

Multi-view Deep Subspace Clustering Networks

Pengfei Zhu, Binyuan Hui, Changqing Zhang, Dawei Du, Longyin Wen, Qinghua Hu

Abstract—Multi-view subspace clustering aims to discover the inherent structure by fusing multi-view complementary information. Most existing methods first extract multiple types of hand-crafted features and then learn a joint affinity matrix for clustering. The disadvantage lies in two aspects: 1) Multi-view relations are not embedded into feature learning. 2) The end-to-end learning manner of deep learning is not well used in multi-view clustering. To address the above issues, we propose a novel multi-view deep subspace clustering network (MvDSCN) by learning a multi-view self-representation matrix in an end-to-end manner. MvDSCN consists of two sub-networks, i.e., diversity network (Dnet) and universality network (Unet). A latent space is built upon deep convolutional auto-encoders and a self-representation matrix is learned in the latent space using a fully connected layer. Dnet learns view-specific self-representation matrices while Unet learns a common self-representation matrix for all views. To exploit the complementarity of multi-view representations, Hilbert Schmidt Independence Criterion (HSIC) is introduced as a diversity regularization, which can capture the non-linear and high-order inter-view relations. As different views share the same label space, the self-representation matrices of each view are aligned to the common one by a universality regularization. Experiments on both multi-feature and multi-modality learning validate the superiority of the proposed multi-view subspace clustering model.

Index Terms—subspace clustering, multi-view learning, self-representation, deep clustering.

I. INTRODUCTION

Subspace clustering aims to segment a set of unlabeled samples drawn from a union of multiple subspaces corresponding to different clusters into several groups. Recently, self-representation based models have achieved superior performance in subspace clustering [1], [2], [3]. It assumes that a sample can be represented by a linear combination of a set of samples:

$$\min_{\mathbf{Z}} L(\mathbf{X}, \mathbf{Z}) + R(\mathbf{Z}), \quad s.t. \quad \mathbf{X} = \mathbf{X}\mathbf{Z}, \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{d \times n}$ and $\mathbf{Z} \in \mathbb{R}^{n \times n}$ denote the training data and self-representation matrix, respectively. $L(\mathbf{X}, \mathbf{Z})$ represents the reconstruction loss and $R(\mathbf{Z})$ is the regularization item. The key differences of self-representation based subspace clustering models lie in the option of the loss function and regularizer. For $R(\mathbf{Z})$, l_0 -norm, l_1 -norm, square of Frobenius norm, elastic net, trace Lasso, and k-block diagonal regularizer have been used under certain subspace assumptions [4]. As hand-crafted features cannot well capture the severe variations,

Pengfei Zhu, Binyuan Hui, Changqing Zhang, and Qinghua Hu are with School of Computer Science and Technology, Tianjin University. Dawei Du is with Computer Science Department, University at Albany, State University of New York. Longyin Wen is with JD Digits. This work was supported by the National Natural Science Foundation of China under Grants 61876127, Natural Science Foundation of Tianjin Under Grants 17JCZDJC30800, 18YFZCGX00390, 18YFZCGX00680, and Young Elite Scientists Sponsorship Program by Tianjin.

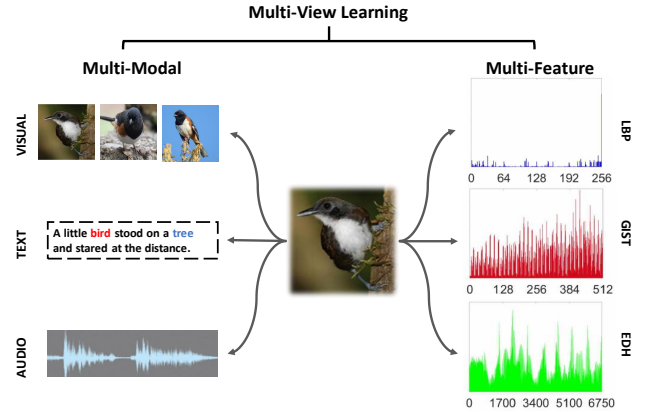


Fig. 1: Examples of multi-view learning. A sample can be represented by different modalities, e.g., image, video and text. Different kinds of features, e.g., SIFT, Gabor, and deep features can be extracted. Generally, multi-view learning covers multi-modal and multi-feature learning.

deep subspace clustering models are developed to jointly learn hierarchical representation and cluster structure [5], [6], [7], [3], [8], [9], [10], [11], [12], [13].

The rapid growth of digital sensors and widespread application of social networks bring about the explosion of multi-modal data in multimedia analysis [14], medical image analysis [15], autonomous driving [16], etc. As shown in Figure 1, different types of data can be collected, including text, image, audio and video, to represent a sample. Even for single-modal data, e.g., images or video sequences, diverse features can be extracted to capture scale, occlusion, illumination, rotation variations for robust recognition [17]. Generally, multi-view learning covers multi-modal and multi-feature learning. Multi-modal and multi-feature information can be fused to boost the performance of subspace clustering.

Multi-view subspace clustering (MVSC) aims to utilize data collected from different modalities or represented by different types of features to discover the underlying clustering structure. Most MVSC methods design multi-view regularizer to characterize the inter-view relationships between several types of hand-crafted features for multi-view clustering [18], [19], [20], [21], [22], [17], [23]. However, their performances are still not satisfactory for two reasons. Firstly, existing methods adopt a two-stage strategy, i.e., first extracting features and then learning the affinity matrix. The features extraction process is irrelevant to the subspace clustering task. Multi-view relationships can only work during the affinity matrix learning process, which ignores the role of inter-view relations in feature learning. Secondly, they consider little about hierarchical representation learning in an end-to-end manner.

In this paper, we propose a multi-view deep subspace

clustering networks (MvDSCN) by learning a multi-view self-representation matrix in an end-to-end manner. MvDSCN is composed of diversity network (Dnet) and universality network (Unet). Dnet learns view-specific self-representation matrices ($\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_v$) while Unet learns a common self-representation matrix \mathbf{Z} . Deep convolutional auto-encoders are learned for each view with either handcrafted features or raw data as the input. Multi-view reconstruction and self-representation losses are simultaneously minimized. To exploit the complementary multi-view information, a diversity regularizer is defined based on Hilbert Schmidt Independence Criterion (HSIC). Additionally, we used a universality regularizer to make view-specific \mathbf{Z}_i close to the common \mathbf{Z} . With diversity and universality regularization, multi-view relations are well embedded in both feature learning and self-representation stages. Experiments on both multi-feature and multi-modal clustering validate the superiority of the proposed model to the state-of-the-art subspace clustering methods.

The structure of this paper is organized as: Section II introduces the related work on multi-view learning, self-representation and auto-encoders. Section III presents the proposed multi-view clustering model. Section IV conducts experiments on both multi-feature and multi-modal tasks. Section V concludes and gives the future work.

II. RELATED WORK

In this section, we give a brief review of multi-view clustering, self-representation, and auto-encoders.

A. Multi-view Clustering

Subspace clustering aims to uncover the inherent cluster structure from data composed of multiple subspaces [3]. In the past few years, most existing subspace clustering methods focus on learning a good affinity matrix and then conduct spectral clustering. Self-representation based subspace clustering methods are essentially based on the hypothesis that a sample can be reconstructed by a linear combination of other samples. Sparse subspace clustering (SSC) imposes l_1 -norm regularization on the representation coefficients to enhance the sparsity [24]. Low rank representation (LRR) explores the multi-block diagonal property of the self-representation matrix to discover the multiple subspace structure [1], [4]. Deep subspace clustering network embeds self-representation into deep convolutional autoencoder by a fully connected layer [2]. To conduct clustering in an end-to-end manner, clustering loss is proposed to output the clustering results directly [8]. Deep adversarial subspace clustering uses GAN-like model to evaluate the clustering performance besides self-representation loss [3].

Multi-view clustering boosts the performance of clustering by exploring the complementary information by modeling inter-view relations or learning a latent representation. Most existing multi-view clustering methods can be considered as an extension of single-view models, including spectral clustering [20], matrix factorization [17], and k-means [22], etc. Multi-view relations can be generally categorized into universality and diversity [19], [20], [22], [21]. Universality emphasizes

on that all views should be similar while diversity focuses on the inter-view complementarity and therefore induces diverse view-specific representation. Some work build multi-view connections by a common latent representation for clustering and model multi-view relations using neural networks [25]. Deep learning has achieved superior performance in many tasks because of the end to end learning manner to a great extent. However, the existing multi-view subspace clustering methods treat multi-view feature extraction and affinity learning as two separate stages. Additionally, due to the view-specific characteristic, it is unreasonable to force the self-representation matrices of all views to be the same.

