

Embedding Graph Auto-Encoder for Graph Clustering

Hongyuan Zhang, Rui Zhang*, *Member, IEEE*, and Xuelong Li, *Fellow, IEEE*

Abstract—Graph clustering, aiming to partition nodes of a graph into various groups via an unsupervised approach, is an attractive topic in recent years. To improve the representative ability, several graph auto-encoder (GAE) models, which are based on semi-supervised graph convolution networks (GCN), have been developed and they achieve good results compared with traditional clustering methods. However, all existing methods either fail to utilize the orthogonal property of the representations generated by GAE, or separate the clustering and the learning of neural networks. We first prove that the relaxed k -means will obtain an optimal partition in the inner-products used space. Driven by theoretical analysis about relaxed k -means, we design a specific GAE-based model for graph clustering to be consistent with the theory, namely Embedding Graph Auto-Encoder (EGAE). Meanwhile, the learned representations are well explainable such that the representations can be also used for other tasks. To further induce the neural network to produce deep features that are appropriate for the specific clustering model, the relaxed k -means and GAE are learned simultaneously. Therefore, the relaxed k -means can be equivalently regarded as a decoder that attempts to learn representations that can be linearly constructed by some centroid vectors. Accordingly, EGAE consists of one encoder and dual decoders. Extensive experiments are conducted to prove the superiority of EGAE and the corresponding theoretical analyses.

Index Terms—Graph clustering, unsupervised representation learning, graph auto-encoder, inner-products space, relaxed k -means.

I. INTRODUCTION

Clustering, which plays important roles in plenty of applications, is one of the most fundamental topics in machine learning [1]–[5]. One specific task, namely graph clustering or node clustering [6], [7], is to group nodes of a given graph, which is common in citation networks, social networks, recommendation systems, *etc.* Some conventional methods (*e.g.*, k -means, DBSCAN [8], AP [9], *etc.*) only utilize the features of nodes but ignore the graph structure. One applicable kind of clustering model is the graph-based model such as spectral clustering [3], [4]. These methods only employ the graph but neglect the features of nodes. Although some models [10] use both features and graph structure, the capacity of the model limits the performance a lot.

With the rise of deep learning, many efforts have been made to promote the capacities of traditional clustering models

via neural networks [11]–[15]. Auto-encoder (AE) [11], as a classical variant of neural network for unsupervised learning, is often employed to perform clustering [13], [16]. Roughly speaking, they employ a multi-layer neural network, namely encoder, to learn non-linear features and reconstruct raw features from the learned features via the decoder. However, most of them separate clustering from training auto-encoder. DEC [17] embeds a self-supervised clustering model into AE and optimizes it by stochastic gradient descent (SGD). Since SGD is applied for clustering, DEC converges slowly. SpectralNet [12], which is not based on AE, intends to extend spectral clustering with neural networks. However, all these methods fail to utilize the structure information provided by graph type data.

A highly related task of graph clustering is network embedding, a fundamental task that aims to learn latent representations (namely embedding) for nodes of a graph. Network embedding [18]–[20] has been applied in diverse applications, such as community networks [21], [22], bioscience, recommendation systems [20], *etc.* Since the network embedding and node clustering are compatible, it is natural to integrate node clustering and network embedding [7], [20], [23]. Specifically, we can utilize the embedding to perform node clustering. Due to the success of CNN, graph convolution neural networks (GCN) have been widely studied for network embedding. A focusing problem is how to extend the convolution operation into irregular data. All existing convolution operations can be classified as spectral based methods [24]–[26] and spatial based methods [27]–[29] according to [30]. On the one hand, spectral-based methods are motivated by the convolution theorem, the Fourier transformation, and the characteristics of the Laplacian operator [31]. In spectral-based models, the spectral domain is regarded as the frequency domain. On the other hand, spatial-based methods do not transform the domain and focus on how to select nodes to perform convolution. For example, PATCH-SAN [27] orders other nodes for a node and chooses top k neighbors to perform convolution. In particular, GCN proposed in [26] combines spectral-based models and spatial-based methods. It employs linear approximation of convolution kernels via Chebyshev polynomials [25]. Although GCNs are usually employed for semi-supervised learning, some graph auto-encoders [6], [7], [32], [33], inspired by auto-encoders, are proposed for unsupervised representation learning and graph clustering. However, all these methods often suffer from overfitting and most of them overlook the crucial characteristics of the generated representations such that some unsuitable clustering methods are applied on embeddings. Besides, they separate the clustering process from the optimization of GAE.

* Rui Zhang is the corresponding author.

Hongyuan Zhang, Rui Zhang, and Xuelong Li are with the School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an 710072, Shaanxi, P. R. China.

E-mail: hyzhang98@gmail.com, ruizhang8633@gmail.com, xuelong_li@nwpu.edu.cn

In this paper, we propose a GAE-based clustering model, Embedding Graph Auto-Encoder (EGAE), for both graph clustering and unsupervised representation learning. The main contributions include:

- 1) We prove that the relaxed k -means can obtain the optimal partition with inner-products if some conditions hold.
- 2) Since the decoder of GAE rebuilds the graph according to inner-products, a specific architecture is designed to cater to the conditions of theoretical analyses. Besides, the learned embedding is well explainable such that EGAE is also a qualified representation learning model.
- 3) Insteading of learning representations and performing clustering separately, the two tasks are processed simultaneously. In particular, the part of clustering can be regarded as a decoder as well. Equivalently, EGAE employs two decoders to capture different information.

II. BACKGROUND

A. Notations

In this paper, all matrices are represented by uppercase words and all vectors are denoted by bold lowercase words. For a matrix M , $\text{tr}(M)$ is the trace of M and $M \geq 0$ means all elements are non-negative. $\text{diag}(\mathbf{m})$ denotes a diagonal matrix whose (i, i) -th entry is m_i . I denotes the identify matrix, $\mathbf{1}_n \in \mathbb{R}^n$ denotes vector whose elements are all 1. If x is a positive scalar, $\text{sign}(x) = 1$. If x is negative, $\text{sign}(x) = -1$. In particular, $\text{sign}(0) = 0$. $\nabla \cdot$ is the gradient operator. In general, we use n , d , and c are used to represent the size of datasets, dimension of data points, and amount of clusters respectively. Given a dataset $\{\mathbf{x}_i\}_{i=1}^n$, it can be denoted by

$$X = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}. \quad (1)$$

The target of clustering is to partition $\{\mathbf{x}_i\}_{i=1}^n$ into c groups, $\{\mathcal{C}_i\}_{i=1}^c$. $|\mathcal{C}_i|$ denotes the amount of samples that are assigned to \mathcal{C}_i . We assume that the graph is stored by an adjacency matrix, A .

B. Convolution on Graph

To apply convolution operation on irregular data, we utilize a graph convolution operation that can be explained as both spectral operators [24] and spatial operators [27], [28]. According to the convolution theorem, the convolution operator can be defined from the frequency domain, which is conventionally named as the spectral domain of graph signals. The adjacency matrix A is employed to represent a graph. Formally, $A_{ij} = 1$ if the i -th point is connected with the j -th one; $A_{ij} = 0$, otherwise. $\mathcal{L} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ is the normalized Laplacian matrix where D is diagonal and $D_{ii} = \sum_{j=1}^n A_{ij}$, while n denotes the amount of nodes in the graph. Formally, a spatial signal of graph $x \in \mathbb{R}^n$ can be transformed into spectral domain by Ux where $\mathcal{L} = U^T \Lambda U$ (eigenvalue decomposition).

