
DEEP EMBEDDED MULTI-VIEW CLUSTERING WITH COLLABORATIVE TRAINING

A PREPRINT

Jie Xu^a, Yazhou Ren^{a,*}, Guofeng Li^a, Lili Pan^a, Ce Zhu^a and Zenglin Xu^{b,c}

^a*University of Electronic Science and Technology of China, Chengdu, China*

^b*Harbin Institute of Technology, Shenzhen, China*

^c*Peng Cheng Lab, Shenzhen, China*

July 28, 2020

ABSTRACT

Multi-view clustering has attracted increasing attentions recently by utilizing information from multiple views. However, existing multi-view clustering methods are either with high computation and space complexities, or lack of representation capability. To address these issues, we propose deep embedded multi-view clustering with collaborative training (DEMVC) in this paper. Firstly, the embedded representations of multiple views are learned individually by deep autoencoders. Then, both consensus and complementary of multiple views are taken into account and a novel collaborative training scheme is proposed. Concretely, the feature representations and cluster assignments of all views are learned collaboratively. A new consistency strategy for cluster centers initialization is further developed to improve the multi-view clustering performance with collaborative training. Experimental results on several popular multi-view datasets show that DEMVC achieves significant improvements over state-of-the-art methods. The code and datasets are available at <https://github.com/JieXuUESTC/DEMVC>.

Keywords Deep embedded clustering; Multi-view clustering; Unsupervised learning; Collaborative training

1 Introduction

Cluster analysis is a fundamental unsupervised learning task in machine learning, which categorizes data samples without labels based on their association with each other [3, 4, 20, 21]. Recently, clustering methods based on deep neural networks (DNN) have achieved impressive clustering performance [32, 35, 38, 23, 19, 22]. However, these methods typically solve a single-view clustering problem.

In many real-world clustering tasks, a data example often has different observable views. For instance, a webpage can be described by both page content and linkage, an object can be captured with different poses, a handwritten digit can be written by different persons. To significantly make use of multiple views' complementary information to enhance the clustering performance, multi-view clustering (MVC) has been proposed. The consensus and complementary principles are two basic concepts in multi-view clustering [9]. On the one hand, since multiple views are exactly multiple maps of the same object, the consensus principle seeks to make multiple predictions of the same object consistent among multiple views. On the other hand, due to the diversity of different views, the complementary principle aims at comprehensively utilizing the complementary information of all views to make better predictions.

The key of multi-view clustering is to effectively mine the information contained in multiple views to achieve better clustering performance. Multi-view information includes common information and complementary information among multiple views. Common information refers to the similar information contained in multiple views. For example, both of two pictures about cats have contour and facial features. The common information of multiple views is

*Corresponding author: yazhou.ren@uestc.edu.cn

helpful to improve the understanding of the commonness of the research objects. Complementary information means that multiple views have specific information about the same object. One view often contains incomplete information, while complementary information can complement each other. For example, one view shows the side of a cat and the other shows the front of the cat, these two views allow for a more complete depiction of the cat. In order to fully depict the cat, all these scattered complementary information in different views is useful. In our study, extracting common information and complementary information are corresponding to the consensus and complementary principles, respectively.

In general, the multi-view clustering approaches can be divided into four categories as below: (1) canonical correlation analysis based MVC, e.g., [12, 10], associates two related views to explore information that is conducive to clustering; (2) subspace clustering based MVC, e.g., [14, 13], explores a shared representation of multiple views to obtain a similarity metric matrix for spectral clustering; (3) matrix factorization based MVC, e.g., [18, 17], decomposes each view into a low-rank matrix with specific constraints and then applies a specific clustering algorithm; (4) graph based MVC, e.g., [15, 11], uses multiple views' information to construct graphs for clustering.

Although existing MVC methods have been successfully applied in various fields, they still have two main disadvantages. (1) Traditional shallow MVC approaches have limited representation capability and are not applicable for many applications, e.g., image clustering. (2) Existing methods usually solve spectral clustering or matrix factorization problems, leading to high computation and space complexities. Thus, these methods can not handle large-scale data clustering tasks.

To address the above mentioned issues, we propose a novel MVC model in this work, namely deep embedded multi-view clustering with collaborative training (DEMVC). Considering a simple clustering of two views, we set one view as the referred view and use its objective to guide the training of itself and the other view. We expect that the view with better clustering performance will serve as a guide to the other view, so as to better mine the common information and complementary information in both views that is beneficial to clustering and enhance the clustering performance. This idea can be extended to clustering of multiple views.

To this end, DEMVC firstly performs feature learning individually for each view by employing deep autoencoders — with good representation capability — to obtain the embedded feature representations. Secondly, DEMVC applies k -means on one view (which is named the referred view) to obtain the initial cluster centers in the embedded space and then computes the auxiliary target distribution of this view. This auxiliary distribution of the referred view is used to refine the deep autoencoders and clustering soft assignments for all views. Each view will become the referred view in sequence to ensure that the multi-view clustering takes full advantage of all views. Performing such collaborative training, all views share the same auxiliary target distribution in every round, and each view can learn from its own view and the other views. Finally, the clustering assignments of all views are summarized to generate the final clustering result. It is verified that the consensus and complementary principles of MVC are guaranteed by the proposed framework.

In summary, the contributions of this paper are three-folds:

- We propose a novel deep embedded multi-view clustering method, which can well utilize the common and complementary information of multiple views by training multiple deep neural networks collaboratively.
- A shared scheme of the auxiliary distribution and a new consistency strategy of cluster centers initialization are developed to improve the performance of MVC.
- The proposed model has good representation capability. In addition, it can be solved efficiently and applied to large-scale datasets. Experiments on several popular datasets demonstrate that DEMVC achieves state-of-the-art performance.

2 Related Work

Deep clustering. In recent years, a number of clustering methods based on deep neural networks have been proposed. In deep embedded clustering (DEC) [32], the cluster assignment and the deep autoencoders are jointly learned. To avoid distortion of the embedded space, [34] proposed an improved version of DEC (IDEC). [33] proposed a deep convolutional embedded clustering algorithm to improve the performance of DEC on image data. Specifically, it uses a convolutional autoencoder to learn the embedded feature space in an end-to-end manner. [35] further improved the performance of DEC by stacking multinomial logistic regression function on top of a multi-layer convolutional autoencoder. [27] introduced a deep neural network with a novel self-expressive layer to improve the traditional subspace clustering. [37] introduced a method that uses deep neural networks to update multiple subspaces and simultaneously combine the reconstruction error to obtain the embedded space, which realizes end-to-end subspace

clustering. [26] applied DEC to semi-supervised clustering. [39] proposed a cooperative subspace clustering method, which can find clusters of data points from a joint low-dimensional subspace.

In this paper, we follow the idea of deep embedded clustering and use the same deep autoencoder as [34, 36, 19], to learn the representations and clustering assignments of samples in low-dimensional embedded space for each view. In addition, we propose a multi-view collaborative training strategy to apply deep embedded clustering to multi-view learning.

