

---

# A Simple Approach to Automated Spectral Clustering

## Appendices

---

Anonymous Author(s)

Affiliation

Address

email

### 1 A More discussion about LSR and KLSR

2 Note that if  $n \ll m$ , using the *push-through identity* [?], we reformulate (??) as  $\mathbf{C} = \mathbf{X}^\top(\lambda\mathbf{I} +$   
3  $\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}$  to reduce the computational cost from  $O(n^3)$  to  $O(mn^2)$ . In (??), when  $n$  is  
4 large (e.g. > 5000), we perform randomized SVD [?] on  $\mathbf{K}$ :  $\mathbf{K} \approx \mathbf{V}_r\boldsymbol{\Sigma}_r\mathbf{V}_r^\top$ . Then  $\mathbf{C} \approx$   
5  $\mathbf{V}_r\boldsymbol{\Sigma}_r^{1/2}(\lambda\mathbf{I} + \boldsymbol{\Sigma})^{-1}\boldsymbol{\Sigma}_r^{1/2}\mathbf{V}_r^\top$ , where  $r = 20k$  works well in practical applications. The time com-  
6 plexity of computing  $\mathbf{C}$  is  $O(r\tau n + rn^2)$ . The computation of the smallest  $k + 1$  eigenvalues of  $\mathbf{L}$   
7 is equivalent to compute the largest  $k + 1$  eigenvalues and eigenvectors of  $\mathbf{D}^{-1/2}\mathbf{A}\mathbf{D}^{-1/2}$ , which  
8 is sparse. The time complexity is  $O(k\tau n)$ . We have the follows.

9 **Proposition A.1.** *Let  $\hat{\mathbf{c}}$  be the optimal solution of minimize  $\frac{1}{2}\|\phi(\mathbf{y}) - \phi(\mathbf{X})\mathbf{c}\|^2 + \frac{\lambda}{2}\|\mathbf{c}\|^2$ ,*  
10 *where  $\phi$  is induced by Gaussian kernel and  $\mathbf{y}$  is arbitrary. Then  $\|\hat{\mathbf{c}}_i - \hat{\mathbf{c}}_j\| \leq$   
11  $\sqrt{2 - 2\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/(2\zeta^2))}$ .*

12 It shows that when two data points in  $\mathbf{X}$ , e.g.  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , are close to each other, the corresponding  
13 two elements in  $\hat{\mathbf{c}}$ , e.g.  $\hat{\mathbf{c}}_i$  and  $\hat{\mathbf{c}}_j$ , have small difference. Hence (??) with Gaussian kernel utilizes  
14 local information to enhance  $\mathbf{C}$ .

In LSR and KLSR, let  $\lambda \in \Lambda$ ,  $\tau \in \mathcal{T}$ , and  $\Theta = \Lambda \times \mathcal{T}$ . The algorithm of AutoSC-GD with only  
LSR and KLSR is shown in Algorithm 1. The total time complexity is

$$O(|\Lambda|(mn^2 + r\bar{\tau}n + rn^2) + 2|\Lambda||\mathcal{T}|k\bar{\tau}n),$$

15 where  $\bar{\tau}$  denotes the mean value in  $\mathcal{T}$ . The time complexity is at most  $O(|\Lambda|(mn^2 + |\mathcal{T}|kmn))$   
16 when  $\tau \leq r \leq m \leq n$ . It is worth noting that Algorithm 1 can be easily implemented parallelly,  
17 which will reduce the time complexity to  $O(\max(m, r)n^2 + kmn)$ . On the contrary, SSC, LRR,  
18 and their variants require iterative optimization and hence their time complexity is about  $O(tmn^2)$ ,  
19 where  $t$  denotes the iteration number and is often larger than 100.

### 20 B The algorithm of AutoSC+NSE

21 See Algorithm 2.

### 22 C More theoretical results

#### 23 C.1 Theoretical guarantee for KLSR

24 **Definition C.1** (Polynomial Deterministic Model). The columns of  $\mathbf{X}_0 \in \mathbb{R}^{m \times n}$  are drawn from a  
25 union of  $k$  different polynomials  $\{g_j : \mathbb{R}^r \rightarrow \mathbb{R}^m, r < m\}_{j=1}^k$  of order at most  $p$  and are further  
26 corrupted by noise, say  $\mathbf{X} = \mathbf{X}_0 + \mathbf{E}$ . Denote the eigenvalue decomposition of the kernel matrix  
27  $\mathbf{K}$  of  $\mathbf{X}$  as  $\mathbf{K} = \mathbf{V}\boldsymbol{\Sigma}\mathbf{V}^\top$ , where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$  and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ . Let  $\gamma = \sigma_{d+1}/\sigma_d$ .

---

**Algorithm 1** AutoSC-GD with Only LSR and KLSR

---

**Input:**  $\mathbf{X}$ ,  $k$ ,  $\mathcal{F}$ ,  $\Lambda$ ,  $\mathcal{T}$ 

- 1: Normalize the columns of  $\mathbf{X}$  to have unit  $\ell_2$  norm.
  - 2: **for**  $f_u$  in  $\mathcal{F}$  **do**
  - 3:   **for**  $\lambda_i$  in  $\Lambda$  **do**
  - 4:     Construct  $\mathbf{C}$  by (??) or (??).
  - 5:     **for**  $\tau_j$  in  $\mathcal{T}$  **do**
  - 6:        $\mathbf{C} \leftarrow |\mathbf{C} \odot (\mathbf{1} - \mathbf{I})|$ .
  - 7:       Truncate  $\mathbf{C}$  with parameter  $\tau_j$ .
  - 8:       For  $j = 1, \dots, n$ , let  $\mathbf{c}_j \leftarrow \mathbf{c}_j / |\mathbf{c}_j|_1$ .
  - 9:        $\mathbf{A} = (\mathbf{C} + \mathbf{C}^\top) / 2$ .
  - 10:        $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ .
  - 11:       Compute  $\sigma_1, \dots, \sigma_{k+1}$  and  $\mathbf{v}_1, \dots, \mathbf{v}_k$ .
  - 12:        $\Delta_{uij} = \text{REG}(\mathbf{L})$ ,  $\mathcal{V}_{uij} = [\mathbf{v}_1, \dots, \mathbf{v}_k]$ .
  - 13:     **end for**
  - 14:   **end for**
  - 15: **end for**
  - 16:  $\mathbf{Z} = \mathcal{V}_{\bar{u}\bar{i}\bar{j}}^\top$ , where  $\{\bar{u}, \bar{i}, \bar{j}\} = \text{argmax}_{u,i,j} \Delta_{uij}$ .
  - 17: Normalize the columns of  $\mathbf{Z}$  to have unit  $\ell_2$  norm.
  - 18: Perform k-means on  $\mathbf{Z}$ .
- Output:**  $k$  clusters:  $C_1, \dots, C_k$ .
- 

