

---

# Image Representation Learning Using Graph Regularized Auto-Encoders

---

Yiyi Liao, Yue Wang, Yong Liu  
State Key Laboratory of Industrial Control and Technology  
Zhejiang University, China  
yongliu@iipc.zju.edu.cn

## Abstract

It is an important task to learn a representation for images which has low dimension and preserve the valuable information in original space. At the perspective of manifold, this is conduct by using a series of local invariant mapping. Inspired by the recent successes of deep architectures, we propose a local invariant deep nonlinear mapping algorithm, called graph regularized auto-encoder (GAE). The local invariant is achieved using a graph regularizer, which preserves the local Euclidean property from original space to the representation space, while the deep nonlinear mapping is based on an unsupervised trained deep auto-encoder. This provides an alternative option to current deep representation learning techniques with its competitive performance compared to these methods, as well as existing local invariant methods.

## 1 Introduction

Although the dense original image can provide intuitive visual representation, it is well known that this representation may cover the hidden semantic patterns which need to be recognized by those image based learning tasks. On the other side, the performance of machine learning methods is also strongly dependent on the corresponding data representation on which they are applied. Thus image representation becomes a fundamental problem in visual analysis. Given an image data matrix  $X \in \mathbb{R}^{m \times n}$ , each column of  $X$  corresponding to an image, the image representation is to find a representation function  $H = f(X) (H \in \mathbb{R}^{l \times n})$ , which can extract useful information from  $X$ . And each column vector of  $H$  is the representation of an image in this concept space.

At the perspective of dimension reduction, the learned representation should have a lower dimension than the original one, i.e.  $l < m$ , and express the property of data better at the same time. The former is directive, while the later, usually be measured by the performance of clustering  $H$ . In this paper, we aim on this dimension reduction problem. One of the usual frameworks to model this problem is an optimization problem minimizing a cost shown as

$$C = \Phi(X, H) + \Psi(H) \quad (1)$$

where the first term measures the approximation of  $H$  to  $X$ , while the second term, constrains the representation space.

In this paper, we propose an implementation of (1) based on deep learning and manifold, called graph regularized auto-encoder (GAE). The choice of  $\Phi$  is graph regularizer, which constrains the  $H$  to have the similar local geometry of original data space  $X$ . This is motivated by the property that a manifold resembles Euclidean space near each point (Wiki). Regard  $X$  as a manifold, then a neighborhood of each  $x$  has Euclidean property, which we want to be kept in  $H$ . However, whether this can be achieved depends on the choice of  $f$ , which maps  $X$  to  $H$ . It should have enough expressive power to map the original space to the constrained representation space. So we choose deep

network to achieve better performance beyond the existing many interesting linear functions with its nonlinearity. It is also expected that many recent successes on deep learning based approaches in supervised tasks [9, 22] can be extended to the context of unsupervised ones.

The remainder of this paper is organized as follows: In section 2, we will give a brief review of auto-encoder based representation learning and the related works; Section 3 will introduce our graph regularized auto-encoder for image representation learning tasks including both unsupervised conditions and semi-supervised conditions. Extensive experimental results on clustering are presented in Section 4. Finally, we provide a conclusion and future works in Section 5.

## 2 Background

Auto-Encoder [7, 8, 3] is a special neural network, whose input is same as the output of the network. Given a data set  $X = \{x_1, \dots, x_n\} \in \mathbb{R}^{m \times n}$ , each column of  $X$  is a sample vector.  $H \in \mathbb{R}^{l \times n}$  is a feature representation of the original data set  $X$  by an encoder function  $H = f_\theta(X)$ . Normally,  $l < m$ , and  $H$  can be regarded as a low dimensional representation (or subspace) of the original data set  $X$ . And another feature mapping function, which is called decoder, maps from feature space back into input space, thus producing a reconstruction  $Q = q_\theta(H)$ . A reconstruction error function  $L(X, Q)$ , which is also called loss function, is defined, and the set of parameters  $\theta$  of the encoder and decoder are learned simultaneously on the task of reconstructing as well as possible the original input, i.e. attempting to incur the lowest possible reconstruction error of  $L(X, Q)$ <sup>1</sup>.

