# EXPLORING A PRINCIPLED FRAMEWORK FOR DEEP SUBSPACE CLUSTERING

#### Xianghan Meng<sup>†</sup>, Zhiyuan Huang<sup>†</sup> & Wei He

Beijing University of Posts and Telecommunications, Beijing 100876, P.R. China {mengxianghan, huangzhiyuan, wei.he}@bupt.edu.cn

#### Xianbiao Qi & Rong Xiao

Intellifusion, Shenzhen, P.R. China

#### Chun-Guang Li\*

Beijing University of Posts and Telecommunications, Beijing 100876, P.R. China lichunguang@bupt.edu.cn

#### ABSTRACT

Subspace clustering is a classical unsupervised learning task, built on a basic assumption that high-dimensional data can be approximated by a union of subspaces (UoS). Nevertheless, the real-world data are often deviating from the UoS assumption. To address this challenge, state-of-the-art deep subspace clustering algorithms attempt to jointly learn UoS representations and self-expressive coefficients. However, the general framework of the existing algorithms suffers from a catastrophic feature collapse and lacks a theoretical guarantee to learn desired UoS representation. In this paper, we present a Principled fRamewOrk for Deep Subspace Clustering (PRO-DSC), which is designed to learn structured representations and self-expressive coefficients in a unified manner. Specifically, in PRO-DSC, we incorporate an effective regularization on the learned representations into the self-expressive model, prove that the regularized self-expressive model is able to prevent feature space collapse, and demonstrate that the learned optimal representations under certain condition lie on a union of orthogonal subspaces. Moreover, we provide a scalable and efficient approach to implement our PRO-DSC and conduct extensive experiments to verify our theoretical findings and demonstrate the superior performance of our proposed deep subspace clustering approach.

### **1** INTRODUCTION

Subspace clustering is an unsupervised learning task, aiming to partition high dimensional data that are approximately lying on a union of subspaces (UoS), and finds wide-ranging applications, such as motion segmentation (Costeira & Kanade, 1998; Vidal et al., 2008; Rao et al., 2010), hybrid system identification (Vidal, 2004; Bako & Vidal, 2008), image representation and clustering (Hong et al., 2006; Lu et al., 2012), genes expression clustering (McWilliams & Montana, 2014) and so on.

Existing subspace clustering algorithms can be roughly divided into four categories: iterative methods (Tseng, 2000; Ho et al., 2003; Zhang et al., 2009), algebraic geometry based methods (Vidal et al., 2005; Tsakiris & Vidal, 2017), statistical methods (Fischler & Bolles, 1981), and spectral clustering-based methods (Chen & Lerman, 2009; Elhamifar & Vidal, 2009; Liu et al., 2010; Lu et al., 2012; You et al., 2016a; Zhang et al., 2021). Among them, spectral clustering based methods gain the most popularity due to the broad theoretical guarantee and superior performance.

The vital component in spectral clustering based methods is a so-called *self-expressive* model (El-hamifar & Vidal, 2009; 2013). Formally, given a dataset  $\mathcal{X} \coloneqq \{x_1, \dots, x_N\}$  where  $x_j \in \mathbb{R}^D$ ,

<sup>\*</sup>Corresponding author. <sup>†</sup> These two authors are equally contributed.

self-expressive model expresses each data point  $x_i$  by a linear combination of other points, i.e.,

$$\boldsymbol{x}_j = \sum_{i \neq j} c_{ij} \boldsymbol{x}_i,\tag{1}$$

where  $c_{ij}$  is the corresponding self-expressive coefficient. The most intriguing merit of the self-expressive model is that the solution of the self-expressive model under proper regularizer on the coefficients  $c_{ij}$  is guaranteed to satisfy a subspace-preserving property, namely,  $c_{ij} \neq 0$  only if  $x_i$  and  $x_j$  are in the same subspace (Elhamifar & Vidal, 2013; Soltanolkotabi & Candes, 2012; Li et al., 2018). Having had the optimal self-expressive coefficients  $\{c_{ij}\}_{i,j=1}^N$ , the data affinity can be induced by  $|c_{ij}| + |c_{ji}|$  for which spectral clustering is applied to yield the partition of the data.

Despite the broad theoretical guarantee, the vanilla self-expressive model still faces great challenges when applied to the complex real-world data that may not well align with the UoS assumption. Earlier works devote to address this deficiency by learning a linear transform of the data (Patel et al., 2013; 2015) or introducing a nonlinear kernel mapping (Patel & Vidal, 2014) under which the representations of the data are supposed to be aligned with the UoS assumption. However, there is a lack of principled mechanism to guide the learning of the linear transforms or the design of the nonlinear kernels to guarantee the representations of the data to form a UoS structure.

To handle complex real-world data, in the past few years, there is a surge of interests in designing deep subspace clustering frameworks, e.g., (Ji et al., 2017; Peng et al., 2018; Zhou et al., 2018; Zhang et al., 2019a; Dang et al., 2020; Peng et al., 2020; Lv et al., 2021; Wang et al., 2023b; Zhao et al., 2024). In these works, usually a deep neural network-based representation learning module is integrated to the self-expressive model, to learn the representations  $Z \in \mathbb{R}^{d \times N}$  and the self-expressive coefficients  $C = \{c_{ij}\}_{i,j=1}^{N}$  in a joint optimization framework. However, as analyzed in (Haeffele et al., 2021) that, the optimal representations Z of these methods tend to catastrophically collapse into subspaces with dimensions much lower than the ambient space, which is detrimental to subspace clustering and there is no evidence that the learned representations form a UoS structure.

In this paper, we attempt to propose a Principled fRamewOrk for Deep Subspace Clustering (PRO-DSC), which is able to simultaneously learn structured representations and self-expressive coefficients. Specifically, in PRO-DSC, we incorporate an effective regularization on the learned representations into the self-expressive model and prove that our PRO-DSC can effectively prevent feature collapse. Moreover, we demonstrate that our PRO-DSC under certain condition can yield structured representations forming a UoS structure and provide a scalable and efficient approach to implement PRO-DSC. We conduct extensive experiments on the synthetic data and six benchmark datasets to verify our theoretical findings and the superior performance of our proposed approach.

Contributions. The contributions of the paper are highlighted as follows.

- 1. We propose a Principled fRamewOrk for Deep Subspace Clustering (PRO-DSC) that learns both structured representations and self-expressive coefficients simultaneously, in which an effective regularization on the learned representations is incorporated to prevent feature space collapse.
- We provide a rigorous analysis for the optimal solution of our PRO-DSC, derive a sufficient condition that guarantees the learned representations to escape from feature collapse, and further demonstrate that our PRO-DSC under certain condition can yield structured representations of a UoS structure.
- 3. We conduct extensive experiments to verify our theoretical findings and to demonstrate the superior performance of the proposed approach.

To the best of our knowledge, this is the first principled framework for deep subspace clustering that is guaranteed to prevent feature collapse problem and is shown to yield the UoS representations.

# 2 DEEP SUBSPACE CLUSTERING: A PRINCIPLED FRAMEWORK, JUSTIFICATION, AND IMPLEMENTATION

In this section, we review the popular framework for deep subspace clustering, called Self-Expressive Deep Subspace Clustering (SEDSC) at first, then present our principled framework for deep subspace clustering and provide a rigorous characterization of the optimal solution and the

property of the learned structured representations. Finally we describe a scalable implementation based on differential programming for the proposed framework. Please refer to Appendix A for the detailed proofs of our theoretical results.

#### 2.1 PREREQUISITE

To apply subspace clustering to complex real-world data that may not well align with the UoS assumption, there has been a surge of interests in exploiting deep neural networks to learn representations and then apply self-expressive model to the learned representations, e.g., (Peng et al., 2016; 2018; Ji et al., 2017; Zhou et al., 2018; Zhang et al., 2019a; Dang et al., 2020; Peng et al., 2020; Lv et al., 2021; Wang et al., 2023b; Zhao et al., 2024).

Formally, the optimization problem of these SEDSC models can be formulated as follows:<sup>1</sup>

$$\min_{\boldsymbol{Z},\boldsymbol{C}} \quad \frac{1}{2} \|\boldsymbol{Z} - \boldsymbol{Z}\boldsymbol{C}\|_F^2 + \beta \cdot r(\boldsymbol{C}) \quad \text{s.t.} \quad \|\boldsymbol{Z}\|_F^2 = N,$$
(2)

where  $Z \in \mathbb{R}^{d \times N}$  denotes the learned representation,  $C \in \mathbb{R}^{N \times N}$  denotes the self-expressive coefficient matrix, and  $\beta > 0$  is a hyper-parameter. The following lemma characterizes the property of the optimal solution Z for problem (2).

**Lemma 1** (Haeffele et al., 2021). The rows of the optimal solution Z for problem (2) are the eigenvectors that associate with the smallest eigenvalues of  $(I - C)(I - C)^{\top}$ .

In other words, the optimal representation Z in SEDSC is restricted to an extremely "narrow" subspace whose dimension is much smaller than d, leading to an undesirable collapsed solution.<sup>2</sup>

#### 2.2 OUR PRINCIPLED FRAMEWORK FOR DEEP SUBSPACE CLUSTERING

In this paper, we attempt to propose a principled framework for deep subspace clustering that provably learns structured representations with maximal intrinsic dimensions.

To be specific, we try to optimize the self-expressive model (2) while preserving the intrinsic dimension of the representation space. Other than using the rank, which is a common measure of the dimension, inspired by (Fazel et al., 2003; Ma et al., 2007; Yu et al., 2020; Liu et al., 2022), we propose to prevent the feature space collapse by incorporating the  $\log \det(\cdot)$ -based concave smooth surrogate which is defined as follows:

$$R(\boldsymbol{Z};\alpha) \coloneqq \log \det(\boldsymbol{I} + \alpha \boldsymbol{Z}^{\top} \boldsymbol{Z}), \tag{3}$$

where  $\alpha > 0$  is the hyper-parameter. Unlike the commonly used nuclear norm, which is a convex surrogate of the rank, the  $\log \det(\cdot)$ -based function is concave and differentiable, offers a tighter approximation and encourages learning subspaces with maximal intrinsic dimensions.<sup>3</sup>

By incorporating the maximization of  $R(\mathbf{Z}; \alpha)$  as a regularizer into the formulation of SEDSC in (2), we have a Principled fRamewOrk for Deep Subspace Clustering (PRO-DSC):

$$\min_{\mathbf{Z},\mathbf{C}} \quad -\frac{1}{2}\log\det\left(\mathbf{I} + \alpha \mathbf{Z}^{\top}\mathbf{Z}\right) + \frac{\gamma}{2}\|\mathbf{Z} - \mathbf{Z}\mathbf{C}\|_{F}^{2} + \beta \cdot r(\mathbf{C}) \quad \text{s.t.} \quad \|\mathbf{Z}\|_{F}^{2} = N, \quad (4)$$

where  $\gamma > 0$  is a hyper-parameter. Now, we will give our theoretical findings for problem (4).

**Theorem 1** (Eigenspace Alignment). Denote the optimal solution of PRO-DSC in (4) as  $(\mathbf{Z}_{\star}, \mathbf{C}_{\star})$ ,  $\mathbf{G}_{\star} := \mathbf{Z}_{\star}^{\top} \mathbf{Z}_{\star}$  and  $\mathbf{M}_{\star} := (\mathbf{I} - \mathbf{C}_{\star})(\mathbf{I} - \mathbf{C}_{\star})^{\top}$ . Then  $\mathbf{G}_{\star}$  and  $\mathbf{M}_{\star}$  share eigenspaces, i.e.,  $\mathbf{G}_{\star}$  and  $\mathbf{M}_{\star}$  can be diagonalized simultaneously by  $\mathbf{U} \in \mathcal{O}(N)$  where  $\mathcal{O}(N)$  is an orthogonal group.

Note that Theorem 1 provides a perspective from eigenspace alignment for analyzing the property of the optimal solution. Figure 1(a) and (b) show empirical evidences to demonstrate that alignment occurs during the training period, where  $G_b = Z_b^{\top} Z_b$ ,  $M_b = (I - C_b)(I - C_b)^{\top}$ ,  $Z_b \in \mathbb{R}^{d \times n_b}$ ,  $C_b \in \mathbb{R}^{n_b \times n_b}$  and  $n_b$  is batch size, are computed in mini-batch training at different epoch.

<sup>&</sup>lt;sup>1</sup>Without loss of generality, we omit the constraint diag(C) = 0 throughout the analysis.

<sup>&</sup>lt;sup>2</sup>The dimension equals to the multiplicity of the smallest eigenvalues of  $(I - C)(I - C)^{\top}$ .

<sup>&</sup>lt;sup>3</sup>Please refer to (Ma et al., 2007) for a packing-ball interpretation.



Figure 1: Empirical Validation to Eigenspace Alignment and Noncollapse Representation in Mini-batch on CIFAR-100. (a): Alignment error curve during the training period. (b): Eigenspace correlation curves measured via  $\langle u_j, \frac{G_b u_j}{\|G_b u_j\|_2} \rangle$  for  $j = 1, \dots, n_b$ . (c) and (d): Eigenvalue curves.



Figure 2: Empirical Validation to Noncollapse Representation on CIFAR-10 and CIFAR-100. Clustering accuracy (ACC%) and subspace-preserving representation error (SRE%) are displayed under varying  $\alpha$  and  $\gamma$ . When collapse occurs, both ACC and SRE dramatically degenerate. The perceivable phase transition phenomenon is consistent with the condition to avoid collapse.

Next, we will analyze problem (4) from the perspective of alternating optimization. When Z is fixed, the optimization problem with respect to (w.r.t.) C reduces to a standard self-expressive model, which has been extensively studied in (Soltanolkotabi & Candes, 2012; Pimentel-Alarcon & Nowak, 2016; Wang & Xu, 2016; Li et al., 2018; Tsakiris & Vidal, 2018). On the other hand, when C is fixed, the optimization problem w.r.t. Z becomes:

$$\min_{\boldsymbol{Z}} \quad -\frac{1}{2} \log \det \left( \boldsymbol{I} + \alpha \boldsymbol{Z}^{\top} \boldsymbol{Z} \right) + \frac{\gamma}{2} \| \boldsymbol{Z} - \boldsymbol{Z} \boldsymbol{C} \|_{F}^{2} \quad \text{s.t.} \quad \| \boldsymbol{Z} \|_{F}^{2} = N,$$
(5)

which is a *non-convex* optimization problem, whose optimal solution remains under-explored.

In light of the fact that G and M converge to share eigenspaces, we decompose G and M to  $U \operatorname{Diag}(\lambda_G^{(1)}, \dots, \lambda_G^{(N)}) U^{\top}$  and  $U \operatorname{Diag}(\lambda_M^{(1)}, \dots, \lambda_M^{(N)}) U^{\top}$ , respectively. Recall that  $G := Z^{\top}Z$ ,  $M := (I - C)(I - C)^{\top}$ , by using the eigenvalue decomposition, we reformulate problem (5) into a *convex* problem w.r.t.  $\{\lambda_G^{(i)}\}_{i=1}^{\min\{d,N\}}$  (See Appendix A) and have the following result. **Theorem 2** (Noncollapse Representation). Suppose that G and M are aligned in the same

eigenspaces and  $\gamma < \frac{1}{\lambda_{\max}(M)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$ . Then we have: a)  $\operatorname{rank}(\mathbf{Z}_{\star}) = \min\{d, N\}$ , and b) the singular values  $\sigma_{\mathbf{Z}_{\star}}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_{M}^{(i)} + \nu_{\star}} - \frac{1}{\alpha}}$  for all  $i = 1, \ldots, \min\{d, N\}$ , where  $\mathbf{Z}_{\star}$  and  $\nu_{\star}$  are the optimal primal solution and dual solution, respectively.

