

# Attentive Multi-View Deep Subspace Clustering Net<sup>★</sup>

Run-kun Lu<sup>a</sup>, Jian-wei Liu<sup>a</sup>, Xin Zuo<sup>a</sup>

<sup>a</sup>Department of Automation, College of Information Science and Engineering, China University of Petroleum, Beijing, 260 Mailbox, Changping District, Beijing 102249, China

## Abstract

In this paper, we propose a novel Attentive Multi-View Deep Subspace Nets (AMVDSN), which deeply explores underlying consistent and view-specific information from multiple views and fuse them by considering each view's dynamic contribution obtained by attention mechanism. Unlike most multi-view subspace learning methods that they directly reconstruct data points on raw data or only consider consistency or complementarity when learning representation in deep or shallow space, our proposed method seeks to find a joint latent representation that explicitly considers both consensus and view-specific information among multiple views, and then performs subspace clustering on learned joint latent representation. Besides, different views contribute differently to representation learning, we therefore introduce attention mechanism to derive dynamic weight for each view, which performs much better than previous fusion methods in the field of multi-view subspace clustering. The proposed algorithm is intuitive and can be easily optimized just by using Stochastic Gradient Descent (SGD) because of the neural network framework, which also provides strong non-linear characterization capability compared with traditional subspace clustering approaches. The experimental results on seven real-world data sets have demonstrated the effectiveness of our proposed algorithm against some state-of-the-art subspace learning approaches.

**Keywords:** multi-view learning, subspace clustering, deep learning, attention

## 1. Introduction

In recent years, multi-view learning has attracted widespread attention from machine learning researchers because it usually considers more complete information than single-view algorithms. Generally, the field is consisting of multi-modal and multi-feature learning, where multi-modal emphasizes that data come from multiple modalities (such as video, audio, and text); multi-feature emphasizes that data obtained by different feature extraction methods. Most existing so called multi-view algorithms are belong to multi-feature learning [1], and this paper is also under the assumption of multi-feature data. On the other hand, in this field, multi-view self-representation subspace learning is one of the most attractive directions, where most existing multi-view subspace clustering algorithms are actually self-representation based methods [2, 3, 4, 5, 6, 7].

Self-representation assumes that each instance in input space can be expressed by linear weighted combination of other instances, and given the instances set  $\mathbf{X} \in \mathbb{R}^{M \times N}$ , the process can be formulated as:

$$\begin{aligned} \min_{\mathbf{C}} R(\mathbf{C}) \\ \text{s.t. } \mathbf{X} = \mathbf{X}\mathbf{C}, \mathbf{C} \in \Omega, \end{aligned} \quad (1)$$

where  $\Omega$  is some matrix set, and  $R(\cdot)$  is the regularization of self-representation matrix which can implement different constraints on optimal problem (usually frobenius,  $L_1$ , nuclear

norms are alternatives). Besides,  $\mathbf{C} \in \mathbb{R}^{N \times N}$  is the self-representation matrix, where  $[\mathbf{C}]_{i,j}$  measures the correlation between instance  $i$  and  $j$ , and then the self-representation matrix can be utilized to construct affinity matrix of spectral clustering according to

$$\frac{|\mathbf{C}| + |\mathbf{C}^T|}{2}.$$

Similar to other multi-view learning algorithms [8, 9], multi-view subspace clustering also focus on multiple views' consistency and complementarity, where consistency represents views' consistent subspace which can be achieved by low-rank regularization [4] or views alignment [5], and complementarity usually focus on the view-specific subspace which can be achieved by frobenius regularization [6].

Subspace clustering works well under the assumption that sample space has relatively separable decision boundaries, when facing highly nonlinear data, only considering the regularization of self-representation matrix can not effectively improve the clustering performance due to the limitation of input representation. Based on this, latent multi-view subspace clustering [4] first computed a latent representation of multiple views, and then applied it in subspace clustering; more recently, multi-view low-rank sparse subspace clustering [5] has proposed kernel-based variants that first project input into high dimensional kernel space. However, the first approach essentially obtained the latent feature through a linear transformation, and the second one encountered the problem of kernel functions selection. To mitigate these problems, neural network is considered more suitable for its great non-linear and

<sup>★</sup>This work was supported by

Email addresses: zsy1rk@gmail.com (Run-kun Lu), liujw@cup.edu.cn (Jian-wei Liu), zuox@cup.edu.cn (Xin Zuo)

adaptable (without the kernel functions selection) ability, such as Deep Subspace Clustering Networks [10] and Deep Multimodal Subspace Clustering Networks (DMSCN) [11], they added additional self-representation layer in fully convolutional auto-encoders to achieve deep subspace learning. However, they are available only for 3-D tensor (we view gray-scale image as one channel but still 3-D tensor, such as for MNIST the dimension is [28, 28, 1]), but for most multi-view data, each view is always a kind of 1-D hand-craft feature of raw data that the fully convolutional network is not suitable for this case; and DMSCN is a multi-modal model but not considers different modal’s contribution in modal fusion process.

In multi-view learning task, different data uses various methods when constructing multiple views, but there is no universal criterion to measure the quality of view representation. Compared to relying on experience, an end-to-end approach is necessary to measure the degree of view importance through learning. Fortunately, with the great progress made by transformer [12] in NLP [13] and CV [14] tasks, attention mechanism has become an essential module for many of the recent neural architectures, which allows the model to dynamically pay attention to only certain parts of the input. Motivated by this, we propose a novel framework called Attentive Multi-View Deep Subspace Net (AMVDSN), which is a fully connected layer based network. Our proposed method takes both multiple views’ consistent and view-specific information into consideration and fuses them with dynamic weights into a joint latent representation through attention, and by utilizing this mechanism, we can effectively exploit the contribution of each view. Based on this, just a simple regularization like frobenius norm of self-representation matrix, which can be optimized automatically in deep learning programming framework, will derive much better clustering performance than traditional self-representation subspace clustering algorithms. On the other hand, AMVDSN is a deep model constructed mainly by fully connected layers, which is easy to encounter the problem of model degradation. Therefore, by introducing shortcut connection [15], we can effectively alleviate model degradation and improve clustering performance.

In summary, the main contributions of this paper are:

- (1) the proposed algorithm is the first approach to apply self-attention mechanism in multi-view subspace clustering, which incorporates multiple views’ dynamic contributions to representation learning and derive more reasonable non-linear joint latent representation consists of consistent and view-specific information;
- (2) compare with traditional subspace learning methods that have complicated regularization on self-representation matrix and complex optimization process, we just select simple frobenius norm as regularization term that can be optimized automatically using SGD and derive much better clustering performance;
- (3) the proposed algorithm is the first multi-view subspace clustering method that introduces shortcut connection to solve the problem of model degradation.
- (4) the proposed method has achieved state-of-the-art clustering performance on some real-world data, especially on

Prokaryotic, which improved by more than 10 percentage points on multiple metrics.

Our paper is organized as follows. In section 2, we introduce some related works on multi-view learning, subspace clustering, and multi-view subspace clustering; In section 3, we give a detailed description on our proposed algorithm; In section 4, we demonstrate the effectiveness of our proposed method through some experiments on real world-data sets; In section 5, we will conclude this paper.

## 2. Related Work

### 2.1. Multi-View Learning

Recently, multi-view learning has become an important research direction of machine learning, and the field usually focus on utilizing more complete information from multiple views compared with single view algorithms. Canonical correlation analysis (CCA) [16] and co-training [17] are always viewed as the early work of multi-view learning, and scholars have used them as a basis to develop many variants in this field. According to Sun’s book [18] (the first book on multi-view learning), the research topics in this area mainly include: multi-view supervised [19, 20] and semi-supervised [21, 22] learning, multi-view subspace learning [23, 24], multi-view clustering [25, 26, 27], multi-view active learning [28, 29], multi-view transfer [30] and multi-task learning [31, 9], multi-view deep learning [32], and view construction [33]. Besides, the field can be widely used in plenty of applications, such as computer vision [34], social network [35], recommendation Systems [36], and medical research [37].

### 2.2. Self-representation Subspace Clustering

In the past decade, subspace clustering has been widely used in many real-world applications [38, 39, 40], and particularly, self-representation based methodologies have achieved promising clustering results. Such methods belong to spectral clustering, which means they all need to learn a affinity matrix to measure the similarities among different instances, and one of the most effective approaches is self-representation clustering. According to equation (1), different regularization terms correspond to various subspace clustering algorithms, such as least squares regression (LSR) [41], sparse subspace clustering (SSC) [42], and low-rank representation (LRR) [43] use frobenius,  $L_1$ , and nuclear norms respectively. Further more, multi-subspace representation (MSR) [44] combines  $L_1$ , and nuclear norms together to utilize the advantages of SSC and LRR. The works discussed above all based on the assumption that self-representation matrix  $C$  follows block diagonal property, but actually most existing subspace clustering algorithms have indirect structure priors, and based on this Lu et al. has proposed block diagonal representation (BDR) to deal with this problem [45]. Although the above research works have achieved promising performance, they still suffered from the problem that each instance should be linear reconstructed, which is at odds with most complex real-world applications involve non-linear relationships. Based on this, some researchers have proposed kernel-based approaches [46, 47, 48, 49], but they still

faced the dilemma of kernel selection. Therefore, neural network has naturally become the better solution for its strong non-linear ability without kernel selection problem, some auto-encoder based subspace clustering methods have been proposed and achieved impressive performance [50, 51].