The most commonly used forms for the decoder and encoder are affine mappings, optimally followed by a non-linearity as:

$$H = f_\theta(X) = s_f(b_H + W_H X) \quad (2)$$

$$Q = q_\theta(H) = s_q(b_Q + W_Q X) \quad (3)$$

$s_f$  and  $s_q$  are the encoder and decoder activation functions, e.g. non-linear functions of sigmoid and hyperbolic tangent or linear identify function etc. Then the set of parameters is  $\theta = \{W_H, b_H, W_Q, b_Q\}$ , and the problem is formally presented as follows:

$$\hat{\theta} = \arg \min_\theta L(X, q_\theta(f_\theta(X))) \quad (4)$$

Formula (4) can be easily solved by the stochastic gradient descent based backpropagation approaches. The auto-encoders are also able to support multiple layers, e.g in Hinton's work [7], which train the encoder network (from  $X$  to  $H_i$ ,  $i$  is the number of the layer) one-by-one using the Restricted Boltzmann Machines and the decoder layers of the network are formed by the inverse of the trained encoder layers, such as  $W_H = (W_Q)^T$  in one layer auto-encoder.

There are also some regularized auto-encoders such as sparse auto-encoders [17, 11, 5, 10], denoising auto-encoders [24, 26, 25] and contractive auto-encoders [19, 18]. It is pointed out that the sparse penalty used in sparse auto-encoder will tend to make only few input configurations can have a low reconstruction error [16], which may hurt the numerical optimization of parameters. The other two kinds of regularized auto-encoders are regarded to make the representation as insensitive as possible with respect to changes in input, which is commonly useful in supervised learning condition, however, it may not provide positive impacts in unsupervised and semi-supervised conditions.

Previous studies have also shown that the locally invariant idea [6] will play an important role in the image representation, especially for those tasks of unsupervised learning and semi-supervised learning. There are many successful manifold learning algorithms, such as Locally Linear Embedding (LLE) [20], ISOMAP [23], and Laplacian Eigenmap [1], which all implement the locally invariant idea that nearby points are likely to have similar embeddings. However, these methods are all linear, which may not provide enough expressive power to find a representation space that can preserve the local geometry.

There are some similar works on graph regularized neural network architecture. [14] proposed a graph regularizer that constrains the similarity between consecutive frames, which shows the human knowledge can be applied in this term. In [13], a graph constrains the data points belonging to the

<sup>1</sup>Normally, the loss function is defined as the Euclidian distance of the two data set, that is  $\|X - Q\|^2$ .

same label is proposed. In this work, the deep network is trained, and then minimizes only the graph regularizer using this network. Both these works use the correct graph regularizer since it is built using the correct supervised or human knowledge information.

Another work similar to us is [6], in which a convolutional neural network is applied to minimize a graph regularizer. In our work, we minimize a combined cost (1) using a fully connected network. The reason is that, the graph in the unsupervised tasks is not completely correct (pair of data point belonging to the different label actually), an introduction of first term in (1) can act as a regularizer avoiding fitting of wrong information. Besides, we do not introduce a mechanism that pulls apart discriminative pairs.

### 3 Graph Regularized Auto-Encoder

The problem is to design the second term  $\Psi(H)$  in (1) to constrain the representation space. In this section, we introduce our geometrical regularization that implement locally invariance in unsupervised and semi-supervised learning.

#### 3.1 Single Layer Auto-Encoder Regularized with Graph

In our Graph regularized Auto-Encoder (GAE), both the decoder and the encoder use sigmoid as their activation functions. Denote sigmoid function as  $S(x) = 1/(1 + e^{-x})$ . Then the encoder and decoder can be presented as follow:

$$H = f_{\theta}(X) = S(W_H X + b_H) \quad (5)$$

$$Q = q_{\theta}(H) = S(W_Q H + b_Q) \quad (6)$$

As the representation should discover the latent structure of the original data, the geometrical structure of the data will be the ideal latent structure in representation learning especially in unsupervised or semi-supervised learning. A reasonable assumption is that if two data points  $x_i, x_j$  are close in the intrinsic geometry of the data distribution, then their corresponding representations,  $h_i, h_j$ , should be also close to each other. This assumption is usually referred to as local invariance assumption [1, 20, 2], which plays an essential role in designing of veracious algorithms, such as dimensionality reduction and semi-supervised learning.

