

EXPLORING A PRINCIPLED FRAMEWORK FOR DEEP SUBSPACE CLUSTERING

Anonymous authors

Paper under double-blind review

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, and prove that the regularized self-expressive model is able to prevent feature space collapse and 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 (Elhamifar & Vidal, 2009; 2013). Formally, given a dataset $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ where $\mathbf{x}_j \in \mathbb{R}^D$, self-expressive model expresses each data point \mathbf{x}_j by a linear combination of other points, i.e.,

$$\mathbf{x}_j = \sum_{i \neq j} c_{ij} \mathbf{x}_i, \quad (1)$$

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 \mathbf{x}_i and \mathbf{x}_j are in the same subspace (Elhamifar & Vidal, 2013; Soltanolkotabi & Candes, 2012). 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.

054 Earlier works devote to address this deficiency by learning a linear transform of the data (Patel
055 et al., 2013; 2015) or introducing a nonlinear kernel mapping (Patel & Vidal, 2014) under which the
056 representations of the data are supposed to be aligned with the UoS assumption. However, there is
057 a lack of principled mechanism to guide the learning of the linear transforms or the design of the
058 nonlinear kernels to guarantee the representations of the data to form a UoS structure.

059 To handle complex real-world data, in the past few years, there is a surge of interests in designing
060 deep subspace clustering frameworks, e.g., (Peng et al., 2016; 2018; Ji et al., 2017; Zhou et al., 2018;
061 Zhang et al., 2019a; Dang et al., 2020; Peng et al., 2020; Lv et al., 2021; Wang et al., 2023b). In
062 these works, usually a deep neural network-based representation learning module is integrated to the
063 self-expressive model, to learn the representations $\mathbf{Z} \in \mathbb{R}^{d \times N}$ and the self-expressive coefficients
064 $\mathbf{C} = \{c_{ij}\}_{i,j=1}^N$ in a joint optimization framework. However, as analyzed in (Haeffele et al., 2021)
065 that, the optimal representations \mathbf{Z} of these methods tend to **catastrophically** collapse into subspaces
066 with dimensions much lower than the ambient space, which is detrimental to subspace clustering
067 and there is no evidence that the learned representations form a UoS structure.

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

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

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

086 To the best of our knowledge, this is the first principled framework for deep subspace clustering that
087 is guaranteed to yield the desired UoS representations.

088 2 DEEP SUBSPACE CLUSTERING: A PRINCIPLED FRAMEWORK, 089 JUSTIFICATION, AND IMPLEMENTATION

090 In this section, we review the deficiency that was suffering in the existing Self-Expressive Deep
091 Subspace Clustering (SEDSC) frameworks at first, then present our principled framework for deep
092 subspace clustering and provide a rigorous characterization of the optimal solution and the property
093 of the learned structured representations. Finally we describe a scalable implementation based on
094 differential programming for the proposed framework. Please refer to Appendix A for the detailed
095 proofs of our theoretical results.

096 2.1 PREREQUISITE

097 To apply subspace clustering to complex real-world data that may not well align with the UoS
098 assumption, there has been a surge of interests in exploiting deep neural networks to learn represen-
099 tations and then apply self-expressive model to the learned representations, e.g., (Peng et al., 2016;
100 2018; Ji et al., 2017; Zhou et al., 2018; Zhang et al., 2019a; Dang et al., 2020; Peng et al., 2020; Lv
101 et al., 2021; Wang et al., 2023b).

Formally, the optimization problem of these SEDSC models can be formulated as:¹

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

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

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

In other words, the optimal representation \mathbf{Z} in SEDSC is restricted to an extremely “narrow” subspace whose dimension is much smaller than d , leading to an undesirable collapsed solution.²

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 space collapse by incorporating the $\log \det(\cdot)$ -based concave smooth surrogate which is defined as follows:

$$R(\mathbf{Z}; \alpha) := \log \det(\mathbf{I} + \alpha \mathbf{Z}^\top \mathbf{Z}) \quad \text{s.t.} \quad \mathbf{Z} \in \mathbb{S}^{d-1}, \quad (3)$$

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

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):

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

where $\gamma > 0$ is a balancing hyper-parameter. Now, we will give our theoretical findings about the optimal solution for problem (4). Note that the Lagrangian of problem (4) is:

$$\mathcal{L}(\mathbf{Z}, \mathbf{C}, \nu) := -\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{Z}^\top \mathbf{Z}) + \frac{\gamma}{2} \|\mathbf{Z} - \mathbf{Z}\mathbf{C}\|_F^2 + \beta r(\mathbf{C}) + \nu(\|\mathbf{Z}\|_F^2 - N), \quad (5)$$

where $\nu \in \mathbb{R}$ is the Lagrange multiplier. By analyzing the gradient flow dynamics of problem (5), we show that the eigenspaces of $\mathbf{Z}^\top \mathbf{Z}$ and $(\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$ align progressively, which is crucial for the subsequent analysis. Precisely, we have a theorem as follows.

Theorem 1 (Informal Statement). *Denote $\mathbf{G} := \mathbf{Z}^\top \mathbf{Z}$, $\mathbf{M} := (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$. Under mild condition, then \mathbf{G} and \mathbf{M} will converge to have eigenspaces aligned, i.e., \mathbf{G} and \mathbf{M} can be diagonalized simultaneously by $\mathbf{U} \in \mathcal{O}^{N \times N}$ where $\mathcal{O}^{N \times N}$ is an $N \times N$ orthogonal matrix group.*

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

$$\min_{\mathbf{Z}} -\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{Z}^\top \mathbf{Z}) + \frac{\gamma}{2} \|\mathbf{Z} - \mathbf{Z}\mathbf{C}\|_F^2 \quad \text{s.t.} \quad \|\mathbf{Z}\|_F^2 = N, \quad (6)$$

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

¹Without loss of generality, we omit the constraint $\text{Diag}(\mathbf{C}) = \mathbf{0}$ throughout the analysis.

²The dimension equals to the multiplicity of the smallest eigenvalues of $(\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$.

³Please refer to (Ma et al., 2007) for packing-ball interpretation.

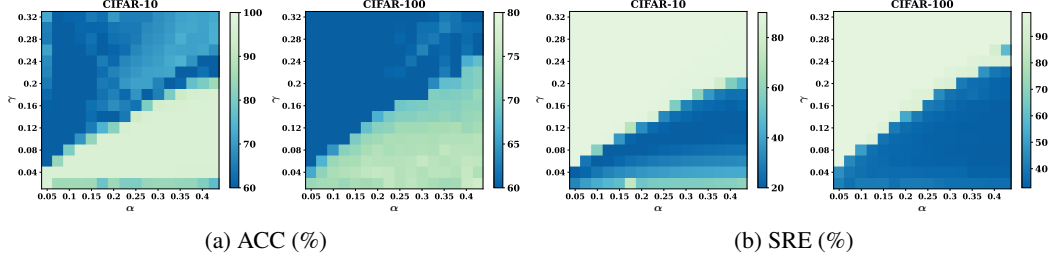


Figure 1: **Empirical validation to Theorem 2.** We train PRO-DSC on CIFAR-10 and CIFAR-100, and report the clustering accuracy (ACC%) and subspace-preserving representation error (SRE%) with varying α and γ . A clear linear phase transition phenomenon can be observed, which is consistent with the condition to avoid collapse: $\gamma < (\alpha - \nu_*)/\lambda_{max}(\mathbf{M})$.

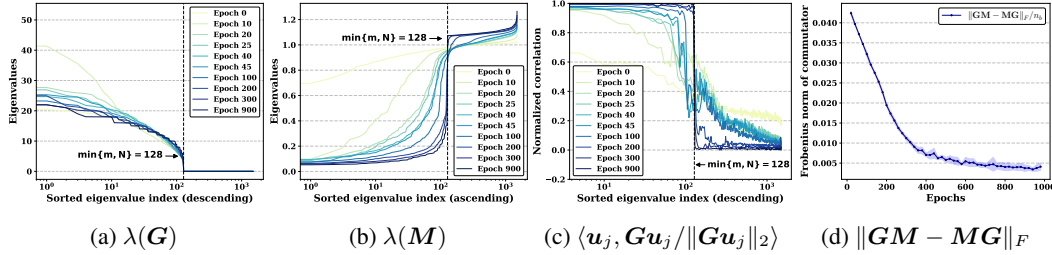


Figure 2: **Empirical validations to Theorem 1 and 2.** We train PRO-DSC on CIFAR-100. (a) Eigenvalues of \mathbf{G} . (b) Eigenvalues of \mathbf{M} . (c) For each eigenvector \mathbf{u}_j of \mathbf{M} , we compute $\langle \mathbf{u}_j, \frac{\mathbf{G}\mathbf{u}_j}{\|\mathbf{G}\mathbf{u}_j\|_2} \rangle$ to measure the eigenspace alignment of \mathbf{G} and \mathbf{M} . (d) We compute $\|\mathbf{L}\|_F$ to measure the eigenspace alignment of \mathbf{G} and \mathbf{M} .

In light of the fact that \mathbf{G} and \mathbf{M} converge to share eigenspaces, we decompose \mathbf{G} and \mathbf{M} to $\mathbf{U} \text{Diag}(\lambda_{\mathbf{G}}^{(1)}, \dots, \lambda_{\mathbf{G}}^{(N)})\mathbf{U}^\top$ and $\mathbf{U} \text{Diag}(\lambda_{\mathbf{M}}^{(1)}, \dots, \lambda_{\mathbf{M}}^{(N)})\mathbf{U}^\top$, respectively. Recall that $\mathbf{G} := \mathbf{Z}^\top \mathbf{Z}$, $\mathbf{M} := (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$, by using the eigenvalue decomposition, we have that problem (6) is transformed to a *convex* problem (detailed in Appendix A):

$$\begin{aligned} \min_{\{\lambda_{\mathbf{G}}^{(i)}\}_{i=1}^{\min\{d, N\}}} & -\frac{1}{2} \sum_{i=1}^{\min\{d, N\}} \log(1 + \alpha \lambda_{\mathbf{G}}^{(i)}) + \frac{\gamma}{2} \lambda_{\mathbf{M}}^{(i)} \lambda_{\mathbf{G}}^{(i)} \\ \text{s.t.} & \sum_{i=1}^{\min\{d, N\}} \lambda_{\mathbf{G}}^{(i)} = N, \lambda_{\mathbf{G}}^{(i)} \geq 0, \quad \text{for all } i = 1, \dots, \min\{d, N\}, \end{aligned} \quad (7)$$

which is a classical reverse water-filling problem (Yu et al., 2004). By solving problem (7), we have the following theorem.

Theorem 2. Suppose that \mathbf{G} and \mathbf{M} have eigenspaces aligned and $\gamma < (\alpha - \nu_*)/\lambda_{max}(\mathbf{M})$, then we have that $\text{rank}(\mathbf{Z}_*) = \min\{d, N\}$ and the singular values $\sigma_{\mathbf{Z}_*}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_{\mathbf{M}}^{(i)} + \nu_*} - \frac{1}{\alpha}}$, for all $i = 1, \dots, \min\{d, N\}$, where \mathbf{Z}_* and ν_* are the primal optimal solution and dual optimal solution.

The above theorem characterizes the optimal solution for problem (6). As shown, the rank of the minimizers for PRO-DSC satisfies that $\text{rank}(\mathbf{Z}_*) = \min\{d, N\}$. On the contrary, SEDSC yields a collapsed solution, where $\text{rank}(\mathbf{Z}_*) \ll \min\{d, N\}$. In Figure 1, we show the subspace clustering accuracy (ACC) and subspace-representation error⁴ (SRE) as a function of the parameter α and γ . The phase transition phenomenon around $\gamma = (\alpha - \nu_*)/\lambda_{max}(\mathbf{M})$ well illustrates the sufficient condition in Theorem 2 to avoid collapse. In Figure 2, we illustrate the curves of the eigenvalues and the alignment of eigenspaces of \mathbf{G} and \mathbf{M} , which demonstrate that the learned representation does no longer collapse and the two eigenspaces are approximately to be aligned.

⁴For each row $\mathbf{c}^{(j)}$ in \mathbf{C} , SRE is computed by $\frac{100}{N} \sum_j (1 - \sum_i w_{ij} \cdot |c_{ij}|) / \|\mathbf{c}^{(j)}\|_1$, where $w_{ij} \in \{0, 1\}$ is the ground-truth affinity.

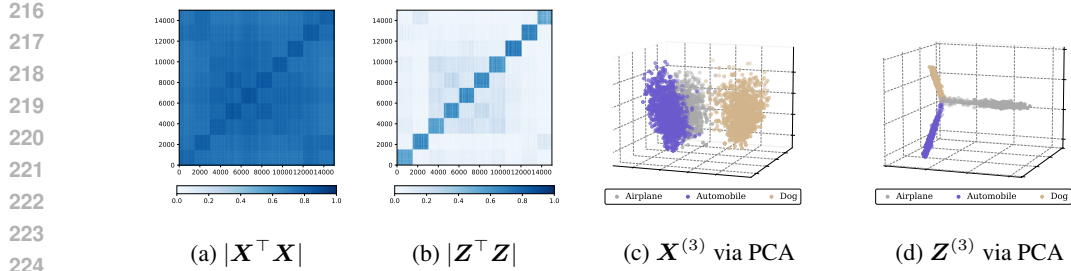


Figure 3: **Empirical validation to Theorem 3.** We visualize the Gram matrices and the dimension reduction results via PCA for CLIP features \mathbf{X} and the learned representations \mathbf{Z} on CIFAR-10. For the clarity of visualization, we apply PCA to the samples from three categories $\mathbf{X}^{(3)}$ and $\mathbf{Z}^{(3)}$.

Furthermore, from the perspective of joint optimizing \mathbf{Z} and \mathbf{C} , the following theorem demonstrates that PRO-DSC yields union-of-orthogonal-subspaces representations \mathbf{Z} and block-diagonal self-expressive matrix \mathbf{C} under certain condition.

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 $\mathbf{\Gamma}$, 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 $\langle (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \rangle \rightarrow 0$, where $\mathbf{G}^* := \text{Diag}(\mathbf{G}_{11}, \mathbf{G}_{22}, \dots, \mathbf{G}_{kk})$ and \mathbf{G}_{jj} is the block Gram matrix corresponding to \mathbf{Z}_j .*

Remark 1. Theorem 3 suggests that our PRO-DSC is able to learn 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 $\langle (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \rangle \rightarrow 0$ is met. We call this condition a *compatibly structured coherence* (CSC), which relates to the properties of the distribution of the representations in \mathbf{Z} and the self-coefficients in \mathbf{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 Figures 3 and 4).⁵

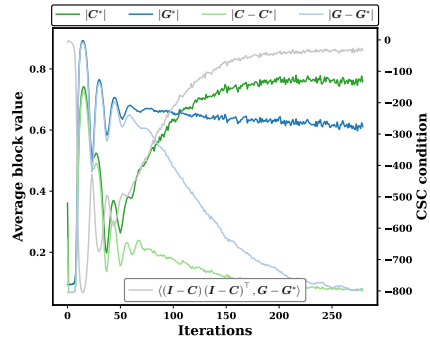


Figure 4: **Empirical validation to Theorem 3 on CIFAR-10.**

2.3 SCALABLE IMPLEMENTATION

Existing SEDSC models typically use autoencoders to learn the representations, and learn the self-expressive matrix \mathbf{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 is 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 \mathbf{x}_i is fed into a neural network $\mathbf{h}(\cdot; \Psi) : \mathbb{R}^D \rightarrow \mathbb{R}^d$ to yield normalized representations, i.e.,

$$\mathbf{y}_i := \mathbf{h}(\mathbf{x}_i; \Psi) / \|\mathbf{h}(\mathbf{x}_i; \Psi)\|_2, \quad (8)$$

⁵Please refer to Appendix B.2 for more details about Figures 1–4.