### 2.3. Multi-View Subspace Clustering

Subspace clustering has become one of the most attractive research directions of multi-view learning in recent years, and most existing algorithms are self-representation based subspace clustering. Similar to single view subspace clustering algorithms, the early works are usually matrix factorization based methods but have different constraints on self-representation matrix. However, as explained in section 1, the regularization terms usually correspond to the underlying physical meanings [5, 6] in multi-view learning, such as low-rank regularization focus on the consistency between views [4]. Comparing with single view approaches, multi-view subspace clustering needs to consider multiple views information and the views fusion methods usually follow two ways: (1) obtain self-representation matrix of all views and fusion them together as one affinity matrix of spectral clustering [2, 3, 5]; (2) first compute a latent representation of multiple views and then obtain its self-representation matrix [6, 4, 52]. Unfortunately, these algorithms also face the limitation of linear views reconstruction process, kernel based methods cannot deal with the issue well and neural network based method must be the better solution as discussed in subsection 2.2.

## 3. Attention-based Multi-View Deep Subspace Network

### 3.1. Notations

For a clearer and more accurate description, we denote that bold uppercase and lowercase characters stand for matrices and column vectors respectively, and other not bold characters represent scalars. Let  $[\mathbf{a}; \mathbf{b}]$  denotes the vertical concatenation between vectors  $\mathbf{a}$  and  $\mathbf{b}$ . And the network structure described in the remaining subsections is shown in Figure 1.

### 3.2. Framework Description

**Encoder embedding layer:** we define  $\mathbf{X}^v \in \mathbb{R}^{M^v \times N}$ ,  $v \in \{1, 2, \dots, V\}$  is the collection of instances for the  $v$ -th view, where  $M^v$  and  $N$  are instance dimension and number of  $v$ -th view respectively. Due to dimensional differences among views, we need first perform embedding operation to ensure that different views characterize as the same dimension, and the encoder of auto-encoder can accomplish this process:

$$\mathbf{h}_n^v = f(\mathbf{W}_{(L,v)} f(\mathbf{W}_{(L-1,v)} \dots f(\mathbf{W}_{(1,v)} [\mathbf{x}_n^v; 1]))) , \quad (2)$$

where  $\mathbf{x}_n^v \in \mathbb{R}^{M^v}$  is  $n$ -th instance of  $\mathbf{X}^v$ , i.e.  $\mathbf{X}^v = \{\mathbf{x}_n^v\}_{n=1}^N$ ;  $\mathbf{W}_{(1,v)} \in \mathbb{R}^{M^h \times (M^v+1)}$  and  $\mathbf{W}_{(l,v)} \in \mathbb{R}^{M^h \times (M^h+1)}$  are the combination of fully connected layer's weights and biases, and  $l \in \{2, \dots, L\}$  stands for the  $l$ -th layer of encoder; in addition,  $f(\cdot)$  is the activation function, and in this paper we use Relu [53].

Specially,  $\mathbf{h}_n^v \in \mathbb{R}^{M^h}$  is the  $v$ -th view's embedding (latent) representation of instance  $n$ , which is viewed as view-specific representation, and we can reformulate it as  $\mathbf{H}^v = \{\mathbf{h}_n^v\}_{n=1}^N \in \mathbb{R}^{M^h \times N}$ .

**Consistent attentive layer:** we try to explore the consistent feature among multiple views by utilizing the following shared weights operation:

$$\mathbf{h}_{c,n}^v = \mathbf{W}_c [\mathbf{h}_n^v; 1], \quad (3)$$

$$\mathbf{h}_{c,n} = \sum_{v=1}^V \alpha_{c,n}^v \mathbf{h}_{c,n}^v, \quad (4)$$

where the weight  $\alpha_{c,n}^v$ , which denotes the contribution for the  $v$ -th view, is calculated according to attention mechanism [12]:

$$\alpha_{c,n}^v = \mathbf{q}_c^T \tanh(\mathbf{K}_c^v [\mathbf{h}_{c,n}^v; 1]), \quad (5)$$

$$\alpha_{c,n}^v = \frac{\exp(a_{c,n}^v)}{\sum_{p=1}^V \exp(a_{c,n}^p)}. \quad (6)$$

Equation (3) implement the process that project each view's embedding vector  $\mathbf{h}_n^v$  onto a same subspace by using a shared weight  $\mathbf{W}_c \in \mathbb{R}^{M^h \times (M^h+1)}$ , and thus the  $v$ -th view's consistent representation is  $\mathbf{H}_c^v = \{\mathbf{h}_{c,n}^v\}_{n=1}^N \in \mathbb{R}^{M^h \times N}$ . Then we fuse multiple views' representation according equation (4), and the updating rule of  $\alpha_{c,n}^v$  is a special kind of attention mechanism [54] shown in equation (5) and (6). Unlike normal alignment operation, we only need one query  $\mathbf{q}_c \in \mathbb{R}^{M^h}$  and corresponding keys  $\mathbf{K}_c^v \in \mathbb{R}^{M^h \times (M^h+1)}$  (they are just viewed as trainable variables) that we can compute  $v$ -th view's contribution weight  $\alpha_{c,n}^v$  and obtain joint consistent representation  $\mathbf{H}_c = \{\mathbf{h}_{c,n}\}_{n=1}^N \in \mathbb{R}^{M^h \times N}$  fused by  $\{\mathbf{H}_c^v\}_{v=1}^V$  according to weight  $\left\{ \left\{ \alpha_{c,n}^v \right\}_{n=1}^N \right\}_{v=1}^V$ .

**Global attentive layer:** through the above steps, we have obtained view-specific representations  $\{\mathbf{H}^1, \mathbf{H}^2, \dots, \mathbf{H}^V\}$ , and consistent representation  $\mathbf{H}_c$ , then we can fuse them to create a joint latent representation consists of complementary and consistent features acquired by following steps:

$$\mathbf{z}_n = \sum_{v=1}^V \alpha_n^v \mathbf{h}_n^v + \alpha_{c,n} \mathbf{h}_{c,n}, \quad (7)$$

$$\alpha_n^v = \mathbf{q}^T \tanh(\mathbf{K}^v [\mathbf{h}_n^v; 1]), \quad (8)$$

$$\alpha_{c,n} = \mathbf{q}^T \tanh(\mathbf{K}_c [\mathbf{h}_{c,n}; 1]), \quad (9)$$

$$\alpha_n^v = \frac{\exp(a_n^v)}{\sum_{p=1}^V \exp(a_n^p) + \exp(a_{c,n})}, \quad (10)$$

$$\alpha_{c,n} = \frac{\exp(a_{c,n})}{\sum_{p=1}^V \exp(a_n^p) + \exp(a_{c,n})}. \quad (11)$$

Among them,  $\mathbf{Z} = \{\mathbf{z}_n\}_{n=1}^N \in \mathbb{R}^{M^h \times N}$  is the joint latent representation fused by view-specific and consistent representations according to equation (7) with the contribution weights

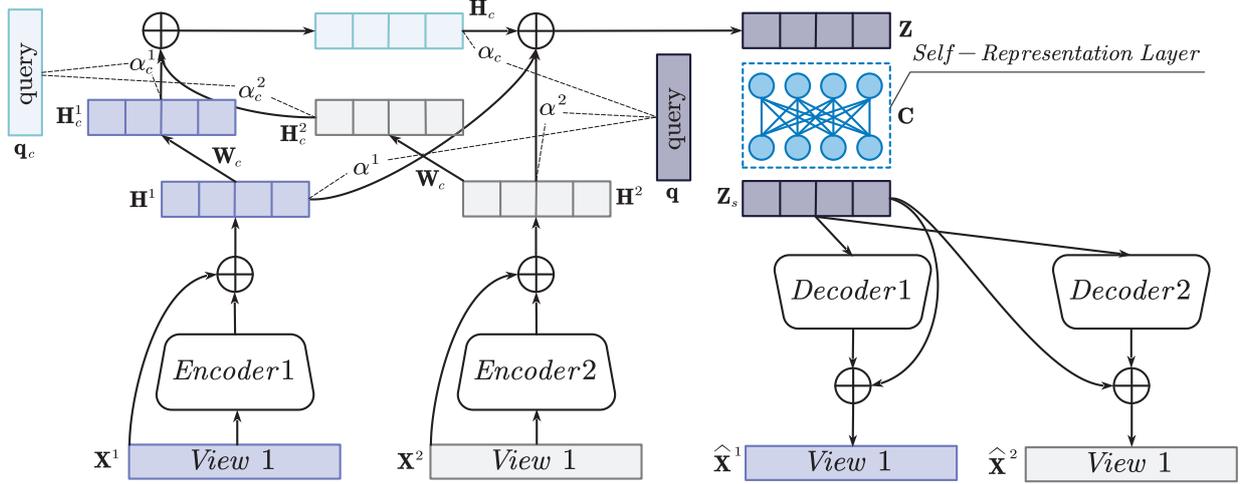


Figure 1: AMVDSN's Framework: Take data with two views as an example.  $\mathbf{X}^1$  and  $\mathbf{X}^2$  (view 1 and 2) are embedded into the same dimension through Encoder 1 and Encoder 2 with structure of shortcut connection, and the corresponding latent representations are  $\mathbf{H}^1$  and  $\mathbf{H}^2$  respectively. Further, consistent features of each view  $\mathbf{H}_c^1$  and  $\mathbf{H}_c^2$  are calculated through share weight  $\mathbf{W}_c$ , and they can be utilized to compute joint consistent feature  $\mathbf{H}_c$  through attention as:  $\mathbf{H}_c = \alpha_c^1 \mathbf{H}_c^1 + \alpha_c^2 \mathbf{H}_c^2$ . Then a global latent representation can also be calculated through attention as:  $\mathbf{Z} = \alpha^1 \mathbf{H}^1 + \alpha^2 \mathbf{H}^2 + \alpha_c \mathbf{H}_c$ .  $\mathbf{Z}_s$  is the reconstruction of  $\mathbf{Z}$  by self-representation layer, and it is used to reconstruct each original view through Decoder 1 and Decoder 2, which means  $\hat{\mathbf{X}}^1$  and  $\hat{\mathbf{X}}^2$  are the reconstructions of  $\mathbf{X}^1$  and  $\mathbf{X}^2$ .