In manifold learning, the local property of the data space are preserved in the reduced representation. But most algorithms considering this problem is linear. In our GAE, we introduce a locality preserved constraint to the nonlinear auto-encoder to better reflect its nature of manifold. Based on formula (4)-(6), we optimize the auto-encoder regularized with graph as follows

$$\hat{\theta} = \arg \min(\|X - Q\|^2 + \lambda \text{tr}(HGH^T)) \quad (7)$$

where  $\lambda$  is a coefficient of the training algorithm,  $\text{tr}(HGH^T)$  is the term of graph regularizer,  $\text{tr}(\cdot)$  denotes the trace of a matrix, and  $G$  is a graph coding the local property of the original data  $X$ . Denote  $v_{ij}$  as the weight for the locality between data sample  $x_i$  and  $x_j$ , and their corresponding representations are  $h_i$  and  $h_j$ . Based on the local invariance assumption, the regularization that requires the points in subspace keeping the same geometrical structure as the original data can be presented as the following weighted formula.

$$\begin{aligned} & \sum_i \sum_j v_{ij} \|h_i - h_j\|^2 \\ &= \sum_i h_i^T \sum_j v_{ij} h_i + \sum_j h_j^T \sum_i v_{ij} h_j - 2 \sum_i \sum_j h_i v_{ij} h_j \\ &= \text{tr}(HD_1 H^T) + \text{tr}(HD_2 H^T) - 2\text{tr}(H V H^T) \\ &= \text{tr}(HGH^T) \end{aligned} \quad (8)$$

where  $v_{ij}$  are entries of  $V$ , and  $V$  is the weight matrix,  $D_1$  and  $D_2$  are diagonal form with  $D_{1,ii} = \sum_j v_{ij}$  and  $D_{2,jj} = \sum_i v_{ij}$ , and  $G = D_1 + D_2 - 2V$ . With this constraint, the local property of the data in original space will be preserved after auto-encoder mapping. The matrix  $G$  has significant expressive power of the structure in the original data space. It is calculated from the weight matrix  $V$ , whose design will be introduced in section 3.3.

Formula (7) can be solved by the Broyden-Fletcher-Goldfarb-Shanno (BFGS) training algorithm. In our GAE, the weight-tying constraint is not used, which means our GAE does not require  $W_H = (W_Q)^T$ .

### 3.2 Multiple-layer Auto-Encoders Regularized with Graph

The auto-encoder was proposed in the context of the neural network, which is later applied to train the deep structure of networks to obtain better expressive power. In data representation, the idea of representing with multiple layers still works. Thus the representation of the original data can be presented with one layer of mapping, as well as multiple-layer mapping. We also implement the locally invariant constraint into the multiple layer auto-encoders by adding the graph regularized terms.

As training all the layers simultaneous in multiple-layer auto-encoders may be stacked, in our multiple-layer GAE, we train the multiple-layer GAE layer-by-layer. We use  $H_i$  to denote the data representation of the  $i$ th layer, and its corresponding decoder is denoted as  $Q_i$ . The input data of the  $i$ th layer is the data representation of  $i - 1$  th layer<sup>2</sup>. That is:

$$H_i = f_{\theta_i}(H_{i-1}) = S(W_{H_i}H_{i-1} + b_{H_i}) \quad (9)$$

$$Q_i = q_{\theta_i}(H_i) = S(W_{Q_i}H_i + b_{Q_i}) \quad (10)$$

Here,  $\theta_i = \{W_{H_i}, W_{Q_i}, b_{H_i}, b_{Q_i}\}$ , and then the objective function of the  $i$ th layer of the GAE is,

$$\hat{\theta}_i = \arg \min(\|H_{i-1} - Q_i\|^2 + \lambda \text{tr}(H_i G_{i-1} H_i^T)) \quad (11)$$

Where  $G_{i-1}$  is the graph regularizer generated from data  $H_{i-1}$ .

Then the Multiple-Layer GAE (ML-GAE) algorithm can be given as follow:

---

#### Algorithm 1: ML-GAE

---

**Input:**  $X$ , total layer number  $j$

**Output:**  $W_{H_1}, b_{H_1}, \dots, W_{H_j}, b_{H_j}$

```

1 for  $i = 1$  to  $j$  do
2   Solve  $\theta_i$  by formula (11), and obtain the  $i$ th layer of data representation  $H_i$ ;
3 end

```

---

### 3.3 Graph Regularizer Design

As mentioned above, the performance of data representation regularized with graph mainly lies on the design of the weight matrix  $V$  since it encodes the local invariance information of the data space. In this section, we will focus on the weight matrix design of supervised learning and unsupervised learning in the context of data representation with auto-encoders.