B. Self-representation

Self-representation reflects the intra-relations among samples, and has been widely used in image processing, clustering, feature selection, and deep learning. In image processing, especially image denoising, non-local mean has been widely used by reconstructed a pixel or image patch using related pixels or patches in the image [26], which inspires many successful image processing models in low-level vision. Besides pixel-level self-representation, a sample can be well reconstructed by a linear combination of bases. Self-representation has been successfully used in clustering in that it can accurately capture the sample relations by embedding sparse, dense, or low-rank priors [1], [2], [4]. To alleviate the curse of dimensionality, feature selection aims to select a subset of features by evaluating the importance of features. Feature-level self-representation assumes that one feature can be reconstructed by all features, and the self-representation coefficients can be used for feature evaluation [27], [28], [29]. Inspired by the success of non-local mean in image denoising, a non-local neural network is proposed to utilize the relations across elements of feature maps, channels, or frames to improve the representation ability of the networks [30].

C. Auto-Encoders

Auto-encoders (AE) extract features of data by mapping the data to a low-dimensional space. With the rapid development of deep learning, deep (or stacked) auto-encoders have become popular for unsupervised learning. Deep auto-encoders have been widely used in dimensionality reduction [31] and image denoising [32]. Recently, deep auto-encoders have been used to initialize deep embedding networks for unsupervised clustering [33]. The work in [34] uses a fully connected deep auto-encoders by incorporating a sparsity prior into the hidden representation learning to preserve the sparse reconstruction relation. By comparison, [2] directly learn the affinities between all data points through a deep auto-encoder network by using a fully connected self-representation layer.

Since convolutional layers have fewer parameters and stronger learning ability than the fully connected layer, convolutional auto-encoders(CAE) that can be trained in an end to end manner are designed for feature learning from unlabeled data. The work in [35] is the first trial to train CAE directly in an end to end manner without pre-training. Convolutional neural networks can be initialized by a CAE stack [36] which

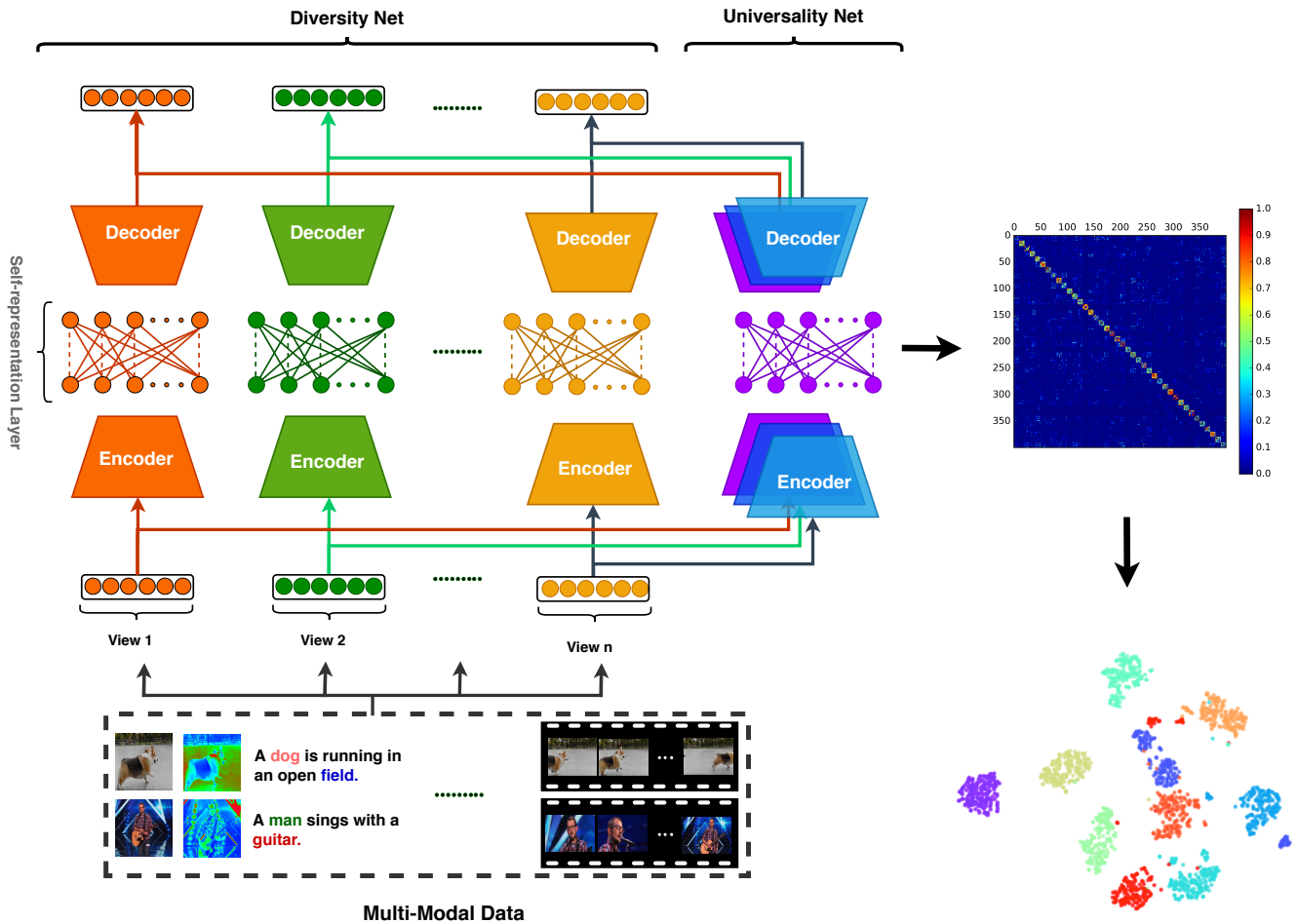


Fig. 2: The network architecture of multi-view deep subspace clustering

is an unsupervised method for hierarchical feature extraction. CAE have been successfully used for generative adversarial networks (GANs). [37] combines a convolutional auto-encoder loss, a GAN loss, and a classification loss defined using a pre-trained classifier. In the field of natural language processing, [38] proposed a general framework for text modeling by embedding a paragraph into a latent representation vector using CAE.

III. MULTI-VIEW DEEP SUBSPACE CLUSTERING

In this section we present the proposed multi-view deep subspace clustering networks (MvDSCN).

A. Network Architecture

Let $\mathbf{X}_1, \dots, \mathbf{X}_i, \dots, \mathbf{X}_v$ denote the inputs of multiple views, where $\mathbf{X}_i \in \mathbb{R}^{n \times d_i}$, v , n and d_i are the number of views, samples, and features in the i^{th} view, respectively. \mathbf{X}_i can be hand-crafted features or raw data, such as image and RGB-D data. The architecture of multi-view deep subspace clustering is shown in Figure 2. The proposed network consists of two parts, i.e., diversity net (Dnet) that learns view-specific representation and universality net (Unet) that learns view-consistent self-representation matrix.

Dnet embeds the input \mathbf{X}_i into the hidden representation \mathbf{F}_i^s by the view-specific encoder for the i^{th} view. Then self-representation is conducted by a fully connected layer without bias and non-linear activations, i.e., $\mathbf{F}_i^s = \mathbf{F}_i^s \mathbf{Z}_i$. Unet uses a common self-representation matrix \mathbf{Z} for all views, which is connected with hidden representation $\mathbf{F}_1^c, \mathbf{F}_2^c, \dots, \mathbf{F}_v^c$ of all views. After self-representation layer, the samples are recovered by the view-specific decoders.

For both Dnet and Unet, we advocate the usage of convolutional auto-encoders with fewer parameters and stronger learning ability rather than the fully connected layer. We use three-layer encoders with [64, 32, 16] channels, and three-layer decoders with [16, 32, 64] channels correspondingly. We adopt a 3×3 kernel and rectified linear unit (ReLU) [39] for the non-linear activations. Notably, no pooling layers are used. The latent features are then back to the space of the same size as the input via the transpose convolution layers.

