

NOISY ℓ^0 -SPARSE SUBSPACE CLUSTERING ON DIMENSIONALITY REDUCED DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

High-dimensional data often lie in or close to low-dimensional subspaces. Sparse subspace clustering methods with sparsity induced by ℓ^0 -norm, such as ℓ^0 -Sparse Subspace Clustering (ℓ^0 -SSC) Yang et al. (2016), are demonstrated to be more effective than its ℓ^1 counterpart such as Sparse Subspace Clustering (SSC) Elhamifar & Vidal (2013). However, these ℓ^0 -norm based subspace clustering methods are restricted to clean data that lie exactly in subspaces. Real data often suffer from noise and they may lie close to subspaces. We propose noisy ℓ^0 -SSC to handle noisy data so as to improve the robustness. We show that the optimal solution to the optimization problem of noisy ℓ^0 -SSC achieves subspace detection property (SDP), a key element with which data from different subspaces are separated, under deterministic and randomized models. Our results provide theoretical guarantee on the correctness of noisy ℓ^0 -SSC in terms of SDP on noisy data. We further propose Noisy-DR- ℓ^0 -SSC which provably recovers the subspaces on dimensionality reduced data. Noisy-DR- ℓ^0 -SSC first projects the data onto a lower dimensional space by linear transformation, then performs noisy ℓ^0 -SSC on the dimensionality reduced data so as to improve the efficiency. The experimental results demonstrate the effectiveness of noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC.

1 INTRODUCTION

Clustering is an important unsupervised learning procedure for analyzing a broad class of scientific data in biology, medicine, psychology and chemistry. On the other hand, high-dimensional data, such as facial images and gene expression data, often lie in low-dimensional subspaces in many cases, and clustering in accordance to the underlying subspace structure is particularly important. For example, the well-known Principal Component Analysis (PCA) works perfectly if the data are distributed around a single subspace. The subspace learning literature develops more general methods that recover multiple subspaces in the original data, and subspace clustering algorithms Vidal (2011) aim to partition the data such that data belonging to the same subspace are identified as one cluster. Among various subspace clustering algorithms, the ones that employ sparsity prior, such as Sparse Subspace Clustering (SSC) Elhamifar & Vidal (2013) and ℓ^0 -Sparse Subspace Clustering (ℓ^0 -SSC) Yang et al. (2016), have been proven to be effective in separating the data in accordance with the subspaces that the data lie in under certain assumptions.

Sparse subspace clustering methods construct the sparse similarity matrix by sparse representation of the data. Subspace detection property (SDP) defined in Section 4.1 ensures that the similarity between data from different subspaces vanishes in the sparse similarity matrix, and applying spectral clustering Ng et al. (2001) on such sparse similarity matrix leads to compelling clustering performance. Elhamifar and Vidal Elhamifar & Vidal (2013) prove that when the subspaces are independent or disjoint, SDP can be satisfied by solving the canonical sparse linear representation problem using data as the dictionary, under certain conditions on the rank, or singular value of the data matrix and the principle angle between the subspaces. SSC has been successfully applied to a novel deep neural network architecture, leading to the first deep sparse subspace clustering method Peng et al. (2016). Under the independence assumption on the subspaces, low rank representation Liu et al. (2010; 2013) is also proposed to recover the subspace structures. Relaxing the assumptions on the subspaces to allowing overlapping subspaces, the Greedy Subspace Clustering Park et al. (2014) and the Low-Rank Sparse Subspace Clustering Wang et al. (2013) achieve subspace detection property with high

probability. The geometric analysis in Soltanolkotabi & Cands (2012) shows the theoretical results on subspace recovery by SSC. In the following text, we use the term SSC or ℓ^1 -SSC exchangeably to indicate the Sparse Subspace Clustering method in Elhamifar & Vidal (2013).

Real data often suffer from noise. Noisy SSC proposed in Wang & Xu (2013) handles noisy data that lie close to disjoint or overlapping subspaces. While ℓ^0 -SSC Yang et al. (2016) has guaranteed clustering correctness via subspace detection property under much milder assumptions than previous subspace clustering methods including SSC, it assumes that the observed data lie in exactly in the subspaces and does not handle noisy data. In this paper, we present noisy ℓ^0 -SSC, which enhances ℓ^0 -SSC by theoretical guarantee on the correctness of clustering on noisy data. It should be emphasized that while ℓ^0 -SSC on clean data Yang et al. (2016) empirically adopts a form of optimization problem robust to noise, it lacks theoretical analysis on the correctness of ℓ^0 -SSC on noisy data. In this paper, the correctness of noisy ℓ^0 -SSC on noisy data in terms of the subspace detection property is established. Our analysis is under both deterministic model and randomized models, which is also the model employed in the geometric analysis of SSC Soltanolkotabi & Cands (2012). Our randomized analysis demonstrates potential advantage of noisy ℓ^0 -SSC over its ℓ^1 counterpart as more general assumption on data distribution can be adopted. Moreover, we present Noisy Dimensionality Reduced ℓ^0 -Sparse Subspace Clustering (Noisy-DR- ℓ^0 -SSC), an efficient version of noisy ℓ^0 -SSC which also enjoys robustness to noise. Noisy-DR- ℓ^0 -SSC first projects the data onto a lower dimensional space by random projection, then performs noisy ℓ^0 -SSC on the dimensionality reduced data. Noisy-DR- ℓ^0 -SSC provably recovers the underlying subspace structure in the original data from the dimensionality reduced data under deterministic model. Experimental results demonstrate the effectiveness of both noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC.

We use bold letters for matrices and vectors, and regular lower letter for scalars throughout this paper. The bold letter with superscript indicates the corresponding column of a matrix, e.g. \mathbf{A}^i is the i -th column of matrix \mathbf{A} , and the bold letter with subscript indicates the corresponding element of a matrix or vector. $\|\cdot\|_F$ and $\|\cdot\|_p$ denote the Frobenius norm and the vector ℓ^p -norm or the matrix p -norm, and $\text{diag}(\cdot)$ indicates the diagonal elements of a matrix. $\mathbf{H}_{\mathbf{T}} \subseteq \mathbb{R}^d$ indicates the subspace spanned by the columns of \mathbf{T} , and $\mathbf{A}_{\mathbf{I}}$ denotes a submatrix of \mathbf{A} whose columns correspond to the nonzero elements of \mathbf{I} (or with indices in \mathbf{I} without confusion). $\sigma_t(\cdot)$ denotes the t -th largest singular value of a matrix, and $\sigma_{\min}(\cdot)$ indicates the smallest singular value of a matrix. $\text{supp}(\cdot)$ is the support of a vector, $\mathbb{P}_{S'}$ is an operator indicating projection onto the subspace S' .

2 PROBLEM SETUP

2.1 NOTATIONS

We hereby introduce the notations for subspace clustering on noisy data considered in this paper. The uncorrupted data matrix is denoted by $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{d \times n}$, where d is the dimensionality and n is the size of the data. The uncorrupted data \mathbf{Y} lie in a union of K distinct subspaces $\{\mathcal{S}_k\}_{k=1}^K$ of dimensions $\{d_k\}_{k=1}^K$. The observed noisy data is $\mathbf{X} = \mathbf{Y} + \mathbf{N}$, where $\mathbf{N} = [\mathbf{n}_1, \dots, \mathbf{n}_n] \in \mathbb{R}^{d \times n}$ is the additive noise. $\mathbf{x}_i = \mathbf{y}_i + \mathbf{n}_i$ is the noisy data point that is corrupted by the noise \mathbf{n}_i .

Let $\mathbf{Y}^{(k)} \in \mathbb{R}^{d \times n_k}$ denote the data belonging to subspace \mathcal{S}_k with $\sum_{k=1}^K n_k = n$, and denote the

corresponding columns in \mathbf{X} by $\mathbf{X}^{(k)}$. The data \mathbf{X} are normalized such that each column has unit ℓ^2 -norm in our deterministic analysis. We consider deterministic noise model where the noise \mathbf{Z} is fixed and $\max \|\mathbf{n}_i\| \leq \delta$. Note that our analysis can be extended to a random noise model which is common and also considered by noisy SSC Wang & Xu (2013), and the random noise model assumes that columns of \mathbf{Z} are sampled i.i.d. and $\max \|\mathbf{n}_i\| \leq \delta$ with high probability. Note that such random noise model does not require spherical symmetric noise as that in Wang & Xu (2013).

2.2 METHOD

ℓ^0 -SSC Yang et al. (2016) proposes to solve the following ℓ^0 sparse representation problem

$$\min_{\mathbf{Z}} \|\mathbf{Z}\|_0 \quad \text{s.t. } \mathbf{X} = \mathbf{X}\mathbf{Z}, \text{diag}(\mathbf{Z}) = \mathbf{0}, \quad (1)$$

and it proves that the subspace detection property defined in Definition 1 is satisfied with the globally optimal solution to (1). We resort to solve the ℓ^0 regularized sparse approximation problem below to

handle noisy data for ℓ^0 -SSC, which is the optimization problem of noisy ℓ^0 -SSC:

$$\min_{\mathbf{Z} \in \mathbb{R}^n \times \mathbb{R}^n, \text{diag}(\mathbf{Z})=\mathbf{0}} L(\mathbf{Z}) = \|\mathbf{X} - \mathbf{X}\mathbf{Z}\|_F^2 + \lambda\|\mathbf{Z}\|_0. \quad (2)$$

The definition of subspace detection property for noisy ℓ^0 -SSC and noiseless ℓ^0 -SSC, i.e. ℓ^0 -SSC on noiseless data, is defined in Definition 1 below.

Definition 1. (*Subspace detection property for noisy and noiseless ℓ^0 -SSC*) Let \mathbf{Z}^* be the optimal solution to (2). The subspaces $\{\mathcal{S}_k\}_{k=1}^K$ and the data \mathbf{X} satisfy subspace detection property for noisy ℓ^0 -SSC if \mathbf{Z}^i is a nonzero vector, and nonzero elements of \mathbf{Z}^i correspond to the columns of \mathbf{X} from the same subspace as \mathbf{y}_i for all $1 \leq i \leq n$.

Similarly, in the noiseless setting where $\mathbf{X} = \mathbf{Y}$, let \mathbf{Z}^* be the optimal solution to (1). The subspaces $\{\mathcal{S}_k\}_{k=1}^K$ and the data \mathbf{X} satisfy the subspace detection property for noiseless ℓ^0 -SSC if \mathbf{Z}^i is a nonzero vector, and nonzero elements of \mathbf{Z}^i correspond to the columns of \mathbf{X} that from the same subspace as \mathbf{y}_i for all $1 \leq i \leq n$.

We say that subspace detection property holds for \mathbf{x}_i if nonzero elements of \mathbf{Z}^{*i} correspond to the data that lie in the same subspace as \mathbf{y}_i , for either noisy ℓ^0 -SSC or noiseless ℓ^0 -SSC.

2.3 MODELS

Similar to Soltanolkotabi & Cands (2012), we introduce the deterministic, semi-random and fully-random models for the analysis of noisy ℓ^0 -SSC.

- **Deterministic Model:** the subspaces and the data in each subspace are fixed.
- **Semi-Random Model:** the subspaces are fixed but the data are independent and identically distributed in each of the subspaces.
- **Fully-Random Model:** both the subspaces and the data of each subspace are independent and identically distributed.

The data in the above definitions refer to clean data without noise. We refer to semi-random model and fully-random model as randomized models in this paper. All the three models are extensively employed to analyze the subspace detection property in the subspace learning literature Soltanolkotabi & Cands (2012); Wang et al. (2013); Wang & Xu (2013); Yining Wang & Singh (2016).

