



## INFORMS Journal on Data Science

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### Compressed Smooth Sparse Decomposition

Shancong Mou, Jianjun Shi

To cite this article:

Shancong Mou, Jianjun Shi (2023) Compressed Smooth Sparse Decomposition. INFORMS Journal on Data Science 2(1):60–80. <https://doi.org/10.1287/ijds.2022.0023>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2022, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

# Compressed Smooth Sparse Decomposition

Shancong Mou,<sup>a</sup> Jianjun Shi<sup>a,\*</sup>

<sup>a</sup>H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, Georgia 30332

\*Corresponding author

Contact: [shancong.mou@gatech.edu](mailto:shancong.mou@gatech.edu) (SM); [jshi33@isye.gatech.edu](mailto:jshi33@isye.gatech.edu),  <https://orcid.org/0000-0002-3774-9176> (JS)

Received: January 11, 2022

Revised: June 20, 2022

Accepted: September 2, 2022

Published Online in Articles in Advance:  
November 7, 2022

<https://doi.org/10.1287/ijds.2022.0023>

Copyright: © 2022 INFORMS

**Abstract.** Image-based anomaly detection systems are of vital importance in various manufacturing applications. The resolution and acquisition rate of such systems are increasing significantly in recent years under the fast development of image sensing technology. This enables the detection of tiny anomalies in real time. However, such a high resolution and a high acquisition rate of image data not only slow down the speed of image processing algorithms but also, increase data storage and transmission cost. To tackle this problem, we propose a fast and data-efficient method with theoretical performance guarantee that is suitable for sparse anomaly detection in images with a smooth background (smooth plus sparse signal). The proposed method, named compressed smooth sparse decomposition (CSSD), is a one-step method that unifies the compressive image acquisition- and decomposition-based image processing techniques. To further enhance its performance in a high-dimensional scenario, a Kronecker compressed smooth sparse decomposition (KronCSSD) method is proposed. Compared with traditional smooth and sparse decomposition algorithms, significant transmission cost reduction and computational speed boost can be achieved with negligible performance loss. Simulation examples and several case studies in various applications illustrate the effectiveness of the proposed framework.

**History:** Kwok-Leung Tsui served as the senior editor for this article.

**Funding:** This work is partially support by the National Science Foundation Division of Engineering Education and Centers [Grant 2052714].

**Data Ethics & Reproducibility Note:** The code capsule is available on Code Ocean at <https://doi.org/10.24433/CO.6352310.v2> and in the e-Companion to this article (available at <https://doi.org/10.1287/ijds.2022.0023>).

**Keywords:** anomaly detection • compressive sensing • image processing • smooth sparse decomposition

## 1. Introduction

High-quality image sensing systems are widely used in manufacturing processes for product quality monitoring and fault diagnosis. The resolution and acquisition rate of such systems increase significantly benefiting from the rapid development of image sensing technology. For example, in a hot rolling process, an in situ image-based sensor can detect a micrometer-sized seam on a rolling bar at a speed of up to 225 miles per hour (Yan et al. 2018). For another example, to monitor solar activity, satellites can capture high-resolution solar images with a high acquisition rate, producing terabytes of data per day (Wang et al. 2018). To achieve real-time inspection, a large volume of high-resolution images needs to be transmitted and processed in real time. Such a large volume of high-resolution image data poses a big challenge not only on the speed of image processing algorithms but also, for the storage and transmission of the data itself.

Matrix decomposition-based image processing techniques are widely used in image-based process monitoring and anomaly detection. They achieve the goal by integrating the prior for background and anomaly

components into the optimization problem. In terms of utilizing the low-rank and sparse property, robust principal component analysis was first proposed by Candès et al. (2011) to decompose a data matrix into low-rank and element-wise sparse components. One of its famous applications is dynamic foreground and static background separation (Bouwmans and Zahzah 2014). Following this approach, numerous algorithm variants have been proposed, including outlier pursuit (Xu et al. 2012), which aims to decompose the data matrix into a low-rank component and a column-wise sparse component, and low-rank plus compressed sparse decomposition (Mardani et al. 2013), which aims to decompose the data matrix into low-rank and compressed sparse components and so on. For utilizing smooth and sparse properties, smooth and sparse decomposition (SSD) methods (Minaee et al. 2015, Yan et al. 2017) are proposed for anomaly detection in images with smooth backgrounds. Following this approach, several explorations have been conducted, including spatiotemporal smooth sparse decomposition (ST-SSD) (Yan et al. 2018), additive tensor decomposition (Mou et al. 2021), and so on. By adopting

the matrix decomposition approach, both the background and anomaly can be captured without detection time delay. However, because of the requirement of storage, transmission, and processing of the whole image signal, it cannot be applied in the scenario with low-transmission bandwidth but high processing speed requirements: for example, the solar flare detection application (Augusto et al. 2011).

To mitigate the data storage and transmission burden and improve sensing efficiency, compressive sensing (CS) (Candes et al. 2006) has been proposed, in which the data are directly collected in a compressed form and then, reconstructed accurately with high probability. More specifically, suppose that the original signal is a *sparse* vector  $y \in \mathbb{R}^n$ , and the main idea is to store and transmit a small set of compressive measurements  $y' = Ay \in \mathbb{R}^p$ , where  $A \in \mathbb{R}^{p \times n}$  is an underdetermined sensing matrix ( $p \ll n$ ) satisfying specific properties. Then, the original signal can be reconstructed from its compressed form  $y'$ , on which assorted image processing algorithms can be applied for defect detection and so on. For a comprehensive review of CS, please refer to Marques et al. (2018) and Rani et al. (2018). Even though promising, the naïve approach that tries to first reconstruct the image from the compressed measurement and then, apply matrix decomposition algorithms for anomaly detection has two issues.

- i. For smooth plus sparse signals, the existence of such a sensing matrix  $A$  satisfying specific properties is unknown.
- ii. The reconstruction process is usually computationally intensive (Marques et al. 2018), which on the other hand, slows down the overall computational speed of the defect detection algorithms.

Recently, to integrate the CS with matrix decomposition algorithms, Waters et al. (2011) proposed an SpaRCS method to recover low-rank and sparse matrices directly from compressive measurements, and Tanner and Vary (2020) gave a rigorous performance discussion. However, those methods do not consider the smooth plus sparse decomposition problems and are not efficient in dealing with high-order data.

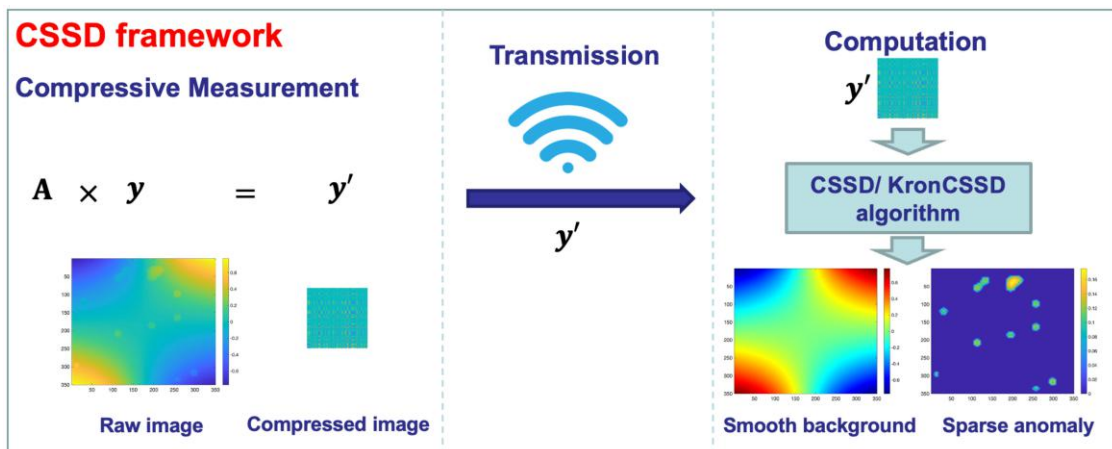
In this paper, we discuss the possibility of adopting compressive data acquisition systems for image-based quality monitoring and fault detection in applications where the background is smooth, and anomalies are sparse. To achieve so, we propose a compressed smooth sparse decomposition (CSSD) framework. In this framework, the signal processing algorithms are directly applied to the compressed data, and no reconstruction step is needed. By doing so, a significant cost reduction in sensing, storage, and transmission as well as a boost in the speed of image processing algorithms can be achieved with negligible performance loss. We also established the theoretical foundation of adopting such a compressive data acquisition system for smooth plus sparse signals as well as the performance guarantee of the proposed algorithm. To further improve its performance in high-order scenarios, a Kronecker compressed smooth sparse sensing (KronCSSD) is proposed.

The remainder of this paper is organized as follows. In Section 2, we present the CSSD framework. In Section 3, we use simulation studies to validate the proposed framework. In Section 4, we demonstrate the proposed framework using several case studies. Finally, Section 5 concludes the paper.

## 2. Compressed Smooth Sparse Decomposition Framework

In this section, we will present the proposed CSSD framework. As mentioned in Section 1, we aim to design a fast and data-efficient method for sparse anomaly (sparse signal component) detection in signals with a smooth background (smooth signal component). For simplicity, we first discuss the methodology for one-dimensional (1D) signals and then, generalize it to  $n$ -dimensional images. Figure 1 provides an overview of the proposed methodology. There are three stages in the proposed methodology. (i) The signal is acquired in its compressed form through compressive measurement. (ii) Then, the compressed data are transmitted to

Figure 1. (Color online) An Overview of the CSSD/KronCSSD Framework



the server. (iii) Finally, a decomposition algorithm will be applied directly to the compressed signal to decompose it into its corresponding smooth and sparse signal components.

Mathematically, let  $\mathbf{y} \in \mathbb{R}^n$  be a smooth plus sparse signal (which will be defined formally in Section 2.1). We aim to store and transmit a small set of compressive measurements  $\mathbf{y}'$  (i.e.,  $\mathbf{y}' = \mathbf{A}\mathbf{y}$ ) and then, reconstruct the smooth and sparse signal components from  $\mathbf{y}'$ . To achieve so, there are three questions to be addressed.

- i. What is a smooth plus sparse signal?
- ii. How do we compress such a signal?
- iii. How do we reconstruct such a signal from compressive measurements?

The remainder of this section is organized as follows to answer those three questions. We start with a formal definition of 1D smooth plus sparse signals in Section 2.1. Based on that, we introduce the compressive measurement method and discuss its theoretical properties for such signals in Section 2.2. In Section 2.3, we present the proposed CSSD framework that can directly reconstruct the smooth and sparse signal components from the compressive measurement by solving the following optimization problem:

$$\begin{aligned} \min_{\boldsymbol{\theta}, \boldsymbol{\theta}_a} \|\boldsymbol{\theta}_a\|_1 \\ \text{s.t. } \|\mathbf{A}(\mathbf{B}\boldsymbol{\theta} + \mathbf{B}_a\boldsymbol{\theta}_a) - \mathbf{y}'\|_2 \leq \epsilon_1, \end{aligned} \quad (1)$$

where  $\mathbf{B}$  and  $\mathbf{B}_a$  are bases,  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}_a$  are corresponding coefficients, and  $\epsilon_1$  is the bound for measurement error. The reconstruction accuracy is also characterized theoretically. Then, we generalize the CSSD algorithm to  $n$  dimensions using Kronecker compressive sensing (KCS) (Duarte and Baraniuk 2011) and propose a KronCSSD formulation in Section 2.4. Finally, in Section 2.5, we present the advantage of the proposed framework and give the strategy of selecting the compressive ratio, tuning parameters, and bases in practice.

## 2.1. The Set of Smooth Plus Sparse Signals

In this section, we define the set of smooth plus sparse signals mathematically. The smooth signals originate from the spline smoothing (De Boor and De Boor 1978, Eilers and Marx 1996), where the raw signal is approximated by a linear combination of a set of spline basis functions for smooth interpolation and denoising, and a spline regression technique is usually utilized. To improve the regression robustness with respect to outliers, outliers are explicitly accounted for in the regression model as a sparse component of

the raw signal (Giannakis et al. 2011, Mateos and Giannakis 2011). This idea was further generalized to incorporate the special structure of outliers (Yan et al. 2017, 2018).

As a summary, a 1D signal  $\mathbf{y} \in \mathbb{R}^n$  is defined as a smooth plus sparse signal if it can be decomposed into two signal components: (i) a smooth signal  $\mathbf{m} \in \mathbb{R}^n$  in a low-dimensional subspace spanned by a set of smooth bases (i.e.,  $\mathbf{m} = \mathbf{B}\boldsymbol{\theta}$ , where  $\mathbf{B} \in \mathbb{R}^{n \times r}$  is a basis matrix with  $r \ll n$ ) and (ii) a sparse signal  $\mathbf{a} \in \mathbb{R}^n$  in a relatively high-dimensional subspace spanned by a set of predefined bases, of which the coefficients admit sparse property (i.e.,  $\mathbf{a} = \mathbf{B}_a\boldsymbol{\theta}_a$ , where  $\mathbf{B}_a \in \mathbb{R}^{n \times q}$  is a basis matrix with  $q \leq n$  and  $\boldsymbol{\theta}_a \in \mathbb{R}^q$  is an  $s$ -sparse vector; i.e.,  $\|\boldsymbol{\theta}_a\|_0 \leq s$ ). Given a smooth plus sparse signal  $\mathbf{y}$ , we define the aforementioned decomposition as SSD (i.e.,  $\mathbf{y} = \mathbf{m} + \mathbf{a}$ ).

To ensure the uniqueness of SSD in a nontrivial case when  $n \leq r + q$ , the following definition is introduced.

**Definition 2.1.** The local support property of  $\mathbf{B}_a$ ,  $\mathbf{B}_a \in \mathbb{R}^{n \times q}$ .  $\mathbf{B}_a$  only has local support such that each column of  $\mathbf{B}_a$  only has nonzero values inside a specific interval. The length of this interval is defined as  $l(\mathbf{B}_a)$ .

Notice that  $l(\mathbf{B}_a) \in \{1, \dots, n\}$ . The local support property with a small  $l$  ensures the sparsity of  $\mathbf{a}$ . For example, the  $B$ -spline basis has a local support property (Unser 1999).

**Definition 2.2.** The incoherence condition of  $\mathbf{B}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times r}$ . Let  $\mathbf{B} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$  be the reduced singular value decomposition (SVD) of  $\mathbf{B}$ , where  $\mathbf{U} \in \mathbb{R}^{n \times r}$ ,  $\boldsymbol{\Sigma} \in \mathbb{R}^{r \times r}$ , and  $\mathbf{V} \in \mathbb{R}^{r \times r}$ . Its incoherence condition parameter  $\mu(\mathbf{B})$  is defined as the smallest value such that

$$\max_{i \in \{1, \dots, r\}} \|\mathbf{U}^T \mathbf{e}_i\|_2 \leq \sqrt{\frac{\mu(\mathbf{B})r}{n}},$$

where  $\mathbf{e}_i$  is the  $i$ th standard basis vector in  $\mathbb{R}^n$ .

Notice that  $\mu(\mathbf{B}) \in [1, \sqrt{n/r}]$ . The incoherent condition with a small  $\mu$  ensures that the  $\mathbf{m}$  is not sparse (Candès et al. 2011).

The following theorem ensures the uniqueness of the SSD decomposition.

**Theorem 2.1.** If  $\mu(\mathbf{B}) < n(2rs)^{-1}$ , then the SSD decomposition is unique with respect to  $\mathbf{m}$  and  $\mathbf{a}$ .

Theorem 2.1 gives the condition that the smooth plus sparse signal can be uniquely decomposed into a smooth part and a sparse part. The proof of Theorem 2.1 is in Appendix A.

Formally, we define the set of smooth plus sparse signals as follows.