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

$$C_\Psi := \mathcal{P}(Y^\top Y), \quad (9)$$

where $Y := [\mathbf{y}_1, \dots, \mathbf{y}_N] \in \mathbb{R}^{d \times N}$ 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).⁶ To enable efficient representation learning, we introduce another learnable mapping $f(\cdot; \Theta) : \mathbb{R}^D \rightarrow \mathbb{R}^d$, for which

$$\mathbf{z}_j := f(\mathbf{x}_j; \Theta) / \|f(\mathbf{x}_j; \Theta)\|_2 \quad (10)$$

is the learned representation for the input \mathbf{x}_j , where Θ denotes the parameters in $f(\cdot)$ to learn the structured representation $Z_\Theta := [\mathbf{z}_1, \dots, \mathbf{z}_N] \in \mathbb{R}^{d \times N}$.

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

$$\min_{\Theta, \Psi} \mathcal{L}(\Theta, \Psi) := -\frac{1}{2} \log \det(\mathbf{I} + \alpha Z_\Theta^\top Z_\Theta) + \frac{\gamma}{2} \|Z_\Theta - Z_\Theta C_\Psi\|_F^2 + \beta \cdot r(C_\Psi). \quad (11)$$

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_Ψ , which is induced by default as $A_\Psi := \frac{1}{2} (|C_\Psi| + |C_\Psi^\top|)$, the block diagonal regularizer is defined as:

$$r(C_\Psi) := \|A_\Psi\|_{\square}, \quad (12)$$

where $\|A_\Psi\|_{\square}$ is the sum of the k smallest eigenvalues of the Laplacian matrix of the affinity A_Ψ .⁷

Consequently, the parameters in Θ and Ψ of reparameterized PRO-DSC can be trained by Stochastic Gradient Descent (SGD) with the loss function $\mathcal{L}(\Theta, \Psi)$ defined in (11). 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 owing to the learned structured representation (Please refer to Table 4 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).

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 α, β, γ , number of iterations T , learning rate η

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

Training:

- 1: **for** $t = 1, \dots, T$ **do**
- 2: Sample a batch $X_b \in \mathbb{R}^{D \times n_b}$ from $\mathcal{X}_{\text{train}}$
 # *Forward propagation*
- 3: Compute self-expressive matrix $C_b \in \mathbb{R}^{n_b \times n_b}$ by Eqs. (8–9)
- 4: Compute representations $Z_b \in \mathbb{R}^{d \times n_b}$ by Eq. (10)
 # *Backward propagation*
- 5: Compute $\nabla_\Psi := \frac{\partial \mathcal{L}}{\partial \Psi}, \nabla_\Theta := \frac{\partial \mathcal{L}}{\partial \Theta}$
- 6: Set $\Psi \leftarrow \Psi - \eta \cdot \nabla_\Psi, \Theta \leftarrow \Theta - \eta \cdot \nabla_\Theta$
- 7: **end for**

Testing:

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

⁶In practice, we set $\text{diag}(C_\Psi) = \mathbf{0}$ to prevent trivial solution $C_\Psi = \mathbf{I}$.

⁷Recall that the number of zero eigenvalues of the Laplacian matrix equals to the number of connected components in the graph (von Luxburg, 2007).

3 EXPERIMENTS

To validate the effectiveness of theoretical findings and to demonstrate the superior 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.

3.1 EXPERIMENTS ON SYNTHETIC DATA

To validate that PRO-DSC addresses the collapse issue in SEDSC and learns representations with a UoS structure, we conduct experiments on synthetic data and visualize the results.

We first follow the procedure in (Ding et al., 2023) to generate synthetic data, as shown in the first row of Figure 5a. To conduct experiments on more complicated scenarios, we randomly select half of the points from the “Arctic” and place them near the “Antarctic” (Figure 5a, second row).

As shown in Figure 5b, the SEDSC models overly compress all the representations to a closed curve on the hypersphere. With increased weights (i.e., $\gamma \uparrow$) of the self-expressive term, the representations collapse to a few points (Figure 5c). The manifold linearizing and clustering method, MLC (Ding et al., 2023) approximately compresses the representations to the orthogonal subspaces. In contrast, our PRO-DSC learns linearized representations lying on orthogonal subspaces in both scenarios, confirming the validity of our theoretical analysis.

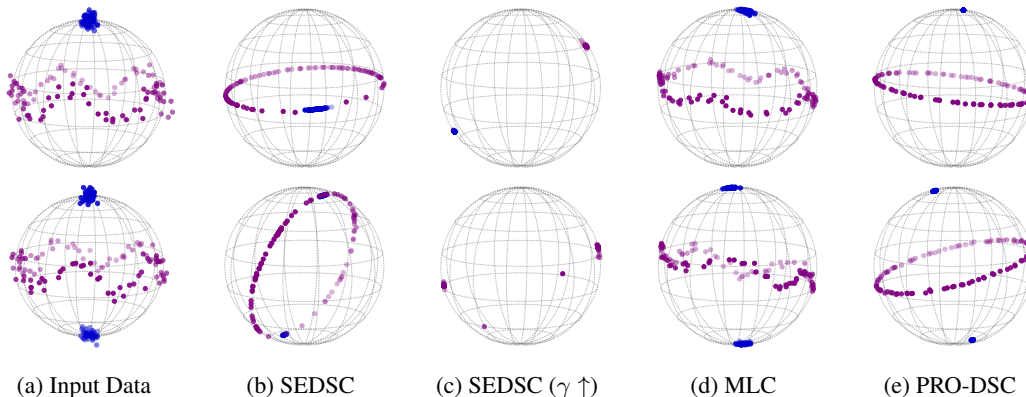


Figure 5: Visualization of the input data and learned representations with different algorithms.

3.2 EXPERIMENTS ON REAL-WORLD DATA

To evaluate the performance of our proposed approach, we conduct experiments on six real-world image datasets: CIFAR-10, CIFAR-20, CIFAR-100, ImageNet-Dogs-15, Tiny-ImageNet-200, and ImageNet-1k. We measure clustering performance using clustering accuracy (ACC) and normalized mutual information (NMI). The results of our PRO-DSC in Tables 1 and 2 are **averaged over 10 trials (with \pm std) and other results of PRO-DSC are averaged over 3 trials.**

Main results. Table 1 compares the clustering performance of our PRO-DSC on CLIP features (Radford et al., 2021) with various baseline methods, including classical clustering algorithms, e.g., k -means (MacQueen, 1967), spectral clustering (Shi & Malik, 2000), subspace clustering algorithm, e.g., EnSC (You et al., 2016a), SENet (Zhang et al., 2021), deep clustering algorithms, e.g., SCAN (Van Gansbeke et al., 2020), TEMI (Adaloglou et al., 2023), CPP (Chu et al., 2024), and deep subspace clustering algorithms, e.g., EDESC (Cai et al., 2022), DSCNet (Ji et al., 2017).⁸ Since that the clustering performance with the CLIP feature is not reported for most baselines, we conduct experiments using the implementations provided by the authors. As shown in Table 1, PRO-DSC significantly outperforms the vanilla subspace clustering algorithms, achieving a 10% improvement on both datasets CIFAR-20 and CIFAR-100. In contrast, the previous deep subspace clustering algorithms perform less competitively, primarily because they do not learn structured representations forming a UoS. Compared to state-of-the-art deep clustering algorithms, our PRO-DSC improves

⁸Please refer to Appendix B.3 for the results on other pre-trained models.

Table 1: **Clustering results of our PRO-DSC on the CLIP features.** The best results are in bold font and the second best results are underlined. “OOM” means out of GPU memory.

Method	CIFAR-10		CIFAR-20		CIFAR-100		TinyImgNet-200		ImgNetDogs-15		ImageNet-1k	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
<i>k</i> -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	<u>64.5</u>	<u>77.7</u>	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 \pm 0.2	92.8 \pm 0.4	71.6 \pm 1.2	73.2 \pm 0.5	77.3 \pm 1.0	82.4 \pm 0.5	69.8 \pm 1.1	80.5 \pm 0.7	84.0 \pm 0.6	81.2 \pm 0.8	65.0 \pm 1.2	83.4 \pm 0.6

Table 2: **Clustering results of our PRO-DSC when training from scratch.** The best results are in bold font and the second best results are underlined. Performance marked with “*” is based on our re-implementation.

Method	CIFAR-10		CIFAR-20		CIFAR-100		TinyImgNet-200		ImgNetDogs-15	
	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI	ACC	NMI
<i>k</i> -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	<u>53.1</u>	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>	<u>85.5</u>	58.3	<u>59.6</u>	49.4	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 \pm 0.1	31.1 \pm 0.3	46.0 \pm 1.0	74.1 \pm 0.5	69.5 \pm 0.6

the accuracy by nearly 4% and 6% on datasets CIFAR-20 and TinyImageNet, respectively, further validating its effectiveness.

To validate the effectiveness of our PRO-DSC without using CLIP features, we conduct a fair comparison with existing deep clustering approaches and report the clustering results with training from scratch in Table 2. By stacking the $f(\cdot)$ and $h(\cdot)$ on a learnable backbone, PRO-DSC can learn representations and self-expressive coefficients directly from raw images. As illustrated in Table 2, 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).

Evaluation on learned representations. To quantitatively evaluate the effectiveness of the learned representations, we conduct experiments to compare the clustering performance with the CLIP features and the representations learned from CPP and our PRO-DSC (additional SEDSC results can be found in Table B.6). Specifically, we use *k*-means (MacQueen, 1967), spectral clustering (Shi & Malik, 2000), and EnSC (You et al., 2016a) to yield the clustering results.

Experimental result are shown in Figure 6. We see 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. We note that the clustering performance is further improved by the learnable mapping $h(\cdot; \Psi)$, suggesting its superior generalization ability.

Sensitivity of hyper-parameters. In Figure 1, we verify that our PRO-DSC yields satisfactory results when the sufficient conditions in Theorem 2 to avoid collapse are met. Moreover, we evaluate the performance sensitivity to hyper-parameters γ and β by experiments on the CLIP features of CIFAR-10, CIFAR-100 and TinyImageNet-200 with varying γ and β . In Figure 7, we observe that the clustering performance maintains satisfactory under a broad range of γ and β .

Time and memory cost. The most time-consuming operations in our PRO-DSC are computing the term involving $\log \det(\cdot)$ and the term $\|\mathbf{A}\|_{\text{F}}$ involving eigenvalue decomposition, respectively. For

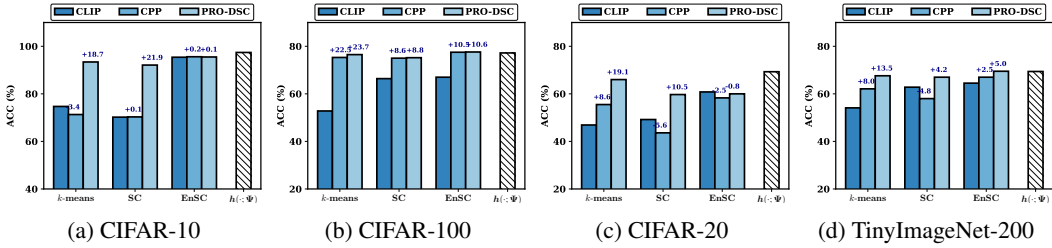


Figure 6: Clustering accuracy with CLIP features and learned representations.

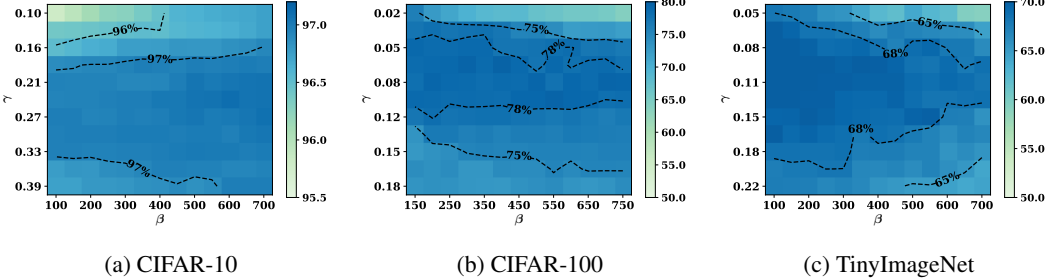


Figure 7: Evaluation on sensitivity to hyper-parameters γ and β on three datasets.

the former, the Gram matrix has a shape of $n_b \times n_b$, getting a time complexity of $\mathcal{O}(n_b^3)$. For the latter, the Laplacian matrix has a shape of $n_b \times n_b$, also resulting in a time complexity of $\mathcal{O}(n_b^3)$. Therefore, the overall time complexity of our PRO-DSC is $\mathcal{O}(n_b^3)$. TEMI (Adaloglou et al., 2023) employs $H = 50$ cluster heads during training, adding further time and memory costs. CPP (Chu et al., 2024) involves computing $\log \det(\cdot)$ for $n_b + 1$ times, leading to complexity $\mathcal{O}(n_b d^3)$. As shown in Table 3, our PRO-DSC significantly reduces the time consumption, particularly for datasets with a large number of clusters. All the experiments are conducted on a single NVIDIA RTX 3090 GPU and Intel Xeon Platinum 8255C CPU.

Table 3: Comparison on time (s) and memory cost (MiB). ‘‘OOM’’ means out of GPU memory.

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

Ablation study. To verify the effectiveness of the loss function in PRO-DSC, we conduct ablation studies on the CLIP features of CIFAR-10, CIFAR-100, and ImageNetDogs-15, and report the results in Table 4, where $\mathcal{L}_1 := -\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{Z}_\Theta^\top \mathbf{Z}_\Theta)$ and $\mathcal{L}_2 := \frac{1}{2} \|\mathbf{Z}_\Theta - \mathbf{Z}_\Theta \mathbf{C}_\Psi\|_F^2$. 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 for learning the self-expressive coefficients. In both cases, clustering performance drops significantly. More interestingly, when we replace the block diagonal regularizer $\|\mathbf{A}\|_{\text{tr}}$ with $\|\mathbf{C}\|_1$, $\|\mathbf{C}\|_*$, and $\|\mathbf{C}\|_F^2$ 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 owing 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

Table 4: Ablation studies on different loss functions and regularizers.

	Loss Term						CIFAR-10		CIFAR-100		ImgNetDogs-15	
	\mathcal{L}_1	\mathcal{L}_2	$\ A\ _{\square}$	$\ C\ _1$	$\ C\ _F^2$	$\ C\ _*$	ACC	NMI	ACC	NMI	ACC	NMI
Ablation		✓	✓				56.9	47.7	54.6	60.9	46.7	37.1
	✓			✓			69.6	56.4	64.7	71.7	10.5	1.7
	✓	✓					97.0	93.0	74.6	80.9	80.9	78.8
Regularizer	✓	✓				✓	97.0	92.6	75.2	81.1	81.3	79.1
	✓	✓			✓		97.0	92.6	75.2	80.9	80.9	78.8
	✓	✓		✓			96.7	91.9	76.4	81.8	81.0	78.8
	✓	✓	✓				97.2	92.8	77.3	82.4	84.0	81.2

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), SEDSC suffers from a catastrophic feature collapse, which is detrimental to subspace clustering. 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., 2020), 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.

Representations 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 (MCR²) is proposed to learn representations that favor the structure of a union of orthogonal subspaces (Wang et al., 2024). More recently, the MCR² 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 MCR² framework combines with contrastive learning to perform manifold clustering and representation learning; in (Ding et al., 2023), the MCR² framework combines with doubly stochastic affinity learning to perform manifold linearizing and clustering; and in (Chu et al., 2024), the performance of (Ding et al., 2023) on large pre-trained model (e.g., CLIP) has been investigated. While the modified MCR² framework has been incorporated into these methods for manifold clustering, none of them provide 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, 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 the theoretical findings and the superior performance of the proposed approach.