B. Loss Function

The losses of the proposed MvDSCN consists of two parts, i.e., reconstruction loss by auto-encoders and self-

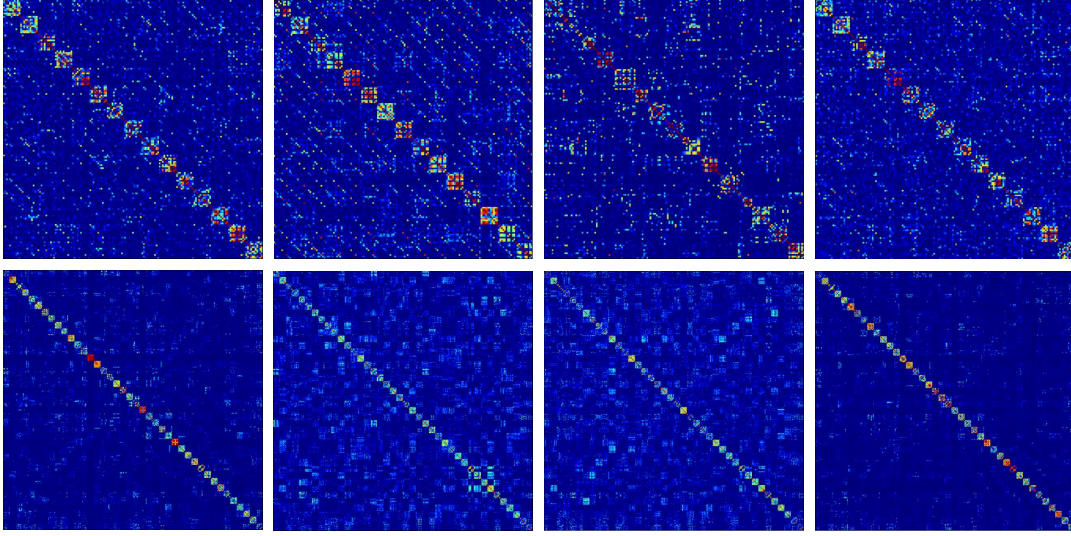


Fig. 3: Visualization of affinity matrices of different views. The first three columns are affinity matrices of view1, view2, view3 learned by DSCN [2]. The last column is the proposed MvDSCN obtained all views. The top row is the result on Yale dataset, and the bottom row is the result on ORL dataset.

representation loss.

$$\min \left\{ \begin{array}{l} \sum_{i=1}^v \left\| \mathbf{X}_i - \hat{\mathbf{X}}_i^s \right\|_F^2 + \left\| \mathbf{X}_i - \hat{\mathbf{X}}_i^c \right\|_F^2 + \\ \sum_{i=1}^v \left\| \mathbf{F}_i^s - \mathbf{F}_i^s \mathbf{Z}_i \right\|_F^2 + \left\| \mathbf{F}_i^c - \mathbf{F}_i^c \mathbf{Z} \right\|_F^2 \end{array} \right\} \quad (2)$$

To embed multi-view relations into feature learning and self-representation, two types of regularizer are used. To exploit the complementary information from multiple views, e.g., RGB and depth information, a diversity regularization is defined based on Hilbert Schmidt Independence Criterion (HSIC). HSIC measures the nonlinear and high-order correlations and has been successfully used in multi-view subspace clustering.

Assuming that there are two variables $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_N]$ and $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_i, \dots, \mathbf{b}_N]$, we define a mapping $\phi(\mathbf{a})$ from $\mathbf{a} \in \mathfrak{A}$ to kernel space \mathfrak{F} , where the inner product of two vectors is defined as $k(\mathbf{a}_1, \mathbf{a}_2) = \langle \phi(\mathbf{a}_1), \phi(\mathbf{a}_2) \rangle$. Then $\varphi(\mathbf{b})$ is defined to map $\mathbf{b} \in \mathfrak{B}$ to kernel space \mathfrak{G} . Similarly, the inner product of two vectors in \mathfrak{G} is defined as $g(\mathbf{b}_1, \mathbf{b}_2) = \langle \phi(\mathbf{b}_1), \phi(\mathbf{b}_2) \rangle$. The empirical version of HSIC is induced as:

Definition 1: Consider a series of N independent observations drawn from $p_{\mathbf{ab}}$, $\mathcal{Z} := \{(\mathbf{a}_1, \mathbf{b}_1), \dots, (\mathbf{a}_N, \mathbf{b}_N)\} \subseteq \mathfrak{A} \times \mathfrak{B}$, an estimator of HSIC, written as $\text{HSIC}(\mathfrak{Z}, \mathfrak{F}, \mathfrak{G})$, is given by:

$$\text{HSIC}(\mathfrak{Z}, \mathfrak{F}, \mathfrak{G}) = (N-1)^{-2} \text{tr}(\mathbf{G}_1 \mathbf{H} \mathbf{G}_2 \mathbf{H}), \quad (3)$$

where $\text{tr}(\cdot)$ is the trace of a square matrix. \mathbf{G}_1 and \mathbf{G}_2 are the Gram matrices with $g_{1,ij} = g_1(\mathbf{a}_i, \mathbf{a}_j)$, $g_{2,ij} = g_2(\mathbf{b}_i, \mathbf{b}_j)$. $h_{ij} = \delta_{ij} - 1/N$ centers the Gram matrix which has zero mean in the feature space. Please refer to [40], [41] for more details about HSIC.

Based on HSIC, the diversity regularizer is defined as

$$R_d(\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_v) = \sum_{ij} \text{HSIC}(\mathbf{Z}_i, \mathbf{Z}_j) \quad (4)$$

The diversity regularizer in Eq. (4) can effectively exploit the complementary information from multiple views. As all views share the same decision space, the view-specific self-representation matrices that reflect sample relations should be aligned with the common self-representation matrix in Unet. We define a centralization regularizer as follows

$$R_c(\mathbf{Z}, \mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_v) = \sum_{i=1}^v \left\| \mathbf{Z} - \mathbf{Z}_i \right\|_F^2 \quad (5)$$

By taking multi-view relations into account, the objective function becomes

$$\min \left\{ \begin{array}{l} \underbrace{\sum_{i=1}^v \left\| \mathbf{X}_i - \hat{\mathbf{X}}_i^s \right\|_F^2 + \left\| \mathbf{X}_i - \hat{\mathbf{X}}_i^c \right\|_F^2}_{\text{auto-encoder loss}} \\ + \lambda_1 \underbrace{\sum_{i=1}^v \left\| \mathbf{F}_i^s - \mathbf{F}_i^s \mathbf{Z}_i \right\|_F^2 + \left\| \mathbf{F}_i^c - \mathbf{F}_i^c \mathbf{Z} \right\|_F^2}_{\text{self-representation loss}} \\ + \lambda_2 \left(\left\| \mathbf{Z} \right\|_p + \sum_{i=1}^v \left\| \mathbf{Z}_i \right\|_p \right) \\ \underbrace{\hspace{10em}}_{\text{lp-norm regularizer}} \\ + \lambda_3 \underbrace{\sum_{i=1}^v \left\| \mathbf{Z} - \mathbf{Z}_i \right\|_F^2}_{\text{universality regularizer}} \\ + \lambda_4 \underbrace{\sum_{ij} \text{HSIC}(\mathbf{Z}_i, \mathbf{Z}_j)}_{\text{diversity regularizer}} \end{array} \right\}, \quad (6)$$

s.t. $\text{diag}(\mathbf{Z}_i) = \mathbf{0}, i = 1, 2, \dots, v, \text{diag}(\mathbf{Z}) = \mathbf{0}$

where $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are positive constants, and $\left\| \mathbf{Z} \right\|_p$ is l_p -norm on \mathbf{Z} . Here we can also consider other types of regularizer, e.g., nuclear norm [1], and block diagonal regularizer [42]. Figure 3 shows the affinity matrix of each view

learned independently by deep subspace clustering network in [2] and the one learned by our proposed MvDSCN on multi-view data. The affinity matrix learned by MvDSCN has better block diagonal property and less noise.

C. Optimization

We use a gradient decent method to solve the problem in Eq. (6). For the back propagation (BP) process, the gradients should be derived for each variable. The encoders and decoders can be updated by standard BP. Here we focus on the updating of self-representation layer. As the optimization problems in Eq. (6) with respect to the view-specific self-representation matrices $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_v$ and view-consistent \mathbf{Z} are convex, we can get the gradients easily.

