

# Matrix Recipes for Hard Thresholding Methods

Anastasios Kyrillidis · Volkan Cevher

Published online: 6 April 2013  
© Springer Science+Business Media New York 2013

**Abstract** In this paper, we present and analyze a new set of low-rank recovery algorithms for linear inverse problems within the class of hard thresholding methods. We provide strategies on how to set up these algorithms via basic ingredients for different configurations to achieve complexity vs. accuracy tradeoffs. Moreover, we study acceleration schemes via memory-based techniques and randomized,  $\epsilon$ -approximate matrix projections to decrease the computational costs in the recovery process. For most of the configurations, we present theoretical analysis that guarantees convergence under mild problem conditions. Simulation results demonstrate notable performance improvements as compared to state-of-the-art algorithms both in terms of reconstruction accuracy and computational complexity.

**Keywords** Affine rank minimization · Hard thresholding ·  $\epsilon$ -approximation schemes · Randomized algorithms

## 1 Introduction

In this work, we consider the general affine rank minimization (ARM) problem, described as follows:

**THE ARM PROBLEM:** Assume  $X^* \in \mathbb{R}^{m \times n}$  is a rank- $k$  matrix of interest ( $k \ll \min\{m, n\}$ ) and let  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  be a known linear operator. Given a set of observations as

$y = \mathcal{A}X^* + \epsilon \in \mathbb{R}^p$ , we desire to recover  $X^*$  from  $y$  in a scalable and robust manner.

The challenge in this problem is to recover the true low-rank matrix in subsampled settings where  $p \ll m \cdot n$ . In such cases, we typically exploit the prior information that  $X^*$  is low-rank and thus, we are interested in finding a matrix  $X$  of rank at most  $k$  that minimizes the data error  $f(X) := \|y - \mathcal{A}X\|_2^2$  as follows:

$$\begin{aligned} & \underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} && f(X) \\ & \text{subject to} && \text{rank}(X) \leq k. \end{aligned} \tag{1}$$

The ARM problem appears in many applications; low dimensional embedding [1], matrix completion [2], image compression [3], function learning [4, 5] just to name a few. We present below important ARM problem cases, as characterized by the nature of the linear operator  $\mathcal{A}$ .

**General linear maps:** In many ARM problem cases,  $\mathcal{A}$  or  $\mathcal{A}^*$  has a dense range, satisfying specific incoherence or restricted isometry properties (discussed later in the paper); here,  $\mathcal{A}^*$  is the adjoint operator of  $\mathcal{A}$ . In Quantum Tomography, [6] studies the Pauli operator, a *compressive* linear map  $\mathcal{A}$  that consists of the Kronecker product of  $2 \times 2$  matrices and obeys restricted isometry properties, defined later in the paper. Furthermore, recent developments indicate connections of ridge function learning [4, 7] and phase retrieval [8] with the ARM problem where  $\mathcal{A}$  is a Bernoulli and a Fourier operator, respectively.

**Matrix Completion (MC):** Let  $\Omega$  be the set of ordered pairs that represent the coordinates of the observable entries in  $X^*$ . Then, the set of observations satisfy  $y = \mathcal{A}_\Omega X^* + \epsilon$  where  $\mathcal{A}_\Omega$  defines a linear mask over the observable entries  $\Omega$ . To solve the MC problem, a potential criterion is given by (1) [2]. As a motivating example, consider the famous

---

A. Kyrillidis (✉) · V. Cevher  
Laboratory for Information and Inference Systems, Ecole  
Polytechnique Federale de Lausanne, Lausanne, Switzerland  
e-mail: [anastasios.kyrillidis@epfl.ch](mailto:anastasios.kyrillidis@epfl.ch)  
url: <http://lions.epfl.ch>

V. Cevher  
e-mail: [volkan.cevher@epfl.ch](mailto:volkan.cevher@epfl.ch)

Netflix problem [9], a recommender system problem where users' movie preferences are inferred by a limited subset of entries in a database.

**Principal Component Analysis:** In Principal Component Analysis (PCA), we are interested in identifying a low rank subspace that best explains the data in the Euclidean sense from the observations  $\mathbf{y} = \mathcal{A}\mathbf{X}^*$  where  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  is an identity linear map that stacks the columns of the matrix  $\mathbf{X}^*$  into a single column vector with  $p = m \cdot n$ . We observe that the PCA problem falls under the ARM criterion in (1). While (1) is generally NP-hard to solve optimally, PCA can be solved in polynomial time using the truncated Singular Value Decomposition (SVD) of  $\mathcal{A}^* \mathbf{y}$ . As an extension to the PCA setting, [10] considers the Robust PCA problem where  $\mathbf{y}$  is further corrupted by gross sparse noise. We extend the framework proposed in this paper for the RPCA case and its generalizations in [11].

For the rest of the paper, we consider only the low rank estimation case in (1). As running test cases to support our claims, we consider the MC setting as well as the general ARM setting where  $\mathcal{A}$  is constituted by permuted subsampled noiselets [12].

### 1.1 Two Camps of Recovery Algorithms

**Convex relaxations:** In [13], the authors study the nuclear norm  $\|\mathbf{X}\|_* := \sum_{i=1}^{\text{rank}(\mathbf{X})} \sigma_i$  as a convex surrogate of  $\text{rank}(\mathbf{X})$  operator so that we can leverage convex optimization approaches, such as interior-point methods—here,  $\sigma_i$  denotes the  $i$ -th singular value of  $\mathbf{X}$ . Under basic incoherence properties of the sensing linear mapping  $\mathcal{A}$ , [13] provides provable guarantees for unique low rank matrix recovery using the nuclear norm.

Once (1) is relaxed to a convex problem, decades of knowledge on convex analysis and optimization can be leveraged. Interior point methods find a solution with fixed precision in polynomial time but their complexity might be prohibitive even for moderate-sized problems [14, 15]. More suitable for large-scale data analysis, first-order methods constitute low-complexity alternatives but most of them introduce complexity vs. accuracy tradeoffs [16–19].

**Non-convex approaches:** In contrast to the convex relaxation approaches, iterative greedy algorithms maintain the nonconvex nature of (1). Unfortunately, solving (1) optimally is in general NP-hard [20]. Due to this computational intractability, the algorithms in this class greedily refine a rank- $k$  solution using only “local” information available at the current iteration [21–23].

### 1.2 Contributions

In this work, we study a special class of iterative greedy algorithms known as hard thresholding methods. Similar re-

sults have been derived for the vector case [24]. Note that the transition from sparse vector approximation to ARM is *non-trivial*; while  $s$ -sparse signals “live” in the union of finite number of subspaces, the set of rank- $k$  matrices expands to infinitely many subspaces. Thus, the selection rules do not generalize in a straightforward way.

Our contributions are the following:

**Ingredients of hard thresholding methods:** We analyze the behaviour and performance of hard thresholding methods from a global perspective. Five building blocks are studied: (i) step size selection  $\mu_i$ , (ii) gradient or least-squares updates over restricted low-rank subspaces (e.g., adaptive block coordinate descent), (iii) memory exploitation, (iv) active low-rank subspace tracking and, (v) low-rank matrix approximations (described next). We highlight the impact of these key pieces on the convergence rate and signal reconstruction performance and provide optimal and/or efficient strategies on how to set up these ingredients under different problem conditions.

**Low-rank matrix approximations in hard thresholding methods:** In [25], the authors show that the solution efficiency can be significantly improved by  $\epsilon$ -approximation algorithms. Based on similar ideas, we analyze the impact of  $\epsilon$ -approximate low rank-revealing schemes in the proposed algorithms with well-characterized time and space complexities. Moreover, we provide extensive analysis to prove convergence using  $\epsilon$ -approximate low-rank projections.

**Hard thresholding-based framework with improved convergence conditions:** We study hard thresholding variants that provide salient computational tradeoffs for the class of greedy methods on low-rank matrix recovery. These methods, as they iterate, exploit the non-convex scaffold of low rank subspaces on which the approximation problem resides. Using simple analysis tools, we derive improved conditions that guarantee convergence, compared to state-of-the-art approaches.

The organization of the paper is as follows. In Sect. 2, we set up the notation and provide some definitions and properties, essential for the rest of the paper. In Sect. 3, we describe the basic algorithmic frameworks in a nutshell, while in Sect. 4 we provide important “ingredients” for the class of hard-thresholding methods; detailed convergence analysis proofs are provided in Sect. 5. The complexity analysis of the proposed algorithms is provided in Sect. 6. We study two acceleration schemes in Sects. 7 and 8, based on memory utilization and  $\epsilon$ -approximate low-rank projections, respectively. We further improve convergence speed by exploiting randomized low rank projections in Sect. 9, based on power iteration-based subspace finder tools [26]. We provide empirical support for our claims through experimental results on synthetic and real data in Sect. 10. Finally, we conclude with future work directions in Sect. 11.