540 ETHICS STATEMENT

541

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

549

550 REPRODUCIBILITY STATEMENT

551

552 To ensure the reproducibility of our work, we have submitted the anonymized source code. The-
 553oretical proofs of the claims made in this paper are provided in Appendix A, and the empirical
 554 validation of these theoretical results is shown in Figures 1– 4, with further detailed explanations in
 555 Appendix B.2. All datasets used in our experiments are publicly available, and we have provided
 556 a comprehensive description of the data processing steps in Appendix B.1. Additionally, detailed
 557 experimental settings and configurations are outlined in Appendix B.1 to facilitate the reproduction
 558 of our results.

559

560 REFERENCES

561

562 Nikolas Adaloglou, Felix Michels, Hamza Kalisch, and Markus Kollmann. Exploring the limits of
 563 deep image clustering using pretrained models. *arXiv preprint arXiv:2303.17896*, 2023.

564 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix
 565 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.

566

567 Laurent Bako and René Vidal. Algebraic identification of MIMO SARX models. In *International
 568 Workshop on Hybrid Systems: Computation and Control*, pp. 43–57. Springer-Verlag, 2008.

569

570 Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. Greedy layer-wise training of
 571 deep networks. In Bernhard Schölkopf, John C. Platt, and Thomas Hofmann (eds.), *Advances in
 572 Neural Information Processing Systems 19, Proceedings of the Twentieth Annual Conference on
 573 Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 4-7,
 574 2006*, pp. 153–160. MIT Press, 2006.

575 Jinyu Cai, Jicong Fan, Wenzhong Guo, Shiping Wang, Yunhe Zhang, and Zhao Zhang. Efficient
 576 deep embedded subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer
 577 Vision and Pattern Recognition*, pp. 1–10, 2022.

578

579 Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for un-
 580 supervised learning of visual features. In *Proceedings of the European conference on computer
 581 vision (ECCV)*, pp. 132–149, 2018.

582 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
 583 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural
 584 information processing systems*, 33:9912–9924, 2020.

585

586 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
 587 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of
 588 the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.

589 Jianlong Chang, Gaofeng Meng, Lingfeng Wang, Shiming Xiang, and Chunhong Pan. Deep self-
 590 evolution clustering. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):
 591 809–823, 2018.

592

593 Guangliang Chen and Gilad Lerman. Spectral curvature clustering (SCC). *International Journal of
 Computer Vision*, 81(3):317–330, 2009. ISSN 0920-5691.

- 594 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
595 contrastive learning of visual representations. In *International conference on machine learning*,
596 pp. 1597–1607. PMLR, 2020.
- 597 Tianzhe Chu, Shengbang Tong, Tianjiao Ding, Xili Dai, Benjamin David Haeffele, René Vidal, and
598 Yi Ma. Image clustering via the principle of rate reduction in the age of pretrained models. In
599 *The Twelfth International Conference on Learning Representations*, 2024.
- 600 Joao Paulo Costeira and Takeo Kanade. A multibody factorization method for independently moving
601 objects. *International Journal of Computer Vision*, 29:159–179, 1998.
- 602 Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural
603 information processing systems*, 26, 2013.
- 604 Zhiyuan Dang, Cheng Deng, Xu Yang, and Heng Huang. Multi-scale fusion subspace clustering
605 using similarity constraint. In *Proceedings of IEEE International Conference on Computer Vision
606 and Pattern Recognition*, pp. 6657–6666, 2020.
- 607 Zhiyuan Dang, Cheng Deng, Xu Yang, Kun Wei, and Heng Huang. Nearest neighbor matching for
608 deep clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
609 Recognition*, pp. 13693–13702, 2021.
- 610 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
611 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
612 pp. 248–255. Ieee, 2009.
- 613 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE
614 Signal Processing Magazine*, 29(6):141–142, 2012.
- 615 Tianjiao Ding, Shengbang Tong, Kwan Ho Ryan Chan, Xili Dai, Yi Ma, and Benjamin D. Haeffele.
616 Unsupervised manifold linearizing and clustering. In *Proceedings of the IEEE/CVF International
617 Conference on Computer Vision (ICCV)*, pp. 5450–5461, October 2023.
- 618 Ehsan Elhamifar and René Vidal. Sparse subspace clustering. In *IEEE Conference on Computer
619 Vision and Pattern Recognition*, pp. 2790–2797, 2009.
- 620 Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications.
621 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2765–2781, 2013.
- 622 Maryam Fazel, Haitham Hindi, and Stephen P Boyd. Log-det heuristic for matrix rank minimization
623 with applications to hankel and euclidean distance matrices. In *Proceedings of the 2003 American
624 Control Conference, 2003.*, volume 3, pp. 2156–2162. IEEE, 2003.
- 625 Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting
626 with applications to image analysis and automated cartography. *Communications of the ACM*, 24
627 (6):381–395, 1981.
- 628 Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena
629 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
630 et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural
631 information processing systems*, 33:21271–21284, 2020.
- 632 Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro.
633 Implicit regularization in matrix factorization. In *Neural Information Processing Systems*, pp.
634 6151–6159, 2017.
- 635 Benjamin D Haeffele, Chong You, and René Vidal. A critique of self-expressive deep subspace
636 clustering. In *The Ninth International Conference on Learning Representations*, 2021.
- 637 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
638 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on
639 computer vision and pattern recognition*, pp. 9729–9738, 2020.

- 648 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
649 toencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*
650 *vision and pattern recognition*, pp. 16000–16009, 2022.
- 651 Jeffrey Ho, Ming-Husang Yang, Jongwoo Lim, Kuang-Chih Lee, and David Kriegman. Clustering
652 appearances of objects under varying illumination conditions. In *Proceedings of IEEE Interna-*
653 *tional Conference on Computer Vision and Pattern Recognition*, pp. 11–18, 2003.
- 654 Wei Hong, John Wright, Kun Huang, and Yi Ma. Multiscale hybrid linear models for lossy image
655 representation. *IEEE Transactions on Image Processing*, 15(12):3655–3671, 2006.
- 656 Zhizhong Huang, Jie Chen, Junping Zhang, and Hongming Shan. Learning representation for clus-
657 tering via prototype scattering and positive sampling. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45
658 (6):7509–7524, 2023.
- 659 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
660 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.
661 pmlr, 2015.
- 662 Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering
663 networks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
664 R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 24–33. Curran
665 Associates, Inc., 2017.
- 666 Yuheng Jia, Jianhong Cheng, Hui Liu, and Junhui Hou. Towards calibrated deep clustering network.
667 *CoRR*, abs/2403.02998, 2024.
- 668 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann
669 LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB,*
670 *Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- 671 José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. OLE: Orthogonal low-rank embed-
672 ding - a plug and play geometric loss for deep learning. In *Proceedings of IEEE International*
673 *Conference on Computer Vision and Pattern Recognition*, pp. 8109–8118, 2018.
- 674 Chun-Guang Li, Chong You, and René Vidal. Structured sparse subspace clustering: A joint affinity
675 learning and subspace clustering framework. *IEEE Transactions on Image Processing*, 26(6):
676 2988–3001, 2017.
- 677 Jun Li, Hongfu Liu, Zhiqiang Tao, Handong Zhao, and Yun Fu. Learnable subspace clustering.
678 *arXiv preprint arXiv:2004.04520*, 2020.
- 679 Yunfan Li, Peng Hu, Zitao Liu, Dezhong Peng, Joey Tianyi Zhou, and Xi Peng. Contrastive cluster-
680 ing. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 8547–8555,
681 2021.
- 682 Zengyi Li, Yubei Chen, Yann LeCun, and Friedrich T Sommer. Neural manifold clustering and
683 embedding. *arXiv preprint arXiv:2201.10000*, 2022.
- 684 Derek Lim, René Vidal, and Benjamin D Haeffele. Doubly stochastic subspace clustering. *arXiv*
685 *preprint arXiv:2011.14859*, 2020.
- 686 Guangcan Liu, Zhouchen Lin, and Yong Yu. Robust subspace segmentation by low-rank represen-
687 tation. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 663–670,
688 2010.
- 689 Xin Liu, Zhongdao Wang, Ya-Li Li, and Shengjin Wang. Self-supervised learning via maximum
690 entropy coding. *Advances in Neural Information Processing Systems*, 35:34091–34105, 2022.
- 691 Canyi Lu, Hai Min, Zhong-Qiu Zhao, Lin Zhu, De-Shuang Huang, and Shuicheng Yan. Robust
692 and efficient subspace segmentation via least squares regression. In *European Conference on*
693 *Computer Vision*, pp. 347–360, 2012.

- 702 Canyi Lu, Jiashi Feng, Zhouchen Lin, Tao Mei, and Shuicheng Yan. Subspace clustering by block
703 diagonal representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
704
- 705 Juncheng Lv, Zhao Kang, Xiao Lu, and Zenglin Xu. Pseudo-supervised deep subspace clustering.
706 *IEEE Transactions on Image Processing*, 30:5252–5263, 2021.
- 707 Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of multivariate mixed data via
708 lossy coding and compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
709 29(9):1546–1562, 2007.
- 710 J. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceed-*
711 *ings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, pp. 281–297,
712 1967.
- 713 Ryan McConville, Raul Santos-Rodriguez, Robert J Piechocki, and Ian Craddock. N2d:(not too)
714 deep clustering via clustering the local manifold of an autoencoded embedding. In *2020 25th*
715 *international conference on pattern recognition (ICPR)*, pp. 5145–5152. IEEE, 2021.
- 716 Brian McWilliams and Giovanni Montana. Subspace clustering of high dimensional data: a predic-
717 tive approach. *Data Mining and Knowledge Discovery*, 28(3):736–772, 2014.
- 718 Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In
719 *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814,
720 2010.
- 721 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number
722 of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pp.
723 722–729. IEEE, 2008.
- 724 Chuang Niu, Hongming Shan, and Ge Wang. SPICE: Semantic pseudo-labeling for image cluster-
725 ing. *IEEE Transactions on Image Processing*, 31:7264–7278, 2022.
- 726 Foivos Ntelemis, Yaochu Jin, and Spencer A Thomas. Information maximization clustering via
727 multi-view self-labelling. *Knowledge-Based Systems*, 250:109042, 2022.
- 728 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
729 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
730 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 731 Sungwon Park, Sungwon Han, Sundong Kim, Danu Kim, Sungkyu Park, Seunghoon Hong, and
732 Meeyoung Cha. Improving unsupervised image clustering with robust learning. In *Proceedings*
733 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12278–12287,
734 2021.
- 735 Vishal M Patel and René Vidal. Kernel sparse subspace clustering. pp. 2849–2853, 2014.
- 736 Vishal M Patel, Hien Van Nguyen, and René Vidal. Latent space sparse subspace clustering. In
737 *Proceedings of the IEEE international conference on computer vision*, pp. 225–232, 2013.
- 738 Vishal M Patel, Hien Van Nguyen, and René Vidal. Latent space sparse and low-rank subspace
739 clustering. *IEEE Journal of Selected Topics in Signal Processing*, 9(4):691–701, 2015.
- 740 Xi Peng, Shijie Xiao, Jiashi Feng, Wei-Yun Yau, and Zhang Yi. Deep subspace clustering with
741 sparsity prior. In *International Joint Conference on Artificial Intelligence*, pp. 1925–1931, 2016.
- 742 Xi Peng, Jiashi Feng, Shijie Xiao, Wei-Yun Yau, Joey Tianyi Zhou, and Songfan Yang. Structured
743 autoencoders for subspace clustering. *IEEE Transactions on Image Processing*, 27(10):5076–
744 5086, 2018.
- 745 Xi Peng, Jiashi Feng, Joey Tianyi Zhou, Yingjie Lei, and Shuicheng Yan. Deep subspace clustering.
746 *IEEE Transactions on Neural Networks and Learning Systems*, 31(12):5509–5521, 2020.
- 747 Daniel Pimentel-Alarcon and Robert Nowak. The information-theoretic requirements of subspace
748 clustering with missing data. In *International Conference on Machine Learning*, pp. 802–810,
749 2016.

- 756 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
757 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
758 models from natural language supervision. In *International conference on machine learning*, pp.
759 8748–8763. PMLR, 2021.
- 760 Shankar Rao, Roberto Tron, René Vidal, and Yi Ma. Motion segmentation in the presence of outly-
761 ing, incomplete, or corrupted trajectories. 32(10):1832–1845, 2010.
- 762
- 763 Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on*
764 *pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- 765
- 766 Mahdi Soltanolkotabi and Emmanuel J Candes. A geometric analysis of subspace clustering with
767 outliers. *Annals of Statistics*, 40(4):2195–2238, 2012.
- 768
- 769 Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics
770 without contrastive pairs. In *International Conference on Machine Learning*, pp. 10268–10278.
771 PMLR, 2021.
- 772
- 773 Manolis Tsakiris and René Vidal. Algebraic clustering of affine subspaces. *IEEE Transactions on*
774 *Pattern Analysis and Machine Intelligence*, 2017.
- 775
- 776 Manolis Tsakiris and René Vidal. Theoretical analysis of sparse subspace clustering with missing
777 entries. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 4975–
778 4984, 2018.
- 779
- 780 Paul Tseng. Nearest q-flat to m points. *Journal of Optimization Theory and Applications*, 105:
781 249–252, 2000.
- 782
- 783 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
784 *learning research*, 9(11), 2008.
- 785
- 786 Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc
787 Van Gool. SCAN: Learning to classify images without labels. In *European Conference on Com-*
788 *puter Vision*, pp. 268–285, 2020.
- 789
- 790 René Vidal. Identification of PWARX hybrid models with unknown and possibly different orders.
791 pp. 547–552, 2004.
- 792
- 793 René Vidal, Yi Ma, and Shankar Sastry. Generalized Principal Component Analysis (GPCA). 27
794 (12):1–15, 2005.
- 795
- 796 René Vidal, Roberto Tron, and Richard Hartley. Multiframe motion segmentation with missing data
797 using PowerFactorization, and GPCA. 79(1):85–105, 2008.
- 798
- 799 Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416,
800 2007.
- 801
- 802 Libin Wang, Yulong Wang, Hao Deng, and Hong Chen. Attention reweighted sparse subspace
803 clustering. *Pattern Recognition*, 139:109438, 2023a.
- 804
- 805 Peng Wang, Huikang Liu, Druv Pai, Yaodong Yu, Zhihui Zhu, Qing Qu, and Yi Ma. A global
806 geometric analysis of maximal coding rate reduction. In *Forty-first International Conference on*
807 *Machine Learning*, 2024.
- 808
- 809 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.