---

**Algorithm 2** AutoSC+NSE

---

**Input:**  $\mathbf{X}$ ,  $k$ ,  $\mathcal{F}$ ,  $\Theta$ ,  $\hat{n}$ .

- 1: Select  $\hat{n}$  landmarks from  $\mathbf{X}$  by k-means to form  $\hat{\mathbf{X}}$ .
  - 2: Apply AutoSC-G or AutoSC-BO to  $\hat{\mathbf{X}}$  with  $\mathcal{F}$  and  $\Theta$ .
  - 2: Get  $\hat{\mathbf{Z}}$  from the best Laplacian matrix given by AutoSC-G or AutoSC-BO.
  - 3: Use mini-batch Adam to solve (14).
  - 4: Compute  $\mathbf{Z}$  by (15).
  - 5: Perform k-means on  $\mathbf{Z}$ .
- Output:**  $k$  clusters:  $C_1, \dots, C_k$ .
- 

28 Denote  $\mathbf{v}_i = (v_{i1}, \dots, v_{in})$  the  $i$ -th row of  $\mathbf{V}$  and let  $\bar{\mathbf{v}}_i = (v_{i1}, \dots, v_{id})$ , where  $d < n$ . Suppose  
29 the following conditions hold: 1) for every  $i \in [n]$ , the  $\bar{\tau}$ -th largest element of  $\{|\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j| : j \in C_{\pi(i)}\}$   
30 is greater than  $\alpha$ ; 2)  $\max_{i \in [n]} \max_{j \in [n] \setminus C_{\pi(i)}} |\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j| \leq \beta$ ; 3)  $\max_{i,j,l} |v_{il} v_{jl}| \leq \mu$ .

31 Here we consider polynomials because they are easy to analyze and can well approximate smooth  
32 functions provided that  $p$  is sufficiently large. Clustering the columns of  $\mathbf{X}$  given by Definition  
33 C.1 according to the polynomials is actually a manifold clustering problem beyond the setting of  
34 subspace clustering. Similar to the subspace detection property, we define

35 **Definition C.2** (Manifold Detection Property). A symmetric affinity matrix  $\mathbf{A}$  obtained from  $\mathbf{X}$  has  
36 manifold detection property if for all  $i$ , the nonzero elements of  $\mathbf{a}_i$  correspond to the columns of  $\mathbf{X}$   
37 lying on the same manifold as  $\mathbf{x}_i$ .

38 The following theorem verifies the effectiveness of (12) followed by the truncation operation in  
39 manifold detection.

40 **Theorem C.3.** Suppose  $\mathbf{X}$  and  $\mathbf{K}$  are given by Definition C.1 and  $\mathbf{C}$  is given by (12), where the  
41 kernel function is a polynomial kernel of order  $q$ ,  $\text{rank}(\mathbf{K}_0) = d$  ( $\mathbf{K}_0$  is from  $\mathbf{X}_0$ ), and

$$\frac{(\rho - \sqrt{\rho^2 - 4(2\mu d - \Delta)(2\mu n - 2\mu d - \Delta)}) \sigma_d^2}{4\mu d - 2\Delta} < \lambda < \frac{(\rho + \sqrt{\rho^2 - 4(2\mu d - \Delta)(2\mu n - 2\mu d - \Delta)}) \sigma_d^2}{4\mu d - 2\Delta} \quad (1)$$

42 where  $\rho = 2\mu n \gamma^2 - \Delta(1 + \gamma^2)$ . Then  $d \leq k \binom{r+pq}{pq}$  and the  $\mathbf{C}$  truncated by  $\tau \leq \bar{\tau}$  has the manifold  
43 detection property.

44 In the theorem,  $\sigma_d$  can be much larger than  $\sigma_{d+1}$  provided that the noise is small enough. Then  
 45 we get a wide range for  $\lambda$ . Compared to Theorem 3.7, Theorem C.3 allows a much larger  $d$ , which  
 46 means the kernel method is able to handle more difficult clustering problems than the linear method.

## 47 C.2 Theoretical analysis for NSE

48 The following proposition shows that a small number of hidden nodes in NSE are sufficient to make  
 49 the clustering succeed.

50 **Proposition C.4.** *Suppose the columns (with unit  $\ell_2$  norm) of  $\mathbf{X}$  are drawn from a union of  $k$   
 51 independent subspaces of dimension  $r$ :  $\sum_{j=1}^k \dim(\mathcal{S}_j) = \dim(\mathcal{S}_1 \cup \dots \cup \mathcal{S}_k) = kr$ . For  $j =$   
 52  $1, \dots, k$ , let  $\mathbf{U}^j$  be the bases of  $\mathcal{S}_j$  and  $\mathbf{x}_i = \mathbf{U}^j \mathbf{v}_i$ , if  $\mathbf{x}_i \in \mathcal{S}_j$ . Suppose  $\max\{\|\mathbf{U}_{:,l}^{i\top} \mathbf{U}^j\| : 1 \leq l \leq$   
 53  $r, 1 \leq i \neq j \leq k\} \leq \mu$ . Suppose that for all  $i = 1, \dots, n$ ,  $\max\{v_{1i}, \dots, v_{ri}\} > \mu$ . Then there exist  
 54  $\mathbf{W}_1 \in \mathbb{R}^{d \times m}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{k \times d}$ ,  $\mathbf{b}_1 \in \mathbb{R}^d$ , and  $\mathbf{b}_2 \in \mathbb{R}^k$  such that performing  $k$ -means on  $\mathbf{Z}$  given by  
 55 (15) identifies the clusters correctly, where  $d = kr$ .*

## 56 D More details about AutoSC-BO

57 In AutoSC-BO, we use Expected Improvement (EI) acquisition function

$$a_{\text{EI}}(\mathbf{s}|\mathcal{D}_t) = \mathbb{E}_p [\max(g_{\min} - g(\mathbf{s}), 0)], \quad (2)$$

58 where  $g_{\min}$  is the best function value known. The closed-form formulation is

$$a_{\text{EI}}(\mathbf{s}|\mathcal{D}_t) = (g_{\min} - \mu) \Phi\left(\frac{g_{\min} - \mu}{\sigma}\right) + \phi\left(\frac{g_{\min} - \mu}{\sigma}\right), \quad (3)$$

59 where  $\mu = \mu(\mathbf{s}|\mathcal{D}_t, \theta_K)$  and  $\sigma = \sigma(\mathbf{s}|\mathcal{D}_t, \theta_K)$  are the mean value and variance of the Gaussian  
 60 process,  $\phi$  and  $\Phi$  are standard Gaussian cumulative density function and probability density func-  
 61 tion respectively, and  $\theta_K$  denotes the hyperparameters of the Gaussian process. For the covariance  
 62 function, we use the automatic relevance determination (ARD) Matérn 5/2 kernel Matérn [2013]:

$$k_{M52}(\mathbf{s}, \mathbf{s}') = \theta_0 \left(1 + \sqrt{5r^2(\mathbf{s}, \mathbf{s}') + \frac{5}{3}r^2(\mathbf{s}, \mathbf{s}')}\right) \times \exp\left(-\sqrt{5r^2(\mathbf{s}, \mathbf{s}')}\right), \quad (4)$$

63 where  $r^2(\mathbf{s}, \mathbf{s}') = \sum_{j=1}^d (s_j - s'_j)^2 / \theta_j^2$ . AutoSC-BO is implemented in MATLAB.