$\{\{\alpha_n^v\}_{n=1}^N\}_{v=1}^V$  and  $\{\alpha_c^n\}_{n=1}^N$ , and equation (8) - (11) are their update formulas; similar to consistent attentive layer,  $\mathbf{q} \in \mathbb{R}^{M^h}$ ,  $\mathbf{K}^v \in \mathbb{R}^{M^h \times (M^h+1)}$ ,  $\mathbf{K}_c \in \mathbb{R}^{M^h \times (M^h+1)}$  represent query, view-specific keys, consistent key respectively.

**Self-representation layer:** the general self-representation subspace learning is usually formulated as follow:

$$\begin{aligned} \min_{\mathbf{C}} R(\mathbf{C}) \\ \text{s.t. } \mathbf{X} = \mathbf{X}\mathbf{C}, \text{diag}(\mathbf{C}) = 0. \end{aligned}$$

In this paper, we concern subspace clustering on joint latent representation  $\mathbf{Z}$ , therefore suppose that  $\mathbf{C} \in \mathbb{R}^{N \times N}$  is the self-representation coefficient matrix and the process can be formulated as follows:

$$\mathbf{Z}_s = \mathbf{Z}\mathbf{C}, \quad (12)$$

where  $\mathbf{Z}_s = \{\mathbf{z}_{s,n}\}_{n=1}^N \in \mathbb{R}^{M^h \times N}$  is the reconstruction representation of  $\mathbf{Z}$  through self-representation process, and the constraint on  $\mathbf{C}$  will be discussed in the following subsections. Note that we cannot constraint the diagonal entries of  $\mathbf{C}$  equal to zero in gradient descent, therefore when programming we can subtract the diagonal elements of  $\mathbf{C}$  in advance before we implement equation (12).

**Decoder layer:** we have obtained multiple views' global joint latent representation  $\mathbf{Z}$  and its reconstruction representation  $\mathbf{Z}_s$ , suppose that  $\mathbf{Z}_s$  can be utilized to reconstruct each view in reverse, and we have:

$$\hat{\mathbf{x}}_n^v = f(\bar{\mathbf{W}}_{(L,v)} f(\bar{\mathbf{W}}_{(L-1,v)} \dots f(\bar{\mathbf{W}}_{(1,v)} [\mathbf{z}_{s,n}; 1]))) \quad (13)$$

Among of them,  $\hat{\mathbf{X}}^v = \{\hat{\mathbf{x}}_n^v\}_{n=1}^N \in \mathbb{R}^{M^v \times N}$  is the reconstruction of the  $v$ -th view;  $\bar{\mathbf{W}}_{(l,v)} \in \mathbb{R}^{M^h \times (M^h+1)}$  and  $\bar{\mathbf{W}}_{(L,v)} \in \mathbb{R}^{M^v \times (M^h+1)}$  are the combination of fully connected layer's weights and biases,

and  $l \in \{1, \dots, L-1\}$  stands for the  $l$ -th layer of decoder; in addition,  $f(\cdot)$  is the activation function, and in this paper we use Relu.

**Shortcut connection:** Based on discussion above, we have proposed a complicated multi-view framework. In practice, the layers number  $L$  of encoder or decoder is usually set from 2 to 4, together with attentive layers and self-representation layer, the depth of our proposed model is relatively deep. Simultaneously, the main body of network structure is fully connected layer, which makes the model easy to degrade in training process. Fortunately, by introducing shortcut connection structure [15], we can effectively deal with the degradation problem and significantly improve learning performance.

In this paper, we apply shortcut connection in the modules that have more than two layers, i.e., encoder, decoder, attentive layers, and therefore we should reformulate equation (2), (7), and (13) respectively as follows:

$$\begin{aligned} \mathbf{h}_n^v &= \mathbf{W}_e^v \mathbf{x}_n^v + f(\mathbf{W}_{(L,v)} f(\mathbf{W}_{(L-1,v)} \dots f(\mathbf{W}_{(1,v)} [\mathbf{x}_n^v; 1]))) \\ \mathbf{z}_n &= \sum_{v=1}^V \alpha^v \mathbf{h}_n^v + \alpha_c \mathbf{h}_{c,n} + \frac{1}{V} \sum_{v=1}^V \mathbf{h}_n^v, \quad (14) \end{aligned}$$

$$\hat{\mathbf{x}}_n^v = \mathbf{W}_d^v \mathbf{z}_n + f(\bar{\mathbf{W}}_{(L,v)} f(\bar{\mathbf{W}}_{(L-1,v)} \dots f(\bar{\mathbf{W}}_{(1,v)} [\mathbf{z}_{s,n}; 1])))$$

Among of them,  $\mathbf{W}_e^v \in \mathbb{R}^{M^h \times (M^v)}$  and  $\mathbf{W}_d^v \in \mathbb{R}^{M^v \times (M^h)}$  are linear projections to match the corresponding dimensions [15]. Note that consistent attentive layer's depth is shallow, but this module associates with global attentive layer, which makes the framework complicated when they combined. As a consequence, we apply shortcut connection only in global attentive layer in

equation (14), and because there are multiple inputs in this sub-module, we need to calculate their mean value.

This module is significant for our proposed algorithm, because the model is a complex hybrid network and will degrade when training without it. The detail discussion on it is in subsection 4.6.

### 3.3. Loss Function

In this paper, we attempt to use deep learning framework to achieve subspace clustering, so the loss function mainly contains two parts. Based on the assumptions of subsection 2.2, the first part of loss function is referred to auto-encoder reconstruction loss:

$$\min_{\Theta} \frac{1}{NV} \sum_{v=1}^V \|\mathbf{X}^v - \hat{\mathbf{X}}^v\|_F^2 + \lambda \Omega(\mathbf{W}, \bar{\mathbf{W}}, \mathbf{W}_c), \quad (15)$$

where  $\Theta = \{\mathbf{C}, \mathbf{W}, \bar{\mathbf{W}}, \mathbf{W}_c, \mathbf{K}^v, \mathbf{K}_c^v, \mathbf{q}, \mathbf{q}_c\}$  is the variable parameters set,  $\Omega(\cdot)$  is the regularization term of model parameters, and  $\lambda$  is the trade-off parameter. In this paper, two alternatives to regularization terms are  $l_1$ -norm and  $l_2$ -norm, which can be viewed as a hyper-parameter. Then, the second part of loss function is referred to self-representation subspace learning:

$$\min_{\mathbf{C}, \mathbf{W}, \mathbf{W}_c, \mathbf{K}^v, \mathbf{K}_c^v, \mathbf{q}, \mathbf{q}_c} R(\mathbf{C}) + \frac{\lambda}{N} \|\mathbf{Z} - \mathbf{Z}_s\|_F^2, \quad (16)$$

where  $R(\cdot)$  is regularization term, and  $\lambda$  is the trade-off parameter. In subspace learning, three alternatives to regularization terms are usually  $l_1$ -norm, nuclear norm, and frobenius norm, which correspond to sparse [42], low-rank [43], least squares regression [41] subspace learning respectively. In this paper, to facilitate optimization, we select frobenius norm as the regularization term, i.e., equation (16) can be reformulated as:

$$\min_{\mathbf{C}, \mathbf{W}, \mathbf{W}_c, \mathbf{K}^v, \mathbf{K}_c^v, \mathbf{q}, \mathbf{q}_c} \|\mathbf{C}\|_F^2 + \frac{\lambda}{N} \|\mathbf{Z} - \mathbf{Z}_s\|_F^2. \quad (17)$$

In summary, we combine equation (15) and (17) to formulate joint object function:

$$\min_{\Theta} \|\mathbf{C}\|_F^2 + \frac{\lambda_1}{N} \|\mathbf{Z} - \mathbf{Z}_s\|_F^2 + \lambda_2 \left[ \frac{1}{NV} \sum_{v=1}^V \|\mathbf{X}^v - \hat{\mathbf{X}}^v\|_F^2 + \lambda_3 \Omega(\mathbf{W}, \bar{\mathbf{W}}, \mathbf{W}_c) \right], \quad (18)$$

where  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are trade-off parameters. As a consequence, by minimizing equation (18) with stochastic gradient descent, we can easily optimize self-representation matrix  $\mathbf{C}$ . In other words, we only need to build forward process of AMVDSN using deep learning framework like Tensorflow and Pytorch when programming without considering the specific optimization process.