Multi-view clustering. Canonical correlation analysis (CCA) [7] is used to find linear projections of two maximally correlated random vectors. [12] explored linearly correlated representation by learning nonlinear transformations of two views with deep canonical correlation analysis (DCCA). [10] proposed deep canonically correlated autoencoders (DCCA), which is an improved version of DCCA. Subspace multi-view clustering approaches assume that multiple views of data come from the same latent space. [14] proposed a diversity-induced multi-view clustering method by extending the traditional subspace clustering. [16] proposed a multi-view clustering method by learning the potential subspace representations of samples. [13] applied deep learning to multi-modal subspace clustering. [28] presented a multi-view sparse subspace clustering method. Combining with convolutional autoencoders and CCA-based self-expressive module, [29] introduced a deep multi-view sparse subspace clustering. [17] proposed to encode collaboratively descriptors of multi-view images into a binary code space. [8] used different kinds of autoencoders to learn multiple deep embedded features and clustering assignments with multi-view fusion mechanism. [11] proposed a multi-view clustering method to firstly learn a connection graph for each view and then minimize the discrepancy of pairwise connection graphs. [15] explored Laplacian rank constrained graph and proposed a self weighted multi-view clustering method. [25] proposed a novel multi-view co-clustering method, which learns optimal weight for bipartite graphs automatically. [24] proposed a three-way multi-view data clustering (uncertain, belong-to and not belong-to) via low-rank matrices. [31, 30] applied self-paced learning in multi-view clustering to address the non-convexity issue.

Based on the consensus principle and complementary principle of multi-view clustering, this paper proposes a novel deep embedded multi-view clustering method. Unlike the above mentioned MVC methods, our approach achieves the consistency of multi-view prediction by a novel collaborative training strategy which shares an auxiliary distribution alternately. Experiments show that this method can mine both the common information and complementary information in multi-view data, and can improve the multi-view clustering performance significantly.

3 Proposed Method

This section presents our deep embedded multi-view clustering with collaborative training (DEMVC) in detail.

3.1 Multi-view Collaborative Training

Consider clustering the dataset $\{\mathbf{x}_i^v \in \mathbb{R}^D\}_{i=1}^N$ into K clusters, where N is the number of samples and D is the dimensionality. Let V be the number of views and represent that there are V subsamples for each object to be clustered. For the v -th view ($v = 1, 2, \dots, V$), let f_Θ^v and g_Ω^v be the encoder and decoder, respectively. Θ and Ω are the corresponding learnable parameters. Based on the nonlinear mapping ability of neural networks, f_Θ^v and g_Ω^v realize:

$$\mathbf{z}_i^v = f_\Theta^v(\mathbf{x}_i^v) \quad (1)$$

and

$$\hat{\mathbf{x}}_i^v = g_\Omega^v(\mathbf{z}_i^v) = g_\Omega^v(f_\Theta^v(\mathbf{x}_i^v)), \quad (2)$$

where $\mathbf{z}_i^v \in \mathbb{R}^d$ is the embedded point, which is encoded by f_Θ^v , of \mathbf{x}_i^v in the low d -dimensional feature space. g_Ω^v decodes \mathbf{z}_i^v and reconstructs the sample as $\hat{\mathbf{x}}_i^v \in \mathbb{R}^D$. For multi-view clustering, we define the loss function as:

$$L = \sum_{v=1}^V L_r^v + \gamma \sum_{v=1}^V L_c^v, \quad (3)$$

where L_r^v and L_c^v are the reconstruction loss and clustering loss of the v -th view, respectively. γ is a trade-off coefficient. In fact, Eq. (3) can be considered as a multi-view generalized version of IDEC [34]. Following IDEC, the trade-off parameter γ is set to 0.1. Combined with Eq. (2), the reconstruction loss is defined as:

$$L_r^v = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i^v - g_\Omega^v(f_\Theta^v(\mathbf{x}_i^v))\|_2^2. \quad (4)$$

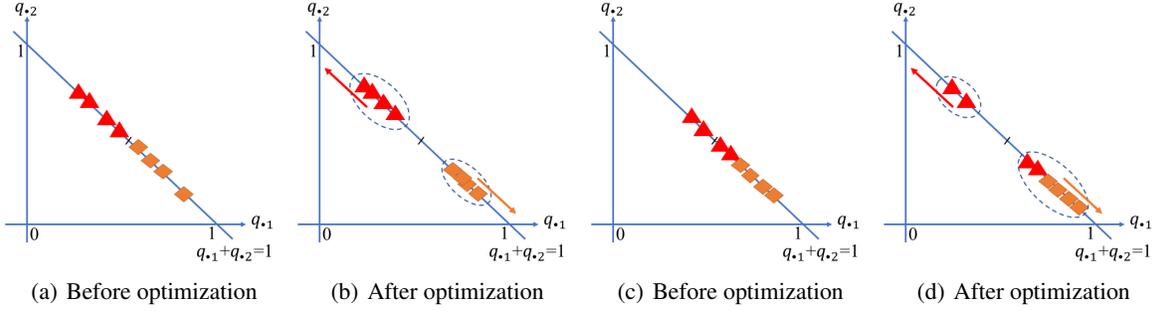


Figure 1: The mechanism of deep embedded clustering to minimize KL divergence. Triangles and rhomboids represent different classes. The two coordinates ($q_{\bullet 1}$ and $q_{\bullet 2}$) of a point represent the probabilities that the corresponding sample belongs to two classes, respectively. Obviously, $q_{\bullet 1} + q_{\bullet 2} = 1$ holds.

Define $\mu_j^v \in \mathbb{R}^d$ as the j -th cluster center of the v -th view. The Kullback-Leibler (KL) divergence is used to define the clustering loss:

$$L_c^v = KL(P^v || Q^v) = \sum_{i=1}^N \sum_{j=1}^K p_{ij}^v \log \frac{p_{ij}^v}{q_{ij}^v}, \quad (5)$$

where $q_{ij}^v \in \mathbb{R}^K$ is the similarity between the embedded point z_i^v and cluster center μ_j^v and is calculated by Student's t -distribution [1] as:

$$q_{ij}^v = \frac{(1 + \|z_i^v - \mu_j^v\|^2)^{-1}}{\sum_j (1 + \|z_i^v - \mu_j^v\|^2)^{-1}}. \quad (6)$$

In a clustering task, q_{ij}^v is treated as the soft label that represents the probability of assigning the i -th sample of the v -th view to the j -th category. With the operation of square and normalization of clustering soft label q_{ij}^v , DEC [32] built the auxiliary target $p_{ij} \in \mathbb{R}^K$ to implement deep single-view clustering. Similarly, p_{ij}^v in our deep multi-view clustering is calculated by:

$$p_{ij}^v = \frac{(q_{ij}^v)^2 / \sum_i q_{ij}^v}{\sum_j ((q_{ij}^v)^2 / \sum_i q_{ij}^v)}. \quad (7)$$

The autoencoders have good representation capability, but the learned representations may not suitable for clustering. We regard the soft label q as a point in the K -dimensional space. By minimizing the KL divergence of Q and P , DEC and IDEC refine autoencoders such that the soft label q is more distinguishing. We take 2-class clustering for example as shown in Figure 1(a) and (b). After minimizing the KL divergence, the two clusters become better separated. Therefore, the refined encoders can acquire more discriminative clustering capabilities, this is why DEC and IDEC can achieve impressive clustering performance. However, they optimize the KL divergence according to a single view. In this way, those hard samples (the samples that are close to other classes, that are fuzzy and hard to distinguish) are prone to misclassification. As shown in Figure 1(c) and (d), some samples which near the classification boundary might be misclassified.

