

# Learning Robust Simplex Sparse Representation Using Alternating Linearized Minimization Method

Zhennan Shi

Department of Computer Science,  
Colorado School of Mines  
Golden, CO, USA  
zhennanshi@mines.edu

Qiushi Wei

Department of Computer Science,  
Colorado School of Mines  
Golden, CO, USA  
qiushiwei@mines.edu

**Abstract**—The construction of graphical representations is crucial for dealing with high-dimensional data in various fields. However, creating such representations often encounters challenges of conventional methods. In this work, we present a novel methodology of a Robust Simplex Sparse Representation (RSSR). Additionally, we address the optimization of the standard simplex by introducing an effective iterative algorithm known as the Proximal Alternating Linearized Minimization method. Our algorithm achieves global convergence and convergence of the objective function and results in a notable reduction in computational expenses while bolstering resilience to outliers and noise. Extensive experiments on both synthetic and real-world datasets validate our method’s effectiveness and advantages over state-of-the-art clustering techniques. Our research offers novel insights into data representation and establishes a foundation for more precise analyses of big data.

**Index Terms**—Graphical representation, Proximal Alternating Linearized Minimization, clustering.

## I. INTRODUCTION

In the ever-evolving landscape of machine learning, one fundamental mathematical discipline emerges as a graphical representation [1]. While machine learning has made significant strides in modeling data, it is often closely related to the data itself. Hence, we focus on datasets that can be readily transformed into similar graphical representations using straightforward methods. Numerous techniques for constructing graphical representations have been developed, such as k-Nearest Neighbors (k-NN) [2], [3], Graph Minimum Spanning Tree (MST) [4], Principal Component Analysis (PCA), t-Distributed Stochastic Neighbor Embedding (t-SNE),  $\epsilon$ -neighborhood Graph [5], and Fully Connected Graph Self-Tuning Graph Construction [6]. Nevertheless, these methods still have problems. For example, PCA results in information loss, especially when the reduced dimensions do not adequately capture the variability in the original data [7]. The  $\epsilon$ -neighbor graph method needs to find the appropriate  $\epsilon$  to establish the edge connections, and the k-nearest neighbor graph relies heavily on k to connect neighbors.

Many existing sparse representation methods, such as cosine similarity, Euclidean distance, and kernel functions, are used to calculate the affinity (similarity or weight) between data points based on their sparse representations [8], [9]. However, constructing sparse graphical representations still faces huge

computational overhead for large datasets and may overlook the global structure of the data. In this work, we focus on Sparse Representation Graphs [10], [11], which is one of the most effective graph representation methods renowned for its resilience to noise and outliers. Importantly, this approach captures similarities among high-dimensional data without requiring scale consistency.

This paper makes a significant contribution to solving a non-convex and non-smooth objective function. To address this challenge, we introduce an effective solution update algorithm, known as the Proximal Alternating Linearized Minimization method [12], [13]. Through rigorous experimentation, our proposed Robust Simplex Sparse Representation (RSSR) surpasses several competing methods and exhibits a significant advantage in handling diverse datasets. The contributions of our work can be summarized as follows:

- 1) The RSSR inherits the advantages of sparse representation and exhibits robustness to issues related to scale inconsistency and outlier noise.
- 2) The proximal alternating linearized minimization method is tailored for the optimization of the objective function, which ensures the convergence of both the objective and the sequence and easily can find the globally optimal solution.
- 3) Our approach integrates a simplex constraint into sparse representation, guaranteeing shift invariance and fostering more compact representations. This contributes to improved accuracy in data interpretation and enhances computational efficiency.
- 4) Extensive experiments comprise two aspects: synthetic datasets and real-world datasets to highlight the remarkable clustering effectiveness of RSSR.

## II. ROBUST SIMPLEX SPARSE REPRESENTATION

In this section, we introduce the background of the RSSR.

### A. Background of Sparse Representation

Given a dataset consisting of  $n$  vectors with a dimension of  $d$  features, it can be arranged into columns of a training sample matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ . Given a new observation  $\mathbf{y}$  which should satisfy  $\mathbf{y} \approx \sum_{i=1}^n \mathbf{x}_i \alpha_i = \mathbf{X}\boldsymbol{\alpha}$ , a representation vector  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^T \in \mathbb{R}^n$  can be thus computed. In

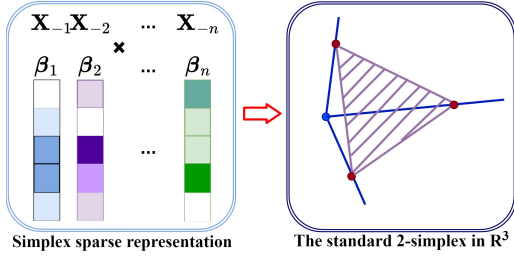


Fig. 1. Framework of the Robust Simplex Sparse Representation Method

order to obtain a sparse solution for  $\alpha$ , it is equivalent to solving the following optimization problem:

$$\min_{\alpha} \|\mathbf{X}\alpha - \mathbf{y}\|_2^2 + \lambda_0 \|\alpha\|_0, \quad (1)$$

where the pseudo-norm  $\ell_0$  quantifies the number of non-zero elements in  $\alpha$ , and  $\lambda_0 > 0$  serves as a control parameter. The sparse solution in Eq. (1) can be effectively approximated by solving the  $\ell_1$  minimization problem. Remarkably, standard linear programming techniques enable the solution of  $\ell_1$  problem in polynomial time.

$$\min_{\alpha} \|\alpha\|_1, \text{ s.t. } \mathbf{X}\alpha = \mathbf{y} \text{ or } \min_{\alpha} \|\mathbf{X}\alpha - \mathbf{y}\|_2^2 + \lambda_1 \|\alpha\|_1. \quad (2)$$

The sparse representation framework has gained recognition for its ability to handle noise and outliers effectively in image classification tasks. Importantly, it is worth mentioning that this framework does not impose any constraints on the scale consistency of the data vectors. In other words, sparse representation holds the potential to address the challenges of scale inconsistency and outliers highlighted in the introduction.

### B. RSSR Construction

The proposed approach can be employed to compute the similarity matrix  $\mathbf{S}$  in this context. Considering a data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ , where each sample is represented as a vector of dimension  $d$ . Specifically, the similarities  $\beta_i \in \mathbb{R}^{n-1}$  between the  $i$ -th sample and other samples are computed using sparse representation through the optimization problem:

$$\min_{\beta_i} \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 + \lambda \|\beta_i\|_1. \quad (3)$$

Here,  $\mathbf{X}_{-i} = [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times (n-1)}$  represents the data matrix without the  $i$  column. To ensure the non-negativity of the similarity matrix, we can impose a non-negative constraint on Eq. (3) to minimize the sum of the residue error in sparse representation.

$$\sum_{i=1}^n \min_{\beta_i \geq 0} \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 + \lambda \|\beta_i\|_1. \quad (4)$$

The parameter  $\lambda$  requires tuning in this approach. Note that Eq. (4) represents the sum of  $n$  independent variables. Therefore, for the purpose of our discussion, we will focus on a single term. Furthermore, when the data is shifted by constants  $\mathbf{b} = [b_1, \dots, b_d]^T$ , i.e.,  $\mathbf{x}_k = \mathbf{x}_k + \mathbf{b}$  for any  $k$ , the computed similarities will be altered. To ensure shift-invariance in the similarities, we introduce the following equation:

$$\|(\mathbf{X}_{-i} + \mathbf{b}\mathbf{1}^T)\beta_i - (\mathbf{x}_i + \mathbf{b})\|_2^2 = \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2. \quad (5)$$

Eq. (5) implies that  $\beta_i^T \mathbf{1} = 1$ . By adding this constraint to problem Eq. (4), we obtain:

$$\min_{\beta_i} \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 + \lambda \|\beta_i\|_1. \quad (6)$$

$$\text{s.t. } \beta_i \geq 0, \beta_i^T \mathbf{1} = 1$$

The constraints in problem Eq. (6) render the second term constant. As a result, problem Eq. (6) can be simplified to:

$$\min_{\beta_i} \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2. \quad (7)$$

$$\text{s.t. } \beta_i \geq 0, \beta_i^T \mathbf{1} = 1$$

The constraints in problem Eq. (7) form a simplex, leading to the term "simplex representation" for this type of sparse representation. In this paper, we use the Alternating Minimization (AM) method to optimize Eq. (7). Details of applying AM method will be introduced in the rest sections.