### 3.4. Pretraining Strategy

Our proposed AMVDSN is a complex hybrid network, which is hard to directly train from scratch and trivial all-zero solution may appear when minimizing the loss function. A very effective and simple solution is to introduce pretraining strategy. We

first train an auto-encoder (without attention, shortcut connection, and self-presentation process) for each view separately to prevent interference from other views. The process is helpful for quickly finding each view's view-specific embedding representation, multiple views encoders' weights construct  $\mathbf{W}_0$ , and decoders' weights are used to construct  $\bar{\mathbf{W}}_0$ . We then set the initial values of AMVDSN's weights  $\mathbf{W}$  and  $\bar{\mathbf{W}}$  as  $\mathbf{W}_0$  and  $\bar{\mathbf{W}}_0$ , and train AMVDSN according to equation (18) with fixed  $\mathbf{W}$ .

Note that we fix  $\mathbf{W}$  because the embedding is good enough to represent the corresponding view. However,  $\bar{\mathbf{W}}$  is need to update because the input of AMVDSN's decoders is  $\mathbf{Z}_s$ , as shown in equation (13), which is different from the inputs of pre-training stage's decoders (in pretraining stage, the input of  $v$ -th view's decoder is  $\mathbf{H}^v$ ). The training process is simple, the configuration is same as the corresponding modules in AMVDSN, and with an early stopping constraint (loss no longer decreases within 200 epochs), each view's net usually converges in hundreds to thousands of epochs.

The pretraining based method is named as AMVDSN-ft, extensive experiments indicate that with the strategy can significantly improve AMVDSN's performance on most data sets, and the detailed experiments will be discussed in section 4.

## 4. Experiment

To verify the performance of our proposed method, we compare AMVDSN with the 5 baseline methods on seven real-world data sets, where one method has additional 3 variants, i.e., there are actually 8 comparison methods. In this section, we will first introduce the data sets and baseline methods, and then discuss the experiments we designed.

### 4.1. Data Sets

In this subsection, we will introduce the data construction method for short, and seven data sets include ORL<sup>1</sup>, Reuters [55], 3-sources<sup>2</sup>, Yale<sup>1</sup>, uci-digit<sup>3</sup>, Prokaryotic [56], and NUS-WIDE-OBJ<sup>4</sup>.

**ORL**: a face data with 400 images of 40 subjects, which means each category involves only 10 pictures, and three types of features (intensity, LBP and Gabor) are utilized to construct three views.

**Yale**: a face data includes 15 different persons, and there are 11 facial expression or configuration for each category. Similar to ORL, three types of features (intensity, LBP and Gabor) are utilized to construct three views.

**uci-digit**: a data set of handwritten digit with 10 categories (0-9), and each digit has 200 samples. The data contains 6 types of features, where 76 Fourier coefficients of the character shapes, 216 profile correlations, and 64 Karhunen-Love coefficients are selected to construct three views.

<sup>1</sup><http://www.cad.zju.edu.cn/home/dengcai/Data/FaceData.html>

<sup>2</sup><http://mlg.ucd.ie/datasets/3sources.html>

<sup>3</sup><http://archive.ics.uci.edu/ml/datasets/Multiple+Features>

<sup>4</sup><https://lms.comp.nus.edu.sg/wp-content/uploads/2019/research/nuswide/NUS-WIDE.html>

**Reuters:** a documents data with 6 categories which originally written in English, and they are also translated into French, German, Spanish and Italian. We view each language as a view therefore we obtain five views. All documents are in the bag-of-words representations, and we sampled 100 documents from each category to build the data set.

**3-sources:** a news data set, where 169 news belong to some dominant topic classes are available in three different online news sources (BBC, Reuters, and The Guardian), the report from each source is viewed as a view. All articles are in the bag-of-words representations.

**Prokaryotic:** a data set describes 551 prokaryotic species with heterogeneous representations, the first is instance’s textual description, which is in the bag-of-words representation and is viewed as the first view. And another two genomic representations: 1) the proteome composition, encoded as relative frequencies of amino acids; 2) the gene repertoire, encoded as presence/absence indicators of gene families in a genome, are viewed as the remaining two views. Then component analysis (PCA) is utilized to reduce the dimensionality on three views respectively, and note that we retain principal components explaining 90% of the variance.

**NUS-WIDE-OBJ:** a subset of NUS-WIDE data set with 30000 samples of 31 objects, which is intended for several object-based tasks. The data involves six types of features, and we select five of them (color histogram, color correlogram, edge direction histogram, wavelet texture, and block-wise color moments) as five views. Since the highly computational complexity of most compared baseline methods, we randomly select 1500 samples from raw data, and it is turn out that there are not too many differences between the clustering performance on the subset and whole set.

In summary, the detail information of each data are summarized in Table 1, where  $V$ ,  $C$ ,  $N$  and  $M^v$  represent views number, clusters number, instance number and each view’s dimension respectively.

Table 1: Description of Data Sets

Data	$V$	$C$	$N$	$M^v$
ORL	3	40	400	[4096, 3304, 6750]
Reuters	5	6	600	[21526, 24892, 34121, 15487, 11539]
3-sources	3	6	169	[3560, 3631, 3068]
Yale	3	15	165	[4096, 3304, 6750]
uci-digit	3	10	2000	[216, 76, 64]
Prokaryotic	3	4	551	[438, 3, 393]
NUS-WIDE-OBJ	5	31	30000	[64, 225, 144, 73, 128]

#### 4.2. Baseline Methods

To demonstrate the efficiency of our proposed algorithm, some state-of-the-art subspace clustering algorithms are selected as the baseline methods:

**LRSC** [43]: a single view subspace clustering method named Low-Rank Subspace Clustering, we conduct the algorithm on all views and select the best result.

**LMSC** [4]: a multi-view subspace clustering method named latent Multi-view Subspace Clustering. The algorithm is based

on the assumption that all views come from a same latent representation through projection of multiple certain orthogonal matrix. And the algorithm applies the latent representation into self-representation subspace clustering with nuclear norm regularization.

**CSMSC** [6]: a multi-view subspace clustering method named Consistent and Specific Multi-View Subspace Clustering. The algorithm is based on the assumption that self-representation matrix includes consistent and specific parts, and uses different norms to regularize them, where nuclear norm and frobenius norm are consistent and specific parts’ regularization terms respectively.

**MLRSSC** [5]: a multi-view subspace clustering method named Multi-view Low-Rank Sparse Subspace Clustering, which let the self-representation matrix sparse and low-rank simultaneously. Besides, the paper provides two alignment methods to ensure the multiple views’ consistent property, the first let each view adjacent to the rest views, which is also MLRSSC’s idea; the second let all views adjacent to a center (a new variable), which is called centroid-based multi-view low-rank sparse subspace clustering (**MLRSSC-C**). On the other hand, to deal with non-linear data, the paper proposed kernel style variants based on two aforementioned assumptions, which are abbreviated as **KMLRSSC** and **MLRSSC-C**.

**MvDSCN** [1]: a neural network based framework named Multi-view Deep Subspace Clustering Networks, which aims to discover the inherent structure by fusing multi-view complementary information. Up until now, only the source code on a multi-modal data set has been released, therefore we only compare with this algorithm on ORL and Yale data sets because MvDSCN was also conducted comparison experiments on them in reference [1].

**CSI** [57]: a multi-view subspace clustering method named Common Subspace Integration. The algorithm conducts self-representation reconstruction on each view and integrates different subspace into common subspace using a graph-based method. Compare with most multi-view subspace algorithms that just calculate the mean of multiple subspace, CSI integrates subspace through learning and works well. Although it gains performance enhancements, but also causes higher computational complexity.

**FCMSC** [58]: a multi-view subspace clustering method named Feature Concatenation Multi-view Subspace Clustering. The algorithm deals with the problem that multiple views have different statistic properties and simply concatenating them directly cannot derive a satisfied clustering performance.

Table 2: Model Configuration

Data	Layers Setting	$\lambda_1, \lambda_2, \lambda_3$	Regularization
ORL	[128, 128, 128, 128]	[0.5, 0.5, 0.1]	$l_2$ -norm
Reuters	[512, 512]	[0.5, 0.5, 0.01]	$l_2$ -norm
3-sources	[512, 512]	[0.5, 0.5, 0.1]	$l_1$ -norm
Yale	[128, 128]	[0.5, 0.5, 0.1]	$l_2$ -norm
uci-digit	[512, 512]	[0.5, 0.5, 0.1]	$l_1$ -norm
Prokaryotic	[128, 128]	[0.5, 0.5, 0.1]	$l_1$ -norm
NUS-WIDE-OBJ	[256, 256, 256]	[0.5, 0.5, 0.01]	$l_2$ -norm

Table 3: Clustering results on ORL, Prokaryotic, and Yale.