If Eq. (3) is optimized directly, the embedded features and clustering assignments of each view will be learned independently. Thus, the complementary information of multiple views is ignored. But, as we mentioned previously, more complementary information should be used and the clustering prediction of multiple views should be consistent as far as possible. Hence, in order to use the common information and complementary information of multiple views, we let each view become the referred view in turn to guide the whole networks to learn the features, which are conducive to clustering. This training idea is called multi-view collaborative training. Specifically, we define P^r as the auxiliary target distribution of the referred view. In Eq. (5), let P^r be the shared auxiliary target distribution for all views. Then, the clustering loss of the v -th view is:

$$L_c^v = KL(P^r || Q^v). \quad (8)$$

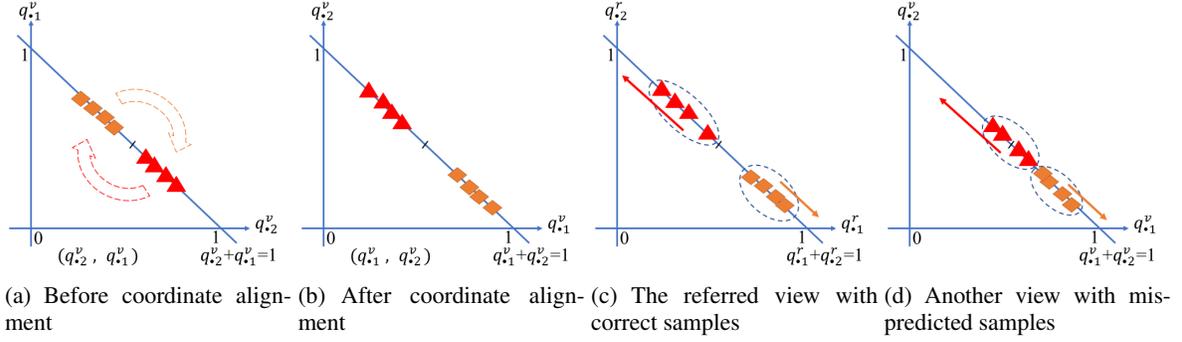


Figure 2: The mechanism by which DEMVC works.

When a view becomes the referred view, all views are collaboratively trained for a certain amount of iterations. Let $(q_{\bullet,1}^r, q_{\bullet,2}^r)$ represent the referred view’s coordinates, as shown in Figure 2(a) and (b). Since the proposed collaborative training shares the same auxiliary target distribution P^r across all views, it can align the coordinates of the other views according to the referred view’s coordinate system. When setting the referred view in multiple views, the accurately predicted samples of the referred view can correct the mispredicted samples in other views, as shown in Figure 2(c) and (d). However, the referred view may also have some mispredicted samples which may mislead other views. Therefore, it is necessary to change the referred view alternately during collaborative training, so that multiple views can supervise each other to obtain more accurate Q^v . In this way, every view can become the referred view in turn to mine the useful complementary information of multiple views to obtain better clustering performance. This corresponds to the complementary principle of multi-view clustering.

After incorporating Eq. (8), the new clustering loss function of collaborative training, Eq. (3) becomes:

$$L = \sum_{v=1}^V \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i^v - g_{\Omega}^v(f_{\Theta}^v(\mathbf{x}_i^v))\|_2^2 + \gamma \sum_{v=1}^V KL(P^r \| Q^v). \quad (9)$$

For the collaborative training — where multiple views in turn become the referred view, it is necessary to optimize the reconstruction loss, as the first term of Eq. (9), to maintain the representation capability of the autoencoders for all views. Otherwise, because of the differences between views, the referred view tend to destroy the capability of the other autoencoders to extract the information from their views. In this way, the autoencoders cannot accurately extract the common information and complementary information of multiple views, resulting in poor clustering performance. Therefore, both clustering loss and reconstruction loss are retained in DEMVC.

The autoencoders, in the beginning, with random network parameters do not have representation capability of the input data. In order to avoid that the referred view guides blindly the training of all views, we first pre-train deep autoencoders of all views via minimizing Eq. (4). Then, we collaboratively fine-tune all autoencoders, cluster assignments and cluster centers of all views by minimizing Eq. (9). Please refer to Section 4.3 for implementation details. The training framework of our model is shown in Figure 3.

When the fine-tuning phase is finished, based on the soft label q_{ij}^v , the clustering prediction of i -th sample of v -th view is computed by:

$$s_i^v = \arg \max_j (q_{ij}^v). \quad (10)$$

In fact, with the influence of shared auxiliary target distribution P^r , the final status of collaborative training is that the distribution of clustering soft labels Q^v of multiple views are similar to each other. For most samples, the predictive soft labels q_{ij}^v of their multiple views are consistent and aligned, corresponding to the consensus principle of multi-view clustering. The final prediction is obtained by averaging the multi-view clustering soft labels as:

$$s_i = \arg \max_j \left(\frac{1}{V} \sum_{v=1}^V q_{ij}^v \right). \quad (11)$$