## 64 E More about the experiments

### 65 E.1 Dataset description

66 The description for the benchmark image datasets considered in this paper are as follows.

- 67 • **Extended Yale B Face** [Kuang-Chih *et al.*, 2005] (Yale B for short): face images  
 68 (192×168) of 38 subjects. Each subject has about 64 images under various illumination  
 69 conditions. We resize the images into 32 × 32.
- 70 • **ORL Face** [Samaria and Harter, 1994]: face images (112×92) of 40 subjects. Each sub-  
 71 ject has 10 images with different poses and facial expressions. We resize the images into  
 72 32×32.
- 73 • **COIL20** [Nene *et al.*, 1996]: images (32 × 32) of 20 objects. Each object has 72 images  
 74 of different poses.
- 75 • **AR Face** [Martínez and Kak, 2001]: face images (165×120) of 50 males and 50 females.  
 76 Each subject has 26 images with different facial expressions, illumination conditions, and  
 77 occlusions. We resize the images into 42 × 30.
- 78 • **MNIST** [LeCun *et al.*, 1998]: 70,000 grey images (28 × 28) of handwritten digits 0 – 9.
- 79 • **MNIST-1k(10k)**: a subset of MNIST containing 1000(10000) samples, 100(1000) ran-  
 80 domly selected samples per class.

- 81 • **Fashion-MNIST** [Xiao *et al.*, 2017]: 70,000 gray images ( $28 \times 28$ ) of 10 types of fashion  
82 product.
- 83 • **Fashion-MNIST-1k(10k)**: a subset of Fashion-MNIST containing 1000(10000) samples,  
84 100(1000) randomly selected samples per class.
- 85 • **MNIST-feature**: following the same procedures of [Chen *et al.*, 2020], we compute a fea-  
86 ture vector of dimension 3,472 using the scattering convolution network Bruna and Mallat  
87 [2013] and then reduce the dimension to 500 using PCA.
- 88 • **Fashion-MNIST-feature**: similar to MNIST-feature.
- 89 • **GTSRB** [Stallkamp *et al.*, 2012]: consisting of 12,390 images of street signs in 14 cate-  
90 gories. Following [Chen *et al.*, 2020], we extract a 1568-dimensional HOG feature, and  
91 reduce the dimension to 500 by PCA.

92 All experiments are conducted in MATLAB on a MacBook Pro with 2.3 GHz Intel i5 Core and 8GB  
93 RAM.

## 94 E.2 Hyperparameter settings for the small datasets and the results on COIL20 and 95 Fashion-MNIST 1k

96 The parameter  $\lambda$  in each of SSC, LRR, and KSSC is chosen from  
97  $\{0.01, 0.02, 0.05, 0.1, 0.2, \dots, 0.5\}$ . The  $\lambda$  in BDR is chosen from  $\{5, 10, 20, \dots, 80\}$ . The  
98  $\gamma$  in BDR-B and BDR-Z is chosen from  $\{0.01, 0.1, 1\}$ . The parameter  $s$  in SSC-OMP is chosen  
99 from  $\{3, 4, \dots, 15\}$ . We report the results of these methods with their best hyperparameters. In  
100 AutoSC, we set  $\Lambda = \{0.01, 0.1, 1\}$  and  $\mathcal{T} = \{5, 6, \dots, 15\}$ . In AutoSC-BO, we consider two  
101 models: 1) Gaussian kernel similarity; 2) KLSR with polynomial kernel; 3) KLSR with Gaussian  
102 kernel, in which the hyperparameters of kernels are optimized adaptively. Then we needn't to  
103 consider LSR explicitly because it is a special case of KLSR with polynomial kernel. See Appendix  
104 E.5.

105 The clustering results on COIL20 and Fashion-MNIST-1k are shown in Table 1.

Table 1: Clustering results on COIL20 and Fashion-MNIST-1k

		SSC	LRR	EDSC	KSSC	SSC-OMP	BDR-Z	BDR-B	AutoSC-GD	AutoSC-BO
COIL20	acc	0.871	0.729	0.759	<b>0.912</b>	0.658(0.030)	0.713	0.791	0.782(0.012)	<b>0.878</b>
	time	61.8	221.2	15.4	100.6	<b>2.5</b>	86.8	86.8	<b>7.6</b>	39.2
Fashion-MNIST-1k	acc	0.553	0.515	0.544	0.548	0.566	0.574	0.563	<b>0.581</b>	<b>0.584</b>
	time	(0.025)	(0.014)	(0.017)	(0.016)	(0.034)	(0.019)	(0.031)	(0.025)	(0.021)
		24.1	68.5	5.1	35.9	<b>1.2</b>	25.7	25.7	<b>2.6</b>	22.7

## 106 E.3 Clustering results in terms of NMI

107 In addition to the clustering accuracy reported in Table 4, here we also compare the normalized  
108 mutual information (NMI) in Table 2. We see that the comparative performance of all methods are  
109 similar to the results in Table 4 and our methods AutoSC-GD and AutoSC-BO outperformed other  
methods in almost all cases.

Table 2: Normalized Mutual Information on the six small datasets

	SSC	LRR	EDSC	KSSC	SSC-OMP	BDR-Z	BDR-B	AutoSC-GD	AutoSC-BO
Yale B	0.817	0.703	0.835	0.730	0.841	0.666	0.743	<b>0.919</b>	<b>0.928</b>
ORL	0.849	0.872	0.856	0.872	0.815	0.875	0.865	<b>0.907</b>	<b>0.903</b>
COIL20	0.954	0.706	0.843	<b>0.983</b>	0.671	0.843	0.873	0.897	<b>0.963</b>
AR	0.818	0.872	0.825	0.809	0.691	0.865	0.861	<b>0.887</b>	<b>0.904</b>
MNIST-1k	0.612	0.538	0.631	0.626	0.546	0.634	0.580	<b>0.667</b>	<b>0.652</b>
Fashion-MNIST-1k	0.616	0.601	0.621	0.621	0.559	0.614	0.605	<b>0.633</b>	<b>0.629</b>

111 **E.4 The stability of AutoSC**

112 Though we have used a relatively compact search space in AutoSC to reduce the highly unnecessary  
 113 computational cost, the search space can be arbitrarily large. Figure 1 shows the clustering accuracy  
 114 and the corresponding relative-eigen-gap. We can see that the region with highest relative-eigen-gap  
 is in accordance with the region with highest clustering accuracy.