Theorem 2 characterizes the optimal solution for problem (5). Recall that SEDSC in (2) yields a collapsed solution, where rank( $Z_{\star}$ )  $\ll \min\{d, N\}$ ; whereas the rank of the minimizers for PRO-DSC in (5) satisfies that rank( $Z_{\star}$ ) =  $\min\{d, N\}$ . In Figure 1(c) and (d), we show the curves of the eigenvalues of  $G_b$  and  $M_b$ , which are computed in the mini-batch training at different epoch, demonstrating that the learned representation does no longer collapse. In Figure 2, we show the subspace clustering accuracy (ACC) and subspace-representation error<sup>4</sup> (SRE) as a function of the parameters  $\alpha$  and  $\gamma$ . The phase transition phenomenon around  $\gamma < \frac{1}{\lambda_{\max}(M)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$  well illustrates the sufficient condition in Theorem 2 to avoid representation collapse.

<sup>&</sup>lt;sup>4</sup>For each column  $c_j$  in C, SRE is computed by  $\frac{100}{N} \sum_j (1 - \sum_i w_{ij} \cdot |c_{ij}|) / ||c_j||_1$ , where  $w_{ij} \in \{0, 1\}$  is the ground-truth affinity.



Figure 3: Empirical Validation to Structured Representation on CIFAR-10. Gram matrices for CLIP features X and learned representations Z are shown in (a) and (b); whereas Data visualization of the samples from three categories  $X_{(3)}$  and  $Z_{(3)}$  via PCA are shown in (c) and (d), respectively.

Furthermore, from the perspective of joint optimizing Z and C, the following theorem demonstrates that PRO-DSC promotes a union-of-orthogonal-subspaces representation Z and block-diagonal self-expressive matrix C under certain condition.

**Theorem 3.** Suppose that the sufficient conditions to prevent feature collapse are satisfied. Without loss of generality, we further assume that the columns in matrix  $\mathbf{Z}$  are arranged into k blocks according to a certain  $N \times N$  permutation matrix  $\Gamma$ , i.e.,  $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_k]$ . Then the condition for that PRO-DSC promotes the optimal solution  $(\mathbf{Z}_*, \mathbf{C}_*)$  to have desired structure (i.e.,  $\mathbf{Z}_*^\top \mathbf{Z}_*$  and  $\mathbf{C}_*$  are both block-diagonal), is that  $\langle (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \rangle \to 0$ , where  $\mathbf{G}^* \coloneqq \text{Diag}(\mathbf{G}_{11}, \mathbf{G}_{22}, \dots, \mathbf{G}_{kk})$  and  $\mathbf{G}_{jj}$  is the block Gram matrix corresponding to  $\mathbf{Z}_j$ .

Theorem 3 suggests that Remark 1. our PRO-DSC is able to promote learning representations and self-expressive matrix with desired structures, i.e., the representations form a union of orthogonal subspaces and accordingly the self-expressive matrix is block-diagonal, when the condition  $\langle (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top}, \boldsymbol{G} - \boldsymbol{G}^* \rangle \to 0$ is met. We call this condition a compatibly structured coherence (CSC), which relates to the properties of the distribution of the representations in Z and the selfcoefficients in C. While it is not possible for us to give a theoretical justification when the CSC condition will be satisfied in general, we do have the empirical evidence that our implementation for PRO-DSC with careful designs does approximately satisfy such a condition and thus yields representations and self-expressive matrix with desired structure (See Figure 3).<sup>5</sup>



Figure 4: Empirical validation to Theorem 3 in Mini-batch on CIFAR-10. The mean curves of the absolute values of the in-block-diagonal entries (thick) and the off-block-diagonal entries (thin) are displayed along with the CSC condition (gray) during training PRO-DSC.

In Figure 4, we show the curves for the compatibly structured coherence (CSC) condition, and for the average values of the entries in  $|G_b^*|$ ,  $|G_b - G_b^*|$ ,  $|C_b^*|$ ,  $|C_b - C_b^*|$  computed in mini-batch during training PRO-DSC on CIFAR-10. As illustrated, the CSC condition is progressively satisfied and consequently the average off-block values  $|G_b - G_b^*|$  and  $|C_b - C_b^*|$  gradually decrease, while the average in-block values  $|G_b^*|$  and  $|C_b^*|$  gradually increase, which empirically validates that PRO-DSC promotes block-diagonal  $G_b$  and  $C_b$ .

<sup>&</sup>lt;sup>5</sup>Please refer to Appendix B.2 for more details about Figures 1-4.

#### 2.3 SCALABLE IMPLEMENTATION

Existing SEDSC models typically use autoencoders to learn the representations and learn the selfexpressive matrix C through an  $N \times N$  fully-connected layer (Peng et al., 2016; 2018; Ji et al., 2017; Zhou et al., 2018; Zhang et al., 2019a). While such implementation is straightforward, there are two major drawbacks: a) since that the number of self-expressive coefficients is quadratic to the number of data points, solving these coefficients suffers from expensive computation burden; b) the learning process is transductive, i.e., the network parameters cannot be generalized to unseen data.

To address these issues, similar to (Zhang et al., 2021), we reparameterize the self-expressive coefficients  $c_{ij}$  by a neural network. Specifically, the input data  $\boldsymbol{x}_i$  is fed into a neural network  $\boldsymbol{h}(\cdot; \boldsymbol{\Psi}) : \mathbb{R}^D \to \mathbb{R}^d$  to yield normalized representations, i.e.,

$$\boldsymbol{y}_i \coloneqq \boldsymbol{h}(\boldsymbol{x}_i; \boldsymbol{\Psi}) / \| \boldsymbol{h}(\boldsymbol{x}_i; \boldsymbol{\Psi}) \|_2, \tag{6}$$

where  $\Psi$  denotes all the parameters in  $h(\cdot)$ . Then, the parameterized self-expressive matrix  $C_{\Psi}$  is generated by:

$$\boldsymbol{C}_{\boldsymbol{\Psi}} \coloneqq \mathcal{P}(\boldsymbol{Y}^{\top}\boldsymbol{Y}), \tag{7}$$

where  $\boldsymbol{Y} \coloneqq [\boldsymbol{y}_1, \dots, \boldsymbol{y}_N] \in \mathbb{R}^{d \times N}$  and  $\mathcal{P}(\cdot)$  is the sinkhorn projection (Cuturi, 2013), which has been widely applied in deep clustering (Caron et al., 2020; Ding et al., 2023).<sup>6</sup> To enable efficient representation learning, we introduce another learnable mapping  $\boldsymbol{f}(\cdot; \boldsymbol{\Theta}) : \mathbb{R}^D \to \mathbb{R}^d$ , for which

$$\boldsymbol{z}_{j} \coloneqq \boldsymbol{f}(\boldsymbol{x}_{j}; \boldsymbol{\Theta}) / \| \boldsymbol{f}(\boldsymbol{x}_{j}; \boldsymbol{\Theta}) \|_{2}$$
 (8)

is the learned representation for the input  $x_j$ , where  $\Theta$  denotes the parameters in  $f(\cdot)$  to learn the structured representation  $Z_{\Theta} \coloneqq [z_1, \ldots, z_N] \in \mathbb{R}^{d \times N}$ .

Therefore, our principled framework for deep subspace clustering (PRO-DSC) in (4) can be reparameterized and reformulated as follows:

$$\min_{\boldsymbol{\Theta},\boldsymbol{\Psi}} \quad \mathcal{L}(\boldsymbol{\Theta},\boldsymbol{\Psi}) \coloneqq -\frac{1}{2} \log \det \left( \boldsymbol{I} + \alpha \boldsymbol{Z}_{\boldsymbol{\Theta}}^{\top} \boldsymbol{Z}_{\boldsymbol{\Theta}} \right) + \frac{\gamma}{2} \left\| \boldsymbol{Z}_{\boldsymbol{\Theta}} - \boldsymbol{Z}_{\boldsymbol{\Theta}} \boldsymbol{C}_{\boldsymbol{\Psi}} \right\|_{F}^{2} + \beta \cdot r\left( \boldsymbol{C}_{\boldsymbol{\Psi}} \right).$$
(9)

To strengthen the block-diagonal structure of self-expressive matrix, we choose the block-diagonal regularizer (Lu et al., 2018) for  $r(C_{\Psi})$ . To be specific, given the data affinity  $A_{\Psi}$ , which is induced by default as  $A_{\Psi} := \frac{1}{2} \left( |C_{\Psi}| + |C_{\Psi}^{\top}| \right)$ , the block diagonal regularizer is defined as:

$$r(\boldsymbol{C}_{\boldsymbol{\Psi}}) \coloneqq \|\boldsymbol{A}_{\boldsymbol{\Psi}}\|_{\boldsymbol{\varepsilon}},\tag{10}$$

where  $\|A_{\Psi}\|_{\overline{k}}$  is the sum of the k smallest eigenvalues of the Laplacian matrix of the affinity  $A_{\Psi}$ .<sup>7</sup>

Consequently, the parameters in  $\Theta$  and  $\Psi$  of reparameterized PRO-DSC can be trained by Stochastic Gradient Descent (SGD) with the loss function  $\mathcal{L}(\Theta, \Psi)$  defined in (9). For clarity, we summarize the procedure for training and testing of our PRO-DSC in Algorithm 1.

**Remark 2.** We note that all the commonly used regularizers with extended block-diagonal property for self-expressive model as discussed in (Lu et al., 2018) can be used to improve the block-diagonal structure of self-expressive matrix. More interestingly, the specific type of the regularizers is not essential owning to the learned structured representation (Please refer to Table 3 for details), and using a specific regularizer or not is also not essential since that the SGD-based optimization also induces some implicit regularization, e.g., low-rank (Gunasekar et al., 2017; Arora et al., 2019).

#### 3 EXPERIMENTS

To validate our theoretical findings and to demonstrate the performance of our proposed framework, we conduct extensive experiments on synthetic data (Sec. 3.1) and real-world data (Sec. 3.2). Implementation details and more results are provided in Appendices B.1 and B.3, respectively.

<sup>&</sup>lt;sup>6</sup>In practice, we set  $\operatorname{diag}(C_{\Psi}) = 0$  to prevent trivial solution  $C_{\Psi} = I$ .

<sup>&</sup>lt;sup>7</sup>Recall that the number of zero eigenvalues of the Laplacian matrix equals to the number of connected components in the graph (von Luxburg, 2007).

Algorithm 1 Scalable & Efficient Implementation of PRO-DSC via Differential Programming

**Input:** Dataset  $\mathcal{X} = \mathcal{X}_{\text{train}} \cup \mathcal{X}_{\text{test}}$ , batch size  $n_b$ , hyper-parameters  $\alpha, \beta, \gamma$ , number of iterations T, learning rate  $\eta$ 

**Initialization:** Random initialize the parameters  $\Psi, \Theta$  in the networks  $h(\cdot; \Psi)$  and  $f(\cdot; \Theta)$ Training:

1: for t = 1, ..., T do

Sample a batch  $X_b \in \mathbb{R}^{D \times n_b}$  from  $\mathcal{X}_{\text{train}}$ 2: # Forward propagation

Compute self-expressive matrix  $C_b \in \mathbb{R}^{n_b \times n_b}$  by Eqs. (6–7) 3:

- Compute representations  $Z_b \in \mathbb{R}^{d \times n_b}$  by Eq. (8) 4:
  - *# Backward propagation*
- 5:
- Compute gradients:  $\nabla_{\Psi} := \frac{\partial \mathcal{L}}{\partial \Psi}, \nabla_{\Theta} := \frac{\partial \mathcal{L}}{\partial \Theta}$ Update  $\Psi$  and  $\Theta$  via:  $\Psi \leftarrow \Psi \eta \cdot \nabla_{\Psi}, \Theta \leftarrow \Theta \eta \cdot \nabla_{\Theta}$ 6:
- 7: end for

#### Testing:

- 8: Compute self-expressive matrix  $C_{\text{test}}$  by Eqs. (6–7) for  $\mathcal{X}_{\text{test}}$
- 9: Apply spectral clustering on the affinity  $A_{\text{test}}$

#### **EXPERIMENTS ON SYNTHETIC DATA** 3.1

To validate whether PRO-DSC resolves the collapse issue in SEDSC and learns representations with a UoS structure, we first follow the procedure in (Ding et al., 2023) to generate two sets of synthetic data, as shown in the first column of Figure 5, and then visualize in Figure 5(b)-(e) the learned representations which are obtained from different methods on these synthetic data.

We observe that the SEDSC model overly compress all the representations to a closed curve on the hypersphere; whereas with increased weights (i.e.,  $\gamma \uparrow$ ) of the self-expressive term, the representations collapse to a few points. Our PRO-DSC yields linearized representations lying on orthogonal subspaces in both cases, confirming the effectiveness of our approach. Nevertheless, MLC (Ding et al., 2023) yields representations approximately on orthogonal subspaces.



Figure 5: Visualization Experiments on Synthetic Data.

#### 3.2 EXPERIMENTS ON REAL-WORLD DATA

To evaluate the performance of our proposed approach, we conduct experiments on six real-world image datasets, including CIFAR-10, CIFAR-20, CIFAR-100, ImageNet-Dogs-15, Tiny-ImageNet-200, and ImageNet-1k, with the pretrained CLIP features<sup>8</sup> (Radford et al., 2021), and compare to several baseline methods, including classical clustering algorithms, e.g., k-means (MacQueen, 1967) and spectral clustering (Shi & Malik, 2000), subspace clustering algorithm, e.g., EnSC (You et al., 2016a) and SENet (Zhang et al., 2021), deep clustering algorithms, e.g., SCAN (Van Gansbeke et al., 2020), TEMI (Adaloglou et al., 2023) and CPP (Chu et al., 2024), and deep subspace clustering algorithms, e.g., DSCNet (Ji et al., 2017) and EDESC (Cai et al., 2022). We measure clustering

<sup>&</sup>lt;sup>8</sup>Please refer to Appendix B.3 for the results on other pre-trained models.

Mathad	CIFA	R-10	CIFA	R-20	CIFA	R-100	TinyImg	gNet-200	ImgNet	Dogs-15	Image	Net-1k
Method	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
k-means	83.5	84.1	46.9	49.4	52.8	66.8	54.1	73.4	52.7	53.6	53.9	79.8
SC	79.8	84.8	53.3	61.6	66.4	77.0	62.8	77.0	48.3	45.7	56.0	81.2
SSCOMP	85.5	83.0	61.4	63.4	55.6	69.7	56.7	72.7	25.6	15.9	44.1	74.4
EnSC	95.4	90.3	61.0	68.7	67.0	77.1	64.5	77.7	57.9	56.0	59.7	83.7
SENet	91.2	82.5	65.3	68.6	67.0	74.7	63.9	76.6	58.7	55.3	53.2	78.1
SCAN	95.1	90.3	60.8	61.8	64.1	70.8	56.5	72.7	70.5	68.2	54.4	76.8
TEMI	<u>96.9</u>	<u>92.6</u>	61.8	64.5	73.7	79.9	-	-	-	-	<u>64.0</u>	-
CPP	96.8	92.3	67.7	70.5	75.4	82.0	63.4	75.5	83.0	81.5	62.0	82.1
EDESC	84.2	79.3	48.7	49.1	53.1	68.6	51.3	68.8	53.3	47.9	46.5	75.5
DSCNet	78.5	73.6	38.6	45.7	39.2	53.4	62.3	68.3	40.5	30.1	OOM	OOM
Our PRO-DSC	97.2±0.2	$92.8 \pm 0.4$	71.6±1.2	$73.2 \pm 0.5$	77.3±1.0	$82.4 \pm 0.5$	69.8±1.1	$80.5 \pm 0.7$	$84.0 \pm 0.6$	$81.2 \pm 0.8$	65.0±1.2	$83.4 \pm 0.6$

Table 1: **Clustering performance Comparison on the CLIP features.** The best results are in bold and the second best results are underlined. "OOM" means out of GPU memory.

performance using clustering accuracy (ACC) and normalized mutual information (NMI), and report the experimental results in Table 1, where the results of our PRO-DSC are averaged over 10 trials (with  $\pm$ std). Since that for most baselines, except for TEMI, the clustering performance with the CLIP feature has not been reported, we conduct experiments using the implementations provided by the authors. For TEMI, we directly cited the results from (Adaloglou et al., 2023).