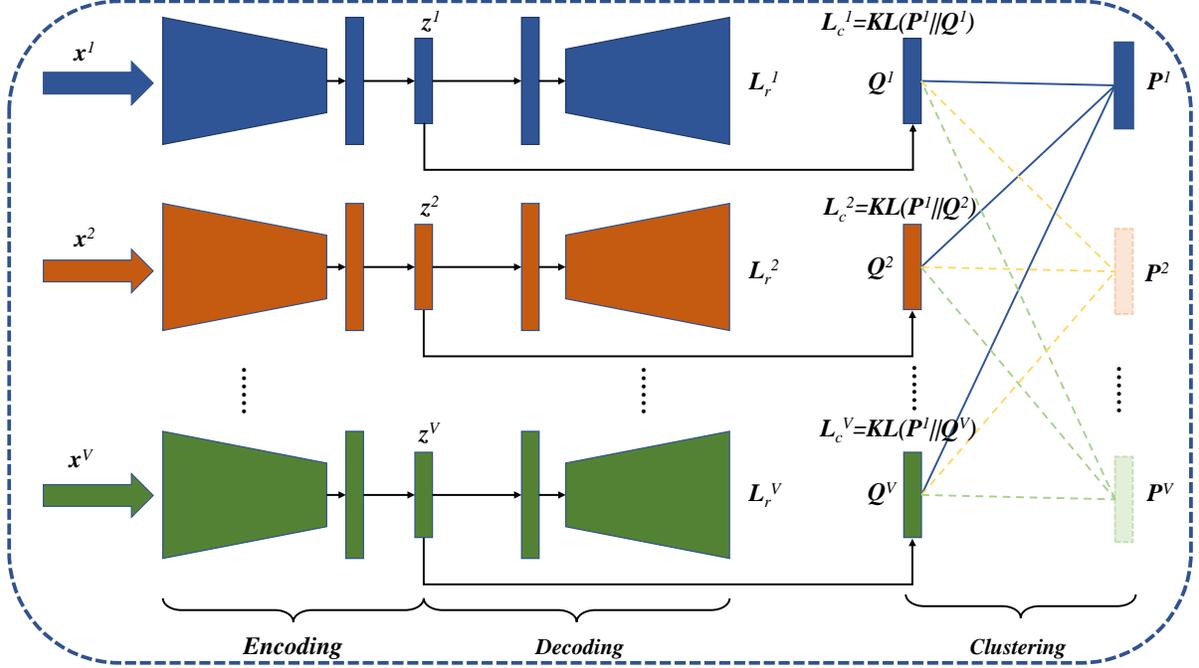


Figure 3: The framework of DEMVC. It consists of V encoders, V decoders, and V clustering layers. The diagram shows that the first view is becoming the referred view as shown with solid lines. The dashed lines represent the corresponding P^i 's do not participate in calculation in this round. The referred view will change sequentially.

According to Eq. (11), for the samples with consistent prediction of multiple views, the final prediction is the same. For those samples whose predictions are inconsistent, the final prediction have the chance to correct the misaligned prediction by integrating the soft labels of multiple views.

3.2 Consistency strategy for cluster centers initialization

According to Eq. (6), q_{ij}^v is calculated by z_i^v and μ_j^v . Before the fine-tuning phase, consider setting the cluster centers μ_j^v of multiple views to be the same to better follow the consensus principle. In this way, multiple views are not limited to their own cluster centers and are easier to accept the guidance of the referred view. We use k -means [3] to initialize the clustering centers in the first referred view (denoted by V_s). The corresponding loss function is:

$$L_{k\text{-means}}^{V_s} = \sum_{i=1}^N \sum_{j=1}^K \left\| \mathbf{z}_i^{V_s} - \mathbf{c}_j^{V_s} \right\|^2, \quad (12)$$

where $\mathbf{c}_j^{V_s} \in \mathbb{R}^d$ ($j = 1, 2, \dots, K$) is the j -th cluster center of the referred view in the low d -dimensional space. Then, we let:

$$\mu_j^v = \mathbf{c}_j^{V_s}, \forall v \in \{1, 2, \dots, V\}. \quad (13)$$

In Eq. (6), the μ^v of all views are initialized with the same cluster centers. Note that the cluster centers of the referred view are shared by other views only in the initialization phase. In the fine-tuning phase, the cluster centers of each view are learned by multi-view collaborative training and only the auxiliary target distribution of the referred view are shared.

The time and space complexities of DEMVC are in the same level with IDEC, which are linear to N . This allows our algorithm to deal with large-scale data clustering problems. The proposed DEMVC model is summarized in Algorithm 1.

Algorithm 1 Deep Embedded Multi-view Clustering with Collaborative Training

Input:Multi-view dataset, number of clusters K **Output:**Multi-view cluster assignment \mathbf{s} Reconstructed samples $\hat{\mathbf{x}}^v$

- 1: *//Initialization phase*
 - 2: Pre-train deep autoencoders by Eq. (4)
 - 3: Initialize cluster centers by Eq. (13)
 - 4: *//Fine-tuning phase*
 - 5: **while** not reaching the maximum iterations **do**
 - 6: **for** V_r in $\{1, 2, \dots, V\}$ **do**
 - 7: Calculate multi-view prediction Q^v by Eq. (6)
 - 8: Calculate P^r of the referred view V_r by Eq. (7)
 - 9: Fine-tune all the deep autoencoders by Eq. (9)
 - 10: **end for**
 - 11: **end while**
 - 12: Output \mathbf{s} and $\hat{\mathbf{x}}^v$ calculated by Eq. (11) and Eq. (2)
-

4 Experimental Setup

4.1 Datasets

NosiyMnist-RotatingMnist (Noisy-Rotating). The Mnist dataset [5] collects 70,000 samples of 28×28 pixel size from 10 classes, i.e., digits 0-9. Following [10], we construct RotatingMnist by randomly rotating the images with angles uniformly sampled from $[-\pi/4, \pi/4]$. To build NosiyMnist, for each sample in RotatingMnist, we randomly select an image with the same label from the Mnist dataset. Then, each pixel is masked with independent random noise uniformly sampled from $[0, 255]$. After that, the values of all pixels are truncated to $[0, 255]$. In multi-view clustering, NosiyMnist and RotatingMnist are two different views corresponding to each other.

Mnist-USPS. USPS is also a handwritten digital dataset, each sample of which is a 16×16 image. Mnist and USPS are treated as two different views of digits. We use the same dataset as [11] did, each view of which contains 5000 digits. Every class provides 500 samples.

Fashion-10K. We use the test set of Fashion dataset [6], which consists of 10,000 28×28 gray images. It contains 10 categories, such as T-shirt, Dress, Coat and is a more challenging dataset. We consider this test set as the first view, and then for each sample, we randomly select a sample with the same label from this set to construct the second/third view. Different views of each sample are different individuals from the same category.

Mnist-10K. The test set of Mnist dataset is used as the first view. The second/third views are constructed in the same way as Fashion-10K.

The input features of each dataset are scaled to $[0, 1]$.

4.2 Comparing Methods

We compare our DEMVC against the following multi-view clustering methods on NosiyMnist-RotatingMnist and Mnist-USPS:

- (1) deep canonical correlation analysis (DCCA) [12].
- (2) deep canonically correlated autoencoders (DCCAE) [10].
- (3) diversity-induced multi-view subspace clustering (DiMSC) [14].
- (4) latent multi-view subspace clustering (LMSC) [16].
- (5) binary multi-view clustering (BMVC) [17].
- (6) multi-view clustering without parameter selection (COMIC) [11].

In order to illustrate the significant improvement of our DEMVC compared to single-view deep clustering approaches, we test several state-of-the-art deep clustering methods on Fashion-10K and Mnist-10K:

Table 1: Comparison of multi-view clustering algorithms. The suffixes ‘-V1’ and ‘-V2’ represent the clustering results of the first view (NoisyMnist or Mnist), and the second view (RotatingMnist or USPS), respectively.