3 THEORETICAL ANALYSIS FOR NOISY ℓ^0 -SSC

The theoretical results on the subspace detection property for noisy ℓ^0 -SSC are presented in this section under deterministic model and randomized models.

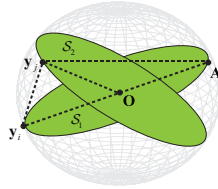


Figure 1: Illustration of an external subspace. All the data \mathbf{Y} are normalized to have unit norm for illustration purpose, so they lie on the surface of the sphere. \mathcal{S}_1 and \mathcal{S}_2 are two subspaces in the three-dimensional ambient space. The subspace spanned by $\mathbf{y}_i \in \mathcal{S}_1$ and $\mathbf{y}_j \in \mathcal{S}_2$ is an external subspace, and the intersection of this external subspace and \mathcal{S}_1 is a dashed line $\mathbf{y}_i\mathbf{O}\mathbf{A}$.

3.1 NOISY ℓ^0 -SSC: DETERMINISTIC ANALYSIS

We introduce the definition of general position and external subspace before our analysis on noisy ℓ^0 -SSC.

Definition 2. (General position) For any $1 \leq k \leq K$, the data $\mathbf{Y}^{(k)}$ are in general position if any subset of $L \leq d_k$ data points (columns) of $\mathbf{Y}^{(k)}$ are linearly independent. \mathbf{Y} are in general position if $\mathbf{Y}^{(k)}$ are in general position for $1 \leq k \leq K$.

The assumption of general condition is rather mild. In fact, if the data points in $\mathbf{X}^{(k)}$ are independently distributed according to any continuous distribution, then they almost surely in general position.

Let the distance between a point $\mathbf{x} \in \mathbb{R}^d$ and a subspace $\mathcal{S} \subseteq \mathbb{R}^d$ be defined as $d(\mathbf{x}, \mathcal{S}) = \inf_{\mathbf{y} \in \mathcal{S}} \|\mathbf{x} - \mathbf{y}\|_2$, the definition of external subspaces is presented as follows. Figure 1 illustrates an example of external subspace.

Definition 3. (External subspace) For a point $\mathbf{y} \in \mathbf{Y}^{(k)}$, a subspace $\mathbf{H}_{\{\mathbf{y}_{i_j}\}_{j=1}^L}$ spanned by a set of linear independent points $\{\mathbf{y}_{i_j}\}_{j=1}^L \subseteq \mathbf{Y}$ is defined to be an external subspace of \mathbf{y} if $\{\mathbf{y}_{i_j}\}_{j=1}^L \not\subseteq \mathbf{Y}^{(k)}$ and $\mathbf{y} \notin \{\mathbf{y}_{i_j}\}_{j=1}^L$. The point \mathbf{y} is said to be away from its external subspaces if $\min_{\mathbf{H} \in \mathcal{H}_{\mathbf{y}, d_k}} d(\mathbf{y}, \mathbf{H}) > 0$, where $\mathcal{H}_{\mathbf{y}, d}$ are the set of all external subspaces of \mathbf{y} of dimension no greater than d for \mathbf{y} , i.e. $\mathcal{H}_{\mathbf{y}, d} = \{\mathbf{H} : \mathbf{H} = \mathbf{H}_{\{\mathbf{y}_{i_j}\}_{j=1}^L}, \dim[\mathbf{H}] = L, L \leq d, \{\mathbf{y}_{i_j}\}_{j=1}^L \not\subseteq \mathbf{Y}^{(k)}, \mathbf{y} \notin \{\mathbf{y}_{i_j}\}_{j=1}^L\}$. All the data points in $\mathbf{Y}^{(k)}$ are said to be away from the external subspaces if each of them is away from the its associated external spaces.

Remark 1. (Subspace detection property holds for noiseless ℓ^0 -SSC under the deterministic model) It can be verified that the following statement is true. Under the deterministic model, suppose data is noiseless, $n_k \geq d_k + 1$, $\mathbf{Y}^{(k)}$ is in general position. If all the data points in $\mathbf{Y}^{(k)}$ are away from the external subspaces for any $1 \leq k \leq K$, then the subspace detection property for ℓ^0 -SSC holds with the optimal solution \mathbf{Z}^* to (1).

To present our theoretical results of the correctness of noisy ℓ^0 -SSC, we also need the definitions of the minimum restricted eigenvalue and the subspace separation margin, which are defined as follows. In the following analysis, we employ β to denote the sparse code of datum \mathbf{x}_i so that a simpler notation other than \mathbf{Z}^i is dedicated to our analysis.

Definition 4. The minimum restricted eigenvalue of the uncorrupted data is defined as

$$\sigma_{\mathbf{Y}, r} \triangleq \min_{\beta: \|\beta\|_0 = r, \text{rank}(\mathbf{Y}\beta) = \|\beta\|_0} \sigma_{\min}(\mathbf{Y}\beta) \quad (3)$$

for $r \geq 1$. In addition, the normalized minimum restricted eigenvalue of the uncorrupted data is defined by

$$\bar{\sigma}_{\mathbf{Y}, r} \triangleq \frac{\sigma_{\mathbf{Y}, r}}{\sqrt{r}} \quad (4)$$

We have the following perturbation bound for the distance between a data point and the subspaces spanned by noisy and noiseless data, which is useful to establish the conditions when the subspace detection property holds for noisy ℓ^0 -SSC.

Lemma 1. Let $\beta \in \mathbb{R}^n$ and $\mathbf{Y}\beta$ has full column rank. Suppose $\delta < \bar{\sigma}_{\mathbf{Y}, r}$ where $r = \|\beta\|_0$, then $\mathbf{X}\beta$ is a full column rank matrix, and

$$|d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}\beta}) - d(\mathbf{x}_i, \mathbf{H}_{\mathbf{Y}\beta})| \leq \frac{\delta}{\bar{\sigma}_{\mathbf{Y}, r} - \delta} \quad (5)$$

for any $1 \leq i \leq n$.

The optimization problem of noisy ℓ^0 -SSC (2) is separable. For each $1 \leq i \leq n$, the optimization problem with respect to the sparse code of i -th data point is

$$\min_{\beta \in \mathbb{R}^n, \beta_i = 0} L(\beta) = \|\mathbf{x}_i - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_0. \quad (6)$$

Lemma 2 shows that the optimal solution to the noisy ℓ^0 -SSC problem (6) is also that to a ℓ^0 -minimization problem with tolerance to noise.

Lemma 2. Let nonzero vector β^* be the optimal solution to the noisy ℓ^0 -SSC problem (6) for point \mathbf{x}_i with $\|\beta^*\|_0 = r^* > 1$. If $\lambda > \tau_0$ where τ_0 is defined as

$$\tau_0 \triangleq \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*} + \tau_1,$$

where

$$\tau_1 \triangleq \frac{\delta}{\bar{\sigma}_{\mathbf{Y}}^* - \delta}, \quad \sigma_{\mathbf{X}}^* \triangleq \sigma_{\min}(\mathbf{X}\beta^*),$$

with $\delta < \bar{\sigma}_{\mathbf{Y}}^*$, and $\bar{\sigma}_{\mathbf{Y}}^*$ is defined as

$$\bar{\sigma}_{\mathbf{Y}}^* \triangleq \min_{1 \leq r < r^*} \bar{\sigma}_{\mathbf{Y},r},$$

then β^* is the optimal solution to the following sparse approximation problem with the uncorrupted data as the dictionary:

$$\min_{\beta} \|\beta\|_0 \quad \text{s.t.} \quad \|\mathbf{x}_i - \mathbf{Y}\beta\|_2 \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}, \quad \beta_i = 0. \quad (7)$$

where $c^* \triangleq \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2$.

Define $\mathbf{B}(\mathbf{x}_i, c_0) = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_i\| \leq c_0\}$ be the ball centered at \mathbf{x}_i with radius c_0 . If $\mathbf{B}(\mathbf{x}_i, c_0)$ is away from the corresponding confusion area, i.e. all the external subspaces in $\mathcal{H}_{\mathbf{y}_i, d_k}$, then subspace detection property holds with the solution to a proper sparse approximation problem where \mathbf{x}_i is approximated by the uncorrupted data, as shown in the following Lemma.

Lemma 3. *Suppose \mathbf{Y} is in general position and $\mathbf{y}_i \in \mathcal{S}_k$ for some $1 \leq k \leq K$. For positive number c_0 such that $c_0 \geq d(\mathbf{x}_i, \mathcal{S}_k)$, suppose $\mathbf{B}(\mathbf{x}_i, c_0) \cap \mathbf{H} = \emptyset$ for any $\mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}$. Then the subspace detection property holds for \mathbf{x}_i with the optimal solution to the following sparse approximation problem, denoted by β^* , i.e. nonzero elements of β^* correspond to the columns of \mathbf{X} from the same subspace as \mathbf{y}_i .*

$$\min_{\beta} \|\beta\|_0 \quad \text{s.t.} \quad \|\mathbf{x}_i - \mathbf{Y}\beta\|_2 \leq c_0, \quad \beta_i = 0. \quad (8)$$

Now we use the above results to present the main result on the correctness of noisy ℓ^0 -SSC.

Theorem 1. (Subspace detection property holds for noisy ℓ^0 -SSC) *Let nonzero vector β^* be the optimal solution to the noisy ℓ^0 -SSC problem (6) for point \mathbf{x}_i with $\|\beta^*\|_0 = r^* > 1$, and $c^* \triangleq \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2$. Suppose \mathbf{Y} is in general position, $\mathbf{y}_i \in \mathcal{S}_k$ for some $1 \leq k \leq K$, $\delta < \bar{\sigma}_{\mathbf{Y}}^*$, $\lambda > \tau_0$, $\mathbf{B}(\mathbf{y}_i, \delta + c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}) \cap \mathbf{H} = \emptyset$ for any $\mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}$. Then the subspace detection property holds for \mathbf{x}_i with β^* . Here τ_0 , τ_1 , $\bar{\sigma}_{\mathbf{Y}}^*$ and $\sigma_{\mathbf{X}}^*$ are defined in Lemma 2.*

Remark 2. *When $\delta = 0$ and there is no noise in the data \mathbf{X} , the conditions for the correctness of noisy ℓ^0 -SSC in Theorem 1 almost reduce to that for noiseless ℓ^0 -SSC. To see this, the conditions are reduced to $\mathbf{B}(\mathbf{y}_i, c^*) \cap \mathbf{H} = \emptyset$, which are exactly the conditions required by noiseless ℓ^0 -SSC, namely data are away from the external subspaces by choosing $\lambda \rightarrow 0$ and it follows that $c^* = 0$.*

While Theorem 1 establishes geometric conditions under which the subspace detection property holds for noisy ℓ^0 -SSC, it can be seen that these conditions are often coupled with the optimal solution β^* to the noisy ℓ^0 -SSC problem (6). In the following theorem, the correctness of noisy ℓ^0 -SSC is guaranteed in terms of λ , the weight for the ℓ^0 regularization term in (6), and the geometric conditions independent of the optimal solution to (6).