## 2 Elementary Definitions and Properties

We reserve lower-case and bold lower-case letters for scalar and vector variable representation, respectively. Bold upper-case letters denote matrices while bold calligraphic upper-case letters represent linear operators. We use calligraphic upper-case letters for set representations. We use  $X(i)$  to represent the matrix estimate at the  $i$ -th iteration.

The rank of  $X$  is denoted as  $\text{rank}(X) \leq \min\{m, n\}$ . The empirical data error is denoted as  $f(X) := \|y - \mathcal{A}X\|_2^2$  with gradient  $\nabla f(X) := -2\mathcal{A}^*(y - \mathcal{A}X)$ , where  $*$  is the adjoint operation over the linear mapping  $\mathcal{A}$ . The inner product between matrices  $A, B \in \mathbb{R}^{m \times n}$  is denoted as  $\langle A, B \rangle = \text{trace}(B^T A)$ , where  $T$  represents the transpose operation.  $\mathbf{I}$  represents an identity matrix with dimensions apparent from the context.

Let  $\mathcal{S}$  be a set of orthonormal, rank-1 matrices that span an arbitrary subspace in  $\mathbb{R}^{m \times n}$ . We reserve  $\text{span}(\mathcal{S})$  to denote the subspace spanned by  $\mathcal{S}$ . With slight abuse of notation, we use:

$$\text{rank}(\text{span}(\mathcal{S})) \equiv \max_X \{\text{rank}(X) : X \in \text{span}(\mathcal{S})\}, \tag{2}$$

to denote the *maximum* rank a matrix  $X \in \mathbb{R}^{m \times n}$  can have such that  $X$  lies in the subspace spanned by the set  $\mathcal{S}$ . Given a finite set  $\mathcal{S}$ ,  $|\mathcal{S}|$  denotes the cardinality of  $\mathcal{S}$ . For any matrix  $X$ , we use  $R(X)$  to denote its range.

We define a *minimum cardinality* set of orthonormal, rank-1 matrices that span the subspace induced by a set of rank-1 (and possibly non-orthogonal) matrices  $\mathcal{S}$  as:

$$\text{ortho}(\mathcal{S}) \in \arg \min_{\mathcal{T}} \{|\mathcal{T}| : \mathcal{T} \subseteq \mathcal{U} \text{ s.t. } \text{span}(\mathcal{T}) = \text{span}(\mathcal{S})\},$$

where  $\mathcal{U}$  denotes the superset that includes all the sets of *orthonormal*, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that  $\langle T_i, T_j \rangle = 0, i \neq j, \forall T_i, T_j \in \mathcal{T}$  and,  $\|T_i\|_F = 1, \forall i$ . In general,  $\text{ortho}(\mathcal{S})$  is not unique.

A well-known lemma used in the convergence rate proofs of this class of greedy hard thresholding algorithms is defined next.

**Lemma 1** [27] *Let  $\mathcal{J} \subseteq \mathbb{R}^{m \times n}$  be a closed convex set and  $f : \mathcal{J} \rightarrow \mathbb{R}$  be a smooth objective function defined over  $\mathcal{J}$ . Let  $X^* \in \mathcal{J}$  be a local minimum of the objective function  $f$  over the set  $\mathcal{J}$ . Then*

$$\langle \nabla f(X^*), X - X^* \rangle \geq 0, \quad \forall X \in \mathcal{J}. \tag{3}$$

### 2.1 Singular Value Decomposition (SVD) and Its Properties

**Definition 1** [SVD] Let  $X \in \mathbb{R}^{m \times n}$  be a rank- $l$  ( $l < \min\{m, n\}$ ) matrix. Then, the SVD of  $X$  is given by:

$$X = U \Sigma V^T = [U_\alpha \ U_\beta] \begin{bmatrix} \tilde{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} V_\alpha^T \\ V_\beta^T \end{bmatrix}, \tag{4}$$

where  $U_\alpha \in \mathbb{R}^{m \times l}, U_\beta \in \mathbb{R}^{m \times (m-l)}, V_\alpha \in \mathbb{R}^{n \times l}, V_\beta \in \mathbb{R}^{n \times (n-l)}$  and  $\tilde{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_l) \in \mathbb{R}^{l \times l}$  for  $\sigma_1, \dots, \sigma_l \in \mathbb{R}_+$ . Here, the columns of  $U, V$  represent the set of left and right singular vectors, respectively, and  $\sigma_1, \dots, \sigma_l$  denote the singular values.

For any matrix  $X \in \mathbb{R}^{m \times n}$  with arbitrary  $\text{rank}(X) \leq \min\{m, n\}$ , its best orthogonal projection  $\mathcal{P}_k(X)$  onto the set of rank- $k$  ( $k < \text{rank}(X)$ ) matrices  $\mathcal{C}_k := \{A \in \mathbb{R}^{m \times n} : \text{rank}(A) \leq k\}$  defines the optimization problem:

$$\mathcal{P}_k(X) \in \arg \min_{Y \in \mathcal{C}_k} \|Y - X\|_F. \tag{5}$$

According to the Eckart-Young theorem [28], the best rank- $k$  approximation of a matrix  $X$  corresponds to its truncated SVD: if  $X = U \Sigma V^T$ , then  $\mathcal{P}_k(X) := U_k \Sigma_k V_k^T$  where  $\Sigma_k \in \mathbb{R}^{k \times k}$  is a diagonal matrix that contains the first  $k$  diagonal entries of  $\Sigma$  and  $U_k, V_k$  contain the corresponding left and right singular vectors, respectively. Moreover, this projection is not always unique. In the case of multiple identical singular values, the lexicographic approach is used to break ties. In any case,  $\|\mathcal{P}_k(X) - X\|_F \leq \|W - X\|_F$  for any rank- $k$   $W \in \mathbb{R}^{m \times n}$ .

### 2.2 Subspace Projections

Given a set of orthonormal, rank-1 matrices  $\mathcal{S}$ , we denote the orthogonal projection operator onto the subspace induced by  $\mathcal{S}$  as  $\mathcal{P}_\mathcal{S}$ <sup>1</sup> which is an idempotent linear transformation; furthermore, we denote the orthogonal projection operator onto the orthogonal subspace of  $\mathcal{S}$  as  $\mathcal{P}_{\mathcal{S}^\perp}$ . We can always decompose a matrix  $X \in \mathbb{R}^{m \times n}$  into two matrix components, as follows:

$$X := \mathcal{P}_\mathcal{S} X + \mathcal{P}_{\mathcal{S}^\perp} X, \quad \text{such that } \langle \mathcal{P}_\mathcal{S} X, \mathcal{P}_{\mathcal{S}^\perp} X \rangle = 0.$$

If  $X \in \text{span}(\mathcal{S})$ , the best projection of  $X$  onto the subspace induced by  $\mathcal{S}$  is the matrix  $X$  itself. Moreover,  $\|\mathcal{P}_\mathcal{S} X\|_F \leq \|X\|_F$  for any  $\mathcal{S}$  and  $X$ .

**Definition 2** [Orthogonal projections using SVD] Let  $X \in \mathbb{R}^{m \times n}$  be a matrix with arbitrary rank and SVD decomposition given by (4). Then,  $\mathcal{S} := \{u_i v_i^T : i = 1, \dots, k\}$  ( $k \leq \text{rank}(X)$ ) constitutes a set of orthonormal, rank-1 matrices that spans the best  $k$ -rank subspace in  $R(X)$  and  $R(X^T)$ ; here,  $u_i$  and  $v_i$  denote the  $i$ -th left and right singular vectors, respectively. The orthogonal projection onto this subspace is given by [2]:

$$\mathcal{P}_\mathcal{S} X = \mathcal{P}_U X + X \mathcal{P}_V - \mathcal{P}_U X \mathcal{P}_V, \tag{6}$$

<sup>1</sup>The distinction between  $\mathcal{P}_\mathcal{S}$  and  $\mathcal{P}_k$  for  $k$  positive integer is apparent from context.

where  $\mathcal{P}_U = U_{:,1:k}U_{:,1:k}^T$  and  $\mathcal{P}_V = V_{:,1:k}V_{:,1:k}^T$  in MATLAB notation. Moreover, the orthogonal projection onto the  $\mathcal{S}^\perp$  is given by:

$$\mathcal{P}_{\mathcal{S}^\perp} \mathbf{X} = \mathbf{X} - \mathcal{P}_S \mathbf{X}. \tag{7}$$

In the algorithmic descriptions, we use  $\mathcal{S} \leftarrow \mathcal{P}_k(\mathbf{X})$  to denote the set of rank-1, orthonormal matrices as outer products of the  $k$  left  $\mathbf{u}_i$  and right  $\mathbf{v}_i$  principal singular vectors of  $\mathbf{X}$  that span the best rank- $k$  subspace of  $\mathbf{X}$ ; e.g.  $\mathcal{S} = \{\mathbf{u}_i \mathbf{v}_i, i = 1, \dots, k\}$ . Moreover,  $\widehat{\mathbf{X}} \leftarrow \mathcal{P}_k(\mathbf{X})$  denotes a/the best rank- $k$  projection matrix of  $\mathbf{X}$ . In some cases, we use  $\{\mathcal{S}, \widehat{\mathbf{X}}\} \leftarrow \mathcal{P}_k(\mathbf{X})$  when we compute both. The distinction between these cases is apparent from the context.