#### 3.3.1 Unsupervised Learning

In unsupervised learning, the label for the data is unavailable. We can only obtain the structure of the data from the local property of the data samples. There are three kinds of weights employed in our GAE, which are introduced as follows,

- **KNN-graph:** It first constructs the  $k$ -nearest neighbor sets for each data sample. If  $x_i$  lies in  $x_j$ 's  $k$ -nearest neighbor set, the weights  $v_{ij}$  is set as the distance between these two data samples, that is  $\exp(-\|x_i - x_j\|)$ , otherwise  $v_{ij}$  is set to zero.
- **$\epsilon$ -graph:** It first constructs the  $\epsilon$ -neighbor sets for each data sample, the data sample  $x_i$ 's  $\epsilon$ -neighbor set contains all the data samples whose distances to  $x_i$  are less than  $\epsilon$ . If  $x_i$  lies in  $x_j$ 's  $\epsilon$ -nearest neighbor set, the weights  $v_{ij}$  is set as the distance between these two data samples, that is  $\exp(-\|x_i - x_j\|)$ , otherwise  $v_{ij}$  is set to zero.

---

<sup>2</sup>If  $i = 1$ , then the input data is the original data set  $X$ , and the training process is back to single layer GAE.

- **$l_1$ -graph:** The weight setting is considered to resolve the following optimization problem,

$$v_{ij} = \arg \min \|x_i - \sum_{j=1\dots n, j \neq i} v_{ij} x_j\| + \lambda \sum_{j=1\dots n, j \neq i} |v_{ij}|$$

### 3.3.2 Semi-supervised Learning

In semi-supervised learning, the data labels are partially available, which brings some ground truth information to the estimation of the data representation. In our GAE, the graph regularizer design for semi-supervised learning task is similar to the unsupervised learning task. We first construct the  $k$ -nearest neighbor sets or  $\epsilon$ -neighbor sets for the whole data set, and set the weight  $v_{ij}$  to zero if  $x_i$  and  $x_j$  are not neighbors. For the condition that  $x_i$  and  $x_j$  are neighbors, the weights are calculated as follow:

$$v_{ij} = \begin{cases} \exp(-\|x_i - x_j\|) & \text{either } x_i \text{ or } x_j \text{ is unlabeled} \\ 1 & x_i, x_j \text{ have the same label} \\ 0 & x_i, x_j \text{ have different labels} \end{cases} \quad (12)$$

As mentioned before, the weights in the graph are computed by  $\exp(-\|x_i - x_j\|)$ , which is a value less than 1 but larger than 0. Since the labeled data will provide ground truth information, the weights of two samples with the same labels are directly set to 1. And the weights between two samples with different labels are directly set to 0.

The graph constraint constructed with formula (12) is called semi-graph regularizer. Apparently, the marginal value 0 and 1 give the most confident level of the similarity since their corresponding are labeled. With this semi-graph regularizer, both labeled and unlabeled data samples are regarded fairly.

## 4 Experimental Results

In this section, comparison experiments are carried out to demonstrate the performance of our proposed method in tasks of both unsupervised learning and semi-supervised learning. To evaluate image representations learned by different methods quantitatively, the  $k$ -means clustering is applied to the representations learned by different methods. Two metrics, the normalized mutual information(MI) and accuracy(AC) are used to measure the clustering performance. For fair comparison, the dimension of learned representation through all algorithms are set to be the same.

The normalized mutual information(MI) is given in [4], which is a normalized measure to evaluate how similar two sets of clusters are. The accuracy (AC) is used to measure the percentage of correct labels compared to the ground truth label provided by the data set. Specifically, given a data sample  $x_i$  with clustered label and ground truth label  $c_i$  and  $g_i$ , the accuracy is defined as

$$AC = \frac{\sum_i \delta(g_i, \text{map}(c_i))}{n}$$

where  $n$  is total number of samples,  $\delta(a, b)$  is delta function, which outputs 1 when  $a = b$  and outputs 0 otherwise,  $\text{map}(c)$  is the permutation mapping function that maps each clustered label  $c_i$  to the best label provided by the data set. This function is implemented using the code published by [4]. For the normalized mutual information, denote  $C$  and  $C'$  as the set of clusters obtained from the ground truth and our algorithm. We first compute the mutual information as follows,

$$\overline{MI}(C, C') = \sum_{c_i \in C, c'_j \in C'} p(c_i, c'_j) \log \frac{p(c_i, c'_j)}{p(c_i)p(c'_j)}$$

where  $p(c_i)$  and  $p(c'_j)$  are the probabilities that a sample selected from the data set that belong to cluster  $c_i$  and  $c'_j$ ,  $p(c_i, c'_j)$  is the probability that a sample selected from the data set that belong to both  $c_i$  and  $c'_j$ . Then, the normalized mutual information can be computed as

$$MI(C, C') = \frac{\overline{MI}(C, C')}{\max(H(C), H(C'))}$$

where  $H(C)$  and  $H(C')$  are the entropy of  $C$  and  $C'$ . When  $MI = 1$ , the two clusters are identical. when  $MI = 0$ , the two clusters are independent.

