

# Discriminatively Boosted Image Clustering with Fully Convolutional Auto-Encoders

Fengfu Li<sup>a</sup>, Hong Qiao<sup>b,c</sup>, Bo Zhang<sup>a</sup>, Xuanyang Xi<sup>b</sup>

<sup>a</sup>*Academy of Mathematics and Systems Science, Chinese Academy of Sciences  
Beijing 100190, China*

*and*

*School of Mathematical Sciences, University of Academy of Sciences  
Beijing 100049, China*

<sup>b</sup>*Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China  
and*

*University of Academy of Sciences, Beijing 100049, China*

<sup>c</sup>*CAS Centre for Excellence in Brain Science and Intelligence Technology  
Shanghai 200031, China*

---

## Abstract

Traditional image clustering methods take a two-step approach, feature learning and clustering, sequentially. However, recent research results demonstrated that combining the separated phases in a unified framework and training them jointly can achieve a better performance. In this paper, we first introduce fully convolutional auto-encoders for image feature learning and then propose a unified clustering framework to learn image representations and cluster centers jointly based on a fully convolutional auto-encoder and soft  $k$ -means scores. At initial stages of the learning procedure, the representations extracted from the auto-encoder may not be very discriminative for latter clustering. We address this issue by adopting a boosted discriminative distribution, where high score assignments are highlighted and low score ones are de-emphasized. With the gradually boosted discrimination, clustering assignment scores are discriminated and cluster purities are enlarged. Experiments on several vision benchmark datasets show that our methods can achieve a state-of-the-art performance.

*Keywords:* image clustering, fully convolutional auto-encoder, representation learning, discriminatively boosted clustering

---

## 1. Introduction

Clustering methods are very important techniques for exploratory data analysis with wide applications ranging from data mining [1, 2], dimension reduction [3], segmentation [4] and so on. Their aim is to partition data points into clusters so that data in the same cluster are similar to each other while data in different clusters are dissimilar. Approaches to achieve this aim include partitional methods such as  $k$ -means and  $k$ -medoids, hierarchical methods like agglomerative clustering and divisive clustering, methods based on density estimation such as DBSCAN [5], and recent methods based on finding density peaks such as CFSFDP [6] and LDPS [7].

Image clustering [8] is a special case of clustering analysis that seeks to find compact, object-level models from many unlabeled images. Its applications include automatic visual concept discovery [9], content-based image retrieval and image annotation. However, image clustering is a hard task mainly owing to the following two reasons: 1) images often are of high dimensionality, which will significantly affect the performance of clustering methods such as  $k$ -means [10], and 2) objects in images usually have two-dimensional or three-dimensional local structures which should not be ignored when exploring the local structure information of the images.

To address these issues, many representation learning methods have been proposed for image feature extractions as a preprocessing step. Traditionally, various hand-crafted features such as SIFT [11], HOG [12], NMF [13], and (geometric) CW-SSIM similarity [14, 15] have been used to encode the visual information. Recently, many approaches have been proposed to combine clustering methods with deep neural networks (DNN), which have shown a remarkable performance improvement over hand-crafted features [16]. Roughly speaking, these methods can be categorized into two groups: 1) sequential methods that apply clustering on the learned DNN representations, and 2) unified approaches that jointly optimize the deep representation learning and clustering objectives.

In the first group, a kind of deep (convolutional) neural networks, such as deep belief network (DBN) [17] and stacked auto-encoders [18], is first trained in an unsupervised manner to approximate the non-linear feature embedding from the raw image space to the embedded feature space (usually being low-dimensional). And then, either  $k$ -means or spectral clustering or agglomerative clustering can be applied to partition the feature space. However, since the feature learning and clustering are separated from each

other, the learned DNN features may not be reliable for clustering.

There are a few recent methods in the second group which take the separation issues into consideration. In [19], the authors proposed deep embedded clustering that simultaneously learns feature representations with stacked auto-encoders and cluster assignments with soft  $k$ -means by minimizing a joint loss function. In [20], joint unsupervised learning was proposed to learn deep convolutional representations and agglomerative clustering jointly using a recurrent framework. In [21], the authors proposed an infinite ensemble clustering framework that integrates deep representation learning and ensemble clustering. The key insight behind these approaches is that good representations are beneficial for clustering and conversely clustering results can provide supervisory signals for representation learning. Thus, two factors, designing a proper representation learning model and designing a suitable unified learning objective will greatly affect the performance of these methods.

In this paper, we follow recent advances to propose a unified clustering method named Discriminatively Boosted Clustering (DBC) for image analysis based on fully convolutional auto-encoders (FCAE). See Fig. 1 for a glance of the overall framework. We first introduce a fully convolutional encoder-decoder network for fast and coarse image feature extraction. We then discard the decoder part and add a soft  $k$ -means model on top of the encoder to make a unified clustering model. The model is jointly trained with gradually boosted discrimination where high score assignments are highlighted and low score ones are de-emphasized. The our main contributions are summarized as follows:

- We propose a fully convolutional auto-encoder (FCAE) for image feature learning. The FCAE is composed of convolution-type layers (convolution and de-convolution layers) and pool-type layers (pooling and un-pooling layers). By adding batch normalization (BN) layers to each of the convolution-type layers, we can train the FCAE in an end-to-end way. This avoids the tedious and time-consuming layer-wise pre-training stage adopted in the traditional stacked (convolutional) auto-encoders. To the best of our knowledge, this is the first attempt to learn a deep auto-encoder in an end-to-end manner.
- We propose a discriminatively boosted clustering (DBC) framework based on the learned FCAE and an additional soft  $k$ -means model.

We train the DBC model in a self-paced learning procedure, where deep representations of raw images and cluster assignments are jointly learned. This overcomes the separation issue of the traditional clustering methods that use features directly learned from auto-encoders.

- We show that the FCAE can learn better features for clustering than raw images on several image datasets include MNIST, USPS, COIL-20 and COIL-100. Besides, with discriminatively boosted learning, the FCAE based DBC can outperform several state-of-the-art analogous methods in terms of  $k$ -means and deep auto-encoder based clustering.

The remaining part of this paper is organized as follows. Some related work including stacked (convolutional) auto-encoders, deconvolutional neural networks, and joint feature learning and clustering are briefly reviewed in Section 2. Detailed descriptions of the proposed FCAE and DBC are presented in Section 3. Experimental results on several real datasets are given in Section 4 to validate the proposed methods. Conclusions and future works are discussed in Section 5.

## 2. Related work

Stacked auto-encoders [22, 23, 24, 17, 25, 26] have been studied in the past years for unsupervised deep feature extraction and nonlinear dimension reduction. Their extensions for dealing with images are convolutional stacked auto-encoders [27, 28]. Most of these methods contain a two-stage training procedure [26]: one is layer-wise pre-training and the other is overall fine-tuning. One of the significant drawbacks of this learning procedure is that the layer-wise pre-training is time-consuming and tedious, especially when the base layer is a Restricted Boltzmann Machine (RBM) rather than a traditional auto-encoder or when the overall network is very deep.