**Performance comparison.** As shown in Table 1, our PRO-DSC significantly outperforms subspace clustering algorithms, e.g., SSCOMP, EnSC and SENet, and deep subspace clustering algorithms, e.g., DSCNet and EDESC. Moreover, our PRO-DSC obtains better performance than the state-of-the-art deep clustering and deep manifold clustering methods, e.g., SCAN, TEMI and CPP.

Validation to the theoretical results. To validate whether the alignment emerges and when representations collapse occurs during the training period, we compute  $G_b = Z_b^{\top} Z_b$  and  $M_b = (I - C_b)(I - C_b)^{\top}$  in mini-batch at different epoch during the training period, and then measure the alignment error via  $||G_bM_b - M_bG_b||_F$  and the eigenspace correlation via  $\langle u_j, \frac{G_bu_j}{||G_bu_j||_2} \rangle$  where  $u_j$  is the *j*-th ending eigenvector<sup>9</sup> of  $M_b$  for  $j = 1, \dots, n_b$ , and plot the eigenvalues of  $G_b$  and  $M_b$ , where  $n_b$  is the sample size per mini-batch. Moreover, we also record empirical performance ACC and SRE on CIFAR-10 and CIFAR-100 under varying hyper-parameters  $\alpha$  and  $\gamma$  to validate the condition in Theorem 2 to avoid collapse. Experimental results are displayed in Figures 1 and 2. We observe that  $G_b$  and  $M_b$  are increasingly aligned and the representations will no longer collapse provided that the parameters are properly set. More details are provided in Section B.2.

**Evaluation on learned representations.** To quantitatively evaluate the effectiveness of the learned representations, we run k-means (MacQueen, 1967), spectral clustering (Shi & Malik, 2000), and EnSC (You et al., 2016a) on four datasets with three different features: a) the CLIP features, b) the representations learned via CPP, and c) the representations learned by our PRO-DSC. Experimental results are shown in Figure 6 (and more results are given in Table B.4 of Appendix B.3). We observe that the representations learned by our PRO-DSC outperform the CLIP features and the CPP representations in most cases across different clustering algorithms and datasets. Notably, the clustering accuracy with the representations learned by our PRO-DSC exceeds 90% on CIFAR-10 and 75% on CIFAR-100, whichever clustering algorithm is used. Besides, the clustering performance is further improved by using the learnable mapping  $h(\cdot; \Psi)$ , indicating a good generalization ability.



Figure 6: Clustering accuracy with CLIP features and learned representations.

<sup>&</sup>lt;sup>9</sup>The eigenvectors are sorted according to eigenvalues of  $M_b$  in ascending order.

Sensitivity to hyper-parameters. In Figure 2, we verify that our PRO-DSC yields satisfactory results when the conditions in Theorem 2 to avoid collapse are met. Moreover, we evaluate the performance sensitivity to hyper-parameters  $\gamma$  and  $\beta$  by experiments on the CLIP features of CIFAR-10, CIFAR-100 and TinyImageNet-200 with varying  $\gamma$  and  $\beta$ . In Figure 7, we observe that the clustering performance maintains satisfactory under a broad range of  $\gamma$  and  $\beta$ .



Figure 7: Evaluation on sensitivity to hyper-parameters  $\gamma$  and  $\beta$  on three datasets.

Time and memory cost. The most time-consuming operations in our PRO-DSC are computing the term involving  $\log \det(\cdot)$  and the term  $\|A\|_{\overline{k}}$  involving eigenvalue decomposition, respectively. The time complexity for  $\log \det(\cdot)$  is  $\mathcal{O}(\min\{n_b^3, d^3\})$  due to the commutative property of  $\log \det(\cdot)$ function (Yu et al., 2020), and the time complexity for  $\|A\|_{\overline{k}}$  is  $\mathcal{O}(kn_b^2)$ .<sup>10</sup> Therefore the overall time complexity of our PRO-DSC is  $\mathcal{O}(kn_b^2 + \min\{n_b^3, d^3\})$ . Note that TEMI (Adaloglou et al., 2023) employs H = 50 cluster heads during training, adding further time and memory costs and CPP (Chu et al., 2024) involves computing  $\log \det(\cdot)$  for  $n_b + 1$  times, leading to complexity  $\mathcal{O}((n_b + 1) \min\{n_b^3, d^3\})$ . The computation time and memory costs are shown in Table 2 for which all the experiments are conducted on a single NVIDIA RTX 3090 GPU and Intel Xeon Platinum 8255C CPU. We read that our PRO-DSC significantly reduces the time consumption, particularly for datasets with a large number of clusters.

Table 2: Comparison on time (s) and memory cost (MiB). "OOM" means out of GPU memory.

Mathada	Complexity	CIFAR-10		CIF	AR-100	ImageNet-1k	
Methods	Complexity	Time	Memory	Time	Memory	Time	Memory
SEDSC	$O(N^2d)$	-	OOM	-	OOM	-	OOM
TEMI	$\mathcal{O}(Hn_b d^2)$	6.9	1,766	5.1	2,394	262.1	2,858
CPP	$\mathcal{O}((n_b+1)\min\{n_b^3, d^3\})$	3.5	3,802	7.1	10,374	1441.2	22,433
PRO-DSC	$\mathcal{O}(kn_b^2 + \min\{n_b^3, d^3\})$	4.5	2,158	4.0	2,328	90.0	2,335

T 11 0	A 1 1 /*	. 1.		1.00	1	c	1	1 .
Inhia 4	Ablation	cfuidiac	nn	dittoront	Ince	tunctione	and	ramilarizare
Table	ADIALION	Studies	VII.	uniterent	1055	runctions	anu	TUPUIALIZEIS.

							CIFAR-10		CIFAR-100		ImgNetDogs-15	
	$\mathcal{L}_1$	$\mathcal{L}_2$	$\ oldsymbol{A}\ _{\kappa}$	$\ m{C}\ _1$	$\ oldsymbol{C}\ _F^2$	$\ m{C}\ _*$	ACC	NMI	ACC	NMI	ACC	NMI
on							56.9	47.7	54.6	60.9	46.7	37.1
lati	$\checkmark$						69.6	56.4	64.7	71.7	10.5	1.7
Ab	$\checkmark$	$\checkmark$					97.0	93.0	74.6	80.9	80.9	78.8
zer	$\checkmark$						97.0	92.6	75.2	81.1	81.3	79.1
ari							97.0	92.6	75.2	80.9	80.9	78.8
gul				$\checkmark$			96.7	91.9	76.4	81.8	81.0	78.8
Re	$\checkmark$	$\checkmark$					97.2	92.8	77.3	82.4	84.0	81.2

Ablation study. To verify the effectiveness of each components in the loss function of our PRO-DSC, we conduct a set of ablation studies with the CLIP features on CIFAR-10, CIFAR-100, and ImageNetDogs-15, and report the results in Table 3, where  $\mathcal{L}_1 \coloneqq -\frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z}_{\Theta}^{\top} \mathbf{Z}_{\Theta})$  and  $\mathcal{L}_2 \coloneqq \frac{1}{2} \| \mathbf{Z}_{\Theta} - \mathbf{Z}_{\Theta} \mathbf{C}_{\Psi} \|_F^2$ . The absence of the term  $\mathcal{L}_1$  leads to catastrophic feature collapse (as demonstrated in Sec. 2.1); whereas without the self-expressive  $\mathcal{L}_2$ , the model lacks a loss function

<sup>&</sup>lt;sup>10</sup>For an  $N \times N$  matrix, the complexity of computing its k eigenvalues by Lanczos algorithm is  $\mathcal{O}(kN^2)$  and the complexity of computing its det(·) is  $\mathcal{O}(N^3)$ .

for learning the self-expressive coefficients. In both cases, clustering performance drops significantly. More interestingly, when we replace the block diagonal regularizer  $\|A\|_{\mathbb{K}}$  with  $\|C\|_1$ ,  $\|C\|_*$ , and  $\|C\|_F^2$  or even drop the explicit regularizer  $r(\cdot)$ , the clustering performance still maintains satisfactory. This confirms that the choice of the regularizer is not essential owning to the structured representations learned by our PRO-DSC.

## 4 RELATED WORK

**Deep subspace clustering.** To tackle with complex real world data, a number of Self-Expressive Deep Subspace Clustering (SEDSC) methods have been developed in the past few years, e.g., (Peng et al., 2016; 2018; Ji et al., 2017; Zhou et al., 2018; Zhang et al., 2019a;b; Dang et al., 2020; Peng et al., 2020; Lv et al., 2021; Cai et al., 2022; Wang et al., 2023b). The key step in SEDSC is to adopt a deep learning module to embed the input data into feature space. For example, deep autoencoder network is adopted in (Peng et al., 2016; 2018), deep convolutional autoencoder network is used in (Ji et al., 2017; Zhou et al., 2018; Zhang et al., 2019a). Unfortunately, as pointed out in (Haeffele et al., 2021), the existing SEDSC methods suffer from a catastrophic feature collapse and there is no evidence that the learned representations align with a UoS structure. To date, however, a principled deep subspace clustering framework has not been proposed.

**Deep clustering.** Recently, most of state-of-the-art deep clustering methods adopt a two-step procedure: at the first step, self-supervised learning based pre-training, e.g., SimCLR (Chen et al., 2020a), MoCo (He et al., 2020), BYOL (Grill et al., 2020) and SwAV (Caron et al., 2020), is adopted to learn the representations; and then deep clustering methods are incorporated to refine the representations, via, e.g., pseudo-labeling (Caron et al., 2018; Van Gansbeke et al., 2020; Park et al., 2021; Niu et al., 2022), cluster-level contrastive learning (Li et al., 2021), local and global neighbor matching (Dang et al., 2021), graph contrastive learning (Zhong et al., 2021), self-distillation (Adaloglou et al., 2023). Though the clustering performance has been improved remarkably, the underlying geometry structure of the learned representations is unclear and ignored.

Representation learning with a UoS structure. The methods for representation learning that favor a UoS structure are pioneered in supervised setting, e.g., (Lezama et al., 2018; Yu et al., 2020). In (Lezama et al., 2018), a nuclear norm based geometric loss is proposed to learn representations that lie on a union of orthogonal subspaces; in (Yu et al., 2020), a principled framework called Maximal Coding Rate Reduction (MCR<sup>2</sup>) is proposed to learn representations that favor the structure of a union of orthogonal subspaces (Wang et al., 2024). More recently, the MCR<sup>2</sup> framework is modified to develop deep manifold clustering methods, e.g., NMCE (Li et al., 2022), MLC (Ding et al., 2023) and CPP (Chu et al., 2024). In (Li et al., 2022), the MCR<sup>2</sup> framework combines with constrastive learning to perform manifold clustering and representation learning; in (Ding et al., 2023), the MCR<sup>2</sup> framework combines with doubly stochastic affinity learning to perform manifold linearizing and clustering; and in (Chu et al., 2024), the features from large pre-trained model (e.g., CLIP) are adopted to evaluate the performance of (Ding et al., 2023). While the MCR<sup>2</sup> framework has been modified in these methods for manifold clustering, none of them provides theoretical justification to yield structured representations. Though our PRO-DSC shares the same regularizer defined in Eq. (3) with MLC (Ding et al., 2023), we are for the first time to adopt it into the SEDSC framework to attack the catastrophic feature collapse issue with theoretical analysis.

# 5 CONCLUSION

We presented a Principled fRamewOrk for Deep Subspace Clustering (PRO-DSC), which jointly learn structured representations and self-expressive coefficients. Specifically, our PRO-DSC incorporates an effective regularization into self-expressive model to prevent the catastrophic representation collapse with theoretical justification. Moreover, we demonstrated that our PRO-DSC is able to learn structured representations that form a desirable UoS structure, and also developed an efficient implementation based on reparameterization and differential programming. We conducted extensive experiments on synthetic data and six benchmark datasets to verify the effectiveness of our proposed approach and validate our theoretical findings.

#### **ACKNOWLEDGMENTS**

The authors would like to thank the constructive comments from anonymous reviewers. This work is supported by the National Natural Science Foundation of China under Grant 61876022.

#### ETHICS STATEMENT

In this work, we aim to extend traditional subspace clustering algorithms by leveraging deep learning techniques to enhance their representation learning capabilities. Our research does not involve any human subjects, and we have carefully ensured that it poses no potential risks or harms. Additionally, there are no conflicts of interest, sponsorship concerns, or issues related to discrimination, bias, or fairness associated with this study. We have taken steps to address privacy and security concerns, and all data used comply with legal and ethical standards. Our work fully adheres to research integrity principles, and no ethical concerns have arisen during the course of this study.

#### **Reproducibility Statement**

To ensure the reproducibility of our work, we have released the source code. Theoretical proofs of the claims made in this paper are provided in Appendix A, and the empirical validation of these theoretical results is shown in Figures 2-4, with further detailed explanations in Appendix B.2. All datasets used in our experiments are publicly available, and we have provided a comprehensive description of the data processing steps in Appendix B.1. Additionally, detailed experimental settings and configurations are outlined in Appendix B.1 to facilitate the reproduction of our results.

#### REFERENCES

- Nikolas Adaloglou, Felix Michels, Hamza Kalisch, and Markus Kollmann. Exploring the limits of deep image clustering using pretrained models. In *British Machine Vision Conference*, pp. 297–299, 2023.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32:7411–7422, 2019.
- Laurent Bako and René Vidal. Algebraic identification of MIMO SARX models. In *International Workshop on Hybrid Systems: Computation and Control*, pp. 43–57, 2008.
- Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of deep networks. *Advances in Neural Information Processing Systems*, pp. 153–160, 2006.
- Jinyu Cai, Jicong Fan, Wenzhong Guo, Shiping Wang, Yunhe Zhang, and Zhao Zhang. Efficient deep embedded subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21–30, 2022.
- Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision*, pp. 132–149, 2018.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9650–9660, 2021.
- Jianlong Chang, Gaofeng Meng, Lingfeng Wang, Shiming Xiang, and Chunhong Pan. Deep selfevolution clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4): 809–823, 2018.

- Guangliang Chen and Gilad Lerman. Spectral curvature clustering (SCC). *International Journal of Computer Vision*, 81(3):317–330, 2009.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pp. 1597–1607, 2020a.
- Ying Chen, Chun-Guang Li, and Chong You. Stochastic sparse subspace clustering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4155–4164, 2020b.
- Tianzhe Chu, Shengbang Tong, Tianjiao Ding, Xili Dai, Benjamin David Haeffele, René Vidal, and Yi Ma. Image clustering via the principle of rate reduction in the age of pretrained models. In *International Conference on Learning Representations*, 2024.
- Joao Paulo Costeira and Takeo Kanade. A multibody factorization method for independently moving objects. *International Journal of Computer Vision*, 29:159–179, 1998.
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. Advances in Neural Information Processing Systems, 26:2292–2300, 2013.
- Zhiyuan Dang, Cheng Deng, Xu Yang, and Heng Huang. Multi-scale fusion subspace clustering using similarity constraint. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6657–6666, 2020.
- Zhiyuan Dang, Cheng Deng, Xu Yang, Kun Wei, and Heng Huang. Nearest neighbor matching for deep clustering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13693–13702, 2021.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE* Signal Processing Magazine, 29(6):141–142, 2012.
- Tianjiao Ding, Shengbang Tong, Kwan Ho Ryan Chan, Xili Dai, Yi Ma, and Benjamin D. Haeffele. Unsupervised manifold linearizing and clustering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5450–5461, October 2023.
- Ehsan Elhamifar and René Vidal. Sparse subspace clustering. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pp. 2790–2797, 2009.
- Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2765–2781, 2013.
- Maryam Fazel, Haitham Hindi, and Stephen P Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. In *American Control Conference*, volume 3, pp. 2156–2162, 2003.
- Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24 (6):381–395, 1981.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, pp. 6151–6159, 2017.