**Definition 2.3.** The set of smooth and sparse signals is defined as  $MS_{r,s,\mu,l}$ :

$$MS_{r,s,\mu,l} = \{\mathbf{y} \in \mathbb{R}^n \mid \mathbf{y} = \mathbf{B}\boldsymbol{\theta} + \mathbf{B}_a\boldsymbol{\theta}_a, \mathbf{B}_a \in \mathbb{R}^{n \times q}, \\ l(\mathbf{B}_a) = l, \boldsymbol{\theta}_a \in \mathbb{R}^q, \|\boldsymbol{\theta}_a\|_0 \leq s, \mathbf{B} \in \mathbb{R}^{n \times r}, \\ \mu(\mathbf{B}) = \mu < n(2rsl)^{-1}, \boldsymbol{\theta} \in \mathbb{R}^r\}.$$

For such a smooth plus sparse signal, how to compress it while ensuring the reconstruction performance will be discussed in the next section.

## 2.2. Compressive Sensing for Smooth Plus Sparse Signals

As mentioned in Section 1, to ensure the reconstruction performance, the sensing matrix  $A$  has to satisfy the so-called restricted isometry property (RIP) (Candes 2008). It has been proven that a random matrix can satisfy the RIP property for the sparse signal (Candes 2008), the low-rank signal (Recht et al. 2010), and the rank plus sparse signal (Tanner and Vary 2020). However, the existence of such a matrix for the smooth plus sparse signal, which is the foundation of adopting compressed data acquisition techniques in applications with smooth background and sparse anomalies, is still unknown. In this section, we will discuss the existence of such a sensing matrix. Before stating the result, we will first present the relevant definitions that are necessary to derive the main result.

**Definition 2.4.** RIP for  $MS_{r,s,\mu,l}$ . Let  $A \in \mathbb{R}^{p \times n}$  be a linear measurement matrix. For every quadruple  $(r, s, \mu, l)$ , define the restricted isometry constant (RIC)  $\delta_{r,s,\mu,l}$  to be the smallest positive constant such that

$$(1 - \delta_{r,s,\mu,l})\|\mathbf{y}\|_2 \leq \|\mathbf{A}\mathbf{y}\|_2 \leq (1 + \delta_{r,s,\mu,l})\|\mathbf{y}\|_2, \quad \forall \mathbf{y} \in MS_{r,s,\mu,l}.$$

If such a  $\delta_{r,s,\mu,l} \in (0, 1)$  exists, we say that  $A$  satisfies the RIP.

**Theorem 2.2.** Suppose that  $\delta_{r,2s,\mu,l} < 1$  for some integer  $r, s, l \geq 1$  and positive numbers  $\mu < n(2rsl)^{-1}$ ; then, there is a  $\mathbf{y}_0$  in the set  $MS_{r,s,\mu,l}$ , which is the only solution for  $\mathbf{A}\mathbf{y}_0 = \mathbf{b}$ .

Theorem 2.2 guarantees the uniqueness of the smooth plus sparse signal that satisfies the sensing equation when  $A$  satisfies the RIP. The proof of Theorem 2.2 is in Appendix B.

Next, we prove that for the set of smooth plus sparse signals,  $MS_{r,s,\mu,l}$ , there exists such a matrix  $A$  satisfying the RIP property with RIC =  $\delta_{r,s,\mu,l}$  with high probability.

Notice that the RIP for a matrix is difficult to verify. A suitable set of random matrices that obey the RIP for the set of sparse vectors with high probability (Recht et al. 2010, Tanner and Vary 2020) is defined as follows.

**Definition 2.5.** Nearly isometric matrices (Baraniuk et al. 2008). Let  $A \in \mathbb{R}^{p \times n}$  be a random variable that takes values in linear maps from  $\mathbb{R}^n$  to  $\mathbb{R}^p$ ; then, for any  $\mathbf{y} \in \mathbb{R}^n$ ,  $A$  is nearly isometrically distributed if

- i.  $\mathbb{E}[\|\mathbf{A}\mathbf{y}\|_2^2] = \|\mathbf{y}\|_2^2$  and
  - ii.  $Pr(\|\|\mathbf{A}\mathbf{y}\|_2^2 - \|\mathbf{y}\|_2^2\| \geq \epsilon\|\mathbf{y}\|_2^2) \leq 2e^{-pc_0(\epsilon)}$ ,  $0 < \epsilon < 1$ ,
- where  $c_0(\epsilon)$  is a constant that only depends on  $\epsilon$ .

The  $p \times n$  matrix with independent, identically distributed (i.i.d.) Gaussian entries satisfies those two properties (Baraniuk et al. 2008) (i.e.,  $A_{ij} \sim \mathcal{N}(0, \frac{1}{p})$ , with  $c_0(\epsilon) = \epsilon^2/4 - \epsilon^3/6$ ). There are also other distributions satisfying the nearly isometric property, such as the  $p \times n$  matrix with i.i.d. Bernoulli entries and their related distribution (Baraniuk et al. 2008).

Then, the following theorem states that the nearly isometric matrices can also serve as the sensing matrix for smooth plus sparse signals and gives the magnitude of the number of linear measurements.

**Theorem 2.3.** Let  $A \in \mathbb{R}^{p \times n}$  be a matrix from the families described in Definition 2.5. Furthermore, assume that  $\mu < n(2rsl)^{-1}$  and the basis matrix  $\mathbf{B}_a$  for the sparse signal component satisfies the RIP with RIC  $\delta_{\mathbf{B}_a,s} \in (0, 1)$ ; that is,  $\delta_{\mathbf{B}_a,s}$  to be the smallest positive constant such that

$$(1 - \delta_{\mathbf{B}_a,s})\|\boldsymbol{\theta}_a\|_2 \leq \|\mathbf{B}_a\boldsymbol{\theta}_a\|_2 \leq (1 + \delta_{\mathbf{B}_a,s})\|\boldsymbol{\theta}_a\|_2, \\ \forall \boldsymbol{\theta}_a \in \{\boldsymbol{\theta}_a \in \mathbb{R}^q, \|\boldsymbol{\theta}_a\|_0 \leq s\}.$$

For a given  $\delta \in (0, 1)$ , there exists constants  $c_1, c_2 > 0$  depending only on  $\delta$ , such that the RIC for  $MS_{r,s,\mu,l}$  is upper bounded by  $\delta$ , with the probability of at least  $1 - \exp(-c_1p)$ , whenever

$$p \geq c_2 \left( \ln 2 + r \ln \frac{24}{\delta} \tau_1 + s \left( 1 + \ln \frac{24}{\delta} \tau_0 + \ln \frac{n}{s} \right) \right),$$

where  $\eta = \sqrt{\frac{\mu r s l}{n}}$ ,  $\tau_0 = \frac{1}{\sqrt{(1-\delta_{\mathbf{B}_a,s})(1-\eta^2)}}$ ,  $\tau_1 = \|\mathbf{B}^\dagger\|_2 \left( 1 + \frac{1}{\sqrt{1-\eta^2}} \right)$ , and  $\mathbf{B}^\dagger = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T$ .

Theorem 2.3 states that if the nearly isometric matrix is selected as the sensing matrix, the RIC for  $MS_{r,s,\mu,l}$  is upper bounded with high probability. In practice,  $\mathbf{B}$  and  $\mathbf{B}_a$  are prespecified based on the understanding/engineering knowledge of the process (please refer to Section 2.5.4 for more detail). Theorem 2.3 also provides a guidance on the magnitude of linear measurements  $p$ , which determines the compressive measurement matrix  $A$ .

The proof of Theorem 2.3 is in Appendix C.

In this section, we answered two fundamental questions. (i) How do we compress the smooth plus sparse signal (design the compressive measurement matrix  $A$ )?

(ii) How many linear measurements are needed to preserve the information in smooth plus sparse signals with high probability?

In the next section, we will discuss the problem formulation to recover the smooth and sparse signal components simultaneously from the compressed signal using the CSSD framework.

### 2.3. Compressed Smooth Sparse Decomposition

As mentioned in Section 1, one way to recover the smooth component and sparse component is first reconstructing the compressed image and then, using the SSD algorithm. However, this will slow down the speed of the defect detection algorithm. Instead, we propose to solve the one-step convex relaxation Problem (1). Natural questions are if we can recover the smooth and sparse signals from the compressed measurement  $\mathbf{y}'$  by solving Problem (1) and what the accuracy is. The following theorem guarantees the recovery performance of the proposed convex relaxation Problem (1).

**Theorem 2.4.** Let  $A \in \mathbb{R}^{p \times n}$  be a matrix from the families described in Definition 2.5. Let the signal  $\mathbf{y} = \mathbf{m}_0 + \mathbf{a}_0 \in MS_{r,s,\mu,l}$ , where  $\mathbf{m}_0 = B\boldsymbol{\theta}_0$  and  $\mathbf{a}_0 = B_a\boldsymbol{\theta}_{a0}$ . Assume that Problem (1) is feasible, and let the optimal solution be  $\mathbf{m}^* = B\boldsymbol{\theta}^*$  and  $\mathbf{a}^* = B_a\boldsymbol{\theta}_a^*$ . Assume that the basis matrix  $B_a$  for sparse signal component satisfies the RIP with RIC  $\delta_{B_a,2s} \in (0, 1)$ . Let  $a = (1 + \alpha_1 + \alpha_2)\gamma^2 + 2\alpha_2 + 2$  and  $c = 1 - \gamma^2\alpha_1\alpha_2 - \alpha_2^2$ , where  $\alpha_1 = \frac{\eta}{1-\eta^2}$ ,  $\alpha_2 = \frac{\sqrt{2}\eta}{1-2\eta^2}$ ,  $\gamma = \sqrt{\frac{1+\delta_{B_a,2s}}{1-\delta_{B_a,2s}}}$ , and  $\eta = \sqrt{\frac{lrs}{n}}$ . Suppose that  $r, s, l \in \mathbb{N}$  and  $\mu < n(2rsl)^{-1}$ , such that  $c > 0$  and  $\delta_{r,3s,\mu,l} \in (0, c/a)$ ; then,

$$\|\mathbf{a}_0 - \mathbf{a}^*\|_2 = \|B_a\boldsymbol{\theta}_{a0} - B_a\boldsymbol{\theta}_a^*\|_2 \leq C_a\epsilon_1$$

and

$$\|\mathbf{m}_0 - \mathbf{m}^*\|_2 = \|B\boldsymbol{\theta}_0 - B\boldsymbol{\theta}^*\|_2 \leq C_m\epsilon_1,$$

where

$$C_a = \frac{(1 + \gamma^2)(1 + \alpha_2)\sqrt{1 + \delta_{r,3s,\mu,l}}}{c - a\delta_{r,3s,\mu,l}}$$

and

$$C_m = \frac{\sqrt{1 + \delta_{r,3s,\mu,l}} + \left(\delta_{r,3s,\mu,l} + \frac{\gamma^2}{1+\gamma^2}\alpha_1 + \frac{1}{1+\gamma^2}\alpha_2\right)C_a}{(1 - \delta_{r,3s,\mu,l})}.$$

Theorem 2.4 gives the conditions that the proposed convex relaxation Problem (1) can recover the true smooth and sparse signal up to a constant times the noise bound.

The proof of Theorem 2.4 is in Appendix D.

Notice that one advantage of the proposed CSSD framework is that it is compatible with existing decomposition

algorithms. For example, for the SSD algorithm proposed by Yan et al. (2017), the problem formulation becomes

$$\min_{\boldsymbol{\theta}, \boldsymbol{\theta}_a} \|\mathbf{y}' - A(B\boldsymbol{\theta} + B_a\boldsymbol{\theta}_a)\|_2^2 + \lambda\|\boldsymbol{\theta}_a\|_1, \quad (2)$$

which can be solved efficiently by using the algorithm proposed by Yan et al. (2017).

### 2.4. Kronecker Compressed Smooth Sparse Decomposition

In the previous section, we discussed the CSSD framework for the 1D signal. In this section, we will generalize the proposed CSSD to KronCSSD framework for high-order tensor data.

Let  $\mathcal{Y} \in \mathbb{R}^{n_1 \times \dots \times n_d}$  be the original signal and  $\mathbf{y} \in \mathbb{R}^N$  be its corresponding vectorized signal (i.e.,  $\mathbf{y} = \text{vec}(\mathcal{Y})$ ) and  $N = \prod_{i=1}^d n_i$ . Let  $A \in \mathbb{R}^{p \times N}$  be a measurement matrix satisfying the RIP. Let  $\mathbf{y}' \in \mathbb{R}^p$  be the compressed data, such that  $\mathbf{y}' = A\mathbf{y}$ .  $B = \otimes_{i=1}^d B_i$  and  $B_a = \otimes_{i=1}^d B_{ai}$  are the known bases for smooth and sparse components, respectively. Notice that problem formulation (1) can still be used for recovering the high-order smooth and sparse signal components from the compressive measurement. However, there are two issues. (i) In practice, the global CS measurements matrix  $A$  is hard to realize using the CS device (Duarte and Baraniuk 2011). (ii) The resulting bases  $B$  and  $B_a$  can be extremely large as the dimension of the data increases, which will not only cause a big challenge in the storage of such a large matrix but also, result in the computational issue in handling such large matrices.

In reality, the high-dimensional tensor data usually have low-rank properties along each mode, which have been extensively exploited in tensor low-rank modeling techniques, such as CANDECOMP/PARAFAC (CP)/Tucker decompositions (Kolda and Bader 2009). This makes it possible for designing a sensing matrix for each mode, which is called Kronecker CS (Duarte and Baraniuk 2011). Inspired by the KCS method, we propose the KronCSSD framework formulation as follows.

Let  $A_i \in \mathbb{R}^{p_i \times n_i}$  be a measurement matrix with RIP for each mode of the tensor; then, we have the following formulation,

$$\begin{aligned} & \min_{\boldsymbol{\Theta}, \boldsymbol{\Theta}_a} \|\text{vec}(\boldsymbol{\Theta}_a)\|_1 \\ & \text{s.t. } \|\text{vec}(\mathcal{Y}_1 - \mathcal{Y}')\|_2 \leq \epsilon_1, \\ & \mathcal{Y}_1 = \boldsymbol{\Theta} \times_1 (A_1 B_1) \times_2 \cdots \times_d (A_d B_d) \\ & \quad + \boldsymbol{\Theta}_a \times_1 (A_1 B_{a1}) \times_2 \cdots \times_d (A_d B_{ad}), \end{aligned} \quad (3)$$

where  $\mathcal{Y}' = \mathcal{Y} \times_1 A_1 \times_2 \cdots \times_d A_d$  is the compressive measurement.  $\boldsymbol{\Theta} \in \mathbb{R}^{r_1 \times \dots \times r_d}$  and  $\boldsymbol{\Theta}_a \in \mathbb{R}^{q_1 \times \dots \times q_d}$  are the basis coefficients for the smooth and sparse signal components, respectively.

The proposed KronCSSD framework is a nontrivial generalization of the CSSD framework for 1D signals. Its performance will be shown empirically in simulation and case studies. The theoretical discussion of the KronSSD framework is left for future work.

For example, for a two-dimensional (2D) image  $Y \in \mathbb{R}^{n_1 \times n_2}$ , when adopting the SSD algorithm (Yan et al. 2017), the problem formulation becomes

$$\min_{\Theta, \Theta_a} \|Y' - A_1 B_1 \Theta B_2^T A_2^T - A_1 B_{a1} \Theta_a B_{a2}^T A_2^T\|_2^2 + \lambda \|\text{vec}(\Theta_a)\|_1, \quad (4)$$

where  $Y'$  is the compressed image (i.e.,  $Y' = A_1 Y A_2^T$ ). It can be solved efficiently by using the algorithm proposed by Yan et al. (2017).

## 2.5. Discussion

**2.5.1. Advantages of the Proposed CSSD/KronCSSD Methods.** In this section, we will give a brief discussion about the advantages of the proposed CSSD/KronCSSD methods. We define the compressive ratio as  $c = \prod_{i=1}^d p_i/n_i$ . The smaller the compressive ratio, the fewer data will be transmitted. The proposed CSSD/KronCSSD methods have the following characteristics.

i. We propose to directly acquire the compressed image, which not only reduces the sensing cost but also, reduces the data transmission and storage cost by  $\prod_{i=1}^d p_i/n_i$  times.

ii. The smooth and sparse signal components can be recovered by solving a smaller-scale convex optimization problem with input data  $\prod_{i=1}^d p_i/n_i$  times smaller than that of the original problem, which significantly boosts the computation.