### 2.3 Restricted Isometry Property

Many conditions have been proposed in the literature to establish solution uniqueness and recovery stability such as null space property [29], exact recovery condition [30], etc. For the matrix case, [13] proposed the *restricted isometry property* (RIP) for the ARM problem.

**Definition 3** [Rank Restricted Isometry Property (R-RIP) for matrix linear operators [13]] A linear operator  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  satisfies the R-RIP with constant  $\delta_k(\mathcal{A}) \in (0, 1)$  if and only if:

$$(1 - \delta_k(\mathcal{A})) \|\mathbf{X}\|_F^2 \leq \|\mathcal{A}\mathbf{X}\|_2^2 \leq (1 + \delta_k(\mathcal{A})) \|\mathbf{X}\|_F^2, \tag{8}$$

$\forall \mathbf{X} \in \mathbb{R}^{m \times n}$  such that  $\text{rank}(\mathbf{X}) \leq k$ . We write  $\delta_k$  to mean  $\delta_k(\mathcal{A})$ , unless otherwise stated.

[6] shows that Pauli operators satisfy the rank-RIP in compressive settings while, in function learning, the linear map  $\mathcal{A}$  is designed specifically to satisfy the rank-RIP [7].

### 2.4 Some Useful Bounds Using R-RIP

In this section, we present some lemmas that are useful in our subsequent developments—these lemmas are consequences of the R-RIP of  $\mathcal{A}$ .

**Lemma 2** [21] Let  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  be a linear operator that satisfies the R-RIP with constant  $\delta_k$ . Then,  $\forall \mathbf{v} \in \mathbb{R}^p$ , the following holds true:

$$\|\mathcal{P}_S(\mathcal{A}^* \mathbf{v})\|_F \leq \sqrt{1 + \delta_k} \|\mathbf{v}\|_2, \tag{9}$$

where  $\mathcal{S}$  is a set of orthonormal, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that  $\text{rank}(\mathcal{P}_S \mathbf{X}) \leq k, \forall \mathbf{X} \in \mathbb{R}^{m \times n}$ .

**Lemma 3** [21] Let  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  be a linear operator that satisfies the R-RIP with constant  $\delta_k$ . Then,  $\forall \mathbf{X} \in \mathbb{R}^{m \times n}$ ,

the following holds true:

$$\begin{aligned} (1 - \delta_k) \|\mathcal{P}_S \mathbf{X}\|_F &\leq \|\mathcal{P}_S \mathcal{A}^* \mathcal{A} \mathcal{P}_S \mathbf{X}\|_F \\ &\leq (1 + \delta_k) \|\mathcal{P}_S \mathbf{X}\|_F, \end{aligned} \tag{10}$$

where  $\mathcal{S}$  is a set of orthonormal, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that  $\text{rank}(\mathcal{P}_S \mathbf{X}) \leq k, \forall \mathbf{X} \in \mathbb{R}^{m \times n}$ .

**Lemma 4** [22] Let  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  be a linear operator that satisfies the R-RIP with constant  $\delta_k$  and  $\mathcal{S}$  be a set of orthonormal, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that  $\text{rank}(\mathcal{P}_S \mathbf{X}) \leq k, \forall \mathbf{X} \in \mathbb{R}^{m \times n}$ . Then, for  $\mu > 0$ ,  $\mathcal{A}$  satisfies:

$$\lambda(\mu \mathcal{P}_S \mathcal{A}^* \mathcal{A} \mathcal{P}_S) \in [\mu(1 - \delta_k), \mu(1 + \delta_k)], \tag{11}$$

where  $\lambda(\mathcal{B})$  represents the range of eigenvalues of the linear operator  $\mathcal{B} : \mathbb{R}^p \rightarrow \mathbb{R}^{m \times n}$ . Moreover,  $\forall \mathbf{X} \in \mathbb{R}^{m \times n}$ , it follows that:

$$\begin{aligned} \|(\mathbf{I} - \mu \mathcal{P}_S \mathcal{A}^* \mathcal{A} \mathcal{P}_S) \mathcal{P}_S \mathbf{X}\|_F \\ \leq \max\{\mu(1 + \delta_k) - 1, 1 - \mu(1 - \delta_k)\} \|\mathcal{P}_S \mathbf{X}\|_F. \end{aligned} \tag{12}$$

**Lemma 5** [22] Let  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$  be a linear operator that satisfies the R-RIP with constant  $\delta_k$  and  $\mathcal{S}_1, \mathcal{S}_2$  be two sets of orthonormal, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that

$$\text{rank}(\mathcal{P}_{\mathcal{S}_1 \cup \mathcal{S}_2} \mathbf{X}) \leq k, \quad \forall \mathbf{X} \in \mathbb{R}^{m \times n}. \tag{13}$$

Then, the following inequality holds:

$$\|\mathcal{P}_{\mathcal{S}_1} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_1^\perp} \mathbf{X}\|_F \leq \delta_k \|\mathcal{P}_{\mathcal{S}_1^\perp} \mathbf{X}\|_F, \quad \forall \mathbf{X} \in \text{span}(\mathcal{S}_2). \tag{14}$$

## 3 Algebraic Pursuits in a Nutshell

Explicit descriptions of the proposed algorithms are provided in Algorithms 1 and 2. Algorithm 1 follows from the ALgebraic PursuitS (ALPS) scheme for the vector case [31]. MATRIX ALPS I provides efficient strategies for adaptive step size selection and additional signal estimate updates at each iteration (these motions are explained in detail in the next subsection). Algorithm 2 (ADMIRA) [21] further improves the performance of Algorithm 1 by introducing least squares optimization steps on restricted subspaces—this technique borrows from a series of vector reconstruction algorithms such as CoSaMP [32], Subspace Pursuit (SP) [33] and Hard Thresholding Pursuit (HTP) [34].

In a nutshell, both algorithms simply seek to improve the subspace selection by iteratively collecting an extended subspace  $\mathcal{S}_i$  with  $\text{rank}(\text{span}(\mathcal{S}_i)) \leq 2k$  and then finding the rank- $k$  matrix that fits the measurements in this restricted subspace using least squares or gradient descent motions.

**Algorithm 1** MATRIX ALPS I

---

**Input:**  $\mathbf{y}$ ,  $\mathcal{A}$ ,  $k$ , Tolerance  $\eta$ , MaxIterations  
**Initialize:**  $\mathbf{X}(0) \leftarrow 0$ ,  $\mathcal{X}_0 \leftarrow \{\emptyset\}$ ,  $i \leftarrow 0$   
**repeat**

1:  $\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i)))$  (Best rank- $k$  subspace orthogonal to  $\mathcal{X}_i$ )  
2:  $\mathcal{S}_i \leftarrow \mathcal{D}_i \cup \mathcal{X}_i$  (Active subspace expansion)  
3:  $\mu_i \leftarrow \arg \min_{\mu} \|\mathbf{y} - \mathcal{A}(\mathbf{X}(i) - \frac{\mu}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)))\|_2^2 = \frac{\|\mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i))\|_F^2}{\|\mathcal{A} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i))\|_2^2}$  (Step size selection)  
4:  $\mathbf{V}(i) \leftarrow \mathbf{X}(i) - \frac{\mu_i}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i))$  (Error norm reduction via gradient descent)  
5:  $\{\mathcal{W}_i, \mathbf{W}(i)\} \leftarrow \mathcal{P}_k(\mathbf{V}(i))$  (Best rank- $k$  subspace selection)  
6:  $\xi_i \leftarrow \arg \min_{\xi} \|\mathbf{y} - \mathcal{A}(\mathbf{W}(i) - \frac{\xi}{2} \mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i)))\|_2^2 = \frac{\|\mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i))\|_F^2}{\|\mathcal{A} \mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i))\|_2^2}$  (Step size selection)  
7:  $\mathbf{X}(i+1) \leftarrow \mathbf{W}(i) - \frac{\xi_i}{2} \mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i))$  with  $\mathcal{X}_{i+1} \leftarrow \mathcal{P}_k(\mathbf{X}(i+1))$  (De-bias using gradient descent)  
     $i \leftarrow i+1$