### III. SIMPLEX CONSTRAINED OPTIMIZATION

To solve the simplex constrained minimization problem, We can construct a Lagrange function, given by:

$$\mathcal{J}(\beta_i, \lambda_i, \mu_i) = \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 - \lambda_i(\beta_i^T \mathbf{1} - 1) - \mu_i^T \beta_i. \quad (8)$$

In Eq. (8), the new variable  $\lambda_i$  acts as a Lagrange multiplier for the constraint and  $\mu_i$  is a Lagrangian coefficient vector. Instead of minimizing  $\mathcal{J}$  concerning all parameters at once, we can use the AM method to minimize  $\mathcal{J}$  concerning one parameter at a time. This approach yields the following update rules when applied to the objective in Eq. (8):

$$\beta_{i(k+1)} = \mathcal{J}(\beta_i, \lambda_{i(k)}, \mu_{i(k)}), \beta_i \in \mathbf{B}_i \quad (9)$$

$$\lambda_{i(k+1)} = \mathcal{J}(\beta_{i(k)}, \lambda_i, \mu_{i(k)}), \lambda_i \in \Lambda_i \quad (10)$$

$$\mu_{i(k+1)} = \mathcal{J}(\beta_{i(k)}, \lambda_{i(k)}, \mu_i), \mu_i \in \mathbf{U}_i \quad (11)$$

To address the problem in Eq. (9), Eq. (10), and Eq. (11) arising from the constraint on each sequence set, we employ the ALM [12], [13] method. This method has been used in similar non-convex and non-smooth optimization problems. To do this, we add a proximal regularization term to the updated equations, resulting in the following formulation:

$$\beta_{i(k+1)} = \arg \min \langle \beta_i - \beta_{i(k)}, \nabla \mathcal{J}(\beta_{i(k)}) \rangle + \frac{\gamma_1}{2} \|\beta_i - \beta_{i(k)}\|_2^2, \quad (12)$$

$$\lambda_{i(k+1)} = \arg \min \langle \lambda_i - \lambda_{i(k)}, \nabla \mathcal{J}(\lambda_{i(k)}) \rangle + \frac{\gamma_2}{2} (\lambda_i - \lambda_{i(k)})^2, \quad (13)$$

$$\mu_{i(k+1)} = \arg \min \langle \mu_i - \mu_{i(k)}, \nabla \mathcal{J}(\mu_{i(k)}) \rangle + \frac{\gamma_3}{2} \|\mu_i - \mu_{i(k)}\|_2^2. \quad (14)$$

Here,  $\nabla \mathcal{J}(\beta_{i(k)})$  represents the gradient of  $\mathcal{J}$  with respect to  $\beta_i$  at iteration  $k$ , while  $\langle \beta_i - \beta_{i(k)}, \nabla \mathcal{J}(\beta_{i(k)}) \rangle$  refers to the dot product between the gradient  $\nabla \mathcal{J}(\beta_{i(k)})$  and the offset  $\beta_i - \beta_{i(k)}$ . Additionally, the regularization terms introduced in the optimization problem help prevent the updated values at each iteration from changing too rapidly. Eq. (12) to Eq. (14) can be solved by setting its derivative to zero. Thus, we obtain the following equations:

$$\beta_{i(k+1)} = \beta_{i(k)} - \frac{1}{\gamma_1} (2\mathbf{X}_{-i}^T (\mathbf{X}_{-i}\beta_{i(k)} - \mathbf{x}_i) - \lambda_i \mathbf{1} - \mu_i), \quad (15)$$

$$\lambda_{i(k+1)} = \lambda_{i(k)} - \frac{1}{\gamma_2} (1 - \beta_{i(k)}^T \mathbf{1}), \quad (16)$$

$$\mu_{i(k+1)} = \mu_{i(k)} - \frac{1}{\gamma_3}(-\beta_i). \quad (17)$$

The pseudocode for updating the three regularization terms is shown in Algorithm 1 below. Here the convergence criteria

---

**Algorithm 1:** Alternating Linearized Minimization

---

**Input:** Data Matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$ ,  $\lambda_i$ ,  $\gamma_1$ ,  $\gamma_2$ ,  $\gamma_3$ , and  $K$   
**1** Initialize  $\beta_i$ ;  
**2** **while**  $k \leq K$  **do**  
**3**     Update  $\beta_{i(k)}$  as Eq. (15);  
**4**     Update  $\lambda_{i(k)}$  as Eq. (16);  
**5**     Update  $\mu_{i(k)}$  as Eq. (17);  
**6** **return**

---

is that the relative change of  $\|\beta\|$ , the Frobenius norm of  $\beta_i$ , which indicates that  $\|\beta_{i(k+1)} - \beta_{i(k)}\|_F$  is less than  $10^{-4}$ . The empirical experiments conducted in this paper show that our algorithm converges fast and always ends within 10 iterations.

From the graphical representation we obtain, the similarity matrix can be further constructed. However, the similarity matrix  $\mathbf{S} = [\hat{\beta}_1, \dots, \hat{\beta}_n]$  is not necessarily symmetric. Thus,  $\hat{\beta}_i$  is the resulting vector of dimension  $n$  of inserting coefficient 0 in the  $i$ -th position of  $\beta_i$ . The coefficient  $\beta_{ij}$  represents the similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  for  $\mathbf{X}_{-i}$ , which does not have to be equal to  $\beta_{ji}$ . To deal with this, the final similarity matrix is obtained by:  $\mathbf{W} = \frac{1}{2}(\mathbf{S} + \mathbf{S}^T)$ . The proof of the correctness of the final similarity matrix is omitted here.

#### IV. CONVERGENCE ANALYSIS

In this section, we focus on analyzing the convergence of our solution algorithm. We start by establishing the Lipschitz continuity of the algorithm's gradient. This analysis provides the stability and convergence properties of the algorithm.

##### A. Lipschitz Continuous Gradient

We first show that  $\mathcal{J}(\beta_i, \lambda_i, \mu_i)$  has a Lipschitz continuous gradient, which assists in later convergence analysis.

**Theorem IV.1.**  $\mathcal{J}(\beta_i, \lambda_i, \mu_i)$  has Lipschitz continuous gradient for all  $\beta_i, \beta_{i(k)} \in \mathbf{B}_i$ ,  $\lambda_i, \lambda_{i(k)} \in \mathbf{\Lambda}_i$  and  $\mu_i, \mu_{i(k)} \in \mathbf{U}_i$ . That is, there exists a constant  $L_c \geq 0$  such that:

$$\|\nabla \mathcal{J}(\beta_i, \lambda_i, \mu_i) - \nabla \mathcal{J}(\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)})\|_F \leq L_c \|(\beta_i, \lambda_i, \mu_i) - (\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)})\|_F. \quad (18)$$

The  $L_c \geq 0$  here is referred to as the Lipschitz constant,  $\mathcal{J}(\beta_i, \lambda_i, \mu_i)$  is non-negative, convex and twice differentiable.