Dataset	Method	ACC	NMI	ARI	Precision	Recall	F-score
ORL	LRSC	0.769(0.015)	0.883(0.015)	0.688(0.037)	0.668(0.037)	0.726(0.036)	0.695(0.036)
	LMSC	0.822(0.037)	0.927(0.013)	0.774(0.043)	0.813(0.038)	0.822(0.037)	0.802(0.041)
	CSMSC	0.868(0.0012)	0.942(0.005)	0.827(0.002)	0.860(0.002)	0.804(0.003)	0.831(0.001)
	MLRSSC	0.637(0.034)	0.813(0.017)	0.524(0.037)	0.640(0.033)	0.637(0.034)	0.623(0.034)
	MLRSSC-Centroid	0.780(0.027)	0.917(0.009)	0.729(0.029)	0.785(0.033)	0.780(0.027)	0.761(0.032)
	KMLRSSC	0.786(0.041)	0.903(0.016)	0.721(0.042)	0.793(0.040)	0.786(0.041)	0.774(0.042)
	KMLRSSC-Centroid	0.783(0.031)	0.907(0.008)	0.721(0.026)	0.793(0.039)	0.783(0.031)	0.772(0.035)
	CSI	0.863(0.023)	0.930(0.009)	0.805(0.025)	0.782(0.029)	0.840(0.020)	0.810(0.025)
	FCSMC	0.847(0.019)	0.928(0.008)	0.782(0.021)	0.746(0.025)	0.833(0.020)	0.787(0.020)
	MvDSCN	0.870(0.006)	0.943(0.002)	0.819(0.001)	\	\	0.834(0.012)
AMVSC	0.887(0.010)	0.936(0.003)	0.831(0.021)	0.885(0.022)	0.887(0.010)	0.887(0.013)	
AMVSC-ft	<b>0.921(0.005)</b>	<b>0.958(0.003)</b>	<b>0.882(0.007)</b>	<b>0.918(0.005)</b>	<b>0.921(0.005)</b>	<b>0.914(0.005)</b>	
Prokaryotic	LRSC	0.582(0.005)	0.079(0.011)	0.038(0.029)	0.407(0.011)	<b>0.964(0.047)</b>	0.571(0.001)
	LMSC	0.709(0.032)	0.418(0.034)	0.394(0.052)	0.669(0.067)	0.585(0.068)	0.618(0.026)
	CSMSC	0.658(0.006)	0.352(0.004)	0.367(0.005)	0.665(0.005)	0.527(0.007)	0.588(0.004)
	MLRSSC	0.646(0.038)	0.309(0.018)	0.324(0.037)	0.529(0.011)	0.531(0.006)	0.498(0.005)
	MLRSSC-Centroid	0.620(0.009)	0.196(0.008)	0.260(0.005)	0.444(0.033)	0.399(0.018)	0.363(0.024)
	KMLRSSC	0.638(0.045)	0.420(0.041)	0.355(0.071)	0.609(0.032)	0.641(0.053)	0.589(0.041)
	KMLRSSC-Centroid	0.655(0.045)	0.405(0.030)	0.350(0.072)	0.606(0.031)	0.657(0.046)	0.597(0.042)
	CSI	0.619(0.000)	0.413(0.000)	0.291(0.000)	0.634(0.000)	0.440(0.000)	0.520(0.000)
	FCSMC	0.674(0.001)	0.443(0.000)	0.394(0.001)	0.688(0.001)	0.537(0.001)	0.603(0.001)
	AMVSC	0.823(0.301)	0.543(0.265)	0.593(0.883)	0.758(0.651)	0.794(0.744)	0.762(0.733)
AMVSC-ft	<b>0.872(0.023)</b>	<b>0.592(0.033)</b>	<b>0.683(0.047)</b>	<b>0.827(0.445)</b>	0.861(0.009)	<b>0.837(0.034)</b>	
Yale	LRSC	0.675(0.027)	0.706(0.019)	0.491(0.032)	0.501(0.030)	0.547(0.030)	0.523(0.030)
	LMSC	0.789(0.018)	0.789(0.019)	0.628(0.030)	0.632(0.029)	0.672(0.027)	0.651(0.028)
	CSMSC	0.752(0.001)	0.784(0.001)	0.615(0.005)	0.673(0.002)	0.610(0.006)	0.794(0.029)
	MLRSSC	0.660(0.050)	0.697(0.033)	0.477(0.054)	0.482(0.056)	0.544(0.045)	0.511(0.050)
	MLRSSC-Centroid	0.627(0.034)	0.702(0.021)	0.458(0.024)	0.711(0.050)	0.627(0.034)	0.648(0.041)
	KMLRSSC	0.649(0.053)	0.689(0.032)	0.485(0.046)	0.679(0.070)	0.649(0.053)	0.653(0.059)
	KMLRSSC-Centroid	0.639(0.043)	0.680(0.032)	0.480(0.044)	0.661(0.046)	0.639(0.043)	0.640(0.045)
	CSI	0.734(0.000)	0.777(0.000)	0.606(0.000)	0.606(0.000)	0.659(0.000)	0.739(0.000)
	FCSMC	0.775(0.025)	0.796(0.019)	0.589(0.041)	0.569(0.046)	0.674(0.028)	0.617(0.037)
	MvDSCN	0.824(0.004)	0.797(0.007)	0.626(0.011)	\	\	0.650(0.010)
AMVSC	0.781(0.022)	0.747(0.011)	0.540(0.074)	0.870(0.063)	0.782(0.022)	0.806(0.022)	
AMVSC-ft	<b>0.867(0.012)</b>	<b>0.844(0.013)</b>	<b>0.685(0.047)</b>	<b>0.921(0.021)</b>	<b>0.867(0.012)</b>	<b>0.877(0.009)</b>	

### 4.3. Model Configuration

Different data sets require different model configurations, and firstly, the initial method of  $\mathbf{W}$ ,  $\bar{\mathbf{W}}$  and  $\mathbf{W}_c$  is Lecun normal initializer [59, 60], which is a function in keras: keras.initializers.lecun\_norm. For optimizer, Adam [61] with the learning rate 0.001 is satisfied, and usually after 100 - 200 epochs, the network converge to optimal results. Besides, the setting of each layer’s node number, regularization, trade-off parameters  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are summarized in Table 2.

### 4.4. Experimental Setting

Before we formally discuss the experimental results, we need to explain some experimental considerations.

(1) Once the network is trained, we can use the self-representation layer’s weight  $\mathbf{C}$  to construct affinity matrix of spectral clustering, and the calculation rule is:

$$\frac{|\mathbf{C}| + |\mathbf{C}^T|}{2}.$$

However, over the years, researchers have developed some heuristics to calculate better affinity matrix, especially for neural network based model, the constraint on  $\mathbf{C}$  is not strict, therefore we utilize the method proposed by [62] to compute the

affinity matrix of the spectral clustering will derive better learning performance.

(2) We select 6 metrics to measure the clustering effect, and they are NMI (Normalized Mutual Information), ARI (Adjusted Rand Index), ACC (Accuracy), Precision, Recall, and F-score. The first two ones are traditional clustering metrics, but the last four ones are actually classification metrics, therefore we need to compute the best map of the predicted labels compared with the ground truth using Kuhn-Munkres algorithm first. Note that higher values indicates better performance for all our selected metrics.

### 4.5. Clustering Result Analysis

We conduct the experiments on seven real-world data sets using 10 baseline methods and our proposed AMVDSN and AMVDSN-ft, and then summarize the clustering results in Table 3 and 4. Obviously, AMVDSN-ft have the best performance on most data sets, especially on Prokaryotic, and it has improved performance by more than 10 percentage points on all metrics; we also find even if the pretraining strategy is not used, AMVDSN is very stable on all data sets, and achieve the best performance on NUS-WIDE-OBJ.

Experimental results indicate that compare with traditional matrix factorization based methods, AMVDSN finds better rep-

Table 4: Clustering results on Reuters, 3-sources, uci-digit and NUS-WIDE-OBJ.