**until**  $\|\mathbf{X}(i) - \mathbf{X}(i-1)\|_2 \leq \eta \|\mathbf{X}(i)\|_2$  or MaxIterations.

---

**Algorithm 2** ADMiRA Instance

---

**Input:**  $\mathbf{y}$ ,  $\mathcal{A}$ ,  $k$ , Tolerance  $\eta$ , MaxIterations  
**Initialize:**  $\mathbf{X}(0) \leftarrow 0$ ,  $\mathcal{X}_0 \leftarrow \{\emptyset\}$ ,  $i \leftarrow 0$   
**repeat**

1:  $\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i)))$  (Best rank- $k$  subspace orthogonal to  $\mathcal{X}_i$ )  
2:  $\mathcal{S}_i \leftarrow \mathcal{D}_i \cup \mathcal{X}_i$  (Active subspace expansion)  
3:  $\mathbf{V}(i) \leftarrow \arg \min_{\mathbf{V}: \mathbf{V} \in \text{span}(\mathcal{S}_i)} \|\mathbf{y} - \mathcal{A} \mathbf{V}\|_2^2$  (Error norm reduction via least-squares optimization)  
4:  $\{\mathcal{X}_{i+1}, \mathbf{X}(i+1)\} \leftarrow \mathcal{P}_k(\mathbf{V}(i))$  (Best rank- $k$  subspace selection)  
     $i \leftarrow i+1$

**until**  $\|\mathbf{X}(i) - \mathbf{X}(i-1)\|_2 \leq \eta \|\mathbf{X}(i)\|_2$  or MaxIterations.

---

At each iteration, the Algorithms 1 and 2 perform motions from the following list:

(1) *Best rank- $k$  subspace orthogonal to  $\mathcal{X}_i$  and active subspace expansion:* We identify the best rank- $k$  subspace of the current gradient  $\nabla f(\mathbf{X}(i))$ , orthogonal to  $\mathcal{X}_i$  and then merge this low-rank subspace with  $\mathcal{X}_i$ . This motion guarantees that, at each iteration, we expand the current rank- $k$  subspace estimate with  $k$  new, rank-1 orthogonal subspaces to explore.

(2a) *Error norm reduction via greedy descent with adaptive step size selection (Algorithm 1):* We decrease the data error by performing a single gradient descent step. This scheme is based on a one-shot step size selection procedure (Step size selection step)—detailed description of this approach is given in Sect. 4.

(2b) *Error norm reduction via least squares optimization (Algorithm 2):* We decrease the data error  $f(\mathbf{X})$  on the active  $O(k)$ -low rank subspace. Assuming  $\mathcal{A}$  is well-conditioned over low-rank subspaces, the main complexity of this operation is dominated by the solution of a symmetric linear system of equations.

(3) *Best rank- $k$  subspace selection:* We project the constrained solution onto the set of rank- $k$  matrices  $\mathcal{C}_k := \{\mathbf{A} \in$

$\mathbb{R}^{m \times n} : \text{rank}(\mathbf{A}) \leq k\}$  to arbitrate the active support set. This step is calculated in polynomial time complexity as a function of  $m \times n$  using SVD or other matrix rank-revealing decomposition algorithms—further discussions about this step and its approximations can be found in Sects. 8 and 9.

(4) *De-bias using gradient descent (Algorithm 1):* We de-bias the current estimate  $\mathbf{W}(i)$  by performing an additional gradient descent step, decreasing the data error. The step size selection procedure follows the same motions as in (2a).

**4 Ingredients for Hard Thresholding Methods**

4.1 Step Size Selection

For the sparse vector approximation problem, recent works on the performance of the IHT algorithm provide strong convergence rate guarantees in terms of RIP constants [35]. However, as a prerequisite to achieve these strong isometry constant bounds, the step size is set  $\mu_i = 1, \forall i$ , given that the sensing matrix satisfies  $\|\Phi\|_2^2 < 1$  where  $\|\cdot\|_2$  denotes the spectral norm [34]; similar analysis can be found in [3] for the matrix case. From a different perspective, [36] proposes a constant step size  $\mu_i = 1/(1 + \delta_{2K}), \forall i$ , based on



a simple but intuitive convergence analysis of the gradient descent method.

Unfortunately, most of the above problem assumptions are not naturally met; the authors in [37] provide an intuitive example where IHT algorithm behaves differently under various scalings of the sensing matrix; similar counterexamples can be devised for the matrix case. Violating these assumptions usually leads to unpredictable signal recovery performance of the class of hard thresholding methods. Therefore, more sophisticated step size selection procedures should be devised to tackle these issues during actual recovery. On the other hand, the computation of R-RIP constants has exponential time complexity for the strategy of [3].

To this end, existing approaches broadly fall into two categories: constant and adaptive step size selection. In this work, we present efficient strategies to adaptively select the step size  $\mu_i$  that implies fast convergence rate, for mild R-RIP assumptions on  $\mathcal{A}$ . Constant step size strategies easily follow from [24] and are not listed in this work.

**Adaptive step size selection.** There is limited work on the adaptive step size selection for hard thresholding methods. To the best of our knowledge, apart from [24, 37, 38] are the only studies that attempt this via line searching for the vector case. At the time of review process, we become aware of [39] which implements ideas presented in [37] for the matrix case.

According to Algorithm 1, let  $\mathbf{X}(i)$  be the current rank- $k$  matrix estimate spanned by the set of orthonormal, rank-1 matrices in  $\mathcal{X}_i$ . Using regular gradient descent motions, the new rank- $k$  estimate  $\mathbf{W}(i)$  can be calculated through:

$$\mathbf{V}_i = \mathbf{X}(i) - \frac{\mu}{2} \nabla f(\mathbf{X}(i)), \quad \{\mathcal{W}_i, \mathbf{W}(i)\} \leftarrow \mathcal{P}_k(\mathbf{V}(i)).$$

We highlight that the rank- $k$  approximate matrix may not be unique. It then holds that the subspace spanned by  $\mathcal{W}_i$  originates: (i) either from the subspace of  $\mathcal{X}_i$ , (ii) or from the best subspace (in terms of the Frobenius norm metric) of the current gradient  $\nabla f(\mathbf{X}(i))$ , *orthogonal to*  $\mathcal{X}_i$ , (iii) or from the combination of orthonormal, rank-1 matrices lying on the union of the above two subspaces. The statements above can be summarized in the following expression:

$$\text{span}(\mathcal{W}_i) \in \text{span}(\mathcal{D}_i \cup \mathcal{X}_i) \tag{15}$$

for any step size  $\mu_i$  and  $\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i)))$ . Since  $\text{rank}(\text{span}(\mathcal{W}_i)) \leq k$ , we easily deduce the following key observation: let  $\mathcal{S}_i \leftarrow \mathcal{D}_i \cup \mathcal{X}_i$  be a set of rank-1, orthonormal matrices where  $\text{rank}(\text{span}(\mathcal{S}_i)) \leq 2k$ . Given  $\mathcal{W}_i$  is unknown before the  $i$ -th iteration,  $\mathcal{S}_i$  spans the smallest subspace that

contains  $\mathcal{W}_i$  such that the following equality

$$\begin{aligned} &\mathcal{P}_k \left( \mathbf{X}(i) - \frac{\mu_i}{2} \nabla f(\mathbf{X}(i)) \right) \\ &= \mathcal{P}_k \left( \mathbf{X}(i) - \frac{\mu_i}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)) \right) \end{aligned} \tag{16}$$

necessarily holds.<sup>2</sup>

To compute step-size  $\mu_i$ , we use:

$$\begin{aligned} \mu_i &= \arg \min_{\mu} \left\| \mathbf{y} - \mathcal{A} \left( \mathbf{X}(i) - \frac{\mu}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)) \right) \right\|_2^2 \\ &= \frac{\| \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)) \|_F^2}{\| \mathcal{A} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)) \|_2^2}, \end{aligned} \tag{17}$$

i.e.,  $\mu_i$  is the minimizer of the objective function, given the current gradient  $\nabla f(\mathbf{X}(i))$ . Note that:

$$1 - \delta_{2k}(\mathcal{A}) \leq \frac{1}{\mu_i} \leq 1 + \delta_{2k}(\mathcal{A}), \tag{18}$$

due to R-RIP—i.e., we select  $2k$  subspaces such that  $\mu_i$  satisfies (18). We can derive similar arguments for the additional step size selection  $\xi_i$  in Step 6 of Algorithm 1.