*Proof.* It is equivalent to showing  $\|\nabla^2 \mathcal{J}\|_F \leq L_c$  for all variables. The objective can be rewritten in the following forms applying the Taylor expansion:

$$\begin{aligned} \mathcal{J}(\beta_{i(k)}, \lambda_{i(k)}) &= \mathcal{J}(\beta_i, \lambda_i) + \Delta\beta_i \nabla \mathcal{J}(\beta_i) \\ &+ \Delta\lambda_i \nabla \mathcal{J}(\lambda_i) + \frac{1}{2}(\Delta\beta_i, \Delta\lambda_i) \nabla^2 \mathcal{J}(\beta_i, \lambda_i) (\Delta\beta_i, \Delta\lambda_i), \end{aligned} \quad (19)$$

and

$$\begin{aligned} \mathcal{J}(\beta_{i(k)}, \mu_{i(k)}) &= \mathcal{J}(\beta_i, \mu_i) + \Delta\beta_i \nabla \mathcal{J}(\beta_i) \\ &+ \Delta\mu_i \nabla \mathcal{J}(\mu_i) + \frac{1}{2}(\Delta\beta_i, \Delta\mu_i) \nabla^2 \mathcal{J}(\beta_i, \mu_i) (\Delta\beta_i, \Delta\mu_i). \end{aligned} \quad (20)$$

Equation (19) and Equation (20) gives us:

$$\begin{aligned} \|\nabla^2 \mathcal{J}(\beta_i, \lambda_i)\|_F &\leq \max_{\|\Delta\|_F=1} |[\nabla^2 \mathcal{J}(\beta_i, \lambda_i)](\Delta, \Delta)| \leq \\ &|\mathcal{J}(\beta_{i(k)}, \lambda_{i(k)}) - \mathcal{J}(\beta_i, \lambda_i) - \Delta\beta_i \nabla \mathcal{J}(\beta_i) \\ &- \Delta\lambda_i \nabla \mathcal{J}(\lambda_i)|. \end{aligned} \quad (21)$$

This can be expanded as:

$$\begin{aligned} \|\nabla^2 \mathcal{J}(\beta_i, \lambda_i)\|_F &\leq \|\mathbf{X}_{-i}(\beta_i + \Delta\beta_i - \mathbf{x}_i)\|_2^2 - \lambda_i((\beta_i + \Delta\beta_i)^T \mathbf{1} - 1) - \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 - \lambda_i(\beta_i^T \mathbf{1} - 1) \\ &- 2(\mathbf{X}_{-i}\beta_i - \mathbf{x}_i)\Delta\beta_i + \lambda_i \mathbf{1}\Delta\beta_i + (\beta_i^T \mathbf{1} - 1)\Delta\lambda_i \\ &:= \mathbf{C}. \end{aligned} \quad (22)$$

Similarly, Equation (20) can be extended as:

$$\begin{aligned} \|\nabla^2 \mathcal{J}(\beta_i, \mu_i)\|_F &\leq \|\mathbf{X}_{-i}(\beta_i + \Delta\beta_i - \mathbf{x}_i)\|_2^2 \\ &- \mu_i^T(\beta_i + \Delta\beta_i) - \|\mathbf{X}_{-i}\beta_i - \mathbf{x}_i\|_2^2 - 2(\mathbf{X}_{-i}\beta_i \\ &- \mathbf{x}_i)\Delta\beta_i + \lambda_i \mathbf{1}\Delta\beta_i + (\beta_i)\Delta\mu_i := \mathbf{C}. \end{aligned} \quad (23)$$

Thus proving that our objective function has a Lipschitz continuous gradient.  $\square$

##### B. Subsequence Convergence

To analyze the convergence, our objective function is:

$$\min \mathcal{J}(\beta_i, \lambda_i, \mu_i) = h(\beta_i, \lambda_i, \mu_i) + \delta_\beta(\beta_i) + \delta_\lambda(\lambda_i) + \delta_\mu(\mu_i), \quad (24)$$

where  $\delta$  is an indicator function. If the element is in the set, the function equals zero, otherwise, it approaches infinity.

**Theorem IV.2.** (Subsequence convergence) The solution subsequence  $S_k$  of Equation (24) is bounded and has three important properties:

1) *Sufficient decrease:*

$$f(S_{k-1}) - f(S_k) \geq \frac{\min(\gamma_1, \gamma_2, \gamma_3) - L_c}{2} \|S_k - S_{k-1}\|_F^2, \quad (25)$$

which implies  $\lim_{k \rightarrow \infty} \|S_{k-1} - S_k\|_F = 0$ .

2) *Convergence:*

$$\lim_{k \rightarrow \infty} \|f(S_k)\| = \text{constant}. \quad (26)$$

The difference between subsequences,

$$\begin{aligned} \lim_{k \rightarrow \infty} \|\beta_{k-1} - \beta_k\|_F &= 0, \quad \lim_{k \rightarrow \infty} \|\lambda_{k-1} - \lambda_k\|_F = 0 \\ \lim_{k \rightarrow \infty} \|\mu_{k-1} - \mu_k\|_F &= 0 \end{aligned} \quad (27)$$

3) For any convergent subsequence  $\{S_{k'}\}$ , the limit point  $S^*$  is a critical point of  $f$  and

$$\lim_{k' \rightarrow \infty} f(S_{k'}) = \lim_{k \rightarrow \infty} f(S_k) = f(S^*) \quad (28)$$

Before showing the proof, we provide necessary definitions to establish a clear understanding of the concepts involved.

**Definition IV.1.** Let  $f: \mathbb{R}^d \rightarrow (-\infty, +\infty]$  be a proper and lower semi-continuous function, given the domain of  $\sigma$  as:

$$\text{dom } f := \{x \in \mathbb{R}^d : \sigma(x) < +\infty\}. \quad (29)$$

Fréchet subdifferential of  $f(\partial f)$  at  $x \in \text{dom } f$  is the set of all vectors  $z \in \mathbb{R}^d$  representing:

$$\partial f(z) = \{z : \liminf_{y \rightarrow x} \frac{f(y) - f(x) - \langle z, y - x \rangle}{\|y - x\|} \geq 0\}. \quad (30)$$

if  $x \notin \text{dom } f$ , then  $\partial f(x) = \emptyset$ . The point  $x$  resulting in a local minimum of a domain,  $0 \in \partial f(x)$ , is called a limiting critical point of  $f$ , or critical point in simple.

*Proof.* According to the definition of the indicator function, we have the function of  $\delta$  equal to zero if the derived solution belongs to its sequence. Based on Equation (24):

$$h_{L_c}(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) \geq h(\beta_{i(k-1)}, \lambda_i, \mu_{i(k-1)}), \quad (31)$$

and by the definition of proximal map:

$$\lambda_{i(k)} = \arg \min_{\lambda_i} \delta_{\lambda_i}(\lambda_i) + \langle \nabla_{\lambda_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}), \Delta_{\lambda_i} \rangle + \frac{\gamma_2}{2} \|\Delta_{\lambda_i}\|_F^2. \quad (32)$$

Since  $\lambda_{i(k)} \leq \lambda_{i(k-1)}$ , we have:

$$\begin{aligned} \delta_{\lambda_i}(\lambda_{i(k)}) + \langle \nabla_{\lambda_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}), \Delta_{\lambda_i} \rangle + \\ \frac{\gamma_2}{2} \|\Delta_{\lambda_i}\|_F^2 \leq \delta_{\lambda_i}(\lambda_{i(k-1)}) + \langle \nabla_{\lambda_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}), \Delta_{\lambda_i} \rangle + \\ \frac{\gamma_2}{2} \|\Delta_{\lambda_i}\|_F^2 = 0. \end{aligned} \quad (33)$$

Applying Equation (33) and Lipschitz continuous gradient definition, and Taylor expansion:

$$\begin{aligned} h(\beta_{i(k-1)}, \lambda_{i(k)}, \mu_{i(k-1)}) &\leq h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) \\ &+ \langle \nabla_{\lambda_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}), \Delta_{\lambda_i} \rangle + \frac{L_c}{2} \|\Delta_{\lambda_i}\|_F^2 \\ &\leq h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) + \frac{L_c}{2} \|\Delta_{\lambda_i}\|_F^2 - \frac{\gamma_2}{2} \|\Delta_{\lambda_i}\|_F^2 \end{aligned} \quad (34)$$

$$\begin{aligned} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) - h(\beta_{i(k-1)}, \lambda_{i(k)}, \mu_{i(k-1)}) \\ \geq \frac{\gamma_2 - L_c}{2} \|\Delta_{\lambda_i}\|_F^2. \end{aligned} \quad (35)$$

Similarly, the other two sequence inequalities can be derived as:

$$\begin{aligned} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) - h(\beta_{i(k)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) \\ \geq \frac{\gamma_1 - L_c}{2} \|\Delta_{\beta_i}\|_F^2, \end{aligned} \quad (36)$$

and

$$\begin{aligned} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) - h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k)}) \\ \geq \frac{\gamma_3 - L_c}{2} \|\Delta_{\mu_i}\|_F^2. \end{aligned} \quad (37)$$

Repeating Equation (35), Equation (36) and Equation (37) together  $\forall k$  will give us Equation (25).