If the convolution kernel θ is constrained as a function of Λ , a spectral convolution operator can be defined as [30]

$$f(\mathbf{x}; \theta) = U^T g(\Lambda; \theta) U \mathbf{x}. \quad (2)$$

Suppose that $g(\Lambda; \theta)$ is diagonal and can be approximated by Chebyshev polynomials [25]. If the linear approximation is utilized [26], the convolution can be defined as $f(\mathbf{x}; \theta) = U^T (\theta_0 - \theta_1 \tilde{\Lambda}) U \mathbf{x} = (\theta_0 I - \theta_1 \tilde{\mathcal{L}}) \mathbf{x}$ where $\tilde{\Lambda} = \frac{2}{\lambda_{max}} \Lambda - I$, λ_{max} represents the maximum eigenvalue of \mathcal{L} , and $\tilde{\mathcal{L}} = U^T \tilde{\Lambda} U$. To reduce the amount of parameters to learn and simplify the graph convolutional network, suppose that $\theta_0 = -\theta_1$ and $\lambda_{max} \approx 2$. Accordingly, the above equation becomes $f(\mathbf{x}; \theta) = \theta_0 (I + D^{-\frac{1}{2}} A D^{-\frac{1}{2}}) \mathbf{x}$. Furthermore, we can renormalize the convolutional matrix as

$$\hat{\mathcal{L}} = \hat{D}^{-\frac{1}{2}} (I + A) \hat{D}^{-\frac{1}{2}}, \quad (3)$$

where $\hat{D}_{ii} = \sum_{j=1}^n (I + A)_{ij}$. Therefore, the signal processed by convolution can be rewritten as $f(\mathbf{x}; \theta) = \theta_0 \hat{\mathcal{L}} \mathbf{x}$. If the graph signal is multiple-dimensional and d' convolution kernels are applied, then we have

$$f(X; W) = \hat{\mathcal{L}} X W, \quad (4)$$

where $W \in \mathbb{R}^{d \times d'}$ is the parameters to learn. From the spatial perspective, $\hat{\mathcal{L}}$ is the normalized Laplacian of $I + A$, which is particularly the adjacency of the original graph A with self-loops. Particularly, $\hat{\mathcal{L}} X$ is equivalent to aggregate the information from neighbors, since $\hat{\mathbf{x}}_i = \sum_{j \in \mathcal{N}_i} \mathcal{L}_{ij} \mathbf{x}_j$ where \mathcal{N}_i represents the neighbors of node \mathbf{x}_i and $\hat{\mathbf{x}}_i$ is the i -th column of $\hat{\mathcal{L}} X$. To further improve the efficiency of convolution, some works [34]–[36] employ the high-order Laplacian to equivalently increase the depth of GCNs.

C. Graph Auto-Encoder

Inspired by the conventional auto-encoders, graph auto-encoder [6] (GAE) employs multiple GCN layers to learn embeddings of nodes. Rather than restoring the inputted features from deep representations, GAE intends to reconstruct the graph since connections of nodes can be regarded as weakly supervised information. Specifically, the decoder computes inner-products between any two nodes and then maps them into a probability space to model similarities via the sigmoid function. In GAEs [6], [7], [36], [37], architectures of the network are often asymmetric while a variant [33] designs the symmetric architecture via graph sharpening. For instance, Adversarial Regularized Graph Auto-Encoder [37] incorporates the adversarial learning into GAE to enhance the robustness. GAE with Adaptive Graph Convolution [36] employs the high-order convolution operator to promote the capacity of GAE.

III. EMBEDDING GRAPH AUTO-ENCODER WITH CLUSTERING

In this section, we will first discuss the drawbacks of existing models that employ graph convolution networks for graph clustering. Then we propose *Embedding Graph Auto-Encoder* to address the mentioned problems. The architecture of EGAE is illustrated in Fig. 1.

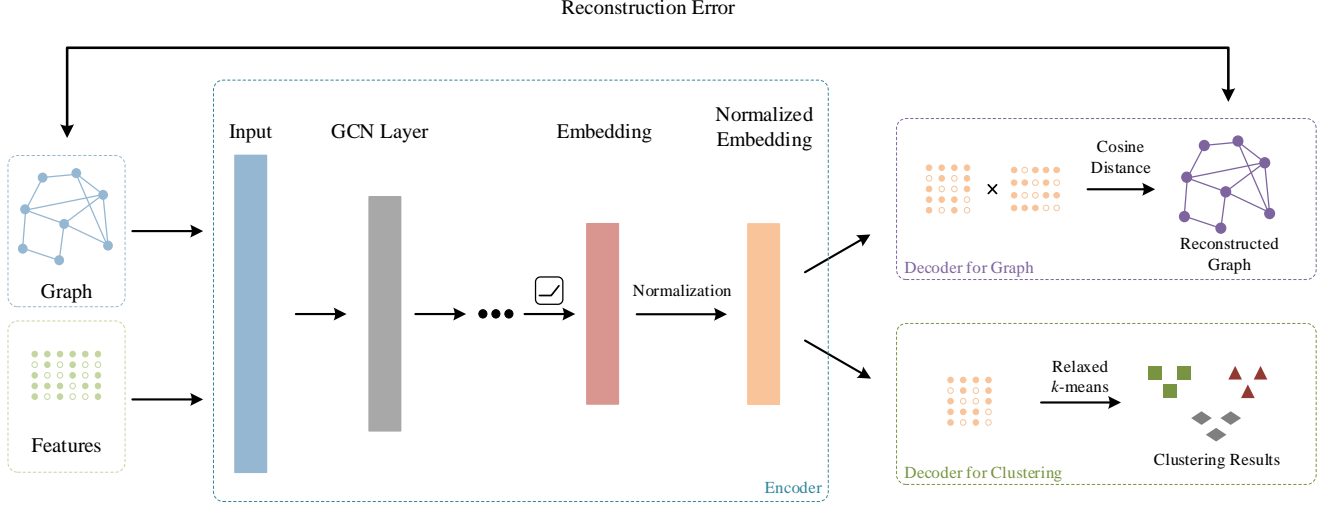


Fig. 1. Conceptual illustration of EGAE. Inputs of EGAE consist of two parts, graph and features. After encoding, data is mapped into a latent feature space which employs inner-products as metric. Note that the last layer of encoder utilizes ReLU as activation function to ensure Assumption 1 holds, and the normalization step is designed due to Assumption 2. Relaxed k -means is embedded into GAE to induce it to generate preferable embeddings. Note that any complicated layers and mechanisms (e.g., attention, pooling, etc.) can be integrated into our framework.

A. Problem Revisit

In network embedding, divergences between nodes are usually measured by inner-product distances (or namely inner-product similarity). Formally,

Definition 1. For nodes i , j , and k , node i is more similar with node j than k under inner-products if $\mathbf{x}_i^T \mathbf{x}_j \geq \mathbf{x}_i^T \mathbf{x}_k$ where \mathbf{x}_i is the representation of node i . $\mathbf{x}_i^T \mathbf{x}_j^T$ is named as the inner-products between \mathbf{x}_i and \mathbf{x}_j .

Compared with the traditional auto-encoders, GAE, which is designed as an unsupervised model for graph data, attempts to reconstruct the graph matrix according to inner-products. However, most GAE-based models, which are used for graph clustering, neglect that the learned embedding scatters in a inner-product space rather than Euclidean space, such that clustering models based on Euclidean distances usually give unsatisfactory results. For instance, k -means are used for node clustering in GAE [6] and AGAE [37] but the adjacency matrix is reconstructed by inner-products. Overall, they pay more attention to how to learn embedding efficiently but fail to focus on the clustering task. Inappropriate clustering methods may provide biased results due to the conflicting metrics.

MGAE [7] and AGC [36] utilize the spectral clustering by constructing a similarity matrix regarding inner-products. Since inner-products may be negative, they simply use absolute values as valid similarities, which also results in a problem. Formally, if the deep representation learned by GAE is denoted by $Z \in \mathbb{R}^{n \times d'}$, the similarity matrix is constructed as

$$S = |ZZ^T|. \quad (5)$$

Note that $s_{ij} = |z_i^T z_j|$. Provided that $z_i^T z_j < 0$ and $z_i^T z_k = 0$, the i -th node is more similar to the k -th one according to the definition of inner-product distance. However, $s_{ij} > s_{ik}$ will

hold according to Eq. (5). The reversed relationships probably mislead the spectral clustering from a theoretical aspect.