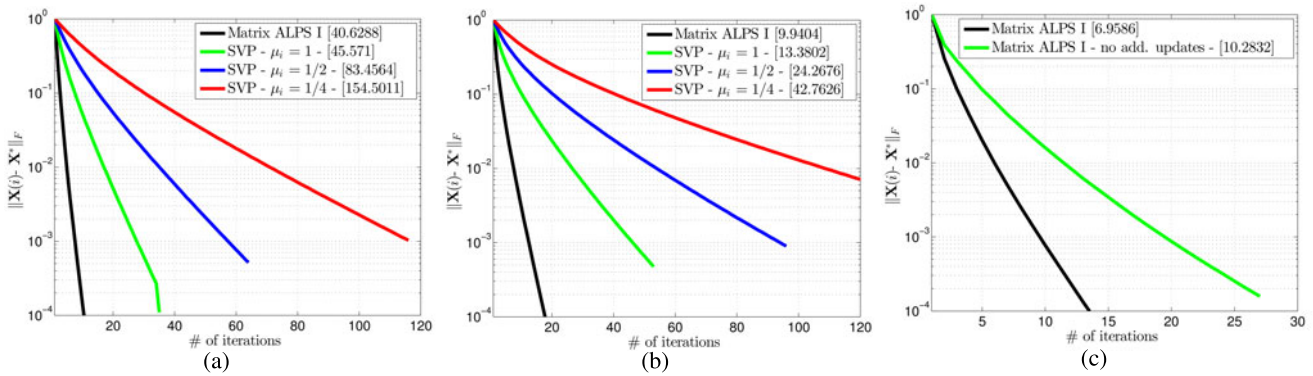
Adaptive  $\mu_i$  scheme results in more restrictive worst-case isometry constants compared to [3, 34, 41], but faster convergence and better stability are empirically observed in general. In [3], the authors present the Singular Value Projection (SVP) algorithm, an iterative hard thresholding algorithm for the ARM problem. According to [3], both constant and iteration dependent (but user-defined) step sizes are considered. Adaptive strategies presented in [3] require the computation of R-RIP constants which has exponential time complexity. Figures 1(a)–(b) illustrate some characteristic examples. The performance varies for different problem configurations. For  $\mu > 1$ , SVP *diverges* for various test cases. We note that, for large fixed matrix dimensions  $m, n$ , adaptive step size selection becomes computationally expensive compared to constant step size selection strategies, as the rank of  $\mathbf{X}^*$  increases.

#### 4.2 Updates on Restricted Subspaces

In Algorithm 1, at each iteration, the new estimate  $\mathbf{W}(i) \leftarrow \mathcal{P}_k(\mathbf{V}(i))$  can be further refined by applying a single or multiple gradient descent updates with line search restricted on  $\mathcal{W}_i$  [34] (Step 7 in Algorithm 1):

$$\mathbf{X}(i + 1) \leftarrow \mathbf{W}(i) - \frac{\xi_i}{2} \mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i)),$$

<sup>2</sup>In the case of multiple identical singular values, any ties are lexicographically dissolved.



**Fig. 1** Median error per iteration for various step size policies and 20 Monte-Carlo repetitions. In brackets, we present the median time consumed for convergence in seconds. (a)  $m = n = 2048$ ,  $p = 0.4n^2$ , and rank  $k = 70$ — $\mathcal{A}$  is formed by permuted and subsampled noiselets [40].

(b)  $n = 2048$ ,  $m = 512$ ,  $p = 0.4n^2$ , and rank  $k = 50$ —we use underdetermined linear map  $\mathcal{A}$  according to the MC problem (c)  $n = 2048$ ,  $m = 512$ ,  $p = 0.4n^2$ , and rank  $k = 40$ —we use underdetermined linear map  $\mathcal{A}$  according to the MC problem

where  $\xi_i = \frac{\|\mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i))\|_F^2}{\|\mathcal{A} \mathcal{P}_{\mathcal{W}_i} \nabla f(\mathbf{W}(i))\|_2^2}$ . In spirit, the gradient step above is the same as block coordinate descent in convex optimization where we find the subspaces adaptively. Figure 1(c) depicts the acceleration achieved by using additional gradient updates over restricted low-rank subspaces for a test case.

### 4.3 Acceleration via Memory-Based Schemes and Low-Rank Matrix Approximations

Memory-based techniques can be used to improve convergence speed. Furthermore, low-rank matrix approximation tools overcome the computational overhead of computing the best low-rank projection by inexactly solving (5). We keep the discussion on memory utilization for Sect. 7 and low-rank matrix approximations for Sects. 8 and 9 where we present new algorithmic frameworks for low-rank matrix recovery.

### 4.4 Active Low-Rank Subspace Tracking

Per iteration of Algorithms 1 and 2, we perform projection operations  $\mathcal{P}_{\mathcal{S}}\mathbf{X}$  and  $\mathcal{P}_{\mathcal{S}^\perp}\mathbf{X}$  where  $\mathbf{X} \in \mathbb{R}^{m \times n}$ , as described by (6) and (7), respectively. Since  $\mathcal{S}$  is constituted by outer products of left and right singular vectors as in Definition 2,  $\mathcal{P}_{\mathcal{S}}\mathbf{X}$  (resp.  $\mathcal{P}_{\mathcal{S}^\perp}\mathbf{X}$ ) projects onto the (resp. complement of the) best low-rank subspace in  $R(\mathbf{X})$  and  $R(\mathbf{X}^T)$ . These operations are highly connected with the adaptive step size selection and the updates on restricted subspaces. Unfortunately, the time-complexity to compute  $\mathcal{P}_{\mathcal{S}}\mathbf{X}$  is dominated by three matrix-matrix multiplications which decelerates the convergence of the proposed schemes in high-dimensional settings. To accelerate the convergence in many test cases, it turns out that we do not have to use the best projection  $\mathcal{P}_{\mathcal{S}}$

in practice.<sup>3</sup> Rather, employing *inexact* projections is sufficient to converge to the optimal solution: either (i)  $\mathcal{P}_{\mathcal{U}}\mathbf{X}$  onto the best low-rank subspace in  $R(\mathbf{X})$  only (if  $m \ll n$ ) or (ii)  $\mathbf{X}\mathcal{P}_{\mathcal{V}}$  onto the best low-rank subspace in  $R(\mathbf{X}^T)$  only (if  $m \gg n$ );<sup>4</sup>  $\mathcal{P}_{\mathcal{U}}$  and  $\mathcal{P}_{\mathcal{V}}$  are defined in Definition 2 and require only one matrix-matrix multiplication.

Figure 2 shows the time overhead due to the exact projection application  $\mathcal{P}_{\mathcal{S}}$  compared to  $\mathcal{P}_{\mathcal{U}}$  for  $m \leq n$ . In Fig. 2(a), we use subsampled and permuted noiselets for linear map  $\mathcal{A}$  and in Figs. 2(b)–(c), we test the MC problem. While in the case  $m = n$  the use of (6)–(7) has a clear advantage over inexact projections using only  $\mathcal{P}_{\mathcal{U}}$ , the latter case converges faster to the desired accuracy  $5 \times 10^{-4}$  when  $m \ll n$  as shown in Figs. 2(a)–(b). In our derivations, we assume  $\mathcal{P}_{\mathcal{S}}$  and  $\mathcal{P}_{\mathcal{S}^\perp}$  as defined in (6) and (7).