$$\begin{aligned} f(\beta_{i(0)}) - f(\beta_{i(N)}) \\ \geq \frac{\min(\gamma_1, \gamma_2, \gamma_3) - L_c}{2} \sum_{k=1}^N (\|\Delta_{\beta_i}\|_F^2 + \|\Delta_{\lambda_i}\|_F^2 + \|\Delta_{\mu_i}\|_F^2) \\ \leq \frac{2f(\beta_{i(0)})}{\min(\gamma_1, \gamma_2, \gamma_3) - L_c}. \end{aligned} \quad (38)$$

The sequence  $\sum_{k=1}^N (\|\Delta_{\beta_i}\|_F^2 + \|\Delta_{\lambda_i}\|_F^2 + \|\Delta_{\mu_i}\|_F^2)$  is also

convergent because it is non-decreasing and upper-bounded.

$$\lim_{i \rightarrow \infty} (\|\Delta_{\beta_i}\|_F^2 + \|\Delta_{\lambda_i}\|_F^2 + \|\Delta_{\mu_i}\|_F^2) = 0. \quad (39)$$

For critical points:

$$\begin{aligned} \lim_{\bar{k} \rightarrow \infty} f(S_{\bar{k}}) &= \lim_{\bar{k} \rightarrow \infty} h(\beta_{i\bar{k}}, \lambda_{i\bar{k}}, \mu_{i\bar{k}}) \\ &+ \delta_{\beta_i}(\beta_{i\bar{k}}) + \delta_{\lambda_i}(\lambda_{i\bar{k}}) + \delta_{\mu_i}(\mu_{i\bar{k}}) \\ &= f(\bar{S}). \end{aligned} \quad (40)$$

As a result convergence of  $\{f(S_k)\}_{k \geq 0}$ , we have

$$\lim_{k \rightarrow \infty} \|S_{(k+1)} - S_{(k)}\|_F = 0. \quad (41)$$

To show that  $\bar{S}$  is a critical point, we can consider the first-order optimal condition of Eq. (32):

$$\begin{aligned} \nabla_{\beta_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) + \gamma_1(\beta_{i(k)} - \beta_{i(k-1)}) \\ + \partial \delta_{\beta_i}(\beta_{i(k)}) = 0. \end{aligned} \quad (42)$$

similarly, given

$$\begin{aligned} \nabla_{\lambda_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) + \gamma_2(\lambda_{i(k)} - \lambda_{i(k-1)}) \\ + \partial \delta_{\lambda_i}(\lambda_{i(k)}) = 0. \end{aligned} \quad (43)$$

as well as

$$\begin{aligned} \nabla_{\mu_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)}) + \gamma_3(\mu_{i(k)} - \mu_{i(k-1)}) \\ + \partial \delta_{\mu_i}(\mu_{i(k)}) = 0. \end{aligned} \quad (44)$$

As  $k$  approaches infinity:

$$\begin{aligned} \|\nabla_{\beta_i} h(\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)}) + \partial \delta_{\beta_i}(\beta_{i(k)})\|_F \leq \\ \|\nabla_{\beta_i} h(\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)}) - \nabla_{\beta_i} h(\beta_{i(k-1)}, \lambda_{i(k-1)}, \mu_{i(k-1)})\|_F \\ + \gamma_1 \|\beta_{i(k)} - \beta_{i(k-1)}\|_F \leq (L_c + \gamma_1) \|S_k - S_{k-1}\| = 0 \end{aligned} \quad (45)$$

and

$$\begin{aligned} \|\nabla_{\lambda_i} h(\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)}) + \partial \delta_{\lambda_i}(\lambda_{i(k)})\|_F &\leq (L_c + \gamma_2) \|S_k - S_{k-1}\| = 0 \\ \|\nabla_{\mu_i} h(\beta_{i(k)}, \lambda_{i(k)}, \mu_{i(k)}) + \partial \delta_{\mu_i}(\mu_{i(k)})\|_F &\leq (L_c + \gamma_3) \|S_k - S_{k-1}\| = 0 \end{aligned} \quad (46)$$

Finally, we have:  $(2L_c + \gamma_1 + \gamma_2 + \gamma_3) \|S_k - S_{k-1}\| \geq \text{dist}(0, \partial f(S_k))$ , thus  $0 \in \partial f(\bar{S})$  since  $\partial f(S_i)$  is closed. Therefore, we can say that  $\bar{S}$  is a critical point of  $f$ .  $\square$

**Theorem IV.3. (Sequence convergence)** The sequence  $\{M_k\}_{k \geq 0}$  generated by Algorithm with a constant step size  $(\gamma_1, \gamma_2, \gamma_3 \geq L_c)$  is a global convergent sequence.

Before Presenting the proof, we will give the following definition:

**Definition IV.2.** (Kurdyka-Lojasiewicz (KL property)) we say a proper semi-continuous function  $h(u)$  satisfies Kurdyka-Lojasiewicz (KL) property, if  $\bar{u}$  is a critical point of  $h(u)$ , then there exists  $\delta > 0, \theta \in [0, 1], C_1 > 0$ , s.t.

$$|h(u) - h(\bar{u})|^\theta \leq C_1 \text{dist}(0, \partial h(u)), \forall u \in B(\bar{u}, \delta). \quad (47)$$

It is known that semi-algebraic functions satisfy the KL inequality. Clearly, the objective function  $f$  is semi-algebraic as  $h, \delta_{\beta_i}, \delta_{\lambda_i}$  and  $\delta_{\mu_i}$  are semi-algebraic. Therefore the objective function satisfies the KL property.

*Proof.* According to the KL property (in the Supplementary), the objective function is subjected to the following:

$$[f(S_k) - f(\bar{S})]^\theta \leq C_2 \text{dist}(0, \partial f(S_k)), \forall k \geq k_0. \quad (48)$$

In order to construct the concavity, we introduce a concave function of  $x^{(1-\theta)}$  for some  $\theta \in [0, 1]$  with domain  $x > 0$ ,

$$x_2^{1-\theta} - x_1^{1-\theta} \geq (1-\theta)x_2^{-\theta}(x_2 - x_1), \forall x_1 > 0, x_2 > 0. \quad (49)$$

Combined with the sufficient decrease property, Equation (49) can be updated by replacing  $x_1$  and  $x_2$  with  $f(S_{k+1}) - f(\bar{S})$

and  $f(S_k) - f(\bar{S})$ , respectively. Then we have the following:

$$\begin{aligned}
& [f(S_k) - f(\bar{S})]^{1-\theta} - [f(S_{k+1}) - f(\bar{S})]^{1-\theta} \\
& \geq (1-\theta) \frac{f(S_k) - f(S_{k+1})}{[f(S_k) - f(\bar{S})]^\theta} \geq \frac{\lambda(1-\theta)}{2C_2} \frac{\|S_k - S_{k+1}\|_F^2}{\text{dist}(0, \partial f(S_k))} \\
& \geq \frac{\lambda(1-\theta)}{2C_2C_3} \frac{\|S_k - S_{k+1}\|_F^2}{\|S_k - S_{k-1}\|_F} \\
& = \kappa \left( \frac{\|S_k - S_{k+1}\|_F^2}{\|S_k - S_{k-1}\|_F} + \|S_k - S_{k-1}\|_F \right) - \kappa \|S_k - S_{k-1}\|_F \\
& \geq \kappa (2\|S_k - S_{k+1}\|_F - \|S_k - S_{k-1}\|_F).
\end{aligned} \tag{50}$$

Accordingly:

$$\begin{aligned}
& 2\|S_k - S_{k+1}\|_F - \|S_k - S_{k-1}\|_F \\
& \leq \alpha \left( [f(S_k) - f(\bar{S})]^{1-\theta} - [f(S_{k+1}) - f(\bar{S})]^{1-\theta} \right).
\end{aligned} \tag{51}$$

with  $C_3 := 2L_c + \gamma_1 + \gamma_2$ ,  $\kappa := \frac{\gamma_2(1-\theta)}{2C_2C_3}$ ,  $\beta := \left( \frac{\gamma_2(1-\theta)}{2C_2C_3} \right)^{-1}$ . Summing the above inequalities up from some  $\tilde{k} > k_0$  to infinity yields:

$$\sum_{k=\tilde{k}}^{\infty} \|S_k - S_{k+1}\|_F \leq \|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F + \alpha [f(S_{\tilde{k}}) - f(\bar{S})]^{1-\theta} \tag{52}$$