Recently, there is an attempt to discard the layer-wise pre-training procedure and train a deep auto-encoder type network in an end-to-end way. In [29], a deep deconvolution network is learned for image segmentation. The input of the architecture is an image and the output is a segmentation mask. The network achieves the state-of-the-art performance compared with analogous methods thanks to three factors: 1) introducing a deconvolution layer and a unpooling layer [30, 31, 32] to recover the original image size of the segmentation mask, 2) applying the batch normalization [33] to each convolution layer and each deconvolution layer to reduce the internal covariate

shifts, which not only makes an end-to-end training procedure possible but also speeds up the process, and 3) adopting a pre-trained encoder on large-scale datasets such as VGG-16 model [34]. The success of the architecture motivates us that it is possible to design an end-to-end training procedure for fully convolutional auto-encoders.

Clustering has also been studied in the past years based on independent features extracted from auto-encoders (see, e.g. [10, 18, 35, 36]). Recently, there are attempts to combine the auto-encoders and clustering in a unified framework [19, 40]. In [19], the authors proposed Deep Embedded Clustering (DEC) that learns deep representations and cluster assignments jointly. DEC uses a deep stacked auto-encoder to initialize the feature extraction model and a Kullback-Leibler divergence loss to fine-tune the unified model. In [40], the authors proposed Deep Clustering Network (DCN), a joint dimensional reduction and  $k$ -means clustering framework. The dimensional reduction model is based on deep neural networks. Although these methods have achieved some success, they are not suitable for dealing with high-dimensional images due to the use of stacked auto-encoders rather than convolutional ones. This motivates us to design a unified clustering framework based on convolutional auto-encoders.

### 3. Proposed methods

In this section, we propose a unified image clustering framework with fully convolutional auto-encoders and a soft  $k$ -means clustering model (see Fig. 1). The framework contains two parts: part I is a fully convolutional auto-encoder (FCAE) for fast and coarse image feature extraction, and part II is a discriminatively boosted clustering (DBC) method which is composed of a fully convolutional encoder and a soft  $k$ -means categorizer. The DBC takes an image as input and exports soft assignments as output. It can be jointly trained with a discriminatively boosted distribution assumption, which makes the learned deep representations more suitable for the top categorizer. Our idea is very similar to self-paces learning [9], where easiest instances are first focused and more complex objects are expanded progressively. In the following subsections, we will explain the detailed implementation of the idea.

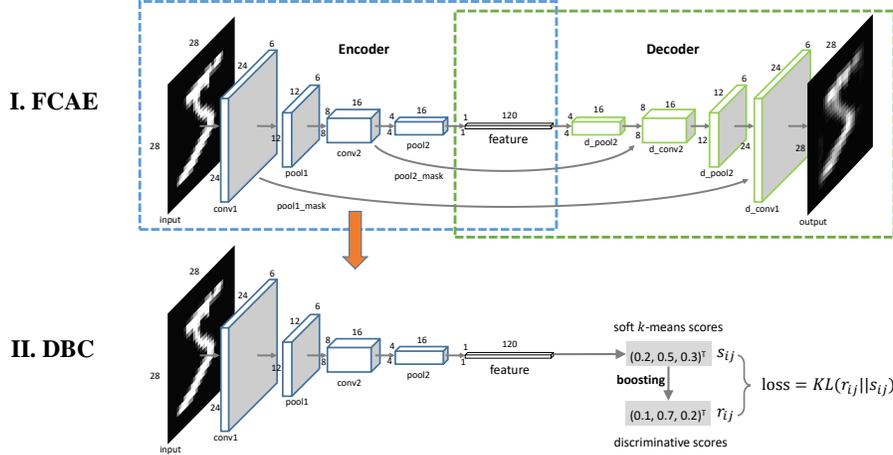


Figure 1: A unified image clustering framework. Part I is a fully convolutional auto-encoder (FCAE), and Part II is a discriminatively boosted clustering (DBC) framework based on the FCAE.

### 3.1. Fully convolutional auto-encoder for image feature extraction

Traditional deep convolutional auto-encoders adopt a greedy layer-wise training procedure for feature transformations. This could be tedious and time-consuming when dealing with very deep neural networks. To address this issue, we propose a fully convolutional auto-encoder architecture which can be trained in an end-to-end manner. Part I of Fig. 1 shows an example of FCAE on the MNIST dataset. It has the following features:

**Fully Convolutional** As pointed out in [27], the max-pooling layers are very crucial for learning biologically plausible features in the convolutional architectures. Thus, we adopt convolution layers along with max-pooling layers to make a fully convolutional encoder (FCE). Since the down-sampling operations in the FCE reduce the size of the output feature maps, we use an unpooling layer introduced in [29] to recover the feature maps. As a result, the unpooling layers along with deconvolution layers (see [29]) are adopted to make a fully convolutional decoder (FCD).

**Symmetric** The overall architecture is symmetric around the feature layer. In practice, it is suggested to design layers of an odd number. Otherwise, it will be ambiguous to define the feature layer. Besides, fully

connected layers (dense layers) should be avoided in the architecture since they destroy the local structure of the feature layer.

**Normalized** The depth of the whole network grows in log-magnitude as the input image size increases. This could make the network very deep if the original image has a very large width or height. To overcome this problem, we adopt the batch normalization (BN) [33] strategy for reducing the internal covariate shift and speeding up the training. The BN operation is performed after each convolutional layer and each deconvolutional layer except for the last output layer. As pointed out in [29], BN is critical to optimize the fully convolutional neural networks.

FCAE utilizes the two-dimensional local structure of the input images and reduces the redundancy in parameters compared with stacked auto-encoders (SAEs). Besides, FCAE differs from conventional SAEs as its weights are shared among all locations within each feature map and thus preserves the spatial locality.

### 3.2. Discriminatively boosted clustering

Once FCAE has been trained, we can extract features with the encoder part to serve as the input of a categorizer. This strategy is used in many clustering methods based on auto-encoders, such as GraphEncoder [18], deep embedding networks [35], and auto-encoder based clustering [36]. These approaches treat the auto-encoder as a preprocessing step which is separately designed from the latter clustering step. However, the representations learned in this way could be amphibolous for clustering, and the clusters may be unclear (see the initial stage in Fig. 2)).

To address this issue, we propose a self-paced approach to make feature learning and clustering in a unified framework (see Part II in Fig. 1). We throw away the decoder of the FCAE and add a soft  $k$ -means model on top of the feature layer. To train the unified model, we trust easier samples first and then gradually utilize new samples with the increasing complexity. Here, an *easier* sample (see the regions labelled 2, 3 and 4 in Fig. 2) is much certain to belong to a specific cluster, and a *harder* sample (see the region 1 in Fig. 2) is very likely to be categorized to multiple clusters. Fig. 2 describes the difference between these samples at a different learning stage of DBC.