Methods	Noisy-Rotating		Mnist-USPS	
	ACC	NMI	ACC	NMI
DCCA (ICML 2013)	97.00 [†]	92.00 [†]	97.42*	93.60*
DCCAE (ICML 2015)	97.50 [†]	93.40 [†]	98.00*	94.70*
DiMSC (CVPR 2015)	-	-	48.34*	36.02*
LMSC (CVPR 2017)	-	-	78.60*	78.49*
BMVC (TPAMI 2018)	85.61	81.48	88.68	89.93
COMIC+SC (ICML 2019)	-	-	97.44*	94.83*
DEMVC-V1	99.51	98.37	99.71	99.24
DEMVC-V2	99.08	97.16	99.81	99.35
DEMVC	99.87	99.53	99.83	99.49

- (1) deep embedded clustering (DEC) [32].
- (2) improved deep embedded clustering (IDEC) [34].
- (3) deep embedded clustering with data augmentation (DEC-DA) [36].
- (4) deep clustering network (DCN) [38].
- (5) k -subspace clustering network (k -SCN) [37].
- (6) neural collaborative subspace clustering (NCSC) [39].

4.3 Implementation Details

Network settings. All the used autoencoders are the same convolutional autoencoder. Following deep single-view clustering methods [36, 19], the structure of encoder is: Input \rightarrow Conv₃₂⁵ \rightarrow Conv₆₄⁵ \rightarrow Conv₁₂₈³ \rightarrow Fc₁₀. That is, the convolution kernel sizes are 5-5-3 and stride of 2 as default, and channels are 32-64-128. The dimensionality is reduced to 10 since the embedded layer Fc₁₀ is a fully connected network, which is made up of 10 neurons. The encoders and decoders of multiple views are symmetric correspondingly. The ReLU is always chosen as the activation function except for the input, embedded, output and clustering layers.

We use Adam and default parameters in Keras¹ to optimize the entire networks. The autoencoders of DEMVC are pre-trained for 500 epochs for each view. In the fine-tuning phase, when a view becomes the referred view, it guides each view (including itself) to train 200 batches in an end-to-end manner. The batch size is 256 and the number of fine-tuning iterations is 20,000. All experiments of DEMVC are performed on Windows PC with Intel (R) Core (TM) i5-9400F CPU @ 2.90GHz, 16.0GB RAM, and GeForce RTX 2060 GPU (6GB caches).

4.4 Evaluation Measures

The quantitative metrics are adjusted rand index (ARI), unsupervised clustering accuracy (ACC), and normalized mutual information (NMI). A larger value of ARI/ACC/NMI indicates a better clustering result. All the reported results (except for those values excerpted from the papers) are the average values of 5 independent runs.

5 Results and Analysis

5.1 Results on Real Data

Table 1 shows the results of comparing methods on NoisyMnist-RotatingMnist and Mnist-USPS. Here, the results marked with ‘*’ and ‘[†]’ are excerpted from [11] and [10], respectively. ‘-’ denotes that the corresponding methods, which are based on subspace clustering or spectral clustering, are with high complexity and can not solve the clustering task of 70,000 data samples. The best three results in each column are highlighted in boldface. It can be seen that

¹<https://github.com/fchollet/keras>

Table 2: Comparison of single-view deep clustering algorithms. The suffixes ‘-2 views’ and ‘-3 views’ indicate that 2 and 3 views are built with the corresponding dataset (Fashion-10K and Mnist-10K), respectively.

Methods	Mnist-10K		Fashion-10K	
	ACC	NMI	ACC	NMI
DEC (ICML 2016)	83.41	79.22	56.70	61.29
IDEC (IJCAI 2017)	84.25	82.77	57.43	61.55
DCN (ICML 2017)	83.31 [‡]	80.86 [‡]	58.67 [‡]	59.40 [‡]
DEC-DA (ACML 2018)	97.93	95.81	53.55	59.91
<i>k</i> -SCN (ACCV 2018)	87.14 [‡]	78.15 [‡]	63.78 [‡]	62.04 [‡]
NCSC (arXiv 2019)	94.09 [‡]	86.12 [‡]	72.14 [‡]	68.60 [‡]
DEMVC-2 views	99.87	99.60	84.75	87.14
DEMVC-3 views	99.99	99.96	78.99	90.88

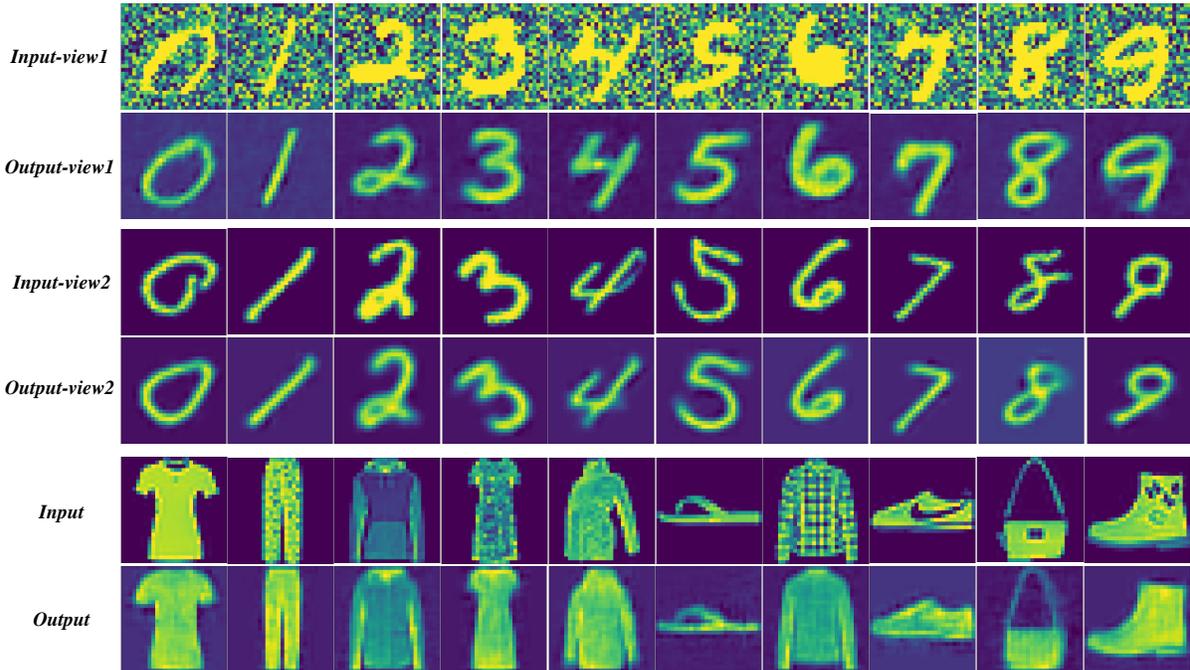


Figure 4: Reconstructed results of multi-view decoders. On NoisyMnist-RotatingMnist, the first and third rows are the corresponding input images of view1 and view2, and the second and fourth rows are the images reconstructed by the corresponding decoders. The fifth and sixth rows are the corresponding input and output images on Fashion-10K. The decoders of DEMVC remove unnecessary noises from the input pictures and make samples look like more standard, indicating their good reconstructed capability.