The experimental results are all average value of multiple times of random experiment. We set that the reduced representation of different dimension reduction techniques share the same dimensions. The data sets employed for the experiments including: ORL[21], Yale and COIL20, whose statistics are shown in table 1.

Table 1: Statistics of data sets employed in the experiments.

Datasets	#Samples	#Classes	#Samples per class
ORL	400	40	10
Yale	165	15	11
COIL20	1440	20	72

#### 4.1 Variants Comparison

Fine-tuning[7] of a pre-trained deep network can sometimes improve the performance in many supervised learning tasks. However, in unsupervised learning, the weight matrix built on Euclidean distance may include some wrong information, i.e. samples with different labels may be connected. We can compute the error rate as the ratio of wrong connections and the total connections.

We construct a 2 layer auto-encoder, and implement layer-wise pre-training based on the method in section3.2. Then we fine-tune the deep auto-encoder with the single graph regularization, i.e. only the second term of (11). The input data is chosen from COIL20. It has 20 classes. In this experiment, we select 8 classes in random for comparison. So the unsupervised weight matrix is also kind of random based on the data set. Experiment result in table 2 shows that when the weight matrix is constructed with no error, the performance will be promoted with fine-tuning. However, when the weight matrix contains wrong connection, then the result turns out to be worse.

We also conduct an experiment that pre-train the GAE without reconstruction error but only with graph regularization which show in the last two rows in table 2. The result is interesting that the deep architecture even cannot learn a meaningful representation. It may give some insights on the reconstruction error term.

As a result, in the next two sections, we train our deep auto-encoders using layer-wise pre-training with both reconstruction error and graph regularization.

Table 2: Variant experiments.

Model	Train Method	MI	AC	Error Rate
GAE	Pre-train	0.8856	0.9167	3.1%
	Fine-tune	0.7666	0.6689	
GAE	Pre-train	0.8020	0.7465	0%
	Fine-tune	0.9026	0.8732	
GAE No Reconstruction Error	Pre-train	0.0156	0.2060	0%
	Fine-tune	0.2129	0.0134	

#### 4.2 Experiments in Unsupervised Learning

For unsupervised learning task, all samples have no labels. So they are directly fed into the algorithm for measure evaluation. The dimension of the reduced representation is set to the number of classes in the input data set. The methods used for comparison including:

- $k$ -means: this is the baseline method which simply performs clustering in the original feature space.
- PCA: Principal Components Analysis. This is the classical linear technique for dimensional reduction.

- GAE: Graph Auto-Encoder (2 layers). It is a contribution of this paper that introduces the graph constraint to the auto-encoder. In the experiments, our GAE employs the KNN graph. There are two coefficients,  $k$  for the number of neighbors in the KNN graph and  $\lambda$  for the intensity of the graph regularizer. They are all selected by the grid based search.
- SAE: Sparse+Auto-encoder(2 layers)[15]. The sparse constraint is equipped to the auto-encoder, which is a very common constraint choice in the field of auto-encoder. The formula is given as follows

$$\hat{\theta} = \arg \min \|X - \hat{X}\|^2 + \eta \sum_j KL(\rho|\rho_j) \quad (13)$$

where  $\rho$  is the user defined sparsity coefficient,  $\rho_j$  is the average response of the  $j$ th hidden unit for the whole dataset. The penalty term can make the hidden response more sparse. The coefficient  $\eta$  is also selected by grid based search.

- GNMF: Graph regularized Nonnegative Matrix Factorization. It is proposed in [4]. It is a combination of nonnegative constraint and locally invariant constraint. The graph parameter settings are similar to the GAE, which employ the KNN graph and the coefficients  $k$  and  $\lambda$  are selected by the optimal grid search.

The results for whole datasets are shown in table 3, 4, 5. To randomize the experiments, we carry out evaluation with different cluster classes. For each given number of the cluster classes, we random choose the cluster classes from the whole datasets for 5 times. One can see that graph regularized auto-encoder achieves the best performance. Although the GNMF and GAE are employed the same encodes on the locally invariant information, the auto-encoder which implement the nonlinear sigmoid scaling on the deep structure will performance better than the nonnegative matrix factorization based approach.