When we use square of Frobenius norm regularization, the gradient for \mathbf{Z}_i is

$$\frac{\partial L_{\mathbf{Z}_i}}{\partial \mathbf{Z}_i} = 2\lambda_1 \mathbf{F}_i^{sT} \mathbf{F}_i^s (\mathbf{Z}_i - \mathbf{I}) - 2\lambda_3 (\mathbf{Z} - \mathbf{Z}_i) + 2\lambda_2 \mathbf{Z} + \lambda_4 (N-1)^{-2} \sum_{j>i} (\mathbf{H}\mathbf{Z}_j \mathbf{H})^T \quad (7)$$

The gradient for \mathbf{Z} is

$$\frac{\partial L_{\mathbf{Z}}}{\partial \mathbf{Z}} = 2\lambda_1 \sum_{i=1}^v \mathbf{F}_i^{eT} \mathbf{F}_i^e (\mathbf{Z} - \mathbf{I}) + 2\lambda_2 \mathbf{Z} + 2\lambda_3 (\mathbf{Z} - \mathbf{Z}_i) \quad (8)$$

Algorithm 1: The algorithm of MvDSCN.

Input: Unlabeled multi-view data $\{\mathbf{X}_1, \dots, \mathbf{X}_v\}$, hyper-parameter $\lambda_1, \lambda_2, \lambda_3$ and λ_4 , pre-trained epochs n , learning rate α_t , initialize parameter $\theta_{\text{Dnet}}, \theta_{\text{Unet}}$ with random values;

for $j = 1$ to n **do**

 Update auto-encoders of θ_{Dnet} ;
 Update auto-encoders of θ_{Unet} ;

while *not converged* **do**

 Compute the gradient of Eq.(6) and update auto-encoders of $\theta_{\text{Dnet}}, \theta_{\text{Unet}}$;
 Optimize $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_v$, and \mathbf{Z} by Eqs.(7) and (8);

Perform spectral clustering using affinity matrix \mathbf{Z} ;

Output: Clustering result C

We first pre-train the deep auto-encoder without the self-representation layer on all multi-view data because the network is difficult to directly train from scratch and avoid the trivial all-zero solution while minimizing the loss function. We then use the pre-trained parameters to initialize the convolutional encoder-decoder layers of both Dnet and Unet. In the fine-tuning stage, we build a big batch using all the data to minimize the loss function. The model is trained with Adam [43] and an initial learning rate of 0.001. For the regularization hyper-parameters of self-representation loss and l_p -norm regularizer, we usually set $\lambda_1 = 1.0 \times 10^{\frac{k}{10}-3}$ where the k is the number of subspaces, $\lambda_2 = 1.0$, $\lambda_3 = 0.1$, $\lambda_4 = 0.1$.

Our network jointly updates Dnet and Unet. Once the network converges, we can use the parameters of the common

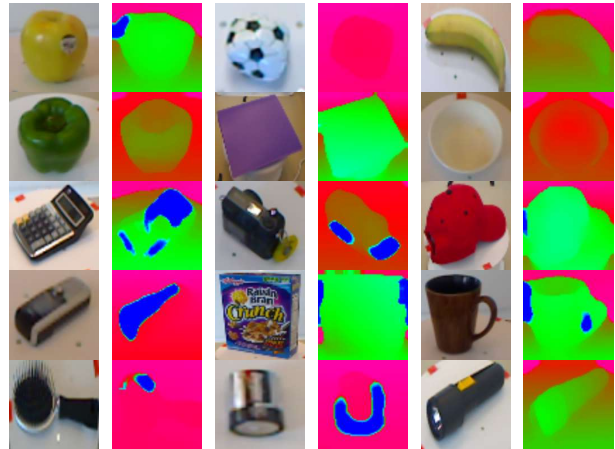


Fig. 4: Sample objects from the RGB-D Object Dataset. RGB image (left) and the corresponding depth image using a recursive median filter(right).

self-expressive layer for all views to construct an affinity matrix $(|\mathbf{Z}| + |\mathbf{Z}|^T) / 2$ for spectral clustering. Similar to [2], since we have no access to labels, our training strategy is unsupervised. Besides, the algorithm for solving MvDSCN is summarized in Alg. 1.

D. Discussions

Recent works on multi-view subspace clustering focus on learning a latent representation across views by dictionary learning or matrix factorization [21], [25], [17]. As shown in Eq. (9), a latent representation \mathbf{C} is learned for \mathbf{X} and then self-representation is conducted on \mathbf{C} . All views can share the same latent representation [21], [25] but view-specific \mathbf{D} .

$$\min_{\{\mathbf{C}, \mathbf{D}, \mathbf{Z}\}} \|\mathbf{X} - \mathbf{C}\mathbf{D}\|_F^2 + \|\mathbf{C} - \mathbf{C}\mathbf{Z}\|_F^2 \quad (9)$$

Compared with latent representation based methods, the proposed MvDSCN also learns a hidden representation \mathbf{F} by auto-encoder ϕ . Auto-encoder can be considered as a mapping function which projects the input to a latent space. MvDSCN has the following two advantages: 1) Compared with the existing shallow models, the hidden representation \mathbf{F} is more informative by using deep convolutional auto-encoders whether the input \mathbf{X} is hand-crafted feature or the raw data. 2) MvDSCN joints feature learning and self-representation together in an end to end manner. Thus, the multi-view relations can guide both affinity matrix learning and feature learning. Hence, our proposed MvDSCN can learn a good affinity matrix and therefore boost the performance of multi-view subspace clustering.

IV. EXPERIMENTS

In this section, extensive experiments are conducted to verify the effectiveness of the proposed clustering model.

TABLE I: Results on four multi-feature datasets (mean \pm standard deviation). Higher value indicates better performance.