, which indicates that  $\sum_{k=\tilde{k}}^{\infty} \|S_k - S_{k+1}\|_F < \infty$ . Following the standard argument, it is obvious that:

$$\lim_{t \rightarrow \infty} \sup \|S_{t_1} - S_{t_2}\|_F = 0, \text{ s.t. } t_1, t_2 \geq t. \tag{53}$$

, which implies that the sequence  $S_k$  is Cauchy, thus convergent. The limit point set  $C(S_0)$  is a singleton  $\bar{S}$ .  $\square$

**Theorem IV.4. (Convergence rate)** *The convergence rate is at least sub-linear.*

*Proof.* Since  $\{S_k\}$  converges to  $\bar{S}$ , i.e.,  $\lim_{k \rightarrow \infty} S^k = \bar{S}$ . From Equation (52) and the triangle inequality:

$$\|S_{\tilde{k}} - \bar{S}\|_F \leq \sum_{i=\tilde{k}}^{\infty} \|S_i - S_{i+1}\|_F \tag{54}$$

$$\leq \|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F + \alpha [f(S_{\tilde{k}}) - f(\bar{S})]^{1-\theta}.$$

which indicates the convergence rate of  $S_{\tilde{k}} \rightarrow \bar{S}$  is at least as fast as the rate that  $\|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F + \alpha [f(S_{\tilde{k}}) - f(\bar{S})]^{1-\theta}$ . In particular, the second term  $\alpha [f(S_{\tilde{k}}) - f(\bar{S})]^{1-\theta}$  can be controlled:

$$\begin{aligned}
& \alpha [f(S_{\tilde{k}}) - f(\bar{S})]^\theta \leq \alpha C_2 \text{dist}(0, \partial f(S_{\tilde{k}})) \\
& \leq \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F) \|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F.
\end{aligned} \tag{55}$$

Thus we have:

$$\sum_{i=\tilde{k}}^{\infty} \|S_i - S_{i+1}\|_F \leq \|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F \tag{56}$$

$$+ \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F) \|S_{\tilde{k}} - S_{\tilde{k}-1}\|_F^{\frac{1-\theta}{\theta}}.$$

We divide the following analysis into three cases based on the value of the KL exponent  $\theta$ :

**Case 1** if  $\theta = 0$ , we set  $Q := k \in \mathbb{N} : S_{k+1} \neq S_k$  and take  $k$  in  $Q$ . When  $k$  is large enough, then we have:

$$\|S_{k+1} - S_k\|_F^2 := C_4 > 0. \tag{57}$$

At the same time, we have:

$$\begin{aligned}
f(S_{k+1}) - f(S_k) & \geq \frac{\min(\gamma_1, \gamma_2, \gamma_3) - L_c}{2} \|S_{k+1} - S_k\|_F^2 \\
& = \frac{\min(\gamma_1, \gamma_2, \gamma_3) - L_c}{2} C_4.
\end{aligned} \tag{58}$$

Since  $f(S_k)$  is known to converge to 0, the above equation shows that  $Q$  is finite and sequence  $S_k$  converges in a finite number of steps.

**Case 2:**  $\theta \in (0, \frac{1}{2}]$ , which corresponds to  $\frac{(1-\theta)}{\theta} \geq 1$ . We define  $P_{\tilde{k}} = \sum_{k=\tilde{k}}^{\infty} \|S_{k+1} - S_k\|_F$ ,

$$P_{\tilde{k}} \leq P_{\tilde{k}-1} - P_{\tilde{k}} + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F) [P_{\tilde{k}-1} - P_{\tilde{k}}]^{\frac{1-\theta}{\theta}}, \tag{59}$$

Since  $P_{\tilde{k}-1} - P_{\tilde{k}} \rightarrow 0$ , there exists a positive integer  $k_1$  such that  $P_{\tilde{k}-1} - P_{\tilde{k}} < 1, \forall \tilde{k} \geq k_1$ . Thus,

$$\begin{aligned}
P_{\tilde{k}} & \leq [1 + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F)] (P_{\tilde{k}-1} - P_{\tilde{k}}), \\
\forall \tilde{k} & \geq \max\{k_0, k_1\}.
\end{aligned} \tag{60}$$

Let  $\rho = \frac{1 + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F)}{2 + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F)} \in (0, 1)$ , and combining with Equation (54) results in

$$\begin{aligned}
P_{\tilde{k}} & \leq [1 + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F)] (P_{\tilde{k}-1} - P_{\tilde{k}}), \\
\text{where } \tilde{k} & = \max\{k_0, k_1\}.
\end{aligned} \tag{61}$$

**Case 3:**  $\theta \in (\frac{1}{2}, 1)$ , which indicates  $\frac{1-\theta}{\theta} \leq 1$ . Based on previous results, we have:

$$\begin{aligned}
P_{\tilde{k}} & \leq [1 + \alpha C_2 (2A_0 + \lambda + \|\mathbf{X}\|_F)] [P_{\tilde{k}-1} - P_{\tilde{k}}]^{\frac{1-\theta}{\theta}}, \\
\forall \tilde{k} & \geq \max\{k_0, k_1\}.
\end{aligned} \tag{62}$$

Similarly, we get:

$$P_{\tilde{k}}^{\frac{1-2\theta}{1-\theta}} - P_{\tilde{k}-1}^{\frac{1-2\theta}{1-\theta}} \geq \zeta, \forall \tilde{k} \geq \tilde{k}, \text{ for some } \zeta > 0. \tag{63}$$

Then repeating and summing up the upper inequality from  $\tilde{k} = \max\{k_0, k_1\}$  to any  $k > \tilde{k}$ , we get:

$$\begin{aligned}
P_{\tilde{k}} & \leq \left[ P_{\tilde{k}-1}^{\frac{1-2\theta}{1-\theta}} + \zeta(\tilde{k} - \tilde{k}) \right]^{\frac{\theta-1}{2\theta-1}} \\
& = \mathcal{O} \left( (\tilde{k} - \tilde{k})^{\frac{\theta-1}{2\theta-1}} \right).
\end{aligned} \tag{64}$$

We thus prove the sub-linear convergence:

$$\|S_k - \bar{S}\|_F \leq \mathcal{O} \left( (\tilde{k} - \tilde{k})^{\frac{\theta-1}{2\theta-1}} \right), \forall k > \tilde{k}. \tag{65}$$

$\square$

TABLE I  
DATA SET DESCRIPTION

Data Set	Size	Dimension	Classes
<i>Wine</i>	178	13	3
<i>Iris</i>	150	4	3
<i>Glass</i>	214	9	6
<i>Breast Cancer</i>	569	30	2
<i>Albalone</i>	4177	8	3
<i>AT&amp;T</i>	400	4096	40
<i>USPS</i>	7291	256	10

## V. EXPERIMENT AND RESULTS

In this section, we apply the RSSR model to both synthetic data and real-world datasets to evaluate the model's performance. Additionally, we compare the clustering performance of the RSSR model with that of the original dataset.

### A. Synthetic Data Experiment

To visually demonstrate the effectiveness of our proposed algorithm, we conduct experiments on a synthetic dataset

TABLE II

CLUSTERING PERFORMANCE OF DIFFERENT METHODS ON REAL-WORLD DATA SETS. THE EVALUATION METRICS ARE AVERAGED WITH 50 RUNS.