Table 3: Results of the unsupervised learning tasks on ORL.

Class	5	10	15	20	25	30	35	40	Average
GAE MI	<b>0.8955</b>	<b>0.8967</b>	<b>0.8694</b>	<b>0.8531</b>	<b>0.8477</b>	<b>0.8315</b>	<b>0.8282</b>	<b>0.8344</b>	<b>0.8571</b>
SAE MI	0.8632	0.8670	0.8438	0.8338	0.8423	0.8054	0.8062	0.7933	0.8319
GNMF MI	0.8464	0.7958	0.7634	0.7548	0.7759	0.7724	0.7627	0.7509	0.7778
PCA MI	0.8436	0.7969	0.7740	0.7492	0.7411	0.7354	0.7505	0.7500	0.7676
Kmeans MI	0.6921	0.7392	0.6907	0.6939	0.6972	0.7088	0.7186	0.7078	0.7060
GAE AC	<b>0.9333</b>	<b>0.8167</b>	0.7733	0.7483	<b>0.7400</b>	<b>0.6878</b>	<b>0.6876</b>	<b>0.6839</b>	<b>0.7589</b>
SAE AC	0.9067	0.8067	<b>0.7800</b>	<b>0.7533</b>	0.7373	0.6867	0.6543	0.6217	0.7433
GNMF AC	0.8840	0.7480	0.6773	0.6310	0.6272	0.6173	0.5714	0.5748	0.6664
PCA AC	0.8800	0.7380	0.6800	0.6190	0.5976	0.5567	0.5777	0.5625	0.6514
Kmeans AC	0.7400	0.6720	0.5933	0.5690	0.5552	0.5300	0.5354	0.5050	0.5875

Table 4: Results of the unsupervised learning tasks on Yale.

Class	3	6	9	12	15	Average
GAE MI	<b>0.7283</b>	<b>0.5748</b>	<b>0.5467</b>	<b>0.5395</b>	<b>0.5689</b>	<b>0.5916</b>
SAE MI	0.6140	0.5154	0.5386	0.5171	0.5350	0.5440
GNMF MI	0.5175	0.5111	0.5130	0.5034	0.4518	0.4994
PCA MI	0.3650	0.3883	0.4686	0.4859	0.5103	0.4436
Kmeans MI	0.4198	0.3082	0.4035	0.4251	0.4532	0.4020
GAE AC	<b>0.8889</b>	<b>0.6515</b>	0.5623	<b>0.5177</b>	<b>0.5313</b>	<b>0.6303</b>
SAE AC	0.8485	0.6313	<b>0.5993</b>	0.4975	0.4828	0.6119
GNMF AC	0.7758	0.5818	0.5414	0.4818	0.4162	0.5594
PCA AC	0.6424	0.4909	0.4970	0.4652	0.4412	0.5073
Kmeans AC	0.6303	0.4303	0.4222	0.4182	0.4048	0.4612

Table 5: Results of the unsupervised learning tasks on COIL20.

Class	6	8	10	12	14	16	20	Average
GAE MI	<b>0.9739</b>	0.8888	<b>0.8762</b>	<b>0.8684</b>	<b>0.8657</b>	<b>0.8552</b>	<b>0.8502</b>	<b>0.8826</b>
SAE MI	0.9264	<b>0.8953</b>	0.8527	0.8280	0.8344	0.8078	0.8013	0.8494
GNMF MI	0.8690	0.8469	0.8696	0.8302	0.8379	0.8460	0.8449	0.8492
PCA MI	0.7350	0.7544	0.7790	0.7988	0.7686	0.7894	0.7817	0.7724
Kmeans MI	0.6550	0.7490	0.7714	0.7720	0.7382	0.7392	0.7354	0.7372
GAE AC	<b>0.9846</b>	<b>0.8779</b>	<b>0.8495</b>	<b>0.8279</b>	<b>0.8056</b>	<b>0.7888</b>	<b>0.7981</b>	<b>0.8475</b>
SAE AC	0.9483	0.8709	0.8352	0.7731	0.7665	0.7237	0.7255	0.8062
GNMF AC	0.8690	0.8097	0.8358	0.7692	0.7579	0.7566	0.7172	0.7879
PCA AC	0.7435	0.7035	0.7553	0.7461	0.6917	0.7148	0.6871	0.7203
Kmeans AC	0.7106	0.7372	0.7472	0.7338	0.6579	0.6646	0.6096	0.6944

Table 6: Results of the semi-supervised learning tasks for ORL.