DEMVC can always achieve the best performance. After sufficient fine-tuning, the soft cluster assignments q_i^1 and q_i^2 are almost the same. In this case, some mispredicted assignments can be corrected to some extent by averaging them. Therefore, DEMVC’s final clustering performance is slightly better than that of the two views (DEMVC-V1 and DEMVC-V2). It is worth noting that the aligned degree of multi-view clustering soft assignments can reflect the progress of collaborative training. When the soft clustering assignments of multiple views are almost consistent, it indicates that DEMVC reaches the consensus principle and meets the stopping criterion.

When comparing with single-view deep clustering methods, we directly apply them on the test set of Mnist and Fashion because they can only deal with one view clustering task. The results are shown in Table 2, where the values with ‘[‡]’ are obtained from [39]. The best two values in each column are highlighted. In the case of two views, the performance of DEMVC is better than all the single-view methods. The performance of DEMVC with three views is better than that with two views in general. This shows that our approach can effectively extract useful feature from multiple views, and also validates the effectiveness of our approach applies in multi-view clustering.

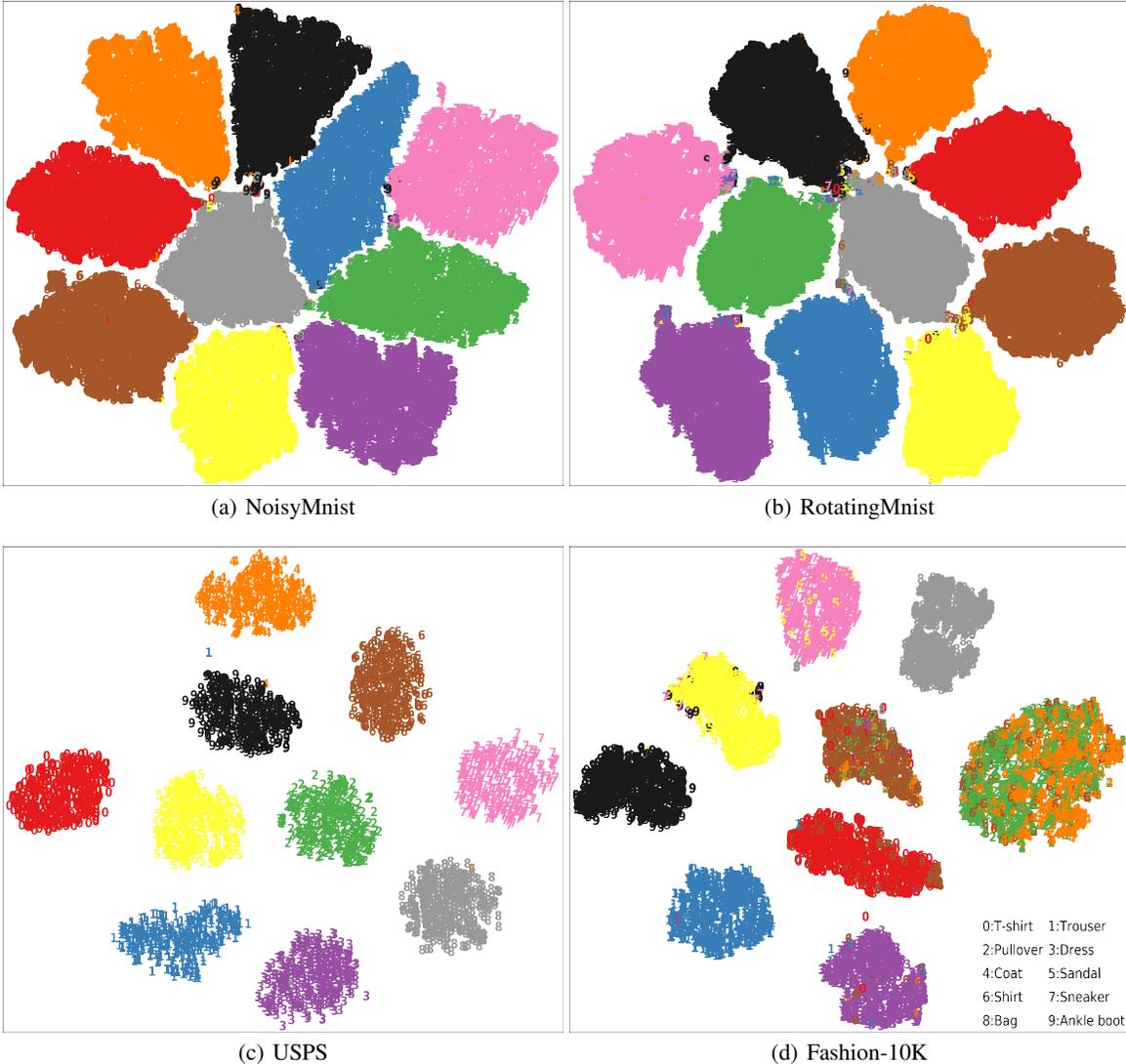


Figure 5: Visualization of the embedded features via t -SNE.

5.2 Visualization of Results

The reconstructed results of DEMVC are shown in Figure 4. Specifically, the image “6”, in the NoisyMnist dataset in the first line, fills the empty space in the hollow position of the digital “6”. After the sample is reconstructed by the decoder, the redundant part of the hollow position of the image “6” is discarded. The handwritten digit “6” has a hollow part, which is the common feature of the images “6” in the training dataset. The autoencoder obtained by DEMVC accurately captures this common information. In addition, the processing of Fashion-10K by the autoencoder focuses on the contour of the extracted object. The reason is that the similarity of clothing products in Fashion-10K mainly lies in their appearance and shape, while the logo or pattern on the clothing products are not so important.

It is shown that, on the noisy digits, rotating digits, and fashionable products, DEMVC’s multi-view decoders can effectively reconstruct the images based on the low-dimensional embedded features of the input samples. The model can even patch up the missing parts while ignoring the unnecessary information of the samples — such as noise, rotation and redundant parts — and finally make the samples look more standard. This indicates its good representation capability of sample features and reconstruction capability, which is the premise to improve clustering performance.