**2.5.2. Compressive Ratio Selection in Practice.** Theorems 2.3 and 2.4 show that the 1D smooth plus sparse signal can be recovered with high probability from a compressive measurement if the compressive ratio is above a specific threshold. However, there are several parameters (such as  $c_1, c_2$ ) that are difficult to obtain when calculating the threshold in practice. We propose a practical procedure of selecting the compressive ratio utilizing the historical data. Suppose some training signals with their real background and anomalies are available; the compressive ratio can be selected with the following guidelines.

i. If there is a requirement for reconstruction accuracy for smooth and sparse signal components, then  $\hat{p}$  is chosen as the smallest value that satisfies such a requirement.

ii. If there is no such requirement, we recommend choosing the compressive ratio corresponding to the sharp change point of the slope of the loss function-compressive ratio curve. This point exists because of the existence of such a threshold, after which the reconstruction with

high probability is guaranteed by Theorems 2.3 and 2.4. We will demonstrate this in the simulation study in Section 3.1.

For high-order tensor data, the selection of the  $p_i$  for each mode can be challenging. We provide some empirical guidelines as follows.

i. If the smoothness of the background is similar along different modes, a unified compressive ratio is recommended.

ii. If the smoothness of the background is different, more sensing budget should be allocated to the mode along which the background is less smooth.

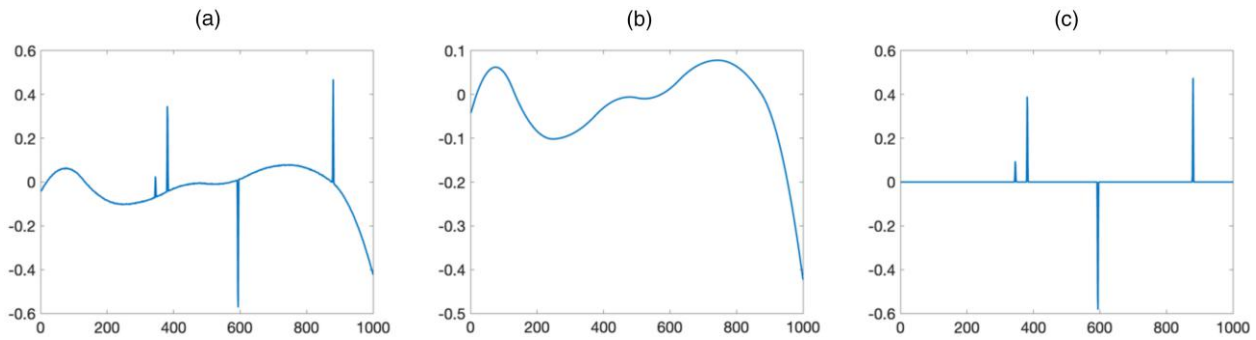
iii. We recommend fixing the ratio among  $p_i$  along different modes and adopting steps (i) and (ii) for 1D signal to determine the compressive ratio.

**2.5.3. Tuning Parameter Selection in Practice.** Notice that in problem formulations (1) and (3), there is a hyperparameter  $\epsilon_1$ , indicating the bound for measurement noise. If the measurement error bound is known from the accuracy of the measurement device, it can be directly used here. Otherwise, a crossvalidation step using historical data is recommended. Similarly, there is a tuning parameter  $\lambda$  in their corresponding Lagrangian form Equations (2) and (4), which controls the sparsity of the decomposed anomaly. Crossvalidation can be used to determine this parameter. For more detail, please refer to Yan et al. (2017).

**2.5.4. Bases Selection in Practice.** The bases  $B$  and  $B_a$  are prespecified based on the understanding/engineering knowledge of the process. The selection of such bases is discussed in detail in section 3.4 of Yan et al. (2017). In general, any smooth basis, such as splines or kernels, can be used for the background. For sparse anomalies, such as small regions scattered over the background or in the form of thin lines, an identity basis is recommended. Linear (quadratic)  $B$  splines are recommended for anomalous regions with sharp corners (curved boundaries). However, to ensure the uniqueness of the smooth sparse decomposition, we do require the bases  $B$  and  $B_a$  to satisfy specific properties, such as those mentioned in Definitions 2.1 and 2.2 and Theorem 2.1, which should be checked for selected bases. We will demonstrate this in the simulation study in Section 3.1.

## 3. Simulation Study

In this section, we will demonstrate the proposed CSSD and KronCSSD framework with simulation studies. First, the CSSD method is applied on 1D signals in Section 3.1, and then, we apply the KronCSSD method on 2D images in Section 3.2.

**Figure 2.** (Color online) A Sample Raw Signal and Its Smooth and Sparse Components

Notes. (a) Raw signal. (b) Smooth signal component  $m$ . (c) Sparse signal component  $a$ .

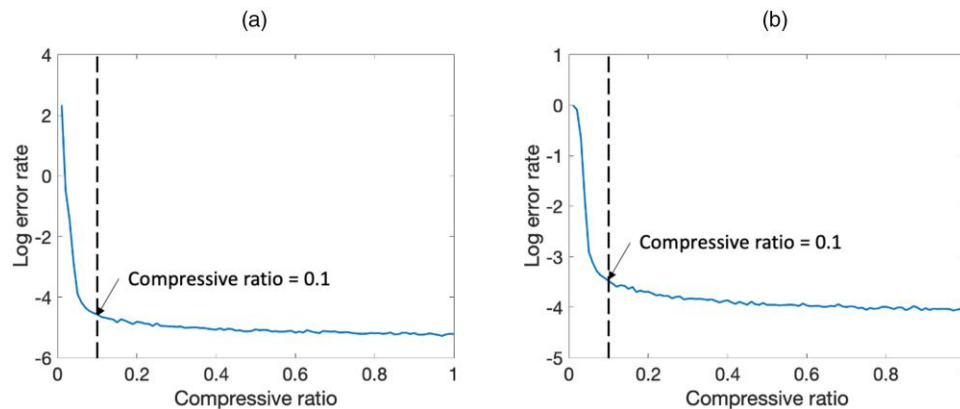
### 3.1. CSSD on a 1D Signal

A 1D signal ( $\mathbf{y} \in \mathbb{R}^n$ ) is assumed to be a superposition of the smooth signal component ( $\mathbf{m} = \mathbf{B}\boldsymbol{\theta}$ ), the sparse signal component ( $\mathbf{a} = \mathbf{B}_a\boldsymbol{\theta}_a$ ), and noise  $\mathbf{e}$  (i.e.,  $\mathbf{y} = \mathbf{m} + \mathbf{a} + \mathbf{e}$ ). To demonstrate the proposed CSSD framework, we will first conduct compressive data acquisition on the raw signal (i.e.,  $\mathbf{y}' = \mathbf{A}\mathbf{y}$ ). Then, the data reconstruction and decomposition are achieved in one step. The performance is evaluated by the relative error between the true signal  $\mathbf{a}$  and the reconstructed one  $\hat{\mathbf{a}}$  (i.e.,  $\|\mathbf{a} - \hat{\mathbf{a}}\|_2 / \|\mathbf{a}\|_2$  and  $\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2$  for the smooth background).

In the simulation study, we generate a 1D signal from  $MS_{r,s,\mu,l}$  and let  $n = 1,000$ . The smooth background is generated from a random linear combination of  $B$ -spline bases with three knots ( $r = 10$ ) (i.e.,  $\mathbf{m} = \mathbf{B}\boldsymbol{\theta}$ , where  $\mathbf{B} \in \mathbb{R}^{n \times r}$  and  $\boldsymbol{\theta} \in \mathbb{R}^r$  is a random vector such that  $\theta_i \sim N(0,1)$ ,  $i \in \{1, \dots, 4\}$ ). The incoherence condition parameter  $\mu = \mu(\mathbf{B}) = 0.82$ . The sparse signal component is generated by a sparse random linear combination of degree 2  $B$ -spline bases with 500 knots ( $q = 500$ ,  $l = 4$ ): that is,  $\mathbf{a} = \mathbf{B}_a\boldsymbol{\theta}_a$ , where  $\mathbf{B}_a \in \mathbb{R}^{n \times q}$  and  $\boldsymbol{\theta}_a \in \mathbb{R}^q$  is a four-sparse random vector ( $s = 4$ ) such that its nonzero

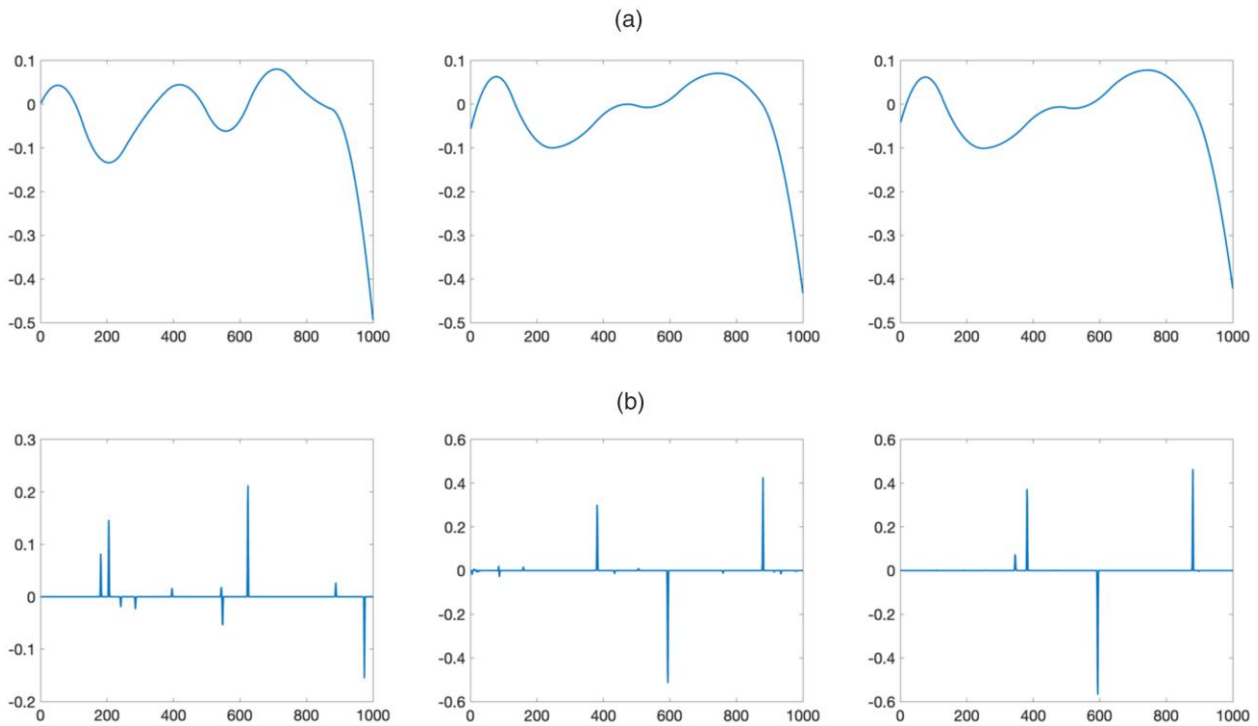
elements follow i.i.d. standard normal distribution. Notice that the RIC,  $\delta_{B_a,s}$  for matrix  $\mathbf{B}_a$ , is hard to calculate in general. However, there is a loose upper bound that can be used, which is  $\delta_{B_a,s} \leq \delta_{B_a,p} = \max\{\lambda_{\max} - 1, 1 - \lambda_{\min}\} = 0.60$ , where  $\lambda_{\max}$  and  $\lambda_{\min}$  are the maximum and minimum singular values of  $\mathbf{B}_a$ , respectively. The noise signal generated as a random vector  $\mathbf{e} \in \mathbb{R}^n$  is a random vector such that  $e_i \sim N(0, 0.001^2)$ ,  $i \in \{1, \dots, n\}$ . Figure 2 shows a sample raw signal and its smooth, sparse components.

Before we state the result, we first check the assumption of Theorems 2.3 and 2.4. For Theorem 2.3, it is easy to check that  $0.82 = \mu \leq \frac{n}{2rs} = \frac{1,000}{2 \times 10 \times 4 \times 4} = 6.25$  and also, that  $\delta_{B_a,s} \leq 0.60 \in (0, 1)$ . For Theorem 2.4,  $\delta_{B_a,2s} \leq 0.60 \in (0, 1)$ ,  $c = 0.37 > 0$ , and  $\delta_{r,3s,\mu,l} \in (0, 0.04)$ . Notice that the range for  $\delta_{r,3s,\mu,l}$  is small in this case because of the loose upper bound for  $\delta_{B_a,2s}$ , which can be further improved. This indicates that the smooth and sparse signal can be recovered by the proposed algorithm with high probability provided that the compressive ratio is above a specific threshold, which is demonstrated by the following observation.

**Figure 3.** (Color online) Average Log Relative Reconstruction Error

Notes. (a) Log relative error for the smooth component. (b) Log relative error for the sparse component.

**Figure 4.** (Color online) Reconstructed Smooth and Sparse Signal Components



Notes. (a) Smooth signal component of compressive ratios 0.02 (left panel), 0.05 (center panel), and 0.1 (right panel). (b) Sparse signal component of compressive ratios 0.02 (left panel), 0.05 (center panel), and 0.1 (right panel).

The simulation is repeated 100 times, and the average log relative error for the background and sparse signal components with respect to compressive ratio are shown in Figure 3. We can observe a large error at the beginning, and it drops very fast with an increase of the compressive ratio. Then, above a threshold (0.1 approximately) of compressive ratio, the error becomes small (below 3%), and the decrease of error becomes less significant, which demonstrates the effectiveness of the reconstruction algorithm. The threshold of 0.1 can be chosen as the compressive ratio mentioned in Section 2.5.2.

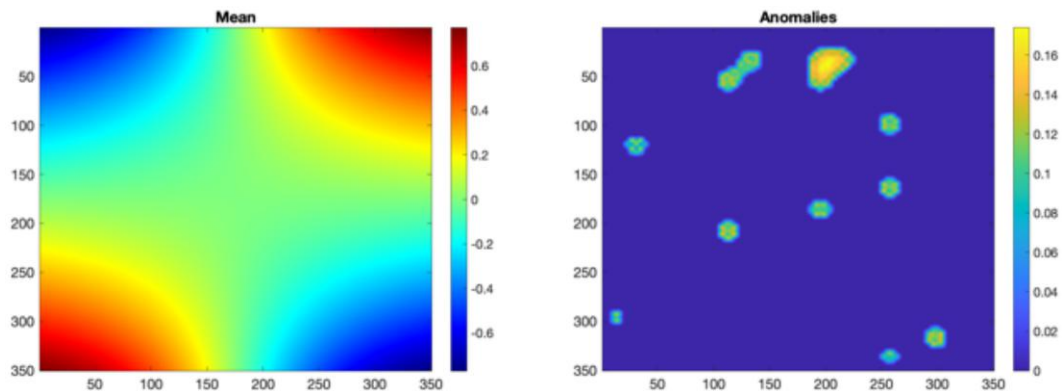
The reconstructed signal components in 1 of the 100 simulations are shown in Figure 4. We can observe that when adopting the compressive ratio of 0.1, both the smooth and sparse signal components can be reconstructed with high accuracy.

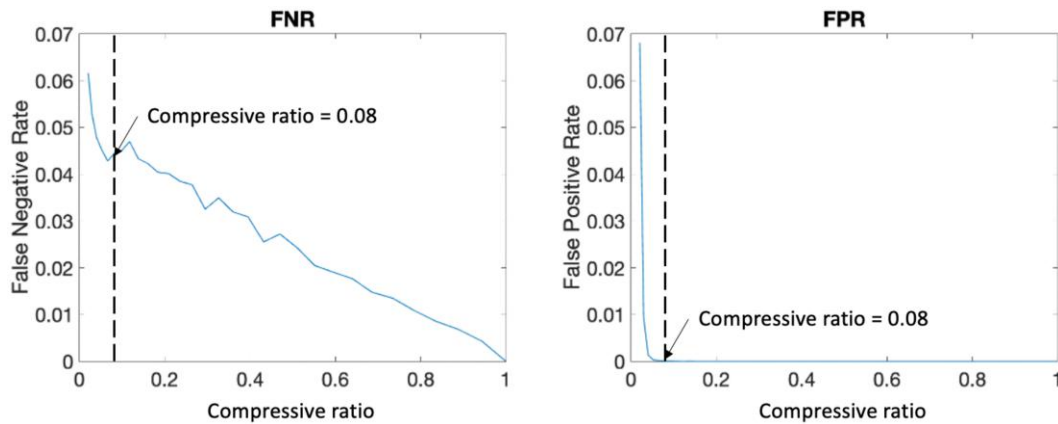
We also vary the magnitude of sparse signals to examine the reconstruction performance of the proposed method. The result and analysis are provided in Appendix H.