B. Motivation: An Approach for Clustering in Inner-Product Space

Before proposing our model formally, we first elaborate the motivation theoretically. K -means, as a fundamental method for clustering, attempts to solve the following problem via a greedy method,

$$\begin{aligned} \min_{\mathbf{f}_j, g_{ij}} \sum_{i=1}^n \sum_{j=1}^c g_{ij} \|\mathbf{x}_i - \mathbf{f}_j\|_2^2 \\ \text{s.t. } g_{ij} \in \{0, 1\}, \sum_{j=1}^c g_{ij} = 1, \end{aligned} \quad (6)$$

where $\{\mathbf{f}_j\}_{j=1}^c$ denotes centroids of c clusters and g_{ij} is the indicator. Specifically speaking, $g_{ij} = 1$ if the i -th point is assigned to the j -th cluster. Otherwise, $g_{ij} = 0$. Clearly, k -means utilize an implicit assumption that Euclidean distance can appropriately depict divergences of data points.

Let g_{ij} be the (i, j) -th entry of matrix G and $F = [\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_c] \in \mathbb{R}^{d \times c}$. Then problem (6) is equivalent to the following problem,

$$\begin{aligned} \min_{F, G} \|X^T - FG^T\|_F^2, \\ \text{s.t. } g_{ij} \in \{0, 1\}, G\mathbf{1}_c = \mathbf{1}_n. \end{aligned} \quad (7)$$

The above problem is hard to solve directly due to the discrete constraint on G . Inspired by the spectral clustering, the indicator can be formulated as

$$\hat{G} = GD_G^{-\frac{1}{2}}, \quad (8)$$

Algorithm 1 Algorithm to optimize problem (14).

Input: Data X .

Calculate c leading left singular vectors, P , of X .

Normalize rows of P .

Perform k -means on normalized P .

Output: Clustering assignments and P .

where D_G is a diagonal matrix which satisfies $(D_G)_{ii} = \sum_{j=1}^n g_{ji} = |\mathcal{C}_i|$. Note that \hat{G} satisfies that $\hat{G}^T \hat{G} = I$. By substituting \hat{G} for G , the objective of k -means can be derived as

$$\begin{aligned} \hat{\mathcal{J}}_c &= \|X^T - F\hat{G}^T\|_F^2 \\ &= \text{tr}(X^T X) - 2\text{tr}(X^T \hat{G} F^T) + \text{tr}(F\hat{G}^T \hat{G} F^T). \end{aligned} \quad (9)$$

Take the derivative of $\hat{\mathcal{J}}_c$ and set it to 0,

$$\nabla_F \hat{\mathcal{J}}_c = -2X^T \hat{G} + 2F\hat{G}^T \hat{G} = 0. \quad (10)$$

Accordingly, we have

$$F = F\hat{G}^T \hat{G} = X^T \hat{G}. \quad (11)$$

Furthermore, \mathcal{J}_c can be written as

$$\hat{\mathcal{J}}_c = \text{tr}(X^T X) - \text{tr}(\hat{G}^T X X^T \hat{G}). \quad (12)$$

Then, problem (6) is converted into the following problem,

$$\begin{aligned} \max_{\hat{G}} \quad & \text{tr}(\hat{G}^T X X^T \hat{G}) \\ \text{s.t.} \quad & \hat{g}_{ij} \in \{0, \frac{1}{|\mathcal{C}_j|}\}, \hat{G}^T \hat{G} = I. \end{aligned} \quad (13)$$

Although the above problem is still intractable as a result of the discrete constraint, it is easy to solve the continuous problem via *singular value decomposition (SVD)*, i.e.,

$$\begin{aligned} \max_P \quad & \text{tr}(P^T X X^T P), \\ \text{s.t.} \quad & P^T P = I, \end{aligned} \quad (14)$$

where P is the continuous indicator matrix. The algorithm is summarized in Algorithm 1.

Although the relaxation seems a common trick from the mathematical perspective, Theorem 1 shows us the relaxed k -means actually assumes that all samples scatter in a inner-product space. Before providing Theorem 1, two assumptions, the basic of our theoretical analysis, are given as follows.

Assumption 1 (Non-Negative Property). *For any two data points \mathbf{x}_i and \mathbf{x}_j , $\mathbf{x}_i^T \mathbf{x}_j \geq 0$ always holds.*

Here we discuss the ideal situation: $\mathbf{x}_i^T \mathbf{x}_j = 0$ if and only if \mathbf{x}_i and \mathbf{x}_j belong to two clusters. Without loss of generality, $Z Z^T$ can be therefore formulated as

$$Q = X X^T = \begin{bmatrix} Q^{(1)} & & & \\ & Q^{(2)} & & \\ & & \ddots & \\ & & & Q^{(c)} \end{bmatrix}, \quad (15)$$

Note that Q can be regarded as a similarity matrix if divergences of data points are measured by inner-products. And $Q^{(i)}$ can

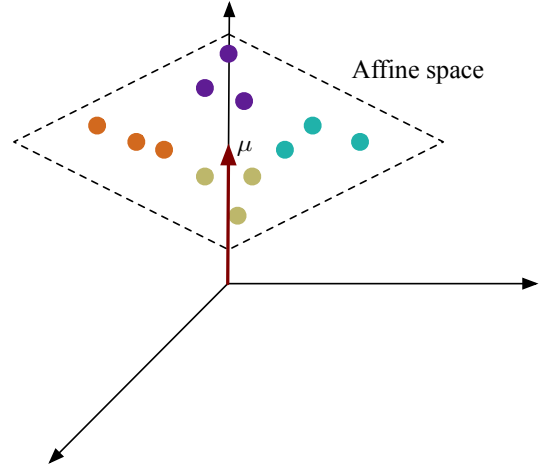


Fig. 2. Illustration of Theorem 2. $\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ denotes the mean vector of samples.

be thus viewed as similarities of samples that belong with the cluster \mathcal{C}_i .

Assumption 2. *Let $\lambda_i^{(m)}$ be the i -th largest eigenvalue of $Q^{(m)}$. For any a and b , $\lambda_1^{(a)} > \lambda_2^{(b)}$ always holds.*

Although Assumption 2 seems too strong and impractical, we will design a GAE-based model in the next subsection to generate qualified embeddings. Based on these assumptions, we state Theorem 1 formally.

Theorem 1. *Suppose that Assumption 1 and 2 hold. Provided that $\mathbf{x}_i^T \mathbf{x}_j = 0$ holds if and only if \mathbf{x}_i and \mathbf{x}_j belong to two different clusters, Algorithm 1 will give an ideal partition.*

Moreover, the following theorem shows the connection between problem (14) and the spectral clustering with normalized cut.

Theorem 2. *Problem (14) is equivalent to the spectral clustering with normalized cut if and only if the mean vector μ is perpendicular to the centralized data. Or equivalently, μ is perpendicular to the affine space that all data points lie in.*

Fig. 2 demonstrates Theorem 2 vividly. In the next subsection, we will show details of EGAE based on the above theoretical knowledge.