Let $M_i > 0$ be the minimum distance between $\mathbf{y}_i \in \mathcal{S}_k$ and its external subspaces when \mathbf{y}_i is away from its external subspaces, i.e.

$$M_i \triangleq \min\{d(\mathbf{y}_i, \mathbf{H}) : \mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}\}, \quad (9)$$

The following two quantities related to the spectrum of clean and noisy data, μ_r and $\sigma_{\mathbf{X},r}$, are defined as follows with $r > 1$ for the analysis in Theorem 2.

$$\mu_r \triangleq \frac{\delta}{\min_{1 \leq r' < r} \bar{\sigma}_{\mathbf{Y},r'} - \delta}, \quad (10)$$

$$\sigma_{\mathbf{X},r} \triangleq \min\{\sigma_{\min}(\mathbf{X}\beta) : 1 \leq \|\beta\|_0 \leq r\} \quad (11)$$

Theorem 2. (Subspace detection property holds for noisy ℓ^0 -SSC under deterministic model, with conditions in terms of λ) *Let nonzero vector β^* be the optimal solution to the noisy ℓ^0 -SSC problem (6) for point \mathbf{x}_i with $\|\beta^*\|_0 = r^*$, $n_k \geq d_k + 1$ for every $1 \leq k \leq K$, and there exists $1 < r_0 \leq d$*

such that $1 < r^* \leq r_0$. Suppose \mathbf{Y} is in general position, $\mathbf{y}_i \in \mathcal{S}_k$ for some $1 \leq k \leq K$, $\delta < \min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y}, r}$, and $M_{i, \delta} \triangleq M_i - \delta$. Suppose

$$M_{i, \delta} > \frac{2\delta}{\sigma_{\mathbf{X}, r_0}}, \quad (12)$$

and

$$\mu_{r_0} < 1 - \frac{2\delta}{\sigma_{\mathbf{X}, r_0}}. \quad (13)$$

Then if

$$\lambda_0 < \lambda < 1, \quad (14)$$

where $\lambda_0 \triangleq \max\{\lambda_1, \lambda_2\}$ and

$$\lambda_1 \triangleq \inf\{0 < \lambda < 1: \sqrt{1 - \lambda} + \frac{2\delta}{\sigma_{\mathbf{X}, r_0} \sqrt{\lambda}} < M_{i, \delta}\}, \quad (15)$$

$$\lambda_2 \triangleq \inf\{0 < \lambda < 1: \lambda - \frac{2\delta}{\sigma_{\mathbf{X}, r_0}} \frac{1}{\sqrt{\lambda}} > \mu_{r_0}\}, \quad (16)$$

the subspace detection property holds for \mathbf{x}_i with β^* . Here M_i , μ_{r_0} and $\sigma_{\mathbf{X}, r_0}$ are defined in (9), (10) and (11) respectively.

Remark 3. The two conditions (12) and (13) are induced by the conditions that $\mathbf{B}(\mathbf{y}_i, \delta + c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}}) \cap \mathbf{H} = \emptyset$ for any $\mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}$, and $\lambda > \tau_0$ in Theorem 1. Note that when (12) and (13) hold, λ_1 and λ_2 can always be chosen in accordance with (15) and (16).

Remark 4. It can be observed from condition (14) that noisy ℓ^0 -SSC encourages sparse solution by a relatively large λ so as to guarantee the subspace detection property. This theoretical finding is consistent with the empirical study shown in the experimental results.

3.2 NOISY ℓ^0 -SSC: RANDOMIZED ANALYSIS

In this subsection, the correctness of noisy ℓ^0 -SSC is analyzed when the clean data in each subspace are distributed at random. We assume that the data in subspace $\mathcal{S}^{(k)}$ are i.i.d. isotropic samples on sphere of radius $\sqrt{d_k}$ centered at the origin according to some continuous distribution, for $1 \leq k \leq K$. A random variable $\mathbf{y} \in \mathcal{S}_k$ is isotropic if $\mathbb{E}[\mathbf{y}\mathbf{y}^\top] = \mathbf{I}_{d_k}$, where \mathbf{I}_{d_k} is the $d_k \times d_k$ identity matrix (with \mathbf{y} represented as a vector in \mathbb{R}^{d_k}). In addition, for each $1 \leq k \leq K$, we assume that the following condition holds:

(a) There exists a constant $M \geq 1$ such that for any $t > 0$, any $\mathbf{y} \in \mathbf{Y}^{(k)}$, and any vector \mathbf{v} with unit ℓ^2 -norm,

$$\Pr[|\langle \mathbf{y}, \mathbf{v} \rangle| > t] \leq \frac{M}{t^4}. \quad (17)$$

Intuitively, condition (a) requires that the projection of any data point onto arbitrary unit vector is bounded from both sides with relatively large probability. This condition is also required in Yaskov (2014) to derive lower bound for the least singular value of a random matrix with independent isotropic columns. In order to meet the conditions in Theorem 2 so as to guarantee the subspace detection property under randomized models, the following lemma is presented and it provides the geometric concentration inequality for the distance between a point $\mathbf{y} \in \mathbf{Y}^{(k)}$ and any of its external subspaces. It renders a lower bound for M_i , namely the minimum distance between $\mathbf{y}_i \in \mathcal{S}_k$ and its external subspaces.

Lemma 4. Under randomized models, given $1 \leq k \leq K$ and $\mathbf{y} \in \mathbf{Y}^{(k)}$, suppose $\mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}$ is any external subspace of \mathbf{y} . Then for any $t > 0$,

$$\Pr[d(\mathbf{y}, \mathbf{H}) \geq 1 - 2t\sqrt{d_k - 1} - t^2] \geq 1 - 8\exp\left(-\frac{d_k t^2}{2}\right). \quad (18)$$

We then have the following results regarding to the subspace detection property of noisy ℓ^0 -SSC under randomized models.

Theorem 3. (Subspace detection property holds for noisy ℓ^0 -SSC under randomized models, with conditions in terms of λ) *Under randomized models, let nonzero vector β^* be the optimal solution to the noisy ℓ^0 -SSC problem (6) for point \mathbf{x}_i with $\|\beta^*\|_0 = r^*$, $n_k \geq d_k + 1$ for every $1 \leq k \leq K$, and there exists $1 < r_0 \leq d$ such that $1 < r^* \leq r_0$. Suppose the data in each subspace are i.i.d. isotropic samples according to some continuous distribution that satisfies condition (a). Let $d_{\max} \triangleq \max_k d_k$, $c \triangleq \frac{1}{\sqrt{r_0}(\sqrt{196Md+1}+14\sqrt{Md})}$. For $t > 0$ such that $1 - 2t\sqrt{d_{\max} - 1} - t^2 > 0$, suppose*

$$\delta < c, \quad (19)$$

$$\delta + \frac{2\delta}{\sqrt{r_0}(c - \delta)} \leq 1 - 2t\sqrt{d_{\max} - 1} - t^2, \quad (20)$$

$$\frac{\delta}{c - \delta} + \frac{2\delta}{\sqrt{r_0}(c - \delta)} < 1, \quad (21)$$

and

$$\lambda'_0 < \lambda < 1, \quad (22)$$

where $\lambda'_0 \triangleq \max\{\lambda'_1, \lambda'_2\}$ and

$$\begin{aligned} \lambda'_1 &\triangleq \inf\{0 < \lambda < 1: \sqrt{1 - \lambda} + \frac{2\delta}{\sqrt{r_0}(c - \delta)\sqrt{\lambda}} \\ &< 1 - 2t\sqrt{d_{\max} - 1} - t^2 - \delta\}, \end{aligned} \quad (23)$$

$$\lambda'_2 \triangleq \inf\{0 < \lambda < 1: \lambda - \frac{2\delta}{\sqrt{r_0}(c - \delta)} \frac{1}{\sqrt{\lambda}} > \frac{\delta}{c - \delta}\}. \quad (24)$$

Then with probability at least $1 - K \exp(-d) - 8 \sum_{k=1}^K n_k \exp(-\frac{d_k t^2}{2})$, the subspace detection property holds for \mathbf{x}_i with β^* .

Remark 5. Note that there is no assumption on the distribution of subspaces in Theorem 3, so it is not required that the subspaces should have uniform distribution, as is required in the geometric analysis of ℓ^1 -SSC Soltanolkotabi & Cands (2012) and its noisy version Wang & Xu (2013). In addition, while Soltanolkotabi & Cands (2012); Wang & Xu (2013) require data in each subspace are i.i.d according to uniform distribution on unit sphere, our randomized result requires data in each subspace are i.i.d. isotropic random vectors on sphere of radius $\sqrt{d_k}$. Note that i.i.d samples uniformly distributed on sphere of radius $\sqrt{d_k}$ centered at the origin are in fact isotropic, our assumption is less restrictive after scaling the data by a factor of $\sqrt{d_k}$.

4 NOISY ℓ^0 -SSC ON DIMENSIONALITY REDUCED DATA: NOISY-DR- ℓ^0 -SSC

Albeit the theoretical guarantee and compelling empirical performance of noisy ℓ^0 -SSC to be shown in the experimental results, the computational cost of noisy ℓ^0 -SSC is high with the high dimensionality of the data. In this section, we propose Noisy Dimensionality Reduced ℓ^0 -SSC (Noisy-DR- ℓ^0 -SSC) which performs noisy ℓ^0 -SSC on dimensionality reduced data. The theoretical guarantee on the correctness of Noisy-DR- ℓ^0 -SSC under deterministic model as well as its empirical performance are presented.

4.1 METHOD

Noisy-DR- ℓ^0 -SSC performs subspace clustering by the following two steps: 1) obtain the dimension reduced data $\tilde{\mathbf{X}} = \mathbf{P}\mathbf{X}$ with a linear transformation $\mathbf{P} \in \mathbb{R}^{p \times d}$ ($p < d$). 2) perform noisy ℓ^0 -SSC on the compressed data $\tilde{\mathbf{X}}$:

$$\min_{\tilde{\beta} \in \mathbb{R}^n, \tilde{\beta}_i = 0} L(\tilde{\beta}) = \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{X}}\tilde{\beta}\|_2^2 + \tilde{\lambda}\|\tilde{\beta}\|_0. \quad (25)$$

If $p < d$, Noisy-DR- ℓ^0 -SSC operates on the compressed data $\tilde{\mathbf{X}}$ rather than on the original data, so that the efficiency is improved.

4.2 ANALYSIS

High-dimensional data often exhibits low-dimensional structures, which often leads to low-rankness of the data matrix. Intuitively, if the data is low rank, then it could be safe to perform noisy ℓ^0 -SSC on its dimensionality reduced version by the linear projection \mathbf{P} , and it is expected that \mathbf{P} can preserve the information of the subspaces contained in the original data as much as possible, while effectively removing uninformative dimensions.

To this end, we propose to choose \mathbf{P} as a random projection induced by randomized low-rank approximation of the data. The key idea is to obtain an approximate low-rank decomposition of the data. Using the random projection induced by such low-rank approximation as the linear transformation \mathbf{P} , the clustering correctness hold for Noisy-DR- ℓ^0 -SSC with a high probability.

Randomized algorithms are efficient and they have been extensively studied in the computer science and numerical linear algebra literature. They have been employed to accelerate various numerical matrix computation and matrix optimization problems, including random projection for matrix decomposition Frieze et al. (2004); Drineas et al. (2004); Sarlos (2006); Drineas et al. (2006; 2008); Mahoney & Drineas (2009); Drineas et al. (2011); Lu et al. (2013).

Formally, a random matrix $\mathbf{T} \in \mathbb{R}^{n \times p}$ is generated such that each element \mathbf{T}_{ij} is sampled independently according to the Gaussian distribution $\mathcal{N}(0, 1)$. QR decomposition is then performed on $\mathbf{X}\mathbf{T}$ to obtain the basis of its column space, namely $\mathbf{X}\mathbf{T} = \mathbf{Q}\mathbf{R}$ where $\mathbf{Q} \in \mathbb{R}^{d \times p}$ is an orthogonal matrix of rank p and $\mathbf{R} \in \mathbb{R}^{p \times p}$ is an upper triangle matrix. The columns of \mathbf{Q} form the orthogonal basis for the sample matrix $\mathbf{X}\mathbf{T}$. An approximation of \mathbf{X} is then obtained by projecting \mathbf{X} onto the column space of $\mathbf{X}\mathbf{T}$: $\mathbf{Q}\mathbf{Q}^\top \mathbf{X} = \mathbf{Q}\mathbf{W} = \hat{\mathbf{X}}$ where $\mathbf{W} = \mathbf{Q}^\top \mathbf{X} \in \mathbb{R}^{p \times n}$. In this manner, a randomized low-rank decomposition of $\hat{\mathbf{X}}$ is achieved as follows:

$$\hat{\mathbf{X}} = \mathbf{Q}\mathbf{W} \quad (26)$$

We present probabilistic result on the correctness of Noisy-DR- ℓ^0 -SSC using the random projection induced by randomized low-rank decomposition of the data \mathbf{X} , namely $\mathbf{P} = \mathbf{Q}^\top$, in Theorem 4. In the sequel, $\tilde{\mathbf{x}} = \mathbf{P}\mathbf{x}$ for any $\mathbf{x} \in \mathbb{R}^n$. To guarantee the subspace detection property on the dimensionality-reduced data $\hat{\mathbf{X}}$, it is crucial to make sure that the conditions, such as (12) and (13) in Theorem 2, still hold after the linear transformation.