- 810 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-
811 ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 812
- 813 Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis.
814 In *International conference on machine learning*, pp. 478–487. PMLR, 2016.
- 815 Jianwei Yang, Devi Parikh, and Dhruv Batra. Joint unsupervised learning of deep representations
816 and image clusters. In *Proceedings of the IEEE conference on computer vision and pattern recog-
817 nition*, pp. 5147–5156, 2016.
- 818
- 819 Chong You, Chun-Guang Li, Daniel Robinson, and René Vidal. Oracle based active set algorithm
820 for scalable elastic net subspace clustering. In *IEEE Conference on Computer Vision and Pattern
821 Recognition*, pp. 3928–3937, 2016a.
- 822 Chong You, Daniel Robinson, and René Vidal. Scalable sparse subspace clustering by orthogonal
823 matching pursuit. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3918–
824 3927, 2016b.
- 825
- 826 Wei Yu, Wonjong Rhee, Stephen Boyd, and John M Cioffi. Iterative water-filling for gaussian vector
827 multiple-access channels. *IEEE Transactions on Information Theory*, 50(1):145–152, 2004.
- 828 Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning diverse
829 and discriminative representations via the principle of maximal coding rate reduction. In *Neural
830 Information Processing Systems (NIPS)*, 2020.
- 831 Pengxin Zeng, Yunfan Li, Peng Hu, Dezhong Peng, Jiancheng Lv, and Xi Peng. Deep fair cluster-
832 ing via maximizing and minimizing mutual information: Theory, algorithm and metric. In *Pro-
833 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23986–
834 23995, 2023.
- 835
- 836 Hongjing Zhang and Ian Davidson. Deep fair discriminative clustering. *arXiv preprint
837 arXiv:2105.14146*, 2021.
- 838 Junjian Zhang, Chun-Guang Li, Chong You, Xianbiao Qi, Honggang Zhang, Jun Guo, and
839 Zhouchen Lin. Self-supervised convolutional subspace clustering network. In *Proceedings of
840 the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5473–5482, 2019a.
- 841
- 842 Shangzhi Zhang, Chong You, René Vidal, and Chun-Guang Li. Learning a self-expressive network
843 for subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
844 Pattern Recognition (CVPR)*, pp. 12393–12403, June 2021.
- 845
- 846 Teng Zhang, Arthur Szlam, and Gilad Lerman. Median k-flats for hybrid linear modeling with
847 many outliers. In *2009 IEEE 12th International Conference on Computer Vision Workshops,
848 ICCV Workshops*, pp. 234–241. IEEE, 2009.
- 849
- 849 Tong Zhang, Pan Ji, Mehrtash Harandi, Wenbing Huang, and Hongdong Li. Neural collaborative
850 subspace clustering. In *International Conference on Machine learning*, pp. 7384–7393, 2019b.
- 851 Huasong Zhong, Jianlong Wu, Chong Chen, Jianqiang Huang, Minghua Deng, Liqiang Nie,
852 Zhouchen Lin, and Xian-Sheng Hua. Graph contrastive clustering. In *Proceedings of the
853 IEEE/CVF international conference on computer vision*, pp. 9224–9233, 2021.
- 854
- 855 Pan Zhou, Yunqing Hou, and Jiashi Feng. Deep adversarial subspace clustering. In *Proceedings of
856 the IEEE conference on computer vision and pattern recognition*, pp. 1596–1604, 2018.
- 857
- 858
- 859
- 860
- 861
- 862
- 863

SUPPLEMENTARY MATERIAL FOR “EXPLORING A PRINCIPLED DEEP SUBSPACE CLUSTERING NETWORK”

The supplementary materials are divided into two 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).

A PROOFS OF MAIN RESULTS

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

Proof. We note that:

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

where $\mathbf{z}^{(i)}$ is the i^{th} row of \mathbf{Z} , thus problem (2) is equivalent to:

$$\begin{aligned} \min_{\{\mathbf{z}^{(i)}\}_{i=1}^d, \mathbf{C}} \quad & \frac{1}{2} \sum_{i=1}^d \mathbf{z}^{(i)} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top \mathbf{z}^{(i)\top} + \beta \cdot r(\mathbf{C}) \\ \text{s.t.} \quad & \|\mathbf{Z}\|_F^2 = N. \end{aligned} \quad (13)$$

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

$$\begin{aligned} \min_{\{\mathbf{z}^{(i)}\}_{i=1}^d, \mathbf{C}} \quad & \frac{1}{2} \sum_{i=1}^d \mathbf{z}^{(i)} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top \mathbf{z}^{(i)\top} + \beta \cdot r(\mathbf{C}) \\ \text{s.t.} \quad & \|\mathbf{z}^{(i)}\|_2^2 = \tau_i, \quad i = 1, \dots, d. \end{aligned} \quad (14)$$

The Lagrangian of problem (14) is:

$$\mathcal{L}(\{\mathbf{z}^{(i)}\}_{i=1}^d, \mathbf{C}, \{\nu_i\}_{i=1}^d) := \frac{1}{2} \sum_{i=1}^d \mathbf{z}^{(i)} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top \mathbf{z}^{(i)\top} + \beta \cdot r(\mathbf{C}) + \frac{1}{2} \sum_{i=1}^d \nu_i (\|\mathbf{z}^{(i)}\|_2^2 - \tau_i),$$

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

$$\begin{cases} \nabla_{\mathbf{z}^{(i)}} \mathcal{L} = \mathbf{z}^{(i)} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top + \nu_i \mathbf{z}^{(i)} = \mathbf{0}, \\ \|\mathbf{z}^{(i)}\|_2^2 = \tau_i, \quad i = 1, \dots, d, \end{cases} \quad (16)$$

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

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

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

We first present two lemmas, which will be used for the proof of Theorem 1.

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

Lemma A2 (Tian et al., 2021). *Let $\mathbf{H}(t) \succ \mathbf{0}$ be a time-varying positive definite matrix whose minimal eigenvalue is bounded away from 0, i.e., $\inf_{t \geq 0} \lambda_{\min}(\mathbf{H}(t)) \geq \lambda_0 > 0$. Then, the following dynamics*

$$\frac{d\mathbf{w}(t)}{dt} = -\mathbf{H}(t)\mathbf{w}(t)$$

satisfies that $\|\mathbf{w}(t)\|_2 \leq e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$, i.e., $\mathbf{w}(t) \rightarrow \mathbf{0}$.

Theorem 1. *Consider the objective function (5) where \mathbf{Z} and \mathbf{C} are optimized by gradient descent with a learning rate $\eta \rightarrow 0$, denote $\mathbf{G} := \mathbf{Z}^\top \mathbf{Z}$, $\mathbf{M} := (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$, $\mathbf{F} := \alpha(\mathbf{I} + \alpha\mathbf{G})^{-1}$, $\hat{\mathbf{M}} := \gamma\mathbf{M} + \nu_*\mathbf{I}$. If $\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) < 0$, and $\frac{8\alpha^2}{27} < \gamma < \frac{\alpha - \nu_*}{\lambda_{\max}(\hat{\mathbf{M}})}$, then \mathbf{G} and \mathbf{M} will converge to have eigenspaces aligned, i.e., \mathbf{G} and \mathbf{M} can be diagonalized simultaneously by $\mathbf{U} \in \mathcal{O}^{N \times N}$ where $\mathcal{O}^{N \times N}$ is an $N \times N$ orthogonal matrix group.*

Proof. As $\mathbf{F} := \alpha(\mathbf{I} + \alpha\mathbf{G})^{-1}$, $\hat{\mathbf{M}} := \gamma\mathbf{M} + \nu_*\mathbf{I}$, we have $\mathbf{G} = \mathbf{Z}^\top \mathbf{Z} = \mathbf{F}^{-1} - \frac{1}{\alpha}\mathbf{I}$ and $\mathbf{M} = \frac{\hat{\mathbf{M}} - \nu_*\mathbf{I}}{\gamma}$. For simplicity, we denote $\mathbf{H} := \mathbf{F} - \hat{\mathbf{M}} = -\gamma\mathbf{M} + \mathbf{F} - \nu_*\mathbf{I}$.

As \mathbf{Z} and \mathbf{C} are optimized by gradient descent with a learning rate $\eta \rightarrow 0$, we have the learning dynamic:

$$\dot{\mathbf{Z}} = \frac{d\mathbf{Z}}{dt} = -\nabla_{\mathbf{Z}}\mathcal{L} = -\gamma\mathbf{Z}\mathbf{M} + \alpha\mathbf{Z}(\mathbf{I} + \alpha\mathbf{G})^{-1} - \nu_*\mathbf{Z} = \mathbf{Z}\mathbf{H}, \quad (17)$$

$$\dot{\mathbf{C}} = \frac{d\mathbf{C}}{dt} = -\nabla_{\mathbf{C}}\mathcal{L} = \gamma\mathbf{G}(\mathbf{I} - \mathbf{C}) = \gamma(\mathbf{F}^{-1} - \frac{1}{\alpha}\mathbf{I})(\mathbf{I} - \mathbf{C}). \quad (18)$$

As \mathbf{F} and $\hat{\mathbf{M}}$ are functions of \mathbf{Z} and \mathbf{C} , respectively, they are updated by:

$$\begin{aligned} \dot{\mathbf{F}} &= \frac{d\mathbf{F}}{dt} = \frac{d}{dt} \left[\alpha(\mathbf{I} + \alpha\mathbf{Z}^\top \mathbf{Z})^{-1} \right] \\ &= -\alpha(\mathbf{I} + \alpha\mathbf{G})^{-1} \frac{d}{dt} (\mathbf{Z}^\top \mathbf{Z}) \alpha(\mathbf{I} + \alpha\mathbf{G})^{-1} \\ &= -\mathbf{F} \left(\dot{\mathbf{Z}}^\top \mathbf{Z} + \mathbf{Z}^\top \dot{\mathbf{Z}} \right) \mathbf{F} \\ &= -\mathbf{F} (\mathbf{H}\mathbf{G} + \mathbf{G}\mathbf{H}) \mathbf{F} \\ &= -\mathbf{F} \left(\mathbf{H}\mathbf{F}^{-1} + \mathbf{F}^{-1}\mathbf{H} - \frac{2}{\alpha}\mathbf{H} \right) \mathbf{F} \\ &= \left(\mathbf{F}\hat{\mathbf{M}} + \hat{\mathbf{M}}\mathbf{F} \right) - 2 \left(\mathbf{F}^2 + \frac{1}{\alpha}\mathbf{F}\hat{\mathbf{M}}\mathbf{F} - \frac{1}{\alpha}\mathbf{F}^3 \right), \end{aligned} \quad (19)$$

$$\begin{aligned} \dot{\hat{\mathbf{M}}} &= \frac{d(\gamma\mathbf{M} + \nu_*\mathbf{I})}{dt} = \gamma \frac{d\mathbf{M}}{dt} = \gamma(-\dot{\mathbf{C}})(\mathbf{I} - \mathbf{C})^\top + \gamma(\mathbf{I} - \mathbf{C})(-\dot{\mathbf{C}})^\top \\ &= -\gamma^2 \left[(\mathbf{F}^{-1} - \frac{1}{\alpha}\mathbf{I})\mathbf{M} + \mathbf{M}(\mathbf{F}^{-1} - \frac{1}{\alpha}\mathbf{I}) \right] \\ &= -\gamma^2 \left(\mathbf{F}^{-1}\mathbf{M} + \mathbf{M}\mathbf{F}^{-1} - \frac{2}{\alpha}\mathbf{M} \right) \\ &= -\gamma \left[(\mathbf{F}^{-1}\hat{\mathbf{M}} + \hat{\mathbf{M}}\mathbf{F}^{-1}) - 2(\nu_*\mathbf{F}^{-1} + \frac{\hat{\mathbf{M}}}{\alpha} - \frac{\nu_*}{\alpha}\mathbf{I}) \right]. \end{aligned} \quad (20)$$

Next, we denote \mathbf{L} as the commutator $[\mathbf{F}, \hat{\mathbf{M}}]$

$$\mathbf{L} := [\mathbf{F}, \hat{\mathbf{M}}] = \mathbf{F}\hat{\mathbf{M}} - \hat{\mathbf{M}}\mathbf{F}, \quad (21)$$

and we notice that:

$$\begin{aligned} \mathbf{F}\hat{\mathbf{M}}^2 - \hat{\mathbf{M}}^2\mathbf{F} &= \hat{\mathbf{M}}\mathbf{L} + \mathbf{L}\hat{\mathbf{M}}, \\ \mathbf{F}^2\hat{\mathbf{M}} - \hat{\mathbf{M}}\mathbf{F}^2 &= \mathbf{F}\mathbf{L} + \mathbf{L}\mathbf{F}, \\ \mathbf{F}^3\hat{\mathbf{M}} - \hat{\mathbf{M}}\mathbf{F}^3 + \hat{\mathbf{M}}\mathbf{F}\hat{\mathbf{M}}\mathbf{F} - \mathbf{F}\hat{\mathbf{M}}\mathbf{F}\hat{\mathbf{M}} &= (\mathbf{F}^2 - \hat{\mathbf{M}}\mathbf{F})\mathbf{L} + \mathbf{L}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) + \mathbf{F}\mathbf{L}\mathbf{F}, \\ \mathbf{F}\hat{\mathbf{M}}\mathbf{F}^{-1} - \mathbf{F}^{-1}\hat{\mathbf{M}}\mathbf{F} &= \mathbf{F}^{-1}\mathbf{L} + \mathbf{L}\mathbf{F}^{-1}. \end{aligned}$$

We compute the dynamic of \mathbf{L} as:

$$\dot{\mathbf{L}} = \frac{d\mathbf{L}}{dt} = \underbrace{(\dot{\mathbf{F}}\hat{\mathbf{M}} - \hat{\mathbf{M}}\dot{\mathbf{F}})}_{\mathbf{L}_1} + \underbrace{(\mathbf{F}\dot{\mathbf{M}} - \dot{\mathbf{M}}\mathbf{F})}_{\mathbf{L}_2}, \text{ where} \quad (22)$$

$$\begin{aligned} \mathbf{L}_1 &= (\mathbf{F}\hat{\mathbf{M}} + \hat{\mathbf{M}}\mathbf{F})\dot{\mathbf{M}} - 2(\mathbf{F}^2 + \frac{1}{\alpha}\mathbf{F}\hat{\mathbf{M}}\mathbf{F} - \frac{1}{\alpha}\mathbf{F}^3)\dot{\mathbf{M}} \\ &\quad - \hat{\mathbf{M}}(\mathbf{F}\dot{\mathbf{M}} + \dot{\mathbf{M}}\mathbf{F}) + 2\hat{\mathbf{M}}(\mathbf{F}^2 + \frac{1}{\alpha}\mathbf{F}\hat{\mathbf{M}}\mathbf{F} - \frac{1}{\alpha}\mathbf{F}^3) \\ &= (\mathbf{F}\hat{\mathbf{M}}^2 - \hat{\mathbf{M}}^2\mathbf{F}) - 2(\mathbf{F}^2\dot{\mathbf{M}} - \dot{\mathbf{M}}\mathbf{F}^2) + \frac{2}{\alpha}[(\mathbf{F}^3\dot{\mathbf{M}} - \dot{\mathbf{M}}\mathbf{F}^3) - (\mathbf{F}\hat{\mathbf{M}}\mathbf{F}\dot{\mathbf{M}} - \dot{\mathbf{M}}\mathbf{F}\hat{\mathbf{M}}\mathbf{F})] \\ &= (\hat{\mathbf{M}}\mathbf{L} + \mathbf{L}\hat{\mathbf{M}}) - 2(\mathbf{F}\mathbf{L} + \mathbf{L}\mathbf{F}) + \frac{2}{\alpha}[(\mathbf{F}^2 - \hat{\mathbf{M}}\mathbf{F})\mathbf{L} + \mathbf{L}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) + \mathbf{F}\mathbf{L}\mathbf{F}], \end{aligned} \quad (23)$$

$$\begin{aligned} \mathbf{L}_2 &= -\gamma\mathbf{F}\left[(\mathbf{F}^{-1}\hat{\mathbf{M}} + \hat{\mathbf{M}}\mathbf{F}^{-1}) - 2(\nu_*\mathbf{F}^{-1} + \frac{\hat{\mathbf{M}}}{\alpha} - \frac{\nu_*}{\alpha}\mathbf{I})\right] \\ &\quad + \gamma\left[(\mathbf{F}^{-1}\hat{\mathbf{M}} + \hat{\mathbf{M}}\mathbf{F}^{-1}) - 2(\nu_*\mathbf{F}^{-1} + \frac{\hat{\mathbf{M}}}{\alpha} - \frac{\nu_*}{\alpha}\mathbf{I})\right]\mathbf{F} \\ &= \gamma(\mathbf{F}^{-1}\hat{\mathbf{M}}\mathbf{F} - \mathbf{F}\hat{\mathbf{M}}\mathbf{F}^{-1}) + \frac{2\gamma}{\alpha}(\mathbf{F}\hat{\mathbf{M}} - \hat{\mathbf{M}}\mathbf{F}) \\ &= -\gamma(\mathbf{F}^{-1}\mathbf{L} + \mathbf{L}\mathbf{F}^{-1}) + \frac{2\gamma}{\alpha}\mathbf{L}. \end{aligned} \quad (24)$$