- Benjamin D Haeffele, Chong You, and René Vidal. A critique of self-expressive deep subspace clustering. In International Conference on Learning Representations, 2021.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.
- Jeffrey Ho, Ming-Husang Yang, Jongwoo Lim, Kuang-Chih Lee, and David Kriegman. Clustering appearances of objects under varying illumination conditions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11–18, 2003.
- Wei Hong, John Wright, Kun Huang, and Yi Ma. Multiscale hybrid linear models for lossy image representation. *IEEE Transactions on Image Processing*, 15(12):3655–3671, 2006.
- Zhizhong Huang, Jie Chen, Junping Zhang, and Hongming Shan. Learning representation for clustering via prototype scattering and positive sampling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(6):7509–7524, 2023.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456, 2015.
- Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering networks. Advances in Neural Information Processing Systems, pp. 24–33, 2017.
- Yuheng Jia, Jianhong Cheng, Hui LIU, and Junhui Hou. Towards calibrated deep clustering network. In *International Conference on Learning Representations*, 2025.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In International Conference on Learning Representations, 2014.
- José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. OLE: Orthogonal low-rank embedding - a plug and play geometric loss for deep learning. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pp. 8109–8118, 2018.
- Chun-Guang Li, Chong You, and René Vidal. Structured sparse subspace clustering: A joint affinity learning and subspace clustering framework. *IEEE Transactions on Image Processing*, 26(6): 2988–3001, 2017.
- Chun-Guang Li, Chong You, and René Vidal. On geometric analysis of affine sparse subspace clustering. *IEEE Journal on Selected Topics in Signal Processing*, 12(6):1520–1533, 2018.
- Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive clustering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 8547–8555, 2021.
- Zengyi Li, Yubei Chen, Yann LeCun, and Friedrich T Sommer. Neural manifold clustering and embedding. *arXiv preprint arXiv:2201.10000*, 2022.
- Derek Lim, René Vidal, and Benjamin D Haeffele. Doubly stochastic subspace clustering. *arXiv* preprint arXiv:2011.14859, 2020.
- Guangcan Liu, Zhouchen Lin, and Yong Yu. Robust subspace segmentation by low-rank representation. In *International Conference on Machine Learning*, pp. 663–670, 2010.
- Xin Liu, Zhongdao Wang, Ya-Li Li, and Shengjin Wang. Self-supervised learning via maximum entropy coding. Advances in Neural Information Processing Systems, 35:34091–34105, 2022.
- Canyi Lu, Hai Min, Zhong-Qiu Zhao, Lin Zhu, De-Shuang Huang, and Shuicheng Yan. Robust and efficient subspace segmentation via least squares regression. In *European Conference on Computer Vision*, pp. 347–360, 2012.

- Canyi Lu, Jiashi Feng, Zhouchen Lin, Tao Mei, and Shuicheng Yan. Subspace clustering by block diagonal representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2): 487–501, 2018.
- Juncheng Lv, Zhao Kang, Xiao Lu, and Zenglin Xu. Pseudo-supervised deep subspace clustering. *IEEE Transactions on Image Processing*, 30:5252–5263, 2021.
- Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of multivariate mixed data via lossy coding and compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(9):1546–1562, 2007.
- J. MacQueen. Some methods for classification and analysis of multivariate observations. In Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, pp. 281–297, 1967.
- Ryan McConville, Raul Santos-Rodriguez, Robert J Piechocki, and Ian Craddock. N2d:(not too) deep clustering via clustering the local manifold of an autoencoded embedding. In *Proceedings* of the International Conference on Pattern Recognition, pp. 5145–5152, 2021.
- Brian McWilliams and Giovanni Montana. Subspace clustering of high dimensional data: a predictive approach. Data Mining and Knowledge Discovery, 28(3):736–772, 2014.
- Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning*, pp. 807–814, 2010.
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In Proceedings of the Indian Conference on Computer Vision, Graphics & Image Processing, pp. 722–729, 2008.
- Chuang Niu, Hongming Shan, and Ge Wang. SPICE: Semantic pseudo-labeling for image clustering. *IEEE Transactions on Image Processing*, 31:7264–7278, 2022.
- Foivos Ntelemis, Yaochu Jin, and Spencer A Thomas. Information maximization clustering via multi-view self-labelling. *Knowledge-Based Systems*, 250:109042, 2022.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024.
- Sungwon Park, Sungwon Han, Sundong Kim, Danu Kim, Sungkyu Park, Seunghoon Hong, and Meeyoung Cha. Improving unsupervised image clustering with robust learning. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12278–12287, 2021.
- Vishal M Patel and René Vidal. Kernel sparse subspace clustering. In *Proceedings of the IEEE International Conference on Image Processing*, pp. 2849–2853, 2014.
- Vishal M Patel, Hien Van Nguyen, and René Vidal. Latent space sparse subspace clustering. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 225–232, 2013.
- Vishal M Patel, Hien Van Nguyen, and René Vidal. Latent space sparse and low-rank subspace clustering. *IEEE Journal of Selected Topics in Signal Processing*, 9(4):691–701, 2015.
- Xi Peng, Shijie Xiao, Jiashi Feng, Wei-Yun Yau, and Zhang Yi. Deep subspace clustering with sparsity prior. In *International Joint Conference on Artificial Intelligence*, pp. 1925–1931, 2016.
- Xi Peng, Jiashi Feng, Shijie Xiao, Wei-Yun Yau, Joey Tianyi Zhou, and Songfan Yang. Structured autoencoders for subspace clustering. *IEEE Transactions on Image Processing*, 27(10):5076–5086, 2018.

- Xi Peng, Jiashi Feng, Joey Tianyi Zhou, Yingjie Lei, and Shuicheng Yan. Deep subspace clustering. *IEEE Transactions on Neural Networks and Learning Systems*, 31(12):5509–5521, 2020.
- Daniel Pimentel-Alarcon and Robert Nowak. The information-theoretic requirements of subspace clustering with missing data. In *International Conference on Machine Learning*, pp. 802–810, 2016.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763, 2021.
- Shankar Rao, Roberto Tron, René Vidal, and Yi Ma. Motion segmentation in the presence of outlying, incomplete, or corrupted trajectories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(10):1832–1845, 2010.
- Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- Mahdi Soltanolkotabi and Emmanuel J Candes. A geometric analysis of subspace clustering with outliers. *Annals of Statistics*, 40(4):2195–2238, 2012.
- Manolis Tsakiris and René Vidal. Algebraic clustering of affine subspaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(2):482–489, 2017.
- Manolis Tsakiris and René Vidal. Theoretical analysis of sparse subspace clustering with missing entries. In *International Conference on Machine Learning*, pp. 4975–4984, 2018.
- Paul Tseng. Nearest *q*-flat to *m* points. *Journal of Optimization Theory and Applications*, 105(1): 249–252, 2000.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of Machine Learning Research, 9(11), 2008.
- Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. SCAN: Learning to classify images without labels. In *European Conference on Computer Vision*, pp. 268–285, 2020.
- René Vidal. Identification of PWARX hybrid models with unknown and possibly different orders. In *Proceedings of the American Control Conference*, pp. 547–552, 2004.
- René Vidal, Yi Ma, and Shankar Sastry. Generalized Principal Component Analysis (GPCA). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(12):1–15, 2005.
- René Vidal, Roberto Tron, and Richard Hartley. Multiframe motion segmentation with missing data using PowerFactorization, and GPCA. *International Journal of Computer Vision*, 79(1):85–105, 2008.
- Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.
- Libin Wang, Yulong Wang, Hao Deng, and Hong Chen. Attention reweighted sparse subspace clustering. *Pattern Recognition*, 139:109438, 2023a.
- Peng Wang, Huikang Liu, Druv Pai, Yaodong Yu, Zhihui Zhu, Qing Qu, and Yi Ma. A global geometric analysis of maximal coding rate reduction. In *International Conference on Machine Learning*, 2024.
- Shiye Wang, Changsheng Li, Yanming Li, Ye Yuan, and Guoren Wang. Self-supervised information bottleneck for deep multi-view subspace clustering. *IEEE Transactions on Image Processing*, 32: 1555–1567, 2023b.
- Yu-Xiang Wang and Huan Xu. Noisy sparse subspace clustering. *Journal of Machine Learning Research*, 17(12):1–41, 2016.

- Lai Wei, Zhengwei Chen, Jun Yin, Changming Zhu, Rigui Zhou, and Jin Liu. Adaptive graph convolutional subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6262–6271, 2023.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International Conference on Machine Learning*, pp. 478–487, 2016.
- Jianwei Yang, Devi Parikh, and Dhruv Batra. Joint unsupervised learning of deep representations and image clusters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5147–5156, 2016.
- Chong You, Chun-Guang Li, Daniel Robinson, and René Vidal. Oracle based active set algorithm for scalable elastic net subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3928–3937, 2016a.
- Chong You, Daniel Robinson, and René Vidal. Scalable sparse subspace clustering by orthogonal matching pursuit. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3918–3927, 2016b.
- Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. Advances in Neural Information Processing Systems, 33:9422–9434, 2020.
- Pengxin Zeng, Yunfan Li, Peng Hu, Dezhong Peng, Jiancheng Lv, and Xi Peng. Deep fair clustering via maximizing and minimizing mutual information: Theory, algorithm and metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23986– 23995, 2023.
- Hongjing Zhang and Ian Davidson. Deep fair discriminative clustering. arXiv preprint arXiv:2105.14146, 2021.
- Junjian Zhang, Chun-Guang Li, Chong You, Xianbiao Qi, Honggang Zhang, Jun Guo, and Zhouchen Lin. Self-supervised convolutional subspace clustering network. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5473–5482, 2019a.
- Shangzhi Zhang, Chong You, René Vidal, and Chun-Guang Li. Learning a self-expressive network for subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12393–12403, 2021.
- Teng Zhang, Arthur Szlam, and Gilad Lerman. Median *k*-flats for hybrid linear modeling with many outliers. In *IEEE/CVF International Conference on Computer Vision Workshops*, pp. 234–241, 2009.
- Tong Zhang, Pan Ji, Mehrtash Harandi, Wenbing Huang, and Hongdong Li. Neural collaborative subspace clustering. In *International Conference on Machine learning*, pp. 7384–7393, 2019b.
- Chen Zhao, Chun-Guang Li, Wei He, and Chong You. Deep self-expressive learning. In *The First Conference on Parsimony and Learning*, volume 234, pp. 228–247, 2024.
- Huasong Zhong, Jianlong Wu, Chong Chen, Jianqiang Huang, Minghua Deng, Liqiang Nie, Zhouchen Lin, and Xian-Sheng Hua. Graph contrastive clustering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9224–9233, 2021.
- Pan Zhou, Yunqing Hou, and Jiashi Feng. Deep adversarial subspace clustering. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1596–1604, 2018.

# SUPPLEMENTARY MATERIAL FOR "EXPLORING A PRINCIPLED FRAMEWORK FOR DEEP SUBSPACE CLUSTERING"

The supplementary materials are divided into three parts. In Section A, we present the proofs of our theoretical results. In Section B, we present the supplementary materials for experiments, including experimental details (Sec. B.1), empirical validation on our theoretical results (Sec. B.2), and more experimental results (Sec. B.3). In Section C, we discuss about the limitations and failure cases of PRO-DSC.

#### A PROOFS OF MAIN RESULTS

As a preliminary, we start by introducing a lemma from (Haeffele et al., 2021) and provide its proof for the convenience of the readers.

**Lemma 1** (Haeffele et al., 2021). The rows of the optimal solution Z for problem (2) are the eigenvectors that associate with the smallest eigenvalues of  $(I - C)(I - C)^{\top}$ .

*Proof.* We note that:

$$\|\boldsymbol{Z} - \boldsymbol{Z}\boldsymbol{C}\|_F^2 = \operatorname{Tr}\left(\boldsymbol{Z}\left(\boldsymbol{I} - \boldsymbol{C}\right)\left(\boldsymbol{I} - \boldsymbol{C}\right)^\top \boldsymbol{Z}^\top\right) = \sum_{i=1}^d \boldsymbol{z}^{(i)}(\boldsymbol{I} - \boldsymbol{C})(\boldsymbol{I} - \boldsymbol{C})^\top \boldsymbol{z}^{(i)\top},$$

where  $z^{(i)}$  is the *i*<sup>th</sup> row of Z, thus problem (2) is reformulated as:

$$\min_{\{\boldsymbol{z}^{(i)}\}_{i=1}^{d}, \boldsymbol{C}} \quad \frac{1}{2} \sum_{i=1}^{d} \boldsymbol{z}^{(i)} (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top} \boldsymbol{z}^{(i)\top} + \beta \cdot r(\boldsymbol{C})$$
s.t.  $\|\boldsymbol{Z}\|_{F}^{2} = N.$ 
(11)

Without loss of generality, the magnitude of each row of Z is assumed to be fixed, i.e.,  $||z^{(i)}||_2^2 = \tau_i$ , i = 1, ..., d, where  $\sum_{i=1}^{d} \tau_i = N$ . Then, the optimization problem becomes:

$$\min_{\{\boldsymbol{z}^{(i)}\}_{i=1}^{d}, \boldsymbol{C}} \quad \frac{1}{2} \sum_{i=1}^{d} \boldsymbol{z}^{(i)} (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top} \boldsymbol{z}^{(i)\top} + \beta \cdot r(\boldsymbol{C})$$
  
s.t.  $\|\boldsymbol{z}^{(i)}\|_{2}^{2} = \tau_{i}, \quad i = 1, \dots, d.$  (12)

The Lagrangian of problem (12) is:

$$\mathcal{L}(\{\boldsymbol{z}^{(i)}\}_{i=1}^{d}, \boldsymbol{C}, \{\nu_{i}\}_{i=1}^{d}) \coloneqq \frac{1}{2} \sum_{i=1}^{d} \boldsymbol{z}^{(i)} (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top} \boldsymbol{z}^{(i)\top} + \beta \cdot \boldsymbol{r}(\boldsymbol{C}) + \frac{1}{2} \sum_{i=1}^{d} \nu_{i} (\|\boldsymbol{z}^{(i)}\|_{2}^{2} - \tau_{i}),$$
(13)

where  $\{\nu_i\}_{i=1}^d$  are the Lagrangian multipliers. The necessary conditions for optimal solution are:

$$\begin{cases} \nabla_{\boldsymbol{z}^{(i)}} \mathcal{L} = \boldsymbol{z}^{(i)} (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top} + \nu_i \boldsymbol{z}^{(i)} = \boldsymbol{0}, \\ \|\boldsymbol{z}^{(i)}\|_2^2 = \tau_i, \quad i = 1, \dots, d, \end{cases}$$
(14)

which implies that the optimal solutions  $z^{(i)}$  are eigenvectors of  $(I - C)(I - C)^{\top}$ .

By further considering the objective functions, the optimal solution  $z^{(i)}$  should be the eigenvectors w.r.t. the *smallest* eigenvalues of  $(I-C)(I-C)^{\top}$  for all  $i \in \{1, \ldots, d\}$ . The corresponding optimal value is  $\frac{1}{2}\lambda_{\min}((I-C)(I-C)^{\top})\sum_{i=1}^{d}\tau_i + \beta \cdot r(C) = \frac{N}{2}\lambda_{\min}((I-C)(I-C)^{\top}) + \beta \cdot r(C)$ , which is irrelevant to  $\{\tau_i\}_{i=1}^{d}$ .