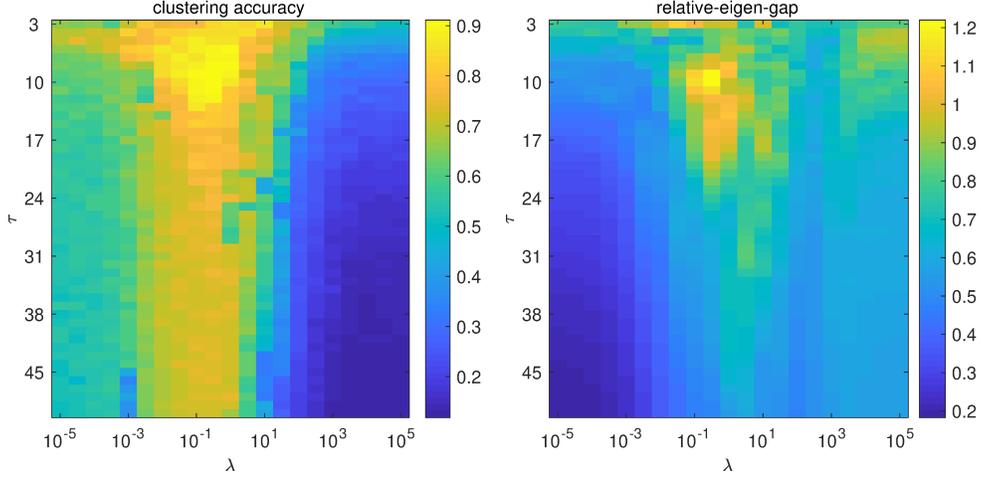


Figure 1: Visualization of the clustering accuracy and the corresponding relative-eigen-gap when a large search space is used.

115

116 **E.5 More about AutoSC-BO in the experiments**

117 For SSC, we consider the following problem

$$\text{minimize}_C \frac{1}{2} \text{Tr}(\mathbf{K} - 2\mathbf{K}\mathbf{C} + \mathbf{C}^\top \mathbf{K}\mathbf{C}) + \lambda \|\mathbf{C}\|_1, \quad (5)$$

118 where  $\mathbf{K}$  is an  $n \times n$  kernel matrix with  $[\mathbf{K}]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . Note that when we use a linear kernel  
 119 function, (5) reduces to the vanilla SSC. We solve the optimization via alternating direction method  
 120 of multipliers (ADMM) Boyd *et al.* [2011], where the Lagrange parameter is 0.1 and max number of  
 121 iterations is 500. In this study, we consider polynomial kernel and Gaussian kernel, and optimize all  
 122 hyperparameters including the order of the polynomial kernel. Particularly, for Gaussian kernel, we  
 123 set  $\varsigma = \frac{\xi}{n^2} \sum_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|$  and optimize  $\xi$ . The search space for the hyperparameters are as follows:  
 124  $10^{-3} \leq \lambda \leq 1, 5 \leq \tau \leq 50, 0 \leq b \leq 10^3, 1 \leq q \leq 5, 0.5 \leq \xi \leq 5$ .

125 In addition to Figure 2 of the main paper, here we report the best hyperparameters of the four models  
 126 found by AutoSC-BO in Table 3. It can be found that the accuracy of KLSR with a linear kernel is  
 127 higher than other models, which is consistent with its highest reg.

Table 3: The best hyperparameters and the corresponding clustering accuracy given by AutoSC<sub>BO</sub> on the first 10 subjects of YaleB Face dataset.

method	hyperparameters	reg	accuracy
KLSR (Polynomial)	$\lambda = 0.207, b = 19.09,$ $q = 1, \tau = 5$	2.379	0.966
KLSR (Gaussian)	$\lambda = 0.013,$ $\xi = 4.92, \tau = 5$	2.217	0.963
KSSC (Polynomial)	$\lambda = 0.519, b = 44.57,$ $q = 2, \tau = 5$	1.388	0.859
KSSC (Gaussian)	$\lambda = 0.0011,$ $\xi = 4.97, \tau = 6$	0.892	0.584

128 **E.6 Hyperparameter settings of large-scale clustering**

129 On MNIST-10k, MNIST, Fashion-MNIST-10k, and Fashion-MNIST, the parameter settings of Chen  
 130 and Cai [2011], SSSC Peng *et al.* [2013], SSC-OMP You *et al.* [2016], and S<sup>5</sup>C Matsushima and  
 131 Brbic [2019], and S<sup>3</sup>COMP-C Chen *et al.* [2020], and AutoSC+NSE are shown in Table 4. These  
 132 hyper parameters have been determined via grid search and the best (as possible) values are used.

Table 4: Hyper-parameter settings of the compared methods on MNIST-10k, MNIST, Fashion-MNIST-10k, and Fashion-MNIST.  $s$  denotes the number of landmark data points. In the optimization (mini-batch Adam) of AutoSC+NSE, the epoch number, batch size, and step size are 200, 128, and  $10^{-3}$  respectively.

LSC-K	$s = 1000, r = 3$
SSSC	$s = 1000, \lambda = 0.01$
SSC-OMP	$K = 10$ (sparsity)
S <sup>5</sup> C	$s = 1000, \lambda = 0.1$ or $0.2$
S <sup>3</sup> COMP-C	$T = 20, \lambda = 0.4, \delta = 0.9$
AutoSC+NSE	$s = 1000, d = 200, \gamma = 10^{-5}$
AutoSC <sub>BO</sub> +NSE	$s = 1000, d = 200, \gamma = 10^{-5}$

133 **E.7 Influence of hyper-parameters in AutoSC+NSE**

134 We investigate the effects of the type of activation function and the number ( $d$ ) of nodes in the  
 135 hidden layer of NSE. For convenience, we used a fixed random seed of MATLAB (rng(1)). Figure 2  
 136 shows the clustering accuracy on MNIST given by AutoSC+NSE with different activation function  
 137 and different  $d$ . We see that ReLU outperformed tanh consistently. The reason is that the nonlinear  
 138 mapping  $g$  from the data space to the eigenspace of the Laplacian matrix is nonsmooth and ReLU  
 139 is more effective than tanh in approximating nonsmooth functions. In addition, when  $d$  increases,  
 140 the clustering accuracy of AutoSC+NSE with ReLU often becomes higher because wider network  
 141 often has higher ability of function approximation.

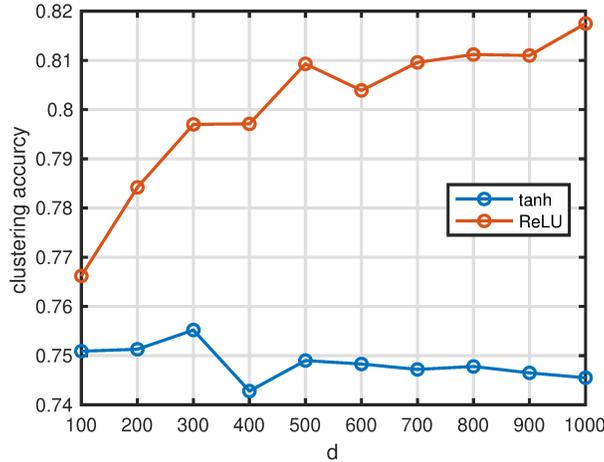


Figure 2: ReLU v.s. tanh (hyperbolic tangent) in the hidden layer of AutoSC+NSE on MNIST. When using ReLU, we set  $\gamma = 10^{-5}$  and  $\alpha = 10^{-3}$  (the step size in Adam). When using tanh, we set  $\gamma = 10^{-3}$  and  $\alpha = 10^{-2}$ , which perform best in this case. Notice that the clustering accuracy when using ReLU is higher than 0.78 in almost all cases, which is higher than the value (say 0.755) we reported in the main paper. The reason is than in the main paper, we reported the mean value of 10 repeated trials but here we reports the value of a single trial.