There are three challenging questions in the learning problem of DBC which will be answered in the following subsections:

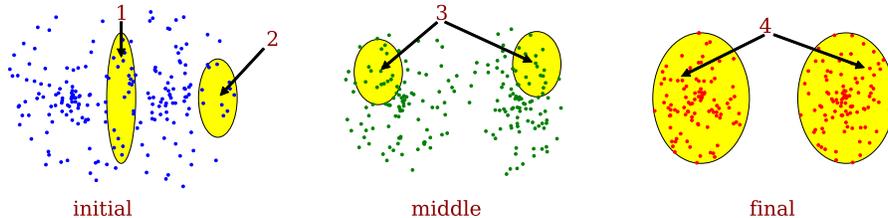


Figure 2: Learning procedure of DBC.

1. How to choose a proper criterion to determine the easiness or hardness of a sample?
2. How to transform harder samples into easier ones?
3. How to learn from easier samples?

### 3.2.1. Easiness measurement with the soft $k$ -means scores

We follow DEC [19] to adopt the  $t$ -distribution-based soft assignment to measure the easiness of a sample. The  $t$ -distribution is investigated in [37] to deal with the crowding problem of low-dimensional data distributions. Under the  $t$ -distribution kernel, the soft score (or similarity) between the feature  $z_i$  ( $i \in 1, 2, \dots, m$ ) and the cluster center  $\mu_j$  ( $j \in 1, 2, \dots, k$ ) is

$$s_{ij} \propto \left(1 + \frac{\|z_i - \mu_j\|^2}{v}\right)^{-\frac{v+1}{2}} \quad (1)$$

$$\text{s.t.} \quad \sum_{j=1}^k s_{ij} = 1$$

Here,  $v$  is the degree of freedom of the  $t$ -distribution and set to be 1 in practice. The most important reason for choosing the  $t$ -distribution kernel is that it has a longer tail than the famous heat kernel (or the Gaussian-distribution kernel). Thus, we do not need to pay much attention to the parameter estimation (see [37]), which is a hard task in unsupervised learning.

### 3.2.2. Boosting easiness with discriminative target distribution

We transform the harder examples to the easier ones by boosting the higher score assignments and, meanwhile, bring down those with lower scores. This can be achieved by constructing an underlying target distribution  $r_{ij}$

from  $s_{ij}$  as follows:

$$\begin{aligned} r_{ij} &\propto s_{ij}^\alpha, \alpha > 1 \\ \text{s.t. } &\sum_{j=1}^k r_{ij} = 1 \end{aligned} \quad (2)$$

Suppose we can ideally learn from the soft scores (denoted as  $S$ ) to the assumptive distribution (denoted as  $R$ ) each time. Then we can generate a learning chain as follows:

$$S^{(0)} \rightarrow R^{(0)} = S^{(1)} \rightarrow R^{(1)} = S^{(2)} \rightarrow \dots$$

The following two properties can be observed from the chain:

**Property 1** If  $s_{ij}^{(0)} = s_{ij'}^{(0)}$  for any  $j$  and  $j'$ , then  $s_{ij}^{(t)} \equiv 1/k$  for all  $j$  and all time step  $t$ .

**Proof** Under the condition, and by (2) we can deduce that  $r_{ij}^{(0)} \equiv r_{ij'}^{(0)}$ . By the chain this is equivalent to the fact that  $s_{ij}^{(1)} \equiv s_{ij'}^{(1)}$ . Thus, the conclusion  $s_{ij}^{(t)} \equiv s_{ij'}^{(t)}$  follows recursively for all  $t$ .  $\square$

**Property 2** If there exists an  $l$  such that  $s_{il}^{(0)} > \max_{j \neq l} s_{ij}^{(0)}$ , then

$$\lim_{t \rightarrow \infty} s_{ij}^{(t)} = \begin{cases} 1 & \text{if } j = l \\ 0 & \text{if } j \neq l \end{cases}$$

**Proof** By (2) we have

$$\frac{s_{ij}^{(t)}}{s_{il}^{(t)}} = \left[ \frac{s_{ij}^{(t-1)}}{s_{il}^{(t-1)}} \right]^\alpha = \left[ \frac{s_{ij}^{(t-2)}}{s_{il}^{(t-2)}} \right]^{\alpha^2} = \dots = \left[ \frac{s_{ij}^{(0)}}{s_{il}^{(0)}} \right]^{\alpha^t}.$$

By the assumption  $s_{il}^{(0)} > \max_{j \neq l} s_{ij}^{(0)}$ , it is seen that  $s_{ij}^{(0)}/s_{il}^{(0)} < 1$  for any  $j \neq l$ .

On the other hand, since  $\alpha > 1$ , we have  $\lim_{t \rightarrow \infty} \alpha^t = \infty$ . Thus,

$$\lim_{t \rightarrow \infty} \frac{s_{ij}^{(t)}}{s_{il}^{(t)}} = \lim_{\alpha^t \rightarrow \infty} \left[ \frac{s_{ij}^{(0)}}{s_{il}^{(0)}} \right]^{\alpha^t} = 0, \quad \forall j \neq l.$$

Since  $s_{il}^{(t)}$  is finite, we have  $\lim_{t \rightarrow \infty} s_{ij}^{(t)} = 0, \forall j \neq l$ . Finally, with the constrains

$\sum_{j=1}^k s_{ij}^{(t)} = 1$ , we obtain

$$\lim_{t \rightarrow \infty} s_{il}^{(t)} = 1 - \sum_{j \neq l} \lim_{t \rightarrow \infty} s_{ij}^{(t)} = 1. \quad \square$$

**Property 1** tells us that the *hardest* sample (which has the equal probability to be assigned to different clusters) would always be the hardest one. However, in practical applications, there can hardly exist such examples. **Property 2** shows that the initial non-discriminative samples could be boosted gradually to be definitely discriminative. As a result, we get the desired features for  $k$ -means clustering.

Note that the boosting factor  $\alpha$  controls the speed of the learning process. A larger  $\alpha$  can make the learning process more quickly than smaller ones. However, it may boost some falsely categorized samples too quickly at initial stages and thus makes their features irrecoverable at later stages.

Besides, it can be helpful to balance the data distribution at different learning stages. In [19], the authors proposed to normalize the boosted assignments to prevent large clusters from distorting the hidden feature space. This problem can be overcome by dividing a normalization factor  $n_j = \sum_{i=1}^m r_{ij}$  for each of the  $r_{ij}$ .

### 3.2.3. Learning with the Kullback-Leibler divergence loss

In the last subsection, it was assumed that we could learn from  $s_{ij}$  to the boosted target distribution  $r_{ij}$ . This aim can be achieved with a joint Kullback-Leibler (KL) divergence loss, that is,

$$(\theta^*, \mu^*) = \arg \min_{\theta, \mu} L = \text{KL}(R||S) = \sum_{i=1}^m \sum_{j=1}^k r_{ij} \log \frac{r_{ij}}{s_{ij}}. \quad (3)$$

Fig. 3 gives an example of the joint loss when  $k = 2$ , where  $L_{ij} = r_{ij} \log(r_{ij}/s_{ij})$  is the loss generated by the sample  $x_i$  with respect to the  $j$ th cluster ( $j = 1$  or 2). Regions marked in Fig. 3 roughly correspond to the regions marked in Fig. 2.