Therefore, we conclude that the rows of optimal solution Z to problem (2) are eigenvectors that associate with the smallest eigenvalues of  $(I - C)(I - C)^{\top}$ .

**Lemma A1.** Suppose that matrices  $A, B \in \mathbb{R}^{n \times n}$  are symmetric, then AB = BA if and only if A and B can be diagonalized simultaneously by  $U \in O(n)$ , where O(n) is an orthogonal group.

Now we present our theorem about the optimal solution of problem PRO-DSC in (4) with its proof.

**Theorem 1.** Denote the optimal solution of PRO-DSC in (4) as  $(\mathbf{Z}_{\star}, \mathbf{C}_{\star})$ , then  $\mathbf{G}_{\star}$  and  $\mathbf{M}_{\star}$  share eigenspaces, where  $\mathbf{G}_{\star} \coloneqq \mathbf{Z}_{\star}^{\top} \mathbf{Z}_{\star}, \mathbf{M}_{\star} \coloneqq (\mathbf{I} - \mathbf{C}_{\star})(\mathbf{I} - \mathbf{C}_{\star})^{\top}$ , i.e.,  $\mathbf{G}_{\star}$  and  $\mathbf{M}_{\star}$  can be diagonalized simultaneously by  $\mathbf{U} \in \mathcal{O}(N)$  where  $\mathcal{O}(N)$  is an orthogonal matrix group.

*Proof.* We first consider the subproblem of PRO-DSC problem in (4) with respect to Z and prove that for all  $C \in \mathbb{R}^{N \times N}$ , the corresponding optimal  $Z_{\star,C}$  satisfies  $G_{\star,C}M = MG_{\star,C}$ , where  $G_{\star,C} = Z_{\star,C}^{\top} Z_{\star,C}$ ,  $M := (I - C)(I - C)^{\top}$ , implying that  $G_{\star,C}$  and M share eigenspace. Then, we will demonstrate that  $G_{\star}$  and  $M_{\star}$  share eigenspace.

The subproblem with respect to  $Z_C$  is reformulated into the following semi-definite program:

$$\min_{\boldsymbol{G}_{\boldsymbol{C}}} -\frac{1}{2} \log \det \left( \boldsymbol{I} + \alpha \boldsymbol{G}_{\boldsymbol{C}} \right) + \frac{\gamma}{2} \operatorname{tr}(\boldsymbol{G}_{\boldsymbol{C}} \boldsymbol{M}) \\
\text{s.t.} \quad \boldsymbol{G}_{\boldsymbol{C}} \succeq \boldsymbol{0}, \ \operatorname{tr}(\boldsymbol{G}_{\boldsymbol{C}}) = N,$$
(15)

which has the Lagrangian as:

$$\mathcal{L}(\boldsymbol{G}_{\boldsymbol{C}},\boldsymbol{\Delta},\nu) \coloneqq -\frac{1}{2}\log\det\left(\boldsymbol{I}+\alpha\boldsymbol{G}_{\boldsymbol{C}}\right) + \frac{\gamma}{2}\operatorname{tr}(\boldsymbol{G}_{\boldsymbol{C}}\boldsymbol{M}) - \operatorname{tr}(\boldsymbol{\Delta}\boldsymbol{G}_{\boldsymbol{C}}) + \frac{\nu}{2}(\operatorname{tr}(\boldsymbol{G}_{\boldsymbol{C}})-N), (16)$$

where scalar  $\nu$  and  $N \times N$  symmetric matrix  $\Delta$  are Lagrange multipliers.

The KKT conditions is:

=

$$\int -\frac{\alpha}{2} (\mathbf{I} + \alpha \mathbf{G}_{\star,\mathbf{C}})^{-1} + \frac{\gamma}{2} \mathbf{M} - \mathbf{\Delta}_{\star} + \frac{\nu_{\star}}{2} \mathbf{I} = \mathbf{0},$$
(17)

$$G_{\star,C} \succeq \mathbf{0},\tag{18}$$

$$\begin{cases} \operatorname{tr}(\boldsymbol{G}_{\star,\boldsymbol{C}}) = N, \\ \boldsymbol{\Delta}_{\star} \succeq \boldsymbol{0}, \end{cases}$$
(19)  
(20)

$$\Delta_* G_{*,C} = 0, \tag{21}$$

which are the sufficient and necessary condition for the global optimality of the solution  $G_{\star,C}$ .

From Eqs. (18),(20) and (21), we have that  $\Delta_{\star}G_{\star,C} - G_{\star,C}\Delta_{\star} = \Delta_{\star}G_{\star,C} - (\Delta_{\star}G_{\star,C})^{\top} = 0$ , implying that  $\Delta_{\star}$  and  $G_{\star,C}$  share eigenspace. By eigenvalue decomposition  $\Delta_{\star} = Q\Lambda_{\Delta_{\star}}Q^{\top}, G_{\star,C} = Q\Lambda_{G_{\star,C}}Q^{\top}$ , where  $\Lambda_{\Delta_{\star}}, \Lambda_{G_{\star,C}}$  are diagonal matrices, we have:

$$2 \cdot \boldsymbol{Q} \boldsymbol{\Lambda}_{\boldsymbol{\Delta}_{\star}} \boldsymbol{Q}^{\top} = -\alpha \boldsymbol{Q} (\boldsymbol{I} + \alpha \boldsymbol{\Lambda}_{\boldsymbol{G}_{\star,\boldsymbol{C}}})^{-1} \boldsymbol{Q}^{\top} + \gamma \boldsymbol{M} + \nu_{\star} \boldsymbol{I}$$
(22)

$$\Rightarrow \quad \gamma \boldsymbol{M} + \nu_{\star} \boldsymbol{I} = \boldsymbol{Q} \left( 2\boldsymbol{\Lambda}_{\boldsymbol{\Delta}_{\star}} + \alpha \left( \boldsymbol{I} + \alpha \boldsymbol{\Lambda}_{\boldsymbol{G}_{\star,\boldsymbol{C}}} \right)^{-1} \right) \boldsymbol{Q}^{\top}, \tag{23}$$

where the first equality is from Eq. (17). Since that  $2\Lambda_{\Delta_{\star}} + \alpha (I + \alpha \Lambda_{G_{\star,C}})^{-1}$  is a diagonal matrix,  $\gamma M + \nu_{\star} I$  can be diagonalized by Q. In other words, for  $\forall M \in \mathbb{S}_{N}^{+}$  in problem (15), M will share eigenspace with the corresponding optimal solution  $G_{\star,C}$ . Next, denote  $(Z_{\star}, C_{\star})$  as the optimal solution of problem (4),  $C := \{Z \mid ||Z||_{F}^{2} = N\}$  as the feasible set and  $f(\cdot, \cdot)$  as the objective function. Since that  $Z_{\star} = \arg \min_{Z \in C} f(Z, C_{\star})$ , otherwise contradicts with the optimality of  $(Z_{\star}, C_{\star})$ , we conclude that  $G_{\star}$  and  $M_{\star}$  share eigenspace, where  $M_{\star} := (I - C_{\star})(I - C_{\star})^{\top}$ .

**Theorem 2.** Suppose that G and M are aligned in the same eigenspaces and  $\gamma < \frac{1}{\lambda_{\max}(M)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$ , then we have that: a)  $\operatorname{rank}(\mathbf{Z}_{\star}) = \min\{d, N\}$ , and b) the singular values  $\sigma_{\mathbf{Z}_{\star}}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_M^{(i)} + \nu_{\star}} - \frac{1}{\alpha}}$  for all  $i = 1, \ldots, \min\{d, N\}$ , where  $\mathbf{Z}_{\star}$  and  $\nu_{\star}$  are the primal optimal solution and dual optimal solution, respectively.

*Proof.* Since  $\|\boldsymbol{Z} - \boldsymbol{Z}\boldsymbol{C}\|_F^2 = \operatorname{Tr}\left(\boldsymbol{Z}^\top \boldsymbol{Z} \left(\boldsymbol{I} - \boldsymbol{C}\right) \left(\boldsymbol{I} - \boldsymbol{C}\right)^\top\right)$  and  $\|\boldsymbol{Z}\|_F^2 = \operatorname{Tr}(\boldsymbol{Z}^\top \boldsymbol{Z})$ , problem (5) is equivalent to:

$$\min_{\boldsymbol{G}} \quad -\frac{1}{2} \log \det \left( \boldsymbol{I} + \alpha \boldsymbol{G} \right) + \frac{\gamma}{2} \operatorname{Tr}(\boldsymbol{G}\boldsymbol{M}) \\
\text{s.t.} \quad \operatorname{Tr}(\boldsymbol{G}) = N, \boldsymbol{G} \succeq \boldsymbol{0},$$
(24)

where  $\boldsymbol{G}\coloneqq \boldsymbol{Z}^{\top}\boldsymbol{Z}$  and  $\boldsymbol{M}\coloneqq (\boldsymbol{I}-\boldsymbol{C})(\boldsymbol{I}-\boldsymbol{C})^{\top}.$ 

Since that G and M have eigenspaces aligned, we can have G and M diagonalized simultaneously by an orthogonal matrix U, i.e.,  $G = U\Lambda_G U^{\top}$ ,  $M = U\Lambda_M U^{\top}$ . Therefore, problem (24) can be transformed into the eigenvalue optimization problem as follows:

$$\min_{\{\lambda_{G}^{(i)}\}_{i=1}^{\min\{d,N\}}} - \frac{1}{2} \sum_{i=1}^{\min\{d,N\}} \log(1 + \alpha \lambda_{G}^{(i)}) + \frac{\gamma}{2} \lambda_{M}^{(i)} \lambda_{G}^{(i)}$$
s.t.
$$\sum_{i=1}^{\min\{d,N\}} \lambda_{G}^{(i)} = N, \quad \lambda_{G}^{(i)} \ge 0, \quad \text{for all } i = 1, \dots, \min\{d,N\},$$
(25)

where  $\{\lambda_M^{(1)}, \dots, \lambda_M^{(\min\{d,N\})}\}\$  are the diagonal entries of  $\Lambda_M$  and  $\{\lambda_G^{(1)}, \dots, \lambda_G^{(\min\{d,N\})}\}\$  are the diagonal entries of  $\Lambda_G$ . Surprisingly, problem (25) is a convex optimization problem. Thus, the KKT condition is sufficient and necessary to guarantee for the global minimizer.

The Lagrangian of problem (25) is:

$$\mathcal{L}\Big(\{\lambda_{G}^{(i)}\}_{i=1}^{\min\{d,N\}}, \{\mu_{i}\}_{i=1}^{\min\{d,N\}}, \nu\Big) \coloneqq -\frac{1}{2} \sum_{i=1}^{\min\{d,N\}} \log(1 + \alpha \lambda_{G}^{(i)}) + \frac{\gamma}{2} \lambda_{M}^{(i)} \lambda_{G}^{(i)} - \mu_{i} \lambda_{G}^{(i)} + \frac{\nu}{2} \Big(\sum_{i=1}^{\min\{d,N\}} \lambda_{G}^{(i)} - N\Big), \quad (26)$$

where  $\mu_i \ge 0, i = 1, ..., \min\{d, N\}$  and  $\nu$  are the Lagrangian multipliers. The KKT conditions are as follows:

$$\begin{cases} \nabla_{\lambda_{G_{\star}}^{(i)}} \mathcal{L} = 0, \qquad \forall i = 1, \dots, \min\{d, N\}, \end{cases}$$
(27)

$$\begin{cases} {}^{(i)}_{G_{\star}} \ge 0, \qquad \qquad \forall i = 1, \dots, \min\{d, N\}, \end{cases}$$

$$(28)$$

$$\sum_{i=1}^{\min\{i,i+1\}} \lambda_{G_{\star}}^{(i)} = N, \tag{29}$$

$$\star \ge 0, \qquad \forall i = 1, \dots, \min\{d, N\}, \tag{30}$$

$$\left(\mu_{i\star}\lambda_{G_{\star}}^{(i)}=0,\qquad \forall i=1,\ldots,\min\{d,N\}.\right.$$
(31)

Then, the stationary condition in (27) is equivalent to:

$$\mu_{i\star} = \frac{1}{2} \Big( \nu_{\star} + \gamma \lambda_{\boldsymbol{M}}^{(i)} - \frac{\alpha}{1 + \alpha \lambda_{\boldsymbol{G}_{\star}}^{(i)}} \Big).$$
(32)

By using Eqs. (28) and (30)-(32), we come up with the following two cases:

$$\int \mu_{i\star} > 0 \Rightarrow \lambda_{G_{\star}}^{(i)} = 0, \frac{1}{\nu_{\star} + \gamma \lambda_{M}^{(i)}} - \frac{1}{\alpha} < 0,$$
(33)

$$\begin{pmatrix}
\mu_{i\star} = 0 \Rightarrow \lambda_{G_{\star}}^{(i)} > 0, \lambda_{G_{\star}}^{(i)} = \frac{1}{\nu_{\star} + \gamma \lambda_{M}^{(i)}} - \frac{1}{\alpha} > 0.
\end{cases}$$
(34)

From the above two cases, we conclude that:

$$\lambda_{\boldsymbol{G}_{\star}}^{(i)} = \max\left\{0, \frac{1}{\nu_{\star} + \gamma \lambda_{\boldsymbol{M}}^{(i)}} - \frac{1}{\alpha}\right\},\tag{35}$$

where  $\nu_{\star}$  satisfies:

$$\sum_{i=1}^{\min\{d,N\}} \max\left\{0, \frac{1}{\nu_{\star} + \gamma \lambda_{M}^{(i)}} - \frac{1}{\alpha}\right\} = N.$$
(36)

Given that  $\gamma < (\alpha - \nu_{\star})/\lambda_{\max}(\boldsymbol{M})$ , we have  $\frac{1}{\nu_{\star} + \gamma \lambda_{\boldsymbol{M}}^{(i)}} - \frac{1}{\alpha} > 0$  for all  $i = 1, \ldots, \min\{d, N\}$ . Therefore, for the optimal solution  $\boldsymbol{Z}_{\star}$  of problem (5), we conclude that:  $\operatorname{rank}(\boldsymbol{Z}_{\star}) = \min\{d, N\}$  and the singular values  $\sigma_{\boldsymbol{Z}_{\star}}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_{\boldsymbol{M}}^{(i)} + \nu_{\star}} - \frac{1}{\alpha}}$ , for all  $i = 1, \ldots, \min\{d, N\}$ .

Note that, the results we established just above rely on a condition  $\gamma \lambda_{\max}(M) < \alpha - \nu_{\star}$  where the  $\nu_{\star}$  is the optimal Lagrangian multiplier, which is set as a fixed value related to  $\alpha, \gamma$  and  $\lambda_{\max}(M)$ . Next, we will develop an upper bound for  $\nu_{\star}$  and justify the fact that we ensure  $\nu_{\star}$  to satisfy the condition  $\gamma \lambda_{\max}(M) < \alpha - \nu_{\star}$  by only adjusting the hyper-parameters  $\alpha$  and  $\gamma$ .