## 5 Convergence Guarantees

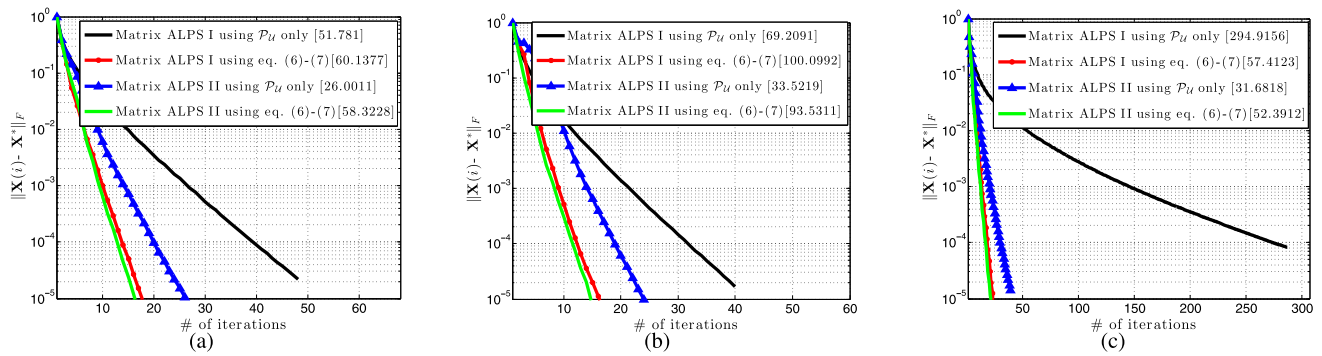
In this section, we present the theoretical convergence guarantees of Algorithms 1 and 2 as functions of R-RIP constants. To characterize the performance of the proposed algorithms, both in terms of convergence rate and noise resilience, we use the following recursive expression:

$$\|\mathbf{X}(i+1) - \mathbf{X}^*\|_F \leq \rho \|\mathbf{X}(i) - \mathbf{X}^*\|_F + \gamma \|\boldsymbol{\epsilon}\|_2. \tag{19}$$

In (19),  $\gamma$  denotes the approximation guarantee and provides insights into algorithm’s reconstruction capabilities when additive noise is present;  $\rho < 1$  expresses the convergence rate towards a region around  $\mathbf{X}^*$ , whose radius is determined

<sup>3</sup>From a different perspective and for a different problem case, similar ideas have been used in [18].

<sup>4</sup>We can move between these two cases by a simple transpose of the problem.



**Fig. 2** Median error per iteration for MATRIX ALPS I and MATRIX ALPS II variants over 10 Monte-Carlo repetitions. In brackets, we present the median time consumed for convergence in seconds. **(a)**  $n =$

2048,  $m = 512$ ,  $p = 0.25n^2$ , and rank  $k = 40$ . **(b)**  $n = 2000$ ,  $m = 1000$ ,  $p = 0.25n^2$ , and rank  $k = 50$ . **(c)**  $n = m = 1000$ ,  $p = 0.25n^2$ , and rank  $k = 50$

by  $\frac{\gamma}{1-\rho} \|\epsilon\|_2$ . In short, (19) characterizes how the distance to the true signal  $X^*$  is decreased and how the noise level affects the accuracy of the solution, at each iteration.

5.1 MATRIX ALPS I

An important lemma for our derivations below is given next:

**Lemma 6** [Active subspace expansion] *Let  $X(i)$  be the matrix estimate at the  $i$ -th iteration and let  $\mathcal{X}_i$  be a set of orthonormal, rank-1 matrices such that  $\mathcal{X}_i \leftarrow \mathcal{P}_k(X(i))$ . Then, at each iteration, the Active Subspace Expansion step in Algorithms 1 and 2 identifies information in  $X^*$ , such that:*

$$\begin{aligned} \|\mathcal{P}_{\mathcal{X}^*} \mathcal{P}_{\mathcal{S}_i^\perp} X^*\|_F &\leq (2\delta_{2k} + 2\delta_{3k}) \|X(i) - X^*\|_F \\ &\quad + \sqrt{2(1 + \delta_{2k})} \|\epsilon\|_2, \end{aligned} \tag{20}$$

where  $\mathcal{S}_i \leftarrow \mathcal{X}_i \cup \mathcal{D}_i$  and  $X^* \leftarrow \mathcal{P}_k(X^*)$ .

Lemma 6 states that, at each iteration, the active subspace expansion step identifies a  $2k$  rank subspace such that the amount of unrecovered energy of  $X^*$ —i.e., the projection of  $X^*$  onto the orthogonal subspace of  $\text{span}(\mathcal{S}_i)$ —is bounded by (20).

Then, Theorem 1 characterizes the iteration invariant of Algorithm 1 for the matrix case:

**Theorem 1** [Iteration invariant for MATRIX ALPS I] *The  $(i + 1)$ -th matrix estimate  $X(i + 1)$  of MATRIX ALPS I satisfies the following recursion:*

$$\|X(i + 1) - X^*\|_F \leq \rho \|X(i) - X^*\|_F + \gamma \|\epsilon\|_2, \tag{21}$$

where  $\rho := \left(\frac{1+2\delta_{2k}}{1-\delta_{2k}}\right) \left(\frac{4\delta_{2k}}{1-\delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1-\delta_{2k}}\right)$  and  $\gamma := \left(\frac{1+2\delta_{2k}}{1-\delta_{2k}}\right) \left(\frac{2\sqrt{1+\delta_{2k}}}{1-\delta_{2k}} + \frac{2\delta_{3k}}{1-\delta_{2k}} \sqrt{2(1 + \delta_{2k})}\right) + \frac{\sqrt{1+\delta_k}}{1-\delta_k}$ . Moreover, when  $\delta_{3k} < 0.1235$ , the iterations are contractive.

To provide some intuition behind this result, assume that  $X^*$  is a rank- $k$  matrix. Then, according to Theorem 1, for  $\rho < 1$ , the approximation parameter  $\gamma$  in (21) satisfies:

$$\gamma < 5.7624, \quad \text{for } \delta_{3k} < 0.1235.$$

Moreover, we derive the following:

$$\rho < \frac{1 + 2\delta_{3k}}{(1 - \delta_{3k})^2} (4\delta_{3k} + 8\delta_{3k}^2) < \frac{1}{2} \Rightarrow \delta_{3k} < 0.079,$$

which is a stronger R-RIP condition assumption compared to state-of-the-art approaches [21]. In the next section, we further improve this guarantee using Algorithm 2.

Unfolding the recursive formula (21), we obtain the following upper bound for  $\|X(i) - X^*\|_F$  at the  $i$ -th iteration:

$$\|X(i) - X^*\|_F \leq \rho^i \|X(0) - X^*\|_F + \frac{\gamma}{1 - \rho} \|\epsilon\|_2. \tag{22}$$

Then, given  $X(0) = \mathbf{0}$ , MATRIX ALPS I finds a rank- $k$  solution  $\hat{X} \in \mathbb{R}^{m \times n}$  such that  $\|\hat{X} - X^*\|_F \leq \frac{\gamma + 1 - \rho}{1 - \rho} \|\epsilon\|_2$  after  $i := \lceil \frac{\log(\|X^*\|_F / \|\epsilon\|_2)}{\log(1/\rho)} \rceil$  iterations.

If we ignore steps 5 and 6 in Algorithm 1, we obtain another projected gradient descent variant for the affine rank minimization problem, for which we obtain the following performance guarantees—the proof follows from the proof of Theorem 1.

**Corollary 1** [MATRIX ALPS I Instance] *In Algorithm 1, we ignore steps 5 and 6 and let  $\{\mathcal{X}_{i+1}, X(i + 1)\} \leftarrow \mathcal{P}_k(V_i)$ . Then, by the same analysis, we observe that the following recursion is satisfied:*

$$\|X(i + 1) - X^*\|_F \leq \rho \|X(i) - X^*\|_F + \gamma \|\epsilon\|_2, \tag{23}$$

for  $\rho := \left(\frac{4\delta_{2k}}{1-\delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1-\delta_{2k}}\right)$  and  $\gamma := \left(\frac{2\sqrt{1+\delta_{2k}}}{1-\delta_{2k}} + \frac{2\delta_{3k}}{1-\delta_{2k}} \sqrt{2(1 + \delta_{2k})}\right)$ . Moreover,  $\rho < 1$  when  $\delta_{3k} < 0.1594$ .



We observe that the absence of the additional estimate update over restricted support sets results in less restrictive isometry constants compared to Theorem 1. In practice, additional updates result in faster convergence, as shown in Fig. 1(c).