142 Figure 3 shows the clustering accuracy on MNIST given by AutoSC+NSE with different  $\gamma$  and  
 143  $\alpha$ . When  $\alpha$  is too small (say  $10^{-4}$ , the clustering accuracy is low, because the training error is  
 144 quite large in 200 epochs. In fact, by increasing the training epochs, the clustering accuracy can  
 145 be improved, which however will increase the time cost. When  $\alpha$  is relatively large, the clustering

146 accuracy is often higher than 0.755. On the other hand, AutoSC+NSE is not sensitive to  $\gamma$  provided  
 147 that it is not too large.

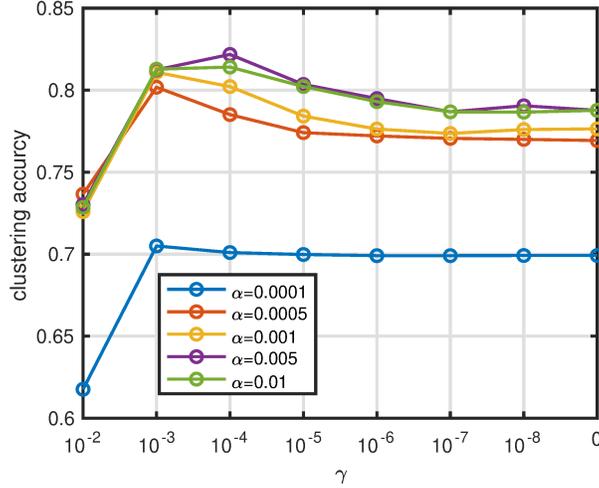


Figure 3: Influence of  $\gamma$  and  $\alpha$  in AutoSC+NSE on MNIST. We set  $d = 200$  and use ReLU.

148 Figure 4 shows the mean value and standard deviation (10 repeated trials) of the clustering accuracy  
 149 on MNIST given by AutoSC+NSE with different number (denoted by  $s$ ) of landmark points. It  
 150 can be found that when the  $s$  increases, the clustering accuracy increases and its standard deviation  
 becomes smaller. When  $s$  is large enough, the improvement is not significant.

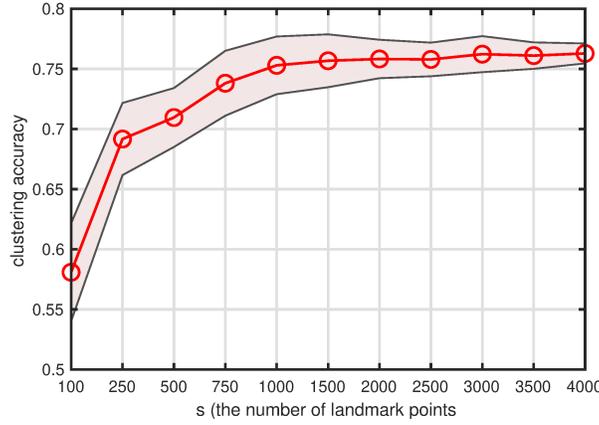


Figure 4: Influence of the number of landmark points in AutoSC+NSE on MNIST. We set  $d = 200$ ,  $\gamma = 10^{-5}$ , and  $\alpha = 10^{-3}$ . The shadow denotes the standard deviation of 10 trials.

151

## 152 **F Proof for the theoretical results**

### 153 **F.1 Proof for Claim 3.2**

154 *Proof.* The stochastic transition matrix of  $G$  is defined as

$$154 \quad \mathbf{P} = \mathbf{D}^{-1} \mathbf{A}. \tag{6}$$

155 In Meila [2001], it was showed that

$$155 \quad \text{MNCut}(\mathcal{C}) \geq k - \sum_{i=1}^k \varrho_i(\mathbf{P}), \tag{7}$$

156 where  $\varrho_i(\mathbf{P})$  denotes the  $i$ -th largest eigenvalue of  $\mathbf{P}$  and  $1 = \varrho_1(\mathbf{P}) \geq \varrho_2(\mathbf{P}) \geq \dots \varrho_k(\mathbf{P})$ .  
 157 According to Lemma 3 of Meila [2001], we have

$$\sigma_i(\mathbf{L}) = 1 - \varrho_i(\mathbf{P}), \quad \forall i = 1, \dots, n. \quad (8)$$

158 Substituting (8) into (7), we have

$$\text{MNCut}(\mathcal{C}) \geq \sum_{i=1}^k \sigma_i(\mathbf{L}). \quad (9)$$

159

□

160 *Remark F.1.*  $\mathcal{C}$  can be any partition of the nodes of  $G$ . Let  $\mathcal{C}^*$  be the optimal partition. Then  
 161  $\text{MNCut}(\mathcal{C}^*) = \sum_{i=1}^k \sigma_i(\mathbf{L})$ . If  $\sum_{i=1}^k \sigma_i(\mathbf{L}) = 0$ , there are no connections (edges) among  
 162  $C_1^*, \dots, C_k^*$ .

### 163 F.2 Proof for Claim 3.3

164 *Proof.* For  $i = 1, \dots, k$ , we aim to partition  $C_i$  into two subsets, denoted by  $C_i^1$  and  $C_i^2$ . Then we  
 165 define

$$\text{MNCut}(C_i) = \frac{\text{Cut}(C_i^1, C_i^2)}{\text{Vol}(C_i^1)} + \frac{\text{Cut}(C_i^2, C_i^1)}{\text{Vol}(C_i^2)}. \quad (10)$$

166 It follows that

$$\text{MNCut}(C_i) \geq \sum_{j=1}^2 \sigma_j(\mathbf{L}_{C_i}) \geq \sigma_2(\mathbf{L}_{C_i}) = ac(C_i), \quad (11)$$

167 where  $\mathbf{L}_{C_i}$  denotes the Laplacian matrix of  $C_i$  an  $i = 1, \dots, k$ . Since  $\sigma_{k+1}(\mathbf{L}) =$   
 168  $\min\{ac(C_1), \dots, ac(C_k)\}$ , we have

$$\min_{1 \leq i \leq k} \text{MNCut}(C_i) \geq \sigma_{k+1}(\mathbf{L}). \quad (12)$$

169 Therefore,  $\sigma_{k+1}(\mathbf{L})$  measures the least connectivity of  $C_1, \dots, C_k$ . This finished the proof. □

170 *Remark F.2.* When  $\sigma_{k+1}(\mathbf{L})$  is large, the connectivity in each of  $C_1, \dots, C_k$  is strong. Otherwise,  
 171 the connectivity in each of  $C_1, \dots, C_k$  is weak. When  $\sigma_{k+1}(\mathbf{L}) = 0$ , at least one of  $C_1, \dots, C_k$   
 172 contains at least two components, which means the nodes of  $G$  can be partitioned into  $k + 1$  or more  
 173 clusters.