Datasets	Methods	NMI	ACC	AR	F-measure
Yale	BestSV	0.654 \pm 0.009	0.616 \pm 0.030	0.440 \pm 0.011	0.475 \pm 0.011
	LRR	0.709 \pm 0.011	0.697 \pm 0.000	0.515 \pm 0.004	0.547 \pm 0.007
	Min-Disagreement	0.645 \pm 0.005	0.615 \pm 0.043	0.433 \pm 0.006	0.470 \pm 0.006
	Co-Reg	0.648 \pm 0.002	0.564 \pm 0.000	0.436 \pm 0.002	0.466 \pm 0.000
	RMSC	0.684 \pm 0.033	0.642 \pm 0.036	0.485 \pm 0.046	0.517 \pm 0.043
	DSCN	0.738 \pm 0.006	0.727 \pm 0.014	0.509 \pm 0.021	0.542 \pm 0.019
	DCSC	0.744 \pm 0.009	0.733 \pm 0.007	0.521 \pm 0.011	0.556 \pm 0.012
	DC	0.756 \pm 0.001	0.766 \pm 0.007	0.553 \pm 0.017	0.579 \pm 0.004
	LMSC	0.702 \pm 0.013	0.670 \pm 0.012	0.472 \pm 0.018	0.506 \pm 0.010
	DMF	0.782 \pm 0.010	0.745 \pm 0.011	0.579 \pm 0.002	0.601 \pm 0.002
	MSCN	0.769 \pm 0.003	0.772 \pm 0.004	0.582 \pm 0.012	0.598 \pm 0.006
	MvDSCN	0.797 \pm 0.007	0.824 \pm 0.004	0.626 \pm 0.011	0.650 \pm 0.010
ORL	BestSV	0.903 \pm 0.016	0.777 \pm 0.033	0.738 \pm 0.001	0.711 \pm 0.043
	LRR	0.895 \pm 0.006	0.773 \pm 0.003	0.724 \pm 0.002	0.731 \pm 0.004
	Min-Disagreement	0.816 \pm 0.001	0.734 \pm 0.040	0.621 \pm 0.003	0.663 \pm 0.003
	Co-Reg	0.853 \pm 0.003	0.715 \pm 0.000	0.602 \pm 0.004	0.615 \pm 0.000
	RMSC	0.872 \pm 0.012	0.723 \pm 0.025	0.645 \pm 0.029	0.654 \pm 0.028
	DSCN	0.883 \pm 0.005	0.801 \pm 0.009	0.704 \pm 0.012	0.711 \pm 0.011
	DCSC	0.893 \pm 0.003	0.811 \pm 0.003	0.709 \pm 0.021	0.718 \pm 0.004
	DC	0.865 \pm 0.011	0.788 \pm 0.002	0.684 \pm 0.007	0.701 \pm 0.008
	LMSC	0.931 \pm 0.011	0.819 \pm 0.017	0.769 \pm 0.044	0.758 \pm 0.009
	DMF	0.933 \pm 0.010	0.823 \pm 0.021	0.783 \pm 0.001	0.773 \pm 0.002
	MSCN	0.928 \pm 0.001	0.833 \pm 0.008	0.790 \pm 0.005	0.787 \pm 0.001
	MvDSCN	0.943 \pm 0.002	0.870 \pm 0.006	0.819 \pm 0.001	0.834 \pm 0.012
Still DB	BestSV	0.104 \pm 0.078	0.297 \pm 0.089	0.063 \pm 0.001	0.221 \pm 0.064
	LRR	0.109 \pm 0.030	0.306 \pm 0.039	0.066 \pm 0.002	0.240 \pm 0.052
	Min-Disagreement	0.097 \pm 0.005	0.336 \pm 0.014	0.103 \pm 0.013	0.223 \pm 0.004
	Co-Reg	0.093 \pm 0.016	0.263 \pm 0.024	0.092 \pm 0.004	0.226 \pm 0.035
	RMSC	0.106 \pm 0.056	0.285 \pm 0.020	0.113 \pm 0.063	0.232 \pm 0.021
	DSCN	0.216 \pm 0.011	0.323 \pm 0.006	0.145 \pm 0.002	0.293 \pm 0.019
	DCSC	0.222 \pm 0.008	0.325 \pm 0.007	0.148 \pm 0.003	0.301 \pm 0.002
	DC	0.199 \pm 0.003	0.315 \pm 0.001	0.131 \pm 0.001	0.280 \pm 0.011
	LMSC	0.137 \pm 0.032	0.328 \pm 0.029	0.088 \pm 0.007	0.269 \pm 0.055
	DMF	0.154 \pm 0.010	0.336 \pm 0.017	0.124 \pm 0.001	0.265 \pm 0.005
	MSCN	0.168 \pm 0.001	0.312 \pm 0.008	0.133 \pm 0.005	0.261 \pm 0.001
	MvDSCN	0.245 \pm 0.020	0.377 \pm 0.023	0.169 \pm 0.003	0.320 \pm 0.015
BBCSport	BestSV	0.715 \pm 0.060	0.836 \pm 0.037	0.659 \pm 0.005	0.768 \pm 0.038
	LRR	0.690 \pm 0.019	0.832 \pm 0.026	0.667 \pm 0.008	0.774 \pm 0.023
	Min-Disagreement	0.776 \pm 0.019	0.797 \pm 0.049	0.783 \pm 0.034	0.260 \pm 0.013
	Co-Reg	0.718 \pm 0.003	0.564 \pm 0.000	0.696 \pm 0.001	0.766 \pm 0.002
	RMSC	0.608 \pm 0.007	0.737 \pm 0.003	0.723 \pm 0.025	0.655 \pm 0.002
	DSCN	0.652 \pm 0.000	0.821 \pm 0.000	0.856 \pm 0.001	0.683 \pm 0.001
	DCSC	0.683 \pm 0.001	0.843 \pm 0.000	0.864 \pm 0.012	0.712 \pm 0.002
	DC	0.556 \pm 0.001	0.724 \pm 0.000	0.781 \pm 0.000	0.492 \pm 0.000
	LMSC	0.826 \pm 0.006	0.900 \pm 0.044	0.893 \pm 0.012	0.887 \pm 0.071
	DMF	0.821 \pm 0.003	0.890 \pm 0.031	0.883 \pm 0.012	0.889 \pm 0.001
	MSCN	0.813 \pm 0.002	0.888 \pm 0.003	0.859 \pm 0.001	0.854 \pm 0.002
	MvDSCN	0.835 \pm 0.000	0.931 \pm 0.001	0.909 \pm 0.001	0.860 \pm 0.000

A. Experiment Setup

Datasets. We extensively evaluate the multi-view clustering performance of the proposed model on benchmark multi-view datasets.

- **Yale** is a widely used face dataset which contains 165 gray scale images, which are composed of 15 individuals with 11 images per person. Variations of the data include center light, with glasses, happy, left light, without glasses, normal, right light, sad, sleepy, surprised and wink.
- **ORL** contains 10 different images of each of 40 distinct subjects. For each subject, the images were taken under varying lighting conditions with different facial expressions (open / closed eyes, smiling / not smiling) and facial details (glasses / no glasses). For the face dataset (Yale and ORL), we adjust the image size to 48×48 and extract three types of features, i.e., intensity

(4,096 dimensions), LBP (3,304 dimensions) and Gabor (6,750 dimensions). The standard LBP features are then extracted from the 72×80 loosely cropped image with a histogram size of 59 over 910 pixels. The Gabor feature is dominated by four directions $\theta = 0^\circ, 45^\circ, 90^\circ, 135^\circ$ and extracted at a scale of $\lambda = 4$. It has a resolution of 25×30 pixels and a loose face cropping. Note that all descriptors except intensity are scaled to have a unit norm.

- **Still DB** consists of 467 images with 6 classes of actions. sift bow, color sift bow and shape context bow are extracted.
- **BBCSport** contains 544 documents from the BBC Sport website of sports news articles, which are related to two viewpoints in five topical areas during 2004-2005. For each sample, there are 3,183 features for the first view and 3,203 features for the second view.

Beyond multi-feature subspace clustering, MvDSCN can be easily extended to multi-modal learning by replacing the input with data with different modalities. We evaluate the proposed deep multi-view subspace clustering methods on real-world RGB-D Object Dataset [44]. It contains visual and depth images of 300 physically distinct objects taken from multiple views and the objects are organized into 51 categories arranged by WordNet hypernym-hyponym relationships (similar to ImageNet). Our experiment datasets is composed of 50 categories randomly selected from RGB-D Object Dataset with each class containing 10 examples. All visual images and depth images are resized to 64×64 pixels. We apply median filter recursively until all missing values are filled to visualize. A subset of RGB and depth images are shown in Figure 4.

Comparison methods. The performance of MvDSCN is compared with the state-of-the-art subspace clustering methods in term of four evaluation metrics. There are six shallow and deep single-view clustering algorithms, and five multi-view clustering algorithms. For all single-view clustering algorithms, the performance of the best view is reported.

- **BestSV** reports the result of the individual view which achieves the best spectral clustering performance with a single view of data [45].
- **LRR** seeks the lowest-rank representation among all the candidates that can represent the data samples as linear combinations of the bases in a given dictionary with the best single view [1].
- **RMSC** has a flavor of low-rank and sparse decomposition [46]. It firstly construct a transition probability matrix from each single view, and then uses these matrices to recover a shared low-rank transition probability matrix as a crucial input to the standard Markov chain method for clustering.
- **DSCN** is a deep auto-encoder framework for subspace clustering with a best single view. [34]
- **DCSC** imposes a self-paced regularizer on the loss and presents a robust deep subspace clustering algorithm [7].
- **DC** proposes a scalable clustering approach for the unsupervised learning of convnets. It iterates between clustering with k-means and updating its weights by predicting the cluster assignments as pseudo-labels in a discriminative loss [47].
- **Min-Disagreement** creates a bipartite graph and is based on the minimizing-disagreement idea [48].
- **Co-Reg SPC** uses spectral clustering objective functions that implicitly combine graphs from multiple views of the data to achieve a better clustering result [49].
- **DMF** learns the hierarchical semantics of multi-view data through the semi-nonnegative matrix factorization [17].
- **LMSC** seeks the underlying latent representation and simultaneously performs data reconstruction based on the learned latent representation [21].
- **MSCN** observes that spatial fusion methods in a deep multimodal subspace clustering task relay on spatial correspondences among the modalities [50].

Evaluation Metrics. Following the experiment setting in [17], [21], four popular metrics are used to evaluate the cluster-

ing quality, including **NMI** (Normalized Mutual Information), **ACC** (Accuracy), **F-Measure**, and **AR** (Adjusted Rand Index) which can comprehensively evaluate the performance.