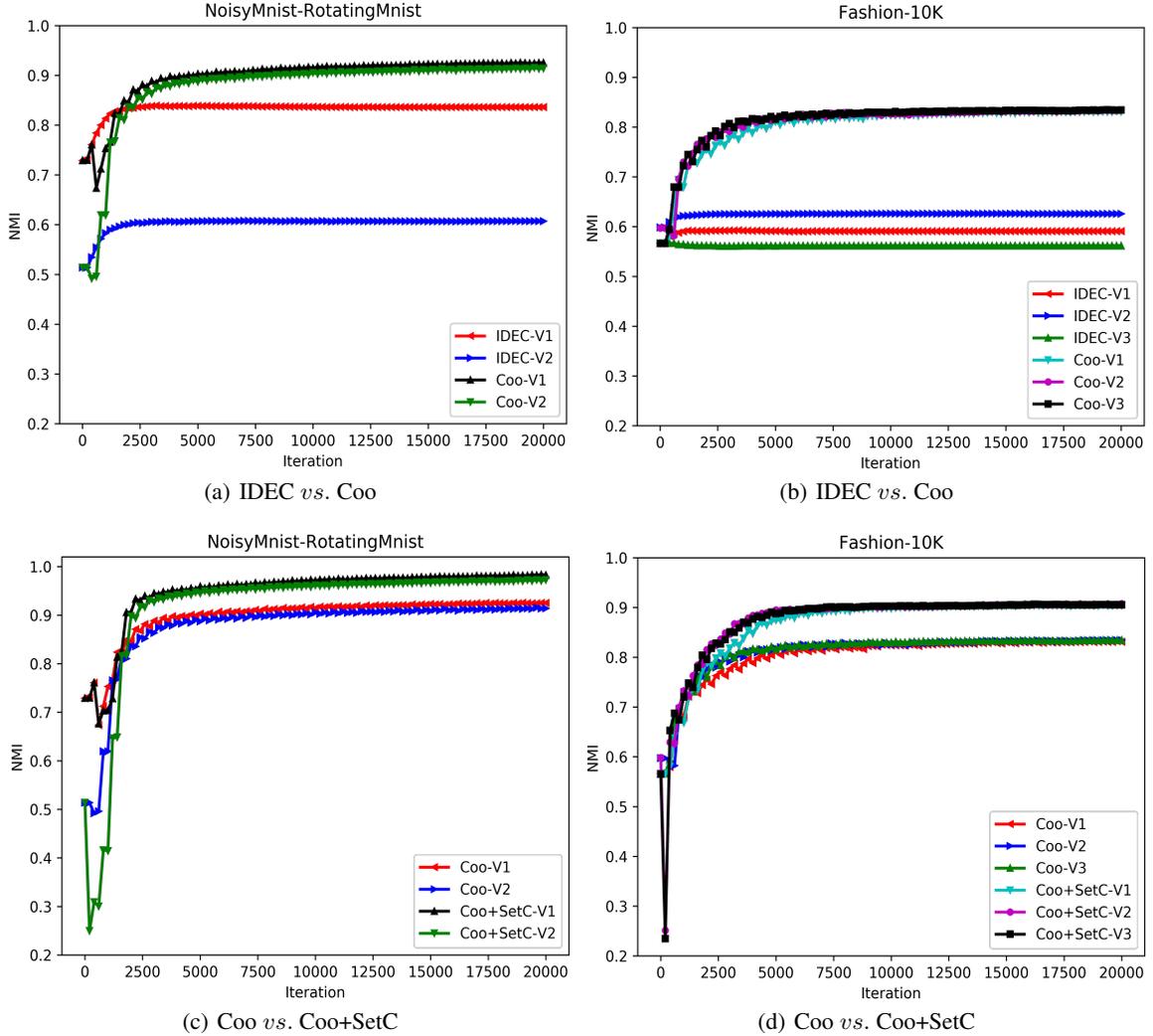


Figure 6: The training process and the comparison of different modules.

We further visualize the embedded features of four datasets in 2-dimensional space via t -SNE [1], as shown in Figure 5. Points of the same class are plotted with the same color. For the sake of demonstration, the ground-truth label of each point is also plotted. We can see that points from the same cluster are highly concentrated and the different clusters are well separated, verifying again the impressive representation capability and clustering performance of DEMVC.

On Fashion-10K, some fashionable products, e.g., pullover, coat, and shirt, look similar in images of 28×28 pixel size. So some points from these classes are closed to each other, as shown in Figure 5(d). Nevertheless, the clustering performance of DEMVC is still significantly better than that of other single-view algorithms. We think that more information (e.g. semi-supervised information) is needed to further separate those samples.

5.3 Module Analysis

This section explores the role of each module of DEMVC. In Figure 6 and Table 3, ‘IDEC-V1’ means IDEC is applied on the first view of the corresponding dataset. ‘Coo’ represents collaborative training without the strategy of sharing cluster centers in the beginning. ‘Coo+SetC’ represents collaborative training method with consistency strategy of cluster centers initialization.

Figure 6(a) and Figure 6(b) show the results of IDEC and Coo on NoisyMnist-RotatingMnist and Fashion-10K (with three views). The proposed collaborative training constantly switches the referred view, so the multiple views

Table 3: The influence of different modules of DEMVC. ‘2 views’ means two views are constructed for the corresponding dataset.

Module	NoisyMnist-RotatingMnist			Mnist-USPS			Mnist-10K(2 views)			Fashion-10K(2 views)		
	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI
IDEC-V1	88.91	83.54	79.67	81.24	78.17	71.36	84.25	82.77	77.56	57.43	61.55	45.31
IDEC-V2	59.22	60.81	44.05	66.20	68.75	55.72	83.58	81.52	76.52	48.84	60.57	40.82
Coo-V1	88.53	93.46	86.51	99.64	98.98	99.20	90.09	96.64	89.69	70.21	81.37	64.10
Coo-V2	87.53	91.76	84.59	99.72	99.23	99.38	90.05	96.54	89.63	70.18	81.25	64.10
Coo+SetC-V1	99.51	98.37	98.92	99.71	99.24	99.25	99.83	99.49	99.62	84.61	86.83	78.49
Coo+SetC-V2	99.08	97.16	97.95	99.81	99.35	99.47	99.74	99.28	99.41	84.73	87.09	78.73
DEMVC	99.87	99.53	99.71	99.83	99.49	99.63	99.87	99.60	99.70	84.75	87.14	78.80

can teach each other such that their clustering performance increases alternately. Overall, for all views, collaborative training significantly outperforms general training, i.e., applying IDEC on each view independently.

Figure 6(c) and Figure 6(d) give the results of Coo and Coo+SetC. It can be observed that the final clustering performance of Coo+SetC is much better than that of Coo, verifying the effectiveness of the developed consistency strategy of cluster centers initialization. For Coo+SetC, except for the referred view, the NMI of other views is low, because the clustering centers of the referred view may not well represent the clustering centers of these views in the beginning. However, in subsequent training, the performance of all views cooperatively increase benefiting from the collaborative training and consistency strategy and get better results.

Results of different modules of DEMVC on four datasets are shown in Table 3. This ablation study demonstrating again that the collaborative training and consistency strategy of initial centers are necessary and useful. Ultimately, DEMVC gets more accurate prediction by averaging the soft labels of multiple views.