We denote by $\tilde{\beta}^*$ the optimal solution to (25). We also define the following quantities in the analysis of the subspace detection property, which correspond to M_i , $\bar{\sigma}_{\mathbf{Y},r}$, $\sigma_{\mathbf{X},r}$ and μ_r used in the analysis on the original data:

$$\tilde{M}_i \triangleq \min\{d(\tilde{\mathbf{y}}_i, \mathbf{H}) : \mathbf{H} \in \mathcal{H}_{\tilde{\mathbf{y}}_i, \tilde{d}_k}\}, \quad (27)$$

where $\mathcal{H}_{\tilde{\mathbf{y}}_i, \tilde{d}_k}$ is all the external subspaces of $\tilde{\mathbf{y}}_i$ with dimension no greater than \tilde{d}_k in the transformed space by \mathbf{P} .

$$\bar{\sigma}_{\tilde{\mathbf{Y}},r} \triangleq \min_{\beta: \|\beta\|_0=r, \text{rank}(\tilde{\mathbf{Y}}\beta)=\|\beta\|_0} \sigma_{\min}(\tilde{\mathbf{Y}}\beta), \quad (28)$$

$$\sigma_{\tilde{\mathbf{X}},r} \triangleq \min\{\sigma_{\min}(\tilde{\mathbf{X}}\beta) : 1 \leq \|\beta\|_0 \leq r\}, \quad (29)$$

$$\tilde{\mu}_r \triangleq \frac{\delta}{\min_{1 \leq r' < r} \bar{\sigma}_{\tilde{\mathbf{Y}},r'} - \delta}. \quad (30)$$

Theorem 4. (Subspace detection property holds for Noisy-DR- ℓ^0 -SSC under deterministic model) *Let nonzero vector β^* be the optimal solution to the noisy ℓ^0 -SSC problem (6) for point \mathbf{x}_i with $\|\beta^*\|_0 = r^*$, $n_k \geq d_k + 1$ for every $1 \leq k \leq K$, and there exists $1 < r_0 \leq d$ such that $1 < r^* \leq r_0$. Suppose \mathbf{Y} is in general position, $\delta < \min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r}$, and $\tilde{M}_{i,\delta} \triangleq \tilde{M}_i - \delta$. Suppose the following conditions hold:*

(i)

$$C_{p,p_0} + 2\delta\sqrt{\tilde{d}_{\max}} < \min_{k=1,\dots,K} \sigma_{\mathbf{Y}}^{(k)}, \quad (31)$$

where $\tilde{d}_{\max} \triangleq \max_k \tilde{d}_k$, $\sigma_{\mathbf{Y}}^{(k)} \triangleq \min\{\sigma_{\min}(\mathbf{A}) : \mathbf{A} \subseteq \mathbf{Y}^{(k)}, \mathbf{A} \in \mathbb{R}^{d \times n'}, n' \leq \tilde{d}_k\}$.

$$(ii) \delta(1 + 2\sqrt{r_0}) < \min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r} - C_{p,p_0},$$

$$(iii) \min_{1 \leq r \leq \tilde{d}_k} \sigma_{\mathbf{Y},r} > C_{p,p_0} - 2\delta\sqrt{\tilde{d}_k} \text{ and}$$

$$\begin{aligned} & M_i - C_{p,p_0} \left(1 + \frac{1}{\min_{1 \leq r \leq \tilde{d}_k} \sigma_{\mathbf{Y},r} - C_{p,p_0} - 2\delta\sqrt{\tilde{d}_k}}\right) \\ & > \delta + \frac{2\delta}{\sigma_{\mathbf{X},r_0} - C_{p,p_0}}, \end{aligned} \quad (32)$$

for all $\mathbf{y}_i \in \mathcal{S}_k$ and $1 \leq k \leq K$.

$$(iv) \min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r_0} > C_{p,p_0} - 2\delta\sqrt{r_0} - \delta \text{ and}$$

$$\begin{aligned} & \frac{\delta}{\min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r_0} - C_{p,p_0} - 2\delta\sqrt{r_0} - \delta} \\ & < 1 - \frac{2\delta}{\sigma_{\mathbf{X},r_0} - C_{p,p_0}}. \end{aligned} \quad (33)$$

If

$$\tilde{\lambda}_0 < \tilde{\lambda} < 1, \quad (34)$$

where $\tilde{\lambda}_0 = \max\{\max\{\tilde{\lambda}_1, \tilde{\lambda}_2, \frac{1}{r_0}\}\}$ and

$$\tilde{\lambda}_1 = \inf\{0 < \tilde{\lambda} < 1: \sqrt{1 - \tilde{\lambda}} + \frac{2\delta}{\sigma_{\tilde{\mathbf{X}},r_0} \sqrt{\tilde{\lambda}}} < \tilde{M}_{i,\delta}\}, \quad (35)$$

$$\tilde{\lambda}_2 = \inf\{0 < \tilde{\lambda} < 1: \tilde{\lambda} - \frac{2\delta}{\sigma_{\mathbf{X},r_0} \sqrt{\tilde{\lambda}}} > \tilde{\mu}_{r_0}\}, \quad (36)$$

then with probability at least $1 - 6e^{-p}$, the subspace detection property holds for $\tilde{\mathbf{x}}_i$ with $\tilde{\boldsymbol{\beta}}^*$. Here \tilde{M}_i , $\tilde{\mu}_r$ and $\tilde{\sigma}_{\tilde{\mathbf{X}},r_0}$ are defined in (27), (30) and (29) respectively.

Table 1: Clustering results on various data sets, with the best two results in bold

Data Set	Measure	KM	SC	Noisy SSC	Noisy DR-SSC	SMCE	SSC-OMP	Noisy ℓ^0 -SSC	Noisy-DR- ℓ^0 -SSC
COIL-20	AC	0.6554	0.4278	0.7854	0.7764	0.7549	0.3389	0.8472	0.8479
	NMI	0.7630	0.6217	0.9148	0.9219	0.8754	0.4853	0.9428	0.9433
COIL-100	AC	0.4996	0.2835	0.5275	0.5013	0.5639	0.1667	0.7683	0.7039
	NMI	0.7539	0.5923	0.8041	0.8019	0.8064	0.3757	0.9182	0.8706
Yale-B	AC	0.0954	0.1077	0.7850	0.7255	0.3293	0.7789	0.8480	0.8231
	NMI	0.1258	0.1485	0.7760	0.7311	0.3812	0.7024	0.8612	0.8533

Table 2: Clustering results on various data sets, with the best two results in bold

Data Set	Measure	KM	SC	Noisy SSC	Noisy DR-SSC	SMCE	SSC-OMP	Noisy ℓ^0 -SSC	Noisy-DR- ℓ^0 -SSC
MPIE S1	AC	0.1164	0.1285	0.5892	0.3588	0.1721	0.1695	0.6741	0.6741
	NMI	0.5049	0.5292	0.7653	0.6806	0.5514	0.3395	0.8622	0.8622
MPIE S2	AC	0.1315	0.1410	0.6994	0.4611	0.1898	0.2093	0.7527	0.7527
	NMI	0.4834	0.5128	0.8149	0.7086	0.5293	0.4292	0.8939	0.7527
MPIE S3	AC	0.1291	0.1459	0.6316	0.4841	0.1856	0.1787	0.7050	0.7050
	NMI	0.4811	0.5185	0.7858	0.7340	0.5155	0.3415	0.8750	0.8750
MPIE S4	AC	0.1308	0.1463	0.6803	0.5511	0.1823	0.1680	0.7246	0.7246
	NMI	0.4866	0.5280	0.8063	0.7955	0.5294	0.3345	0.8837	0.8837

5 OPTIMIZATION OF NOISY ℓ^0 -SSC AND NOISY-DR- ℓ^0 -SSC

We employ Proximal Gradient Descent (PGD) to optimize the objective function of noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC. For example, in the k -th iteration of PGD for problem (6), the variable $\boldsymbol{\beta}$ is updated according to

$$\boldsymbol{\beta}^{(k+1)} = T_{\sqrt{2\lambda s}}(\boldsymbol{\beta}^{(k)} - s\nabla g(\boldsymbol{\beta}^{(k)})), \quad (37)$$

where $g(\beta) \triangleq \|\mathbf{x}_i - \mathbf{X}\beta\|_2^2$, T_θ is an element-wise hard thresholding operator:

$$[T_\theta(\mathbf{u})]_j = \begin{cases} 0 & : |\mathbf{u}_j| \leq \theta \\ \mathbf{u}_j & : \text{otherwise} \end{cases}, \quad 1 \leq j \leq n.$$

It is proved in Yang & Yu (2019) that the sequence $\{\beta^{(k)}\}$ generated by PGD converges to a critical point of (6), denoted by $\hat{\beta}$. Let β^* be the optimal solution to (6). Theorem 5 in Yang & Yu (2019) to problem (6) shows that the $\|\beta^* - \hat{\beta}\|_2$ is bounded. Theorem 5 establishes the conditions under which $\hat{\beta}$ is also the optimal solution to (6).

Define $\mathbf{S}^* \triangleq \text{supp}(\beta^*)$, $H^* \triangleq \max_{1 \leq j \leq n} \text{dist}(\mathbf{x}_j, \mathbf{H}_{\mathbf{X}_{\mathbf{S}^* \setminus \{j\}}})$, $\mu \triangleq \max\{H^* + \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2, 2\|\mathbf{x}_i - \mathbf{X}\hat{\beta}\|_2, 2\|\mathbf{x}_i - \mathbf{X}\beta^*\|_2\}$, $\kappa_0 \triangleq \sigma_{\min}(\mathbf{X}_{\mathbf{S} \cup \mathbf{S}^*}) > 0$ where $\mathbf{S} \triangleq \text{supp}(\beta^{(0)})$. The following theorem demonstrates that $\hat{\beta} = \beta^*$ if λ is two-side bounded and $\hat{\beta}_{\min} = \min_{t: \hat{\beta}_t \neq 0} |\hat{\beta}_t|$ is sufficiently large.

Theorem 5. (Conditions that the sub-optimal solution by PGD is also globally optimal) *If*

$$\hat{\beta}_{\min} \geq \frac{\mu}{\kappa_0^2} \quad (38)$$

and

$$\frac{\mu^2}{2\kappa_0^2} \leq \lambda \leq (\hat{\beta}_{\min} - \frac{\mu}{2\kappa_0^2})\mu, \quad (39)$$

then $\hat{\beta} = \beta^*$.