ACC	Representation	NMF	LSTM	K-means	NCut	STSC	Hierarchical	EgoNetSplitter	MNMF	SymmNMF	GEMSEC
Wine	Original	0.882	0.360	0.376	0.438	0.354	0.337	0.158	<b>0.539</b>	0.163	0.333
	RSSR	<b>0.888</b>	<b>0.764</b>	<b>0.972</b>	<b>0.618</b>	<b>0.371</b>	<b>0.376</b>	<b>0.202</b>	0.494	<b>0.484</b>	<b>0.473</b>
Iris	Original	0.587	0.127	0.833	0.847	<b>0.387</b>	0.687	0.073	0.500	0.087	0.487
	RSSR	<b>0.647</b>	<b>0.667</b>	<b>0.867</b>	<b>0.887</b>	0.320	<b>0.740</b>	<b>0.133</b>	<b>0.533</b>	<b>0.093</b>	<b>0.522</b>
Glass	Original	0.256	0.126	0.304	<b>0.372</b>	0.140	<b>0.336</b>	0.137	0.145	0.112	0.383
	RSSR	<b>0.435</b>	<b>0.136</b>	<b>0.388</b>	0.341	<b>0.154</b>	0.327	<b>0.210</b>	<b>0.383</b>	<b>0.136</b>	<b>0.458</b>
Breast Cancer	Original	0.580	0.362	0.634	0.376	0.496	0.374	0.011	0.509	0.069	0.524
	RSSR	<b>0.647</b>	0.359	<b>0.645</b>	<b>0.622</b>	<b>0.503</b>	<b>0.376</b>	<b>0.019</b>	<b>0.536</b>	<b>0.076</b>	<b>0.634</b>
NMI	Representation	NMF	LSTM	K-means	NCut	STSC	Hierarchical	EgoNetSplitter	MNMF	SymmNMF	GEMSEC
Wine	Original	<b>0.708</b>	0.248	0.876	<b>0.572</b>	0.002	0.018	0.480	0.089	0.375	0.266
	RSSR	0.701	<b>0.924</b>	<b>0.893</b>	0.562	<b>0.005</b>	<b>0.549</b>	<b>0.508</b>	<b>0.252</b>	<b>0.385</b>	<b>0.276</b>
Iris	Original	0.692	0.734	0.657	0.690	<b>0.290</b>	<b>0.713</b>	0.491	0.214	0.386	<b>0.227</b>
	RSSR	<b>0.720</b>	<b>0.771</b>	<b>0.696</b>	<b>0.730</b>	0.010	0.704	<b>0.522</b>	<b>0.258</b>	<b>0.404</b>	<b>0.253</b>
Glass	Original	0.424	0.486	0.471	<b>0.310</b>	0.132	0.101	0.506	0.286	0.466	0.326
	RSSR	<b>0.480</b>	<b>0.531</b>	<b>0.530</b>	0.100	<b>0.156</b>	<b>0.345</b>	<b>0.512</b>	<b>0.319</b>	<b>0.492</b>	<b>0.387</b>
Breast Cancer	Original	0.014	0.038	0.024	0.005	0.290	0.002	0.057	0.214	0.022	0.007
	RSSR	<b>0.027</b>	<b>0.041</b>	<b>0.037</b>	<b>0.730</b>	<b>0.503</b>	<b>0.005</b>	<b>0.058</b>	<b>0.714</b>	<b>0.028</b>	<b>0.015</b>
Purity	Representation	NMF	LSTM	K-means	NCut	STSC	Hierarchical	EgoNetSplitter	MNMF	SymmNMF	GEMSEC
Wine	Original	0.719	0.785	0.966	<b>0.635</b>	0.399	0.399	0.927	0.497	0.817	0.567
	RSSR	<b>0.899</b>	<b>0.978</b>	<b>0.972</b>	0.629	<b>0.404</b>	<b>0.629</b>	<b>0.944</b>	<b>0.685</b>	<b>0.865</b>	<b>0.579</b>
Iris	Original	0.687	0.667	0.831	0.847	<b>0.413</b>	0.687	0.920	0.533	0.880	0.513
	RSSR	<b>0.707</b>	<b>0.873</b>	<b>0.867</b>	<b>0.887</b>	0.393	<b>0.740</b>	<b>0.974</b>	<b>0.633</b>	<b>0.907</b>	<b>0.520</b>
Glass	Original	0.631	0.640	0.603	0.379	0.463	0.379	0.781	0.549	0.745	0.580
	RSSR	<b>0.673</b>	<b>0.696</b>	<b>0.738</b>	<b>0.477</b>	<b>0.463</b>	<b>0.486</b>	<b>0.790</b>	<b>0.673</b>	<b>0.832</b>	<b>0.664</b>
Breast Cancer	Original	0.627	0.638	0.634	0.847	0.627	0.543	<b>0.809</b>	0.533	0.663	0.628
	RSSR	<b>0.647</b>	<b>0.641</b>	<b>0.645</b>	<b>0.887</b>	<b>0.629</b>	<b>0.627</b>	0.714	<b>0.627</b>	<b>0.675</b>	<b>0.638</b>

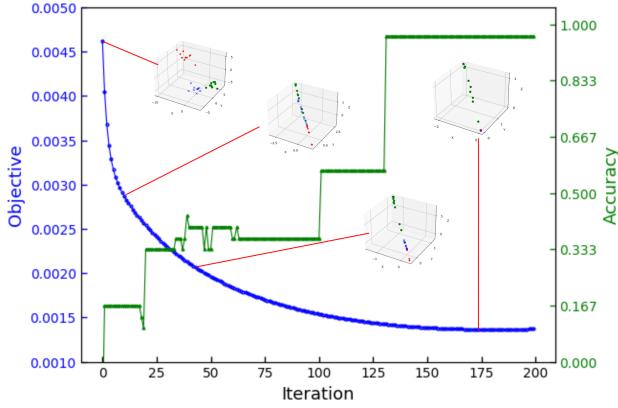
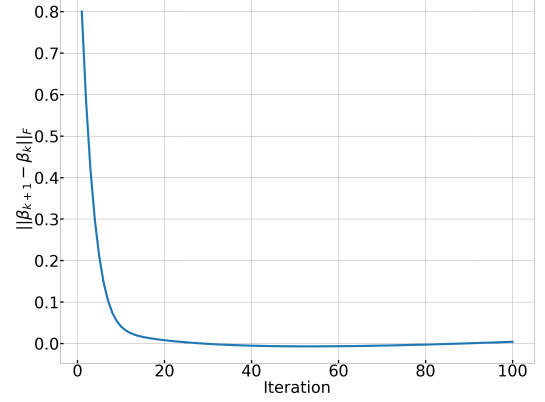


Fig. 2. Objective function with the accuracy on the noisy synthetic dataset

consisting of three centroids in a 20-dimensional space, each surrounded by 30 data points. Three centroids are normal distributed vectors with a dimension of 20 are constructed. Each centroid has a different mean value. The goal is to do the clustering for these three centroids of vectors. In addition, noises (mean is set as 0, standard deviation is set as 0.5) are added to the distributed vectors to avoid trivial testing. Experiments are run 100 times with different random seeds. Three standard metrics for clustering are deployed for comparing the clustering performances among various clustering methods.

In Fig. 2, our objective function monotonically decreases with each iteration, indicating the convergence of our algorithm. For the purpose of visualization, we embed the data into a 3-dimensional space ( $r = 3$ ), where the  $x$ ,  $y$ , and  $z$  axes correspond to the first, second, and third columns of the original matrix  $X$ . In addition, the accuracy stabilizes after 100 iterations and the accuracy reaches 0.96 once stabilized. From the figure, it is evident that the clustering performance on the simplex representation improves as the algorithm it-

Fig. 3.  $\|\beta_{k+1} - \beta_k\|_F$  plotted as RSSR iterates on the noisy synthetic dataset

erates as the accuracy gets higher. This experiment not only demonstrates the effectiveness of our proposed method but also empirically confirms the sub-linear convergence of our objective function and the convergence of the solution. As shown in Fig. 3, the  $\beta$  variance  $\|\beta_{k+1} - \beta_k\|_F$  converges to 0 in less than 10 iterations and stabilizes in 40 iterations.

### B. Real Data Experiment

Table I summarizes the datasets used in the real experiment section. We extracted these datasets from the UCI machine learning repository. In this section, we compare our proposed new method with the NMF [14], LTSM [15], K-means [16], Normalized cut (NCut) [17], Self-Tuning Spectral Clustering (STSC) [18], DBSCAN [19], EgoNetSplitter [20], MNMF [21], SymmNMF [22], GEMSEM [23] methods. Results will be compared including clustering Accuracy (Acc.), Normalized Mutual Information (NMI) and Purity.

The parameter setting for the real dataset experiment is as follows: We split each real dataset into a training and a test set with a proportion of 80% and 20% upon the total number

of samples. During each iteration step, our model learns an updated similarity matrix  $\mathbf{S}$  which will be applied to be used as the simplex representation for the dataset. Once the  $\mathbf{W}$  (recall that  $\mathbf{W} = \frac{1}{2}(\mathbf{S} + \mathbf{S}^T)$ ) has been learned, we run different clustering methods in the test set with the original representation  $\mathbf{X}$  and the simplex representation  $\mathbf{XW}$  and compare these two based on the three metrics: Acc, NMI, and purity. We use a 5-fold cross-validation within the training set to identify the best hyperparameters. The validation set in this inner cross-validation is chosen to be 20% of the training set.