C. Embedding Graph Auto-Encoder

In practice, Assumption 1 and 2 do not hold in most cases, especially on the original features. However, we can construct a representative model that can map the original features into valid representations. Based on the theoretical analyses in the previous subsection, we propose a theory-driven model, *Embedding Graph Auto-Encoder (EGAE)*, formally.

a) *Encoder:* Encoder aims to learn embedding Z of raw data via multiple graph convolution layers. The encoder can be constituted by any valid layers and mechanisms such as graph-attention, graph-pooling, etc. To keep simplicity, only GCN layers are employed in our experiments. Specifically speaking,

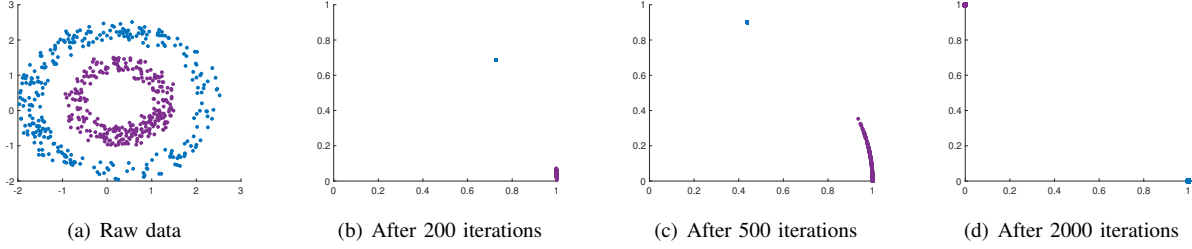


Fig. 3. Illustration of EGAE with $\alpha = 0$ on the two-rings synthetic dataset. After sufficient training, EGAE projects nodes into orthogonal representations.

deep representations provided by the i -th layer, denoted by H_L , is defined as

$$H_i = \varphi_i(\hat{\mathcal{L}}H_{i-1}W_i), \quad (16)$$

where $\varphi_i(\cdot)$ is the activation function of the i -th layer. Let L represent the total amount of layers in the encoder. Due to Assumption 1, the learned embedding Z should satisfy that $ZZ^T \geq 0$. To satisfy this assumption, EGAE employs a simple but effective method by setting

$$\varphi_L(\cdot) = \text{ReLU}(\cdot). \quad (17)$$

Accordingly, the generated embeddings satisfy Assumption 1. The key of EGAE is how to design a model to cater to Assumption 2. The following theorem gives us a feasible scheme.

Theorem 3. Suppose that $\forall i, \|z_i\|_2 = 1$ and $z_i^T z_j = 0$ if and only if the i -th data point and the j -th one belongs to different clusters. Let $\frac{1}{\epsilon}$ ($\epsilon \geq 1$) denote the lower-bound of non-zero entries, i.e.,

$$\begin{cases} z_i^T z_j = 0, & z_i \text{ and } z_j \text{ belong to different clusters,} \\ z_i^T z_j \geq \frac{1}{\epsilon}, & \text{otherwise.} \end{cases} \quad (18)$$

Then, for any valid a and b , $\lambda_1^{(a)} > \lambda_2^{(b)}$ holds if

$$\epsilon < \frac{|\mathcal{C}_{min}|}{|\mathcal{C}_{max}| - 2} + 1, \quad (19)$$

where \mathcal{C}_{min} and \mathcal{C}_{max} represent the largest cluster and smallest one, respectively.

Before continuing the discussion of the encoder, we provide more insights about Theorem 3. Larger $\frac{1}{\epsilon}$ means better representations since samples in the same cluster have high similarities. When $\epsilon \rightarrow 1$, data points of the same cluster are projected into the identical representation. Besides, Eq. 19 demonstrates that the inequality becomes hard to conform to if there is a large gap between the size of the largest cluster and smallest one. Specifically speaking, when $|\mathcal{C}_{max}| \gg |\mathcal{C}_{min}|$, $\frac{|\mathcal{C}_{min}|}{|\mathcal{C}_{max}| - 2} + 1 \rightarrow 1$, which results in a too strict restriction. Provided that $|\mathcal{C}_{min}| = |\mathcal{C}_{max}|$, Eq. 19 becomes

$$\epsilon < 2 < \frac{1}{1 - \frac{2c}{n}} + 1. \quad (20)$$

In other words, the condition of Theorem 3 will hold only if the similarity of two nodes in the same partition exceeds 0.5 in the category-balanced case.

Based on the above discussion, we can process the output of the L -th layer via a normalization step. Formally, the final representation of embedding is defined as

$$z_i = \frac{h_i^{(L)}}{\|h_i^{(L)}\|_2}, \quad (21)$$

where $h_i^{(L)}$ is the i -th column of H_L^T . In summary, the encoder maps nodes into the non-negative part of a hypersphere nonlinearly.

b) Decoder: Decoder intends to reconstruct the adjacency A from the embedding and the reconstruction error is the objective of GAE. The decoder of GAE utilizes a sigmoid function to map $(-\infty, +\infty)$ into a probability space. According to Eq. (17) and (21), it is easy to verify that $z_i^T z_j \in [0, 1]$. In other words, the sigmoid function is no more required. The reconstructed adjacency can be defined as

$$\hat{A} = ZZ^T. \quad (22)$$

The objective to reconstruct the adjacency can be formulated as

$$\mathcal{J}_r = KL(A||\hat{A}). \quad (23)$$

c) Embed Clustering Model as Another Decoder: Instead of separating clustering and embedding learning, EGAE attempts to obtain appropriate representations for the specific model defined in problem (14) via optimizing it and the network simultaneously. Through replacing \hat{G} in $\hat{\mathcal{J}}_c$ (defined in Eq. (9)) by P , the objective of relaxed k -means is formulated as

$$\mathcal{J}_c = \text{tr}(ZZ^T) - \text{tr}(P^T ZZ^T P), \quad (24)$$

and the loss of EGAE is defined as

$$\min_{W_i, P^T P = I} \mathcal{J} = \mathcal{J}_r + \alpha \mathcal{J}_c, \quad (25)$$

where α denotes a tradeoff hyper-parameter. It should be emphasized that \mathcal{J}_c can be employed as an individual decoder. Since one core idea of unsupervised neural networks is to define a loss via reconstruction for training, \mathcal{J}_c provides a novel approach to train neural networks unsupervisedly. On the one hand, \mathcal{J}_c induces the network to produce representations that are preferable for the relaxed k -means. On the other hand, with fixed P and F , \mathcal{J}_c aims to produce embedding to approach FP^T . Therefore, it can be viewed as a decoder as well. Overall, EGAE can be also viewed as a variant of GAE with dual decoders.

Due to the constrained problem in Eq. (25), the simple gradient descent can not be applied directly. As only P

Algorithm 2 Algorithm to optimize problem (25).

Input: Tradeoff parameters α and the adjacency matrix A .
repeat
 repeat
 Calculate the gradients of Eq. (25) w.r.t. $\{W_i\}_{i=1}^L$.
 Update $\{W_i\}_{i=1}^L$ by the gradient descent.
 until convergence or exceeding maximum iterations.
 Compute P and obtain partition via Algorithm 1.
until convergence or exceeding maximum outer-iterations.
Output: Clustering assignments, indicator matrix G and parameters $\{W_i\}_{i=1}^L$.

has a constraint, the alternative method, which consists of gradient descent for $\{W_i\}_{i=1}^L$ and closed-form solution for P , is employed. Details of the optimization are summarized in Algorithm 2.

Remark 1: As we all know, one severe problem that traditional auto-encoders suffer from is that the representations may be inappropriate for specific tasks such as clustering, classification, etc. A reason is the uncertain metric on representations. In other words, we cannot obtain representations under a specific distance metric with the symmetric architecture, which is widely used in auto-encoders. Compared with AE, GAE utilizes the inner-products of deep features to reconstruct the graph via an asymmetric structure. As we focus on graph data in this paper, GAE is a natural choice to fit our theoretical analyses.

Remark 2: Another merit of the normalization defined in Eq. (21) is to unify the inner-products and Euclidean distance. More formally,

$$\|z_i - z_j\|_2^2 = \|z_i\|_2^2 + \|z_j\|_2^2 - 2z_i^T z_j = 2 - 2z_i^T z_j. \quad (26)$$

Therefore, the Euclidean distance of two nodes decreases monotonously with the increasing of inner-products. This property demonstrates the favorable explainability of representations generated by EGAE. Besides node clustering, the embedding is thereby suitable for other tasks such as classification, link prediction, etc.

D. Intuitive Illustration

To understand the motivation of EGAE, we give a visual interpretation via experiment on a toy dataset. Roughly speaking, EGAE can map data into a hypersphere non-linearly. Similarities of data points are measured by inner-products in the hypersphere. Fig. 3 illustrates the effect of EGAE on 2-rings data. In the experiment, we employ a sampling method to generate the graph. Specifically, we connect two nodes that belong to an identical cluster with 90% probability, and there is no link between two different clusters. From the figure, we realize that with the number of iterations increasing, 1) connected samples become more cohesive; 2) samples of different clusters go as orthogonal as possible. In other words, we can obtain orthogonal embedding when the prior adjacency contains sufficient information.