### 3.2. CSSD on a 2D Image

In this simulation study, we aim to decompose an image into the smooth background, sparse anomalies,

**Figure 5.** (Color online) True Background and Anomalies



**Figure 6.** (Color online) Average FNR and FPR

and noise. A  $350 \times 350$  image with smooth background and the sparse anomaly is generated similar to Yan et al. (2017) (i.e.,  $Y = M + A + E$ , where  $M$  is the smooth background,  $A$  is the sparse anomalies, and  $E$  is i.i.d. Gaussian noise such that  $E_i \sim NID(0, \sigma^2)$ ). The smooth background is generated from a linear combination of  $B$ -spline bases with  $3 \times 3$  knots, and the anomalies are generated from a sparse linear combination of  $B$ -spline bases with  $88 \times 88$  knots. The background and anomalies are shown in Figure 5.

We can see that the anomaly size covers a large range. The mean absolute value of the background plus anomaly is  $\mu = 0.21$ . In this simulation, we first study the reconstruction performance under a noise-free scenario. Then, we increase the noise level to test the robustness of the algorithm.

**3.2.1. Noise-Free Case.** In this section, we vary the compressive ratio from 4% to 100%. For each compressive ratio, we simulate 100 times and the average false-negative rate (FNR) and false-positive rate (FPR) are reported in Figure 6. The FPR is defined as the portion of normal pixels predicted as anomaly,

$$FPR = \frac{\sum_{i,j} (1 - I_{A \neq 0}\{A(i,j)\}) I_{\hat{A} \neq 0}\{\hat{A}(i,j)\}}{\sum_{i,j} (1 - I_{A \neq 0}\{A(i,j)\})},$$

and the FNR is defined as the portion of anomalous pixels predicted as normal background,

$$FNR = \frac{\sum_{i,j} I_{A \neq 0}\{A(i,j)\} (1 - I_{\hat{A} \neq 0}\{\hat{A}(i,j)\})}{\sum_{i,j} I_{A \neq 0}\{A(i,j)\}},$$

where  $A$  is the true anomaly,  $\hat{A}$  is the predicted anomaly, and  $I_{\Omega}(\cdot)$  is the indicator function: that is,

$$I_{\Omega}(x) = \begin{cases} 1, & \text{if } x \in \Omega \\ 0, & \text{otherwise.} \end{cases}$$

From Figure 6, we can see that both the FPR and FNR ratios decrease as the compressive ratio increases. The

FPR drops so fast that when the compressive ratio achieves 8%, there is no false alarm, which is desired for the anomaly detection algorithm. The FNR drops slower, and less than 5% of the anomaly pixels are ignored when the compressive ratio achieves 8%. However, it does not miss any cluster of the anomalies, even though some of them are small.

The recovered sparse signal components are shown in Figure 7 for compressive ratios 4% (Figure 7(a)), 8% (Figure 7(b)), 33% (Figure 7(c)), and 73% (Figure 7(d)). For comparison, we apply the SSD algorithm (Yan et al. 2017), of which the FNR and the FPR are zero and computation time is 0.16 seconds.

We record the computation time in Figure 8. A significant boosting of the computation is observed when the compressive ratio is 8%, which speeds up the SSD algorithm by 4.3 times.

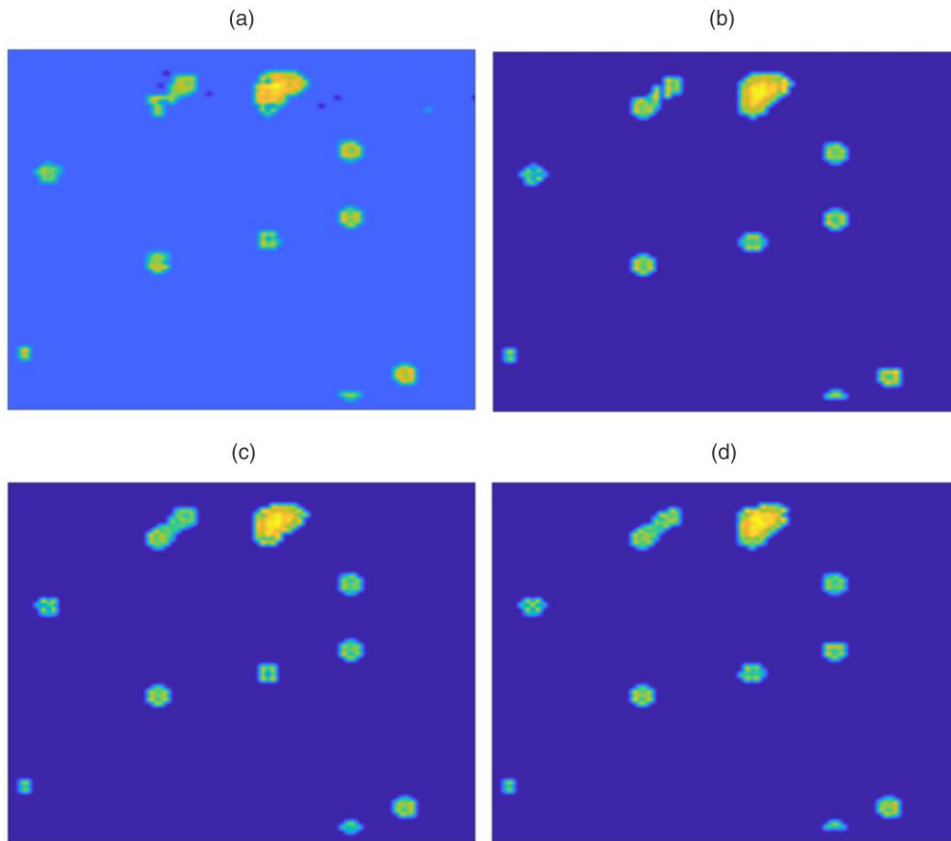
**3.2.2. Noisy Case.** The signal-to-noise ratio  $\mu/\epsilon \in [4, 40]$  is studied to evaluate the robustness of the algorithm. The three-dimensional plot of FPR, the signal-to-noise ratio, and the compressive ratio are shown in Figure 9.

The purple line indicates the equipotential line of  $FPR = 0.01$ . We can see that with the increase of the signal-to-noise ratio, we are allowed to use a less compressive ratio to achieve satisfactory decomposition results. The majority of the area lies below the equipotential line of  $FPR = 0.01$ , which means that the proposed algorithm is robust.

## 4. Case Study

In this section, we use three real cases to demonstrate the effectiveness of the proposed CSSD/KronCSSD framework. For comparison, we also apply the SSD method in each case study. The compressive ratio ( $c$ ) and the average computational time ( $t$ ) for a single image are reported in Table 1. A significant transmission bandwidth reduction and computation boost can be observed with negligible performance degradation (Figures 10–12), compared with the vanilla SSD algorithm.

**Figure 7.** (Color online) The Recovered Sparse Signal Components

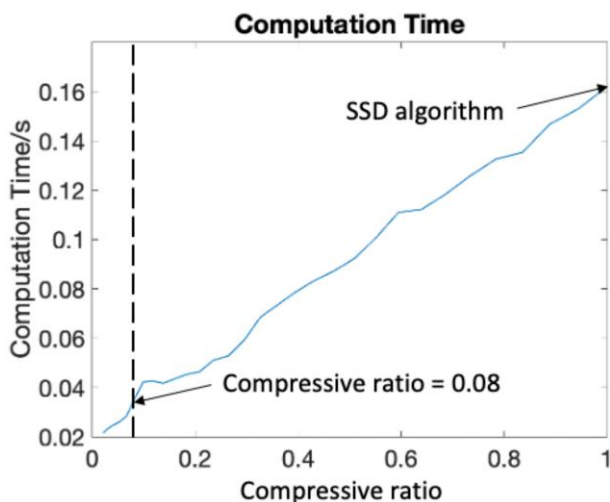


Notes. (a) Compressive ratio 4%. (b) Compressive ratio 8%. (c) Compressive ratio 33%. (d) Compressive ratio 73%.

#### 4.1. Surface Defect Detection in a Steel Rolling Process

As mentioned in Section 1, vision sensors collect high-resolution images of the product surface with a high data acquisition rate in the rolling processes. This poses a

**Figure 8.** (Color online) The Computation Time

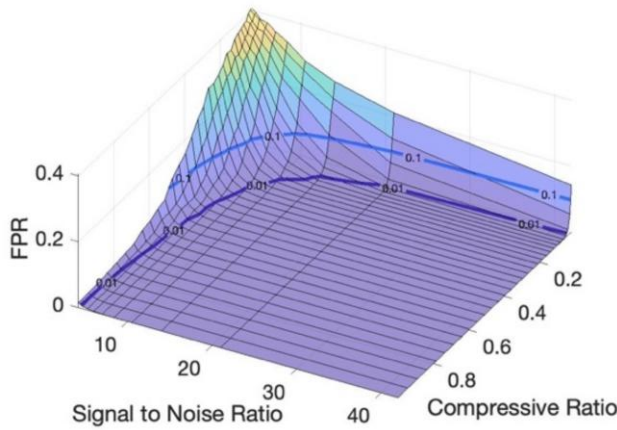


challenge in data storage, transmission, and processing. One sample image of size  $128 \times 512$  with typical anomalies is shown in Figure 10(a). The black scratches shown in the red block are surface anomalies. For a detailed description, the readers are encouraged to refer to Yan et al. (2018). The detected anomalies are shown in Figure 10, (b) and (c) by using the proposed KronCSSD algorithm and the SSD algorithm, respectively. The data set has 100 images, and more example results can be found in Appendix I.

The KronCSSD method achieves a similar anomaly detection performance with 54% bandwidth and is 1.8 times faster. By adopting the KronCSSD method, (i) we can achieve a faster anomaly detection and thus, reduce the loss through a timely intervention of the manufacturing process, and (ii) we can keep the manufacturing inspection information for a longer time period with the same storage capability, which is important for root cause analysis.

#### 4.2. Solar Flare Detection

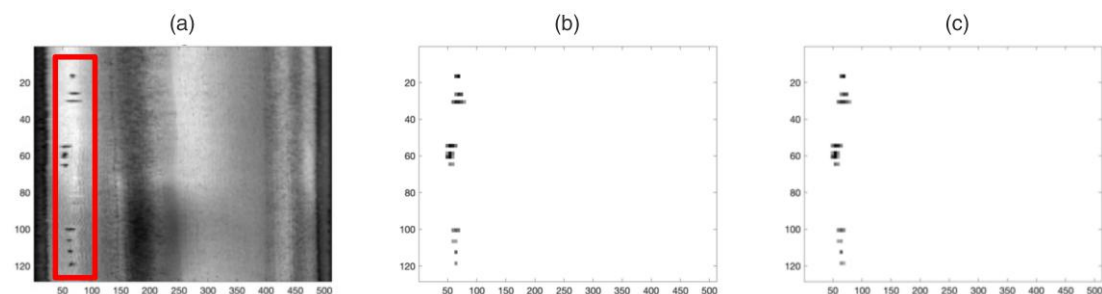
Another important application is solar flare detection from satellite images. A solar flare is defined as a sudden, transient, and intense variation in brightness over

**Figure 9.** (Color online) The Contour Plot of FPR, Signal-to-Noise Ratio, and Compressive Ratio

the sun's surface. It has a significant influence on radio communication on Earth. Each second, thousands of high-resolution images are captured by a satellite, which poses a big challenge to real-time data transmission and processing (Yan et al. 2018). The data set has 300 images, and more example images can be found in Appendix I.

One sample image of size  $232 \times 292$  with a typical solar flare is shown in Figure 11(a), where the yellowish bright region is the solar flare. The detected anomalies are shown in Figure 11, (b) and (c) by using the proposed KronCSSD and SSD algorithms, respectively. The KronCSSD method achieves a similar solar flare detection performance as the SSD method but with 22% bandwidth, and it is 2.5 times faster. Notice that the decomposed images can also be used for downstream tasks, such as control charts and so on, which are beyond the scope of this paper.

By adopting the KronCSSD method, we can improve the transmission rate under the same transmission bandwidth and thus, achieve almost five times faster solar flare detection, which is of vital importance for protecting radio communications, power grids, and navigation systems.

**Figure 10.** (Color online) Steel Rolling Images and Detected Anomalies

Notes. (a) Raw image. (b) KronCSSD result. (c) SSD result.

**Table 1.** Comparison Between KronCSSD and SSD

	Surface defect		Solar flare		Indentation	
	<i>c</i>	<i>t/s</i>	<i>c</i>	<i>t/s</i>	<i>c</i>	<i>t/s</i>
KronCSSD	54%	0.073	22%	0.034	48%	0.053
SSD	1	0.135	1	0.086	1	0.094

### 4.3. Silicon Surface Indentation Detection

The stress map of size  $90 \times 550$  of a silicon surface laminate with surface indentation is shown in Figure 12, where clusters of high-stress areas indicate the surface indentation (Yan et al. 2017). We aim to detect those high-stress areas. The detected anomalies are shown in Figure 12, (b) and (c). We can see that the KronCSSD method achieves a similar detection performance with the SSD method but with 40% bandwidth, and it is 1.8 times faster, which significantly reduces the storage and transmission cost and improves the processing speed.

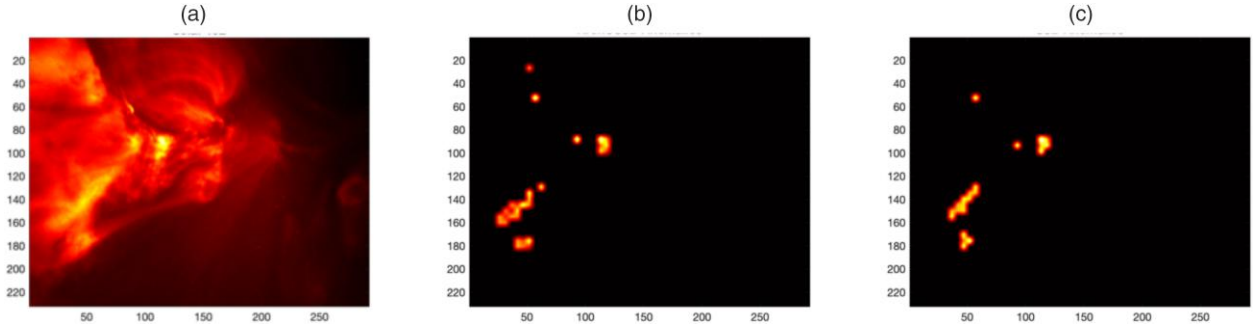
## 5. Conclusion

In this paper, we proposed a CSSD framework for efficient data acquisition, transmission, and processing for sparse anomaly detection in smooth backgrounds. To further enhance its computational efficiency, a KronCSSD framework is proposed for tensor data.

The contributions of this work are twofold. (i) Theoretically, we showed the feasibility of combining compressive sensing and smooth sparse decomposition. This enables the adoption of a compressive data acquisition approach. (ii) Practically, the proposed framework is compatible with many existing decomposition-based anomaly detection algorithms, such as SSD, ST-SSD, and so on, which achieve both a significant cost reduction in sensing, storage, and transmission and a boost in speed but with negligible loss in their performance.

In this article, we use a simulation study to demonstrate the effectiveness and robustness of the proposed CSSD/KronCSSD framework. Three case studies across different applications demonstrate the versatility of the proposed framework. The authors believe that

**Figure 11.** (Color online) Solar Activity Images and Detected Solar Flare



Notes. (a) Raw image. (b) KronCSSD result. (c) SSD result.

the CSSD/KronCSSD framework can be applied in a wider range of applications toward more efficient data acquisition, transmission, and processing.

Further studies on the theoretical properties of the proposed KronCSSD can be a future direction.