6 Conclusion

In this work, we present a novel deep embedded multi-view clustering algorithm (DEMVC). Through multi-view collaborative training, each view can guide all views, in turn, to learn the embedded features. DEMVC also uses a new consistency strategy for cluster centers initialization and follows the consensus and complementary principles of multi-view clustering. The proposed framework can make use of the multi-view common and complementary information to enhance clustering performance. Extensive experiments demonstrate the effectiveness of the proposed model. In addition, DEMVC is of $O(N)$ complexity and can be used for large-scale data clustering.

Acknowledgments

This work was supported in part by the National Key Research and Development Program of China (No. 2018AAA0100204), the National Natural Science Foundation of China (No. 61806043), and the China Postdoctoral Science Foundation (No. 2016M602674).

References

- [1] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 9:2579–2605, 2008.
- [2] François Chollet. keras, github. *GitHub repository*, <https://github.com/fchollet/keras>, 2015.
- [3] James MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, pages 281–297, 1967.
- [4] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, pages 849–856, 2002.
- [5] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [6] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [7] Theodore W Anderson. An introduction to multivariate statistical analysis. Technical report, 1958.

- [8] Bingqian Lin, Yuan Xie, Yanyun Qu, Cuihua Li, and Xiaodan Liang. Jointly deep multi-view learning for clustering analysis. *arXiv preprint arXiv:1808.06220*, 2018.
- [9] Chang Xu, Dacheng Tao, and Chao Xu. A survey on multi-view learning. *arXiv preprint arXiv:1304.5634*, 2013.
- [10] Weiran Wang, Raman Arora, Karen Livescu, and Jeff Bilmes. On deep multi-view representation learning. In *ICML*, pages 1083–1092, 2015.
- [11] Xi Peng, Zhenyu Huang, Jiancheng Lv, Hongyuan Zhu, and Joey Tianyi Zhou. COMIC: Multi-view clustering without parameter selection. In *ICML*, pages 5092–5101, 2019.
- [12] Galen Andrew, Raman Arora, Jeff Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *ICML*, pages 1247–1255, 2013.
- [13] Mahdi Abavisani and Vishal M Patel. Deep multimodal subspace clustering networks. *IEEE Journal of Selected Topics in Signal Processing*, 12(6):1601–1614, 2018.
- [14] Xiaochun Cao, Changqing Zhang, Huazhu Fu, Si Liu, and Hua Zhang. Diversity-induced multi-view subspace clustering. In *CVPR*, pages 586–594, 2015.
- [15] Feiping Nie, Jing Li, and Xuelong Li. Self-weighted multiview clustering with multiple graphs. In *IJCAI*, pages 2564–2570, 2017.
- [16] Changqing Zhang, Qinghua Hu, Huazhu Fu, Pengfei Zhu, and Xiaochun Cao. Latent multi-view subspace clustering. In *CVPR*, pages 4279–4287, 2017.
- [17] Zheng Zhang, Li Liu, Fumin Shen, Heng Tao Shen, and Ling Shao. Binary multi-view clustering. *TPAMI*, 41(7):1774–1782, 2018.
- [18] Handong Zhao, Zhengming Ding, and Yun Fu. Multi-view clustering via deep matrix factorization. In *AAAI*, 2017.
- [19] Yazhou Ren, Ni Wang, Mingxia Li, and Zenglin Xu. Deep density-based image clustering. *Knowledge-Based Systems*, 197:105841, 2020.
- [20] Guo Zhong and Chi-Man Pun. Nonnegative self-representation with a fixed rank constraint for subspace clustering. *Information Sciences*, 2020.
- [21] Maryam Abdolali and Mohammad Rahmati. Neither global nor local: A hierarchical robust subspace clustering for image data. *Information Sciences*, 514:333–353, 2020.
- [22] Nairouz Mrabah, Mohamed Bouguessa, and Riadh Ksantini. Adversarial deep embedded clustering: on a better trade-off between feature randomness and feature drift. *arXiv preprint arXiv:1909.11832*, 2019.
- [23] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep adaptive image clustering. In *Proceedings of the IEEE international conference on computer vision*, pages 5879–5887, 2017.
- [24] Hong Yu, Xincheng Wang, Guoyin Wang, and Xianhua Zeng. An active three-way clustering method via low-rank matrices for multi-view data. *Information Sciences*, 507:823–839, 2020.
- [25] Shudong Huang, Zenglin Xu, Ivor W Tsang, and Zhao Kang. Auto-weighted multi-view co-clustering with bipartite graphs. *Information Sciences*, 512:18–30, 2020.
- [26] Yazhou Ren, Kangrong Hu, Xinyi Dai, Lili Pan, Steven CH Hoi, and Zenglin Xu. Semi-supervised deep embedded clustering. *Neurocomputing*, 325:121–130, 2019.
- [27] Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering networks. In *Advances in Neural Information Processing Systems*, pages 24–33, 2017.
- [28] Maria Brbić and Ivica Kopriva. Multi-view low-rank sparse subspace clustering. *Pattern Recognition*, 73:247–258, 2018.
- [29] Xiaoliang Tang, Xuan Tang, Wanli Wang, Li Fang, and Xian Wei. Deep multi-view sparse subspace clustering. In *Proceedings of the 2018 VII International Conference on Network, Communication and Computing*, pages 115–119, 2018.
- [30] Yazhou Ren, Shudong Huang, Peng Zhao, Minghao Han, and Zenglin Xu. Self-paced and auto-weighted multi-view clustering. *Neurocomputing*, 2019.
- [31] Chang Xu, Dacheng Tao, and Chao Xu. Multi-view self-paced learning for clustering. In *IJCAI*, pages 3974–3980, 2015.
- [32] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *ICML*, pages 478–487, 2016.

- [33] Xifeng Guo, Xinwang Liu, En Zhu, and Jianping Yin. Deep clustering with convolutional autoencoders. In *ICONIP*, pages 373–382. Springer, 2017.
- [34] Xifeng Guo, Long Gao, Xinwang Liu, and Jianping Yin. Improved deep embedded clustering with local structure preservation. In *IJCAI*, pages 1753–1759, 2017.
- [35] Kamran Ghasedi Dizaji, Amirhossein Herandi, Cheng Deng, Weidong Cai, and Heng Huang. Deep clustering via joint convolutional autoencoder embedding and relative entropy minimization. In *ICCV*, pages 5736–5745, 2017.
- [36] Xifeng Guo, En Zhu, Xinwang Liu, and Jianping Yin. Deep embedded clustering with data augmentation. In *ACML*, pages 550–565, 2018.
- [37] Tong Zhang, Pan Ji, Mehrtash Harandi, Richard Hartley, and Ian Reid. Scalable deep k-subspace clustering. In *ACCV*, pages 466–481. Springer, 2018.
- [38] Bo Yang, Xiao Fu, Nicholas D Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *ICML*, pages 3861–3870, 2017.
- [39] Tong Zhang, Pan Ji, Mehrtash Harandi, Wenbing Huang, and Hongdong Li. Neural collaborative subspace clustering. *arXiv preprint arXiv:1904.10596*, 2019.