on NVIDIA v100 GPU. PubMed and Wiki’s α are large in alpha parameters because their nodes are associated with tf-idf weighted word vectors

Table 6: Hyperparameter settings for node clustering and link prediction.

	#Nodes	#Edges	#Features	#Classes	ν_{latent}	α	Q_p
CORA	2,708	5,429	1,433	7	0.001	1.0	50
CiteSeer	3,327	4,732	3,703	6	0.003	0.5	80
PubMed	19,717	44,338	500	3	0.003	60	20
Wiki	2,405	17,981	4,973	17	0.02	150	70

C Comparison of visualization of all datasets

To evaluate our proposed method, we visualize the distribution of learned latent representations compared to each node’s input features in two-dimensional space using UMAP [19]. We show our method’s visualization result compared with other methods on the all datasets in Fig. 6. We generated the embedding results of the baseline method with the official code. The visualization results show that the DMAGE method can produce more exact category boundaries while preserving the interrelationships between clusters.

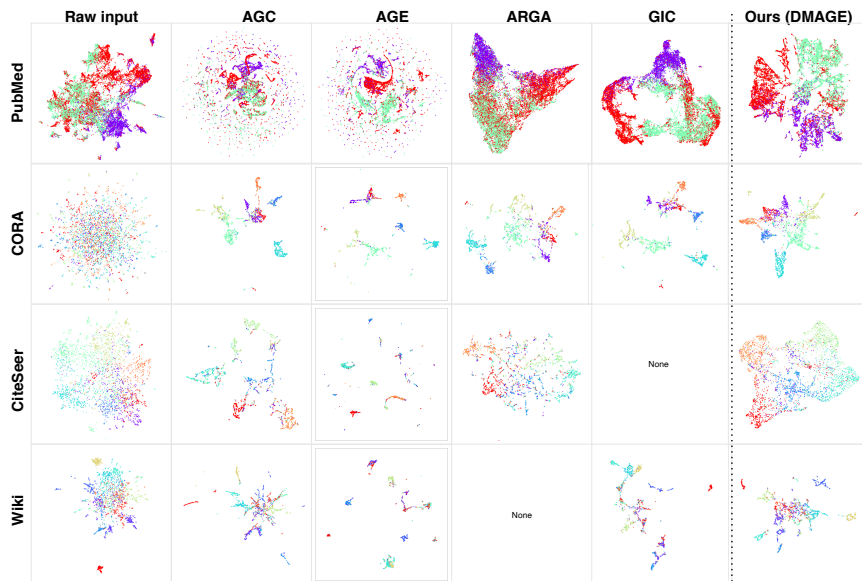


Figure 6: 2D visualization of node representations on Pubmed and Cora using UMAP. The different colors represent different classes. In the two displayed datasets, our method has less blending in the embedding result clusters, and all other methods have more color overlap. In addition, our method does not sever the connection between clusters while ensuring the subclustering, which indicates that the method retains richer information.