## Appendix A

**Proof of Theorem 2.1.** We prove Theorem 2.1 by contradiction. Assume there exists two different decompositions for the same smooth plus sparse signal  $\mathbf{y}$  (i.e.,  $\mathbf{y} = \mathbf{m}_1 + \mathbf{a}_1 = \mathbf{m}_2 + \mathbf{a}_2$ , where  $\mathbf{m}_1 \neq \mathbf{m}_2$  and  $\mathbf{a}_1 \neq \mathbf{a}_2$ ). Then,

$$\mathbf{m}_1 - \mathbf{m}_2 = -\mathbf{a}_1 + \mathbf{a}_2. \quad (\text{A.1})$$

Because  $\mathbf{m}_1 \neq \mathbf{m}_2$ , we can normalize both side by  $\|\mathbf{m}_1 - \mathbf{m}_2\|_2$  and denote  $\tilde{\mathbf{m}} = (\mathbf{m}_1 - \mathbf{m}_2)/\|\mathbf{m}_1 - \mathbf{m}_2\|_2$  and  $\tilde{\mathbf{a}} = (-\mathbf{a}_1 + \mathbf{a}_2)/\|\mathbf{m}_1 - \mathbf{m}_2\|_2$ . Notice that  $\tilde{\mathbf{m}}$  is in the column space of  $\mathbf{B}$ , which is spanned by columns of the  $\mathbf{U}$  (recall that  $\mathbf{B} = \mathbf{U}\Sigma\mathbf{V}^T$ ; i.e.,  $\tilde{\mathbf{m}} = \mathbf{U}\mathbf{x}$ , where  $\mathbf{x}$  is the coefficient vector,  $\mathbf{x} \in \mathbb{R}^r$ ). Because  $\|\tilde{\mathbf{m}}\|_2 = 1$ , we have  $\|\mathbf{x}\|_2 = 1$ . We can bound each element in  $\tilde{\mathbf{m}}$  as follows:

$$|\tilde{m}_i| = \mathbf{e}_i^T \mathbf{U}\mathbf{x} \leq \|\mathbf{e}_i^T \mathbf{U}\|_2 \|\mathbf{x}\|_2 \leq \max_{i \in \{1, \dots, r\}} \|\mathbf{U}^T \mathbf{e}_i\|_2 \leq \sqrt{\frac{\mu(\mathbf{B})r}{n}} < \sqrt{\frac{1}{2ls}}, \quad \forall i \in \{1, \dots, n\}.$$

According to Definition 2.1,  $\|\mathbf{a}_1\|_0 \leq ls$  and  $\|\mathbf{a}_2\|_0 \leq ls$ . Therefore,  $\|\tilde{\mathbf{a}}\|_0 \leq 2ls$ . Moreover, according to Equation (A.1), we

conclude that  $\|\tilde{\mathbf{m}}\|_0 \leq 2ls$ . Denote the support of  $\tilde{\mathbf{m}}$  as  $\mathbf{T}_{\tilde{\mathbf{m}}}$ ; we have  $\|\tilde{\mathbf{m}}\|_2 = \sqrt{\sum_{i \in \mathbf{T}_{\tilde{\mathbf{m}}}} |\tilde{m}_i|^2} < 1$ , which is a contradiction.  $\square$

## Appendix B

**Proof of Theorem 2.2.** We prove Theorem 2.2 with contradiction. Assume there exists another vector  $\mathbf{y} = \mathbf{B}\boldsymbol{\theta} + \mathbf{B}_a\boldsymbol{\theta}_a \in MS_{r,s,\mu,l}$  such that  $\mathbf{A}\mathbf{y} = \mathbf{b}$  and  $\mathbf{y} \neq \mathbf{y}_0$ . Then,  $\mathbf{z} = \mathbf{y} - \mathbf{y}_0 = \mathbf{B}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \mathbf{B}_a(\boldsymbol{\theta}_a - \boldsymbol{\theta}_{a0})$  is a nonzero vector. Because  $\|\boldsymbol{\theta}_a - \boldsymbol{\theta}_{a0}\|_1 \leq \|\boldsymbol{\theta}_a\|_1 + \|\boldsymbol{\theta}_{a0}\|_1 \leq 2s$ , by Definition 2.3, we have that  $\mathbf{z} \in MS_{r,2s,\mu,l}$ . Therefore,  $0 = \|\mathbf{A}\mathbf{z}\|_2 \geq (1 - \delta_{r,2s,\mu,l})\|\mathbf{z}\|_2 > 0$ , which is a contradiction. Notice that the proof is inspired by the proof of lemma 3.1 in Candes and Tao (2005).

## Appendix C

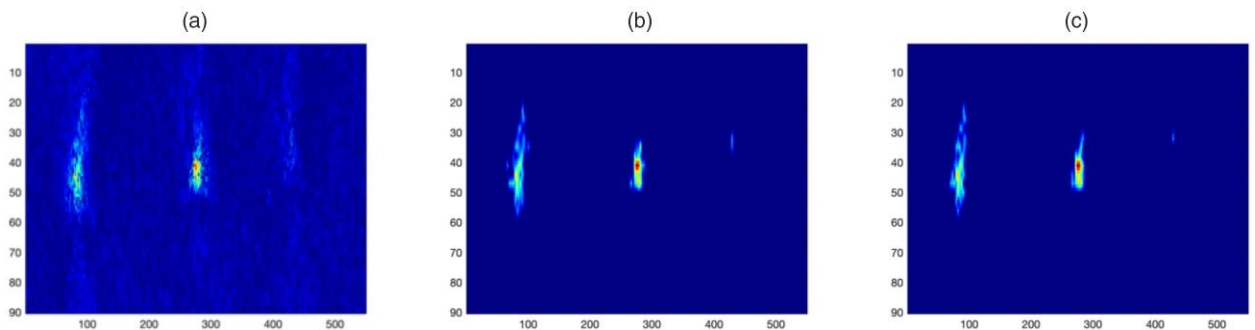
**Proof of Theorem 2.3.** This proof is inspired by Baraniuk et al. (2008) and Tanner and Vary (2020).

We will first derive the RIC for a fixed subspace  $MS_{r,\mathbf{T},\mu,l}$  of  $MS_{r,s,\mu,l}$  when  $\boldsymbol{\theta}_a$  is restricted in a fixed subspace  $\mathbf{T}$  with the fixed support such that the number of nonzero elements is:

$$MS_{r,\mathbf{T},\mu,l} = \{\mathbf{y} \in \mathbb{R}^n \mid \mathbf{y} = \mathbf{B}\boldsymbol{\theta} + \mathbf{B}_a\boldsymbol{\theta}_a, \mathbf{B}_a \in \mathbb{R}^{n \times q}, l(\mathbf{B}_a) = l, \boldsymbol{\theta}_a \in \mathbf{T}, \mathbf{B} \in \mathbb{R}^{n \times r}, \mu(\mathbf{B}) = \mu < n(2rsl)^{-1}, \boldsymbol{\theta} \in \mathbb{R}^r\}.$$

Then, we use a covering argument that counts over all possible sparse subspaces  $\mathbf{T}$  with support less than or equal to  $s$ . Finally, we can derive the RIC for  $MS_{r,s,\mu,l}$ .

**Figure 12.** (Color online) Silicon Stress Map and Detected Indentation



Notes. (a) Raw image. (b) KronCSSD result. (c) SSD result.

The following lemma describes the RIC for a fixed subspace  $MS_{r,T,\mu,l}$  and is proved in Appendix E.

**Lemma C.1.** *RIC for a fixed subspace  $MS_{r,T,\mu,l}$ . Let  $A \in \mathbb{R}^{p \times n}$  be a matrix from the families described in Definition 2.5. Furthermore, assume that  $\mu < n(2rsl)^{-1}$  and that the basis matrix  $B_a$  for the sparse signal component satisfies the RIP with RIC  $\delta_{B_a,s} \in (0, 1)$ : that is,  $\delta_{B_a,s}$  is the smallest positive constant such that*

$$(1 - \delta_{B_a,s})\|\theta_a\|_2 \leq \|B_a\theta_a\|_2 \leq (1 + \delta_{B_a,s})\|\theta_a\|_2, \\ \forall \theta_a \in \{\theta_a \in \mathbb{R}^q \mid \|\theta_a\|_0 \leq s\}.$$

For a given  $\delta \in (0, 1)$ , there exists a constant  $c_0 > 0$  depending only on  $\delta$ , such that the RIC for  $MS_{r,T,\mu,l}$  is upper bounded by  $\delta$  with the probability of at least  $1 - 2\left(\frac{24}{\delta}\tau_1\right)^r \left(\frac{24}{\delta}\tau_0\right)^s e^{-pc_0(\delta/2)}$ , where  $\eta = \sqrt{\frac{ursl}{n}}$ ,  $\tau_0 = \frac{1}{\sqrt{(1-\delta_{B_a,s})(1-\eta^2)}}$ ,  $\tau_1 = \|B^{\dagger}\|_2 \left(1 + \frac{1}{\sqrt{1-\eta^2}}\right)$ , and  $B^{\dagger} = (B^T B)^{-1} B^T$ .

Notice that for a fixed subspace  $MS_{r,T,\mu,l}$ , the RIP will fail with probability less than or equal to  $2\left(\frac{24}{\delta}\tau_1\right)^r \left(\frac{24}{\delta}\tau_0\right)^s e^{-pc_0(\delta/2)}$ . Because there are  $\binom{n}{s} \leq \left(\frac{en}{s}\right)^s$  such subspaces, the probability to fail for  $MS_{r,s,\mu,l}$ , which is a combination of those  $\binom{n}{s}$  subspaces, will be less than or equal to

$$\binom{n}{s} 2\left(\frac{24}{\delta}\tau_1\right)^r \left(\frac{24}{\delta}\tau_0\right)^s e^{-pc_0(\frac{\delta}{2})} \\ \leq \exp\left(-c_0\left(\frac{\delta}{2}\right)p + \ln 2 + r \ln \frac{24}{\delta}\tau_1 + s\left(1 + \ln \frac{24}{\delta}\tau_0 + \ln \frac{n}{s}\right)\right).$$

Then, for any give  $\delta$ , there exist  $c_1, c_2 > 0$ , such that the probability to fail for  $MS_{r,s,\mu,l}$  is less than or equal to  $\exp(-c_1 p)$ , provided that  $p \geq c_2(\ln 2 + r \ln \frac{24}{\delta}\tau_1 + s(1 + \ln \frac{24}{\delta}\tau_0 + \ln \frac{n}{s}))$ , where  $c_2 = [c_0(\frac{\delta}{2}) - c_1]^{-1}$ . This finishes the proof.

## Appendix D

**Proof of Theorem 2.4.** The proof of Theorem 2.4 is inspired by Candes et al. (2006) and Tanner and Vary (2020). Assume that in Problem (1),  $\epsilon_1$  is properly chosen such that Problem (1) is feasible. In the following discussion, we will use  $(\cdot)^*$  to denote the optimal solution of Problem (1) and  $(\cdot)_0$  to denote the signal we wish to recover. Let  $R = X^* - X_0 = R^m + R^a$ , where  $R^m = m - m_0 = B\theta^* - B\theta_0$  and  $R^a = B_a\theta_a^* - B_a\theta_{a0}$  are the residuals of the smooth and sparse signal components, respectively.

Let  $h = \theta_a^* - \theta_{a0} = h_{T_0} + h_{T_0^c}$ , where  $T_0$  is the support of  $\theta_{a0}$ ,  $h_{T_0}$  denotes the projection of  $h$  onto  $T_0$  such that

$$h_{T_0}(t) = \begin{cases} t, & \text{if } t \in T_0 \\ 0, & \text{otherwise,} \end{cases}$$

and  $T_0^c$  denotes the complementary set of  $T_0$ .

Because  $\theta_{a0}$  is feasible and  $\theta_a^*$  is the optimal solution of Problem (1), we must have  $\|\theta_a^*\|_1 \leq \|\theta_{a0}\|_1$ , which is equivalent to  $\|\theta_{a0} + h_{T_0} + h_{T_0^c}\|_1 \leq \|\theta_{a0}\|_1$ . Because  $T_0$  and  $T_0^c$  are complementary to each other, we have  $\|\theta_{a0} + h_{T_0}\|_1 + \|h_{T_0^c}\|_1 \leq \|\theta_{a0}\|_1$ . Because  $\|\theta_{a0} + h_{T_0}\|_1 \geq \|\theta_{a0}\|_1 - \|h_{T_0}\|_1$ , we have  $\|\theta_{a0}\|_1 - \|h_{T_0}\|_1$

+  $\|h_{T_0^c}\|_1 \leq \|\theta_{a0}\|_1$ . Hence,  $\|h_{T_0^c}\|_1 \leq \|h_{T_0}\|_1$ . Because  $\|h_{T_0}\|_1 \leq \sqrt{s}\|h_{T_0}\|_2$ , we have

$$\|h_{T_0^c}\|_1 \leq \sqrt{s}\|h_{T_0}\|_2. \quad (D.1)$$

Similar to Candes et al. (2006), we order the elements of  $T_0^c$  in decreasing order of their magnitude and enumerate  $T_0^c$  as  $v_1, \dots, v_{n-|T_0|}$ . Then,  $T_0^c$  is divided into subsets  $T_i^c$  of size  $M$ , where

$$T_i^c = \{v_j : (i-1)M \leq j \leq iM\}.$$

Let  $h_{T_i^c}$  be the projection of  $h$  onto  $T_i^c$ ; we have

$$\|h_{T_i^c}\|_0 \leq M, \quad \forall i \geq 1 \\ T_i^c \cap T_j^c = \emptyset, \quad \forall i \neq j \\ \|h_{T_{i+1}^c}\|_2 \leq \frac{1}{\sqrt{M}}\|h_{T_i^c}\|_1, \quad \forall i \geq 1, \quad (D.2)$$

where the last inequality comes from the fact that  $T_0^c$  is in decreasing order, such that

$$\left|h_{T_{i+1}^c}\right|_{(v)} \leq \frac{1}{M} \sum_{j \in T_i^c} \left|h_{T_j^c}\right|_{(j)}. \quad \forall v \in T_{i+1}^c.$$

Define  $R_{T_i^c}^a = B_a h_{T_i^c}$  and  $R_{T_0}^a = B_a h_{T_0}$ , and combine Equations (D.1) and (D.2). Then, we have

$$\sum_{j \geq 2} \|R_{T_j^c}^a\|_2 \leq \sum_{j \geq 2} \sqrt{1 + \delta_{B_a, M}} \|h_{T_j^c}\|_2 \\ \leq (a) \sum_{j \geq 1} \frac{\sqrt{1 + \delta_{B_a, M}} \|h_{T_j^c}\|_1}{\sqrt{M}} = \frac{\sqrt{1 + \delta_{B_a, M}} \|h_{T_0^c}\|_1}{\sqrt{M}} \\ \leq (b) \frac{\sqrt{s} \sqrt{1 + \delta_{B_a, M}} \|h_{T_0}\|_2}{\sqrt{M}} \leq \frac{\sqrt{s} \sqrt{1 + \delta_{B_a, M}} \|R_{T_0}^a\|_2}{\sqrt{M} \sqrt{1 - \delta_{B_a, M}}},$$

where (a) follows Equation (D.2), (b) follows Equation (D.1), and the RIP property of  $B_a$  is used because  $\|h_{T_j^c}\|_0 \leq M$ .

Denote  $\gamma = \sqrt{\frac{1 + \delta_{B_a, M+s}}{1 - \delta_{B_a, M+s}}}$ . (A tighter bound can be achieved by using  $\sqrt{\frac{1 + \delta_{B_a, M}}{1 - \delta_{B_a, M}}}$ . However, we adopt  $\delta_{B_a, M+s}$  instead of  $\delta_{B_a, M}$  for simplicity in the following proof.) Then, we have

$$\sum_{j \geq 2} \|R_{T_j^c}^a\|_2 \leq \sqrt{\frac{s}{M}} \gamma \|R_{T_0}^a\|_2. \quad (D.3)$$

Next, we derive the bounds for  $R^m$  and  $R^a$ .