### 174 F.3 Proof for Theorem 3.4

175 *Proof.* According to Theorem 1 of Meila *et al.* [2005], we have

$$\text{dist}(\mathcal{C}, \mathcal{C}') < \frac{3\delta}{\sigma_{k+1}(\mathbf{L}) - \sigma_k(\mathbf{L})}. \quad (13)$$

176 Since  $\text{reg}(\mathbf{L}) = \frac{\sigma_{k+1}(\mathbf{L}) - \frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L})}{\frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L}) + \epsilon}$ , we have

$$\sigma_{k+1}(\mathbf{L}) - \sigma_k(\mathbf{L}) = \text{reg}(\mathbf{L})(\bar{\sigma} + \epsilon) + \bar{\sigma} - \sigma_k(\mathbf{L}), \quad (14)$$

177 where  $\bar{\sigma} = \frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L}) \geq \epsilon$ . Invoking (14) into (13), we arrive at

$$\begin{aligned} \text{dist}(\mathcal{C}, \mathcal{C}') &< \frac{3\delta}{\text{reg}(\mathbf{L})(\bar{\sigma} + \epsilon) + \bar{\sigma} - \sigma_k(\mathbf{L})} \\ &\leq \frac{3\delta}{2\epsilon \text{reg}(\mathbf{L}) + \bar{\sigma} - k\bar{\sigma}} \\ &\leq \frac{3\delta}{2\epsilon \text{reg}(\mathbf{L}) + (1-k)\eta\epsilon} \\ &\leq \frac{1.5\delta\epsilon^{-1}}{\text{reg}(\mathbf{L}) + (1-k)\eta/2}. \end{aligned}$$

178 This finished the proof. □

179 **F.4 Proof for Proposition A.1**

*Proof.* Since  $\hat{\mathbf{c}}$  is the optimal solution, we have

$$\begin{aligned}\phi(\mathbf{x}_i)^\top (\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}) + \lambda\hat{c}_i &= 0, \\ \phi(\mathbf{x}_j)^\top (\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}) + \lambda\hat{c}_j &= 0.\end{aligned}$$

180 It follows that

$$\begin{aligned}\|\hat{c}_i - \hat{c}_j\| &= \|(\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))^\top (\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}})\| \\ &\leq \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\| \|\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}\| \\ &= \sqrt{k(\mathbf{x}_i, \mathbf{x}_i) - 2k(\mathbf{x}_i, \mathbf{x}_j) + k(\mathbf{x}_j, \mathbf{x}_j)} \\ &\quad \times \|\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}\| \\ &= \sqrt{2 - 2k(\mathbf{x}_i, \mathbf{x}_j)} \|\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}\| \\ &\leq \sqrt{2 - 2k(\mathbf{x}_i, \mathbf{x}_j)} \|\phi(\mathbf{y})\| \\ &= \sqrt{2 - 2 \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\zeta^2}\right)}.\end{aligned}\tag{15}$$

181 In the second and last equalities, we used the fact that  $\|\phi(\mathbf{y})\| = \|\phi(\mathbf{x})\| = 1$ . In the second  
182 inequality, we used the fact that  $\frac{1}{2}\|\phi(\mathbf{y}) - \phi(\mathbf{X})\hat{\mathbf{c}}\|^2 + \frac{\lambda}{2}\|\hat{\mathbf{c}}\|^2 \leq \frac{1}{2}\|\phi(\mathbf{y}) - \phi(\mathbf{X})\mathbf{0}\|^2 + \frac{\lambda}{2}\|\mathbf{0}\|^2 =$   
183  $\frac{1}{2}\|\phi(\mathbf{y})\|^2$  because  $\hat{\mathbf{c}}$  is the optimal solution.  $\square$

184 **F.5 Proof for Theorem 3.7**

185 *Proof.* Invoking the SVD of  $\mathbf{X}$  into the closed-form solution of LSR, we get

$$\mathbf{C} = \mathbf{V} \text{diag}\left(\frac{\sigma_1^2}{\sigma_1^2 + \lambda}, \dots, \frac{\sigma_n^2}{\sigma_n^2 + \lambda}\right) \mathbf{V}^\top.\tag{16}$$

186 It means

$$\begin{aligned}c_{it} &= \sum_{l=1}^n \frac{v_{il}v_{jl}\sigma_l^2}{\sigma_l^2 + \lambda} \\ &= \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_t - \sum_{l=1}^d \frac{v_{il}v_{tl}\lambda}{\sigma_l^2 + \lambda} + \sum_{l=d+1}^n \frac{v_{il}v_{tl}\sigma_l^2}{\sigma_l^2 + \lambda}.\end{aligned}\tag{17}$$

187 Suppose  $j \in C_{\pi(i)}$  and  $k \in [n] \setminus C_{\pi(i)}$ . We have

$$\begin{aligned}&|c_{ij}| - |c_{ik}| \\ &= \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j - \sum_{l=1}^d \frac{v_{il}v_{jl}\lambda}{\sigma_l^2 + \lambda} + \sum_{l=d+1}^n \frac{v_{il}v_{jl}\sigma_l^2}{\sigma_l^2 + \lambda} \right| \\ &\quad - \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_k - \sum_{l=1}^d \frac{v_{il}v_{kl}\lambda}{\sigma_l^2 + \lambda} + \sum_{l=d+1}^n \frac{v_{il}v_{kl}\sigma_l^2}{\sigma_l^2 + \lambda} \right| \\ &\geq \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j \right| - \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_k \right| - \left| \sum_{l=1}^d \frac{v_{il}v_{jl}\lambda}{\sigma_l^2 + \lambda} \right| - \left| \sum_{l=1}^d \frac{v_{il}v_{kl}\lambda}{\sigma_l^2 + \lambda} \right| \\ &\quad - \left| \sum_{l=d+1}^n \frac{v_{il}v_{jl}\sigma_l^2}{\sigma_l^2 + \lambda} \right| - \left| \sum_{l=d+1}^n \frac{v_{il}v_{kl}\sigma_l^2}{\sigma_l^2 + \lambda} \right| \\ &\geq \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j \right| - \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_k \right| - 2\mu \sum_{l=1}^d \frac{\lambda}{\sigma_l^2 + \lambda} - 2\mu \sum_{l=d+1}^n \frac{\sigma_l^2}{\sigma_l^2 + \lambda} \\ &\geq \left| \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j \right| - \beta - \frac{2\mu d\lambda}{\sigma_d^2 + \lambda} - \frac{2\mu a\sigma_{d+1}^2}{\sigma_{d+1}^2 + \lambda},\end{aligned}\tag{18}$$

188 where  $a = \min(m, n) - d = m - d$ .