Therefore, we vectorize $\dot{\mathbf{L}}$ by Kronecker product:

$$\text{Vec}(\dot{\mathbf{L}}) = \text{Vec}(\mathbf{L}_1) + \text{Vec}(\mathbf{L}_2) = -(\mathbf{K}_1 + \mathbf{K}_2)\text{Vec}(\mathbf{L}), \quad (25)$$

$$\text{where } \mathbf{K}_1 := -\hat{\mathbf{M}} \oplus \hat{\mathbf{M}} + 2\mathbf{F} \oplus \mathbf{F} - \frac{2}{\alpha}[(\mathbf{F}^2 - \hat{\mathbf{M}}\mathbf{F}) \oplus (\mathbf{F}^2 - \hat{\mathbf{M}}\mathbf{F}) + \mathbf{F} \otimes \mathbf{F}], \quad (26)$$

$$\mathbf{K}_2 := \gamma\mathbf{F}^{-1} \oplus \mathbf{F}^{-1} - \frac{2\gamma}{\alpha}\mathbf{I}, \text{ and denote } \mathbf{W}_1 \oplus \mathbf{W}_2 := \mathbf{W}_1 \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{W}_2. \quad (27)$$

The smallest eigenvalue of $\mathbf{K} := \mathbf{K}_1 + \mathbf{K}_2$ satisfies:

$$\lambda_{\min}(\mathbf{K}) \geq -2\lambda_{\max}(\hat{\mathbf{M}}) + 4\lambda_{\min}(\mathbf{F} - \frac{1}{2\alpha}\mathbf{F}^2 + \frac{\gamma}{2}\mathbf{F}^{-1}) - \frac{4}{\alpha}\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) - \frac{2\gamma}{\alpha}. \quad (28)$$

Given that $\lambda(\mathbf{F}) = \lambda(\alpha(\mathbf{I} + \alpha\mathbf{G})^{-1}) = \frac{\alpha}{1+\alpha\lambda(\mathbf{G})} \in (0, \alpha]$, we denote $f(\lambda) := \lambda - \frac{\lambda^2}{2\alpha} + \frac{\gamma}{2\lambda}$, $\lambda \in (0, \alpha]$, and its derivative $f'(\lambda) = \frac{1}{\lambda^2}(2\lambda^2 - \frac{2}{\alpha}\lambda^3 - \gamma)$. We notice that when $\gamma > \frac{8\alpha^2}{27}$, $f'(\lambda) < 0$ for all λ and α , which means $f(\lambda)$ is monotonically decreasing in $\lambda \in [0, \alpha]$. Therefore,

$$\lambda_{\min}(\mathbf{K}) \geq 2\left[-\lambda_{\max}(\hat{\mathbf{M}}) + 2\lambda_{\max}(\mathbf{F}) - \frac{1}{\alpha}\lambda_{\max}^2(\mathbf{F}) + \frac{\gamma}{\lambda_{\max}(\mathbf{F})} - \frac{2}{\alpha}\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) - \frac{\gamma}{\alpha}\right]. \quad (29)$$

Given that $\lambda_{\max}(\hat{\mathbf{M}}) = \lambda_{\max}(\gamma\mathbf{M} + \nu_*\mathbf{I}) = \gamma\lambda_{\max}(\mathbf{M}) + \nu_*$, we have:

$$\begin{aligned} \lambda_{\min}(\mathbf{K}) &\geq 2\left[-\lambda_{\max}(\hat{\mathbf{M}}) + 2\lambda_{\max}(\mathbf{F}) - \frac{1}{\alpha}\lambda_{\max}^2(\mathbf{F}) + \frac{\gamma}{\lambda_{\max}(\mathbf{F})} - \frac{2}{\alpha}\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) - \frac{\gamma}{\alpha}\right] \\ &= 2\left[-\gamma\lambda_{\max}(\mathbf{M}) - \nu_* + 2\alpha - \alpha + \frac{\gamma}{\alpha} - \frac{\gamma}{\alpha} - \frac{2}{\alpha}\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}})\right] \\ &= 2\left[-\gamma\lambda_{\max}(\mathbf{M}) - \nu_* + \alpha - \frac{2}{\alpha}\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}})\right]. \end{aligned} \quad (30)$$

Given that $\nu_* < \alpha - \gamma\lambda_{\max}(\mathbf{M})$, $\lambda_{\max}(\mathbf{F}^2 - \mathbf{F}\hat{\mathbf{M}}) < 0$, we have $\lambda_{\min}(\mathbf{K}) > 0$.

Therefore, we have:

$$\text{Vec}(\dot{\mathbf{L}}) = \frac{d\text{Vec}(\mathbf{L})}{dt} = -\mathbf{K}\text{Vec}(\mathbf{L}), \text{ where } \mathbf{K} \succ \mathbf{0}. \quad (31)$$

1026 According to Lemma A2, this implies:

$$1027 \lim_{t \rightarrow \infty} \mathbf{L} = \lim_{t \rightarrow \infty} \mathbf{F} \hat{\mathbf{M}} - \hat{\mathbf{M}} \mathbf{F} = \mathbf{0}. \quad (32)$$

1028 By Lemma A1, we find that $\mathbf{F} = \alpha(\mathbf{I} + \alpha\mathbf{G})^{-1}$ and $\hat{\mathbf{M}} = \gamma\mathbf{M} + \nu_*\mathbf{I}$ can be diagonalized
1029 simultaneously by $\mathbf{U} \in \mathcal{O}^{N \times N}$. At the same time, \mathbf{F} with \mathbf{G} , $\hat{\mathbf{M}}$ with \mathbf{M} share an eigenspace,
1030 which means that $\mathbf{G} = \mathbf{Z}^\top \mathbf{Z}$ and $\mathbf{M} = (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$ can be diagonalized simultaneously by
1031 $\mathbf{U} \in \mathcal{O}^{N \times N}$.

1032 \square

1033 To justify the fact that actually we can ensure that ν_* satisfies the condition $\gamma < \frac{1}{\lambda_{\max}(\mathbf{M}_*)}(\alpha - \nu_*)$
1034 by adjusting the hyper-parameters α and γ , we will derive the optimal Lagrangian multiplier ν_* from
1035 the optimality condition of the problem in Eq.(5).

1036 We begin with analyzing the KKT condition of the problem in Eq.(5):

$$1037 \begin{cases} \nabla_{\mathbf{Z}} \mathcal{L}(\mathbf{Z}_*, \mathbf{C}_*, \nu_*) = \gamma \mathbf{Z}_* \mathbf{M}_* - \alpha \mathbf{Z}_* (\mathbf{I} + \alpha \mathbf{G}_*)^{-1} + \nu_* \mathbf{Z}_* = \mathbf{0}, & (33) \\ \nabla_{\mathbf{C}} \mathcal{L}(\mathbf{Z}_*, \mathbf{C}_*, \nu_*) = -\gamma \mathbf{G}_* (\mathbf{I} - \mathbf{C}_*) = \mathbf{0}, & (34) \\ \nabla_{\nu} \mathcal{L}(\mathbf{Z}_*, \mathbf{C}_*, \nu_*) = \|\mathbf{Z}_*\|_F^2 - N = \text{Tr}(\mathbf{G}_*) - N = 0, & (35) \end{cases}$$

1038 where $\mathbf{G}_* = \mathbf{Z}_*^\top \mathbf{Z}_*$ and $\mathbf{M}_* = (\mathbf{I} - \mathbf{C}_*)(\mathbf{I} - \mathbf{C}_*)^\top$. The condition can be rewritten as:

$$1039 \begin{cases} \gamma \mathbf{G}_* \mathbf{M}_* - \alpha \mathbf{G}_* (\mathbf{I} + \alpha \mathbf{G}_*)^{-1} + \nu_* \mathbf{G}_* = \mathbf{0} & (36) \\ \gamma \mathbf{G}_* \mathbf{M}_* = \mathbf{0} & (37) \\ \text{Tr}(\mathbf{G}_*) = N & (38) \end{cases}$$

1040 Then, we add the trace operator to (36) and simplify the formula by (37) as:

$$1041 \text{Tr}(\nu_* \mathbf{G}_* - \alpha \mathbf{G}_* (\mathbf{I} + \alpha \mathbf{G}_*)^{-1}) = 0 \quad (39)$$

1042 Since $\alpha \mathbf{G}_* (\mathbf{I} + \alpha \mathbf{G}_*)^{-1} = \mathbf{I} - (\mathbf{I} + \alpha \mathbf{G}_*)^{-1}$ and $\mathbf{Z} \in \mathbb{R}^{d \times N}$, we discuss two cases between N
1043 and d .

1044 **For case 1.** Suppose $N > d$, then (39) will be rewritten as:

$$1045 N\nu_* = \text{Tr}(\mathbf{I} - (\mathbf{I} + \alpha \mathbf{G}_*)^{-1}) \quad (40)$$

$$1046 \Leftrightarrow N\nu_* = N - (N - d) - \sum_{i=1}^d \frac{1}{1 + \alpha \lambda_{\mathbf{G}_*}^{(i)}}$$

$$1047 \Leftrightarrow \nu_* = \frac{d}{N} - \frac{1}{N} \sum_{i=1}^d \frac{1}{1 + \alpha \lambda_{\mathbf{G}_*}^{(i)}}, \quad (41)$$

1048 where $\lambda_{\mathbf{G}_*}^{(i)} \geq 0, i = 1, \dots, d$, $\lambda_{\mathbf{G}_*}^{(i)} = 0, i = d + 1, \dots, N$ and $\sum_{i=1}^d \lambda_{\mathbf{G}_*}^{(i)} = N$. To ensure that ν_*
1049 satisfies the condition of Theorem 1, i.e., $\nu_* < \alpha - \gamma \lambda_{\max}(\mathbf{M})$ for any $\{\lambda_{\mathbf{G}_*}^{(i)}\}_{i=1}^N$, we have:

$$1050 \nu_* = \frac{d}{N} - \frac{1}{N} \sum_{i=1}^d \frac{1}{1 + \alpha \lambda_{\mathbf{G}_*}^{(i)}} \leq \frac{\alpha}{1 + \alpha \cdot \frac{N}{d}} < \alpha - \gamma \lambda_{\max}(\mathbf{M}), \quad (42)$$

1051 where the equality holds when $\alpha > 0, \lambda_{\mathbf{G}_*}^{(1)} = \dots = \lambda_{\mathbf{G}_*}^{(d)} = \frac{N}{d}$. Therefore, ν_* will satisfy the
1052 condition in Theorem 1 if the last inequality holds, which is equivalent to:

$$1053 \gamma \lambda_{\max}(\mathbf{M}) < \alpha - \frac{\alpha}{1 + \alpha \cdot \frac{N}{d}} = \frac{\alpha^2}{\frac{d}{N} + \alpha}. \quad (43)$$

1054 **For case 2.** Suppose $d \geq N$, then with the similar process, we estimate the upper bound of ν_* as:

$$1055 \nu_* = 1 - \frac{1}{N} \sum_{i=1}^N \frac{1}{1 + \alpha \lambda_{\mathbf{G}_*}^{(i)}} \leq \frac{\alpha}{1 + \alpha} < \alpha - \gamma \lambda_{\max}(\mathbf{M}) \quad (44)$$

where the equality holds when $\alpha > 0, \lambda_{\mathbf{G}_*}^{(1)} = \dots = \lambda_{\mathbf{G}_*}^{(N)} = 1$. Therefore, ν_* will satisfy the condition in Theorem 2 if the last inequality holds, which is equivalent to:

$$\gamma \lambda_{\max}(\mathbf{M}) < \alpha - \frac{\alpha}{1 + \alpha} = \frac{\alpha^2}{1 + \alpha}. \quad (45)$$

This means that when γ and α satisfy Eq. (43) or (45) in each cases, the optimal ν_* will satisfy Theorem 2's condition. Therefore, we may not need to concern whether ν_* satisfies the condition, but focus on the hyper-parameters γ and α .

Theorem 2. *Suppose that \mathbf{G} and \mathbf{M} have eigenspaces aligned and $\gamma < (\alpha - \nu_*)/\lambda_{\max}(\mathbf{M})$, then for the optimal solution \mathbf{Z}_* , we have that $\text{rank}(\mathbf{Z}) = \min\{d, N\}$ and the singular values $\sigma_{\mathbf{Z}_*}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_{\mathbf{M}}^{(i)} + \nu_*} - \frac{1}{\alpha}}$, for all $i = 1, \dots, \min\{d, N\}$, where ν_* is the dual optimal solution.*

Proof. Since $\|\mathbf{Z} - \mathbf{Z}\mathbf{C}\|_F^2 = \text{Tr}(\mathbf{Z}^\top \mathbf{Z}(\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top)$ and $\|\mathbf{Z}\|_F^2 = \text{Tr}(\mathbf{Z}^\top \mathbf{Z})$, problem (6) is equivalent to:

$$\begin{aligned} \min_{\mathbf{G}} \quad & -\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{G}) + \frac{\gamma}{2} \text{Tr}(\mathbf{G}\mathbf{M}) \\ \text{s.t.} \quad & \text{Tr}(\mathbf{G}) = N, \mathbf{G} \succeq \mathbf{0}, \end{aligned} \quad (46)$$

where $\mathbf{G} := \mathbf{Z}^\top \mathbf{Z}$, $\mathbf{M} := (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top$.

Since that \mathbf{G} and \mathbf{M} have eigenspaces aligned, we have $\mathbf{G} = \mathbf{U}\mathbf{\Lambda}_{\mathbf{G}}\mathbf{U}^\top$, $\mathbf{M} = \mathbf{U}\mathbf{\Lambda}_{\mathbf{M}}\mathbf{U}^\top$. Therefore, problem (46) is transformed into:

$$\begin{aligned} \min_{\mathbf{\Lambda}_{\mathbf{G}}} \quad & -\frac{1}{2} \text{Tr} \log(\mathbf{I} + \alpha \mathbf{\Lambda}_{\mathbf{G}}) + \frac{\gamma}{2} \text{Tr}(\mathbf{\Lambda}_{\mathbf{G}}\mathbf{\Lambda}_{\mathbf{M}}) \\ \text{s.t.} \quad & \text{Tr}(\mathbf{\Lambda}_{\mathbf{G}}) = N, \mathbf{\Lambda}_{\mathbf{G}} \succeq \mathbf{0}, \end{aligned} \quad (47)$$

which is further equivalent to the eigenvalue optimization problem below:

$$\begin{aligned} \min_{\{\lambda_{\mathbf{G}}^{(i)}\}_{i=1}^{\min\{d, N\}}} \quad & -\frac{1}{2} \sum_{i=1}^{\min\{d, N\}} \log(1 + \alpha \lambda_{\mathbf{G}}^{(i)}) + \frac{\gamma}{2} \lambda_{\mathbf{M}}^{(i)} \lambda_{\mathbf{G}}^{(i)} \\ \text{s.t.} \quad & \sum_{i=1}^{\min\{d, N\}} \lambda_{\mathbf{G}}^{(i)} = N, \\ & \lambda_{\mathbf{G}}^{(i)} \geq 0, \quad \text{for all } i = 1, \dots, \min\{d, N\}. \end{aligned} \quad (48)$$

It is noteworthy that problem (48) is a convex problem. Thus, the KKT condition is the sufficient and necessary condition for the minimizer.