In Eq. (36), we can easily find an upper bound of  $\nu_{\star}$  as:

$$N = \sum_{i=1}^{\min\{d,N\}} \max\left\{0, \frac{1}{\nu_{\star} + \gamma\lambda_{\boldsymbol{M}}^{(i)}} - \frac{1}{\alpha}\right\} \le \frac{\min\{d,N\}}{\nu_{\star} + \gamma\lambda_{\min}(\boldsymbol{M})} - \frac{\min\{d,N\}}{\alpha}, \quad (37)$$

$$\Rightarrow \nu_{\star} \leq \frac{1}{\frac{N}{\min\{d,N\}} + \frac{1}{\alpha}} - \gamma \lambda_{\min}(\boldsymbol{M}), \tag{38}$$

Therefore, we can find a tighten bound between  $\alpha - \nu_{\star}$  and  $\gamma \lambda_{\max}(M)$  as:

$$\gamma \lambda_{\max}(\boldsymbol{M}) < \frac{\alpha^2}{\alpha + \min\left\{\frac{d}{N}, 1\right\}} < \frac{\alpha^2}{\alpha + \min\left\{\frac{d}{N}, 1\right\}} + \gamma \lambda_{\min}(\boldsymbol{M}) \le \alpha - \nu_{\star}, \tag{39}$$

which means that the condition of  $\gamma \lambda_{\max}(M) < \alpha - \nu_{\star}$  can be reformed as:

$$\gamma < \frac{1}{\lambda_{\max}(\boldsymbol{M})} \frac{\alpha^2}{\alpha + \min\left\{\frac{d}{N}, 1\right\}}$$
(40)

**Remark 3.** We notice that (25) is a reverse water-filling problem, where the water level is controlled by  $1/\alpha$ , as shown in Figure A.1. When G and M have eigenspaces aligned and  $\gamma < (\alpha - \nu_{\star})/\lambda_{\max}(M)$ , we have  $\operatorname{rank}(\mathbf{Z}_{\star}) = \min\{d, N\}$  and  $\lambda_{G_{\star}}^{(i)} \neq 0$  for all  $i \leq \min\{d, N\}$ . When  $\gamma \geq (\alpha - \nu_{\star})/\lambda_{\max}(M)$ , non-zero  $\lambda_{G}^{(i)}$  first disappears for the larger  $\lambda_{M}^{(i)}$ .



Figure A.1: **Illustration of the optimal solution for problem (25)**. The primal problem can be transformed into a classical reverse water-filling problem.

**Theorem 3.** Suppose that the sufficient conditions to prevent catastrophic feature collapse are satisfied. Without loss of generality, we further assume that the columns in matrix  $\mathbf{Z}$  are arranged into k blocks according to a certain  $N \times N$  permutation matrix  $\Gamma$ , i.e.,  $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \cdots, \mathbf{Z}_k]$ . Then the

condition for that PRO-DSC promotes the optimal solution  $(\mathbf{Z}_{\star}, \mathbf{C}_{\star})$  to be desired structure (i.e.,  $\mathbf{Z}_{\star}^{\top}\mathbf{Z}_{\star}$  and  $\mathbf{C}_{\star}$  are block-diagonal), is that  $\langle (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^{\top}, \mathbf{G} - \mathbf{G}^{*} \rangle \rightarrow 0$ , where

and  $G_{jj}$  is the block Gram matrix corresponding to  $Z_j$ .

*Proof.* We begin with the analysis to the first two terms of the loss function  $\tilde{\mathcal{L}} \coloneqq \mathcal{L}_1 + \gamma \mathcal{L}_2$ , where

$$\mathcal{L}_{1} \coloneqq -\frac{1}{2} \log \det \left( \boldsymbol{I} + \alpha (\boldsymbol{Z} \boldsymbol{\Gamma})^{\top} (\boldsymbol{Z} \boldsymbol{\Gamma}) \right) = -\frac{1}{2} \log \det (\boldsymbol{I} + \alpha \boldsymbol{G}),$$
  
$$\mathcal{L}_{2} \coloneqq \frac{1}{2} \| \boldsymbol{Z} \boldsymbol{\Gamma} - \boldsymbol{Z} \boldsymbol{\Gamma} \boldsymbol{\Gamma}^{\top} \boldsymbol{C} \boldsymbol{\Gamma} \|_{F}^{2} = \frac{1}{2} \| \boldsymbol{Z} - \boldsymbol{Z} \boldsymbol{C} \|_{F}^{2} = \frac{1}{2} \operatorname{Tr} \left( \boldsymbol{G} \left( \boldsymbol{I} - \boldsymbol{C} \right) \left( \boldsymbol{I} - \boldsymbol{C} \right)^{\top} \right),$$

since that  $\Gamma^{\top}\Gamma = \Gamma\Gamma^{\top} = I$ . Thus, we have:

$$\tilde{\mathcal{L}}(\boldsymbol{G},\boldsymbol{C}) = \frac{\gamma}{2} \operatorname{Tr} \left( \boldsymbol{G} \left( \boldsymbol{I} - \boldsymbol{C} \right) \left( \boldsymbol{I} - \boldsymbol{C} \right)^{\top} \right) - \frac{1}{2} \log \det(\boldsymbol{I} + \alpha \boldsymbol{G}), \tag{41}$$

which is a convex function with respect to (w.r.t) G and C, separately. By the property of convex function w.r.t. C, we have:

$$\begin{split} \tilde{\mathcal{L}}(\boldsymbol{G},\boldsymbol{C}) &\geq \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*) + \left\langle \nabla_{\boldsymbol{C}} \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*), \boldsymbol{C} - \boldsymbol{C}^* \right\rangle + \left\langle \frac{\gamma}{2} \left( \boldsymbol{I} - \boldsymbol{C} \right) \left( \boldsymbol{I} - \boldsymbol{C} \right)^\top, \boldsymbol{G} - \boldsymbol{G}^* \right\rangle \\ &= \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*) + \left\langle -\gamma \boldsymbol{G}^* (\boldsymbol{I} - \boldsymbol{C}^*), \boldsymbol{C} - \boldsymbol{C}^* \right\rangle + \left\langle \frac{\gamma}{2} \left( \boldsymbol{I} - \boldsymbol{C} \right) \left( \boldsymbol{I} - \boldsymbol{C} \right)^\top, \boldsymbol{G} - \boldsymbol{G}^* \right\rangle, \end{split}$$

where  $C^* = \text{Diag}(C_{11}, C_{22}, \cdots, C_{kk}) = \begin{bmatrix} C_{11} & & \\ & \ddots & \\ & & C_{kk} \end{bmatrix}$  with the blocks associating to the

partition of  $Z = [Z_1, Z_2, \dots, Z_k]$ . Since that  $\langle G^*(I - C^*), C - C^* \rangle = 0$  due to the complementary between  $G^*$  and  $I - C^*$ , we have:

$$\tilde{\mathcal{L}}(\boldsymbol{G},\boldsymbol{C}) \geq \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*) + \left\langle \frac{\gamma}{2} \left( \boldsymbol{I} - \boldsymbol{C} \right) \left( \boldsymbol{I} - \boldsymbol{C} \right)^{\top}, \boldsymbol{G} - \boldsymbol{G}^* \right\rangle.$$

It is easy to see that if  $\langle (\boldsymbol{I} - \boldsymbol{C}) (\boldsymbol{I} - \boldsymbol{C})^{\top}, \boldsymbol{G} - \boldsymbol{G}^* \rangle \rightarrow 0$ , then we will have:

$$\tilde{\mathcal{L}}(\boldsymbol{G},\boldsymbol{C}) \ge \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*), \tag{42}$$

where the equality holds only when  $G = G^*$ ,  $C = C^*$ . Furthermore, if the regularizer  $r(\cdot)$  satisfies the extended block diagonal condition as defined in (Lu et al., 2018), then we have that  $r(C) \ge r(C^*)$ , where the equality holds if and only if  $C = C^*$ . Therefore, we have:

$$\mathcal{L}(\boldsymbol{G},\boldsymbol{C}) = \tilde{\mathcal{L}}(\boldsymbol{G},\boldsymbol{C}) + \beta \cdot r(\boldsymbol{C}) \ge \tilde{\mathcal{L}}(\boldsymbol{G}^*,\boldsymbol{C}^*) + \beta \cdot r(\boldsymbol{C}^*) = \mathcal{L}(\boldsymbol{G}^*,\boldsymbol{C}^*).$$
(43)

Thus we conclude that minimizing the loss function  $\mathcal{L}(G, C) = \tilde{\mathcal{L}}(G, C) + \beta \cdot r(C)$  promotes the optimal solution  $(G_{\star}, C_{\star})$  to have block diagonal structure.

We note that the Gram matrix being block-diagonal, i.e.,  $G_{\star} = G^{\star}$ , implies that  $Z_{\star,j_1}^{\top} Z_{\star,j_2} = 0$  for all  $1 \leq j_1 < j_2 \leq k$ , which is corresponding to the subspaces associated to the blocks  $Z_{\star,j}$ 's are orthogonal to each other.

#### **B** EXPERIMENTAL SUPPLEMENTARY MATERIAL

#### **B.1** EXPERIMENTAL DETAILS

#### **B.1.1** SYNTHETIC DATA

As shown in Figure 5a (top row), data points are generated from two manifolds. The first manifold (colored in purple) is generated by sampling 100 data points from

$$\boldsymbol{x} = \begin{bmatrix} \cos\left(\frac{1}{5}\sin\left(5\varphi\right)\right)\cos\varphi\\ \cos\left(\frac{1}{5}\sin\left(5\varphi\right)\right)\sin\varphi\\ \sin\left(\frac{1}{5}\sin\left(5\varphi\right)\right) \end{bmatrix} + \boldsymbol{\epsilon}, \tag{44}$$

where  $\varphi$  is taken uniformly from  $[0, 2\pi]$  and  $\epsilon \sim \mathcal{N}(\mathbf{0}, 0.05\mathbf{I}_3)$  is the additive noise. The second manifold (colored in blue) is generated by sampling 100 data points from a Gaussian distribution  $\mathcal{N}([0, 0, 1]^{\top}, 0.05\mathbf{I}_3)$ . To further test more complicated cases, we generate the second manifold by sampling 50 data points from a Gaussian distribution  $\mathcal{N}([0, 0, 1]^{\top}, 0.05\mathbf{I}_3)$  and 50 data points from another Gaussian distribution  $\mathcal{N}([0, 0, -1]^{\top}, 0.05\mathbf{I}_3)$ , as shown in Figure 5a (bottom row).

In PRO-DSC, the learnable mappings  $h(\cdot; \Psi)$  and  $f(\cdot; \Theta)$  are implemented with two MLPs with Rectified Linear Units (ReLU) (Nair & Hinton, 2010) as the activation function. The hidden dimension and output dimension of the MLPs are set to 100 and 3, respectively. We train PRO-DSC with batch-size  $n_b = 200$ , learning rate  $\eta = 5 \times 10^{-3}$  for 1000 epochs. We set  $\gamma = 0.5, \beta = 1000$ , and  $\alpha = \frac{3}{0.1 \cdot 200}$ .

We use DSCNet (Ji et al., 2017) as the representative of the SEDSC methods. In Figure 5b, we set  $\gamma = 1$  for both cases, whereas in Figure 5c,  $\gamma$  is set to 5 and 100 for the two cases, respectively. The encoder and decoder of DSCNet are MLPs of two hidden layers, with the hidden dimensions being set to 100 and 3, respectively. We train DSCNet with batch-size  $n_b = 200$ , learning rate  $\eta = 1 \times 10^{-4}$  for 1000 epochs.

#### **B.1.2 REAL-WORLD DATASETS**

**Datasets description.** CIFAR-10 and CIFAR-100 are classic image datasets consisting of 50,000 images for training and 10,000 images for testing. They are split into 10 and 100 classes, respectively. CIFAR-20 shares the same images with CIFAR-100 while taking 20 super-classes as labels. ImageNet-Dogs consists of 19,500 images of 15 different dog species. Tiny-ImageNet consists of 100,000 images from 200 different classes. ImageNet-1k is the superset of the two datasets, containing more than 1,280,000 real world images from 1000 classes. For all the datasets except for ImageNet-Dogs, we train the network to implement PRO-DSC on the train set and test it on the test set to validate the generalization of the learned model. For ImageNet-Dogs dataset which does not have a test set, we train the network to implement PRO-DSC on the train set and report the clustering performance on the training set. For a direct comparison, we conclude the basic information of these datasets in Table B.1.

To leverage the CLIP features for training, the input images are first resized to 224 with respect to the smaller edge, then center-cropped to  $224 \times 224$  and fed into the CLIP pre-trained image encoder to obtain fixed features.<sup>11</sup> The subsequent training of PRO-DSC takes the extracted features as input, instead of loading the entire CLIP pre-trained model.

Network architecture and hyper-parameters. The learnable mappings  $h(\cdot; \Psi)$  and  $f(\cdot; \Theta)$  are two fully-connected layers with the same output dimension *d*. Following (Chu et al., 2024), for the experiments on real-world data, we stack a pre-feature layer before the learnable mappings, which is composed of two fully-connected layers with ReLU and batch-norm (Ioffe & Szegedy, 2015).

We train the network by the SGD optimizer with the learning rate set to  $\eta = 10^{-4}$ , and the weight decay parameters of  $f(\cdot; \Theta)$  and  $h(\cdot; \Psi)$  are set to  $10^{-4}$  and  $5 \times 10^{-3}$ , respectively.

 $<sup>^{11}</sup>We$  use the ViT L/14 pre-trained model provided by <code>https://github.com/openai/CLIP</code> for 768-dimensional features.

Datasets	# Train	# Test	# Classes
CIFAR-10	50,000	10,000	10
CIFAR-20	50,000	10,000	20
CIFAR-100	50,000	10,000	100
ImageNet-Dogs	19,500	N/A	15
TinyImageNet	100,000	10,000	200
ImageNet	1.281.167	50.000	1000

Table B.1: **Basic statistical information of datasets.** We summarize the information in terms of the data with both the train and test split, as well as the number of classes involved.

Following by (Chu et al., 2024), we warm up training  $f(\cdot; \Theta)$  by diversifying the features with  $\mathcal{L}_1 = -\log \det(I + \alpha Z_{\Theta}^{\top} Z_{\Theta})$  for a few iterations and share the weights to  $h(\cdot; \Psi)$ . We set  $\alpha = \frac{d}{0.1 \cdot n_b}$  for all the experiments. We summarize the hyper-parameters for training the network to implement PRO-DSC in Table B.2.

Table B.2: Hyper-parameters configuration for training the network to implement PRO-DSC with CLIP pre-trained features. Here  $\eta$  is the learning rate,  $d_{pre}$  is the hidden and output dimension of pre-feature layer, m is the output dimension of h and f,  $n_b$  is the batch size for training, and "# warm-up" is the number of iterations of warm-up stage.

	$\eta$	$d_{pre}$	d	#epochs	$n_b$	#warm-up	$\gamma$	β
CIFAR-10	$1 \times 10^{-4}$	4096	128	10	1024	200	$300/n_{b}$	600
CIFAR-20	$1 \times 10^{-4}$	4096	256	50	1500	0	$600/n_{b}$	300
CIFAR-100	$1 \times 10^{-4}$	4096	128	100	1500	200	$150/n_{b}$	500
ImageNet-Dogs	$1 \times 10^{-4}$	4096	128	200	1024	0	$300/n_{b}$	400
TinyImageNet	$1 \times 10^{-4}$	4096	256	100	1500	0	$200/n_{b}$	400
ImageNet	$1 \times 10^{-4}$	4096	1024	100	2048	2000	$800/n_{b}$	400
MNIST	$1 \times 10^{-4}$	4096	128	100	1024	200	$700/n_{b}$	400
F-MNIST	$1 \times 10^{-4}$	1024	128	200	1024	400	$50/n_{b}$	100
Flowers	$1 \times 10^{-4}$	1024	256	200	1024	200	$400/n_{b}$	200

**Running other algorithms.** Since that k-means (MacQueen, 1967), spectral clustering (Shi & Malik, 2000), EnSC (You et al., 2016a), SSCOMP (You et al., 2016b), and DSCNet (Ji et al., 2017) are based on transductive learning, we evaluate these models directly on the test set for all the experiments.