Class	5	10	15	20	25	30	35	40	Average
GAE MI	<b>0.9242</b>	<b>0.9335</b>	<b>0.8996</b>	<b>0.8872</b>	<b>0.8772</b>	<b>0.8787</b>	<b>0.8681</b>	<b>0.8694</b>	<b>0.8923</b>
CNMF MI	0.8257	0.7998	0.8275	0.8161	0.8007	0.7943	0.7885	0.7873	0.8050
GAE AC	<b>0.9533</b>	<b>0.9167</b>	<b>0.8400</b>	<b>0.7783</b>	<b>0.7800</b>	<b>0.7756</b>	<b>0.7543</b>	<b>0.7325</b>	<b>0.8163</b>
CNMF AC	0.8680	0.7580	0.7520	0.7290	0.6824	0.6620	0.6354	0.6250	0.7140

### 4.3 Experiments in Semi-supervised Learning

For semi-supervised learning task, a small part of the samples are labeled. In this experiment, these labeled samples are selected in random. For COIL20, 10% samples are labeled in each class, so there are 7 labeled samples in each class. For ORL and Yale, 20% are labeled, then 2 samples are labeled in each class<sup>3</sup>. Still referring to the unsupervised learning experiments, the dimension of the learned representation is equal to the number of classes in the data set. The comparison methods used in this experiments including:

- CNMF: constrained NMF. It is proposed in [12]. In their framework, the samples with the same label are required to have the same representation in the reduced space. There is no user defined parameters either.
- SGAE: Semi-Graph regularized Auto-Encoder (2 layers). It is a representation learning algorithm proposed in this paper consisting of the auto-encoder regularized by the semi-graph regularizer presented in section 3.3.2. The parameters include the intensity of the graph constraint,  $\lambda$  and the number of neighbors in the KNN graph,  $k$ . Similarly to the GAE in unsupervised learning experiment, these two parameters are selected by optimal grid search.

The clustering results on all classes of the datasets are shown in table 6, 7, 8. Similarly to the randomize experiment of the unsupervised learning, we also conduct the randomize experiment

<sup>3</sup>Here one labeled sample is meaningless to both CNMF and SGAE, so we label 20% of the samples in each class.

Table 7: Results of the semi-supervised learning tasks for Yale.

Class	3	6	9	12	15	Average
GAE MI	<b>0.6400</b>	<b>0.7322</b>	<b>0.6336</b>	<b>0.6126</b>	<b>0.6550</b>	<b>0.6547</b>
CNMF MI	0.5915	0.5733	0.5620	0.5710	0.5543	0.5704
GAE AC	<b>0.8384</b>	<b>0.7172</b>	<b>0.6604</b>	<b>0.6162</b>	<b>0.5818</b>	<b>0.6828</b>
CNMF AC	0.7636	0.6636	0.5939	0.5439	0.4949	0.6120



Table 8: Results of the semi-supervised learning tasks for COIL20.

Class	6	8	10	12	14	16	20	Average
GAE MI	<b>0.9182</b>	<b>0.9406</b>	<b>0.9281</b>	<b>0.9131</b>	<b>0.9215</b>	<b>0.8908</b>	<b>0.8642</b>	<b>0.9109</b>
CNMF MI	0.7861	0.8129	0.7580	0.7606	0.7576	0.7904	0.7440	0.7728
GAE AC	<b>0.9529</b>	<b>0.9630</b>	<b>0.9000</b>	<b>0.9028</b>	<b>0.8991</b>	<b>0.8313</b>	<b>0.8236</b>	<b>0.8961</b>
CNMF AC	0.7755	0.7951	0.7396	0.7043	0.7131	0.7109	0.6505	0.7270

for semi-supervised learning on different number of classes. The classes for the experiment are also randomly sampled from the whole datasets with 5 times, and the average clustering results are shown in rightmost column of the table. It can be found that the semi-graph regularized auto-encoder gives a significant improvement of performance compared to the constrained NMF. The reason may be that the CNMF only utilizes the labeled data while ignoring the geometric structure hidden in the unlabeled data. When it comes to the semi-graph regularized auto-encoder, all the information from labeled and unlabeled data are all considered. As we expected, semi-graph regularized auto-encoder achieves better performance compared to the unsupervised clustering results in table 3, 4, 5.