The Lagrangian of problem (48) is:

$$\begin{aligned} \mathcal{L}(\{\lambda_{\mathbf{G}}^{(i)}\}_{i=1}^{\min\{d, N\}}, \{\mu_i\}_{i=1}^{\min\{d, N\}}, \nu) := \\ -\frac{1}{2} \sum_{i=1}^{\min\{d, N\}} \log(1 + \alpha \lambda_{\mathbf{G}}^{(i)}) + \frac{\gamma}{2} \lambda_{\mathbf{M}}^{(i)} \lambda_{\mathbf{G}}^{(i)} - \mu_i \lambda_{\mathbf{G}}^{(i)} + \frac{\nu}{2} \left(\sum_{i=1}^{\min\{d, N\}} \lambda_{\mathbf{G}}^{(i)} - N \right), \end{aligned} \quad (49)$$

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

$$\begin{cases} \nabla_{\lambda_{\mathbf{G}_*}^{(i)}} \mathcal{L} = 0, & \forall i = 1, \dots, \min\{d, N\}, \end{cases} \quad (50)$$

$$\begin{cases} \lambda_{\mathbf{G}_*}^{(i)} \geq 0, & \forall i = 1, \dots, \min\{d, N\}, \end{cases} \quad (51)$$

$$\begin{cases} \sum_{i=1}^{\min\{d, N\}} \lambda_{\mathbf{G}_*}^{(i)} = N, \end{cases} \quad (52)$$

$$\begin{cases} \mu_{i_*} \geq 0, & \forall i = 1, \dots, \min\{d, N\}, \end{cases} \quad (53)$$

$$\begin{cases} \mu_{i_*} \lambda_{\mathbf{G}_*}^{(i)} = 0, & \forall i = 1, \dots, \min\{d, N\}. \end{cases} \quad (54)$$

Then, (50) is equivalent to:

$$\mu_{i_*} = \frac{1}{2} \left(\nu_* + \gamma \lambda_M^{(i)} - \frac{\alpha}{1 + \alpha \lambda_{\mathbf{G}_*}^{(i)}} \right). \quad (55)$$

By Eqs. (51) and (53-54)(55), we come up with the following two cases:

$$\begin{cases} \mu_{i_*} > 0 \Rightarrow \lambda_{\mathbf{G}_*}^{(i)} = 0, \frac{1}{\nu_* + \gamma \lambda_M^{(i)}} - \frac{1}{\alpha} < 0, \\ \mu_{i_*} = 0 \Rightarrow \lambda_{\mathbf{G}_*}^{(i)} > 0, \lambda_{\mathbf{G}_*}^{(i)} = \frac{1}{\nu_* + \gamma \lambda_M^{(i)}} - \frac{1}{\alpha} > 0. \end{cases} \quad (56)$$

From the above two cases, we conclude that:

$$\lambda_{\mathbf{G}_*}^{(i)} = \max \left\{ 0, \frac{1}{\nu_* + \gamma \lambda_M^{(i)}} - \frac{1}{\alpha} \right\}, \quad (58)$$

where ν_* satisfies:

$$\sum_{i=1}^{\min\{d, N\}} \max \left\{ 0, \frac{1}{\nu_* + \gamma \lambda_M^{(i)}} - \frac{1}{\alpha} \right\} = N. \quad (59)$$

Given that $\gamma < (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$, we have $\frac{1}{\nu_* + \gamma \lambda_M^{(i)}} - \frac{1}{\alpha} > 0$ for all $i = 1, \dots, \min\{d, N\}$. Therefore, for the optimal \mathbf{Z}_* in problem (6), we have $\text{rank}(\mathbf{Z}_*) = \min\{d, N\}$ and the singular values $\sigma_{\mathbf{Z}_*}^{(i)} = \sqrt{\frac{1}{\gamma \lambda_M^{(i)} + \nu_*} - \frac{1}{\alpha}}$, for all $i = 1, \dots, \min\{d, N\}$.

Remark 3. We notice that (48) is a reverse water-filling problem, where the water level is controlled by $1/\alpha$, as shown in Figure A.1. When \mathbf{G} and \mathbf{M} have eigenspaces aligned and $\gamma < (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$, we have $\text{rank}(\mathbf{Z}_*) = \min\{d, N\}$ and $\lambda_{\mathbf{G}_*}^{(i)} \neq 0$ for all $i \leq \min\{d, N\}$. When $\gamma \geq (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$, non-zero $\lambda_{\mathbf{G}}^{(i)}$ first disappears for the larger $\lambda_M^{(i)}$. \square

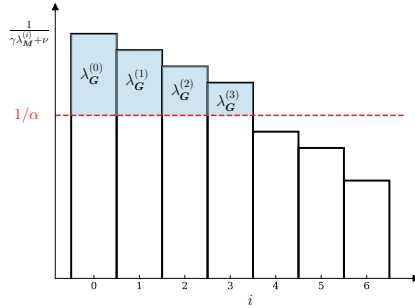


Figure A.1: **Illustration of the optimal solution for problem (6).** The primal problem can be transformed into a classical reverse water-filling problem.

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 $\mathbf{\Gamma}$, 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 be desired structure, i.e., $\mathbf{Z}_*^\top \mathbf{Z}_*$ and \mathbf{C}_* are block-diagonal, is that $\langle (\mathbf{I} - \mathbf{C})(\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \rangle \rightarrow 0$, where

$$\mathbf{G}^* := \text{Diag}(\mathbf{G}_{11}, \mathbf{G}_{22}, \dots, \mathbf{G}_{kk}) = \begin{bmatrix} \mathbf{G}_{11} & & & \\ & \ddots & & \\ & & \mathbf{G}_{kk} & \\ & & & \end{bmatrix},$$

and \mathbf{G}_{jj} is the block Gram matrix corresponding to \mathbf{Z}_j .

1188 *Proof.* We begin with the analysis to the first two terms of the loss function $\tilde{\mathcal{L}} := \mathcal{L}_1 + \gamma\mathcal{L}_2$, where

$$\begin{aligned} 1189 \mathcal{L}_1 &:= -\frac{1}{2} \log \det (\mathbf{I} + \alpha(\mathbf{Z}\mathbf{\Gamma})^\top (\mathbf{Z}\mathbf{\Gamma})) = -\frac{1}{2} \log \det (\mathbf{I} + \alpha\mathbf{G}), \\ 1190 \\ 1191 \\ 1192 \mathcal{L}_2 &:= \frac{1}{2} \|\mathbf{Z}\mathbf{\Gamma} - \mathbf{Z}\mathbf{\Gamma}\mathbf{\Gamma}^\top \mathbf{C}\mathbf{\Gamma}\|_F^2 = \frac{1}{2} \|\mathbf{Z} - \mathbf{Z}\mathbf{C}\|_F^2 = \frac{1}{2} \text{Tr} \left(\mathbf{G} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top \right), \\ 1193 \end{aligned}$$

1194 since that $\mathbf{\Gamma}^\top \mathbf{\Gamma} = \mathbf{\Gamma}\mathbf{\Gamma}^\top = \mathbf{I}$. Thus, we have:

$$1195 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) = \frac{\gamma}{2} \text{Tr} \left(\mathbf{G} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top \right) - \frac{1}{2} \log \det (\mathbf{I} + \alpha\mathbf{G}), \quad (60)$$

1196 which is a convex function with respect to (w.r.t) \mathbf{G} and \mathbf{C} , separately. By the property of convex

1197 function w.r.t. \mathbf{G} , we have:

$$\begin{aligned} 1201 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) &\geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}) + \left\langle \nabla_{\mathbf{G}} \tilde{\mathcal{L}}|_{[\mathbf{G}^*, \mathbf{C}]}, \mathbf{G} - \mathbf{G}^* \right\rangle \\ 1202 &= \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}) + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top - \frac{\alpha}{2} (\mathbf{I} + \alpha\mathbf{G}^*)^{-1}, \mathbf{G} - \mathbf{G}^* \right\rangle \\ 1203 &= \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}) + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle - \left\langle \frac{\alpha}{2} (\mathbf{I} + \alpha\mathbf{G}^*)^{-1}, \mathbf{G} - \mathbf{G}^* \right\rangle. \\ 1204 \\ 1205 \\ 1206 \end{aligned}$$

1207 Since that $\left\langle (\mathbf{I} + \alpha\mathbf{G}^*)^{-1}, \mathbf{G} - \mathbf{G}^* \right\rangle = 0$, we have:

$$1208 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) \geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}) + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle.$$

1209 By the property of convex function w.r.t. \mathbf{C} , we have:

$$\begin{aligned} 1210 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) &\geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}^*) + \left\langle \nabla_{\mathbf{C}} \tilde{\mathcal{L}}|_{[\mathbf{G}^*, \mathbf{C}^*]}, \mathbf{C} - \mathbf{C}^* \right\rangle + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle \\ 1211 &= \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}^*) + \left\langle -\gamma\mathbf{G}^* (\mathbf{I} - \mathbf{C}^*), \mathbf{C} - \mathbf{C}^* \right\rangle + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle. \\ 1212 \\ 1213 \end{aligned}$$

1214 Since that $\left\langle \mathbf{G}^* (\mathbf{I} - \mathbf{C}^*), \mathbf{C} - \mathbf{C}^* \right\rangle = 0$, we have:

$$1215 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) \geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}^*) + \left\langle \frac{\gamma}{2} (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle.$$

1216 It is easy to see that if $\left\langle (\mathbf{I} - \mathbf{C}) (\mathbf{I} - \mathbf{C})^\top, \mathbf{G} - \mathbf{G}^* \right\rangle \rightarrow 0$, then we will have:

$$1217 \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) \geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}^*), \quad (61)$$

1218 where the equality holds *only when* $\mathbf{G} = \mathbf{G}^*, \mathbf{C} = \mathbf{C}^*$, in which $\mathbf{C}^* =$

1219 $\text{Diag}(\mathbf{C}_{11}, \mathbf{C}_{22}, \dots, \mathbf{C}_{kk}) = \begin{bmatrix} \mathbf{C}_{11} & & \\ & \ddots & \\ & & \mathbf{C}_{kk} \end{bmatrix}$. Furthermore, if the regularizer $r(\cdot)$ satisfies the

1220 extended block diagonal condition as defined in (Lu et al., 2018), then we have that $r(\mathbf{C}) \geq r(\mathbf{C}^*)$,

1221 where the equality holds if and only if $\mathbf{C} = \mathbf{C}^*$. Therefore, we have:

$$1222 \mathcal{L}(\mathbf{G}, \mathbf{C}) = \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) + \beta \cdot r(\mathbf{C}) \geq \tilde{\mathcal{L}}(\mathbf{G}^*, \mathbf{C}^*) + \beta \cdot r(\mathbf{C}^*) = \mathcal{L}(\mathbf{G}^*, \mathbf{C}^*). \quad (62)$$

1223 Thus we conclude that minimizing the loss function $\mathcal{L}(\mathbf{G}, \mathbf{C}) = \tilde{\mathcal{L}}(\mathbf{G}, \mathbf{C}) + \beta \cdot r(\mathbf{C})$ promotes the

1224 optimal solution $(\mathbf{G}_*, \mathbf{C}_*)$ to have block diagonal structure.

1225 We note that the Gram matrix being block-diagonal, i.e., $\mathbf{G}_* = \mathbf{G}^*$, implies that $\mathbf{Z}_{*,j_1}^\top \mathbf{Z}_{*,j_2} = \mathbf{0}$ for

1226 all $1 \leq j_1 < j_2 \leq k$, which is corresponding to the subspaces associated to the blocks $\mathbf{Z}_{*,j}$'s are

1227 orthogonal to each other. \square

B EXPERIMENTAL SUPPLEMENTARY MATERIAL

B.1 EXPERIMENTAL DETAILS

B.1.1 SYNTHETIC DATA

As shown in Figure 5a (line 1), the data points are generated from two manifolds. The first manifold (colored in purple) is generated by sampling 100 data points from

$$\mathbf{x} = \begin{bmatrix} \cos\left(\frac{1}{5}\sin(5\varphi)\right)\cos\varphi \\ \cos\left(\frac{1}{5}\sin(5\varphi)\right)\sin\varphi \\ \sin\left(\frac{1}{5}\sin(5\varphi)\right) \end{bmatrix} + \epsilon, \quad (63)$$

where φ 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 $\mathcal{N}\left([0, 0, 1]^\top, 0.05\mathbf{I}_3\right)$. To further test more complicated cases, we remove 50 data points from the second manifold and replace them with 50 data points sampled from another Gaussian distribution $\mathcal{N}\left([0, 0, -1]^\top, 0.05\mathbf{I}_3\right)$ (Figure 5a line 2).

In PRO-DSC, the learnable mappings $\mathbf{h}(\cdot; \Psi)$ and $\mathbf{f}(\cdot; \Theta)$ 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 MLP is 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 $\alpha = 3/0.1 \cdot 200$.

For SEDSC methods, we use DSCNet (Ji et al., 2017) as the representative. In Figure 5b, we set $\gamma = 1$ for both cases, whereas in Figure 5c, γ is set to 5 and 100 for the two cases, respectively. The encoder and decoder of DSCNet are MLPs with two hidden layers, with the hidden dimensions set to 100 and 3. 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 and 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 pre-trained features for training, the input images are first resized to 224 with respect to the smaller edge, then center-cropped to 224×224 and fed into the CLIP pre-trained image encoder to obtain fixed features.⁹ The subsequent training of PRO-DSC takes the fixed extracted features as input, instead of loading the entire CLIP pre-trained model.

Network architecture and hyper-parameters. The learnable mappings $\mathbf{h}(\cdot; \Psi)$ and $\mathbf{f}(\cdot; \Theta)$ 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 SGD optimizer with the learning rate set to $\eta = 10^{-4}$, and the weight decay parameters of $\mathbf{f}(\cdot; \Theta)$ and $\mathbf{h}(\cdot; \Psi)$ are set to 10^{-4} and 5×10^{-3} , respectively. Following (Chu et al., 2024), we warm up training $\mathbf{f}(\cdot; \Theta)$ by diversifying the features with $\mathcal{L}_1 = -\log \det(\mathbf{I} +$

⁹We use the ViT L/14 pre-trained model provided by <https://github.com/openai/CLIP> for 768-dimensional features.

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.

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

$\alpha \mathbf{Z}_\Theta^\top \mathbf{Z}_\Theta$) for a few iterations and duplicate the weights to $\mathbf{h}(\cdot; \Psi)$. We set $\alpha = d/0.1 \cdot n_b$ 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**, where η is the learning rate, d_{pre} is the hidden and output dimension of pre-feature layer, m is the output dimension of \mathbf{h} and \mathbf{f} , n_b is the batch size for training, and “# warm-up” is the number of iterations of warm-up stage.

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

Running other algorithms. Since 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 train and test these models directly on the test set for all the experiments.

(1) For EnSC, we tune the hyper-parameter in front of the self-expressive term $\gamma \in \{1, 2, 5, 10, 20, 50, 100, 200, 400, 800, 1600, 3200\}$ and tune the hyper-parameter to balance the ℓ_1 and ℓ_2 norms $\tau \in \{0.9, 0.95, 1\}$ to report the best clustering result.

(2) For SSCOMP, we tune the hyper-parameter to control the sparsity $k_{max} \in \{1, 2, 5, 10, 20, 50, 100, 200\}$ and the residual $\epsilon \in \{10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}\}$ to report the best clustering result.

(3) To apply DSCNet to the CLIP features, we substitute MLPs with two hidden layers for the convolutional encoder and decoder. The hidden dimension of the MLPs are 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 test data set.

(4) 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 <https://github.com/LeslieTrue/CLIP/blob/main/main.py> and https://github.com/LeslieTrue/CLIP/blob/main/main_efficient.py), 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.

(5) 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.

(6) For SENet, SCAN and EDESC, we adjust the hyper-parameters and repeat experiments for three times, with only the best results are reported.

(7) For TEMI, we directly cited the results from the paper Adaloglou et al. (2023).