Intuitively, the loss has the following main features:

- For an ambiguous (or hard) sample (i.e.,  $s_{ij} \approx s_{il}, \forall j, l$ ), its loss  $L_i = \sum_{j=1}^k r_{ij} \log(r_{ij}/s_{ij}) \approx 0$  according to **Property 1**. Therefore, it will not be seriously treated in the learning process. (Region 1)
- For a good categorized sample (i.e., there exists an  $l$  such that  $1 \gg s_{il} > \max_{j \neq l} s_{ij}$ ), its loss will be much greater than zero, and thus it will be treated more seriously. (Regions 2 and 3)

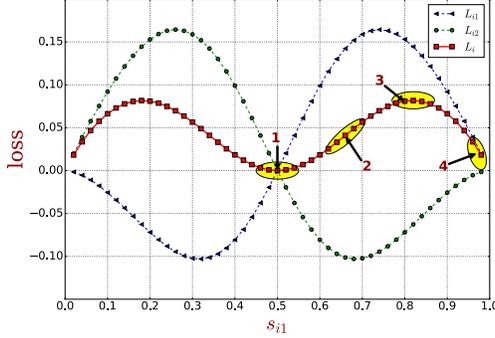


Figure 3: KL divergence loss with respect to the soft scores assigned to the first cluster. Here we assume that there are 2 clusters, so 0.5 is a random guess probability.

- For a definitely well categorized sample (i.e., there exists an  $l$  such that  $1 \approx s_{il} \gg \max_{j \neq l} s_{ij}$ ), its loss will be near zero. This means that its features do not need to be changed much more. (Region 4)

By (1)-(3), the gradients of the KL divergence loss w.r.t.  $z_i$  and  $\mu_j$  can be deduced as follows:

$$\frac{\partial L}{\partial z_i} = \frac{1+v}{v} \sum_{j=1}^k (r_{ij} - s_{ij}) \frac{z_i - \mu_j}{1 + \|z_i - \mu_j\|^2/v} \quad (4)$$

and

$$\frac{\partial L}{\partial \mu_j} = \frac{1+v}{v} \sum_{i=1}^m (r_{ij} - s_{ij}) \frac{\mu_j - z_i}{1 + \|\mu_j - z_i\|^2/v} \quad (5)$$

The derivation of (4) and (5) can be found in the appendix.

#### 3.2.4. Training algorithm

In this section, we summarize the overall training procedure of the proposed method in Algorithm 1 and Algorithm 2. They implement the framework showed in Fig. 1. Here  $T$  is the maximum learning epochs,  $B$  is the maximum updating iterations in each epoch and  $m_b$  is the mini-batch size. The encoder part of FCAE is  $f: x \rightarrow z$ , which is parameterized by  $\theta_e$  and the decoder part of FCAE is  $g: z \rightarrow \hat{x}$ , which is parameterized by  $\theta_d$ .

---

**Algorithm 1** Discriminatively Boosted Clustering (DBC)

---

**Require:**  $X, T, B, m_b, \alpha, k$

**Ensure:**  $\theta$  and  $\mu$

//**Stage I: Train a FCAE and clustering with its features**

1: Train a deep fully convolutional auto-encoder

$$x_i \xrightarrow{\theta_e} z_i(\text{features}) \xrightarrow{\theta_d} \hat{x}_i \quad (\text{M1})$$

with the Euclidian loss

$$(\theta_e^*, \theta_d^*) = \arg \min_{\theta_e, \theta_d} \sum_{i=1}^m \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^m \|x_i - g_{\theta_d}(f_{\theta_e}(x_i))\|_2^2$$

by using the traditional error back-propagation algorithm.

2: Extract features:  $Z \leftarrow f_{\theta_e^*}(X)$

3: Clustering with the features:  $\mu_z \leftarrow k$ -means centers

---

## 4. Experiments

In this section, we present experimental results on several real datasets to evaluate the proposed methods by comparing with several state-of-the-art methods. To this end, we first introduce several evaluation benchmarks and then present visualization results of the inner features, the learned FCAE weights, the frequency hist of soft assignments during the learning process and the features embedded in a low-dimensional space. We will also give some ablation studies with respect to the boosting factor  $\alpha$ , the normalization factor  $n_j$  and the FCAE initializations.

### 4.1. Evaluation benchmarks

**Datasets** We evaluate the proposed FCAE and DBC methods on two hand-written digit image datasets (MNIST <sup>1</sup> and USPS <sup>2</sup>) and two multi-view object image datasets (COIL-20 <sup>3</sup> and COIL-100 <sup>4</sup>). The size of the datasets,

---

<sup>1</sup><http://yann.lecun.com/exdb/mnist/>

<sup>2</sup><http://www.cs.nyu.edu/~roweis/data.html>

<sup>3</sup><http://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php>

<sup>4</sup><http://www.cs.columbia.edu/CAVE/software/softlib/coil-100.php>

---

**Algorithm 2** DBC (Continued)

---

//**Stage II: Jointly learn the FCE and cluster centers**

- 4: Construct a unified clustering model with encoder parameters
- $\theta$
- and cluster centers
- $\mu$

$$x_i \xrightarrow{\theta} z_i \xrightarrow{\mu_j} s_{ij} \quad (\text{M2})$$

- 5: Initialization:
- $\theta \leftarrow \theta_e^*$
- ,
- $\mu \leftarrow \mu_z$

6: **for**  $t = 1$  to  $T$  **do**

- 7:     Forward propagate (M2) and update the
- soft*
- assignments

$$s_{ij} \leftarrow \frac{(1 + \|z_i - \mu_j\|^2)^{-1}}{\sum_{j=1}^k (1 + \|z_i - \mu_j\|^2)^{-1}}, \text{ where } z_i = f_\theta(x_i).$$

- 8:     Update the target distribution

$$r_{ij} \leftarrow \frac{s_{ij}^\alpha / n_j}{\sum_{j=1}^k s_{ij}^\alpha / n_j}, \text{ where } n_j = \sum_{i=1}^m s_{ij}^\alpha.$$

9:     **for**  $b = 1$  to  $B$  **do**

- 10:         Forward propagate (M2) with a mini-batch of
- $m_b$
- samples.

- 11:         Backward propagate (M2) from (4) and (5) to get
- $\partial L / \partial \theta$
- and
- $\partial L / \partial \mu$
- .

- 12:         Update
- $\theta$
- and
- $\mu$
- with the gradients.

13:     **end for**

- 14:     Stop if
- hard*
- assignments remain unchanged.

15: **end for**

---

the number of categories, the image sizes and the number of channels are summarized in Table 1.

Table 1: Datasets used in our experiments.