The NMI calculates the normalized measure of similarity between two labels of the same data as follows:

$$NMI = \frac{I(l; c)}{\max\{H(l), H(c)\}}, \quad (10)$$

where $I(l; c)$ denotes the mutual information between l and c , and H represents their entropy. Result of NMI do not change by permutations of clusters (classes), and they are normalized to the range of $[0, 1]$, with 0 meaning no correlation and 1 exhibiting perfect correlation.

The ACC score is calculated as:

$$ACC = \max_m \frac{\sum_{i=1}^n \{l_i = m(c_i)\}}{n}, \quad (11)$$

where l_i is the ground-truth label, c_i is the cluster assignment produced by the model. $m(c_i)$ is the permutation map function, which maps the cluster labels into class labels. n is the number of samples. The best map can be obtained by the Kuhn-Munkres algorithm.

The F-measure can be interpreted as a weighted average of the precision and recall, where an F-measure reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F-measure are equal. Its formulation is:

$$F_{measure} = \left(\frac{\text{recall}^{-1} + \text{precision}^{-1}}{2} \right)^{-1} \quad (12)$$

The adjusted rand index is the corrected-for-chance version of the rand index [51].

Note that lower values indicate better performance for average entropy, and higher values indicate better performance for the other metrics. We optimize all the parameters to achieve the best performance of the comparison method. Especially, we run each method 30 times and report the average performance and standard deviation.

B. Results of Multi-feature Subspace Clustering

The multi-view clustering performance of different methods is given in Table I. Our proposed method significantly outperforms other methods on Yale, ORL and Still DB datasets, and shows very competitive performance on BBCSport dataset. For Yale, we raise the performance bar by around 7.9% in ACC, 4.7% in AR, 4.9% in F-measure. In addition, for ORL, we raise the performance bar by around 4.7% in ACC, 3.6% in AR, 6.1% in F-measure. On StillDB our method gains significant improvements around 9.1%, 3.7%, 4.5%, 5.5%, over the second best method in terms of NMI, ACC, AR, F-measure, respectively. For BBCSport, the performance of the proposed method is better than DMF in terms of three evaluation metrics. This demonstrates the effectiveness of the proposed MvDSCN on multi-feature subspace clustering task. The performance improvement owns to two aspects, i.e., the end to end manner in learning the affinity matrix, and multi-view relations embedded into both feature learning and self-representation.

TABLE II: Results on four multi-feature datasets (mean \pm standard deviation). Higher value indicates better performance.

Datasets	Views	NMI	ACC	AR	F-measure
Yale	View1	0.738 \pm 0.006	0.727 \pm 0.014	0.509 \pm 0.021	0.542 \pm 0.019
	View2	0.613 \pm 0.007	0.598 \pm 0.008	0.401 \pm 0.008	0.439 \pm 0.008
	View3	0.545 \pm 0.009	0.522 \pm 0.012	0.267 \pm 0.011	0.311 \pm 0.011
	ALL	0.797 \pm 0.007	0.824 \pm 0.004	0.626 \pm 0.011	0.650 \pm 0.010
ORL	View1	0.883 \pm 0.005	0.801 \pm 0.009	0.704 \pm 0.012	0.711 \pm 0.011
	View2	0.793 \pm 0.011	0.627 \pm 0.024	0.504 \pm 0.023	0.516 \pm 0.023
	View3	0.764 \pm 0.009	0.580 \pm 0.024	0.458 \pm 0.024	0.471 \pm 0.023
	ALL	0.943 \pm 0.002	0.870 \pm 0.006	0.819 \pm 0.001	0.834 \pm 0.012
Still DB	View1	0.113 \pm 0.001	0.329 \pm 0.004	0.083 \pm 0.003	0.243 \pm 0.014
	View2	0.216 \pm 0.011	0.323 \pm 0.006	0.145 \pm 0.002	0.293 \pm 0.019
	View3	0.211 \pm 0.002	0.313 \pm 0.012	0.142 \pm 0.014	0.289 \pm 0.007
	ALL	0.245 \pm 0.020	0.377 \pm 0.023	0.169 \pm 0.003	0.320 \pm 0.015
BBCSport	View1	0.617 \pm 0.000	0.801 \pm 0.000	0.847 \pm 0.001	0.559 \pm 0.000
	View2	0.652 \pm 0.000	0.821 \pm 0.000	0.856 \pm 0.001	0.683 \pm 0.001
	ALL	0.835 \pm 0.000	0.931 \pm 0.001	0.909 \pm 0.001	0.860 \pm 0.000

TABLE III: Results on RGB-D Object datasets (mean \pm standard deviation). Higher value indicates better performance.

Datasets	Methods	NMI	ACC	AR	F-measure
RGB-D Object	BestSV	0.554 \pm 0.006	0.278 \pm 0.001	0.106 \pm 0.006	0.125 \pm 0.006
	LRR	0.589 \pm 0.002	0.299 \pm 0.010	0.143 \pm 0.002	0.156 \pm 0.001
	Min-Disagreement	0.605 \pm 0.008	0.332 \pm 0.002	0.160 \pm 0.013	0.177 \pm 0.011
	Co-Reg	0.602 \pm 0.007	0.268 \pm 0.003	0.155 \pm 0.020	0.175 \pm 0.018
	RMSC	0.603 \pm 0.006	0.341 \pm 0.015	0.162 \pm 0.010	0.178 \pm 0.010
	DSCN	0.589 \pm 0.004	0.339 \pm 0.006	0.163 \pm 0.004	0.179 \pm 0.004
	DCSC	0.591 \pm 0.002	0.340 \pm 0.002	0.170 \pm 0.001	0.182 \pm 0.003
	DC	0.594 \pm 0.003	0.340 \pm 0.002	0.177 \pm 0.004	0.184 \pm 0.004
	LMSC	0.593 \pm 0.030	0.335 \pm 0.028	0.151 \pm 0.035	0.167 \pm 0.034
	DMF	0.549 \pm 0.004	0.286 \pm 0.006	0.107 \pm 0.002	0.123 \pm 0.001
	MSCN	0.608 \pm 0.001	0.354 \pm 0.003	0.190 \pm 0.002	0.203 \pm 0.004
MvDSCN	0.639 \pm 0.003	0.388 \pm 0.005	0.210 \pm 0.004	0.225 \pm 0.004	

TABLE IV: Results on RGB-D Object dataset (mean \pm standard deviation). Higher value indicates better performance.

Datasets	Views	NMI	ACC	AR	F-measure
RGB-D Object	RGB	0.589 \pm 0.004	0.339 \pm 0.006	0.163 \pm 0.004	0.179 \pm 0.004
	Depth	0.576 \pm 0.004	0.300 \pm 0.005	0.131 \pm 0.004	0.147 \pm 0.004
	RGB+Depth	0.639 \pm 0.003	0.388 \pm 0.005	0.210 \pm 0.004	0.225 \pm 0.004

TABLE V: Ablation study on RGB-D Object dataset (mean \pm standard deviation). Higher value indicates better performance.

Methods	NMI	ACC	AR	F-measure
D-MvDSCN	0.594 \pm 0.004	0.343 \pm 0.007	0.190 \pm 0.006	0.205 \pm 0.006
U-MvDSCN	0.593 \pm 0.005	0.350 \pm 0.006	0.192 \pm 0.006	0.197 \pm 0.006
MvDSCN	0.639 \pm 0.003	0.388 \pm 0.005	0.210 \pm 0.004	0.225 \pm 0.004

Single View versus Multiple Views. To further investigate the improvement of our method, we compare our method with deep subspace clustering networks (DSCN) [2] with only single-view data. Figure 6 shows the detailed results on different datasets. According to Table II, the clustering performance with multiple views consistently outperforms that of each single view, which empirically proves that clustering with multiple views is more robust than that with single view.

C. Results of Multi-modal Subspace Clustering

For multi-modal experiments, we use the pre-trained deep auto-encoders to extract features for the comparison shallow methods. There are 4,096 features for both RGB image and depth image. The experimental results on the RGB-D dataset are presented in Table III. Our method outperforms all the other competitors. We raise the performance bar by around 3.4% in NMI, 5.3% in ACC, 4.8% in AR, 5.8% in F-measure

compared with the second best method. Compared with the shallow models that first extract deep features and then conduct subspace clustering, our proposed MvDSCN jointly feature learning and subspace clustering together and multi-view relations affect on both parts. Hence, MvDSCN outperforms the state-of-the-art subspace clustering algorithms.