E. Extension: EGAE with Adjacency Sharing

As is shown in Eq. (22), ZZ^T is the reconstructed adjacency. Meanwhile, ZZ^T can be viewed as a similarity matrix in relaxed k -means according to Eq. (14). If the quality of reconstruction is poor, then the performance of relaxed k -means may be affected severely. To alleviate the bias caused by incomplete training, a considerable manner is to employ the prior information provided by the graph to rectify the clustering result.

A substantial clustering model that can explore clustering information of graphs is the spectral clustering. Let $L = D - A$ be the unnormalized Laplacian matrix. Then the rectified model that incorporates the prior graph information can be defined as

$$\min_{F, P^T P = I} \|X^T - FP^T\|_F^2 + \beta \text{tr}(P^T LP), \quad (27)$$

where β is the tradeoff coefficient to control the importance of the prior knowledge. Note that the latter term is the spectral clustering with ratio cut. To solve the above problem, it is equivalent to

$$\min_{P^T P = I} \text{tr}(P^T (\beta L - XX^T) P). \quad (28)$$

Therefore, the optimal solution of problem (27) is given by eigenvalue decomposition of $XX^T + \beta L$. Accordingly, the loss of this extension of EGAE is formulated as

$$\min_{W_i, P^T P = I} \mathcal{J}_r + \alpha (\|Z^T - FP^T\|_F^2 + \beta \text{tr}(P^T LP)). \quad (29)$$

In particular, the adjacency matrix is shared by both GAE and the clustering part. This mechanism helps EGAE to avoid overfitting by the given graph structure. In this paper, we only conduct experiments on EGAE without this extension since it will slow down the optimization. The details can be found in the next subsection.

F. Computational Complexity

Let n_e and d_i be number of edges in the graph and dimension of the i -th layer's output, respectively. To keep notations uncluttered, $d_0 = d$ is used as the dimension of raw features. Since computation of gradient for each node needs all connected samples in each propagation, every step to calculate the gradient of \mathcal{J}_r w.r.t. $\{W_i\}_{i=1}^L$ requires $O(n_e \sum_{i=0}^k (d_i d_{i+1}))$ time. To compute $\nabla_{W_i} \mathcal{J}_c$, extra $O(nd_L c)$ time is needed. The optimization of clustering model needs $O(nd_L c)$ to compute the c leading singular of Z . It should be pointed out that if the extended model is employed, then the optimization of clustering model requires $O(n^2 c)$ at least which extremely slows down the optimization.

Let T_i be the inner-iterations to update parameters of the network with fixed P and T_o be the outer-iterations. Then, the time complexity of EGAE is $O(T_o(T_i(nd_L c + n_e \sum_{i=0}^k (d_i d_{i+1})) + nd_L c))$. Apparently, the embedded clustering model does not increase the complexity of GAE as c is usually small compared with d_i .

IV. PROOFS

In this section, we will provide proofs of the mentioned theorem respectively.

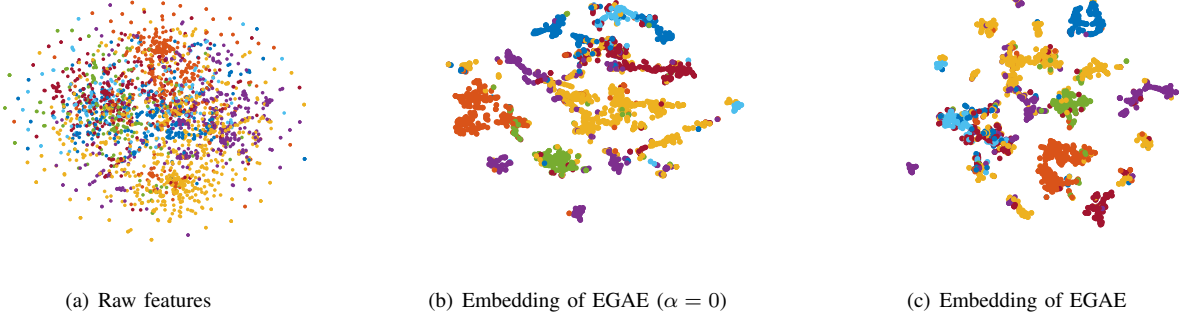


Fig. 4. t-SNE visualization of EGAE on Cora.

A. Proof of Theorem 1

To prove Theorem 1, the following lemma is required.

Lemma 1. For any positive and symmetric matrix K , the most principle component γ satisfies that all elements are not zero and have the same sign. More formally,

$$\forall i, j, \text{sign}(\gamma_i) = \text{sign}(\gamma_j) \quad (30)$$

Proof. According to the definition, we have

$$\begin{cases} \gamma = \arg \max_{\mathbf{x}} \frac{\mathbf{x}^T K \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \\ \lambda_{max} = \max_{\mathbf{x}} \frac{\mathbf{x}^T K \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \end{cases} \quad (31)$$

Suppose that $\exists i, j, \text{sign}(\gamma_i) \neq \text{sign}(\gamma_j)$. We can construct $\hat{\gamma} = |\gamma|$. Note that

$$\mathbf{x}^T K \mathbf{x} = \sum_{i,j} k_{ij} x_i x_j \quad (32)$$

As $K > 0$, we have

$$\hat{\gamma}^T K \hat{\gamma} = \sum_{i,j} k_{ij} |\gamma_i \gamma_j| > \sum_{i,j} k_{ij} \gamma_i \gamma_j = \gamma^T K \gamma \quad (33)$$

Moreover, $\hat{\gamma}^T \hat{\gamma} = \gamma^T \gamma$. Hence, we have

$$\frac{\hat{\gamma}^T K \hat{\gamma}}{\hat{\gamma}^T \hat{\gamma}} = \frac{\hat{\gamma}^T K \hat{\gamma}}{\gamma^T \gamma} > \frac{\gamma^T K \gamma}{\gamma^T \gamma} \quad (34)$$

which leads to a contradiction. Therefore, $\gamma \geq 0$ or $\gamma \leq 0$.

Due to $K > 0$, $\text{tr}(K) = \sum_i k_{ii} > 0$. Clearly, we have $\lambda_{max} > 0$. If $\gamma_k = 0$, then we have

$$\eta = K\gamma = \sum_{i \neq k} \mathbf{k}_i \gamma_i \quad (35)$$

Note that

$$K\gamma = \lambda_{max} \gamma \quad (36)$$

It is not hard to verify that $\eta_k = \lambda_{max} \gamma_k$ if and only if $\gamma = 0$, which leads to a contradiction. Therefore, $\gamma < 0$ or $\gamma > 0$.

Hence, the lemma is proved. \square

Here, we give the complete proof of Theorem 1 based on the above lemma.

Proof of Theorem 1. If samples from different clusters are orthogonal and inner-products of samples from the same cluster are positive, then we have

$$Q_z = ZZ^T = \begin{bmatrix} Q_z^{(1)} & & & \\ & Q_z^{(2)} & & \\ & & \ddots & \\ & & & Q_z^{(c)} \end{bmatrix} \quad (37)$$

where $Q_z^{(i)} = Z_i Z_i^T$ and $Z_i \in \mathbb{R}^{|C_i| \times d}$ consists of samples from the i -th cluster. According to our assumption, we have $Q_z^{(i)} > 0$. With the help of Lemma 1, there exists γ_i which satisfies that

$$\begin{cases} \gamma_i > 0 \\ Q_z^{(i)} \gamma_i = \lambda_{max}^{(i)} \gamma_i \end{cases} \quad (38)$$

Furthermore,

$$Q_z \begin{bmatrix} 0 \\ \gamma_i \\ 0 \end{bmatrix} = \lambda_{max}^{(i)} \begin{bmatrix} 0 \\ \gamma_i \\ 0 \end{bmatrix}. \quad (39)$$

Therefore, a valid relaxed G_0 can be given as

$$P_0 = \begin{bmatrix} \frac{\gamma_1}{\|\gamma_1\|_2} & & & \\ & \frac{\gamma_2}{\|\gamma_2\|_2} & & \\ & & \ddots & \\ & & & \frac{\gamma_c}{\|\gamma_c\|_2} \end{bmatrix}. \quad (40)$$

However, given a orthogonal matrix R , any P which satisfies $P = P_0 R$ is a valid solution. If we normalize rows of P as \hat{P} , then we have

$$\hat{P} = \begin{bmatrix} \pm \mathbf{1}_{|C_1|} r^1 \\ \pm \mathbf{1}_{|C_2|} r^2 \\ \vdots \\ \pm \mathbf{1}_{|C_c|} r^c \end{bmatrix}, \quad (41)$$

where r^i represents the i -th row vector of R .