6 EXPERIMENTAL RESULTS

We demonstrate the performance of noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC, with comparison to other competing clustering methods including K-means (KM), Spectral Clustering (SC), noisy SSC, Sparse Manifold Clustering and Embedding (SMCE) Elhamifar & Vidal (2011) and SSC-OMP Dyer et al. (2013). With the coefficient matrix \mathbf{Z} obtained by the optimization of noisy ℓ^0 -SSC or Noisy-DR- ℓ^0 -SSC, a sparse similarity matrix is built by $\mathbf{W} = \frac{|\mathbf{Z}| + |\mathbf{Z}^\top|}{2}$, and spectral clustering is performed on \mathbf{W} to obtain the clustering results. Two measures are used to evaluate the performance of different clustering methods, i.e. the Accuracy (AC) and the Normalized Mutual Information (NMI) Zheng et al. (2004).

We use randomized rank- p decomposition of the data matrix in Noisy-DR- ℓ^0 -SSC with $p = \frac{\min\{d, n\}}{10}$. It can be observed that noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC always achieve better performance than other methods in Table 1, including the noisy SSC on dimensionality reduced data (Noisy DR-SSC) Wang et al. (2015). Throughout all the experiments we find that the best clustering accuracy is achieved whenever λ is chosen by $0.5 < \lambda < 0.95$, justifying our theoretical finding claimed in Remark 4 and (39) in Theorem 5. More experimental results on the CMU Multi-PIE data are shown in Table 2. For all the methods that involve random projection, we conduct the experiments for 30 times and report the average performance. Note that the cluster accuracy of SSC-OMP on the extended Yale-B data set is reported according to You et al. (2016). The time complexity of running PGD for noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC are $\mathcal{O}(Tnd)$ and $\mathcal{O}(Tpd)$ respectively, where T is the maximum iteration number. The actual running time of both algorithms confirms such time complexity, and we observe that Noisy-DR- ℓ^0 -SSC is always more than 8.7 times faster than noisy ℓ^0 -SSC with the same number of iterations.

7 CONCLUSION

We present provable noisy ℓ^0 -SSC that recovers subspaces from noisy data through ℓ^0 -induced sparsity in a robust manner, with the theoretical guarantee on its correctness in terms of subspace detection property under both deterministic and randomized models. Experimental results shows the superior performance of noisy ℓ^0 -SSC. We also propose Noisy-DR- ℓ^0 -SSC which performs noisy ℓ^0 -SSC on dimensionality reduced data and still provably recovers the subspaces in the original data. Experiment results demonstrate the effectiveness of both noisy ℓ^0 -SSC and Noisy-DR- ℓ^0 -SSC.

REFERENCES

- G. Aubrun and S.J. Szarek. *Alice and Bob Meet Banach: The Interface of Asymptotic Geometric Analysis and Quantum Information Theory*. Mathematical Surveys and Monographs. American Mathematical Society, 2017.
- Yan Mei Chen, Xiao Shan Chen, and Wen Li. On perturbation bounds for orthogonal projections. *Numerical Algorithms*, 73(2):433–444, Oct 2016.
- P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Machine Learning*, 56(1):9–33, 2004.
- Petros Drineas, Ravi Kannan, and Michael W. Mahoney. Fast monte carlo algorithms for matrices ii: Computing a low-rank approximation to a matrix. *SIAM Journal on Computing*, 36(1):158–183, 2006.
- Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. Relative-error Scurl matrix decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30(2):844–881, 2008.
- Petros Drineas, Michael W. Mahoney, S. Muthukrishnan, and Tamás Sarlós. Faster least squares approximation. *Numerische Mathematik*, 117(2):219–249, 2011.
- Eva L. Dyer, Aswin C. Sankaranarayanan, and Richard G. Baraniuk. Greedy feature selection for subspace clustering. *Journal of Machine Learning Research*, 14:2487–2517, 2013.
- Ehsan Elhamifar and René Vidal. Sparse manifold clustering and embedding. In *NIPS*, pp. 55–63, 2011.
- Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(11):2765–2781, 2013.
- Alan Frieze, Ravi Kannan, and Santosh Vempala. Fast monte-carlo algorithms for finding low-rank approximations. *J. ACM*, 51(6):1025–1041, November 2004.
- N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, May 2011. ISSN 0036-1445.
- Guangcan Liu, Zhouchen Lin, and Yong Yu. Robust subspace segmentation by low-rank representation. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pp. 663–670, 2010.
- Guangcan Liu, Zhouchen Lin, Shuicheng Yan, Ju Sun, Yong Yu, and Yi Ma. Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(1):171–184, January 2013.
- Yichao Lu, Paramveer S. Dhillon, Dean Foster, and Lyle Ungar. Faster ridge regression via the subsampled randomized hadamard transform. In *Proceedings of the 26th International Conference on Neural Information Processing Systems, NIPS’13*, pp. 369–377, USA, 2013. Curran Associates Inc.
- Michael W. Mahoney and Petros Drineas. Cur matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106(3):697–702, 2009.
- Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, pp. 849–856, 2001.
- Dohyung Park, Constantine Caramanis, and Sujay Sanghavi. Greedy subspace clustering. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pp. 2753–2761, 2014.
- Xi Peng, Shijie Xiao, Jiashi Feng, Wei-Yun Yau, and Zhang Yi. Deep subspace clustering with sparsity prior. In *Proceedings of the 25 International Joint Conference on Artificial Intelligence*, pp. 1925–1931, New York, NY, USA, 9-15 July 2016.
- T. Sarlos. Improved approximation algorithms for large matrices via random projections. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS’06)*, pp. 143–152, Oct 2006.
- Mahdi Soltanolkotabi and Emmanuel J. Cands. A geometric analysis of subspace clustering with outliers. *Ann. Statist.*, 40(4):2195–2238, 08 2012.
- G. W. Stewart. On the perturbation of pseudo-inverses, projections and linear least squares problems. *SIAM Review*, 19(4):634–662, 1977. ISSN 00361445.
- R. Vidal. Subspace clustering. *Signal Processing Magazine, IEEE*, 28(2):52–68, March 2011.

- Yining Wang, Yu-Xiang Wang, and Aarti Singh. A deterministic analysis of noisy sparse subspace clustering for dimensionality-reduced data. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15*, pp. 1422–1431. JMLR.org, 2015.
- Yu-Xiang Wang and Huan Xu. Noisy sparse subspace clustering. In *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*, pp. 89–97, 2013.
- Yu-Xiang Wang, Huan Xu, and Chenlei Leng. Provable subspace clustering: When lrr meets ssc. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 26*, pp. 64–72. Curran Associates, Inc., 2013.
- H. Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71:441–479, 1912.
- Yingzhen Yang and Jiahui Yu. Fast proximal gradient descent for A class of non-convex and non-smooth sparse learning problems. In *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2019, Tel Aviv, Israel, July 22-25, 2019*, pp. 508, 2019.
- Yingzhen Yang, Jiashi Feng, Nebojsa Jojic, Jianchao Yang, and Thomas S. Huang. L0-sparse subspace clustering. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II*, pp. 731–747, 2016.
- Pavel Yaskov. Lower bounds on the smallest eigenvalue of a sample covariance matrix. *Electron. Commun. Probab.*, 19:10 pp., 2014.
- Yu-Xiang Wang Yining Wang and Aarti Singh. Parameter estimation of generalized linear models without assuming their link function. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2016, Cadiz, Spain, May 9-11, 2016*, 2016.
- Chong You, Daniel P. Robinson, and René Vidal. Scalable sparse subspace clustering by orthogonal matching pursuit. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 3918–3927, 2016.
- Xin Zheng, Deng Cai, Xiaofei He, Wei-Ying Ma, and Xueyin Lin. Locality preserving clustering for image database. In *Proceedings of the 12th Annual ACM International Conference on Multimedia, MULTIMEDIA '04*, pp. 885–891, New York, NY, USA, 2004. ACM.

A APPENDIX

We provide proofs to the lemmas and theorems in the paper in this appendix.

A.1 PROOF OF REMARK 1

Lemma A. (Subspace detection property holds for ℓ^0 -SSC under the deterministic model) *Under the deterministic model, suppose data is noiseless, $n_k \geq d_k + 1$, $\mathbf{Y}^{(k)}$ is in general position. If all the data points in $\mathbf{Y}^{(k)}$ are away from the external subspaces for any $1 \leq k \leq K$, then the subspace detection property for ℓ^0 -SSC holds with the optimal solution \mathbf{Z}^* to (1).*

Proof. Let $\mathbf{x}_i \in \mathcal{S}_k$. Note that \mathbf{Z}^{*i} is the optimal solution to the following ℓ^0 sparse representation problem

$$\min_{\mathbf{Z}^i} \|\mathbf{Z}^i\|_0 \quad \text{s.t. } \mathbf{x}_i = [\mathbf{X}^{(k)} \setminus \mathbf{x}_i \quad \mathbf{X}^{(-k)}] \mathbf{Z}^i, \quad \mathbf{Z}_{ii} = 0, \quad (40)$$

where $\mathbf{X}^{(-k)}$ denotes the data that lie in all subspaces except \mathcal{S}_k . Let $\mathbf{Z}^{*i} = \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix}$ where $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are sparse codes corresponding to $\mathbf{X}^{(k)} \setminus \mathbf{x}_i$ and $\mathbf{X}^{(-k)}$ respectively.

Suppose $\boldsymbol{\beta} \neq \mathbf{0}$, then \mathbf{x}_i belongs to a subspace $\mathcal{S}' = \mathbf{H}_{\mathbf{X}_{\mathbf{Z}^{*i}}}$ spanned by the projected data points corresponding to nonzero elements of \mathbf{Z}^{*i} , and $\mathcal{S}' \neq \mathcal{S}_k$, $\dim[\mathcal{S}'] \leq d_k$. To see this, if $\mathcal{S}' = \mathcal{S}_k$, then the data corresponding to nonzero elements of $\boldsymbol{\beta}$ belong to \mathcal{S}_k , which is contrary to the definition of $\mathbf{X}^{(-k)}$. Also, if $\dim[\mathcal{S}'] > d_k$, then any d_k points in $\mathbf{X}^{(k)}$ can be used to linearly represent \mathbf{x}_i by the condition of general position, contradicting with the optimality of \mathbf{Z}^{*i} . Since the data points (or columns) in $\mathbf{X}_{\mathbf{Z}^{*i}}$ are linearly independent, it follows that \mathbf{x}_i lies in an external subspace $\mathbf{H}_{\mathbf{X}_{\mathbf{Z}^{*i}}}$ spanned by linearly independent points in $\mathbf{X}_{\mathbf{Z}^{*i}}$, and $\dim[\mathbf{H}_{\mathbf{X}_{\mathbf{Z}^{*i}}}] = \dim[\mathcal{S}'] \leq d_k$. This contradicts with the assumption that \mathbf{x}_i is away from the external subspaces. Therefore,

$\beta = \mathbf{0}$. Perform the above analysis for all $1 \leq i \leq n$, we can prove that the subspace detection property holds for all $1 \leq i \leq n$. \square

A.2 PROOF OF LEMMA 1

The following proposition is used for proving Lemma 1.