Single Modal versus Multiple Modal. As shown in Table IV, our methods significantly outperform subspace clustering with only single modal data, which further demonstrates the superiority of multi-modal fusion. Overall, the RGB modality achieves better performance than Depth modality. We improve the clustering performance by 5.0% in NMI, 4.9% in ACC, 4.7% in AR and 4.6% in F-measure when we fuse them by MvDSCN.

D. Convergence Analysis

To empirically analyze the convergence of MvDSCN, in Figure 7, we show the relationship between the loss of the

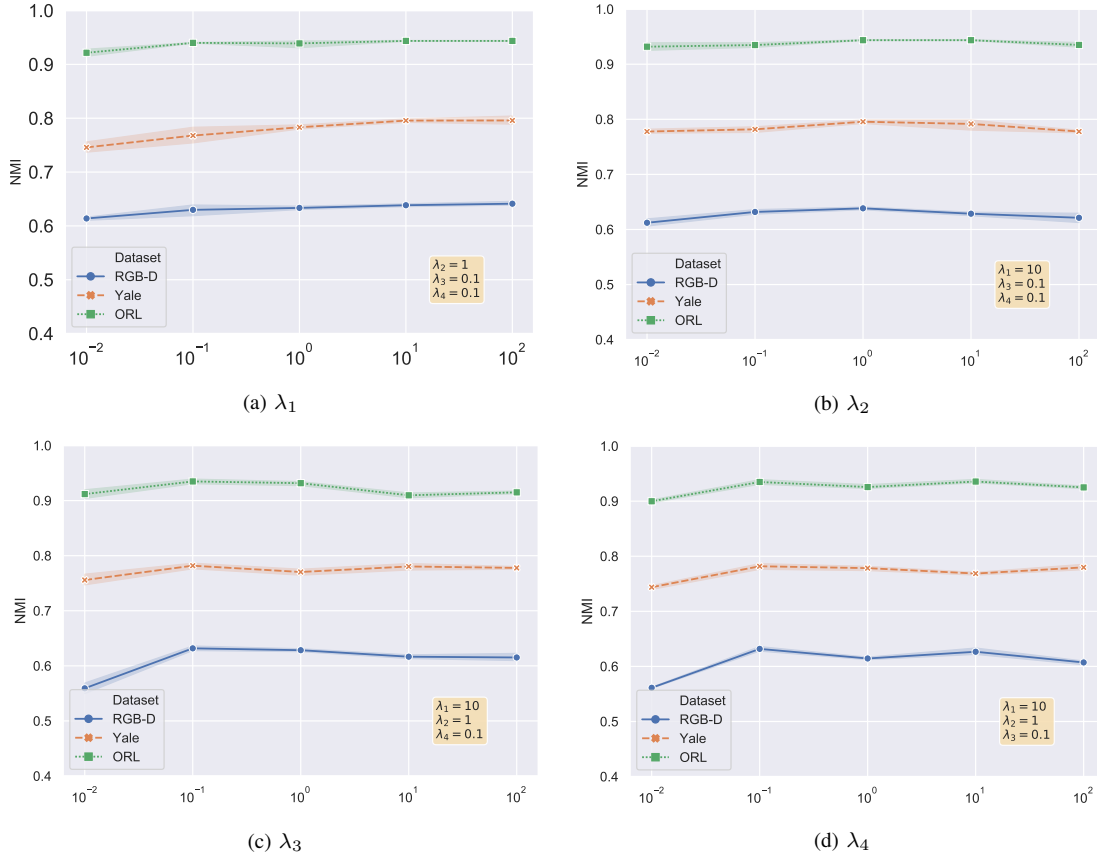


Fig. 5: The effect of different parameters on MvDSCN learning.

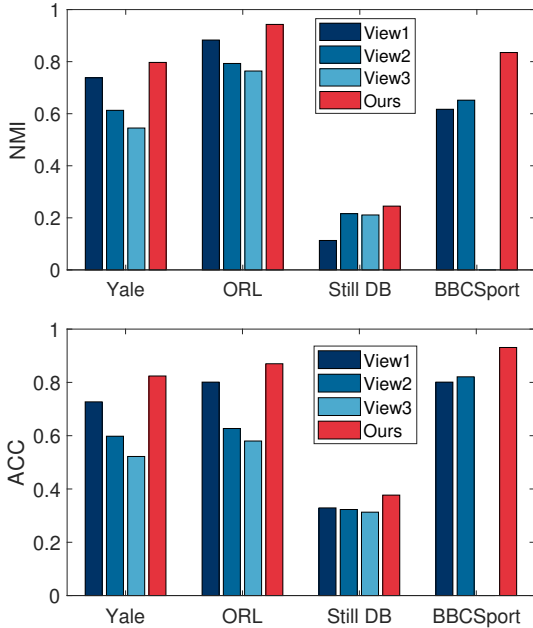


Fig. 6: Performance comparison between DSCN [2] with each view and MvDSCN versus NMI and ACC.

MvDSCN and the clustering performance on the ORL dataset. The reported values in this figure are normalized between zero and one. As can be seen from the figure, the loss decreases rapidly in a few epochs. The clustering performance increases significantly in the first few epochs and then grows slowly. Similar results can be observed on other datasets.

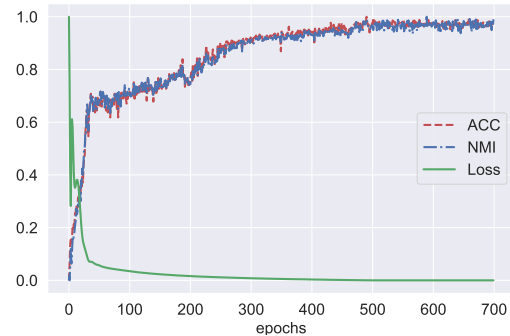


Fig. 7: The loss, and clustering performance (NMI and ACC) of MvDSCN with training epochs.

E. Parameter Sensitivity

To analyze the impact the parameters on the clustering performance of MvDSCN, we plot the performance of MvD-

SCN with different parameters in Figure 5. There are four parameters, i.e., λ_1 , λ_2 , λ_3 , and λ_4 . λ_1 seeks the balance between self-representation loss and reconstruction loss of auto-encoders. λ_2 reflects the impact of the l_p -norm regularization. λ_3 controls the degree of the universality regularizer, while λ_4 controls the degree of the diversity regularizer. We fix the other three parameters and analyze the impact of the rest parameter. The results show that the clustering performance grows with the value of λ_1 and varies little when λ_1 is above 10. For λ_1 , the best performance is achieved across different datasets when λ_1 is set as 1. Similarly, we can observe that when λ_3 and λ_4 are set as 0.1, our proposed MvDSCN achieves superior performance. For all datasets, we fix the values of four parameters as 10, 1, 0.1, 0.1.

F. Ablation Study

To verify the effectiveness of diversity and universality regularization, we conduct the ablation study with respect to the proposed model. D-MvDSCN represents the proposed model without diversity regularization while U-MvDSCN refers to the one without universality regularization. Note that U-MvDSCN can be considered as the case when only the Unet part is kept. As presented in Table V, MvDSCN substantially outperforms D-MvDSCN, which numerically indicates that we cannot ignore diversity regularization that can enhance multi-view complementary information. Besides, our proposed method performs better than U-MvDSCN. Universality regularization forces the view-specific representation to be centralized to the common representation. We conduct a sensitivity test for the regularizer parameter of diversity (λ_3) and universality (λ_4) by varying from 0.001 to 1. Figure 8 shows the influence of different parameter values with respect to NMI on RGB-D dataset. Moreover, our method performs much stable when λ_3 and λ_4 becomes larger. In sum, both diversity and universality regularization contribute to the enhancement of the proposed model in terms of multi-view clustering performance.

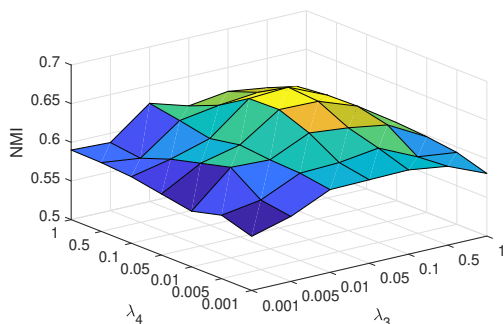


Fig. 8: Sensitivity test on λ_3 and λ_4 versus NMI on RGB-D.