Accordingly, it will obtain ideal partition if the k -means is performed on rows of \hat{P} . \square

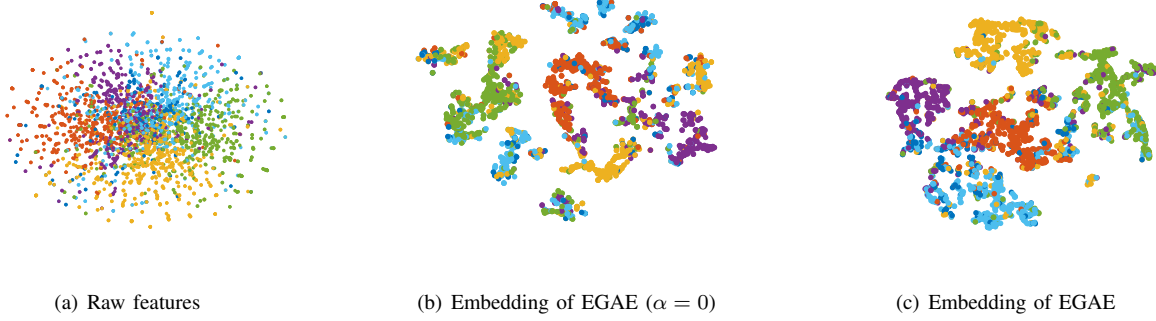


Fig. 5. t-SNE visualization of EGAE on Citeseer.

B. Proof of Theorem 2

Proof. The objective function of normalized cut spectral clustering is given as

$$\begin{aligned} & \min_{P^T P = I} \text{tr}(P^T \mathcal{L} P) \\ \Rightarrow & \min_{P^T P = I} \text{tr}(I - P^T D^{-\frac{1}{2}} A D^{-\frac{1}{2}} P) \\ \Rightarrow & \max_{P^T P = I} \text{tr}(P^T D^{-\frac{1}{2}} A D^{-\frac{1}{2}} P) \end{aligned} \quad (42)$$

Note that $D = \text{diag}(ZZ^T \mathbf{1}) = n \times \text{diag}(Z\boldsymbol{\mu})$ where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$. Let $H = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ be a centralized matrix. Then the centralized data can be represented as $\hat{Z} = HZ = Z - \mathbf{1}_n \boldsymbol{\mu}^T$.

On the one hand, if $\boldsymbol{\mu}$ is perpendicular to the centralized data, then we have

$$\hat{Z}\boldsymbol{\mu} = Z\boldsymbol{\mu} - \mathbf{1}_n \boldsymbol{\mu}^T \boldsymbol{\mu} = Z\boldsymbol{\mu} - \|\boldsymbol{\mu}\|_2^2 \mathbf{1}_n = 0 \quad (43)$$

which means

$$D = n \times \text{diag}(\|\boldsymbol{\mu}\|_2^2 \mathbf{1}_n) = n \|\boldsymbol{\mu}\|_2^2 I_n \quad (44)$$

Hence, problem (42) has the same solution with problem

$$\max \text{tr}(P^T Z Z^T P). \quad (45)$$

On the other hand, if problem (45) is equivalent to problem (42), $D = kI_n$. Therefore, we have

$$HZ^T Z \mathbf{1}_n = HZ^T \cdot n\boldsymbol{\mu} = 0 \quad (46)$$

Accordingly, $\text{rank}(H^T Z) < d$ which means that data points lie in an affine space, and the centralized data is perpendicular to $\boldsymbol{\mu}$.

Hence, the theorem is proved. \square

C. Proof of Theorem 3

Proof. Since $\|\mathbf{z}_i\|_2 = 1$, $(Q_z)_{ii} = 1$ always holds. For each $Q^{(a)}$, we have

$$\sum_i \lambda_i^{(a)} = \text{tr}(Q_z^{(a)}) = |\mathcal{C}_a|. \quad (47)$$

Define

$$\boldsymbol{\gamma} = \frac{1}{\sqrt{|\mathcal{C}_a|}} \mathbf{1}_{|\mathcal{C}_a|}, \quad (48)$$

and we have

$$\begin{aligned} \frac{\boldsymbol{\gamma}^T Q_z^{(a)} \boldsymbol{\gamma}}{\boldsymbol{\gamma}^T \boldsymbol{\gamma}} &= \frac{1}{|\mathcal{C}_a|} \left(\frac{|\mathcal{C}_a|^2}{\epsilon} + |\mathcal{C}_a| \left(1 - \frac{1}{\epsilon}\right) \right) \\ &= \frac{|\mathcal{C}_a|}{\epsilon} + \left(1 - \frac{1}{\epsilon}\right). \end{aligned} \quad (49)$$

The above equation indicates that

$$\lambda_1^{(a)} \geq \frac{|\mathcal{C}_a|}{\epsilon} + \left(1 - \frac{1}{\epsilon}\right). \quad (50)$$

Combine with Eq. (47),

$$\lambda_2^{(a)} \leq (|\mathcal{C}_a| - 1) \left(1 - \frac{1}{\epsilon}\right). \quad (51)$$

Here we get an upper bound for $\lambda_2^{(a)}$ and lower bound for $\lambda_1^{(a)}$. To simplify the discussion, let

$$\begin{cases} \mathcal{B}_\epsilon(|\mathcal{C}_a|) = \frac{|\mathcal{C}_a|}{\epsilon} + \left(1 - \frac{1}{\epsilon}\right), \\ \mathcal{U}_\epsilon(|\mathcal{C}_a|) = (|\mathcal{C}_a| - 1) \left(1 - \frac{1}{\epsilon}\right). \end{cases} \quad (52)$$

It should be pointed out that both $\mathcal{B}_\epsilon(|\mathcal{C}_a|)$ and $\mathcal{U}_\epsilon(|\mathcal{C}_a|)$ increase with $|\mathcal{C}_a|$ monotonously. If the following inequality,

$$1 \leq \epsilon < \frac{|\mathcal{C}_{min}|}{|\mathcal{C}_{max}| - 2} + 1, \quad (53)$$

holds, then we have the following derivation,

$$\begin{aligned} & (\epsilon - 1)(|\mathcal{C}_{max}| - 2) < |\mathcal{C}_{min}| \\ \Leftrightarrow & \left(1 - \frac{1}{\epsilon}\right)(|\mathcal{C}_{max}| - 2) < \frac{|\mathcal{C}_{min}|}{\epsilon} \\ \Leftrightarrow & \left(1 - \frac{1}{\epsilon}\right)(|\mathcal{C}_{max}| - 1) < \frac{|\mathcal{C}_{min}|}{\epsilon} + \left(1 - \frac{1}{\epsilon}\right) \\ \Leftrightarrow & \mathcal{U}_\epsilon(|\mathcal{C}_{max}|) < \mathcal{B}_\epsilon(|\mathcal{C}_{min}|). \end{aligned} \quad (54)$$

Clearly, we have

$$\begin{aligned} \lambda_1^{(a)} &\geq \mathcal{B}_\epsilon(|\mathcal{C}_a|) \geq \mathcal{B}_\epsilon(|\mathcal{C}_{min}|) \\ &> \mathcal{U}_\epsilon(|\mathcal{C}_{max}|) \geq \mathcal{U}_\epsilon(|\mathcal{C}_b|) \geq \lambda_2^{(b)}. \end{aligned} \quad (55)$$

Hence, the theorem is proved. \square

V. EXPERIMENT

In this section, we elaborate on the crucial details of experiments and clustering results of EGAE. To be more intuitive, the t-SNE [38] visualization is shown in Fig. 4, 5, and 6.