- For EnSC, we tune the hyper-parameter  $\gamma \in \{1, 2, 5, 10, 20, 50, 100, 200, 400, 800, 1600, 3200\}$ and the hyper-parameter  $\tau$  in  $\tau \| \cdot \|_1 + \frac{1-\tau}{2} \| \cdot \|_2^2$  to balance the  $\ell_1$  and  $\ell_2$  norms in  $\{0.9, 0.95, 1\}$  and report the best clustering result.
- For SSCOMP, we tune the hyper-parameter to control the sparsity  $k_{\text{max}} \in \{1, 2, 5, 10, 20, 50, 100, 200\}$  and the residual  $\epsilon \in \{10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}\}$  and report the best clustering result.
- To apply DSCNet to the CLIP features, we use MLPs with two hidden layers to replace the convolutional encoder and decoder. The hidden dimension of the MLPs are set to 128. We tune the balancing hyper-parameters  $\gamma \in \{1, 2, 3, 4\}$  and  $\beta \in \{1, 5, 25, 50, 75, 100\}$ and train the model for 100 epochs with learning rate  $\eta = 1 \times 10^{-4}$  and batch-size  $n_b$ equivalent to number of samples in the test data set.
- As the performance of CPP is evaluated by averaging the ACC and NMI metrics tested on each batch, we reproduce the results by their open-source implementation and report the results on the entire test set. The authors provide two implementations (see <a href="https://github.com/LeslieTrue/CPP/blob/main/main\_main.py">https://github.com/LeslieTrue/CPP/blob/main/ main.py</a> and <a href="https://github.com/LeslieTrue/CPP/blob/main/main\_efficient.py">https://github.com/LeslieTrue/CPP/blob/main/ main.py</a> and <a href="https://github.com/LeslieTrue/CPP/blob/main/main\_efficient.py">https://github.com/LeslieTrue/CPP/blob/main/main\_ efficient.py</a>), where one optimizes the cluster head and the feature head separately and the other shares weights between the two heads. In this paper, we test both cases and report the better results.

- For k-means and spectral clustering (including when spectral clustering is used as the final step in subspace clustering), we repeat the clustering 10 times with different random initializations (by setting  $n_{init} = 10$  in scikit-learn) and report the best results.
- For SENet, SCAN and EDESC, we adjust the hyper-parameters and repeat experiments for three times, with only the best results are reported.

#### **B.2** EMPIRICAL VALIDATION ON THEORETICAL RESULTS

**Empirical Validation on Theorem 1.** To validate Theorem 1 empirically, we conduct experiments on CIFAR-100 with the same training configurations as described in Section B.1.2 but change the training period to 1000 epochs. For each epoch, we compute  $G_b = Z_b^{\top} Z_b$  and  $M_b = (I - C_b)(I - C_b)^{\top}$  with the learned representations  $Z_b$  and self-expressive matrix  $C_b$  in mini-batch of size  $n_b$ after the last iteration of different epoch. Then, to quantify the eigenspace alignment of  $G_b$  and  $M_b$ , we directly plot the alignment error which is computed via the Frobenius norm of the commutator  $L := ||G_b M_b - M_b G_b||_F$  in mini-batch of size  $n_b$  during the training period in Figure 1a. We also show the standard deviation with shaded region after repeating the experiments for 5 random seeds. As can be read, the alignment error decreases monotonically during the training period, implying that the eigenspaces are progressively aligned. Moreover, we find the eigenvector  $\{u_j\}$  of  $M_b$  by eigenvalue decomposition, where  $u_j$  denotes the *j*-th eigenvector which are sorted according to the eigenvalues in the ascending order, and calculate the normalized correlation coefficient which is defined as  $\langle u_j, G_b u_j/||G_b u_j||_2 \rangle$ . Note that when the eigenspace alignment holds, one can verify that:

$$\langle \boldsymbol{u}_j, \frac{\boldsymbol{G}_b \boldsymbol{u}_j}{\|\boldsymbol{G}_b \boldsymbol{u}_j\|_2} \rangle = \begin{cases} 1, & \lambda_{\boldsymbol{G}_b}^{(j)} \neq 0\\ 0, & \lambda_{\boldsymbol{G}_b}^{(j)} = 0 \end{cases} \quad \text{for all } j = 1, 2, \dots, n_b.$$

$$\tag{45}$$

As shown in Figure 1b, the normalized correlation curves associated to the first d = 128 eigenvectors converge to 1, whereas the rest converge to 0, implying the progressively alignment between  $G_b$  and  $M_b$ .

Empirical Validation on Theorem 2. To verify Theorem 2, we conduct experiments on CIFAR-10 and CIFAR-100. The experimental setup keeps the same as described in Section B.1.2. In each epoch, we compute  $G_b = Z_b^{\top} Z_b$  and  $M_b = (I - C_b)(I - C_b)^{\top}$  from  $Z_b$  and  $C_b$  in mini-batch after the last iteration, respectively, and then find the eigenvalues of  $G_b$  and  $M_b$ . We display the eigenvalue curves in Figure 1c and 1d, respectively. To enhance the clarity of the visualization, the eigenvalues of  $G_b$  and  $M_b$  are sorted in descending and ascending order, respectively. As can be observed, there are min  $\{d, n_b\} = 128$  non-zero eigenvalues in  $G_b$ , approximately being inversely proportional to the smallest 128 eigenvalues of  $M_b$ . This results empirically demonstrate that rank $(Z_{\star}) = \min\{d, N\}$  and  $\lambda_{G_{\star}}^{(i)} = \frac{1}{\gamma \lambda_{M}^{(i)} + \nu_{\star}} - \frac{1}{\alpha}$  for minimizers.

Furthermore, to verify the sufficient condition of PRO-DSC to prevent feature space collapse, we conduct experiments on CIFAR-10 and CIFAR-100 with varying  $\alpha$  and  $\gamma$ , keeping all the other hyper-parameters consistent with Table B.2. As can be read in Figure 2, Theorem 2 is verified since that  $\gamma < \frac{1}{\lambda_{\max}(M_b)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$  yields satisfactory clustering accuracy (ACC%) and subspace-preserving representation error (SRE%). The satisfactory ACC and SRE confirm that PRO-DSC avoids catastrophic collapse when  $\gamma < \frac{1}{\lambda_{\max}(M)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$  holds. When  $\gamma \geq \frac{1}{\lambda_{\max}(M_b)} \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$ , PRO-DSC yields significantly worse ACC and SRE. There is a phase tran-

sition phenomenon that corresponds to the sufficient condition to prevent collapse.12

**Empirical Validation on Theorem 3.** To intuitively visualize the structured representations learned by PRO-DSC, we visualize the Gram matrices  $|Z^{\top}Z|$  and Principal Component Analysis (PCA) results for both CLIP features and learned representations on CIFAR-10. The experimental setup also keeps the same as described in Section B.1.2.

The Gram matrix shows the similarities between representations within the same class (indicated by in-block diagonal values) and across different classes (indicated by off-block diagonal values).

<sup>12</sup>In experiments, we estimate that  $\lambda_{\max}(M_b) = 1$  and thus the condition reduces to  $\gamma < \frac{\alpha^2}{\alpha + \min\{\frac{d}{N}, 1\}}$ .

Moreover, we display the dimensionality reduction results via PCA for the CLIP features and the learned representation of samples from three categories in CIFAR-10. We use PCA for dimensionality reduction as it performs a linear projection, well preserving the underlying structure.

As can be observed in Figure 3, the CLIP features from three classes approximately lie on different subspaces. Despite of the structured nature of the features, the underlying subspaces are not orthogonal. In the Gram matrix of the CLIP feature, the average similarity between features from different classes is greater than 0.6, resulting in an unclear block diagonal structure. After training with PRO-DSC, the spanned subspaces of the learned representations become orthogonal.<sup>13</sup> Additionally, the off-block diagonal values of the Gram matrix decrease significantly, revealing a clear block diagonal structure. These visualization results qualitatively verify that PRO-DSC aligns the representations with a union of orthogonal subspaces.<sup>14</sup>

#### **B.3 MORE EXPERIMENTAL RESULTS**

#### B.3.1 MORE RESULTS OF PRO-DSC ON SYNTHETIC DATA

To explore the learning ability of our PRO-DSC, we prepare experiments on synthetic data with adding an additional subspace, as presented in Figure B.1.

In case 1, we sample 100 points from Gaussian distribution  $\boldsymbol{x} \sim \mathcal{N}([\frac{1}{\sqrt{2}}, 0, \frac{1}{\sqrt{2}}]^{\top}, 0.05\boldsymbol{I}_3)$  and 100 points from  $\boldsymbol{x} \sim \mathcal{N}([-\frac{1}{\sqrt{2}}, 0, \frac{1}{\sqrt{2}}]^{\top}, 0.05\boldsymbol{I}_3)$ , respectively. We train PRO-DSC with batch-size  $n_b = 300$ , learning rate  $\eta = 5 \times 10^{-3}$  for 5000 epochs and set  $\gamma = 1.3, \beta = 500, \alpha = \frac{3}{0.1\cdot 300}$ . We observe that our PRO-DSC successfully eliminates the nonlinearity in representations and maximally separates the different subspaces.

In case 2, we add a vertical curve

$$\boldsymbol{x} = \begin{bmatrix} \cos\left(\frac{1}{5}\sin\left(5\varphi\right)\right)\cos\varphi\\ \sin\left(\frac{1}{5}\cos\left(5\varphi\right)\right)\\ \cos\left(\frac{1}{5}\sin\left(5\varphi\right)\right)\sin\varphi \end{bmatrix} + \boldsymbol{\epsilon}, \tag{46}$$

from which 100 points are sampled, where  $\epsilon \sim \mathcal{N}(\mathbf{0}, 0.05I_3)$ . We use  $\sin(\frac{1}{5}\cos(5\varphi))$  to avoid overlap in the intersection of the two curves. We train PRO-DSC with batch-size  $n_b = 200$ , learning rate  $\eta = 5 \times 10^{-3}$  for 8000 epochs and set  $\gamma = 0.5, \beta = 500, \alpha = \frac{3}{0.1 \cdot 200}$ . We observe that PRO-DSC finds difficulties to learn representations of data which are located at the intersections of the subspaces. However, for those data points which are away from the intersections are linearized well.



Figure B.1: Additional results on synthetic data.

#### **B.3.2 EXPERIMENTS WITH BYOL PRE-TRAINING**

To validate the effectiveness of our PRO-DSC without using CLIP features, we conduct a fair comparison with existing deep clustering approaches. Similar to most deep clustering algorithms, we

<sup>&</sup>lt;sup>13</sup>The dimension of each subspace is much greater than one (see Figure B.4). The 1-dimensional subspaces observed in the PCA results are a consequence of dimensionality reduction.

<sup>&</sup>lt;sup>14</sup>Please refer to Figure B.3 and B.7 for the results on other datasets and the visualization of the bases of each subspace.

divide the training process into two steps. We begin with pre-training the parameters of the backbone with BYOL (Grill et al., 2020). Then, we leverage the parameters pre-trained in the first stage and fine-tune the model by the proposed PRO-DSC loss function. Specifically, we set the learning rate  $\eta = 0.05$  and the batch size  $n_b = 256$ . The output feature dimension *d* is consistent with the setting for training with the CLIP features. Following (Li et al., 2021; Huang et al., 2023), We use ResNet-18 as the backbone for the experiments on CIFAR-10 and CIFAR-20, and use ResNet-34 as the backbone for the experiments on other datasets, and use a convolution filter of size  $3 \times 3$ and stride 1 to replace the first convolution filter. We use the commonly used data augmentations methods to the input images, which are listed as follows:

```
transforms.RandomResizedCrop(size=img_size, scale=(0.08, 1)),
transforms.RandomHorizontalFlip(),
transforms.RandomApply([transforms.ColorJitter(0.4, 0.4, 0.2,
0.1)], p=0.8),
transforms.RandomGrayscale(p=0.2),
transforms.RandomApply([transforms.GaussianBlur(kernel_size=23,
sigma=(0.1, 2.0))], p=1.0).
```

When re-implementing other baselines, we use the code provided by the respective authors and report the best performance after fine-tuning the hyper-parameters.

We report the clustering results based on BYOL pre-training in Table B.3. As can be read from Table B.3, our PRO-DSC outperforms all the deep clustering baselines, including CC (Li et al., 2021), GCC (Zhong et al., 2021), NNM (Dang et al., 2021), SCAN (Van Gansbeke et al., 2020), NMCE (Li et al., 2022), IMC-SwAV (Ntelemis et al., 2022), and MLC (Ding et al., 2023).

Table B.3: **Clustering performance comparison on BYOL pre-training.** The best results are in bold font and the second best results are underlined. Performance marked with "\*" is based on our re-implementation.

Mathad	CIFA	R-10	CIFA	R-20	CIFA	AR-100	TinyImg	gNet-200	ImgNet	Dogs-15
Wiethou	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
k-means	22.9	8.7	13.0	8.4	9.2	23.0	2.5	6.5	10.5	5.5
SC	24.7	10.3	13.6	9.0	7.0	17.0	2.2	6.3	11.1	3.8
CC	79.0	70.5	42.9	43.1	26.9*	48.1*	14.0	34.0	42.9	44.5
GCC	85.6	76.4	47.2	47.2	28.2*	49.9*	13.8	34.7	52.6	49.0
NNM	84.3	74.8	47.7	48.4	41.2	55.1	-	-	31.1*	34.3*
SCAN	88.3	79.7	50.7	48.6	34.3	55.7	-	-	29.6*	30.3*
NMCE	89.1	81.2	53.1	52.4	40.0*	53.9*	21.6*	40.0*	39.8	39.3
IMC-SwAV	89.7	81.8	51.9	52.7	45.1	<u>67.5</u>	28.2	52.6	-	-
MLC	<u>92.2</u>	85.5	58.3	<u>59.6</u>	<u>49.4</u>	68.3	<u>28.7</u> *	<u>52.2</u> *	<u>71.0</u> *	<u>68.3</u> *
Our PRO-DSC	$93.0 \pm 0.6$	$86.5 \pm 0.2$	$58.3 \pm 0.9$	$60.1 \pm 0.6$	$56.3 \pm 0.6$	$66.7.0 \pm 1.0$	$31.1 \pm 0.3$	$46.0 \pm 1.0$	$74.1 \pm 0.5$	$69.5 \pm 0.6$

#### B.3.3 MORE EXPERIMENTS ON CLIP, DINO AND MAE PRE-TRAINED FEATURES

**Clustering on learned representations.** To quantitatively validate the effectiveness of the structured representations learned by PRO-DSC, we illustrate the clustering accuracy of representations learned by various algorithms in Figure 6. Here, to compared with the representations learned from SEDSC methods, we additionally conduct experiments on DSCNet (Ji et al., 2017) and report the performance in Table B.4. To apply DSCNet on CLIP features, we use MLPs with two hidden layers to replace the stacked convolutional encoder and decoder. As demonstrated in Sec. B.1, we report the best clustering results after the tuning of hyper-parameters. As analyzed in (Haeffele et al., 2021) and Section 2.1, DSCNet overly compresses the representations and yields unsatisfactory clustering results.

**Out of domain datasets.** We evaluate the capability to refine features by training PRO-DSC with pre-trained CLIP features on out-of-domain datasets, namely, MNIST (Deng, 2012), Fashion MNIST (Xiao et al., 2017) and Oxford flowers (Nilsback & Zisserman, 2008). As shown in Table B.5, CPP (Chu et al., 2024) refines the CLIP features and yields better clustering performance comparing with spectral clustering (Shi & Malik, 2000) and EnSC (You et al., 2016a). Our PRO-DSC further demonstrates the best performance on all benchmarks, validating its effectiveness in refining input features.