Lemma B. (Perturbation of distance to subspaces) *Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ are two matrices and $\text{rank}(\mathbf{A}) = r$, $\text{rank}(\mathbf{B}) = s$. Also, $\mathbf{E} = \mathbf{A} - \mathbf{B}$ and $\|\mathbf{E}\|_2 \leq C$, where $\|\cdot\|_2$ indicates the spectral norm. Then for any point $\mathbf{x} \in \mathbb{R}^m$, the difference of the distance of \mathbf{x} to the column space of \mathbf{A} and \mathbf{B} , i.e. $|d(\mathbf{x}, \mathbf{H}_{\mathbf{A}}) - d(\mathbf{x}, \mathbf{H}_{\mathbf{B}})|$, is bounded by*

$$|d(\mathbf{x}, \mathbf{H}_{\mathbf{A}}) - d(\mathbf{x}, \mathbf{H}_{\mathbf{B}})| \leq \frac{C\|\mathbf{x}\|_2}{\min\{\sigma_r(\mathbf{A}), \sigma_s(\mathbf{B})\}}. \quad (41)$$

Proof. Note that the projection of \mathbf{x} onto the subspace $\mathbf{H}_{\mathbf{A}}$ is $\mathbf{A}\mathbf{A}^+\mathbf{x}$ where \mathbf{A}^+ is the Moore-Penrose pseudo-inverse of the matrix \mathbf{A} , so $d(\mathbf{x}, \mathbf{H}_{\mathbf{A}})$ equals to the distance between \mathbf{x} and its projection, namely $d(\mathbf{x}, \mathbf{H}_{\mathbf{A}}) = \|\mathbf{x} - \mathbf{A}\mathbf{A}^+\mathbf{x}\|_2$. Similarly, $d(\mathbf{x}, \mathbf{H}_{\mathbf{B}}) = \|\mathbf{x} - \mathbf{B}\mathbf{B}^+\mathbf{x}\|_2$.

It follows that

$$\begin{aligned} |d(\mathbf{x}, \mathbf{H}_{\mathbf{A}}) - d(\mathbf{x}, \mathbf{H}_{\mathbf{B}})| &= \|\mathbf{x} - \mathbf{A}\mathbf{A}^+\mathbf{x}\|_2 - \|\mathbf{x} - \mathbf{B}\mathbf{B}^+\mathbf{x}\|_2 \\ &\leq \|\mathbf{A}\mathbf{A}^+\mathbf{x} - \mathbf{B}\mathbf{B}^+\mathbf{x}\|_2 \leq \|\mathbf{A}\mathbf{A}^+ - \mathbf{B}\mathbf{B}^+\|_2 \|\mathbf{x}\|_2. \end{aligned} \quad (42)$$

According to the perturbation bound on the orthogonal projection in Chen et al. (2016); Stewart (1977),

$$\|\mathbf{A}\mathbf{A}^+ - \mathbf{B}\mathbf{B}^+\|_2 \leq \max\{\|\mathbf{E}\mathbf{A}^+\|_2, \|\mathbf{E}\mathbf{B}^+\|_2\}. \quad (43)$$

Since $\|\mathbf{E}\mathbf{A}^+\|_2 \leq \|\mathbf{E}\|_2 \|\mathbf{A}^+\|_2 \leq \frac{C}{\sigma_r(\mathbf{A})}$, $\|\mathbf{E}\mathbf{B}^+\|_2 \leq \|\mathbf{E}\|_2 \|\mathbf{B}^+\|_2 \leq \frac{C}{\sigma_s(\mathbf{B})}$, combining (42) and (43), we have

$$\begin{aligned} |d(\mathbf{x}, \mathbf{H}_{\mathbf{A}}) - d(\mathbf{x}, \mathbf{H}_{\mathbf{B}})| &\leq \max\left\{\frac{C}{\sigma_r(\mathbf{A})}, \frac{C}{\sigma_s(\mathbf{B})}\right\} \|\mathbf{x}\|_2 \\ &= \frac{C\|\mathbf{x}\|_2}{\min\{\sigma_r(\mathbf{A}), \sigma_s(\mathbf{B})\}}. \end{aligned} \quad (44)$$

So that (5) is proved. \square

Proof of Lemma 1. We have $\mathbf{y}_i = \mathbf{x}_i - \mathbf{n}_i$, and $\sigma_{\min}(\mathbf{Y}_{\beta}^{\top} \mathbf{Y}_{\beta}) = (\sigma_{\min}(\mathbf{Y}_{\beta}))^2 \geq \sigma_{\mathbf{Y},r}^2$.

By Weyl (1912), $|\sigma_i(\mathbf{Y}_{\beta}) - \sigma_i(\mathbf{X}_{\beta})| \leq \|\mathbf{N}_{\beta}\|_2 \leq \|\mathbf{N}_{\beta}\|_F \leq \sqrt{r}\delta$. Since $\sqrt{r}\delta < \sigma_{\mathbf{Y},r} \leq \sigma_{\min}(\mathbf{Y}_{\beta}) \leq \sigma_i(\mathbf{Y}_{\beta})$, $\sigma_i(\mathbf{X}_{\beta}) \geq \sigma_i(\mathbf{Y}_{\beta}) - \sqrt{r}\delta \geq \sigma_{\mathbf{Y},r} - \sqrt{r}\delta > 0$ for $1 \leq i \leq \min\{d, r\}$. It follows that $\sigma_{\min}(\mathbf{X}_{\beta}) \geq \sigma_{\mathbf{Y},r} - \sqrt{r}\delta > 0$ and \mathbf{X}_{β} has full column rank.

Also, $\|\mathbf{X}_{\beta} - \mathbf{Y}_{\beta}\|_2 \leq \|\mathbf{X}_{\beta} - \mathbf{Y}_{\beta}\|_F \leq \sqrt{r}\delta$. According to Lemma B,

$$\begin{aligned} |d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta}}) - d(\mathbf{x}_i, \mathbf{H}_{\mathbf{Y}_{\beta}})| &\leq \frac{\sqrt{r}\delta}{\min\{\sigma_{\min}(\mathbf{X}_{\beta}), \sigma_{\min}(\mathbf{Y}_{\beta})\}} \\ &\leq \frac{\sqrt{r}\delta}{\sigma_{\mathbf{Y},r} - \sqrt{r}\delta} = \frac{\delta}{\bar{\sigma}_{\mathbf{Y},r} - \delta}. \end{aligned} \quad (45)$$

\square

A.3 PROOF OF LEMMA 2

Proof.

$$\begin{aligned} \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2^2 + \lambda\|\beta^*\|_0 &\leq \|\mathbf{x}_i - \mathbf{X}\mathbf{0}\|_2^2 + \lambda\|\mathbf{0}\|_0 = 1 \\ \Rightarrow c^* = \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2 &< 1. \end{aligned}$$

We first prove that β^* is the optimal solution to the sparse approximation problem

$$\min_{\beta} \|\beta\|_0 \quad s.t. \quad \|\mathbf{x}_i - \mathbf{X}\beta\|_2 \leq c^*, \beta_i = 0. \quad (46)$$

To see this, suppose there is a vector β' such that $\|\mathbf{x}_i - \mathbf{X}\beta'\|_2 \leq c^*$ and $\|\beta'\|_0 < \|\beta^*\|_0$, then $L(\beta') < c^* + \lambda\|\beta^*\|_0 = L(\beta^*)$, contradicting the fact that β^* is the optimal solution to (6).

Note that \mathbf{X}_{β^*} is a full column rank matrix, otherwise a sparser solution to (6) can be obtained as vector whose support corresponds to the maximal linear independent set of columns of \mathbf{X}_{β^*} .

Also, the distance between \mathbf{x}_i and the subspace spanned by columns of \mathbf{X}_{β^*} equals to c^* , i.e. $d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) = c^*$.

To see this, it is clear that $d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) \leq c^*$. If there is a vector $\mathbf{y} = \mathbf{X}\tilde{\beta}$ in $\mathbf{H}_{\mathbf{X}_{\beta^*}}$ with $\text{supp}(\tilde{\beta}) \subseteq \text{supp}(\beta^*)$, and $\|\mathbf{x}_i - \mathbf{y}\|_2 < c^*$, then $L(\tilde{\beta}) < L(\beta^*)$ which contradicts the optimality of β^* . Therefore, $d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) \geq c^*$, and it follows that $d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) = c^*$.

To prove that the subspace separation margin $\text{HS}(\mathbf{x}_i, \mathbf{X}, \beta^*) > 0$, suppose $\text{HS}(\mathbf{x}_i, \mathbf{X}, \beta^*) \leq 0$, so there exists β' such that $\|\beta'\|_0 < r^*$, $\text{rank}(\mathbf{X}_{\beta'}) = \|\beta'\|_0$ and $d(\mathbf{y}_i, \mathbf{H}_{\mathbf{X}_{\beta'}}) \leq d(\mathbf{y}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) \leq c^*$. Then β' is sparser than β^* and it satisfies the constraint of problem (46), contradicting the optimality of β^* .

Since $\|\mathbf{x}_i - \mathbf{X}\beta^*\|_2 \leq 1$, $\|\mathbf{X}\beta^*\|_2 \leq 2$. Also,

$$\sigma_{\min}(\mathbf{X}_{\beta^*}^\top \mathbf{X}_{\beta^*}) \|\beta^*\|_2^2 \leq \|\mathbf{X}\beta^*\|_2^2 \leq 4,$$

it follows that $\|\beta^*\|_2^2 \leq \frac{4}{\sigma_{\mathbf{X}}^{*2}}$. By Cauchy-Schwarz inequality, $\|\beta^*\|_1 \leq \frac{2\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}$ and $\|\mathbf{N}\beta^*\|_2 \leq \|\beta^*\|_1 \delta \leq \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}$. Therefore,

$$\begin{aligned} \|\mathbf{x}_i - \mathbf{Y}\beta^*\|_2 &= \|\mathbf{x}_i - \mathbf{X}\beta^* + \mathbf{N}\beta^*\|_2 \\ &\leq \|\mathbf{x}_i - \mathbf{X}\beta^*\|_2 + \|\mathbf{N}\beta^*\|_2 \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}, \end{aligned}$$

so that β^* is a feasible for problem (7). To prove that β^* is also the optimal solution to (7), suppose this is not the case, and the optimal solution to (7) is a vector β' such that $\|\mathbf{x}_i - \mathbf{Y}\beta'\|_2 \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}$ and $\|\beta'\|_0 = r < r^*$. $\mathbf{Y}_{\beta'}$ is a full column rank matrix, otherwise a sparser solution can be obtained as vector whose support corresponds to the maximal linear independent set of columns of $\mathbf{Y}_{\beta'}$. We have

$$d(\mathbf{x}_i, \mathbf{H}_{\mathbf{Y}_{\beta'}}) \leq \|\mathbf{x}_i - \mathbf{Y}\beta'\|_2 \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}.$$

According to Lemma 1, we have

$$\begin{aligned} |d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta'}}) - d(\mathbf{x}_i, \mathbf{H}_{\mathbf{Y}_{\beta'}})| &\leq \frac{\sqrt{r}\delta}{\sigma_{\mathbf{Y},r} - \sqrt{r}\delta} \\ &= \frac{\delta}{\bar{\sigma}_{\mathbf{Y},r} - \delta} \leq \frac{\delta}{\bar{\sigma}_{\mathbf{Y}} - \delta} \\ \Rightarrow d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta'}}) &\leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*} + \frac{\delta}{\bar{\sigma}_{\mathbf{Y}} - \delta} = c^* + \tau_0. \end{aligned}$$

However, according to the optimality of β^* in the noisy ℓ^0 -SSC problem (6), we have

$$\begin{aligned} d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta'}}) - c^* &= d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta'}}) - d(\mathbf{x}_i, \mathbf{H}_{\mathbf{X}_{\beta^*}}) \\ &\geq (r^* - r)\lambda \geq \lambda > \tau_0 \end{aligned}$$

This contradiction shows that β^* is the optimal solution to (7). \square

A.4 PROOF OF LEMMA 3

Proof. (8) is equivalent to the following problem

$$\min_{\beta} \|\beta\|_0 \quad s.t. \quad \mathbf{y} = \mathbf{Y}\beta, \|\mathbf{x}_i - \mathbf{y}\|_2 \leq c_0, \beta_i = 0. \quad (47)$$

We show that the points (columns) of \mathbf{Y}_{β^*} must come from subspace \mathcal{S}_k . To see this, suppose some columns of \mathbf{Y}_{β^*} come from different subspaces. We first have $\|\beta^*\|_0 \leq d_k$. To see this, we can choose some $\mathbf{y}' \in \mathcal{S}_k$ such

that $\|\mathbf{y}' - \mathbf{x}_i\|_2 \leq c_0$ since $c_0 \geq d(\mathbf{x}_i, \mathcal{S}_k)$. Also, d_k points in $\mathbf{Y}^{(k)}$ can linearly represent \mathbf{y}' since $\mathbf{Y}^{(k)}$ is in general position, and it follows that $\|\beta^*\|_0 \leq d_k$ due to the optimality of β^* .