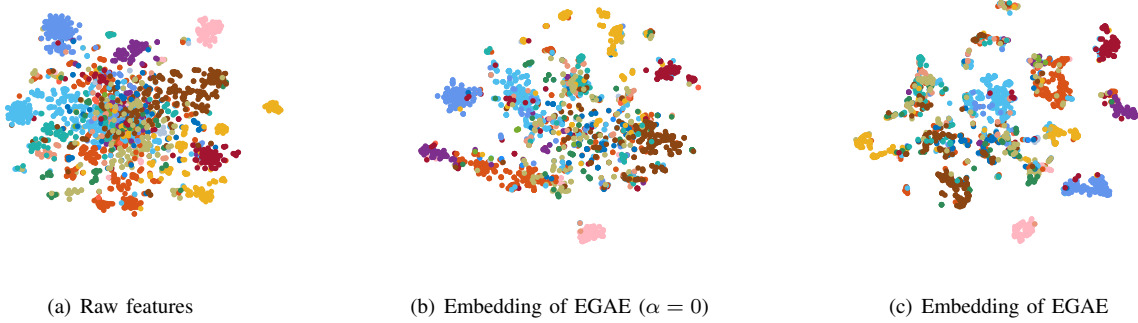


Fig. 6. t-SNE visualization of EGAE on Wiki.

TABLE I
CONCRETE INFORMATION OF DATASETS

Dataset	# Nodes	# Links	# Features	# Classes
Cora	2,708	5,429	1,433	7
Citeseer	3,312	4,732	3,703	6
Wiki	2,405	17,981	4,973	17

A. Benchmark Datasets

To verify the effectiveness of EGAE, experiments are conducted on 3 graph datasets, including *Cora*¹ [41], *Citeseer*¹ [41], and *Wiki*² [23]. These 3 datasets are citation networks where each node represents a publication and link denote a citation between two publications. Features of nodes are usually word vectors which may consist of keywords of papers (or pages). Cora, which contains 7 categories of publications, has 2708 nodes and 5429 links. Citeseer, which contains 6 categories of publications, has 3312 nodes and 4732 links. Wiki has 2405 nodes and 4973 links while all nodes come from 17 classes.

B. Baseline Methods and Evaluation Metrics

Totally ten algorithms are compared in the experiments: *k*-means, Spectral Clustering (SC), Graph-Encoder [39], Deep Walk [20], DNGR [40], TADW [23], GAE (GAE) [6], Adversarial Regularized Graph Autoencoder (ARGE) [37], Adversarial Regularized Variational Graph Autoencoder (ARVGE) [37], and Adaptive Graph Convolution (AGC) [36]. The experimental results of Graph-Encoder, Deep Walk, DNGR, TADW, ARGE, ARVGE, and AGC are reported from original papers. In this paper, totally 3 metrics are employed to testify the performance of various models, including the clustering accuracy (ACC), normalized mutual information (NMI), and adjusted rand index (ARI).

C. Experimental Settings

In our experiments, the encoder is a two-layer GCN. Both activation functions of the two layers are ReLU. To have a stable training process, we employ a LASSO regularization for

parameters of neural networks, and the corresponding tradeoff coefficient is set as 10^{-3} for Cora and Citeseer and 10^{-4} for Wiki. In EGAE, α is an important hyper-parameter that are searched from $\{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$. To avoid the biased indicator matrix P misleading the training of neural networks, we employ EGAE with $\alpha = 0$ as the pretrain for EGAE with a non-zero α . The learning rate of EGAE with $\alpha = 0$ is set as 10^{-3} and the maximum iteration for training is set as 200. After pretraining, the learning rate of fine tuning is set as 10^{-4} . The maximum iteration to update P is set as 30 and the maximum inner iteration to update the neural network is set as 5. For Cora, Citeseer, and Wiki, the encoder is composed of 256-neurons hidden layer and 128-neurons embedding layer. Although the number of clusters in all datasets is not too large and the embedding dimension just needs to equal with c in theory, too small embedding dimension may lead to the slow convergence and difficulty to train EGAE. Accordingly, the embedding layer is set as a 128-neuron graph convolution layer for all datasets. To test the effectiveness of the dual *decoders* EGAE, we also report the performance of EGAE with $\alpha = 0$, the pretrain model. When $\alpha = 0$, the model only has a unique decoder like other GAE-based models. Since clustering results of most compared methods and EGAE depend on *k*-means, all methods are performed 10 times and the means are reported. All codes are implemented by torch-1.3.1 on a Win 10 PC with an NVIDIA GeForce GTX 1660 GPU.

D. Experimental Results

The clustering results of Cora, Citeseer, and Wiki are summarized in Table II. The best results are highlighted by boldface. From this table, we get some conclusions as follows:

- On all datasets, GAE and its extensions outperform other network embedding models, especially two traditional clustering models *k*-means and spectral clustering. Specifically speaking, GAE increases ACC by more than 10 percent compared with Deep Walk.
- Similar to traditional auto-encoder, GAE suffers from overfitting. ARGE and ARVGE absorb adversarial learning into GAE to address this issue and thus promote the performance such that they achieve better performance.
- In particular, EGAE obtains remarkable results on all metrics with the same depth of encoder. Compared with GAE, EGAE is a highly theory-driven model where the

¹lincs.soe.ucsc.edu/data

²github.com/thunlp/TADW

TABLE II
CLUSTERING RESULTS (%)

Methods	Cora			Citeseer			Wiki		
	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
k -means	49.18	32.05	22.81	53.97	30.50	27.81	40.43	42.91	15.03
SC [4]	36.72	12.67	3.11	23.89	5.57	1.00	22.04	18.17	1.46
Graph-Encoder [39]	32.49	10.93	0.55	22.52	3.30	1.00	20.67	12.07	0.49
DeepWalk [20]	48.40	32.70	24.27	33.65	8.78	9.22	38.46	32.38	17.03
DNGR [40]	41.91	31.84	14.22	32.59	18.02	4.29	37.58	35.85	17.97
TADW [23]	56.03	44.11	33.20	45.48	29.14	22.81	30.96	27.13	4.54
GAE [6]	59.61	42.89	34.83	40.84	17.55	18.72	32.85	29.02	7.80
ARGE [37]	<u>64.00</u>	44.90	35.20	<u>57.30</u>	<u>35.00</u>	34.10	38.05	34.45	11.22
ARVGE [37]	63.80	<u>45.00</u>	37.40	54.40	26.10	24.50	38.67	33.88	10.69
AGC [36]	68.92	53.68	—	67.00	41.13	—	47.65	45.28	—
EGAE ($\alpha = 0$)	69.31	51.16	44.81	58.38	33.46	30.89	44.47	41.92	25.92
EGAE	72.42	53.96	47.22	67.42	41.18	43.18	51.52	48.03	33.07

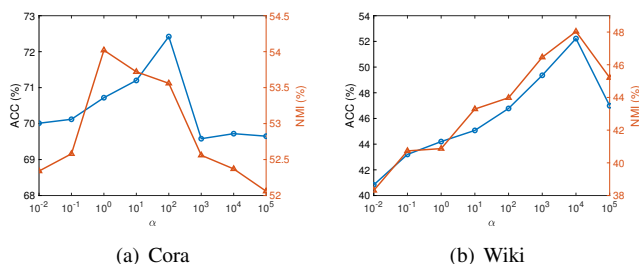


Fig. 7. Impact of α on Cora and Wiki. ACCs and NMIs of EGAE with different α are shown.

extracted deep features are more appropriate for relaxed k -means. On Cora, EGAE improves ACC by 4% and ARI by 9% compared to the second-best method. On Wiki, EGAE improves ACC by 4%, NMI by 3%, and ARI by 16% compared to the second-best method.