## 5 Conclusion

In this paper, we proposed a novel graph regularized auto-encoder, which can learn a locally invariant representation of the images for both unsupervised and semi-supervised learning tasks. In unsupervised learning, our approach trains the image representation by an multiple-layer auto-encoder regularized with the graph, which encodes the locally neighborhood relationships of the original data. And in semi-supervised learning, the graph regularizer used in our auto-encoders is extended to semi-graph regularizer, which adds the penalty and reward obtained from the labeled data points to the locally neighborhood weight matrix. Experimental results on image clustering show our method provides better performance comparing with the stat-of-the-art approaches. The further work may focus on investigating the affections of the parameter settings in GAE and the impacts of the deep structure is also a possible future work.

## References

- [1] M. Belkin and P. Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*, pages 585–591, 2001.
- [2] M. Belkin, P. Niyogi, and V. Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research*, 7:2399–2434, 2006.
- [3] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- [4] D. Cai, X. He, J. Han, and T. S. Huang. Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1548–1560, 2011.
- [5] I. Goodfellow, Q. Le, A. Saxe, H. Lee, and A. Y. Ng. Measuring invariances in deep networks. In *Advances in Neural Information Processing Systems 22*, pages 646–654. 2009.
- [6] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 2*, CVPR '06, pages 1735–1742, Washington, DC, USA, 2006. IEEE Computer Society.
- [7] G. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504 – 507, 2006.
- [8] G. E. Hinton and R. S. Zemel. Autoencoders, minimum description length and helmholtz free energy. In *NIPS*, pages 3–10, 1993.
- [9] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1106–1114. 2012.
- [10] H. Larochelle and Y. Bengio. Classification using discriminative restricted boltzmann machines. In *ICML*, pages 536–543, 2008.
- [11] H. Lee, C. Ekanadham, and A. Y. Ng. Sparse deep belief net model for visual area v2. In *NIPS*, 2007.
- [12] H. Liu, Z. Wu, D. Cai, and T. S. Huang. Constrained nonnegative matrix factorization for image representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1299–1311, 2012.

- [13] M. R. Min, L. Maaten, Z. Yuan, A. J. Bonner, and Z. Zhang. Deep supervised t-distributed embedding. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 791–798, 2010.
- [14] H. Mobahi, R. Collobert, and J. Weston. Deep learning from temporal coherence in video. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 737–744. ACM, 2009.
- [15] C. Poultney, S. Chopra, Y. L. Cun, et al. Efficient learning of sparse representations with an energy-based model. In *Advances in neural information processing systems*, pages 1137–1144, 2006.
- [16] M. Ranzato, Y.-L. Boureau, and Y. LeCun. Sparse feature learning for deep belief networks. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1185–1192. MIT Press, Cambridge, MA, 2008.
- [17] M. Ranzato, C. S. Poultney, S. Chopra, and Y. LeCun. Efficient learning of sparse representations with an energy-based model. In *NIPS*, pages 1137–1144, 2006.
- [18] S. Rifai, G. Mesnil, P. Vincent, X. Muller, Y. Bengio, Y. Dauphin, and X. Glorot. Higher order contractive auto-encoder. In *Proceedings of the 2011 European conference on Machine learning and knowledge discovery in databases - Volume Part II, ECML PKDD’11*, pages 645–660, Berlin, Heidelberg, 2011. Springer-Verlag.
- [19] S. Rifai, P. Vincent, X. Muller, X. Glorot, and Y. Bengio. Contractive Auto-Encoders: Explicit invariance during feature extraction. In *ICML*, 2011.
- [20] S. T. Roweis and L. K. Saul. Nonlinear Dimensionality Reduction by Locally Linear Embedding. *Science*, 290(5500):2323–2326, 2000.
- [21] F. S. Samaria and A. C. Harter. Parameterisation of a stochastic model for human face identification. In *Applications of Computer Vision, 1994., Proceedings of the Second IEEE Workshop on*, pages 138–142. IEEE, 1994.
- [22] R. Socher, B. Huval, B. P. Bath, C. D. Manning, and A. Y. Ng. Convolutional-recursive deep learning for 3d object classification. In *NIPS*, pages 665–673, 2012.
- [23] J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–23, Dec 2000.
- [24] P. Vincent. A connection between score matching and denoising autoencoders. *Neural Comput.*, 23(7):1661–1674, July 2011.
- [25] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning, ICML ’08*, pages 1096–1103, New York, NY, USA, 2008. ACM.
- [26] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *J. Mach. Learn. Res.*, 11:3371–3408, Dec. 2010.