	CI	CIFAR-10		CIFAR-100		CIFAR-20			TinyImgNet-200			
	k-means	SC	EnSC	k-means	SC	EnSC	k-means	SC	EnSC	k-means	SC	EnSC
CLIP	74.7	70.2	95.4	52.8	66.4	67.0	46.9	49.2	60.8	54.1	62.8	64.5
SEDSC	16.4	18.9	16.9	5.4	4.9	5.3	11.7	10.6	12.8	5.7	3.9	7.2
CPP	71.3	70.3	95.6	75.3	75.0	77.5	55.5	43.6	58.3	62.1	58.0	67.0
PRO-DSC	93.4	92.1	<u>95.5</u>	76.5	75.2	77.6	66.0	59.7	<u>60.0</u>	67.6	67.0	69.5

Table B.4: Clustering accuracy of CLIP features and learned representations. We apply *k*-means, spectral clustering, and EnSC to cluster the representations.

Table B.5: Experiments on out-of-domain datasets.

Mathada	MN	IST	F-MI	NIST	Flowers	
Methods	ACC	NMI	ACC	NMI	ACC	NMI
Spectral Clustering (Shi & Malik, 2000)	74.5	67.0	64.3	56.8	85.6	94.6
EnSC (You et al., 2016a)	91.0	85.3	69.1	65.1	90.0	95.9
CPP (Chu et al., 2024)	<u>95.7</u>	<u>90.4</u>	<u>70.9</u>	<u>68.8</u>	<u>91.3</u>	<u>96.4</u>
PRO-DSC	96.1	90.9	71.3	70.3	92.0	97.4

**Experiments on block diagonal regularizer with different** k. To test the robustness of block diagonal regularizer  $||A||_{\mathbb{F}}$  to different k, we vary k and report the clustering performance in Table B.6. As illustrated, k does not necessarily equal to the number of clusters. There exists an interval within which the regularizer works effectively.

Table B.6: Clustering performance with different k in block diagonal regularizer.

	k	2	5	10	15	20	25	30
CIFAR-10	ACC NMI	97.2 93.2	97.2 93.2	97.4 93.5	96.3 92.0	96.3 92.0	95.4 90.7	94.0 88.6
	k	25	50	75	100	125	150	200
CIFAR-100	ACC NMI	74.3 80.9	76.7 82.3	78.1 <b>83.2</b>	78.2 82.9	78.9 83.2	76.4 82.2	74.8 81.5

But if k is significantly smaller than the number of clusters, the effect of block diagonal regularizer will be subtle. Therefore, the performance of PRO-DSC will be similar to that of PRO-DSC without a regularizer (see ablation studies in Section 3). In contrary, if k is significantly larger than the number of clusters, over-segmentation will occur to the affinity matrix, which has negative impact on the subsequent clustering performance.

**Clustering on ImageNet-1k with DINO and MAE.** To test the performance of PRO-DSC based on more pre-trained features other than CLIP (Radford et al., 2021), we further conduct experiments on ImageNet-1k (Deng et al., 2009) pre-trained by DINO (Caron et al., 2021) and MAE (He et al., 2022) (see Table B.7).

DINO and MAE are pre-trained on ImageNet-1k without leveraging external training data, thus their performance on PRO-DSC is lower than CLIP. Similar to the observations in CPP (Chu et al., 2024), DINO initializes PRO-DSC well, yet MAE fails, which is attributed to the fact that features from MAE prefer fine-tuning with labels, while they are less suitable for learning inter-cluster discriminative representations (Oquab et al., 2024). We further extract features from the validation set of ImageNet-1k and visualize through *t*-SNE (Van der Maaten & Hinton, 2008) to validate the hypothesis (see Figure B.2).

**B.3.4** EXPERIMENTS WITHOUT USING PRE-TRAINED MODELS

**Experiments on Reuters and UCI HAR.** During the rebuttal, we conducted extra experiments on datasets Reuters and UCI HAR. The dataset Reuters-10k consists of four text classes, containing 10,000 samples of 2,000 dimension. The UCI HAR is a time-series dataset, consisting of six classes, 10,299 samples of 561 dimension. We take EDESC (Cai et al., 2022) as the baseline method for

Method	Backbone	PRO ACC	-DSC NMI	k-m ACC	eans NMI
MAE (He et al., 2022) DINO (Caron et al., 2021) DINO (Caron et al., 2021)	ViT L/16 ViT B/16 ViT B/8	9.0 57.3 59.7	49.1 79.3 80.8	9.4 52.2 <b>54.6</b>	49.3 79.2 <b>80.5</b>
CLIP (Radford et al., 2021)	ViT L/14	65.1	83.6	52.5	79.7

Table B.7: Clustering Performance of PRO-DSC based on DINO and CLIP pre-trained features on ImageNet-1k.



Figure B.2: The *t*-SNE visualization of CLIP and MAE features on the validation set of ImageNet-1k.

deep subspace clustering on Reuters-10k, and take N2D (McConville et al., 2021) and FCMI (Zeng et al., 2023) as the baseline methods for UCI HAR, in which the results are directly cited from the respective papers. We conducted experiments with PRO-DSC on Reuters and UCI HAR following the same protocol for data processing as the baseline methods. We train and test PRO-DSC on the entire dataset and report the results over 10 trials. Experimental results are provided in Table B.8. The hyper-parameters used for PRO-DSC is listed in Table B.9.

Detect	REUTE	RS-10k	UCI	HAR
Dataset	ACC	NMI	ACC	NMI
k-means (MacQueen, 1967)	52.4	31.2	59.9	58.8
SC (Shi & Malik, 2000)	40.2	37.5	53.8	74.1
AE (Bengio et al., 2006)	59.7	32.3	66.3	60.7
VAE (Kingma & Welling, 2014)	62.5	32.9	-	-
JULE (Yang et al., 2016)	62.6	40.5	-	-
DEC (Xie et al., 2016)	75.6	<u>68.3</u>	57.1	65.5
DSEC (Chang et al., 2018)	78.3	70.8	-	-
EDESC (Cai et al., 2022)	82.5	61.1	-	-
DFDC (Zhang & Davidson, 2021)	-	-	86.2	84.5
N2D (McConville et al., 2021)	-	-	82.8	71.7
FCMI (Zeng et al., 2023)	-	-	88.2	80.7
PRO-DSC	$\textbf{85.7} \pm 1.3$	$64.6\pm1.3$	$\underline{87.1}\pm0.4$	$\underline{80.9}\pm1.2$

Table B.8: Experimental Results on Datasets Reuters and UCI HAR with 10 trials. The results of other methods are cited from the respective papers.

Dataset	$\eta$	$d_{pre}$	d	#epochs	$n_b$	#warm-up	$\gamma$	$\beta$
REUTERS-10k	$10^{-4}$	1024	128	100	1024	50	50	200
UCI HAR	$10^{-4}$	1024	128	100	2048	20	100	300
EYale-B	$10^{-4}$	1080	256	10000	2432	100	200	50
ORL	$10^{-4}$	80	64	5000	400	100	75	10
COIL-100	$10^{-4}$	12800	100	10000	7200	100	200	100

Table B.9: Configuration of hyper-parameters for experiments on Reuters, UCI HAR, EYale-B, ORL and COIL-100.

**Comparison to AGCSC and SAGSC on Extended Yale B, ORL, and COIL-100.** During the rebuttal, we conducted more experiments on two state-of-the-art subspace clustering methods AGCSC (Wei et al., 2023) and ARSSC (Wang et al., 2023a). Since that both of the two methods cannot handle the datasets used for evaluating our PRO-DSC, we conducted experiments on the datasets: Extended Yale B (EYaleB), ORL, and COIL-100. We set the architecture of pre-feature layer in PRO-DSC as the same to the encoder of DSCNet (Ji et al., 2017). The hyper-parameters configuration for training PRO-DSC is summarized in Table B.9. We repeated experiments for 10 trails and report the average with standard deviation in Table B.10. As baseline methods, we use EnSC (You et al., 2016a), SSCOMP (You et al., 2016b), S<sup>3</sup>COMP (Chen et al., 2020b), DSCNet, DSSC (Lim et al., 2020) and DELVE Zhao et al. (2024). The results of these methods, except for S<sup>3</sup>COMP and DELVE, are directly cited them from DSSC (Lim et al., 2020), and the results of S<sup>3</sup>COMP and DELVE are cited from their own papers.

- Comparison to AGCSC. Our method surpasses AGCSC on the Extended Yale B dataset and achieves comparable results on the ORL dataset. However, AGCSC cannot yield the result on COIL-100 in 24 hours.
- Comparison to ARSSC. ARSSC employs three different non-convex regularizers:  $\ell_{\gamma}$  norm Penalty (LP), Log-Sum Penalty (LSP), and Minimax Concave Penalty (MCP). While ARSSC-MCP performs the best on Extended Yale B, our PRO-DSC outperforms ARSSC-MCP on ORL. While AGCSC performs the best on ORL, but it yields inferior results on Extended Yale B and it cannot yield the results on COIL-100 in 24 hours. Thus, we did not report the results of AGCSC on COIL-100 and marked it as Out of Time (OOT). Our PRO-DSC performs the second best results on Extended Yale B, ORL and the best results on COIL-100. Since that we have not found the open-source code for ARSSC, we are unable to have their results on COIL-100. This comparison also confirms the scalablity of our PRO-DSC which is due to the re-parametrization (similar to SENet).

	EYale-B		ORL		COIL-100	
	ACC	NMI	ACC	NMI	ACC	NMI
EnSC	65.2	73.4	77.4	90.3	68.0	90.1
SSCOMP	78.0	84.4	66.4	83.2	31.3	58.8
S <sup>3</sup> COMP-C (Chen et al., 2020b)	87.4	-	-	-	78.9	-
DSCNet	69.1	74.6	75.8	87.8	49.3	75.2
DELVE (Zhao et al., 2024)	89.8	90.1	-	-	79.0	93.9
J-DSSC (Lim et al., 2020)	92.4	95.2	78.5	90.6	79.6	94.3
A-DSSC (Lim et al., 2020)	91.7	94.7	79.0	91.0	82.4	94.6
AGCSC (Wei et al., 2023)	92.3	94.0	86.3	92.8	OOT	OOT
ARSSC-LP (Wang et al., 2023a)	95.7	-	75.5	-	-	-
ARSSC-LSP (Wang et al., 2023a)	95.9	-	71.3	-	-	-
ARSSC-MCP (Wang et al., 2023a)	99.3	-	72.0	-	-	-
PRO-DSC	$\underline{96.0}\pm0.3$	$\textbf{95.7}\pm0.8$	$\underline{83.2}\pm2.2$	$\underline{92.7}\pm0.6$	$\textbf{82.8}\pm0.9$	$\textbf{95.0}\pm0.6$

Table B.10: Experiments on Extended Yale B, ORL and COIL-100.

#### **B.4 MORE VISUALIZATION RESULTS**

Gram matrices and PCA visualizations. To qualitatively validate that PRO-DSC learns representations aligning with a union-of-orthogonal-subspaces distribution, we visualize the Gram matrices and PCA dimension reduction results of CLIP features and learned representations from PRO-DSC for each dataset. As shown in Figure B.3, the off-block diagonal values decrease significantly, implying the orthogonality between representations from different classes. The orthogonal between subspaces can also be observed from the PCA dimension reduction results.

**Singular values visualization.** To show the intrinsic dimension of CLIP features and the representations of PRO-DSC, We plot the singular values of CLIP features and PRO-DSC's representations in Figure B.4. Specifically, the singular values of features from all the samples are illustrated on the left and the singular values of features within each class are illustrated on the middle and right. As can be seen, the singular values of PRO-DSC decrease much slower than that of CLIP, implying that the features of PRO-DSC enjoy a higher intrinsic dimension and more isotropic structure in the ambient space.

**Learning curves.** We plot the learning curves with respect to loss values and performance of PRO-DSC on CIFAR-100, CIFAR-20 and ImageNet-1k in Figure B.5a, Figure B.5b and Figure B.5c, respectively. Recall that  $\mathcal{L}_1 := -\frac{1}{2} \log \det \left( I + \alpha Z_{\Theta}^{\top} Z_{\Theta} \right), \mathcal{L}_2 := \frac{1}{2} ||Z_{\Theta} - Z_{\Theta} C_{\Psi}||_F^2$ , and  $\mathcal{L}_3 := ||A_{\Psi}||_{E}$ . Since  $\mathcal{L}_1$  is the only loss function used in the warm-up stage, we plot all the curves starting from the iteration when warm-up ends. As illustrated, the clustering performance of PRO-DSC steadily increase as the loss values gradually decrease, which shows the effectiveness of the proposed loss functions in PRO-DSC.

*t*-**SNE visualization of learned representations.** We visualize the CLIP features and cluster representations learned by PRO-DSC leveraging *t*-SNE (Van der Maaten & Hinton, 2008) in Figure B.6. As illustrated, the learned cluster representations are significantly more compact compared with the CLIP features, which contributes to the improved clustering performance.

**Subspace visualization.** We visualize the principal components of subspaces learned by PRO-DSC in Figure B.7. For each cluster in the dataset, we apply Principal Component Analysis (PCA) to the learned representations. We select the top eight principal components to represent the learned subspaces. Then, for each principal component, we display eight images whose representations are most closely aligned with that principal component.

Interestingly, we can observe specific semantic meanings from the principal components learned by PRO-DSC. For instance, the third row of Figure B.7a consists of stealth fighters, whereas the fifth row shows airliners. The second row of Figure B.7c consists of birds standing and resting, while the sixth row shows flying eagles. While Figure B.7j consists of all kinds of trucks, the first row shows fire trucks.

# C LIMITATIONS AND FAILURE CASES

**Limitations:** In this paper, we explore an effective framework for deep subspace clustering with theoretical justification. However, it is not clear how to develop the geometric guarantee for our PRO-DSC framework to yield correctly subspace-preserving solution. Moreover, it is an unsupervised learning framework, we left the extension to semi-supervised setting as future work.

**Failure Cases:** In this paper, we evaluate our PRO-DSC framework on four scenarios of synthetic data (Fig. 5 and B.1), six benchmark datasets with CLIP features (Table 1), five benchmark datasets with BYOL pre-trained features (Table B.3), three out-of-domain datasets (Table B.5), using four different regularization terms (Table 3), using different feature extractor (Table B.7) and varying hyper-parameters (Fig. 7 and Table B.6). We also conduct experiments on two face image datasets (Table B.10), text and temporal dataset (Table B.8). However, as demonstrated in Fig. 1, our PRO-DSC will fail if the sufficient condition to prevent catastrophic collapse is not satisfied by using improper hyper-parameters  $\gamma$  and  $\alpha$ .

**Extensibility:** As a general framework for self-expressive model based deep subspace clustering, our PRO-DSC is reasonable, scalable and flexible to miscellaneous extensions. For example, rather than using  $\log \det(\cdot)$ , there are other methods to solve the feature collapse issue, e.g., the nuclear norm. In addition, it is also worthwhile to incorporate the supervision information from the pseudo-label, e.g., (Huang et al., 2023; Jia et al., 2025; Li et al., 2017), for further improving the performance of our PRO-DSC.



(e) ImageNet-1k

Figure B.3: Visualization of the union-of-orthogonal-subspaces structure of the learned representations via Gram matrix and PCA dimension reduction on three categories. Left:  $|X^{\top}X|$ . Mid-left:  $|Z^{\top}Z|$ . Mid-right:  $X_{(3)}$  via PCA. Right:  $Z_{(3)}$  via PCA.



Figure B.4: Singular values of features from all samples (left) and features from each class (mid and right). For the better clarity, we plot the singular values for the first ten classes.



Figure B.5: The learning curves w.r.t. loss values and evaluation performance of PRO-DSC on CIFAR-20, CIFAR-100 and ImageNet-1k dataset.



Figure B.6: *t*-SNE visualization of CLIP features and PRO-DSC's learned representations. The experiments are conducted on CIFAR-10 and CIFAR-100 dataset.



(a) Cluster 1



(d) Cluster 4



(g) Cluster 7









(h) Cluster 8



(c) Cluster 3



(f) Cluster 6

(i) Cluster 9



Figure B.7: Visualization of the principal components in CIFAR-10 dataset. For each cluster, we display the most similar images to its principal components.