- The effect of simultaneously learning GAE and clustering model is apparent, especially on Citeseer and Wiki. For instance, the joint learning improves ACC, NMI, and ARI by 9, 8, and 13 percentage on Citeseer compared with EGAE with $\alpha = 0$. ACC, NMI, and ARI increase by 7, 7, and 8 percentage on Wiki.

E. Parameter Study

In our experiments, we focus on the impact of the tradeoff coefficient α , which plays an important role in balancing two decoders. α varies in $\{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$. Fig. 7 demonstrates different performance with different α on Cora and Wiki. Not surprisingly, too large α may lead to generate biased embedding such that EGAE obtains undesirable results. Combining the results shown in Table II and the two figures, we find that any individual decoder will not lead to the best results.

VI. CONCLUSION

In this paper, we propose a novel GAE-based clustering model, Embedding Graph Auto-Encoder (EGAE) for graph

clustering, driven by rigorous theoretical analyses about the relaxed k -means on inner-products space. Each component of EGAE is supported by theorems. We prove that Algorithm 2 will return an optimal solution in the ideal case. Besides, experiments on synthetic datasets illustrate the theorems vividly. Ablation experiments (EGAE and EGAE with $\alpha = 0$) also testify the effectiveness of the joint learning of embedding and clustering. Or equivalently, the results verify that collaborative training of the two decoders outperforms any single one. The analysis of computational complexity shows that the extra decoder will not result in expensive cost. An extension that tries to rectify the biased indicator via sharing the adjacency matrix is further proposed. As the objective is a constrained problem, the optimization of EGAE is to update GAE by gradient descent and perform clustering alternatively. Extensive experiments prove the superiority of EGAE.

REFERENCES

- [1] J. Dunn, "A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters," vol. 3, no. 3, pp. 32–57, 1973.
- [2] F. Nie, X. Wang, and H. Huang, "Clustering and projected clustering with adaptive neighbors," in *Proc. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2014, pp. 977–986.
- [3] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 8, pp. 888–905, 2000.
- [4] A. Ng, M. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Advances in neural information processing systems*, 2002, pp. 849–856.
- [5] R. Zhang, H. Zhang, and X. Li, "Robust multi-task learning with flexible manifold constraint," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.
- [6] Thomas N Kipf and Max Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016.
- [7] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jiang Jing, "Mgae: Marginalized graph autoencoder for graph clustering," in *the 2017 ACM*, 2017.
- [8] M. Ester, H. Kriegel, J. Sander, and X. Xu, "Density-based algorithm for discovering clusters in large spatial databases with noise," in *KDD*, 1996, pp. 226–231.
- [9] B. Frey and D. Dueck, "Clustering by passing messages between data points," *Science*, vol. 315, no. 5814, pp. 972–976, 2007.
- [10] R. Zhang, F. Nie, M. Guo, X. Wei, and X. Li, "Joint learning of fuzzy k -means and nonnegative spectral clustering with side information," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2152–2162, 2018.
- [11] G. E. Hinton and R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

- [12] Uri Shaham, Kelly Stanton, Henry Li, Ronen Basri, Boaz Nadler, and Yuval Kluger, "Spectralnet: Spectral clustering using deep neural networks," in *International Conference on Learning Representations*, 2018.
- [13] F. Tian, B. Gao, Q. Cui, E. Chen, and T. Liu, "Learning deep representations for graph clustering," in *Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014, pp. 1293–1299.
- [14] R. Zhang, X. Li, H. Zhang, and F. Nie, "Deep fuzzy k-means with adaptive loss and entropy regularization," *IEEE Transactions on Fuzzy Systems*, vol. 28, no. 11, pp. 2814–2824, 2020.
- [15] X. Li, R. Zhang, Q. Wang, and H. Zhang, "Autoencoder constrained clustering with adaptive neighbors," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–7, 2020.
- [16] X. Peng, J. Feng, S. Xiao, W. Yau, T. Zhou, and S. Yang, "Structured autoencoders for subspace clustering," *IEEE Transactions on Image Processing*, vol. 27, no. 10, pp. 5076–5086, 2018.
- [17] J. Xie, R. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," in *International Conference on Machine Learning*, 2016, pp. 478–487.
- [18] Aleksandar Bojchevski, Oleksandr Shchur, Daniel Zügner, and Stephan Günnemann, "Netgan: Generating graphs via random walks," *arXiv preprint arXiv:1803.00816*, 2018.
- [19] Hongwei Wang, Jia Wang, Jialin Wang, Miao Zhao, Weinan Zhang, Fuzheng Zhang, Xing Xie, and Minyi Guo, "Graphgan: Graph representation learning with generative adversarial nets," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018, pp. 2508–2515.
- [20] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena, "Deepwalk: Online learning of social representations," in *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 701–710.
- [21] Xiao Wang, Peng Cui, Jing Wang, Jian Pei, Wenwu Zhu, and Shiqiang Yang, "Community preserving network embedding," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017, pp. 203–209.
- [22] Liyuan Liu, Linli Xu, Zhen Wangy, and Enhong Chen, "Community detection based on structure and content: A content propagation perspective," in *2015 IEEE International Conference on Data Mining*. IEEE, 2015, pp. 271–280.
- [23] Cheng Yang, Zhiyuan Liu, Deli Zhao, Maosong Sun, and Edward Chang, "Network representation learning with rich text information," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015, pp. 2111–2117.
- [24] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- [25] Michael Defferrard, Xavier Bresson, and Pierre Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Advances in neural information processing systems*, 2016, pp. 3844–3852.
- [26] Thomas N Kipf and Max Welling, "Semi-supervised classification with graph convolutional networks," *Proc. of ICLR*, 2017.
- [27] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov, "Learning convolutional neural networks for graphs," in *International conference on machine learning*, 2016, pp. 2014–2023.
- [28] Alessio Micheli, "Neural network for graphs: A contextual constructive approach," *IEEE Transactions on Neural Networks*, vol. 20, no. 3, pp. 498–511, 2009.
- [29] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [30] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu, "A comprehensive survey on graph neural networks," *arXiv preprint arXiv:1901.00596*, 2019.
- [31] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 83–98, 2013.
- [32] Xuelong Li, Hongyuan Zhang, and Rui Zhang, "Adaptive graph auto-encoder for general data clustering," *arXiv preprint arXiv:2002.08648*, 2020.
- [33] Jiwoong Park, Minsik Lee, Hyung Jin Chang, Kyuewang Lee, and Jin Young Choi, "Symmetric graph convolutional autoencoder for unsupervised graph representation learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6519–6528.
- [34] Sami Abu-El-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr Harutyunyan, Greg Ver Steeg, and Aram Galstyan, "Mixhop: Higher-order graph convolutional architectures via sparsified neighborhood mixing," in *International Conference on Machine Learning*, 2019, pp. 21–29.
- [35] Felix Wu, Tianyi Zhang, Amaur Holanda de Souza, Christopher Fifty, Tao Yu, and Kilian Q Weinberger, "Simplifying graph convolutional networks," *Proceedings of Machine Learning Research*, 2019.
- [36] Xiaotong Zhang, Han Liu, Qimai Li, and Xiao-Ming Wu, "Attributed graph clustering via adaptive graph convolution," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 2019, pp. 4327–4333.
- [37] Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang, "Adversarially regularized graph autoencoder for graph embedding," in *IJCAI*, 2018, pp. 2609–2615.
- [38] Van Der Maaten Laurens and Geoffrey Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, no. 2605, pp. 2579–2605, 2008.
- [39] Fei Tian, Bin Gao, Qing Cui, Enhong Chen, and Tie-Yan Liu, "Learning deep representations for graph clustering," in *Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014, pp. 1293–1299.
- [40] Shaosheng Cao, Wei Lu, and Qiongkai Xu, "Deep neural networks for learning graph representations," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [41] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad, "Collective classification in network data," *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.