### C. Original vs Representations

We run several methods for clustering after the new representation is learned. In Table II, each experiment repeats 50 times, and the average is taken. The clustering result is then reported with the original representation versus the simplex representation. We see that RSSR always achieves higher clustering accuracy, normalized mutual information, and purity in analysis. We also compare our method with other methods on large datasets as Table III demonstrates, where since the dataset is large, we split the whole data into small batches due to the high computation cost of the counterpart algorithms and report the average result. Although our method may not consistently outperform others on small datasets, the results remain promising and generally reliable. Overall, the experiments demonstrate strong performance across all datasets, highlighting its potential for broader applications.

TABLE III  
CLUSTERING PERFORMANCE OF DIFFERENT METHODS ON LARGE  
REAL-WORLD DATA SETS, AVERAGED OVER 50 RUNS.

ACC	Representation	NMF	LSTM	K-means	NCut	STSC
Albalone	Original	0.409	0.145	0.431	0.368	0.362
	RSSR	<b>0.483</b>	<b>0.164</b>	<b>0.561</b>	<b>0.382</b>	<b>0.463</b>
AT&T	Original	0.217	0.253	0.264	<b>0.323</b>	<b>0.412</b>
	RSSR	<b>0.376</b>	<b>0.266</b>	<b>0.298</b>	0.312	0.380
USPS	Original	0.165	0.088	0.173	<b>0.171</b>	0.170
	RSSR	<b>0.173</b>	<b>0.164</b>	<b>0.179</b>	0.165	<b>0.182</b>
NMI	Representation	NMF	LSTM	K-means	NCut	STSC
Albalone	Original	0.168	0.323	0.175	<b>0.088</b>	0.104
	RSSR	<b>0.213</b>	<b>0.353</b>	<b>0.181</b>	0.062	<b>0.115</b>
AT&T	Original	0.484	0.443	0.464	0.524	<b>0.556</b>
	RSSR	<b>0.570</b>	<b>0.518</b>	<b>0.581</b>	<b>0.546</b>	0.437
USPS	Original	0.267	0.253	0.257	0.281	0.113
	RSSR	<b>0.294</b>	<b>0.330</b>	<b>0.324</b>	<b>0.302</b>	<b>0.196</b>
Purity	Representation	NMF	LSTM	K-means	NCut	STSC
Albalone	Original	0.326	0.124	0.358	<b>0.422</b>	0.475
	RSSR	<b>0.334</b>	<b>0.141</b>	<b>0.374</b>	0.343	<b>0.484</b>
AT&T	Original	0.333	0.363	0.383	0.378	<b>0.345</b>
	RSSR	<b>0.415</b>	<b>0.428</b>	<b>0.411</b>	<b>0.407</b>	0.326
USPS	Original	0.145	0.160	0.137	0.147	0.096
	RSSR	<b>0.164</b>	<b>0.257</b>	<b>0.231</b>	<b>0.177</b>	<b>0.107</b>

## VI. CONCLUSION

In this work, our proposed Robust Simplex Sparse Representation methodology, coupled with the Proximal Alternating Linearized Minimization algorithm, addresses challenges in constructing graphical representations for high-dimensional data. This approach reduces computation, improves robustness to outliers, and outperforms state-of-the-art clustering methods on diverse datasets. Our research provides valuable insights into data representation challenges for big data.

## REFERENCES

- [1] Z. Peng, W. Huang, M. Luo, Q. Zheng, Y. Rong, T. Xu, and J. Huang, "Graph representation learning via graphical mutual information maximization," in *Proceedings of The Web Conference*, 2020, pp. 259–270.
- [2] J. Laaksonen and E. Oja, "Classification with learning k-nearest neighbors," in *Proceedings of international conference on neural networks (ICNN'96)*, vol. 3. IEEE, 1996, pp. 1480–1483.
- [3] L. Kozma, "k nearest neighbors algorithm (knn)," *Helsinki University of Technology*, vol. 32, 2008.
- [4] Y. Xu, V. Olman, and D. Xu, "Clustering gene expression data using a graph-theoretic approach: an application of minimum spanning trees," *Bioinformatics*, vol. 18, no. 4, pp. 536–545, 2002.
- [5] M. H. Rohban and H. R. Rabiee, "Supervised neighborhood graph construction for semi-supervised classification," *Pattern Recognition*, vol. 45, no. 4, pp. 1363–1372, 2012.
- [6] G. Wen, "Robust self-tuning spectral clustering," *Neurocomputing*, vol. 391, pp. 243–248, 2020.
- [7] J. J. Yanez-Borjas, J. M. Machorro-Lopez, D. Camarena-Martinez, M. Valtierra-Rodriguez, J. P. Amezcua-Sanchez, F. J. Carrion-Viramontes, and J. A. Quintana-Rodriguez, "A new damage index based on statistical features, pca, and mahalanobis distance for detecting and locating cables loss in a cable-stayed bridge," *International Journal of Structural Stability and Dynamics*, vol. 21, no. 09, p. 2150127, 2021.
- [8] Y. Liu, Q. Gao, J. Han, and S. Wang, "Euler sparse representation for image classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [9] A. Shrivastava, V. M. Patel, and R. Chellappa, "Multiple kernel learning for sparse representation-based classification," *IEEE Transactions on Image Processing*, vol. 23, no. 7, pp. 3013–3024, 2014.
- [10] K. Lu, Z. Ding, and S. Ge, "Sparse-representation-based graph embedding for traffic sign recognition," *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 4, pp. 1515–1524, 2012.
- [11] A. Kyrillidis, S. Becker, V. Cevher, and C. Koch, "Sparse projections onto the simplex," in *International Conference on Machine Learning*. PMLR, 2013, pp. 235–243.
- [12] J. Bolte, S. Sabach, and M. Teboulle, "Proximal alternating linearized minimization for nonconvex and nonsmooth problems," *Mathematical Programming*, vol. 146, no. 1-2, pp. 459–494, 2014.
- [13] T. Pock and S. Sabach, "Inertial proximal alternating linearized minimization (ipalm) for nonconvex and nonsmooth problems," *SIAM journal on imaging sciences*, vol. 9, no. 4, pp. 1756–1787, 2016.
- [14] D. Guillamet, M. Bressan, and J. Vitria, "A weighted non-negative matrix factorization for local representations," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1. IEEE, 2001, pp. I–I.
- [15] Q. Wang, C. Downey, L. Wan, P. A. Mansfield, and I. L. Moreno, "Speaker Diarization with LSTM," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5239–5243.
- [16] A. Likas, N. Vlassis, and J. J. Verbeek, "The global k-means clustering algorithm," *Pattern recognition*, vol. 36, no. 2, pp. 451–461, 2003.
- [17] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 8, pp. 888–905, 2000.
- [18] L. Zelnik-Manor and P. Perona, "Self-tuning spectral clustering," *Advances in neural information processing systems*, vol. 17, 2004.
- [19] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *kdd*, vol. 96, no. 34, 1996, pp. 226–231.
- [20] A. Epasto, S. Lattanzi, and R. Paes Leme, "Ego-splitting framework: From non-overlapping to overlapping clusters," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017, pp. 145–154.
- [21] X. Wang, P. Cui, J. Wang, J. Pei, W. Zhu, and S. Yang, "Community preserving network embedding," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017.
- [22] D. Kuang, C. Ding, and H. Park, "Symmetric nonnegative matrix factorization for graph clustering," in *Proceedings of the 2012 SIAM international conference on data mining*. SIAM, 2012, pp. 106–117.
- [23] B. Rozemberczki, R. Davies, R. Sarkar, and C. Sutton, "Gemsec: Graph embedding with self clustering," in *Proceedings of the 2019 IEEE/ACM international conference on advances in social networks analysis and mining*, 2019, pp. 65–72.

## Supplementary Material

### A. Derivation details

In this section, we provide the derivation of the solution for each variable discussed in Section. We take the derivative of the equation Eq. (12) to Eq. (14) with respect to each parameter and set it to zero to earn the update for each parameter.