Training from scratch. Similar to most deep clustering algorithms, we divide the training process into two steps. We begin with pre-training the parameters of the backbone with BYOL (Grill et al., 2020). Then, leveraging the parameters pre-trained in the first stage, we 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. 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, following (Li et al., 2021; Huang et al., 2023). We use a convolution filter of size 3×3 and stride 1 to replace the first convolution filter, following (Huang et al., 2023; Li et al., 2020). The data augmentation strategy is as follows:

Augmentation 1 Augmentation for training from scratch

```
1: transforms.RandomResizedCrop(size=img_size, scale=(0.08, 1)),
2: transforms.RandomHorizontalFlip(),
3: transforms.RandomApply([transforms.ColorJitter(0.4, 0.4, 0.2,
0.1)], p=0.8),
4: transforms.RandomGrayscale(p=0.2),
5: 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.

B.2 EMPIRICAL VALIDATION ON THEORETICAL RESULTS

Empirical Validation on Theorem 1. To measure the eigenspace alignment of \mathbf{G} and \mathbf{M} (Theorem 1), we plot $\langle \mathbf{u}_j, \mathbf{G}\mathbf{u}_j / \|\mathbf{G}\mathbf{u}_j\|_2 \rangle$ in Figure 2c, where \mathbf{u}_j is the j^{th} eigenvector of \mathbf{M} . When the eigenspace alignment holds, one can verify that:

$$\langle \mathbf{u}_j, \frac{\mathbf{G}\mathbf{u}_j}{\|\mathbf{G}\mathbf{u}_j\|_2} \rangle = \begin{cases} 1, & \lambda_{\mathbf{G}}^{(j)} \neq 0 \\ 0, & \lambda_{\mathbf{G}}^{(j)} = 0 \end{cases} \quad \text{for all } j = 1, 2, \dots, N. \quad (64)$$

As shown in Figure 2c, the first $d = 128$ normalized correlation values converge to 1, while the rest converge to 0, implying the progressively alignment between \mathbf{G} and \mathbf{M} . In addition, we plot the Frobenius norm of the commutator \mathbf{L} during training in Figure 2d. The commutator decreases monotonically during the network training, implying the eigenspace alignment by Lemma A1.

Empirical Validation on Theorem 2. To verify Theorem 2, we decompose and plot the eigenvalues of \mathbf{G} , \mathbf{M} in Figure 2a and 2b, respectively. As shown, there are $\min\{d, N\} = 128$ non-zero eigenvalues of \mathbf{G} , approximately being inversely proportional to the smallest 128 eigenvalues of \mathbf{M} . This results empirically demonstrate that $\text{rank}(\mathbf{Z}_*) = \min\{d, N\}$ and $\lambda_{\mathbf{G}_*}^{(i)} = \frac{1}{\gamma \lambda_{\mathbf{M}}^{(i)} + \nu_*} - \frac{1}{\alpha}$ for minimizers. Furthermore, the condition of Theorem 2 is verified in Figure 1, where $\gamma < (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$ yields satisfactory clustering accuracy (ACC%) and subspace-preserving representation error (SRE%). The satisfactory ACC and SRE confirm that PRO-DSC avoids catastrophic collapse when $\gamma < (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$ holds. When $\gamma \geq (\alpha - \nu_*) / \lambda_{\max}(\mathbf{M})$, PRO-DSC yields significantly worse ACC and SRE. There is a phase transition phenomenon that corresponds to the sufficient condition to prevent collapse.

Empirical Validation on Theorem 3. To intuitively visualize the structured representations learned by PRO-DSC, we visualize the Gram matrices $|\mathbf{Z}^T \mathbf{Z}|$ for both CLIP features and learned representations on CIFAR-10. The Gram matrix shows the similarities between representations within the same class (indicated by block diagonal values) and across different classes (indicated by off-block diagonal values). Moreover, we display the dimensionality reduction results via Principal Component Analysis (PCA) for the CLIP features and the learned representation of samples from three

1404 categories in CIFAR-10. We use PCA for dimensionality reduction as it performs a linear projec-
 1405 tion, well preserving the underlying structure. As shown in Figure 3, the CLIP features from three
 1406 classes approximately lie on different subspaces. Despite of the structured nature of the features,
 1407 the underlying subspaces are not orthogonal. In the Gram matrix of the CLIP fature, the average
 1408 similarity between features from different classes is greater than 0.6, resulting in an unclear block
 1409 diagonal structure. After training with PRO-DSC, the spanned subspaces of the learned representa-
 1410 tions become orthogonal.¹⁰ Additionally, the off-block diagonal values of the Gram matrix decrease
 1411 significantly, revealing a clear block diagonal structure. These visualization results qualitatively
 1412 verify that PRO-DSC aligns the representations with a union of orthogonal subspaces.¹¹

1413 In Figure 4, we plot the curves for the compatibly structured coherence (CSC) condition and the
 1414 average values of $|\mathbf{G}^*|$, $|\mathbf{G} - \mathbf{G}^*|$, $|\mathbf{C}^*|$, $|\mathbf{C} - \mathbf{C}^*|$ during the training of PRO-DSC on CIFAR-10.
 1415 As illustrated, the CSC condition progressively satisfies. Consequently, the average off-block values
 1416 $|\mathbf{G} - \mathbf{G}^*|$ and $|\mathbf{C} - \mathbf{C}^*|$ gradually decrease while the average block values $|\mathbf{G}^*|$ and $|\mathbf{C}^*|$ gradually
 1417 increase, which empirically validates that PRO-DSC promotes block-diagonal structure in \mathbf{G} and
 1418 \mathbf{C} .

1419 B.3 MORE EXPERIMENTAL RESULTS

1420 **More results on synthetic data.** The synthetic experiments of adding an additional subspace are
 1421 presented in Figure B.1.

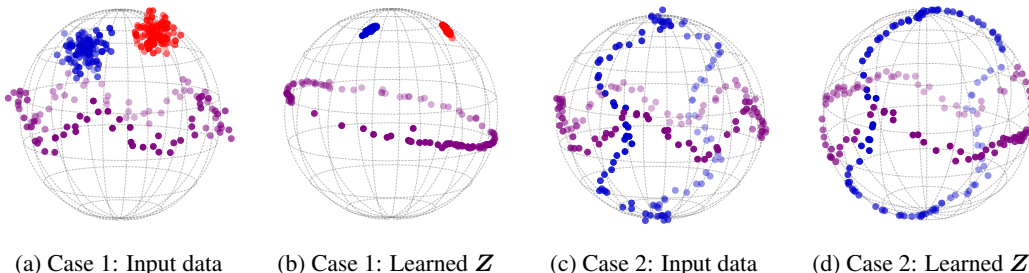
1422 In case 1, we implement two sets with 100 points in each cluster sampled from Gaussian distribution
 1423 $\mathbf{x} \sim \mathcal{N}([\frac{1}{\sqrt{2}}, 0, \sqrt{2}]^\top, 0.05\mathbf{I}_3)$ and $\mathbf{x} \sim \mathcal{N}([-\frac{1}{\sqrt{2}}, 0, \sqrt{2}]^\top, 0.05\mathbf{I}_3)$ in the same side of the sphere.
 1424 PRO-DSC eliminates the nonlinearity in representations and maximally separates the different sub-
 1425 spaces.

1426 In case 2, we add a vertical curve with 100 points sampled by:

$$1427 \mathbf{x} = \begin{bmatrix} \cos(\frac{1}{5} \sin(5\varphi)) \cos \varphi \\ \sin(\frac{1}{5} \cos(5\varphi)) \\ \cos(\frac{1}{5} \sin(5\varphi)) \sin \varphi \end{bmatrix} + \epsilon, \quad (65)$$

1428 where $\epsilon \sim \mathcal{N}(\mathbf{0}, 0.05\mathbf{I}_3)$ and use $\sin(\frac{1}{5} \cos(5\varphi))$ to avoid overlap in the intersection of the two
 1429 curves. PRO-DSC finds difficulties to learn representations of data which located at the intersection
 1430 of subspaces. However, those which are away from the intersection are linearized well.

1431 For the experiments on synthetic data, the learnable mappings $h(\cdot; \Psi)$ and $f(\cdot; \Theta)$ are implemented
 1432 with two MLPs with Rectified Linear Units (ReLU) (Nair & Hinton, 2010) as the activation function.
 1433 The hidden dimension and output dimension of the MLP is set to 100 and 3, respectively. In case
 1434 1, we train PRO-DSC with batch-size $n_b = 300$, learning rate $\eta = 5 \times 10^{-3}$ for 5000 epochs and
 1435 set $\gamma = 1.3, \beta = 500, \alpha = 3/0.1 \cdot 300$. In case 2, we train PRO-DSC with batch-size $n_b = 200$,
 1436 learning rate $\eta = 5 \times 10^{-3}$ for 8000 epochs and set $\gamma = 0.5, \beta = 500, \alpha = 3/0.1 \cdot 200$.



1443 (a) Case 1: Input data (b) Case 1: Learned \mathbf{Z} (c) Case 2: Input data (d) Case 2: Learned \mathbf{Z}

1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Figure B.1: Additional results on synthetic data.

¹⁰The dimension of each subspace is much greater than one (see Figure B.3). The 1-dimensional subspaces observed in the PCA results are a consequence of dimensionality reduction.

¹¹Please refer to Figure B.2 and B.7 for the results on other datasets and the visualization of the bases of each subspace.

Experimental Results 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 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.4. The hyper-parameters used for PRO-DSC is summarized in Table B.3.

Table B.3: Hyper-parameter setting for experiments on Reuters, UCI HAR, EYale-B, ORL and COIL-100.

Dataset	η	d_{pre}	d	#epochs	n_b	#warm-up	γ	β
REUTERS-10k	1e-4	1024	128	100	1024	50	50	200
UCI HAR	1e-4	1024	128	100	2048	20	100	300
EYale-B	1e-4	1080	256	10000	2432	100	200	50
ORL	1e-4	80	64	5000	400	100	75	10
COIL-100	1e-4	12800	100	10000	7200	100	200	100

Table B.4: Experimental Results on Datasets Reuters and UCI HAR with 10 trials. The results of other methods are cited from the respective papers.

Dataset	REUTERS-10k		UCI HAR	
	ACC	NMI	ACC	NMI
<i>k</i> -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	68.3	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	<u>80.7</u>
PRO-DSC	85.7 \pm 1.3	64.6 \pm 1.3	87.1 \pm 0.4	80.9 \pm 1.2