Dataset	#Samples	#Categories	Image Size	#Channels
MNIST	70000	10	28×28	1
USPS	11000	10	16×16	1
COIL-20	1440	20	128×128	1
COIL-100	7200	100	128×128	3

**Evaluation metrics** Two standard metrics are used to evaluate the experiment results explained as follows.

- Accuracy (ACC) [19]. Given the ground truth labels  $\{c_i | 1 \leq i \leq m\}$  and the predicted assignments  $\{\hat{c}_i | 1 \leq i \leq m\}$ , ACC measures the average accuracy:

$$\text{ACC}(\hat{c}, c) = \max_g \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{c_i = g(\hat{c}_i)\}$$

where  $g$  ranges over all possible one-to-one mappings between the labels of the predicted clusters and the ground truth labels. The optimal mapping can be efficiently computed using the Hungarian algorithm [38].

- Normalized mutual information (NMI) [39]. From the information theory point of view NMI can be interpreted as

$$\text{NMI}(\hat{c}, c) = \frac{\text{MI}(\hat{c}, c)}{\max(\text{H}(\hat{c}), \text{H}(c))}$$

where  $\text{H}(c)$  is the entropy of  $c$  and  $\text{NMI}(\hat{c}, c)$  is the mutual information of  $\hat{c}$  and  $c$ .

**Network architectures** Table 2 shows the network architecture of the encoder parts with respect to different datasets. The decoder parts are totally reversed by the encoder parts. We use max-pooling in all the experiments. The size of all the feature layers is  $1 \times 1$ . No padding is used in the convolutional layers except for the USPS dataset whose padding size is 1.

**The comparing methods** To validate the effectiveness of FCAE and DBC, we compare them with the following state-of-the-art methods in terms of the  $k$ -means and deep auto-encoders based clustering.

- **KMS** is the baseline method that applies the  $k$ -means algorithm on raw images.
- **DAE-KMS** [19] uses deep auto-encoders for feature extraction and then applies  $k$ -means for later clustering.
- **AEC** [36] is a variant of DAE-KMS that simultaneously optimizes the data reconstruction error and representation compactness.

Table 2: Detailed configuration of the network architecture of the convolutional encoder. The first rows are the filter sizes of the corresponding layer (filter size or pooling size, #filters). The second rows are the output sizes (feature map size, #channels).

datasets	conv1	pool1	conv2	pool2	conv3	pool3	conv4	pool4	features
MNIST	5, 6	2, -	5, 16	2, -	-	-	-	-	4, 120
	24, 6	12, 6	8, 16	4, 16	-	-	-	-	1, 120
USPS	3, 20	2, -	3, 20	2, -	-	-	-	-	4, 160
	16,20	8, 20	8, 20	4, 20	-	-	-	-	1, 160
COIL	9, 20	2, -	5, 20	2, -	5, 20	2, -	5,40	2, -	4, 320
	120, 20	60, 20	56, 20	28, 20	24, 20	12, 20	8, 40	4, 40	1, 320

- **IEC** [21] incorporates the deep representation learning and ensemble clustering.
- **DEC** [19] simultaneously learns the feature representations and cluster centers using deep auto-encoders and soft  $k$ -means, respectively.
- **DEN** [35] learns the clustering-oriented representations by utilizing deep auto-encoders and manifold constraints.
- **DCN** [40] jointly applies dimensionality reduction and  $k$ -means clustering.
- **FCAE-KMS** (our algorithm) adopts FCAE for feature extraction and applies  $k$ -means for the latter clustering.
- **DBC** (our algorithm) uses Algorithm ?? for training a unified clustering method.

Table 3: Clustering performance on MNIST.

Metric	KMS	DAE-KMS	AEC	IEC	DCN	DEC	FCAE-KMS	DBC
ACC	0.535	0.818	0.760	0.609	0.58 / 0.93 <sup>5</sup>	0.843	0.794	<b>0.964</b>
NMI	0.531	-	0.669	0.542	0.63 / 0.85	-	0.698	<b>0.917</b>

<sup>5</sup>DCN with processed MNIST.

**Results and analysis** Table 3 summarizes the benchmark results on the MNIST dataset. The  $k$ -means method performs badly on raw images. However, based on the end-to-end trained FACE features,  $k$ -means can achieve comparative results compared with DAE-KMS which uses greedily layer-wise trained deep auto-encoder features. Moreover, with an additional joint training, DBC outperforms FCAE-KMS and beats all the other comparing methods in terms of ACC and NMI.

Tables 4-6 show the benchmarks on USPS, COIL-20 and COIL-100, respectively. Similarly to the observations on the MNIST hand-written digits dataset, DBC outperforms FCAE-KMS by a large margin on the USPS hand-written digits dataset. On the COIL sets, DBC obtained a little better results than FCAE-KMS did.

On the hand-written digits datasets, the number of samples is much larger than the number of categories. This results in the distribution of the FCAE features to be closely related, and lots of ambiguous samples may occur. As a result, discriminatively boosting makes sense on these datasets. Thus, there is no doubt that DBC performs much better than FCAE-KMS. On the COIL sets, DBC takes little advantage of the discriminatively boosting procedure since the FCAE features are very helpful for clustering. Thus, there are very few ambiguous samples whose easiness needs to be boosted.

Table 4: Clustering performance on USPS.

Metric	KMS	AEC	IEC	FCAE-KMS	DBC
ACC	0.535	0.715	<b>0.767</b>	0.667	0.743
NMI	0.531	0.651	0.641	0.645	<b>0.724</b>

Table 5: Clustering performance on COIL-20.

Metric	KMS	DEN	FCAE-KMS	DBC
ACC	0.592	0.725	0.787	<b>0.793</b>
NMI	0.767	0.870	0.882	<b>0.895</b>

Table 6: Clustering performance on COIL-100.

Metric	KMS	IEC	FCAE-KMS	DBC
ACC	0.506	0.546	0.766	<b>0.775</b>
NMI	0.772	0.787	0.897	<b>0.905</b>

#### 4.2. Visualization

One of the advantages of fully convolutional neural networks is that we can naturally visualize the inner activations (or features) and the trained weights (or filters) in a two-dimensional space [27]. Besides, we can monitor the learning process of DBC by drawing frequency hists of assignment scores. In addition,  $t$ -SNE can be applied to the embedded features to visualize the manifold structures in a low-dimensional space. Finally, we show some typical falsely categorized samples generated by our algorithm.

##### 4.2.1. Visualization of the inner activations and learned filters

In Fig. 4, we visualize the inner activations of FCAE on the MNIST dataset with three digits: 1, 5, and 9. As shown in the figure, the activations in the feature layer are very sparse. Besides, the deconvolutional layer gradually recovers details of the pooled feature maps and finally gives a rough description of the original image. This indicates that FCAE can learn clustering-friendly features and keep the key information for image reconstruction.