Dataset	Method	ACC	NMI	ARI	Precision	Recall	F-score
Reuters	LRSC	0.209(0.007)	0.131(0.007)	0.015(0.001)	0.172(0.000)	<b>0.835(0.013)</b>	0.285(0.001)
	LMSC	0.495(0.001)	0.355(0.001)	0.223(0.010)	0.312(0.009)	0.497(0.008)	0.383(0.006)
	CSMSC	0.513(0.001)	0.358(0.001)	0.247(0.001)	0.343(0.001)	0.444(0.001)	0.387(0.001)
	MLRSSC	0.530(0.033)	0.382(0.015)	0.282(0.026)	0.501(0.055)	0.530(0.033)	0.483(0.051)
	MLRSSC-Centroid	0.514(0.035)	0.370(0.011)	0.268(0.026)	0.484(0.040)	0.514(0.035)	0.459(0.046)
	KMLRSSC	0.571(0.023)	0.374(0.019)	0.304(0.020)	0.618(0.037)	0.571(0.023)	0.567(0.030)
	KMLRSSC-Centroid	0.551(0.024)	0.357(0.016)	0.294(0.016)	0.591(0.046)	0.551(0.024)	0.540(0.027)
	CSI	0.527(0.009)	0.394(0.009)	0.288(0.007)	0.364(0.006)	0.511(0.008)	0.425(0.007)
	FCSMC	0.518(0.001)	0.369(0.001)	0.268(0.000)	0.357(0.000)	0.470(0.001)	0.405(0.000)
	AMVSC	0.592(0.005)	0.396(0.006)	0.315(0.006)	0.642(0.005)	0.592(0.005)	0.583(0.006)
AMVSC-ft	<b>0.650(0.006)</b>	<b>0.459(0.009)</b>	<b>0.380(0.010)</b>	<b>0.690(0.006)</b>	0.650(0.006)	<b>0.648(0.006)</b>	
3-sources	LRSC	0.596(0.026)	0.462(0.009)	0.351(0.020)	0.471(0.037)	0.585(0.054)	0.518(0.004)
	LMSC	0.734(0.020)	0.676(0.018)	0.579(0.017)	0.729(0.024)	0.621(0.047)	0.669(0.018)
	CSMSC	0.798(0.002)	0.747(0.005)	0.673(0.004)	0.730(0.005)	0.773(0.001)	0.751(0.003)
	MLRSSC	0.677(0.060)	0.593(0.026)	0.550(0.067)	0.535(0.048)	0.596(0.060)	0.544(0.059)
	MLRSSC-Centroid	0.666(0.051)	0.590(0.019)	0.534(0.053)	0.512(0.045)	0.551(0.068)	0.513(0.056)
	KMLRSSC	0.609(0.038)	0.526(0.027)	0.423(0.035)	0.536(0.59)	0.564(0.079)	0.517(0.060)
	KMLRSSC-Centroid	0.621(0.025)	0.531(0.020)	0.457(0.031)	0.544(0.057)	0.562(0.051)	0.519(0.046)
	CSI	0.763(0.000)	0.747(0.000)	0.622(0.000)	0.717(0.000)	0.702(0.000)	0.709(0.000)
	FCSMC	<b>0.908(0.003)</b>	<b>0.795(0.006)</b>	<b>0.793(0.008)</b>	<b>0.887(0.004)</b>	<b>0.795(0.007)</b>	<b>0.839(0.006)</b>
	AMVSC	0.823(0.014)	0.749(0.024)	0.735(0.054)	0.685(0.019)	0.758(0.006)	0.710(0.011)
AMVSC-ft	0.831(0.010)	0.741(0.016)	0.718(0.020)	0.727(0.039)	0.791(0.033)	0.746(0.035)	
uci-digit	LRSC	0.763(0.003)	0.697(0.001)	0.626(0.003)	0.650(0.003)	0.678(0.002)	0.664(0.003)
	LMSC	0.905(0.004)	0.833(0.006)	0.812(0.008)	0.835(0.007)	0.840(0.007)	0.838(0.007)
	CSMSC	0.903(0.000)	0.835(0.001)	0.810(0.000)	0.823(0.000)	0.834(0.004)	0.828(0.000)
	MLRSSC	0.881(0.061)	0.852(0.022)	0.814(0.050)	0.877(0.070)	0.881(0.061)	0.873(0.070)
	MLRSSC-Centroid	0.893(0.055)	0.853(0.023)	0.818(0.050)	0.891(0.063)	0.893(0.055)	0.887(0.063)
	KMLRSSC	0.894(0.049)	0.861(0.018)	0.826(0.044)	0.888(0.060)	0.894(0.049)	0.886(0.058)
	KMLRSSC-Centroid	0.910(0.046)	0.865(0.018)	0.840(0.039)	0.910(0.051)	0.910(0.046)	0.906(0.053)
	CSI	0.863(0.000)	<b>0.915(0.000)</b>	0.846(0.000)	0.805(0.000)	0.929(0.000)	0.862(0.000)
	FCSMC	0.916(0.001)	0.848(0.001)	0.827(0.001)	0.840(0.001)	0.849(0.001)	0.844(0.001)
	AMVSC	0.932(0.001)	0.876(0.005)	0.859(0.004)	0.935(0.001)	0.932(0.001)	0.932(0.001)
AMVSC-ft	<b>0.958(0.001)</b>	0.909(0.002)	<b>0.909(0.002)</b>	<b>0.958(0.001)</b>	<b>0.958(0.001)</b>	<b>0.958(0.001)</b>	
NUS-WIDE-OBJ	LRSC	0.139(0.002)	0.174(0.002)	0.027(0.002)	0.097(0.002)	0.055(0.001)	0.070(0.001)
	LMSC	0.177(0.004)	0.211(0.004)	0.056(0.003)	0.138(0.003)	0.074(0.002)	0.096(0.003)
	CSMSC	0.169(0.004)	0.213(0.003)	0.050(0.003)	0.128(0.004)	0.071(0.002)	0.091(0.002)
	MLRSSC	0.160(0.007)	0.209(0.005)	0.047(0.003)	0.125(0.004)	0.067(0.002)	0.088(0.003)
	MLRSSC-Centroid	0.158(0.010)	0.205(0.006)	0.044(0.005)	0.121(0.006)	0.065(0.003)	0.085(0.004)
	KMLRSSC	0.173(0.006)	0.183(0.004)	0.043(0.001)	0.095(0.002)	0.146(0.017)	0.114(0.005)
	KMLRSSC-Centroid	0.168(0.008)	0.204(0.006)	0.048(0.004)	0.125(0.005)	0.070(0.003)	0.090(0.004)
	CSI	0.171(0.002)	0.167(0.003)	0.032(0.002)	0.086(0.002)	0.137(0.001)	0.105(0.001)
	FCSMC	0.149(0.003)	0.188(0.003)	0.039(0.001)	0.115(0.002)	0.062(0.001)	0.080(0.001)
	AMVSC	<b>0.197(0.003)</b>	<b>0.229(0.004)</b>	<b>0.067(0.004)</b>	<b>0.178(0.006)</b>	<b>0.188(0.008)</b>	<b>0.156(0.004)</b>
AMVSC-ft	0.193(0.002)	0.227(0.006)	0.065(0.004)	0.176(0.006)	0.184(0.009)	0.153(0.003)	

resentation so that despite self-representation matrix is obtained through a simple strategy (the strategy is formulated as equation (17)), we still get much better clustering results. For example, [42, 43, 6] don't have additional feature extraction process; [4] computes a latent representation of multiple views but the process is actually a linear transformation; the kernel style variants of [5] need first map the original input into a high dimensional feature space. However, their strategies perform not well, and compare with them, our proposed method is based on neural network which has good non-linear ability, and attention provides AMVDSN with better views fusion performance.

On the other hand, compare with multi-view deep learning method MvDSCN we still have two advantages: (1) AMVSC-ft performs better than MvDSCN on ORL and Yale (we only list the results on such two data sets, and reason is discussed in subsection 4.2); (2) our pretraining strategy is much easier to train than MvDSCN. In pretraining stage, MvDSCN train mul-

iple auto-encoders simultaneously, the loss function is all net's joint reconstruction error, which converges slowly and often not converges after tens of thousands of epochs. But our strategy is to train multiple views' auto-encoders respectively, with a simple early stopping method (loss no longer decreases within 200 rounds), each view's auto-encoder usually converges in only hundreds to thousands of epochs.

#### 4.6. Ablation Study

In this subsection, we design some ablation experiments to verify the roles of shortcut connection and consistent attentive layer modules. AMVDSN is set as the baseline, we remove shortcut connection layers and consistent attentive layers respectively and compare their clustering accuracy, and the comparisons of them are shown in Figure 2. The orange bars indicate that training from scratch without the shortcut connection

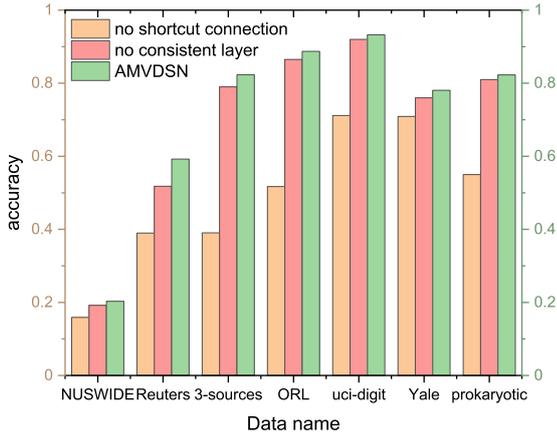


Figure 2: Ablation Study: the orange, red, and green bars show the clustering performance of three types of AMVDSN: no shortcut connection, no consistent attentive layers, and the standard AMVDSN.

module is hard or even infeasible on some data set, the reason is AMVDSN’s structure is deep and complex such that model degrades when training. Fortunately, shortcut connection deals with the problem well and even makes AMVDSN could train from scratch. When using residual module and remove consistent layer, the red bars indicate that the model can be well-trained, but the clustering performance can still be improved by adding the consistent attentive layers. The module works because we use the shared weight to extract multiple embeddings’ consensus information and fuse them as new embedding, and finally the fusion of multiple meaningful features is helpful for constructing a more complete representation.