### 5.2 ADMiRA Instance

In MATRIX ALPS I, the gradient descent steps constitute a first-order approximation to least-squares minimization problems. Replacing Step 4 in Algorithm 1 with the following optimization problem:

$$V(i) \leftarrow \arg \min_{V: V \in \text{span}(\mathcal{S}_i)} \|y - \mathcal{A}V\|_2^2, \tag{24}$$

we obtain ADMiRA (furthermore, we remove the de-bias step in Algorithm 1). Assuming that the linear operator  $\mathcal{A}$ , restricted on sufficiently low-rank subspaces, is well conditioned in terms of the R-RIP assumption, the optimization problem (24) has a unique optimal minimizer. By exploiting the optimality condition in Lemma 1, ADMiRA instance in Algorithm 2 features the following guarantee:

**Theorem 2** [Iteration invariant for ADMiRA instance] *The  $(i + 1)$ -th matrix estimate  $X(i + 1)$  of ADMiRA answers the following recursive expression:*

$$\|X(i + 1) - X^*\|_F \leq \rho \|X(i) - X^*\|_F + \gamma \|\epsilon\|_F,$$

$$\rho := (2\delta_{2k} + 2\delta_{3k}) \sqrt{\frac{1+3\delta_{3k}^2}{1-\delta_{3k}^2}}, \text{ and } \gamma := \sqrt{\frac{1+3\delta_{3k}^2}{1-\delta_{3k}^2}} \sqrt{2(1 + \delta_{3k})} + \left(\frac{\sqrt{1+3\delta_{3k}^2}}{1-\delta_{3k}} + \sqrt{3}\right) \sqrt{1 + \delta_{2k}}.$$

Moreover, when  $\delta_{3k} < 0.2267$ , the iterations are contractive.

Similarly to MATRIX ALPS I analysis, the parameter  $\gamma$  in Theorem 2 satisfies:

$$\gamma < 5.1848, \quad \text{for } \delta_{3k} < 0.2267.$$

Furthermore, to compare the approximation guarantees of Theorem 2 with [21], we further observe:

$$\delta_{3k} < 0.1214, \quad \text{for } \rho < 1/2.$$

We remind that [21] provides convergence guarantees for ADMiRA with  $\delta_{4k} < 0.04$  for  $\rho = 1/2$ .

## 6 Complexity Analysis

In each iteration, computational requirements of the proposed hard thresholding methods mainly depend on the total number of linear mapping operations  $\mathcal{A}$ , gradient descent steps, least-squares optimizations, projection operations and

matrix decompositions for low rank approximation. Different algorithmic configurations (e.g. removing steps 6 and 7 in Algorithm 1) lead to hard thresholding variants with less computational complexity per iteration and better R-RIP conditions for convergence but a degraded performance in terms of stability and convergence speed is observed in practice. On the other hand, these additional processing steps increase the required time-complexity per iteration; hence, low iteration counts are desired to tradeoff these operations.

A non-exhaustive list of linear map examples includes the identity operator (Principal component analysis (PCA) problem), Fourier/Wavelets/Noiselets transformations and the famous Matrix Completion problem where  $\mathcal{A}$  is a mask operator such that only a fraction of elements in  $X$  is observed. Assuming the most demanding case where  $\mathcal{A}$  and  $\mathcal{A}^*$  are dense linear maps with no structure, the computation of the gradient  $\nabla f(X(i))$  at each iteration requires  $O(pkmn)$  arithmetic operations.

Given a set  $\mathcal{S}$  of orthonormal, rank-1 matrices, the projection  $\mathcal{P}_{\mathcal{S}}X$  for any matrix  $X \in \mathbb{R}^{m \times n}$  requires time complexity  $O(\max\{m^2n, mn^2\})$  as a sequence of matrix-matrix multiplication operations.<sup>5</sup> In MATRIX ALPS I, the adaptive step size selection steps require  $O(\max\{pkmn, m^2n\})$  time complexity for the calculation of  $\mu_i$  and  $\xi_i$  quantities. In ADMiRA solving a least-squares system restricted on rank- $2k$  and rank- $k$  subspaces requires  $O(pk^2)$  complexity; according to [21, 32], the complexity of this step can be further reduced using iterative techniques such as the Richardson method or conjugate gradients algorithm.

Using the Lanczos method, we require  $O(kmn)$  arithmetic operations to compute a rank- $k$  matrix approximation for a given constant accuracy; a prohibitive time-complexity that does not scale well for many practical applications. Sections 8 and 9 describe approximate low rank matrix projections and how they affect the convergence guarantees of the proposed algorithms.

Overall, the operation that dominates per iteration requires  $O(\max\{pkmn, m^2n, mn^2\})$  time complexity in the proposed schemes.

## 7 Memory-Based Acceleration

Iterative algorithms can use memory to gain momentum in convergence. Based on Nesterov’s optimal gradient methods [42], we propose a hard thresholding variant, described in Algorithm 3 where an additional update on  $X(i + 1)$  with momentum step size  $\tau_i$  is performed using previous matrix estimates.

<sup>5</sup>While such operation has  $O(\max\{m^2n, mn^2\})$  complexity, each application of  $\mathcal{P}_{\mathcal{S}}X$  requires three matrix-matrix multiplications. To reduce such computational cost, we relax this operation in Sect. 10 where in practice we use only  $\mathcal{P}_{\mathcal{L}}$  that needs one matrix-matrix multiplication.

**Algorithm 3** MATRIX ALPS II

**Input:**  $\mathbf{y}$ ,  $\mathcal{A}$ ,  $k$ , Tolerance  $\eta$ , MaxIterations

**Initialize:**  $\mathbf{X}(0) \leftarrow 0$ ,  $\mathcal{X}_0 \leftarrow \{\emptyset\}$ ,  $\mathbf{Q}(0) \leftarrow 0$ ,  $\mathcal{Q}_0 \leftarrow \{\emptyset\}$ ,  $\tau_i \forall i, i \leftarrow 0$

**repeat**

- 1:  $\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{Q}_i^\perp} \nabla f(\mathbf{Q}(i)))$  (Best rank- $k$  subspace orthogonal to  $\mathcal{Q}_i$ )
  - 2:  $\mathcal{S}_i \leftarrow \mathcal{D}_i \cup \mathcal{Q}_i$  (Active subspace expansion)
  - 3:  $\mu_i \leftarrow \arg \min_{\mu} \|\mathbf{y} - \mathcal{A}(\mathbf{Q}(i) - \frac{\mu}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i)))\|_2^2 = \frac{\|\mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))\|_F^2}{\|\mathcal{A} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))\|_2^2}$  (Step size selection)
  - 4:  $\mathbf{V}(i) \leftarrow \mathbf{Q}(i) - \frac{\mu_i}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))$  (Error norm reduction via gradient descent)
  - 5:  $\{\mathcal{X}_{i+1}, \mathbf{X}(i+1)\} \leftarrow \mathcal{P}_k(\mathbf{V}(i))$  (Best rank- $k$  subspace selection)
  - 6:  $\mathbf{Q}(i+1) \leftarrow \mathbf{X}(i+1) + \tau_i(\mathbf{X}(i+1) - \mathbf{X}(i))$  (Momentum update)
  - 7:  $\mathcal{Q}_{i+1} \leftarrow \text{ortho}(\mathcal{X}_i \cup \mathcal{X}_{i+1})$   
 $i \leftarrow i + 1$
- until**  $\|\mathbf{X}(i) - \mathbf{X}(i-1)\|_2 \leq \eta \|\mathbf{X}(i)\|_2$  or MaxIterations.

Similar to  $\mu_i$  strategies,  $\tau_i$  can be preset as constant or adaptively computed at each iteration. Constant momentum step size selection has no additional computational cost but convergence rate acceleration is not guaranteed for some problem formulations in practice. On the other hand, empirical evidence has shown that adaptive  $\tau_i$  selection strategies result in faster convergence compared to zero-memory methods with *similar complexity*.

For the case of strongly convex objective functions, Nesterov [43] proposed the following constant momentum step size selection scheme:  $\tau_i = \frac{\alpha_i(1-\alpha_i)}{\alpha_i^2 + \alpha_{i+1}}$ , where  $\alpha_0 \in (0, 1)$  and  $\alpha_{i+1}$  is computed as the root  $\in (0, 1)$  of

$$\alpha_{i+1}^2 = (1 - \alpha_{i+1})\alpha_i^2 + q\alpha_{i+1}, \quad \text{for } q \triangleq \frac{1}{\kappa^2(\mathcal{A})}, \quad (25)$$

where  $\kappa(\mathcal{A})$  denotes the condition number of  $\mathcal{A}$ . In this scheme, exact calculation of  $q$  parameter is computationally expensive for large-scale data problems and approximation schemes are leveraged to compensate this complexity bottleneck.