189 To ensure that there exist at least  $\bar{\tau}$  elements of  $\{|c_{ij}| : j \in C_{\pi(i)}\}$  greater than  $|c_{ik}|$  for all  $k \in$   
 190  $[n] \setminus C_{\pi(i)}$ , we need

$$|\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j| - \beta - \frac{2\mu d\lambda}{\sigma_d^2 + \lambda} - \frac{2\mu\alpha\sigma_{d+1}^2}{\sigma_{d+1}^2 + \lambda} > 0 \quad (19)$$

191 holds at least for  $\bar{\tau}$  different  $j$ , where  $j \in C_{\pi(i)}$ . It is equivalent to ensure that

$$\alpha - \beta - \frac{2\mu d\lambda}{\sigma_d^2 + \lambda} - \frac{2\mu\alpha\sigma_{d+1}^2}{\sigma_{d+1}^2 + \lambda} > 0. \quad (20)$$

192 We rewrite (20) as

$$u_1\lambda^2 + u_2\lambda + u_3 > 0, \quad (21)$$

193 where  $u_1 = \alpha - \beta - 2\mu d$ ,  $u_2 = (\alpha - \beta)(\sigma_d^2 + \sigma_{d+1}^2) - 2\mu(d+a)\sigma_{d+1}^2$ , and  $u_3 = (\alpha - \beta - 2\mu\alpha)\sigma_d^2\sigma_{d+1}^2$ .

194 The definition of  $\mu$ ,  $\alpha$ , and  $\beta$  imply  $u_1 < 0$ . Then we solve (21) and obtain

$$\begin{cases} \lambda > \frac{2\mu m\sigma_{d+1}^2 - (\alpha - \beta)(\sigma_d^2 + \sigma_{d+1}^2) - \sqrt{w}}{2(2\mu d - (\alpha - \beta))} \\ \lambda < \frac{2\mu m\sigma_{d+1}^2 - (\alpha - \beta)(\sigma_d^2 + \sigma_{d+1}^2) + \sqrt{w}}{2(2\mu d - (\alpha - \beta))} \end{cases} \quad (22)$$

195 where  $w = u_2^2 - 4u_1u_3$ . To simplify the notations, we let  $\Delta = \alpha - \beta$ ,  $\sigma_{d+1} = \gamma\sigma_d$  and get

$$\begin{cases} \lambda > \frac{(2\mu m\gamma^2 - \Delta(1 + \gamma^2) - \sqrt{(\Delta(1 + \gamma^2) - 2\mu m\gamma^2)^2 - 4(\Delta - 2\mu d)(\Delta - 2\mu m + 2\mu d)})\sigma_d^2}{4\mu d - 2\Delta} \\ \lambda < \frac{(2\mu m\gamma^2 - \Delta(1 + \gamma^2) + \sqrt{(\Delta(1 + \gamma^2) - 2\mu m\gamma^2)^2 - 4(\Delta - 2\mu d)(\Delta - 2\mu m + 2\mu d)})\sigma_d^2}{4\mu d - 2\Delta} \end{cases} \quad (23)$$

196 Further, let  $\rho = 2\mu m\gamma^2 - \Delta(1 + \gamma^2)$ , we arrive at

$$\begin{cases} \lambda > \frac{(\rho - \sqrt{\rho^2 - 4(\Delta - 2\mu d)(\Delta - 2\mu m + 2\mu d)})\sigma_d^2}{4\mu d - 2\Delta} \\ \lambda < \frac{(\rho + \sqrt{\rho^2 - 4(\Delta - 2\mu d)(\Delta - 2\mu m + 2\mu d)})\sigma_d^2}{4\mu d - 2\Delta} \end{cases} \quad (24)$$

197 That means, if (24) holds, for every  $i$ , the indices of the largest  $\bar{\tau}$  absolute elements in the  $i$ -th  
 198 column of  $C$  are in  $C_{\pi(i)}$ . Therefore, the truncation operation with parameter  $\tau \leq \bar{\tau}$  ensures the  
 199 subspace detection property. This finished the proof.

200 □

## 201 F.6 Proof for Proposition 3.8

202 *Proof.* The condition of reg means

$$\frac{\sigma_{k+1}(\mathbf{L}) - \frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L})}{\frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L}) + \epsilon} = \frac{\sigma_{k+1}(\mathbf{L})}{\epsilon} > 0.$$

203 For convenience, denote  $\vartheta = \frac{1}{k} \sum_{i=1}^k \sigma_i(\mathbf{L})$ . We have

$$-\vartheta\epsilon = \vartheta\sigma_{k+1}.$$

204 It indicates  $\vartheta = 0$  and  $\sigma_{k+1} \neq 0$ . Therefore the graph has exactly  $k$  connected components. Since  
 205 the subspace or manifold detection property hold for  $\mathbf{A}$ , each component of  $G$  is composed of the  
 206 columns of  $\mathbf{X}$  in the same subspace or manifold. Thus, all the columns of  $\mathbf{X}$  in the same subspace  
 207 or manifold must be in the same component. Otherwise, the number of connected components is  
 208 larger than  $k$ . □

## 209 F.7 Proof for Theorem C.3

210 The proof is nearly the same as that for Theorem 3.7, except that  $d < n$  and  $\text{rank}(\mathbf{K}_0) \leq k \binom{r+pq}{pq}$ ,  
 211 where  $\mathbf{K}_0 = \phi(\mathbf{X}_0)^\top \phi(\mathbf{X}_0)$ . In this case,  $\mathbf{K}$  can be well approximately by a low-rank matrix of  
 212 rank at most  $k \binom{r+pq}{pq}$  provided that the noise is small enough. More details about  $\mathbf{K}_0$  can be found  
 213 in Fan *et al.* [2020].

214 **F.8 Proof for Proposition C.4**

215 *Proof.* We only need to provide an example of  $\mathbf{W}_1 \in \mathbb{R}^{d \times m}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{k \times d}$ ,  $\mathbf{b}_1 \in \mathbb{R}^d$ , and  $\mathbf{b}_2 \in \mathbb{R}^k$ ,  
 216 where  $d = kr$ , such that the clusters can be recognized by k-means.