V. CONCLUSIONS

In this paper, we proposed a new multi-view deep subspace clustering network (MvDSCN). The proposed method learns multi-view self-representation in an end-to-end manner by combining convolutional auto-encoder and self-representation

together. It consists of diversity net (Dnet) and universality net (Unet), which are connected by diversity and universality regularizer. Experiments on both multi-feature and multi-modal tasks validate the superiority of our method compared with the state-of-the-arts.

REFERENCES

- [1] L. Guangcan, L. Zhouchen, Y. Shuicheng, S. Ju, Y. Yong, and M. Yi, "Robust recovery of subspace structures by low-rank representation," *TPAMI*, vol. 35, no. 1, pp. 171–184, 2013.
- [2] P. Ji, T. Zhang, H. Li, M. Salzmann, and I. Reid, "Deep subspace clustering networks," in *NIPS*, 2017, pp. 24–33.
- [3] P. Zhou, Y. Hou, and J. Feng, "Deep adversarial subspace clustering," in *CVPR*, 2018.
- [4] C. Lu, J. Feng, Z. Lin, M. Tao, and S. Yan, "Subspace clustering by block diagonal representation," *TPAMI*, vol. PP, no. 99, pp. 1–1, 2018.
- [5] D. Chen, J. Lv, and Y. Zhang, "Unsupervised multi-manifold clustering by learning deep representation," in *AAAI Workshops*, 2017.
- [6] K. Tian, S. Zhou, and J. Guan, "Deepcluster: A general clustering framework based on deep learning," in *ECML/PKDD*, 2017.
- [7] Y. Jiang, Z. Yang, Q. Xu, X. Cao, and Q. Huang, "When to learn what: Deep cognitive subspace clustering," in *ACMMM*, 2018.
- [8] X. Peng, J. Feng, S. Xiao, W.-Y. Yau, J. T. Zhou, and S. Yang, "Structured autoencoders for subspace clustering," *TIP*, vol. 27, pp. 5076–5086, 2018.
- [9] X. Guo, L. Gao, X. Liu, and J. Yin, "Improved deep embedded clustering with local structure preservation," in *IJCAI*, 2017.
- [10] J. Li, H. Liu, H. Zhao, and Y. Fu, "Projective low-rank subspace clustering via learning deep encoder," in *IJCAI*, 2017.
- [11] R. Chellappa, "Deep density clustering of unconstrained faces," in *CVPR*, 2018.
- [12] J. Liang, J. Yang, H.-Y. Lee, K. Wang, and M.-H. Yang, "Sub-gan: An unsupervised generative model via subspaces," in *ECCV*, 2018.
- [13] J. Lezama, Q. Qiu, P. Musé, and G. Sapiro, "Ole: Orthogonal low-rank embedding, a plug and play geometric loss for deep learning," in *CVPR*, 2018.
- [14] M. Luo, X. Chang, Z. Li, L. Nie, A. G. Hauptmann, and Q. Zheng, "Simple to complex cross-modal learning to rank," *Computer Vision and Image Understanding*, vol. 163, pp. 67–77, 2017.
- [15] S. Kim, D. Min, B. Ham, S. Ryu, M. N. Do, and K. Sohn, "Dasc: Dense adaptive self-correlation descriptor for multi-modal and multi-spectral correspondence," in *ICCV*, 2015, pp. 2103–2112.
- [16] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," in *CVPR*, 2017, pp. 1907–1915.
- [17] H. Zhao, Z. Ding, and Y. Fu, "Multi-view clustering via deep matrix factorization," in *AAAI*, 2017.
- [18] Y. Cui, X. Z. Fern, and J. G. Dy, "Non-redundant multi-view clustering via orthogonalization," in *ICDM*, 2007.
- [19] S. Günnemann, I. Färber, and T. Seidl, "Multi-view clustering using mixture models in subspace projections," in *KDD*, 2012.
- [20] Y. Wang, X. Lin, L. Wu, W. Zhang, Q. Zhang, and X. Huang, "Robust subspace clustering for multi-view data by exploiting correlation consensus," *TIP*, vol. 24, pp. 3939–3949, 2015.
- [21] C. Zhang, Q. Hu, H. Fu, P. Zhu, and X. Cao, "Latent multi-view subspace clustering," in *CVPR*, 2017.
- [22] J. Xu, J. Han, and F. Nie, "Discriminatively embedded k-means for multi-view clustering," in *CVPR*, 2016.
- [23] Q. Wang, Z. Ding, Z. Tao, Q. Gao, and Y. Fu, "Partial multi-view clustering via consistent gan," in *ICDM*, 2018.
- [24] E. Elhamifar and R. Vidal, "Sparse subspace clustering," in *CVPR*, 2009.
- [25] C. Zhang, H. Fu, Q. Hu, X. Cao, Y. Xie, D. Tao, and D. Xu, "Generalized latent multi-view subspace clustering," *TPAMI*, vol. PP, no. 99, pp. 1–1.
- [26] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2. IEEE, 2005, pp. 60–65.
- [27] P. Zhu, W. Zuo, L. Zhang, Q. Hu, and S. C. Shiu, "Unsupervised feature selection by regularized self-representation," *Pattern Recognition*, vol. 48, no. 2, pp. 438–446, 2015.
- [28] P. Zhu, Q. Hu, C. Zhang, and W. Zuo, "Coupled dictionary learning for unsupervised feature selection," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

- [29] P. Zhu, W. Zhu, Q. Hu, C. Zhang, and W. Zuo, "Subspace clustering guided unsupervised feature selection," *Pattern Recognition*, vol. 66, pp. 364–374, 2017.
- [30] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803.
- [31] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313 5786, pp. 504–7, 2006.
- [32] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *Journal of Machine Learning Research*, vol. 11, pp. 3371–3408, 2010.
- [33] J. Xie, R. B. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," in *ICML*, 2016.
- [34] X. Peng, S. Xiao, J. Feng, W.-Y. Yau, and Z. Yi, "Deep subspace clustering with sparsity prior," in *IJCAI*, 2016.
- [35] F. Li, H. Qiao, and B. Zhang, "Discriminatively boosted image clustering with fully convolutional auto-encoders," *PR*, vol. 83, pp. 161–173, 2018.
- [36] J. Masci, U. Meier, D. C. Ciresan, and J. Schmidhuber, "Stacked convolutional auto-encoders for hierarchical feature extraction," in *ICANN*, 2011.
- [37] A. T. L. Nguyen, J. Yosinski, Y. Bengio, A. Dosovitskiy, and J. Clune, "Plug & play generative networks: Conditional iterative generation of images in latent space," *CVPR*, pp. 3510–3520.
- [38] Y. Zhang, D. Shen, G. Wang, Z. Gan, R. Henaio, and L. Carin, "Deconvolutional paragraph representation learning," in *NIPS*, 2017.
- [39] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010.
- [40] A. Gretton, O. Bousquet, A. J. Smola, and B. Scholkopf, "Measuring statistical dependence with hilbert-schmidt norms," *Algorithmic Learning Theory*, pp. 63–77, 2005.
- [41] X. Cao, C. Zhang, H. Fu, S. Liu, and H. Zhang, "Diversity-induced multi-view subspace clustering," *CVPR*, 2015.
- [42] C. Lu, J. Feng, Z. Lin, T. Mei, and S. Yan, "Subspace clustering by block diagonal representation," *TPAMI*, vol. 41, no. 2, pp. 487–501, 2019.
- [43] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2015.
- [44] K. Lai, L. Bo, X. Ren, and D. Fox, "A large-scale hierarchical multi-view rgb-d object dataset," *ICRA*, pp. 1817–1824, 2011.
- [45] A. Y. Ng, M. I. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *NIPS*, 2001.
- [46] R. Xia, Y. Pan, L. Du, and J. Yin, "Robust multi-view spectral clustering via low-rank and sparse decomposition," in *AAAI*, 2014.
- [47] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, 2018.
- [48] V. R. de Sa, "Spectral clustering with two views," in *ICML*, 2005.
- [49] A. Kumar, P. Rai, and H. Daumé, "Co-regularized multi-view spectral clustering," in *NIPS*, 2011.
- [50] M. Abavisani and V. M. Patel, "Deep multimodal subspace clustering networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, pp. 1601–1614, 2018.
- [51] W. M. Rand, "Objective criteria for the evaluation of clustering methods," *Journal of the American Statistical association*, vol. 66, no. 336, pp. 846–850, 1971.