Also, \mathbf{Y}_{β^*} has full column rank, so that subspace $\mathbf{H}_{\mathbf{Y}_{\beta^*}} \in \mathcal{H}_{\mathbf{y}_i, d_k}$. Let $\mathbf{y}^* = \mathbf{Y}\beta^*$, then $\mathbf{y}^* \in \mathbf{H}_{\mathbf{Y}_{\beta^*}} \cap \mathbf{B}(\mathbf{x}_i, c_0)$ which contradicts the fact that $\mathbf{B}(\mathbf{x}_i, c_0) \cap \mathbf{H} = \emptyset$ for any $\mathbf{H} \in \mathcal{H}_{\mathbf{y}_i, d_k}$. Therefore, columns of \mathbf{Y}_{β^*} must come from \mathcal{S}_k . \square

A.5 PROOF OF THEOREM 1

Proof. We first show that $d(\mathbf{x}_i, \mathcal{S}_k) \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}$. To see this, $\sigma_{\mathbf{X}}^* = \sigma_{\min}(\mathbf{X}_{\beta^*}) \leq 1$ as the columns of \mathbf{X} have unit ℓ^2 -norm. It follows that

$$c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*} \geq 2\delta\sqrt{r^*} \leq 2\delta > \|\mathbf{x}_i - \mathbf{y}_i\| \leq d(\mathbf{x}_i, \mathcal{S}_k) \quad (48)$$

By Lemma 2, it can be verified that β^* is the optimal solution to the following problem

$$\min_{\beta} \|\beta\|_0 \quad s.t. \quad \|\mathbf{x}_i - \mathbf{Y}\beta\|_2 \leq c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}, \quad \beta_i = 0. \quad (49)$$

The subspace detection property holds which follows from applying Lemma 3 with $c_0 = c^* + \frac{2\delta\sqrt{r^*}}{\sigma_{\mathbf{X}}^*}$. \square

A.6 PROOF OF THEOREM 2

Proof. This theorem can be proved by checking that the conditions in Theorem 1 are satisfied. \square

A.7 PROOF OF LEMMA 4

Proof. Let \mathbf{H} be a fixed subspace of dimension $d_e \leq d_k$, and $\mathbf{y} \notin \mathbf{H}$. Since $\mathbf{y} \in \mathcal{S}_k$ and $\mathbf{y} \notin \mathbf{H}$. Let $\mathbf{U}_{\mathcal{S}_k} = \begin{bmatrix} \mathbf{I}_{d_k} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{d \times d_k}$ be the orthonormal basis of \mathcal{S}_k under which the isotropic random vector \mathbf{y} in \mathcal{S}_k satisfies $\mathbb{E}[\mathbf{y}\mathbf{y}^\top] = \begin{bmatrix} \mathbf{I}_{d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$. It follows that more columns vectors can be added to $\mathbf{U}_{\mathcal{S}_k}$ to form a orthonormal basis $\mathbf{U} \in \mathbb{R}^{d \times d'}$ for the minimum subspace that contains \mathcal{S}_k and \mathbf{H} . It can be verified that $d_k + 1 \leq d' \leq \min\{d_k + d_e, d\}$ because $\mathbf{H} \neq \mathcal{S}_k$. Note that \mathbf{U} can be represented as a block matrix as $\mathbf{U} = \begin{bmatrix} \mathbf{I}_{d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}' \end{bmatrix}$ where $\mathbf{U}' \in \mathbb{R}^{(d-d_k) \times (d'-d_k)}$ has orthonormal columns. It can be verified that the basis of \mathbf{H} can be represented as $\mathbf{U}_{\mathbf{H}} = \begin{bmatrix} \mathbf{I}_{d_e-d'+d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}' \end{bmatrix}$. Note that if $d_e - d' + d_k = 0$, $\mathbf{U}_{\mathbf{H}} = \begin{bmatrix} \mathbf{0} \\ \mathbf{U}' \end{bmatrix}$. Then $\mathbb{P}_{\mathbf{H}}(\mathbf{y}) = \mathbf{U}_{\mathbf{H}}\mathbf{U}_{\mathbf{H}}^\top\mathbf{y}$, and we have

$$\begin{aligned} \mathbb{E}[\|\mathbb{P}_{\mathbf{H}}(\mathbf{y})\|_2^2] &= \mathbb{E}[\mathbf{y}^\top \mathbf{U}_{\mathbf{H}} \mathbf{U}_{\mathbf{H}}^\top \mathbf{U}_{\mathbf{H}} \mathbf{U}_{\mathbf{H}}^\top \mathbf{y}] \\ &= \mathbb{E}[\text{Tr}(\mathbf{y}^\top \mathbf{U}_{\mathbf{H}} \mathbf{U}_{\mathbf{H}}^\top \mathbf{y})] \\ &= \mathbb{E}[\text{Tr}(\mathbf{U}_{\mathbf{H}}^\top \mathbf{y} \mathbf{y}^\top \mathbf{U}_{\mathbf{H}})] \\ &= \text{Tr}(\mathbf{U}_{\mathbf{H}}^\top \mathbb{E}[\mathbf{y} \mathbf{y}^\top] \mathbf{U}_{\mathbf{H}}) \\ &= \text{Tr} \left(\begin{bmatrix} \mathbf{I}_{d_e-d'+d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}' \end{bmatrix}^\top \begin{bmatrix} \mathbf{I}_{d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{d_e-d'+d_k} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}' \end{bmatrix} \right) \\ &= d_e - d' + d_k \leq d_e - 1 \leq d_k - 1 \end{aligned} \quad (50)$$

According to the concentration inequality in section 5.2 of Aubrun & Szarek (2017), for any $t > 0$,

$$\Pr[\|\mathbb{P}_{\mathbf{H}}(\mathbf{y})\|_2 - \sqrt{d_e - d' + d_k} \geq t] \leq 8 \exp\left(-\frac{d_k t^2}{2}\right) \quad (51)$$

Now let \mathbf{H} be spanned by data from \mathbf{Y} , i.e. $\mathbf{H} = \mathbf{H}_{\{\mathbf{y}_{i_j}\}_{j=1}^{d_e}}$, where $\{\mathbf{y}_{i_j}\}_{j=1}^{d_e}$ are any d_e linearly independent points that does not contain \mathbf{y} . For any fixed points $\{\mathbf{y}_{i_j}\}_{j=1}^{d_e}$, (51) holds. Let A be the event that $\|\mathbb{P}_{\mathbf{H}}(\mathbf{y}) - \sqrt{d_e - d' + d_k}\| \geq t$, we aim to integrate the indicator function $\mathbf{1}_A$ with respect to the random vectors, i.e. \mathbf{y} and $\{\mathbf{y}_{i_j}\}_{j=1}^{d_e}$, to obtain the probability that A happens over these random vectors. Let $\mathbf{y} = \mathbf{y}_i$, using Fubini theorem, we have

$$\begin{aligned}
\Pr[A] &= \int_{\times_{j=1}^n \mathcal{S}^{(j)}} \mathbf{1}_A \otimes_{j=1}^n d\mu^{(j)} \\
&= \int_{\times_{j \neq i} \mathcal{S}^{(j)}} \Pr[A|\{\mathbf{y}_j\}_{j \neq i}] \otimes_{j \neq i} d\mu^{(j)} \\
&\leq \int_{\times_{j \neq i} \mathcal{S}^{(j)}} 8 \exp\left(-\frac{d_k t^2}{2}\right) \otimes_{j \neq i} d\mu^{(j)} = 8 \exp\left(-\frac{d_k t^2}{2}\right)
\end{aligned} \tag{52}$$

where $\mathcal{S}^{(j)} \in \{\mathcal{S}_k\}_{k=1}^K$ is the subspace that \mathbf{y}_j lies in, and $\mu^{(j)}$ is the probabilistic measure of the distribution in $\mathcal{S}^{(j)}$. The last inequality is due to (51).

Note that for any \mathbf{y} 's external subspace $\mathbf{H} = \mathbf{H}_{\{\mathbf{y}_{i_j}\}_{j=1}^{d_e}}^{d_e}$, $d(\mathbf{y}, \mathbf{H}) = \sqrt{\|\mathbf{y}\|_2^2 - \|\mathbb{P}_{\mathbf{H}}(\mathbf{y})\|_2^2} = \sqrt{d_k - \|\mathbb{P}_{\mathbf{H}}(\mathbf{y})\|_2^2}$. According to (52), we have

$$\Pr[d(\mathbf{y}, \mathbf{H}) \geq 1 - 2t\sqrt{d_k - 1} - t^2] \geq 1 - 8 \exp\left(-\frac{d_k t^2}{2}\right). \tag{53}$$

□

A.8 PROOF OF THEOREM 3

Proof. According to Yaskov (2014) and condition (a), with probability at least $1 - \exp(-d)$, $\sigma_{\min}(\mathbf{Y}_\beta \geq \sqrt{196Md+1} - 14\sqrt{Md}$ for any $\beta \in \mathbb{R}^n$ such that $\|\beta\|_0 = r \leq d$, $\text{rank}(\mathbf{Y}_\beta) = \|\beta\|_0$. It follows that $\sigma_{\mathbf{Y}, r} \geq \sqrt{196Md+1} - 14\sqrt{Md}$. By Weyl Weyl (1912), $|\sigma_{\min}(\mathbf{X}_\beta) - \sigma_{\min}(\mathbf{Y}_\beta)| \leq \|\mathbf{N}_\beta\|_2 \leq \delta\sqrt{r_0}$. Therefore, $\sigma_{\min}(\mathbf{X}_\beta) \geq \sqrt{196Md+1} - 14\sqrt{Md} - \delta\sqrt{r_0} > 0$ if $\delta < \frac{\sqrt{196Md+1} - 14\sqrt{Md}}{\sqrt{r_0}} = c$. It can be verified that (20), (21) and (22) guarantee (12), (13) and (14) in Theorem 2 respectively, therefore, the conclusion holds. □

A.9 PROOF OF THEOREM 4

It is proved that the low rank approximation $\tilde{\mathbf{X}}$ is close to \mathbf{X} in terms of the spectral norm Halko et al. (2011):

Lemma C. (Corollary 10.9 in Halko et al. (2011)) *Let $p_0 \geq 2$ be an integer and $p' = p - p_0 \geq 4$, then with probability at least $1 - 6e^{-p}$, the spectral norm of $\mathbf{X} - \tilde{\mathbf{X}}$ is bounded by*

$$\|\mathbf{X} - \tilde{\mathbf{X}}\|_2 \leq C_{p, p_0} \tag{54}$$

where

$$C_{p, p_0} \triangleq \left(1 + 17\sqrt{1 + \frac{p_0}{p'}}\right) \sigma_{p_0+1} + \frac{8\sqrt{p}}{p'+1} \left(\sum_{j>p_0} \sigma_j^2\right)^{\frac{1}{2}} \tag{55}$$

and $\sigma_1 \geq \sigma_2 \geq \dots$ are the singular values of \mathbf{X} .

Before proving Theorem 4, we present the following lemma on the perturbation bound for the distance between a data point and a subspace before and after the projection \mathbf{P} . Each subspace \mathcal{S}_k is transformed into $\tilde{\mathcal{S}}_k = \mathbf{P}(\mathcal{S}_k)$ with dimension \tilde{d}_k .