Bound for  $R^m$ :

$$\|AR^m\|_2^2 = |\langle AR^m, A(R - R^a) \rangle| \\ = |\langle AR^m, A(R - R^a) \rangle| \\ = |\langle AR^m, AR \rangle + \langle AR^m, -AR^a \rangle| \\ \leq |\langle AR^m, AR \rangle| + |\langle AR^m, -AR^a \rangle| \\ = |\langle AR^m, AR \rangle| + \left| \left\langle AR^m, -A \left( R_{T_0}^a + R_{T_1}^a + \sum_{j \geq 2} R_{T_j}^a \right) \right\rangle \right| \\ \leq |\langle AR^m, AR \rangle| + \left| \left\langle AR^m, A \left( R_{T_0}^a + R_{T_1}^a \right) \right\rangle \right| \\ + \sum_{j \geq 2} \left| \left\langle AR^m, AR_{T_j}^a \right\rangle \right|. \quad (D.4)$$

In the following discussion, we will bound those terms. According to the Cauchy-Schwarz inequality, the first term

can be bounded as

$$|\langle \mathbf{AR}^m, \mathbf{AR} \rangle| \leq \|\mathbf{AR}^m\|_2 \|\mathbf{AR}\|_2 \leq \sqrt{1 + \delta_{r,s,\mu,l}} \epsilon_1 \|\mathbf{R}^m\|_2, \quad (\text{D.5})$$

where the last inequality comes from the RIP property and the first constraint in Problem (1).

The third term can be bounded as follows. Denoting  $\mathbf{z}_1 = \mathbf{R}^m / \|\mathbf{R}^m\|_2$  and  $\mathbf{z}_2 = \mathbf{R}_{T_j^c}^a / \|\mathbf{R}_{T_j^c}^a\|_2$ , we have

$$\begin{aligned} & \frac{|\langle \mathbf{AR}^m, \mathbf{AR}_{T_j^c}^a \rangle|}{\|\mathbf{R}^m\|_2 \|\mathbf{R}_{T_j^c}^a\|_2} \\ &= |\langle \mathbf{Az}_1, \mathbf{Az}_2 \rangle| = \frac{1}{4} \left| \|\mathbf{A}(\mathbf{z}_1 + \mathbf{z}_2)\|_2^2 - \|\mathbf{A}(\mathbf{z}_1 - \mathbf{z}_2)\|_2^2 \right| \\ &\leq_{(a)} \frac{1}{4} \max \left\{ \begin{aligned} & |(1 + \delta_{r,M,\mu,l})\|\mathbf{z}_1 + \mathbf{z}_2\|_2^2 - (1 - \delta_{r,M,\mu,l})\|\mathbf{z}_1 - \mathbf{z}_2\|_2^2|, \\ & |(1 + \delta_{r,M,\mu,l})\|\mathbf{z}_1 - \mathbf{z}_2\|_2^2 - (1 - \delta_{r,M,\mu,l})\|\mathbf{z}_1 + \mathbf{z}_2\|_2^2| \end{aligned} \right\} \\ &= |\delta_{r,M,\mu,l} + \langle \mathbf{z}_1, \mathbf{z}_2 \rangle| \leq_{(b)} \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2}, \end{aligned}$$

where  $\eta_1 = \sqrt{\frac{\mu r M}{n}}$ , inequality (a) follows the RIP property because  $\mathbf{z}_1 + \mathbf{z}_2 \in MS_{r,M,\mu,l}$  and  $\mathbf{z}_1 - \mathbf{z}_2 \in MS_{r,M,\mu,l}$ ; and inequality (b) follows Equation (F.1) in the proof of Lemma E.1 and  $\|\mathbf{z}_1\|_2 = \|\mathbf{z}_2\|_2 = 1$ ,  $\langle \mathbf{z}_1, \mathbf{z}_2 \rangle \leq \frac{\eta_1}{1 - \eta_1^2} \|\mathbf{z}_1\|_2 \|\mathbf{z}_2\|_2 = \frac{\eta_1}{1 - \eta_1^2}$ , provided that  $M \leq 2s$ .

Therefore,

$$|\langle \mathbf{AR}^m, \mathbf{AR}_{T_j^c}^a \rangle| \leq \left( \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \|\mathbf{R}^m\|_2 \|\mathbf{R}_{T_j^c}^a\|_2. \quad (\text{D.6})$$

Similarly, the second term can be bounded as

$$\begin{aligned} & |\langle \mathbf{AR}^m, \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a) \rangle| \\ &\leq \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2, \end{aligned} \quad (\text{D.7})$$

where  $\eta_2 = \sqrt{\frac{\mu r (M+s)l}{n}}$ , provided that  $M \leq s$ .

Plugging Equations (D.5)–(D.7) into Equation (D.4), we have

$$\begin{aligned} & \|\mathbf{AR}^m\|_2^2 \\ &\leq \|\mathbf{R}^m\|_2 \left( \sqrt{1 + \delta_{r,s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2 \right. \\ &\quad \left. + \left( \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \|\mathbf{R}_{T_j^c}^a\|_2 \right) \\ &\leq_{(a)} \|\mathbf{R}^m\|_2 \left( \sqrt{1 + \delta_{r,s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2 \right. \\ &\quad \left. + \left( \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \sqrt{\frac{s}{M}} \gamma \|\mathbf{R}_{T_0}^a\|_2 \right), \end{aligned}$$

where inequality (a) follows Equation (D.3).

According to the RIP property, we have

$$\begin{aligned} & (1 - \delta_{r,s,\mu,l}) \|\mathbf{R}^m\|_2^2 \\ &\leq \|\mathbf{R}^m\|_2 \left( \sqrt{1 + \delta_{r,s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2 \right. \\ &\quad \left. + \left( \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \sqrt{\frac{s}{M}} \gamma \|\mathbf{R}_{T_0}^a\|_2 \right). \end{aligned}$$

Consequently, we have

$$\|\mathbf{R}^m\|_2 \leq \frac{\sqrt{1 + \delta_{r,s,\mu,l}} \epsilon_1 + \left( \delta_{r,M,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \sqrt{\frac{s}{M}} \gamma^2 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2}{(1 - \delta_{r,s,\mu,l})}, \quad (\text{D.8})$$

where the inequality follows from

$$\begin{aligned} \|\mathbf{R}_{T_0}^a\|_2 &= \|\mathbf{B}_a \mathbf{h}_{T_0}\|_2 \leq_{(a)} \sqrt{1 + \delta_{B_a,s}} \|\mathbf{h}_{T_0}\|_2 \leq_{(b)} \sqrt{1 + \delta_{B_a,s}} \|\mathbf{h}_{T_0} + \mathbf{h}_{T_1^c}\|_2 \\ &\leq_{(c)} \sqrt{\frac{1 + \delta_{B_a,s}}{1 - \delta_{B_a,M+s}}} \|\mathbf{B}_a (\mathbf{h}_{T_0} + \mathbf{h}_{T_1^c})\|_2 \leq \gamma \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2, \end{aligned} \quad (\text{D.9})$$

inequalities (a) and (c) follow the RIP property of  $\mathbf{B}_a$ , and inequality (b) follows that  $\mathbf{T}_0 \cap \mathbf{T}_1^c = \emptyset$ .

Bound for  $\mathbf{R}^a$ :

$$\begin{aligned} & \|\mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a)\|_2^2 = |\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a - \mathbf{R} + \mathbf{R}) \rangle| \\ &= |\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), \mathbf{AR} \rangle| + \left| \left\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), -\mathbf{A} \left( \mathbf{R}^m + \sum_{j \geq 2} \mathbf{R}_{T_j^c}^a \right) \right\rangle \right| \\ &\leq |\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), \mathbf{AR} \rangle| + |\langle \mathbf{AR}^m, \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a) \rangle| \\ &\quad + \sum_{j \geq 2} |\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), \mathbf{AR}_{T_j^c}^a \rangle|. \end{aligned} \quad (\text{D.10})$$

In the following discussion, we will bound those terms. According to the Cauchy–Schwarz inequality, the first term can be bounded as

$$\begin{aligned} & |\langle \mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a), \mathbf{AR} \rangle| \leq \|\mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a)\|_2 \|\mathbf{AR}\|_2 \\ &\leq \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2, \end{aligned} \quad (\text{D.11})$$

where the last inequality comes from the RIP property and the first constraint in Problem (1).

The third term can be bounded as follows. Denoting  $\mathbf{z}_2 = \mathbf{R}_{T_j^c}^a / \|\mathbf{R}_{T_j^c}^a\|_2$ ,  $j \geq 2$  and  $\mathbf{z}_3 = (\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a) / \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1^c}^a\|_2$ , we

have

$$\begin{aligned} & \left| \frac{\langle A(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a), \mathbf{A}\mathbf{R}_{T_1}^a \rangle}{\|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \|\mathbf{R}_{T_1}^a\|_2} \right| = |\langle \mathbf{A}\mathbf{z}_3, \mathbf{A}\mathbf{z}_2 \rangle| \\ &= \frac{1}{4} \left| \|\mathbf{A}(\mathbf{z}_3 + \mathbf{z}_2)\|_2^2 - \|\mathbf{A}(\mathbf{z}_3 - \mathbf{z}_2)\|_2^2 \right| \\ &\stackrel{(a)}{\leq} \frac{1}{4} \max \left\{ \left| (1 + \delta_{r,2M+s,\mu,l}) \|\mathbf{z}_3 + \mathbf{z}_2\|_2^2 - (1 - \delta_{r,2M+s,\mu,l}) \|\mathbf{z}_3 - \mathbf{z}_2\|_2^2 \right|, \right. \\ &\quad \left. \left| (1 + \delta_{r,2M+s,\mu,l}) \|\mathbf{z}_3 - \mathbf{z}_2\|_2^2 - (1 - \delta_{r,2M+s,\mu,l}) \|\mathbf{z}_3 + \mathbf{z}_2\|_2^2 \right| \right\} \\ &= |\delta_{r,2M+s,\mu,l} \langle \mathbf{z}_3, \mathbf{z}_2 \rangle| \\ &\stackrel{(b)}{\leq} \delta_{r,2M+s,\mu,l} \end{aligned}$$

where inequality (a) follows the RIP property because  $\mathbf{z}_3 + \mathbf{z}_2 \in MS_{r,2M+s,\mu,l}$  and  $\mathbf{z}_3 - \mathbf{z}_2 \in MS_{r,2M+s,\mu,l}$ . Inequality (b) comes from that  $\mathbf{T}_i^c \cap \mathbf{T}_j^c = \emptyset$ ,  $\forall i \neq j$  and  $\mathbf{T}_0 \cap \mathbf{T}_i^c = \emptyset$ ,  $\forall i$ .

Therefore,

$$\left| \langle \mathbf{A}\mathbf{R}^m, \mathbf{A}\mathbf{R}_{T_1}^a \rangle \right| \leq \delta_{r,2M+s,\mu,l} \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \|\mathbf{R}_{T_1}^a\|_2. \quad (\text{D.12})$$

Plugging Equations (D.7), (D.11), and (D.12) into Equation (D.10), we have

$$\begin{aligned} & \|\mathbf{A}(\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a)\|_2^2 \\ &\leq \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \left( \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \right. \\ &\quad \left. + \delta_{r,2M+s,\mu,l} \sum_{j \geq 2} \|\mathbf{R}_{T_j}^a\|_2 \right) \\ &\stackrel{(a)}{\leq} \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \left( \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \right. \\ &\quad \left. + \delta_{r,2M+s,\mu,l} \sqrt{\frac{s}{M}} \gamma \|\mathbf{R}_{T_0}^a\|_2 \right) \\ &\stackrel{(b)}{\leq} \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \left( \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \right. \\ &\quad \left. + \delta_{r,2M+s,\mu,l} \sqrt{\frac{s}{M}} \gamma^2 \|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2 \right), \end{aligned}$$

where inequalities (a) and (b) follows the same argument as deriving Equation (D.8).

According to the RIP property, we have

$$\begin{aligned} & (1 - \delta_{r,M+s,\mu,l}) \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2^2 \\ &\leq \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2 \left( \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \right. \\ &\quad \left. + \delta_{r,2M+s,\mu,l} \sqrt{\frac{s}{M}} \gamma^2 \|\mathbf{R}_{T_0}^a\|_2 \right). \end{aligned}$$

Consequently, we have

$$\begin{aligned} & \|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2 \\ &\quad \left( \sqrt{1 + \delta_{r,M+s,\mu,l}} \epsilon_1 + \left( \delta_{r,M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right) \|\mathbf{R}^m\|_2 \right. \\ &\quad \left. + \delta_{r,2M+s,\mu,l} \sqrt{\frac{s}{M}} \gamma^2 \|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2 \right) \\ &\leq \frac{\|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2}{(1 - \delta_{r,M+s,\mu,l})}. \quad (\text{D.13}) \end{aligned}$$

Notice that Equations (D.8) and (D.13) still hold if we relax  $\delta_{r,M+s,\mu,l}$ ,  $\delta_{r,s,\mu,l}$  to  $\delta_{r,2M+s,\mu,l}$ . For simplicity, here we replace  $\delta_{r,M+s,\mu,l}$ ,  $\delta_{r,s,\mu,l}$  with  $\delta_{r,2M+s,\mu,l}$  in the following derivation.

Plugging Equation (D.8) into Equation (D.13) and letting  $x \equiv \|\mathbf{R}_{T_0}^a + \mathbf{R}_{T_1}^a\|_2$ ,  $y \equiv \|\mathbf{R}^m\|_2$ , we have

$$\left( D_1 - \frac{B_1 B_2}{D_2} - C_1 \right) x \leq A_1 \epsilon_1 + \frac{B_1}{D_2} A_2 \epsilon_1, \quad (\text{D.14})$$

where

$$A_1 = \sqrt{1 + \delta_{r,2M+s,\mu,l}}, \quad B_1 = \left( \delta_{r,2M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right),$$

$$C_1 = \delta_{r,2M+s,\mu,l} \sqrt{\frac{s}{M}} \gamma^2, \quad D_1 = (1 - \delta_{r,2M+s,\mu,l}),$$

and

$$A_2 = \sqrt{1 + \delta_{r,2M+s,\mu,l}},$$

$$B_2 = \left( \delta_{r,2M+s,\mu,l} + \frac{\eta_1}{1 - \eta_1^2} \right) \sqrt{\frac{s}{M}} \gamma^2 + \left( \delta_{r,2M+s,\mu,l} + \frac{\eta_2}{1 - \eta_2^2} \right),$$

$$D_2 = (1 - \delta_{r,2M+s,\mu,l}).$$

Here, we require  $D_1 - \frac{B_1 B_2}{D_2} - C_1 > 0$  in Equation (D.14) and let  $M = s$ , which is

$$1 - \gamma^2 \alpha_1 \alpha_2 - \alpha_2^2 - ((1 + \alpha_1 + \alpha_2) \gamma^2 + 2\alpha_2 + 2) \delta_{r,3s,\mu,l} > 0. \quad (\text{D.15})$$

Let  $a = (1 + \alpha_1 + \alpha_2) \gamma^2 + 2\alpha_2 + 2$  and  $c = 1 - \gamma^2 \alpha_1 \alpha_2 - \alpha_2^2$ , where  $\alpha_1 = \frac{\eta}{1 - \eta^2}$ ,  $\alpha_2 = \frac{\sqrt{2}\eta}{1 - 2\eta^2}$ ,  $\gamma = \sqrt{\frac{1 + \delta_{R_a,2s}}{1 - \delta_{R_a,2s'}}$ , and  $\eta = \sqrt{\frac{\mu r s l}{n}}$ . If  $c > 0$ , then there exist a  $\delta_{r,3s,\mu,l} > 0$  such that Equation (D.15) is valid. Then, the denominator  $-a\delta_{r,3s,\mu,l} + c > 0 \forall \delta_{r,3s,\mu,l} \in (0, c/a)$ . Consequently,

$$\|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2 \leq \frac{(1 + \alpha_2) \sqrt{1 + \delta_{r,3s,\mu,l}}}{c - a\delta_{r,3s,\mu,l}} \epsilon_1.$$

Notice that from Equations (D.1) and (D.9), we have

$$\sum_{j \geq 2} \|\mathbf{R}_{T_j}^a\|_2 \leq \gamma \|\mathbf{R}_{T_0}^a\|_2 \leq \gamma^2 \|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2.$$

Therefore,  $\|\mathbf{R}^a\|_2 \leq \|\mathbf{R}_{T_0}^a\|_2 + \|\mathbf{R}_{T_1}^a\|_2 + \sum_{j \geq 2} \|\mathbf{R}_{T_j}^a\|_2 \leq C_a \epsilon_1$ , where

$$C_a = \frac{(1 + \gamma^2)(1 + \alpha_2) \sqrt{1 + \delta_{r,3s,\mu,l}}}{c - a\delta_{r,3s,\mu,l}}.$$

Similarly, we can bound  $\|\mathbf{R}^m\|_2$  as  $\|\mathbf{R}^m\|_2 \leq C_m \epsilon_1$ , where

$$C_m = \frac{\sqrt{1 + \delta_{r,3s,\mu,l}} + \left( \delta_{r,3s,\mu,l} + \frac{\gamma^2}{1 + \gamma^2} \alpha_1 + \frac{1}{1 + \gamma^2} \alpha_2 \right) C_a}{(1 - \delta_{r,3s,\mu,l})}. \quad \square$$

## Appendix E

**Proof of Lemma C.1.** In this section, we will provide the proof for Lemma C.1. By linearity of the measurement matrix  $\mathbf{A}$ , without loss of generality, it is enough to prove this lemma when  $\|\mathbf{y}_2\| = 1$ . The proof mainly has two steps. First, the bounds for  $\theta_a$  and  $\theta$  are derived, and a finite set of points to approximate the set  $MS_{r,\mathbf{T},\mu,l}$  to any accuracy in

norm 2 sense can be found. Then, the concentration inequality can be applied through a union bound. This is a common approach in compressive sensing literature (Baraniuk et al. 2008, Tanner and Vary 2020).