Fig. 5 visualizes the learned filters of FCAE on the MNIST dataset. It is observed in [27] that the stacked convolutional auto-encoders trained on noisy inputs (30% binomial noise) and a max-pooling layer can learn localized biologically plausible filters. However, even without adding noise, the learned deconvolutional filters in our architectures are non-trivial Gabor-like filters which are visually the nicest shapes. This is due to the use of max-pooling and unpooling operations. As discussed in [27], the max-pooling layers are elegant way of enforcing sparse codes which are required to deal with the over-complete representations of convolutional architectures.

##### 4.2.2. Monitoring the learning process

We use frequency hist of the soft assignment scores to monitor the learning process of DBC. Fig. 6 shows the hists of scores on the MNIST test dataset

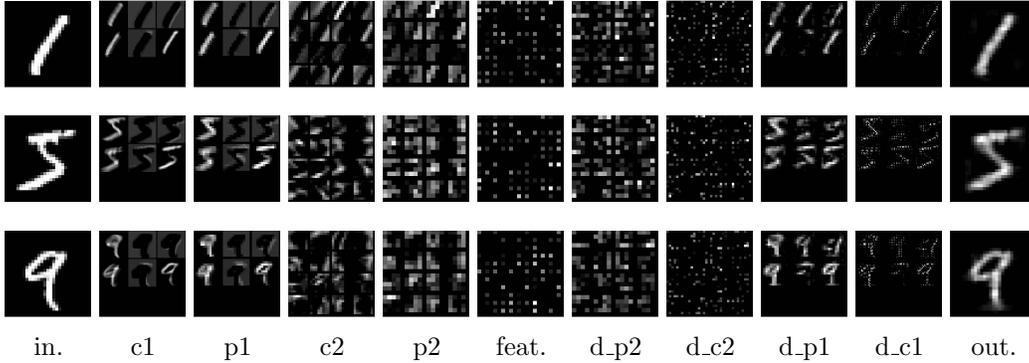


Figure 4: Visualization of the inner activations with respect to digits 1, 5 and 9.

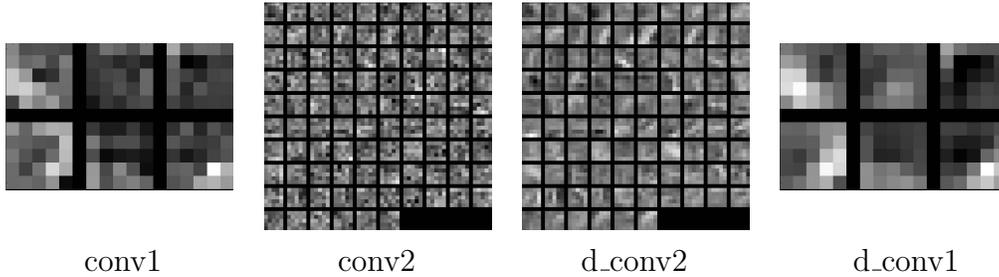


Figure 5: Visualization of the learned FCAE filters.

(a subset of the MNIST dataset with 10000 samples). The scores are assigned to the first cluster at different learning epochs. At early epochs ( $t \leq 4$ ), most of the scores are near 0.1. This is a random guess probability because there are 10 clusters. As the learning procedure goes on, some higher score samples are discriminatively boosted and their scores become larger than others. As a result, the cluster tends to “believe” in these higher score samples and consequently make scores of the others to be smaller (approximating zero). Finally, the scores assigned to the cluster become two-side polarized. Samples with very high scores ( $s_{ij} \approx 0.8$ ) are thought to definitely belong to the first cluster and the others with very small scores ( $s_{ij} \approx 0.02$ ) should belong to other clusters.

#### 4.2.3. Embedding learned features in a low dimensional space

We visualize the distribution of the learned features in a two-dimensional space with  $t$ -SNE [37]. Fig. 7 shows the embedded features of the MNIST test dataset at different epochs. At the initial epoch, the features learned

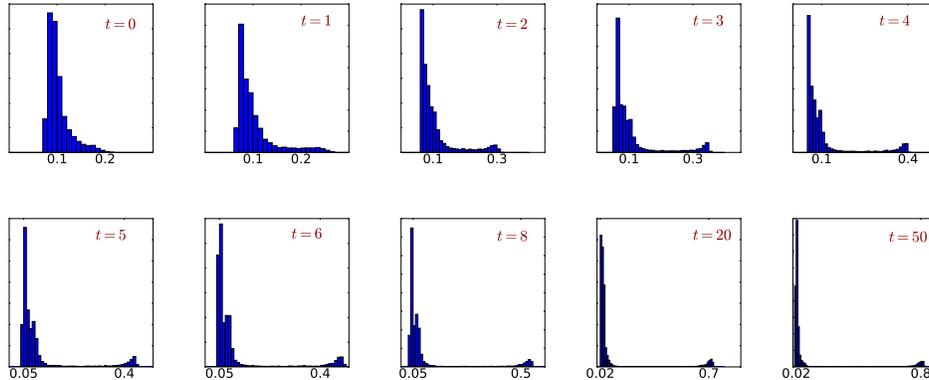


Figure 6: Visualization of the soft score frequency histograms with respect to the first cluster at different learning stages.

with FCAE are not very discriminative for clustering. As shown in Fig. 7(a), the features of digits 3, 5, and 8 are closely related. The same thing happened with digits 4, 7, and 9. At the second epoch, the distribution of the learned features becomes much compact locally. Besides, the features of digit 7 become far away from those of digits 4 and 9. Similarly, the features of digit 8 get far away from those of digits 3 and 5. As the learning procedure goes on, the hardest digits (4 v.s. 9, 3 v.s. 5) for categorization are mostly well categorized after enough discriminative boosting epochs. The observation is consistent with the results showed in Subsection 4.2.2.

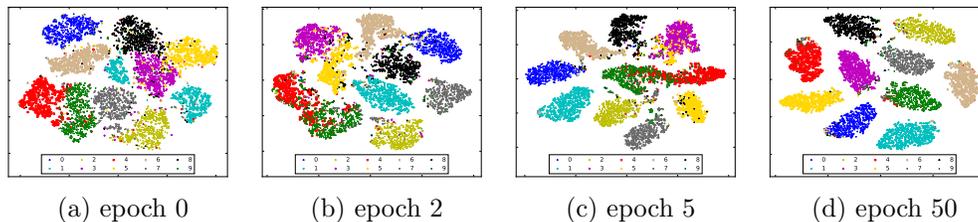


Figure 7: Visualization of the embedded features in a two-dimensional space with  $t$ -SNE.

#### 4.2.4. Visualization of falsely categorized examples

In Fig. 8, we show the top 100 falsely categorized examples whose maximum soft assignment scores are over 0.6. It can be observed that it is very hard to distinguish between some ground truth digits 4, 7 and 9 even with

human experience. Lots of digits 7 are written with transverse lines in their middle space and would be thought to be ambiguous for the clustering algorithm. Besides, some ground truth images are themselves confusing, such as those showed with the gray background.