217 We organize the rows of  $\mathbf{W}_1$  into  $k$  groups:  $\mathbf{W}_1^j \in \mathbb{R}^{r \times m}$ ,  $j = 1, \dots, k$ . Let  $\mathbf{W}_1^j = \mathbf{U}^{j\top}$ ,  
 218  $j = 1, \dots, k$ . Let  $\mathbf{W}_1 \mathbf{x}_i = \boldsymbol{\alpha}_i = (\alpha_i^1, \dots, \alpha_i^r)$ . When  $\mathbf{x}_i \in \mathcal{S}_j$ , we have

$$\boldsymbol{\alpha}_i^j = \mathbf{U}^{j\top} \mathbf{x}_i = \mathbf{U}^{j\top} \mathbf{U}^j \mathbf{v}_i = \mathbf{v}_i. \quad (25)$$

219 It follows from the assumption that

$$\max_p \alpha_{pi}^j > \mu. \quad (26)$$

220 Let  $\mathbf{b}_1 = [\mathbf{b}_1^1; \dots; \mathbf{b}_1^k] = -\mu \mathbf{1}$ . Then  $\mathbf{h}_i^j = \text{ReLU}(\boldsymbol{\alpha}_i^j + \mathbf{b}_1^j)$  has at least one positive element. On  
 221 the other hand, since

$$\boldsymbol{\alpha}_i^l = \mathbf{U}^{l\top} \mathbf{x}_i = \mathbf{U}^{l\top} \mathbf{U}^j \mathbf{v}_i \quad l \neq j, \quad (27)$$

222 using the assumption of  $\mu$ , we have

$$|\alpha_{pi}^l| = |\mathbf{U}_{:p}^{l\top} \mathbf{U}^j \mathbf{v}_i| \leq \|\mathbf{U}_{:p}^{l\top} \mathbf{U}^j\| \|\mathbf{v}_i\| \leq \mu, \quad (28)$$

where we have used the fact  $\|\mathbf{v}_i\| = 1$  because  $\|\mathbf{x}_i\| = 1$ . It follows that

$$\mathbf{h}_i^l = \text{ReLU}(\boldsymbol{\alpha}_i^l + \mathbf{b}_1^l) = \mathbf{0}, \quad l \neq j.$$

223 Now we formulate  $\mathbf{W}_2$  as

$$\mathbf{W}_2 = \begin{bmatrix} \mathbf{q}_{11} & \mathbf{q}_{12} & \dots & \mathbf{q}_{1k} \\ \mathbf{q}_{21} & \mathbf{q}_{22} & \dots & \mathbf{q}_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_{k1} & \mathbf{q}_{k2} & \dots & \mathbf{q}_{kk} \end{bmatrix}, \quad (29)$$

where  $\mathbf{q}_{lj} \in \mathbb{R}^{1 \times r}$ ,  $l, j = 1, \dots, k$ . We have

$$z_{ji} = \mathbf{q}_{j1} \mathbf{h}_i^1 + \mathbf{q}_{j2} \mathbf{h}_i^2 \dots + \mathbf{q}_{jk} \mathbf{h}_i^k = \mathbf{q}_{jj} \mathbf{h}_i^j.$$

and

$$z_{li} = \mathbf{q}_{l1} \mathbf{h}_i^1 + \mathbf{q}_{l2} \mathbf{h}_i^2 \dots + \mathbf{q}_{lk} \mathbf{h}_i^k = \mathbf{q}_{lj} \mathbf{h}_i^j.$$

Here we have let  $\mathbf{b}_2 = \mathbf{0}$ . Let  $\mathbf{q}_{jj} \geq \mathbf{0}$  and  $\mathbf{q}_{lj} = \mathbf{0}$ , we have

$$z_{ji} > z_{li} = 0.$$

224 Therefore, if  $\mathbf{x}_i \in \mathcal{S}_j$ , we have  $z_{ji} > 0$  and  $z_{li} = 0 \forall 1 \leq j \neq l \leq k$ . Now performing k-means on  
 225  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n]$  can identify the clusters trivially.

226

□

227 **References**

- 228 Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.*, 3(1):1–122, 2011.
- 229
- 230
- 231 Joan Bruna and Stéphane Mallat. Invariant scattering convolution networks. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1872–1886, 2013.
- 232
- 233 Xinlei Chen and Deng Cai. Large scale spectral clustering with landmark-based representation. In *Twenty-fifth AAAI conference on artificial intelligence*. Citeseer, 2011.
- 234
- 235 Ying Chen, Chun-Guang Li, and Chong You. Stochastic sparse subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4155–4164, 2020.
- 236
- 237
- 238 Jicong Fan, Yuqian Zhang, and Madeleine Udell. Polynomial matrix completion for missing data imputation and transductive learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3842–3849, 2020.
- 239
- 240

- 241 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
242 *arXiv:1412.6980*, 2014.
- 243 Lee Kuang-Chih, J. Ho, and D. J. Kriegman. Acquiring linear subspaces for face recognition under  
244 variable lighting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(5):684–  
245 698, 2005.
- 246 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
247 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 248 Jun Li, Hongfu Liu, Zhiqiang Tao, Handong Zhao, and Yun Fu. Learnable subspace clustering.  
249 *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–15, 2020.
- 250 Aleix M Martínez and Avinash C Kak. PCA versus LDA. *IEEE transactions on pattern analysis*  
251 *and machine intelligence*, 23(2):228–233, 2001.
- 252 Bertil Matérn. *Spatial variation*, volume 36. Springer Science & Business Media, 2013.
- 253 Shin Matsushima and Maria Brbic. Selective sampling-based scalable sparse subspace clustering.  
254 In *Advances in Neural Information Processing Systems*, pages 12416–12425, 2019.
- 255 Marian Meila and Susan Shortreed. Regularized spectral learning. *Journal of machine learning*  
256 *research*, 1(1):1–20, 2006.
- 257 Marian Meila, Susan Shortreed, and Liang Xu. Regularized spectral learning. In *AISTATS*. PMLR,  
258 2005.
- 259 Marina Meila. The multicut lemma. *UW Statistics Technical Report*, 417, 2001.
- 260 S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (coil-20). Report, Columbia  
261 University, 1996.
- 262 Xi Peng, Lei Zhang, and Zhang Yi. Scalable sparse subspace clustering. In *Proceedings of the IEEE*  
263 *conference on computer vision and pattern recognition*, pages 430–437, 2013.
- 264 F. S. Samaria and A. C. Harter. Parameterisation of a stochastic model for human face identification.  
265 In *Proceedings of 1994 IEEE Workshop on Applications of Computer Vision*, pages 138–142, Dec  
266 1994.
- 267 Sho Sonoda and Noboru Murata. Neural network with unbounded activation functions is universal  
268 approximator. *Applied and Computational Harmonic Analysis*, 43(2):233 – 268, 2017.
- 269 Johannes Stalkamp, Marc Schlipf, Jan Salmen, and Christian Igel. Man vs. computer: Bench-  
270 marking machine learning algorithms for traffic sign recognition. *Neural networks*, 32:323–332,  
271 2012.
- 272 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-  
273 ing machine learning algorithms, 2017.
- 274 Chong You, Daniel Robinson, and René Vidal. Scalable sparse subspace clustering by orthogo-  
275 nal matching pursuit. In *Proceedings of the IEEE conference on computer vision and pattern*  
276 *recognition*, pages 3918–3927, 2016.