#### Step 1: Updating $\beta_{i(k+1)}$

$$\nabla \mathcal{J}(\beta_{i(k)}) = \frac{\partial \mathcal{J}(\beta_{i(k)}, \lambda_i, \mu_i)}{\partial (\beta_{i(k)})} \quad (66)$$

$$= 2\mathbf{X}_{-i}^T(\mathbf{X}_{-i}\beta_{i(k)} - \mathbf{x}_i) - \lambda_i \mathbf{1} - \mu_i$$

Set derivative Eq. (12) to 0:

$$\begin{aligned} & \frac{\partial(\langle \beta_i - \beta_{i(k)}, \nabla \mathcal{J}(\beta_{i(k)}) \rangle + \frac{\gamma_1}{2} \|\beta_i - \beta_{i(k)}\|_2^2)}{\partial (\beta_i)} \\ &= \frac{\partial(\langle \beta_i - \beta_{i(k)}, \nabla \mathcal{J}(\beta_{i(k)}) \rangle)}{\partial (\beta_i)} + \frac{\partial(\frac{\gamma_1}{2} \|\beta_i - \beta_{i(k)}\|_2^2)}{\partial (\beta_i)} \\ &= \frac{\partial(\beta_i - \beta_{i(k)})}{\partial (\beta_i)} \nabla \mathcal{J}(\beta_{i(k)}) + \frac{\partial(\nabla \mathcal{J}(\beta_{i(k)}))}{\partial (\beta_i)} (\beta_i - \beta_{i(k)}) \\ & \quad + \frac{\partial(\frac{\gamma_1}{2} \|\beta_i - \beta_{i(k)}\|_2^2)}{\partial (\beta_i)} \\ &= \mathbf{I} \nabla \mathcal{J}(\beta_{i(k)}) + 0 + \gamma_1 (\beta_i - \beta_{i(k)}) \\ &= 0 \end{aligned} \quad (67)$$

Plug in Eq. (66), this gives us:

$$2\mathbf{X}_{-i}^T(\mathbf{X}_{-i}\beta_{i(k)} - \mathbf{x}_i) - \lambda_i \mathbf{1} - \mu_i + \gamma_1 (\beta_i - \beta_{i(k)}) = 0 \quad (68)$$

Which gives the updated equation (15).

#### Step 2: Updating $\lambda_{i(k+1)}$

$$\nabla \mathcal{J}(\lambda_{i(k)}) = \frac{\partial \mathcal{J}(\beta_i, \lambda_{i(k)}, \mu_i)}{\partial (\lambda_{i(k)})} = 1 - \beta_i^T \mathbf{1} \quad (69)$$

Set derivative Eq. (13) to 0:

$$\begin{aligned} & \frac{\partial(\langle \lambda_i - \lambda_{i(k)}, \nabla \mathcal{J}(\lambda_{i(k)}) \rangle + \frac{\gamma_2}{2} (\lambda_i - \lambda_{i(k)})^2)}{\partial \lambda_i} \\ &= \frac{\partial(\langle \lambda_i - \lambda_{i(k)}, \nabla \mathcal{J}(\lambda_{i(k)}) \rangle)}{\partial \lambda_i} + \frac{\partial(\frac{\gamma_2}{2} (\lambda_i - \lambda_{i(k)})^2)}{\partial \lambda_i} \\ &= \frac{\partial(\lambda_i - \lambda_{i(k)})}{\partial \lambda_i} \nabla \mathcal{J}(\lambda_{i(k)}) + \frac{\partial \nabla \mathcal{J}(\lambda_{i(k)})}{\partial \lambda_i} \\ & \quad + \frac{\partial(\frac{\gamma_2}{2} (\lambda_i - \lambda_{i(k)})^2)}{\partial \lambda_i} \\ &= \nabla \mathcal{J}(\lambda_{i(k)}) + 0 + \gamma_2 (\lambda_i - \lambda_{i(k)}) \\ &= 1 - \beta_{i(k)}^T \mathbf{1} + \gamma_2 (\lambda_i - \lambda_{i(k)}) = 0 \end{aligned} \quad (70)$$

Which gives us the updated equation (16).

#### Step 3: Updating $\mu_{i(k+1)}$

$$\nabla \mathcal{J}(\mu_{i(k)}) = \frac{\partial \mathcal{J}(\beta_i, \lambda_i, \mu_{i(k)})}{\partial (\mu_{i(k)})} = -\beta_i \quad (71)$$

Set derivative Eq. (14) to 0:

$$\begin{aligned} & \frac{\partial(\langle \mu_i - \mu_{i(k)}, \nabla \mathcal{J}(\mu_{i(k)}) \rangle + \frac{\gamma_3}{2} \|\mu_i - \mu_{i(k)}\|_2^2)}{\partial \mu_i} \\ &= \frac{\partial(\mu_i - \mu_{i(k)})}{\partial \mu_i} \nabla \mathcal{J}(\mu_{i(k)}) + \frac{\partial(\nabla \mathcal{J}(\mu_{i(k)}))}{\partial \mu_i} (\mu_i - \mu_{i(k)}) \\ & \quad + \frac{\partial(\frac{\gamma_3}{2} \|\mu_i - \mu_{i(k)}\|_2^2)}{\partial \mu_i} \\ &= \mathbf{I} \nabla \mathcal{J}(\mu_{i(k)}) + 0 + \gamma_3 (\mu_i - \mu_{i(k)}) \\ &= -\beta_i + \gamma_3 (\mu_i - \mu_{i(k)}) = 0 \end{aligned} \quad (72)$$

Which gives us the updated equation (17).

### B. Convergence Analysis Related Lemma

**Lemma VI.1.** [Uniform KL Property] The objective function satisfies the KL property. There exists  $\delta_0 > 0$ ,  $\theta_{KL} \in [0, 1]$ ,  $C_{KL} > 0$  such that as long as  $\text{dist}((S), \mathbb{C}(S_0)) \leq \delta_0$ :  
 $|f(S) - \bar{f}|^{\theta_{KL}} \leq C_{KL} \text{dist}(0, \partial f(S))$ , (73)  
 with  $\bar{f}$  representing the limiting function value defined in assertion (i) of Theorem (IV.2).

*Proof.* First, we recognize that the union  $\bigcup_k B(\overline{S_k}, \delta_k)$  forms an open cover of  $\mathbb{C}(S_0)$  with  $\overline{S_k}$  denoting all points in  $\mathbb{C}(S_0)$  and  $\delta_k$  to be chosen so that the KL property of  $f$  at  $\overline{S_k} \in \mathbb{C}(S_0)$  holds:

$|f(S) - \bar{f}|^{\theta_k} \leq C_i \text{dist}(0, \partial f(S)), \forall (S) \in B(\overline{S_k}, \delta_k)$ , (74)  
 in which we have used all  $f(\overline{S_k}) = \bar{f}$  by assertion (iii) of Theorem (IV.2). Then since  $\mathbb{C}(S_0)$  is compact, it has a finite subcover  $\bigcup_{k=1}^p B(\overline{S_{s_i}}, \delta_{s_i})$  for some positive integer  $p$ . Considering all these facts, we have  $\forall S \in \bigcup_{k=1}^p B(\overline{S_{s_i}}, \delta_{s_i})$ :  
 $|f(S) - \bar{f}|^{\theta_{KL}} \leq C_{KL} \text{dist}(0, \partial f(S))$ , (75)  
 with  $\theta_{KL} = \max_{k=1}^p \{\theta_{s_k}\}$  and  $C_{KL} = \max_{k=1}^p \{C_{s_k}\}$ . In the end, since  $\bigcup_{k=1}^p B(\overline{S_{s_i}}, \delta_{s_i})$  is an open cover of  $\mathbb{C}(S_0)$ , there should exist a sufficiently small number  $\delta_0$ , so that

$$\{(S) : \text{dist}(S, \mathbb{C}(S_0)) \leq \delta_0\} \subset \bigcup_{k=1}^p B(\overline{S_k}, \delta_{s_k}). \quad (76)$$

Thus Equation (75) holds when  $\text{dist}(S, \mathbb{C}(S_0)) \leq \delta_0$ .  $\square$