Comparison to AGCSC and SAGSC. 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.3. We repeated experiments for 10 trails and report the average with standard deviation in Table B.5. For the results of other methods in Table B.5, we directly cited the results from DSSC (Lim et al., 2020).

1) 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.

2) ARSSC. ARSSC employs three different non-convex regularizers: ℓ_γ 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 scalability of our PRO-DSC which is due to the re-parametrization (similar to SENet).

It is noteworthy that the goal of our work is to provide a reasonable (or a principled) framework for deep subspace clustering based on self-expressive (SE) model. We demonstrate theoretically and empirically that adding the $\log \det(\cdot)$ term into the SE model can prevent the catastrophic feature collapse, and show our PRO-DSC promotes to produce representations that are aligned with a union of orthogonal subspaces. Therefore, our PRO-DSC provides a general framework for self-expressive model based deep subspace clustering.

Since that both AGCSC and ARSSC improve the self-expressive model by incorporating GCN modules and self-adaptive attention mechanisms, respectively. Thus, it will be attempting to employ them into the deep self-expressive models under our PRO-DSC framework. Note that our PRO-DSC with a vanilla self-expressive model has already demonstrated highly competitive performance. We believe that extending PRO-DSC with current SOTA self-expressive models (such as AGCSC and ARSSC) would be more promising.

As a general framework for self-expressive model based deep subspace clustering, our PRO-DSC is reasonable, scalable and flexible to miscellaneous extensions.

Table B.5: Experiments on Extended Yale B, ORL and COIL-100. Note that the results of EnSC, SSCOMP, S3COMP, DSCNet, J-DSSC and ADSDSC are cited from (Lim et al., 2020).

	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
S3COMP	87.4	-	-	-	78.9	-
DSCNet	69.1	74.6	75.8	87.8	49.3	75.2
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	<u>82.4</u>	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	<u>96.0</u> \pm 0.3	95.7 \pm 0.8	<u>83.2</u> \pm 2.2	<u>92.7</u> \pm 0.6	82.8 \pm 0.9	95.0 \pm 0.6

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.2, 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.3. 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.4a, Figure B.4b and Figure B.4c, respectively. Recall that $\mathcal{L}_1 := -\frac{1}{2} \log \det(\mathbf{I} + \alpha \mathbf{Z}_\Theta^\top \mathbf{Z}_\Theta)$, $\mathcal{L}_2 := \frac{1}{2} \|\mathbf{Z}_\Theta - \mathbf{Z}_\Theta \mathbf{C}_\Psi\|_F^2$, and $\mathcal{L}_3 := \|\mathbf{A}_\Psi\|_{\overline{\mathcal{K}}}$. 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.

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.6. To apply DSCNet on CLIP features, we substitute MLPs with two hidden layers for 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, SEDSC overly compresses the representations and yields unsatisfactory clustering results.

Table B.6: **Clustering accuracy of CLIP features and learned representations.** We apply k -means, spectral clustering, and EnSC to cluster the representations.

	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	<u>74.7</u>	70.2	95.4	52.8	66.4	67.0	46.9	<u>49.2</u>	60.8	54.1	<u>62.8</u>	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	<u>70.3</u>	95.6	<u>75.3</u>	<u>75.0</u>	<u>77.5</u>	55.5	43.6	58.3	<u>62.1</u>	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

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).

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

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

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., 2023). 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.5).

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.8, 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.

Experiments on block diagonal regularizer with different k . To test the robustness of block diagonal regularizer $\|\mathbf{A}\|_{\square}$ to different k , we vary k and report the clustering performance in Table B.9. As illustrated, k does not necessarily equal to the number of clusters. There exists an interval within which the regularizer works effectively.

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

Table B.8: Experiments on out-of-domain datasets.

Methods	MNIST		F-MNIST		Flowers	
	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

number of clusters, over-segmentation will occur to the affinity matrix, which has negative impact on the subsequent clustering performance.

Table B.9: Clustering performance with different k in block diagonal regularizer.

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

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 subspace-preserving (correct) 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 two scenarios of synthetic data (Fig. 5), six benchmark datasets with CLIP features (Table 1), five benchmark datasets which are for training from scratch (Table 2), three out-of-domain datasets (Table B.8), using four different regularization terms (Table 4), using different feature extractor (Table B.7) and varying hyper-parameters (Fig. 7 and Table B.9). During the rebuttal, to reply Reviewer JbyQ, we also conduct experiments on two face image datasets (Extended Yale b and ORL), text dataset (REUTERS) and temporal dataset (UCI HAR). Currently, we did not find significant failure cases. However, as demonstrated in Fig. 2, our PRO-DSC will fail if the sufficient condition to prevent catastrophic collapse is not satisfied by using improper hyper-parameters γ and α .

Extensibility: As a general framework for self-expressive model based deep subspace clustering, our PRO-DSC is scalable and flexible to miscellaneous extensions. For example, AGCSC (Wei et al., 2023) and ARSSC (Wang et al., 2023a) improve the self-expressive model by incorporating GCN

1674 modules and self-adaptive attention mechanisms, respectively. Thus, it will be attempting to employ
1675 them into the deep self-expressive models under our PRO-DSC framework. Moreover, rather than
1676 using $\log \det(\cdot)$, there are other methods to solve the feature collapse issue, e.g., the nuclear norm.
1677 In addition, it is also worthwhile to incorporate the supervision information from the pseudo-label,
1678 e.g., (Huang et al., 2023; Jia et al., 2024; Li et al., 2017), for further improving the performance of
1679 our PRO-DSC.

1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

1741

1742

1743

1744

1745

1746

1747

1748

1749

1750

1751

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

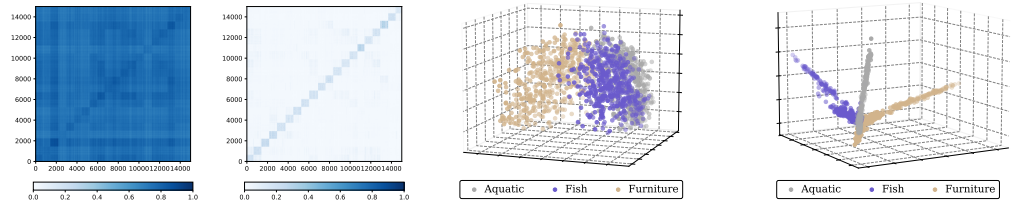
1777

1778

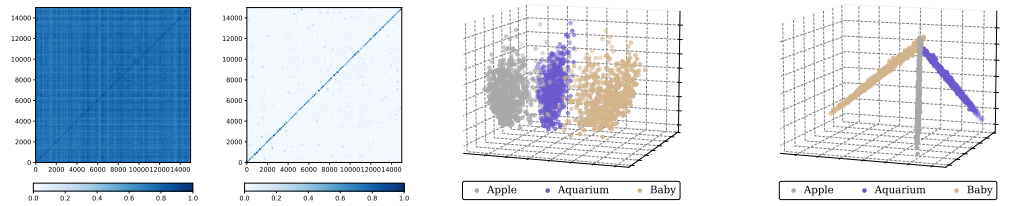
1779

1780

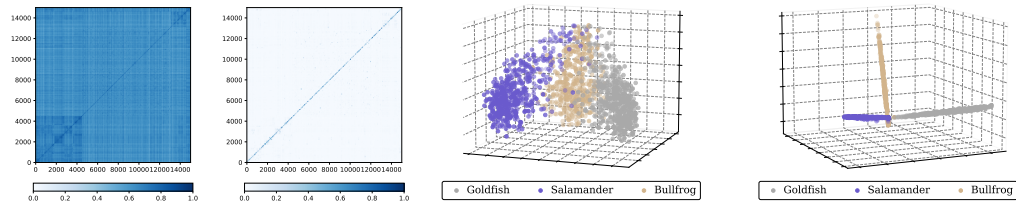
1781



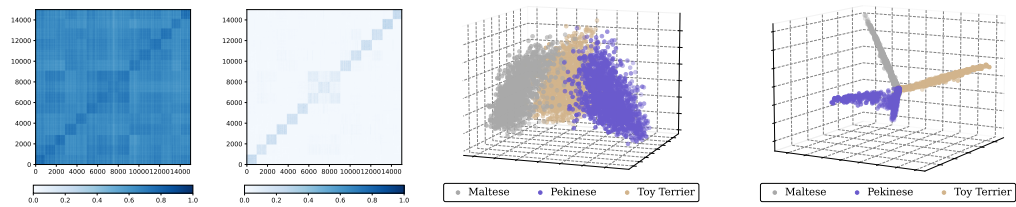
(a) CIFAR-20



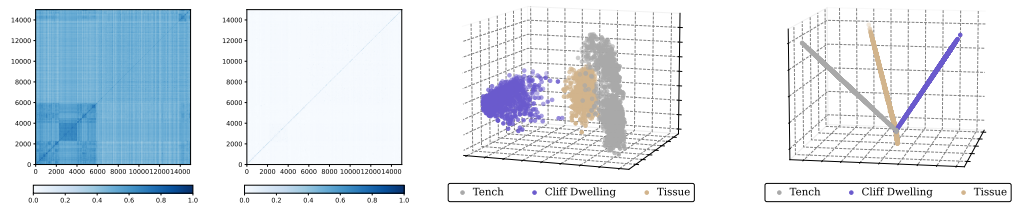
(b) CIFAR-100



(c) TinyImageNet-200



(d) ImageNet-Dogs-15



(e) ImageNet-1k

Figure B.2: Visualization of the union-of-orthogonal-subspaces structure of the learned representations via Gram matrix and PCA dimension reduction on three categories. Left: $|X^T X|$. Mid-left: $|Z^T Z|$. Mid-right: $X^{(3)}$ via PCA. Right: $Z^{(3)}$ via PCA.

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

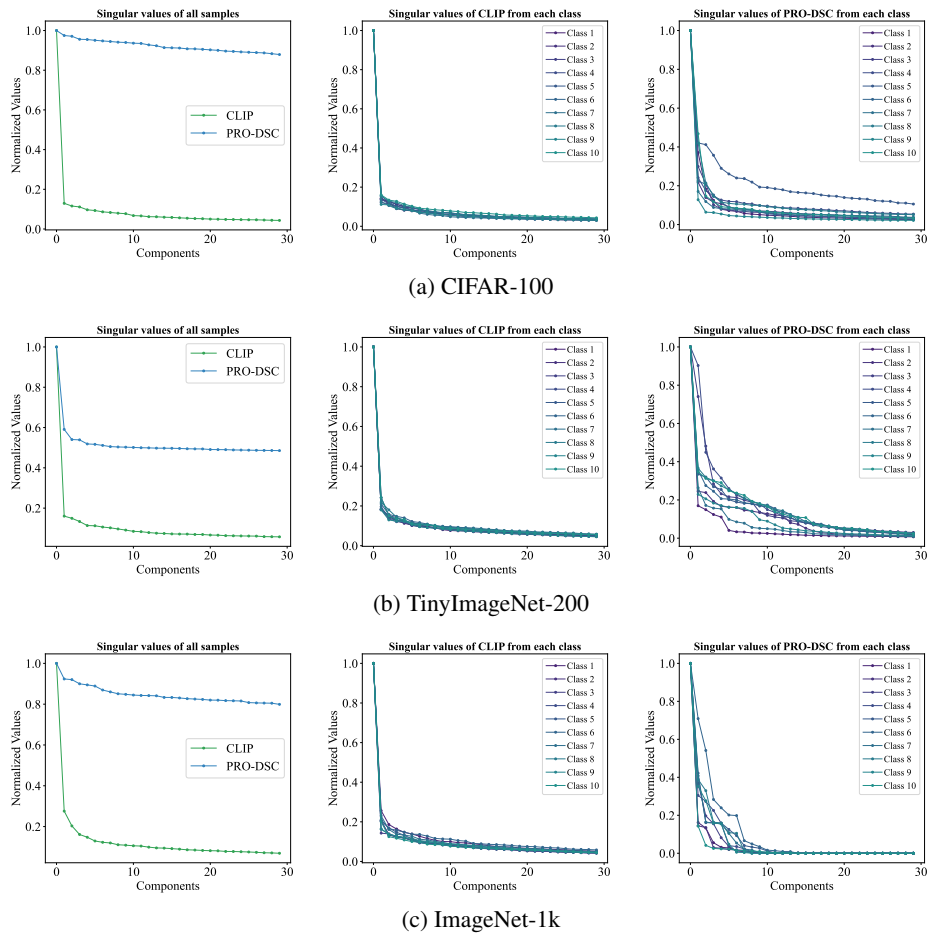


Figure B.3: **Singular values of features from all samples (left) and features from each class (mid and right).** For the better clarity, we plot the singular values for the first ten classes.

1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889

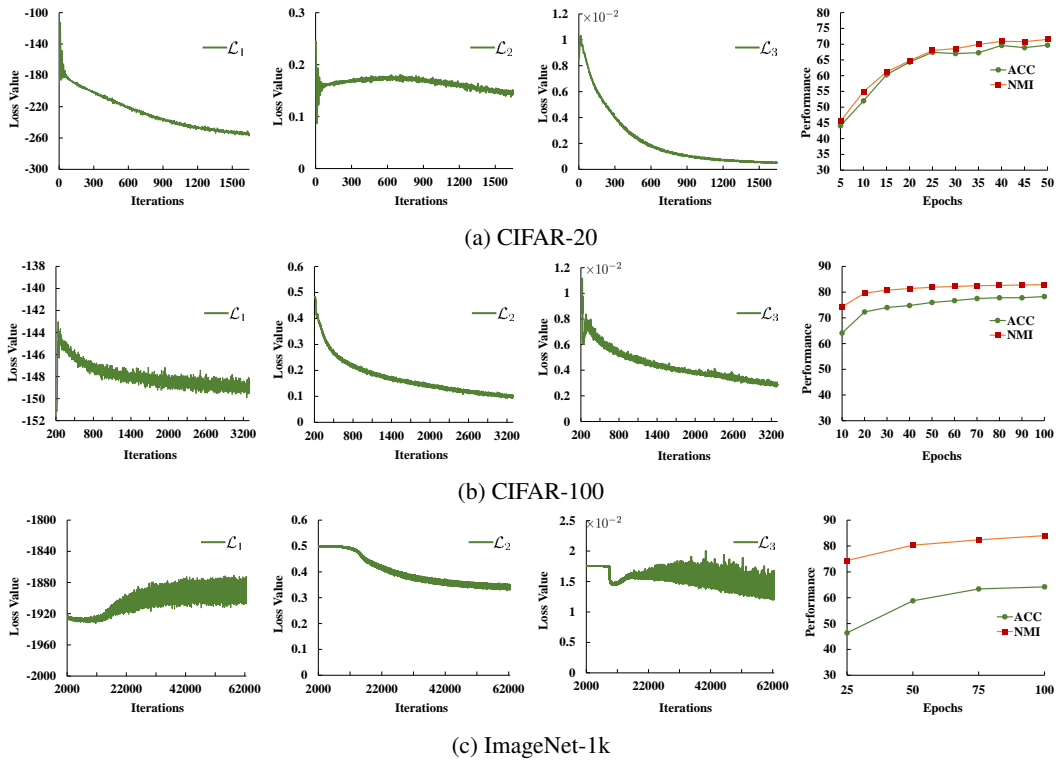


Figure B.4: The learning curves w.r.t. loss values and evaluation performance of PRO-DSC on CIFAR-20, CIFAR-100 and ImageNet-1k dataset.

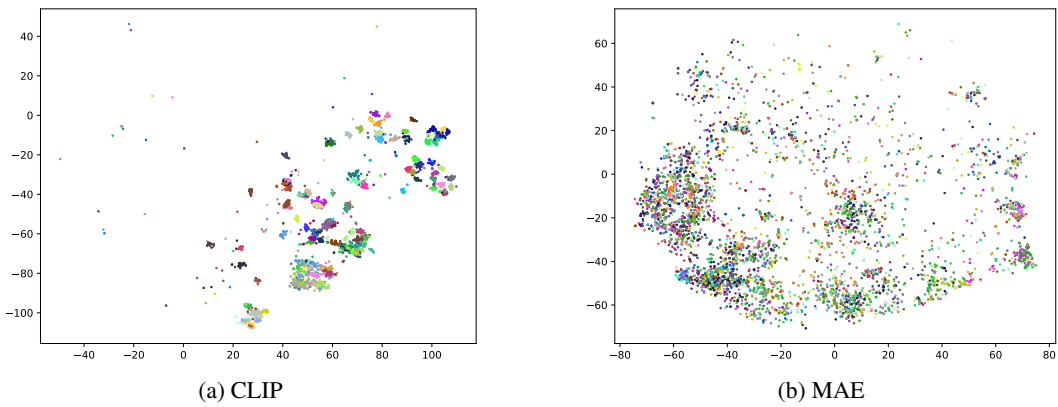


Figure B.5: The t -SNE visualization of CLIP and MAE features on the validation set of ImageNet-1k.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

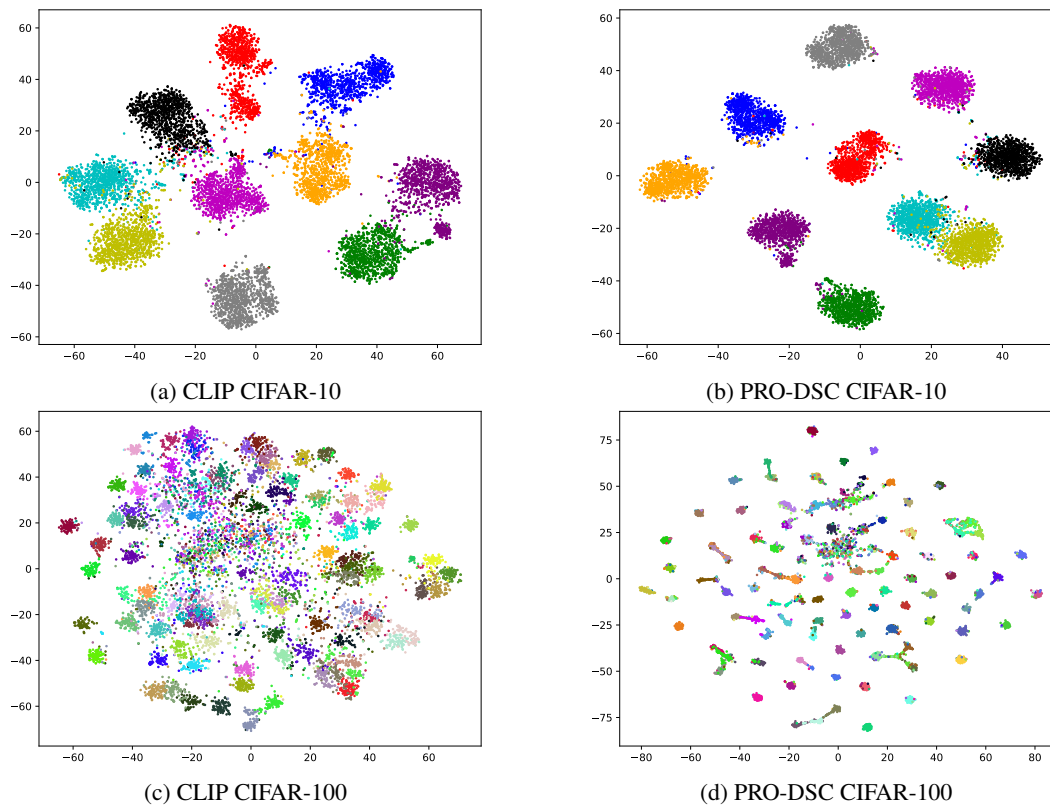
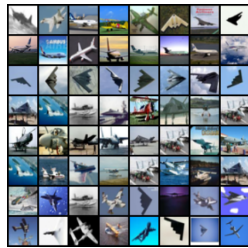
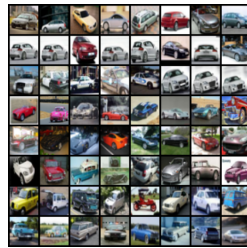


Figure B.6: *t*-SNE visualization of CLIP features and PRO-DSC's learned representations. The experiments are conducted on CIFAR-10 and CIFAR-100 dataset.

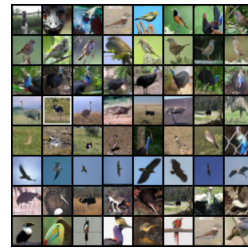
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997



(a) Cluster 1



(b) Cluster 2



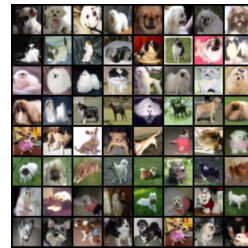
(c) Cluster 3



(d) Cluster 4



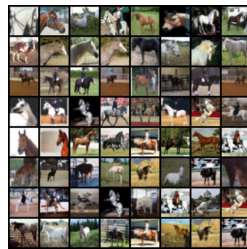
(e) Cluster 5



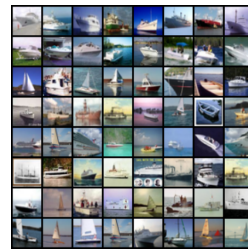
(f) Cluster 6



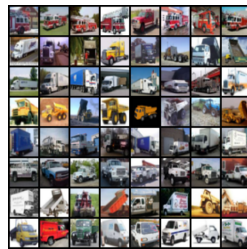
(g) Cluster 7



(h) Cluster 8



(i) Cluster 9



(j) Cluster 10

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.