#### 4.7. Advanced Analysis

To show the difference between representation computed by AMVDSN and each view’s original representation, we reduce  $\mathbf{Z}$  and  $\{\mathbf{X}^v\}_{v=1}^V$  of Prokaryotic to 2-dimensional data by t-sne and visualizes them in Figure 3. We find that the joint representation has a better decision boundary, especially for cluster 2. Note that for original views, view 3 is the best representation, and if we want to divide them into 4 clusters, at least 3 decision boundaries are needed. And for our obtained joint latent representation, we only need two decision boundaries, and the distinction between clusters is clearer.

Furthermore, according to the theory proposed by [45], the self-representation matrix needs to block diagonal so that it can have good clustering performance. Therefore, we visualize the self-representation  $\mathbf{C}$  of all data sets in Figure 4, which shows the matrix is block diagonal according to categories. In Figure 4, lighter colors represent stronger relevance, which means on all data sets we selected, each sample is most closely related to the samples that they belong to the same clusters, and this is the reason that self-representation subspace clustering works well.

## 5. Conclusion and Future Works

We proposed a novel multi-view subspace learning framework based on auto-encoder in this paper. By introducing self-attention mechanism, we can dynamically update views’ contribution during training process, and it also makes the fusion process more reasonable. With the joint latent representation computed by neural network, we can achieve state-of-the-art clustering results only using least squares regression subspace learning, which indicates a better feature space will derive better clustering performance in subspace learning. Besides, by introducing shortcut connection, we solve the problem of model degradation. Finally, extensive experiments on real-world data sets demonstrate the effectiveness of our proposed AMVDSN.

In the future, more advanced regularization strategies on self-representation matrix  $\mathbf{C}$  like low rank and sparse constraints will be introduced to give more reasonable meanings for deep subspace clustering. On the other hand, refer to the success of the transformer, we can attempt to treat the data of multiple views as a sequence and introduce self-representation learning into the structure of encoder-decoder, which will enhance each view’s expression ability by the help of other views and eliminate the corruptions of some disturbed views. Finally, we find that it’s hard to constrain the diagonal entries of  $\mathbf{C}$  equal to zero, we will focus on novel optimal method to handle the problem instead of using some tricks. In general, multi-view subspace clustering is a promising research direction of machine learning, and we will focus on this direction and continue to work for it.

## References

- [1] P. Zhu, B. Hui, C. Zhang, D. Du, L. Wen, and Q. Hu, Multi-view Deep Subspace Clustering Networks, (2019), arXiv preprint arXiv:1908.0.
- [2] C. Zhang, H. Fu, S. Liu, G. Liu, X. Cao, Zhang, C., Fu, H., Liu, S., Liu, G., Low-Rank Tensor Constrained Multiview Subspace Clustering., in: 2015 IEEE International Conference on Computer Vision, 2015, pp. 1582–1590.
- [3] X. Cao, C. Zhang, H. Fu, Si Liu, Hua Zhang, Diversity-induced Multi-view Subspace Clustering, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 586–594.
- [4] C. Zhang, Q. Hu, H. Fu, P. Zhu, X. Cao, Latent Multi-view Subspace Clustering, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 4333–4341.
- [5] M. Brbić, I. Kopriva, Multi-view low-rank sparse subspace clustering, Pattern Recognition 73 (2018) 247–258.
- [6] S. Luo, C. Zhang, W. Zhang, X. Cao, Consistent and specific multi-view subspace clustering, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018, pp. 3730–3737.
- [7] T. Zhou, C. Zhang, X. Peng, H. Bhaskar, J. Yang, Dual Shared-Specific Multiview Subspace Clustering, IEEE Transactions on Cybernetics 50 (8) (2019) 3517–3530.
- [8] J. Liu, Y. Jiang, Z. Li, Z. H. Zhou, H. Lu, Partially shared latent factor learning with multiview data, IEEE Transactions on Neural Networks and Learning Systems 26 (6) (2015) 1233–1246.
- [9] R. Lu, J. Liu, S. Lian, and X. Zuo, Multi-view representation learning in multi-task scene, Neural Computing and Applications 32 (14) (2020) 10403–10422.
- [10] P. Ji, T. Zhang, H. Li, M. Salzmann, I. D. Reid, Deep Subspace Clustering Networks, in: Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems, 2017, pp. 24–33.

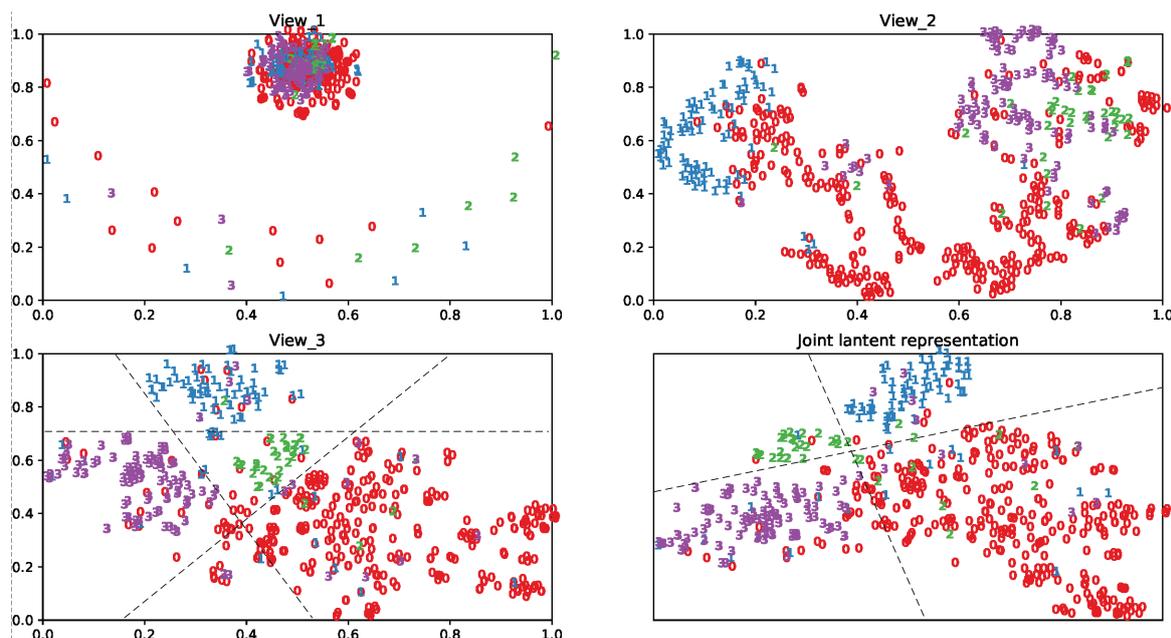


Figure 3: Visualization of original views and fused representation on Prokaryotic: 4 panels represent 2-dim representations of 3 original views and obtained joint latent representation respectively. The data set have 4 clusters, and we label them as 0 - 3, we can find that view 1's original representation is the most chaotic, and view 3's performs best, and it needs 3 decision boundaries to distinct 4 clusters. And for our computed joint latent representation, we only need 2 decision boundaries that can distinct different targets well.

[11] M. Abavisani, V. M. Patel, Deep Multimodal Subspace Clustering Networks, *IEEE Journal of Selected Topics in Signal Processing* 12 (6) (2018) 1601–1614.

[12] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is All you Need, in: *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.

[13] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krüger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, Language models are few-shot learners, (2020), arXiv preprint arXiv:2005.14165.

[14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, An image is worth 16x16 words: Transformers for image recognition at scale, (2020), arXiv preprint arXiv:2010.11929.

[15] K. He, X. Zhang, S. Ren, J. Sun, Deep Residual Learning for Image Recognition, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[16] H. HOTELLING, RELATIONS BETWEEN TWO SETS OF VARIATES, *Biometrika* 28 (34) (1936) 321–377.

[17] A. Blum and T. Mitchell, Combining Labeled and Unlabeled Data with Co-training, in: *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*, 1998, pp. 92–100.

[18] S. Sun, L. Mao, Z. Dong, and L. Wu, *Multiview Machine Learning*, 2019, Springer Singapore.

[19] Q. Liu and S. Sun, Multi-view regularized gaussian processes, in: *Advances in Knowledge Discovery and Data Mining*, 2017, pp. 655–667.

[20] J. Tang, Y. Tian, P. Zhang, and X. Liu, Multiview privileged support vector machines, *IEEE Transactions on Neural Networks and Learning Systems* 29 (8) (2018) 3463–3477.

[21] Z. Qi, Y. Tian, and Y. Shi, Laplacian twin support vector machine for semi-supervised classification, *Neural Networks* 35 (2012) 46–53.

[22] G. Chao and S. Sun, Semi-supervised multi-view maximum entropy discrimination with expectation laplacian regularization, *Information Fusion* 45 (2019) 296–306.

[23] C. Xu, D. Tao, and C. Xu, Multi-view intact space learning, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37 (12) (2015) 2531–2544.

[24] Y. Luo, D. Tao, K. Ramamohanarao, C. Xu, and Y. Wen, Tensor canonical correlation analysis for multi-view dimension reduction, *IEEE Transactions on Knowledge and Data Engineering* 27 (11) (2015) 3111–3124.

[25] Y. Yang, H. Wang, Y. Yang, and H. Wang, Multi-view clustering: A survey, *Big Data Mining and Analytics* 1 (2) (2018) 83–107.

[26] X. Peng, Z. Huang, J. Lv, H. Zhu, and J. T. Zhou, COMIC: Multi-view Clustering Without Parameter Selection, in: *Proceedings of the 36th International Conference on Machine Learning*, 2019, pp. 5092–5101.