Figure 8: Visualization of falsely categorized high score samples (top-100). The number in the top-left corner is the clustering label which is generated by the Hungarian algorithm.

### 4.3. Discussions

In this section, we make some ablation studies on the learning process with respect to different boosting factors ( $\alpha$ ), different normalization methods ( $n_j$ ) and different initialization models generated by FCAE.

#### 4.3.1. Impact of the boosting factor $\alpha$

Fig. 9(a) shows the ACC and NMI curves, where  $\alpha$  equals to 1.5, 2, 4. With a small  $\alpha$  ( $\alpha = 1.5$ ), the learning process is very slow and takes very long time to terminate. On the contrary, when the factor is set to be very large ( $\alpha = 4$ ), the learning process is very fast at the initial stages. However,

this could result in falsely boosting some scores of the ambiguous samples. As a consequence, the model learned too much from some false information so the performance is not so satisfactory. With a moderate boosting factor ( $\alpha = 2$ ), the ACC and NMI curves grow reasonably and progressively.

#### 4.3.2. Impact of the balance normalization

In DEC [19], the authors pointed out that the balance normalization plays an important role in preventing large clusters from distorting the hidden feature space. To address this issue, we compare three normalization strategies: 1) the constant normalization for comparison, that is,  $n_j = 1$ , 2) the normalization by dividing the sum of the original soft assignment score per cluster, that is,  $n_j = \sum_i s_{ij}$ , which is adopted in DEC, and 3) the normalization by dividing the sum of the boosted soft assignment score per cluster, that is,  $n_j = \sum_i s_{ij}^\alpha$ . Fig. 9(b) presents the value curves of ACC and NMI against the epoch with these settings. Initially, the normalization does not affect ACC and NMI very much. However, the constant normalization can easily get stuck at early stages. The normalization by dividing  $n_j = \sum_i s_{ij}$  has certain power of preventing the distortion. Our normalization strategy gives the best performance compared with the previous methods. This is because our normalization directly reflects the changes of the boosted scores.

#### 4.3.3. Impact of the FCAE initialization

To investigate the impact of the FCAE initialization on DBC, we compare the performance of DBC with three different initialization models: 1) the random initialization, 2) the initialization with a half-trained FCAE model, and 3) the initialization with a sufficiently trained FCAE model. The comparison results are shown in Fig. 9(c). As illustrated in the figure, DBC performs greatly based on all the models even when the initialization model is randomly distributed. However, if the FCAE model is not sufficiently trained, the resultant DBC model will be suboptimal.

## 5. Conclusions and future works

In this paper, we proposed FCAE and DBC to deal with image representation learning and image clustering, respectively. Benchmarks on several visual datasets show that our methods can achieve superior performance than the analogous methods. Besides, the visualization shows that the proposed learning algorithm can implement the idea proposed in Section 3.2. Some

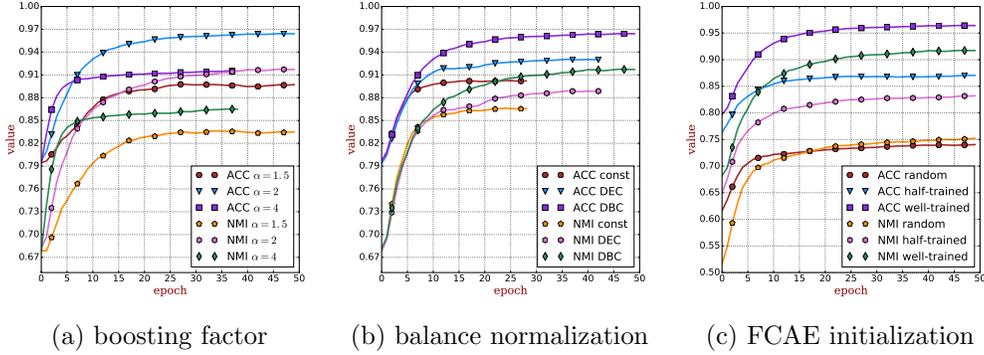


Figure 9: Ablation studies with respect to the boosting factor  $\alpha$ , the balance normalization factor  $n_j$  and the FCAE initialization models.

issues to be considered in the future include: 1) adding suitable constraints on FCAE to deal with natural images, and 2) scaling the algorithm to deal with large-scale datasets such as the ImageNet dataset.

## Acknowledgement

This work was supported in part by NNSF of China under grants 61379093, 61602483 and 61603389. We thank Shuguang Ding, Xuanyang Xi, Lu Qi and Yanfeng Lu for valuable discussions.

## Appendix

### A. Derivation of (4).

We use the chain rule for the deduction. First, we set

$$u_{ij} = 1 + \frac{\|z_i - \mu_j\|^2}{v}. \quad (.1)$$

Then it follows that

$$\frac{\partial u_{ij}}{\partial z_i} = \frac{2}{v}(z_i - \mu_j). \quad (.2)$$

Now set

$$q_{ij} = u_{ij}^{-\frac{1+v}{2}}, \quad (.3)$$

so

$$\begin{aligned}
\frac{\partial q_{ij}}{\partial z_i} &= \frac{\partial q_{ij}}{\partial u_{ij}} \cdot \frac{\partial u_{ij}}{\partial z_i} \\
&= \left(-\frac{1+v}{2}\right) u_{ij}^{-\frac{3+v}{2}} \cdot \frac{2}{v} (z_i - \mu_j) \\
&= \left(-\frac{1+v}{v}\right) u_{ij}^{-1} q_{ij} \cdot (z_i - \mu_j). \tag{.4}
\end{aligned}$$

Further, let

$$s_{ij} = \frac{q_{ij}}{\sum_{j'} q_{ij'}}. \tag{.5}$$

Then we have

$$\begin{aligned}
\frac{\partial s_{ij}}{\partial z_i} &= \frac{\frac{\partial q_{ij}}{\partial z_i} \cdot \sum_{j'} q_{ij'} - q_{ij} \cdot \sum_{j'} \frac{\partial q_{ij'}}{\partial z_i}}{(\sum_{j'} q_{ij'})^2} \\
&= \frac{\left(-\frac{1+v}{v}\right) u_{ij}^{-1} q_{ij} \cdot (z_i - \mu_j) \cdot \sum_{j'} q_{ij'} - q_{ij} \cdot \sum_{j'} \left(-\frac{1+v}{v}\right) u_{ij'}^{-1} q_{ij'} \cdot (z_i - \mu_{j'})}{(\sum_{j'} q_{ij'})^2} \\
&= \left(-\frac{1+v}{v}\right) \frac{q_{ij}}{\sum_{j'} q_{ij'}} \frac{u_{ij}^{-1} \cdot (z_i - \mu_j) \cdot \sum_{j'} q_{ij'} - \sum_{j'} u_{ij'}^{-1} q_{ij'} \cdot (z_i - \mu_{j'})}{\sum_{j'} q_{ij'}} \\
&= \left(-\frac{1+v}{v}\right) s_{ij} (u_{ij}^{-1} \cdot (z_i - \mu_j) - \sum_{j'} u_{ij'}^{-1} s_{ij'} \cdot (z_i - \mu_{j'})). \tag{.6}
\end{aligned}$$

Combine the above expressions to get the required result

$$\frac{\partial L}{\partial z_i} = -\sum_j \frac{r_{ij}}{s_{ij}} \cdot \frac{\partial s_{ij}}{\partial z_i} \tag{.7}$$

$$\begin{aligned}
&= \frac{1+v}{v} \sum_j \frac{r_{ij} - s_{ij}}{u_{ij}} (z_i - \mu_j) \\
&= \frac{1+v}{v} \sum_j (r_{ij} - s_{ij}) \frac{z_i - \mu_j}{1 + \|z_i - \mu_j\|^2/v} \tag{.8}
\end{aligned}$$

## B. Derivation of (5).

(5) can be derived similarly by exchanging  $\mu$  and  $z$  in the above derivations of (4).