Lemma D. *Let $\beta \in \mathbb{R}^n$, $\tilde{\mathbf{y}}_i = \mathbf{P}\mathbf{y}_i$, $\mathbf{H}_{\mathbf{Y}_\beta}$ is an external subspace of \mathbf{y}_i , $\tilde{\mathbf{Y}}_\beta = \mathbf{P}(\mathbf{Y}_\beta)$ and $\tilde{\mathbf{Y}}_\beta$ has full column rank. Then*

$$\begin{aligned}
&|d(\mathbf{y}_i, \mathbf{H}_{\mathbf{Y}_\beta}) - d(\tilde{\mathbf{y}}_i, \mathbf{H}_{\tilde{\mathbf{Y}}_\beta})| \\
&\leq C_{p, p_0} \left(1 + \frac{1}{\min_{1 \leq r \leq \tilde{d}_k} \sigma_{\mathbf{Y}, r} - C_{p, p_0} - 2\delta\sqrt{\tilde{d}_k}}\right)
\end{aligned} \tag{56}$$

for any $1 \leq i \leq n$ and $\mathbf{y}_i \in \mathcal{S}_k$.

Proof. This lemma can be proved by applying Lemma B. □

Proof of Theorem 4. For any matrix $\mathbf{A} \in \mathbb{R}^{p \times q}$, we first show that multiplying \mathbf{Q} to the left of \mathbf{A} would not change its spectrum. To see this, let the singular value decomposition of \mathbf{A} be $\mathbf{A} = \mathbf{U}_A \Sigma \mathbf{V}_A^\top$ where \mathbf{U}_A

and \mathbf{V}_A have orthonormal columns with $\mathbf{U}_A^\top \mathbf{U}_A = \mathbf{V}_A^\top \mathbf{V}_A = \mathbf{I}$. Then $\mathbf{Q}_A = \mathbf{U}_{Q_A} \Sigma \mathbf{V}_{Q_A}$ is the singular value decomposition of \mathbf{Q}_A with $\mathbf{U}_{Q_A} = \mathbf{Q} \mathbf{U}_A$ and $\mathbf{V}_{Q_A} = \mathbf{V}_A$. This is because the columns of \mathbf{U}_{Q_A} are orthonormal since the columns of \mathbf{Q} are orthonormal: $\mathbf{U}_{Q_A}^\top \mathbf{U}_{Q_A} = \mathbf{U}_A^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{U}_A = \mathbf{I}$, and Σ is a diagonal matrix with nonnegative diagonal elements. It follows that $\sigma_{\min}(\mathbf{Q}_A) = \sigma_{\min}(\mathbf{A})$ for any $\mathbf{A} \in \mathbb{R}^{p \times q}$.

For a point $\mathbf{x}_i = \mathbf{y}_i + \mathbf{n}_i$, after projection via \mathbf{P} , we have the projected noise $\tilde{\mathbf{n}}_i = \mathbf{P} \mathbf{n}_i$. Because

$$\|\tilde{\mathbf{n}}_i\|_2 = \|\mathbf{P} \mathbf{n}_i\|_2 = \|\mathbf{Q}^\top \mathbf{n}_i\|_2 \leq \|\mathbf{Q}\|_2 \|\mathbf{n}_i\|_2 \leq \|\mathbf{n}_i\|_2 \leq \delta, \quad (57)$$

the magnitude of the noise in the projected data is also bounded by δ . Also,

$$\|\tilde{\mathbf{x}}_i\|_2 = \|\mathbf{Q}^\top \mathbf{x}_i\|_2 \leq \|\mathbf{x}_i\|_2 \leq 1, \quad (58)$$

Let $\beta \in \mathbb{R}^n$, $\tilde{\mathbf{Y}}_\beta = \mathbf{P} \mathbf{Y}_\beta$ with $\|\beta\|_0 = r$. Then $\sigma_{\min}(\mathbf{Q} \tilde{\mathbf{Y}}_\beta) = \sigma_{\min}(\tilde{\mathbf{Y}}_\beta)$. Since

$$\begin{aligned} |\sigma_{\min}(\tilde{\mathbf{Y}}_\beta) - \sigma_{\min}(\mathbf{Y}_\beta)| &= |\sigma_{\min}(\mathbf{Q} \tilde{\mathbf{Y}}_\beta) - \sigma_{\min}(\mathbf{Y}_\beta)| \\ &\leq \|\mathbf{Q} \tilde{\mathbf{Y}}_\beta - \mathbf{Y}_\beta\|_2 \\ &= \|\mathbf{Q} \mathbf{Q}^\top \mathbf{Y}_\beta - \mathbf{Y}_\beta\|_2 \\ &= \|\mathbf{Q} \mathbf{Q}^\top \mathbf{X}_\beta - \mathbf{X}_\beta + \mathbf{N}_\beta - \mathbf{Q} \mathbf{Q}^\top \mathbf{N}_\beta\|_2 \\ &\leq C_{p,p_0} + \|\mathbf{N}_\beta\|_F + \|\mathbf{Q} \mathbf{Q}^\top \mathbf{N}_\beta\|_F \\ &\leq C_{p,p_0} + 2\delta\sqrt{r} \end{aligned} \quad (59)$$

Therefore, it follows from (59) that if

$$C_{p,p_0} + 2\delta\sqrt{\tilde{d}_{\max}} < \min_{k=1,\dots,K} \sigma_{\tilde{\mathbf{Y}}}^{(k)}, \quad (60)$$

then $\tilde{\mathbf{Y}}$ is also in general position.

In addition, since $\lambda \geq \frac{1}{r_0}$, we have $\lambda \|\tilde{\beta}^*\|_0 \leq L(\mathbf{0}) \leq 1$, and it follows that $\|\tilde{\beta}^*\|_0 \leq \frac{1}{\lambda} \leq r_0$.

Based on (59) we have

$$|\bar{\sigma}_{\tilde{\mathbf{Y}},r} - \bar{\sigma}_{\mathbf{Y},r}| \leq C_{p,p_0} + 2\delta\sqrt{r_0}, \quad (61)$$

it follows that $\delta < \min_{1 \leq r < r_0} \bar{\sigma}_{\tilde{\mathbf{Y}},r}$ because $\delta < \min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r} - C_{p,p_0} - 2\delta\sqrt{r_0}$.

Again, for $\beta \in \mathbb{R}^n$ with $\|\beta\|_0 = r \leq r_0$, we have

$$\begin{aligned} |\sigma_{\min}(\tilde{\mathbf{X}}_\beta) - \sigma_{\min}(\mathbf{X}_\beta)| &= |\sigma_{\min}(\mathbf{Q} \tilde{\mathbf{X}}_\beta) - \sigma_{\min}(\mathbf{X}_\beta)| \\ &\leq \|\mathbf{Q} \tilde{\mathbf{X}}_\beta - \mathbf{X}_\beta\|_2 \\ &= \|\mathbf{Q} \mathbf{Q}^\top \mathbf{X}_\beta - \mathbf{X}_\beta\|_2 = \|\hat{\mathbf{X}} - \mathbf{X}_\beta\|_2 \\ &\leq C_{p,p_0} \end{aligned} \quad (62)$$

It can be verified that

$$|\sigma_{\tilde{\mathbf{X}},r} - \sigma_{\mathbf{X},r}| \leq C_{p,p_0} \quad (63)$$

Combining (63) and Lemma D, noting that $\sigma_{\mathbf{X},r_0} - C_{p,p_0}$, since

$$\begin{aligned} M_i - C_{p,p_0} \left(1 + \frac{1}{\min_{1 \leq r \leq \tilde{d}_k} \sigma_{\mathbf{Y},r} - C_{p,p_0} - 2\delta\sqrt{\tilde{d}_k}}\right) \\ > \delta + \frac{2\delta}{\sigma_{\mathbf{X},r_0} - C_{p,p_0}}, \end{aligned} \quad (64)$$

we have

$$\tilde{M}_{i,\delta} \triangleq \tilde{M}_i - \delta > \frac{2\delta}{\bar{\sigma}_{\tilde{\mathbf{X}},r_0}}, \quad (65)$$

where $\mathbf{y}_i \in \mathcal{S}_k$.

Based on (61) and (63), we have

$$\tilde{\mu}_{r_0} < 1 - \frac{2\delta}{\sigma_{\mathbf{X},r_0}}, \quad (66)$$

because

$$\frac{\delta}{\min_{1 \leq r < r_0} \bar{\sigma}_{\mathbf{Y},r} - C_{p,p_0} - 2\delta\sqrt{r_0} - \delta} < 1 - \frac{2\delta}{\sigma_{\mathbf{X},r_0} - C_{p,p_0}} \quad (67)$$

□

A.10 SKETCH OF PROOF OF THEOREM 5

We first present Theorem 5 in Yang & Yu (2019). Let $g(\mathbf{x}) = \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2$, $\mathbf{y} \in \mathbb{R}^d$, \mathbf{D} is the design matrix of dimension $d \times n$. Let \mathbf{x}^* be the globally optimal solution to

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0, \quad (68)$$

$\mathbf{S}^* = \text{supp}(\mathbf{x}^*)$, $\hat{\mathbf{x}}$ be the suboptimal solution to (68) obtained by Proximal Gradient Descent (PGD), $\hat{\mathbf{S}} = \text{supp}(\hat{\mathbf{x}})$. The following theorem presents the bound between $\hat{\mathbf{x}}$ and \mathbf{x}^* .

Theorem A. (Theorem 5 in Yang & Yu (2019)) *Suppose $\mathbf{D}_{\mathbf{S} \cup \mathbf{S}^*}$ has full column rank with $\kappa_0 \triangleq \sigma_{\min}(\mathbf{D}_{\mathbf{S} \cup \mathbf{S}^*}) > 0$ where \mathbf{S} is the support of the initialization for PGD on problem (68). Let $\kappa > 0$ such that $2\kappa_0^2 > \kappa$ and b is chosen according to (69) as below:*

$$0 < b < \min\left\{\min_{j \in \hat{\mathbf{S}}} |\hat{\mathbf{x}}_j|, \frac{\lambda}{\max_{j \notin \hat{\mathbf{S}}} \left| \frac{\partial g}{\partial \mathbf{x}_j} \Big|_{\mathbf{x}=\hat{\mathbf{x}}} \right|}, \min_{j \in \mathbf{S}^*} |\mathbf{x}_j^*|, \frac{\lambda}{\max_{j \notin \mathbf{S}^*} \left| \frac{\partial g}{\partial \mathbf{x}_j} \Big|_{\mathbf{x}=\mathbf{x}^*} \right|}\right\}. \quad (69)$$

Let $\mathbf{F} = (\hat{\mathbf{S}} \setminus \mathbf{S}^*) \cup (\mathbf{S}^* \setminus \hat{\mathbf{S}})$ be the symmetric difference between $\hat{\mathbf{S}}$ and \mathbf{S}^* , then

$$\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 \leq \frac{1}{2\kappa_0^2 - \kappa} \left(\sum_{j \in \mathbf{F} \cap \hat{\mathbf{S}}} (\max\{0, \frac{\lambda}{b} - \kappa |\hat{\mathbf{x}}_j - b|\})^2 + \sum_{j \in \mathbf{F} \setminus \hat{\mathbf{S}}} (\max\{0, \frac{\lambda}{b} - \kappa b\})^2 \right)^{\frac{1}{2}} \quad (70)$$

Sketch of Proof of Theorem 5. It can be verified that $\max\{0, \frac{\lambda}{b} - \kappa |\hat{\beta}_j - b|\} = 0$ and $\max\{0, \frac{\lambda}{b} - \kappa b\} = 0$ under the conditions (38) and (39), therefore, $\hat{\beta} = \beta^*$ by applying Theorem A. \square