[27] S. Huang, Z. Kang, I.W. Tsang, and Z. Xu, Auto-weighted multi-view clustering via kernelized graph learning, *Pattern Recognition* 88 (2019) 174–184.

[28] I. Muslea, S. Minton, and C.A. Knobloc, Active learning with multiple views, *Journal of Artificial Intelligence Research* 27 (2006) 203–233.

[29] Q. Zhang and S. Sun, Multiple-view multiple-learner active learning, *Pattern Recognition* 43 (9) (2010) 3113–3119.

[30] B. Tan, E. Zhong, E. Xiang, and Q. Yang, Multi-transfer: Transfer learning with multiple views and multiple sources, *Statistical Analysis and Data Mining: The ASA Data Science Journal* 7 (4) (2014) 282–293.

[31] X Zhang, X Zhang, H Liu, and X Liu, Multi-Task Multi-View Clustering, *IEEE Transactions on Knowledge and Data Engineering* 28 (12) (2016) 3324–3338.

[32] G. Andrew, R. Arora, J.A. Bilmes, and K. Livescu, Deep canonical correlation analysis, in: *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *JMLR Workshop and Conference Proceedings*, 2013, pp. 1247–1255.

[33] S. Sun, F. Jin, and W. Tu, View construction for multi-view semi-supervised learning, in: *8th International Symposium on Neural Networks*, 2011, pp. 595–601.

[34] Qi Wang, M. Chen, F. Nie, and X. Li, Detecting Coherent Groups in Crowd Scenes by Multiview Clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (1) (2019) 46–58.

[35] C. He, X. Fei, H. Li, Y. Tang, H. Liu, and Q. Chen, A Multi-View Clustering Method for Community Discovery Integrating Links and Tags, in: *IEEE 14th International Conference on e-Business Engineering*, 2017, pp.

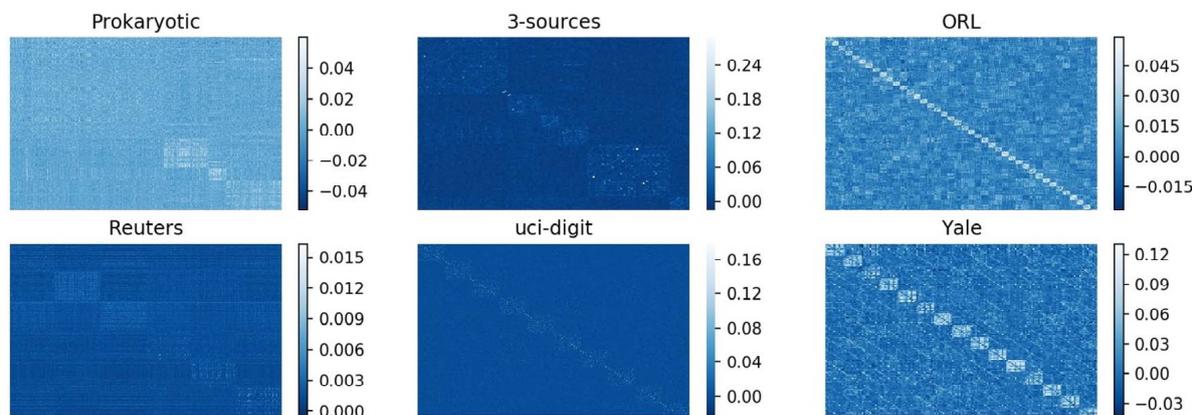


Figure 4: Correlation Diagram of Self-Representation Matrix: we illustrate the correlation heat maps on seven data sets, where lighter color indicates higher relevance. The block numbers represent the number of clusters, and the diagram indicates that each sample has higher relationship with those samples belong to same cluster.

- 23–30.
- [36] A.M. Elkahky, Y. Song, and X. He, A Multi-View Deep Learning Approach for Cross Domain User Modeling in Recommendation Systems, in: *International World Wide Web Conferences*, 2015, pp. 278–288.
- [37] C. Zhang, E. Adely, Z. Wu, G. Li, W. Lin, and D. Shen, Infant brain development prediction with latent partial multi-view representation learning, in: *2018 IEEE 15th International Symposium on Biomedical Imaging*, 2018, pp. 1048–1051.
- [38] S.R. Rao, R. Tron, R. Vidal, and Y. Ma, Motion segmentation via robust subspace separation in the presence of outlying, incomplete, or corrupted trajectories, in: *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2008, pp. 24–26.
- [39] J. Feng, Z. Lin, H. Xu, and S. Yan, Robust subspace segmentation with block-diagonal prior, in: *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3818–3825.
- [40] P. Ji, M. Salzmann, and H. Li, Shape interaction matrix revisited and robustified: Efficient subspace clustering with corrupted and incomplete data, in: *IEEE International Conference on Computer Vision*, 2015, pp. 4687–4695.
- [41] C. Lu, H. Min, Z. Zhao, L. Zhu, D. Huang, S. Yan, Robust and Efficient Subspace Segmentation via Least Squares Regression, in: *12th European Conference on Computer Vision*, 2012, pp. 347–360.
- [42] E. Elhamifar, R. Vidal, Sparse Subspace Clustering: Algorithm, Theory, and Applications, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (11) (2013) 2765–2781.
- [43] R. Vidal, P. Favaro, Low rank subspace clustering (LRSC), *Pattern Recognition Letters* 43 (2014) 47–61.
- [44] D. Luo, F. Nie, C. Ding, and H. Huang, Multi-Subspace Representation and Discovery, in: *Machine Learning and Knowledge Discovery in Databases - European Conference*, 2011, pp. 405–420.
- [45] C. Lu, J. Feng, Z. Lin, T. Mei, and S. Yan, Subspace Clustering by Block Diagonal Representation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41 (2) (2019) 487–501.
- [46] V.M. Patel, H.V. Nguyen, and R. Vidal, Latent Space Sparse Subspace Clustering, in: *IEEE International Conference on Computer Vision*, 2013, pp. 225–232.
- [47] V.M. Patel and R. Vidal, Kernel sparse subspace clustering, in: *IEEE International Conference on Image Processing*, 2014, pp. 2849–2853.
- [48] M. Yin, Y. Guo, J. Gao, Z. He, and S. Xie, Kernel Sparse Subspace Clustering on Symmetric Positive Definite Manifolds, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, 5157–5164.
- [49] S. Xiao, M. Tan, D. Xu, and Z.Y. Dong, Robust Kernel Low-Rank Representation, *IEEE Transactions on Neural Networks and Learning Systems* 27 (11) (2016) 2268–2281.
- [50] X. Peng, J. Feng, J.T. Zhou, Y. Lei, and S. Yan, Deep Subspace Clustering, *IEEE Transactions on Neural Networks and Learning Systems* (2020) 1–13.
- [51] X. Peng, J. Feng, S. Xiao, W. Yau, J.T. Zhou, and S. Yang, Structured AutoEncoders for Subspace Clustering, *IEEE Transactions on Image Processing* 27 (10) (2018) 5076–5086.
- [52] C. Zhang, H. Fu, Q. Hu, X. Cao, Y. Xie, D. Tao, and D. Xu, Generalized Latent Multi-View Subspace Clustering, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42 (1) (2020) 86–99.
- [53] X. Glorot, A. Bordes, Y. Bengio, Deep sparse rectifier neural networks, in: *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 315–323.
- [54] C. Wu, F. Wu, J. Liu, S. He, Y. Huang, X. Xie, Neural demographic prediction using search query, in: *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, 2019, pp. 654–662.
- [55] D.D. Lewis, Y. Yang, T. G. Rose, F. Li, Rcv1, A new benchmark collection for text categorization research, *Journal of Machine Learning Research* 5 (2004) 361–397.

- [56] M. Brbić, M. Piškorec, V. Vidulin, A. Kriško, T. Šmuc, F. Supek, The landscape of microbial phenotypic traits and associated genes, *Nucleic Acids Research* 44 (21) (2016) 10074–10090.
- [57] W. Weng, W. Zhou, J. Chen, H. Peng, and H. Cai, Enhancing multi-view clustering through common subspace integration by considering both global similarities and local structures, *Neurocomputing* (378) (2020) 375–386.
- [58] Q. Zheng, J. Zhu, Z. Li, S. Pang, J. Wang, and Y. Li, Feature concatenation multi-view subspace clustering, *Neurocomputing* (379) (2020) 89–102.
- [59] Y. Lecun, L. Bottou, G. Orr, and K. Muller, Efficient backdrop, *Lecture Notes in Computer Science* (1998) 9–50.
- [60] G. Klambauer, T. Unterthiner, A. Mayr, S. Hochreiter, Self-normalizing neural networks, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Self-Normalizing Neural Networks*, in: *Advances in Neural Information Processing Systems*, 2017, pp. 971–980.
- [61] D.P. Kingma and J. Ba, Adam: A method for stochastic optimization, in: *3rd International Conference on Learning Representations*, 2015.
- [62] P. Ji, M. Salzmann, and H. Li, Efficient dense subspace clustering, in: *IEEE Winter Conference on Applications of Computer Vision*, 2014, pp. 461–468.