## References

## References

- [1] J. Han, J. Pei, M. Kamber. Data Mining: Concepts and Techniques. *Elsevier*, 2011.
- [2] P. Berkhin. A survey of clustering data mining techniques. In: *Grouping Multidimensional Data*, pp. 25-71, 2006.
- [3] C. Boutsidis, A. Zouzias, M. Mahaoney, and P. Drineas. Randomized dimensionality reduction for  $k$ -means clustering. *IEEE Transactions on Information Theory*, vol. 61, no. 2, pp. 1045-1062, 2015.
- [4] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888-905, 2000.
- [5] M. Ester, H.P. Kriegel, J. Sander and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. *KDD*, vol. 96, no. 34, pp. 226-231, 1996.
- [6] A. Rodriguez and A. Laio. Clustering by fast search and find of density peaks. *Science*, vol. 344, no. 6191, pp. 1492-1496, 2014.
- [7] F. Li, H. Qiao and B. Zhang. Effective deterministic initialization for  $k$ -means-like methods via local density peaks searching. *arXiv:1611.06777*, 2016.
- [8] N. Ahmed. Recent review on image clustering. *IET Image Processing*, vol. 9, no. 11, pp. 1020-1032, 2015.
- [9] Y.J. Lee and K. Grauman. Learning the easy things first: Self-paced visual category discovery. *IEEE Conference on CVPR*, pp. 1721-1728, 2011.
- [10] C. Ding, T. Li. Adaptive dimension reduction using discriminant analysis and  $k$ -means clustering. *Proc. 24th International Conference on Machine learning*, pp. 521-528, 2007.
- [11] D. Lowe. Object recognition from local scale-invariant features. *Proc. 7th International Conference on Computer Vision*, vol. 2, pp. 1150-1157, 1999.

- [12] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. *Proc. Computer Vision and Pattern Recognition*, vol. 1, pp. 886-893, 2005.
- [13] S. Hong, J. Choi, J. Feyereisl, B. Han and L. S. Davis. Joint image clustering and labeling by matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 7, pp. 1411-1424, 2016.
- [14] M. Sampat, Z. Wang, S. Gupta, A. Bovik and M. Markey. Complex wavelet structural similarity: A new image similarity index. *IEEE Transactions on Image Processing*, vol. 18, no. 1, pp. 2385-2401, 2009.
- [15] F. Li, X. Huang, H. Qiao and B. Zhang. A new manifold distance for visual object categorization. *The 12th World Congress on Intelligent Control and Automation*, pp. 2232-2236, 2016.
- [16] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. *Advanced Neural Information Processing Systems*, vol. 24, pp. 1097-1105, 2012.
- [17] G. Hinton, S. Osindero, Y. Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, vol. 18, no. 7, pp. 1527-1554, 2006.
- [18] F. Tian, B. Gao, Q. Cui, E. Chen and T. Liu. Learning deep representations for graph clustering. *AAAI*, pp. 1293-1299, 2014.
- [19] J. Xie, R. Girshick and A. Farhadi. Unsupervised deep embedding for clustering analysis. *Proc. 33rd International Conference on Machine Learning*, pp. 478-487, 2016.
- [20] J. Yang, D. Parikh, D. Batra. Joint unsupervised learning of deep representations and image clusters. *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [21] H. Liu, M. Shao, S. Li and Y. Fu. Infinite ensemble for image clustering. *Proc. 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 1745-1754, 2016.

- [22] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio and P. Manzagol. Stacked denoising auto-encoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, vol. 11, pp. 3371-3408, Dec. 2010.
- [23] P. Baldi. Autoencoders, unsupervised learning, and deep architectures. *ICML Workshop on Unsupervised and Transfer Learning*, vol. 27, pp. 37-50, 2012.
- [24] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1789-1828, 2013.
- [25] G. Hinton, R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, vol. 313, no. 5786, pp. 504-507, 2006.
- [26] Y. Bengio, P. Lamblin, D. Popovici and H. Larochelle. Greedy layer-wise training of deep networks. *Advances in Neural Information Processing Systems*, vol. 19, pp. 153, 2007.
- [27] J. Masci, U. Meier, D. ciresan and J. Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. *International Conference on Artificial Neural Networks*, pp. 52-59, 2011.
- [28] H. Lee, R. Grosse, R. Ranganath and A. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. *Proc. 26th Annual International Conference on Machine Learning*, pp. 609-616, 2009.
- [29] H. Noh, S. Hong and B. Han. Learning deconvolution network for semantic segmentation. *Proc. IEEE International Conference on Computer Vision*, pp. 1520-1528, 2015.
- [30] M. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. *European Conference on Computer Vision*, pp. 818-833, 2014.
- [31] R. Mohan. Deep deconvolutional networks for scene parsing. *arXiv:1411.4101*, 2014.
- [32] M. Zeiler, G. Taylor, R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. *2011 International Conference on Computer Vision*, pp. 2018-2025, 2011.

- [33] S. Ioffe, and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv:1502.03167*, 2015.
- [34] K. Simonyan, A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv:1409.1556*, 2014.
- [35] P. Huang, Y. Huang, W. Wang and L. Wang. Deep embedding network for clustering. *ICPR*, pp. 1532-1537, 2014.
- [36] C. Song, F. Liu, Y. Huang, et al. Auto-encoder based data clustering. *Iberoamerican Congress on Pattern Recognition*, pp. 117-124, 2013.
- [37] L. Maaten, G. Hinton. Visualizing data using *t*-SNE. *Journal of Machine Learning Research*, vol. 9, pp. 2579-2605, 2008.
- [38] H. Kuhn. The Hungarian method for the assignment problem. *50 Years of Integer Programming 1958-2008*, pp. 29-47, 2010.
- [39] D. Cai, X. He, J. Han. Locally consistent concept factorization for document clustering. *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 6, pp. 902-913, 2011.
- [40] B. Yang, X. Fu, ND Sidiropoulos, and M Hong. Towards k-means-friendly spaces: simultaneous deep learning and clustering. *arXiv:1610.04794*, 2016.