To derive the upper bounds for  $\theta_a$  and  $\theta$ , we first derive the upper bounds for the signals  $m = B\theta$  and  $a = B_a\theta_a$ , which are given in the following lemma.

**Lemma E.1.** *The smooth signal component  $m$  and sparse signal component  $a$  of the signal  $y$  in  $MS_{r,T,\mu,l}$  with  $\mu < \frac{n}{2rs_l}$  can be bounded as follows:*

$$\|m\|_2 = \|B\theta\|_2 \leq \frac{\|y\|_2}{\sqrt{1-\eta^2}}, \quad (\text{E.1})$$

$$\|a\|_2 = \|B_a\theta_a\|_2 \leq \frac{\|y\|_2}{\sqrt{1-\eta^2}}, \quad (\text{E.2})$$

where  $\eta = \sqrt{\frac{\mu r s_l}{n}}$ .

The proof is presented in Appendix F. According to the RIP for  $B_a$ , we have

$$\sqrt{(1-\delta_{B_a,s})}\|\theta_a\|_2 \leq \|B_a\theta_a\|_2 \leq \sqrt{(1+\delta_{B_a,s})}\|\theta_a\|_2. \quad (\text{E.3})$$

Combining Equations (E.2) and (E.3), we have

$$\begin{aligned} \|\theta_a\|_2 &\leq \frac{1}{\sqrt{(1-\delta_{B_a,s})}}\|B_a\theta_a\|_2 \leq \frac{\|y\|_2}{\sqrt{(1-\delta_{B_a,s})(1-\eta^2)}} \\ &= \frac{1}{\sqrt{(1-\delta_{B_a,s})(1-\eta^2)}}. \end{aligned}$$

Denoting  $\tau_0 = \frac{1}{\sqrt{(1-\delta_{B_a,s})(1-\eta^2)}}$ , we have  $\|\theta_a\|_2 \leq \tau_0$ . Recall that

$y = B\theta + B_a\theta_a$ ; we have  $\theta = B^\dagger(y - B_a\theta_a)$ , where  $B^\dagger = (B^T B)^{-1} B^T$ . Therefore, according to the triangle inequality and the Cauchy–Schwarz inequality, we have

$$\|\theta\|_2 \leq \|B^\dagger\|_2(\|y\|_2 + \|B_a\theta_a\|_2) \leq \|B^\dagger\|_2 \left(1 + \frac{1}{\sqrt{1-\eta^2}}\right).$$

Denoting  $\tau_1 = \|B^\dagger\|_2 \left(1 + \frac{1}{\sqrt{1-\eta^2}}\right)$ , we have  $\|\theta\|_2 \leq \tau_1$ .

Because we have derived the bounds for  $\theta_a$  and  $\theta$ , the covering number of  $MS_{r,T,\mu,l}$  is given by the following lemma whose proof is in Appendix G.

**Lemma E.2.** *There exists a set  $Q \in MS_{r,T,\mu,l}$ , such that for all  $y \in MS_{r,T,\mu,l}$  with  $\|y\|_2 = 1$  we have  $\min_{q \in Q} \|q - y\|_2 \leq \frac{\delta}{4}$  and  $|Q| \leq \left(\frac{24}{\delta}\tau_1\right)^r \left(\frac{24}{\delta}\tau_0\right)^s$ , where  $|Q|$  is its cardinality.*

Next, we will prove the main result by applying the concentration inequality (Definition 2.5(ii)) with union bound. Let  $\epsilon = \delta/2$ ,

$$\left(1 - \frac{\delta}{2}\right)\|q\|_2^2 \leq \|Aq\|_2^2 \leq \left(1 + \frac{\delta}{2}\right)\|q\|_2^2 \quad \forall q \in Q, \quad (\text{E.4})$$

with probability greater than  $1 - 2|Q|e^{-pc_0(\delta/2)}$ .

Because  $\delta \in (0,1)$ , we have  $1 - \frac{\delta}{2} \leq \sqrt{1 - \frac{\delta}{2}}$  and  $\sqrt{1 + \frac{\delta}{2}} \leq 1 + \frac{\delta}{2}$ . Then, Equation (E.4) can be written as

$$\left(1 - \frac{\delta}{2}\right)\|q\|_2 \leq \|Aq\|_2 \leq \left(1 + \frac{\delta}{2}\right)\|q\|_2 \quad \forall q \in Q, \quad (\text{E.5})$$

with probability greater than  $1 - 2|Q|e^{-pc_0(\delta/2)}$ .

By the triangle inequality, we have

$$\|Ay\|_2 \leq \|A(y - q)\|_2 + \|Aq\|_2. \quad (\text{E.6})$$

Define

$$U = \max_{y \in MS_{r,T,\mu,l}, \|y\|_2=1} \|Ay\|_2, \quad (\text{E.7})$$

which is attainable because  $MS_{r,T,\mu,l}$  is closed.

Combining Equations (E.5) and (E.6), we have  $\forall y \in MS_{r,T,\mu,l}$  with  $\|y\|_2 = 1$ ; there exists a  $q \in MS_{r,T,\mu,l}$ , such that  $\|Ay\|_2 \leq \|A(q - y)\|_2 + (1 + \frac{\delta}{2})\|q\|_2$ , with probability greater than  $1 - 2|Q|e^{-pc_0(\delta/2)}$ .

Because  $q - y \in MS_{r,T,\mu,l}$ , if  $q - y = 0$ , we have  $\|Ay\|_2 \leq (1 + \frac{\delta}{2})\|q\|_2 = 1 + \frac{\delta}{2}$ .

If  $q - y \neq 0$ , we have  $\|Ay\|_2 \leq \left\|A \frac{(q-y)}{\|q-y\|_2}\right\|_2 \|q - y\|_2 + (1 + \frac{\delta}{2})\|q\|_2 \leq \frac{\delta}{4}U + 1 + \frac{\delta}{2}$ .

Notice that the second inequality comes from Equation (E.7) combined with  $Q$  being a  $\frac{\delta}{4}$  covering of  $MS_{r,T,\mu,l}$ . In summary, we have  $\|Ay\|_2 \leq \frac{\delta}{4}U + 1 + \frac{\delta}{2}$ .

Because  $U$  is attainable, according to Equation (E.7), we have  $U \leq \frac{\delta}{4}U + 1 + \frac{\delta}{2}$ . Consequently, we have  $U \leq 1 + \frac{3}{4-\delta}\delta \leq 1 + \delta$  because  $\delta < 1$ . Therefore,  $\|Ay\|_2 \leq 1 + \delta$ , with probability greater than  $1 - 2|Q|e^{-pc_0(\delta/2)}$ .

Similarly, we can prove that  $\|Ay\|_2 \geq 1 - \delta$  with probability greater than  $1 - 2|Q|e^{-pc_0(\delta/2)}$ .

Finally, according to Lemma E.2, we have that

$$1 - 2|Q|e^{-pc_0(\frac{\delta}{2})} \geq 1 - 2\left(\frac{24}{\delta}\tau_1\right)^r \left(\frac{24}{\delta}\tau_0\right)^s e^{-pc_0(\frac{\delta}{2})}.$$

This finishes the proof.

## Appendix F

**Proof of Lemma E.1.** To prove the result, we first derive a nontrivial upper bound for the inner produce between  $m$  and  $a$ . Let  $B = U\Sigma V^T$  be the reduced SVD of  $B$ ; then,

$$\begin{aligned} |m^T a| &= |\theta^T B^T a| = |\theta^T V \Sigma U^T a| = \left| \theta^T V \Sigma U^T \sum_i^n a_i e_i \right| \\ &= \left| \theta^T V \Sigma \sum_i^n a_i U^T e_i \right| \\ &\leq \|\theta^T V \Sigma\|_2 \left\| \sum_i^n a_i U^T e_i \right\|_2 \\ &\leq \|\theta^T V \Sigma\|_2 \sum_i^n |a_i| \|U^T e_i\|_2 \\ &\leq \|\theta^T V \Sigma\|_2 \sum_i^n |a_i| \max_{j \in \{1, \dots, r\}} \|U^T e_j\|_2 \\ &\stackrel{(a)}{\leq} \|\theta^T V \Sigma U^T\|_2 \|a\|_1 \sqrt{\frac{\mu r}{n}} \\ &\stackrel{(b)}{\leq} \sqrt{\frac{\mu r s_l}{n}} \|m\|_2 \|a\|_2, \end{aligned} \quad (\text{F.1})$$

where inequality (a) follows  $\|\theta^T V \Sigma\|_2 = \|\theta^T V \Sigma U^T\|_2$  because  $U^T U = I$  and Definition 2.2. Inequality (b) follows from  $\|a\|_1 \leq \sqrt{s} \|a\|_2$ .

Let  $\eta = \sqrt{\frac{\mu r s l}{n}}$  because  $\mu < \frac{n}{r s l}$ , we have  $\eta < 1$  and

$$|\mathbf{m}^T \mathbf{a}| = \frac{|\|\mathbf{y}\|_2^2 - \|\mathbf{m}\|_2^2 - \|\mathbf{a}\|_2^2|}{2} \leq \eta \|\mathbf{m}\|_2 a_2.$$

Therefore, we have

$$\|\mathbf{m}\|_2^2 + \|\mathbf{a}\|_2^2 - \|\mathbf{y}\|_2^2 \leq 2\eta \|\mathbf{m}\|_2 \|\mathbf{a}\|_2.$$

By completing the square, we have

$$(\|\mathbf{m}\|_2 + \eta \|\mathbf{a}\|_2)^2 + (1 - \eta^2) \|\mathbf{a}\|_2^2 - \|\mathbf{y}\|_2^2 \leq 0.$$

Because  $(\|\mathbf{m}\|_2 + \eta \|\mathbf{a}\|_2)^2 \geq 0$ , we have

$$\|\mathbf{a}\|_2 \leq \frac{1}{\sqrt{(1 - \eta^2)}} \|\mathbf{y}\|_2.$$

Similarly, we can derive that

$$\|\mathbf{m}\|_2 \leq \frac{1}{\sqrt{(1 - \eta^2)}} \|\mathbf{y}\|_2.$$

This finishes the proof.

## Appendix G

**Proof of Lemma E.2.** We first state results for the covering number of a set (Vershynin 2018). The covering number of a smallest  $\epsilon$  net for a unit  $l_2$  norm ball in  $d$ -dimensional space is  $(3/\epsilon)^d$ .

Let  $\mathbf{M} = \{\mathbf{m} \in \mathbb{R}^n \mid \mathbf{m} = \mathbf{B}\boldsymbol{\theta}, \boldsymbol{\theta} \in \mathbb{R}^r, \|\boldsymbol{\theta}\|_2 \leq \tau_1, \mu(\mathbf{B}) = \mu\}$  and  $\mathbf{S} = \{\mathbf{a} \in \mathbb{R}^n \mid \mathbf{a} = \mathbf{B}_a \boldsymbol{\theta}_a, \boldsymbol{\theta}_a \in \mathbb{T}, \|\boldsymbol{\theta}_a\|_2 \leq \tau_0, l(\mathbf{B}_a) = l\}$ . There exist two finite  $\frac{\delta}{8}$  covering sets of  $\mathbf{M}$  and  $\mathbf{S}$ , which are  $\mathbf{Q}_M \subseteq \mathbf{M}$  and  $\mathbf{Q}_S \subseteq \mathbf{S}$ .

For all  $\mathbf{q}_M \in \mathbf{Q}_M$  and for all  $\mathbf{m} \in \mathbf{M}$ , we have

$$\min_{\mathbf{q}_M \in \mathbf{Q}_M} \|\mathbf{m} - \mathbf{q}_M\|_2 \leq \frac{\delta}{8};$$

for all  $\mathbf{q}_S \in \mathbf{Q}_S$  and for all  $\mathbf{a} \in \mathbf{S}$ , we have

$$\min_{\mathbf{q}_S \in \mathbf{Q}_S} \|\mathbf{a} - \mathbf{q}_S\|_2 \leq \frac{\delta}{8}.$$

Therefore, we have  $|\mathbf{Q}_M| \leq (\frac{24}{\delta} \tau_1)^r$  and  $|\mathbf{Q}_S| \leq (\frac{24}{\delta} \tau_0)^s$ .

Define  $\mathbf{Q}_{MS} = \{\mathbf{q}_M + \mathbf{q}_S \mid \mathbf{q}_M \in \mathbf{Q}_M, \mathbf{q}_S \in \mathbf{Q}_S\} \subseteq MS_{r,T,\mu,l}$ . Then,  $\forall \mathbf{y} \in MS_{r,T,\mu,l}$ , there exists a pair  $\mathbf{q}_{MS} = \mathbf{q}_M + \mathbf{q}_S \in MS_{r,T,\mu,l}$ , such that

$$\|\mathbf{q}_{MS} - \mathbf{y}\|_2 = \|\mathbf{q}_M - \mathbf{m} + \mathbf{q}_S - \mathbf{a}\|_2 \leq \|\mathbf{q}_M - \mathbf{m}\|_2 + \|\mathbf{q}_S - \mathbf{a}\|_2 \leq \frac{\delta}{4}.$$

Therefore,  $\mathbf{Q}_{MS}$  is a  $\delta/4$  covering of  $MS_{r,T,\mu,l}$  and  $|\mathbf{Q}_{MS}| \leq (\frac{24}{\delta} \tau_1)^r (\frac{24}{\delta} \tau_0)^s$ . This finishes the proof.

## Appendix H. Simulation Study by Varying the Magnitude of Sparse Signals

We adopt the same simulation data generation procedure as in Section 3.1 while changing the distribution of elements in  $\boldsymbol{\theta}_a \in \mathbb{R}^q$  such that the nonzero elements follow i.i.d. normal distribution with mean zero and standard deviation  $\sigma_s$  in the range of  $\{0.065, 0.125, 0.25, 0.5\}$ . The example signals and reconstruction performance are shown in Figure H.1.

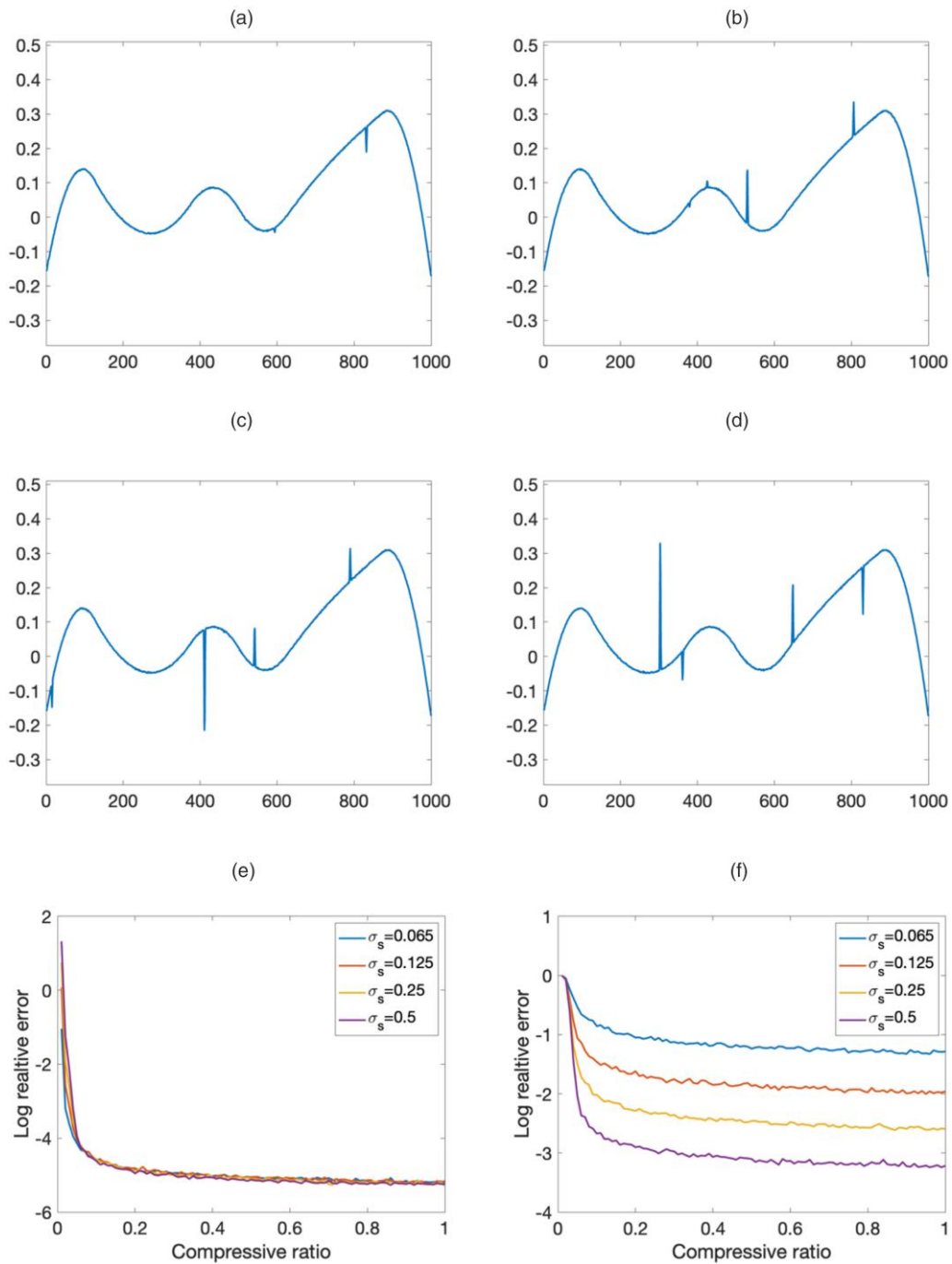
There are several observations.

1. The log relative error of the smooth component does not change with different magnitudes of  $\sigma_s$ . This agrees with the theoretical result in Theorem 2.4, where the reconstruction error is bounded by a constant times the noise bound because the noise level is kept the same in all simulations.

2. The log relative error of the sparse component decreases as  $\sigma_s$  increases. This also agrees with the theoretical result in Theorem 2.4. Because the log relative error of the sparse component is defined as  $\log \|\mathbf{a} - \hat{\mathbf{a}}\|_2 / \|\mathbf{a}\|_2$ , according to Theorem 2.4, the reconstruction error term can be approximated by  $C_a \epsilon_1$  (i.e.,  $\|\mathbf{a} - \hat{\mathbf{a}}\|_2 \sim C_a \epsilon_1$ , where  $C_a$  is independent of  $\boldsymbol{\theta}_a$ ).  $\|\mathbf{a}\|_2 = \|\mathbf{B}_a \boldsymbol{\theta}_a\|_2$  increases as the magnitude of elements in  $\boldsymbol{\theta}_a$  increases. Therefore,  $\log \|\mathbf{a} - \hat{\mathbf{a}}\|_2 / \|\mathbf{a}\|_2$  will decrease as  $\sigma_s$  increases.

3. The 0.1 threshold (approximately) of the compressive ratio still holds with different magnitudes of sparse signal component.

Figure H.1. (Color online) Simulation Study by Varying  $\sigma_s$



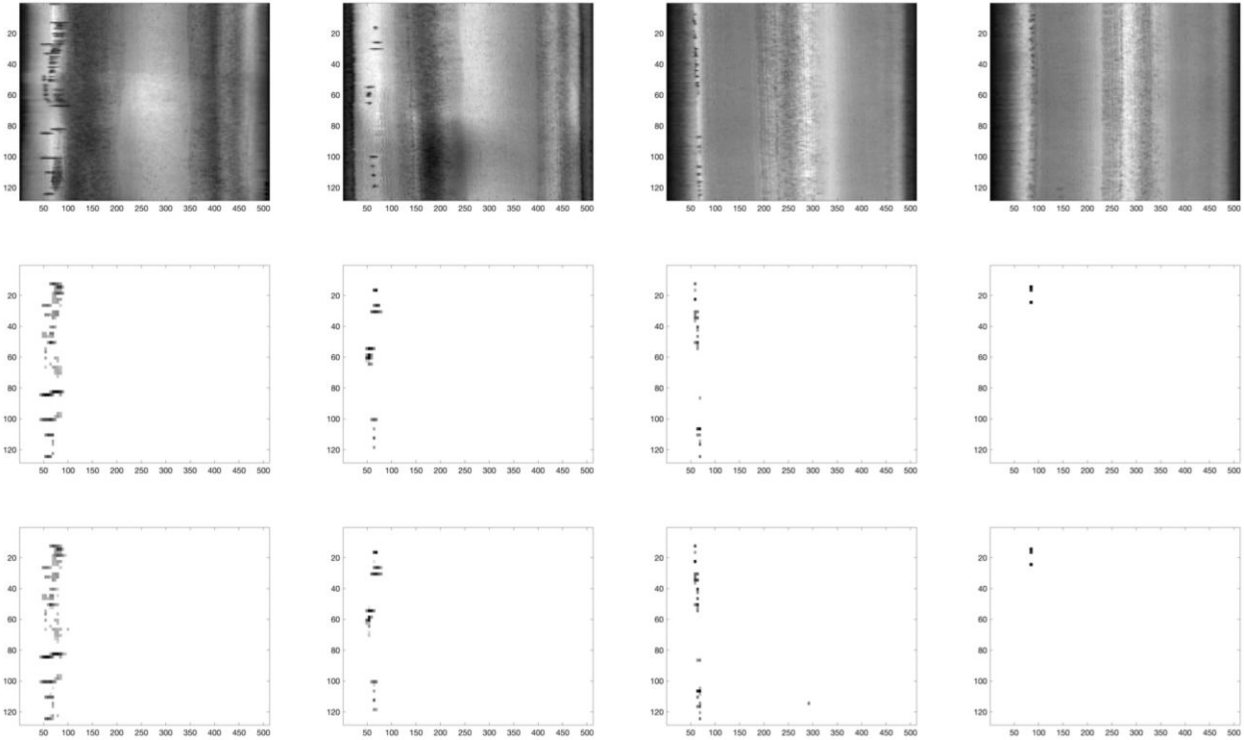
Notes. (a)–(d) Example signals corresponding to different  $\sigma_s$ , (e) log relative error for the smooth component, and (f) log relative error for the sparse component. (a)  $\sigma_s = 0.065$ . (b)  $\sigma_s = 0.125$ . (c)  $\sigma_s = 0.25$ . (d)  $\sigma_s = 0.5$ . (e) Log relative error for the smooth component. (f) Log relative error for the sparse component

## Appendix I. Sample Case Study Images

### Appendix I.1. Sample Images of Case Study 4.1

See Figure I.1.

**Figure I.1.** Sample Images of Surface Defect Detection in the Steel Rolling Process

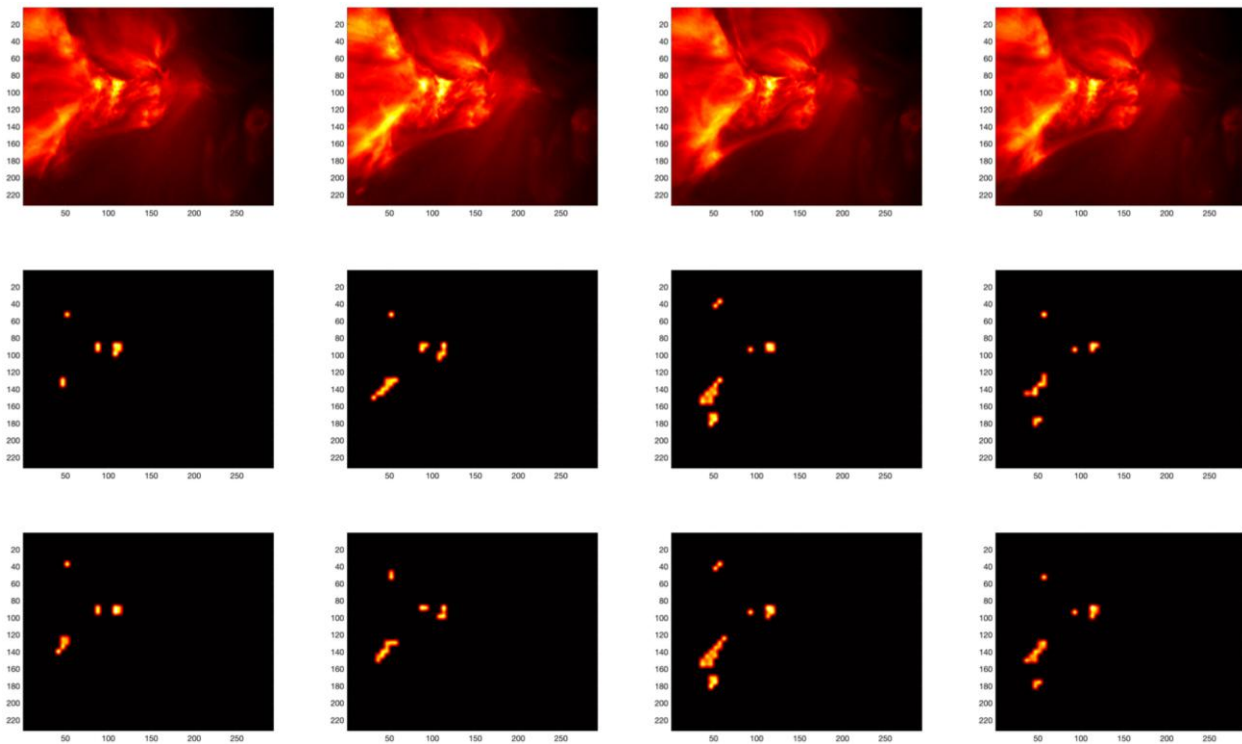


Notes. The top row shows raw images. The middle row shows corresponding SSD results. The bottom row shows KronCSSD results.

## Appendix I.2. Sample Images of Case Study 4.2

See Figure I.2.

Figure I.2. (Color online) Sample Images of Solar Flare Detection



Notes. The top row shows raw images. The middle row shows corresponding SSD results. The bottom row shows KronCSSD results.

## References

- Augusto CRA, Fauth AC, Navia CE, Shigeouka H, Tsui KH (2011) Connection among spacecrafts and ground level observations of small solar transient events. *Experiment. Astronomy* 31(2):177–197.
- Baraniuk R, Davenport M, DeVore R, Wakin M (2008) A simple proof of the restricted isometry property for random matrices. *Constructive Approximation* 28(3):253–263.
- Bouwmans T, Zahzah EH (2014) Robust PCA via principal component pursuit: A review for a comparative evaluation in video surveillance. *Comput. Vision Image Understanding* 122:22–34.
- Candes EJ (2008) The restricted isometry property and its implications for compressed sensing. *Competus Rendus Math.* 346(9–10):589–592.
- Candes EJ, Tao T (2005) Decoding by linear programming. *IEEE Trans. Inform. Theory* 51(12):4203–4215.
- Candes EJ, Romberg JK, Tao T (2006) Stable signal recovery from incomplete and inaccurate measurements. *Comm. Pure Appl. Math.* 59(8):1207–1223.
- Candès EJ, Li X, Ma Y, Wright J (2011) Robust principal component analysis? *J. ACM* 58(3):1–37.
- De Boor C, De Boor C (1978) *A Practical Guide to Splines*, vol. 27 (Springer-Verlag, New York).
- Duarte MF, Baraniuk RG (2011) Kronecker compressive sensing. *IEEE Trans. Image Processing* 21(2):494–504.
- Eilers PH, Marx BD (1996) Flexible smoothing with B-splines and penalties. *Statist. Sci.* 11(2):89–121.
- Giannakis GB, Mateos G, Farahmand S, Kekatos V, Zhu H (2011) USPA-COR: Universal sparsity-controlling outlier rejection. Tichavsky P, Cernocky H, Prochazka A, eds. *2011 IEEE Internat. Conf. Acoustics Speech Signal Processing* (IEEE, New York), 1952–1955.
- Kolda TG, Bader BW (2009) Tensor decompositions and applications. *SIAM Rev.* 51(3):455–500.
- Mardani M, Mateos G, Giannakis GB (2013) Recovery of low-rank plus compressed sparse matrices with application to unveiling traffic anomalies. *IEEE Trans. Inform. Theory* 59(8):5186–5205.
- Marques EC, Maciel N, Naviner L, Cai H, Yang J (2018) A review of sparse recovery algorithms. *IEEE Access* 7:1300–1322.
- Mateos G, Giannakis GB (2011) Robust nonparametric regression by controlling sparsity. *2011 IEEE Internat. Conf. Acoustics Speech Signal Processing (ICASSP)*.
- Minaee S, Abdolrashidi A, Wang Y (2015) Screen content image segmentation using sparse-smooth decomposition. *2015 49th Asilomar Conf. Signals Systems Comput.*
- Mou S, Wang A, Zhang C, Shi J (2021) Additive tensor decomposition considering structural data information. *IEEE Trans. Automation Sci. Engrg.* 19(4):2904–2917.
- Rani M, Dhok SB, Deshmukh RB (2018) A systematic review of compressive sensing: Concepts, implementations and applications. *IEEE Access* 6:4875–4894.
- Recht B, Fazel M, Parrilo PA (2010) Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.* 52(3):471–501.
- Tanner J, Vary S (2020) Compressed sensing of low-rank plus sparse matrices. Preprint, submitted July 18, <https://arxiv.org/abs/2007.09457v1>.

- Unser M (1999) Splines: A perfect fit for signal and image processing. *IEEE Signal Processing Magazine* 16(6):22–38.
- Vershynin R (2018) *High-Dimensional Probability: An Introduction with Applications in Data Science*, vol. 47 (Cambridge University Press, Cambridge, United Kingdom).
- Wang A, Xian X, Tsung F, Liu K (2018) A spatial-adaptive sampling procedure for online monitoring of big data streams. *J. Quality Tech.* 50(4):329–343.
- Waters AE, Sankaranarayanan AC, Baraniuk RG (2011) SpaRCS: Recovering low-rank and sparse matrices from compressive measurements. Shawe-Taylor J, Zemel RS, Bartlett PL, Pereira F, Weinberger KQ, eds. *Conf. Neural Inform. Processing Systems* (Curran Associates Inc., Red Hook, NY), 1089–1097.
- Xu H, Caramanis C, Sanghavi S (2012) Robust PCA via outlier pursuit. *IEEE Trans. Inform. Theory* 58(5):3047–3064.
- Yan H, Paynabar K, Shi J (2017) Anomaly detection in images with smooth background via smooth-sparse decomposition. *Technometrics* 59(1):102–114.
- Yan H, Paynabar K, Shi J (2018) Real-time monitoring of high-dimensional functional data streams via spatio-temporal smooth sparse decomposition. *Technometrics* 60(2):181–197.