Based upon adaptive  $\mu_i$  selection, we propose to select  $\tau_i$  as the minimizer of the objective function:

$$\begin{aligned} \tau_i &= \arg \min_{\tau} \|\mathbf{y} - \mathcal{A}\mathbf{Q}(i+1)\|_2^2 \\ &= \frac{\langle \mathbf{y} - \mathcal{A}\mathbf{X}(i), \mathcal{A}\mathbf{X}(i) - \mathcal{A}\mathbf{X}(i-1) \rangle}{\|\mathcal{A}\mathbf{X}(i) - \mathcal{A}\mathbf{X}(i-1)\|_2^2}, \end{aligned} \quad (26)$$

where  $\mathcal{A}\mathbf{X}(i)$ ,  $\mathcal{A}\mathbf{X}(i-1)$  are already *pre-computed* at each iteration. According to (26),  $\tau_i$  is dominated by the calculation of a vector inner product, a computationally cheaper process than  $q$  calculation.

Theorem 3 characterizes Algorithm 3 for *constant* momentum step size selection. To keep the main ideas simple, we ignore the additional gradient updates in Algorithm 3. In addition, we only consider the noiseless case for clarity. The convergence rate proof for these cases is provided in the [Appendix](#).

**Theorem 3** [Iteration invariant for MATRIX ALPS II] *Let  $\mathbf{y} = \mathcal{A}\mathbf{X}^*$  be a noiseless set of observations. To recover  $\mathbf{X}^*$  from  $\mathbf{y}$  and  $\mathcal{A}$ , the  $(i+1)$ -th matrix estimate  $\mathbf{X}(i+1)$  of MATRIX ALPS II satisfies the following recursion:*

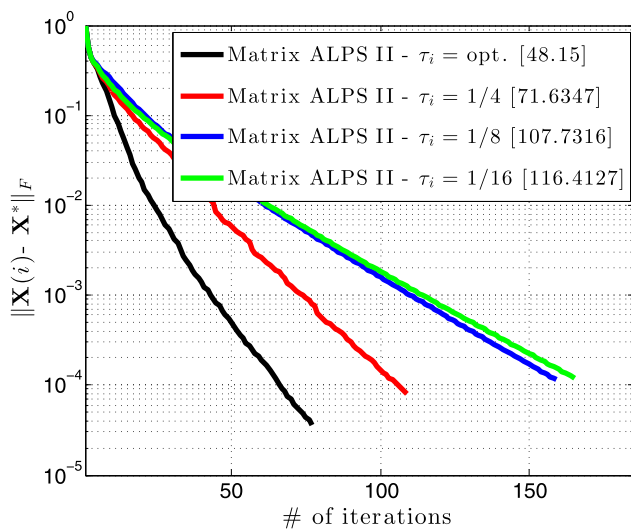
$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F &\leq \alpha(1 + \tau_i) \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\ &\quad + \alpha\tau_i \|\mathbf{X}(i-1) - \mathbf{X}^*\|_F, \end{aligned} \quad (27)$$

where  $\alpha := \frac{4\delta_{3k}}{1-\delta_{3k}} + (2\delta_{3k} + 2\delta_{4k})\frac{2\delta_{3k}}{1-\delta_{3k}}$ . Moreover, solving the above second-order recurrence, the following inequality holds true:

$$\|\mathbf{X}(i+1) - \mathbf{X}^*\|_F \leq \rho^{i+1} \|\mathbf{X}(0) - \mathbf{X}^*\|_F, \quad (28)$$

for  $\rho := \frac{\alpha(1+\tau_i) + \sqrt{\alpha^2(1+\tau_i)^2 + 4\alpha\tau_i}}{2}$ .

Theorem 3 provides convergence rate behaviour proof for the case where  $\tau_i$  is constant  $\forall i$ . The more elaborate case where  $\tau_i$  follows the policy described in (26) is left as an open question for future work. To provide some insight for (28), for  $\tau_i = 1/4, \forall i$  and  $\tau_i = 1/2, \forall i$ ,  $\delta_{4k} < 0.1187$  and  $\delta_{4k} < 0.095$  guarantee convergence in Algorithm 3, respectively. While the RIP requirements for memory-based MATRIX ALPS II are more stringent than the schemes proposed in the previous section, it outperforms Algorithms 1 and 2. Figure 2 shows the acceleration achieved in MATRIX ALPS II by using inexact projections  $\mathcal{P}_{\mathcal{U}}$ . Using the proper projections (6)–(7), Fig. 3 shows acceleration in practice when using the adaptive momentum step size strategy: while a wide range of constant momentum step sizes leads to convergence, providing flexibility to select an appropriate  $\tau_i$ , adaptive  $\tau_i$  avoids this arbitrary  $\tau_i$  selection while further decreases the number of iterations needed for convergence in most cases.



**Fig. 3** Median error per iteration for various momentum step size policies and 10 Monte-Carlo repetitions. Here,  $n = 1024$ ,  $m = 256$ ,  $p = 0.25n^2$ , and rank  $k = 40$ . We use permuted and subsampled noiselets for the linear map  $\mathcal{A}$ . In brackets, we present the median time for convergence in seconds

### 8 Accelerating MATRIX ALPS: $\epsilon$ -Approximation of SVD via Column Subset Selection

A time-complexity bottleneck in the proposed schemes is the computation of the singular value decomposition to find subspaces that describe the unexplored information in matrix  $X^*$ . Unfortunately, the computational cost of regular SVD for best subspace tracking is prohibitive for many applications.

Based on [44, 45], we can obtain randomized SVD approximations of a matrix  $X$  using *column subset selection* ideas: we compute a leverage score for each column that represents its “significance”. In particular, we define a probability distribution that weights each column depending on the amount of information they contain; usually, the distribution is related to the  $\ell_2$ -norm of the columns. The main idea of this approach is to compute a surrogate rank- $k$  matrix  $\mathcal{P}_k^\epsilon(X)$  by subsampling the columns according to this distribution. It turns out that the total number of sampled columns is a function of the parameter  $\epsilon$ . Moreover, [46, 47] proved that, given a target rank  $k$  and an approximation parameter  $\epsilon$ , we can compute an  $\epsilon$ -approximate rank- $k$  matrix  $\mathcal{P}_k^\epsilon(X)$  according to the following definition.

**Definition 4** [ $\epsilon$ -approximate low-rank projection] Let  $X$  be an arbitrary matrix. Then,  $\mathcal{P}_k^\epsilon(X)$  projection provides a rank- $k$  matrix approximation to  $X$  such that:

$$\|\mathcal{P}_k^\epsilon(X) - X\|_F^2 \leq (1 + \epsilon) \|\mathcal{P}_k(X) - X\|_F^2, \tag{29}$$

where  $\mathcal{P}_k(X) \in \arg \min_{Y: \text{rank}(Y) \leq k} \|X - Y\|_F$ .

For the following theoretical results, we assume the following condition on the sensing operator  $\mathcal{A}$ :  $\|\mathcal{A}^* \beta\|_F \leq \lambda$ ,  $\forall \beta \in \mathbb{R}^p$  where  $\lambda > 0$ . Using  $\epsilon$ -approximation schemes to perform the Active subspace selection step, the following upper bound holds. The proof is provided in the Appendix:

**Lemma 7** [ $\epsilon$ -approximate active subspace expansion] Let  $X(i)$  be the matrix estimate at the  $i$ -th iteration and let  $\mathcal{X}_i$  be a set of orthonormal, rank-1 matrices in  $\mathbb{R}^{m \times n}$  such that  $\mathcal{X}_i \leftarrow \mathcal{P}_k(X(i))$ . Furthermore, let

$$\mathcal{D}_i^\epsilon \leftarrow \mathcal{P}_k^\epsilon(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(X(i))),$$

be a set of orthonormal, rank-1 matrices that span rank- $k$  subspace such that (29) is satisfied for  $X := \mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(X(i))$ . Then, at each iteration, the Active Subspace Expansion step in Algorithms 1 and 2 captures information contained in the true matrix  $X^*$ , such that:

$$\begin{aligned} & \|\mathcal{P}_{\mathcal{X}^*} \mathcal{P}_{\mathcal{S}_i^\perp} X^*\|_F \\ & \leq (2\delta_{2k} + 2\delta_{3k}) \|X(i) - X^*\|_F + \sqrt{2(1 + \delta_{2k})} \|\epsilon\|_2 \\ & \quad + 2\lambda\sqrt{\epsilon}, \end{aligned} \tag{30}$$

where  $\mathcal{S}_i \leftarrow \mathcal{X}_i \cup \mathcal{D}_i^\epsilon$  and  $\mathcal{X}^* \leftarrow \mathcal{P}_k(X^*)$ .

Furthermore, to prove the following theorems, we extend Lemma 10, provided in the Appendix, as follows. The proof easily follows from the proof of Lemma 10, using Definition 4:

**Lemma 8** [ $\epsilon$ -approximation rank- $k$  subspace selection] Let  $V(i)$  be a rank- $2k$  proxy matrix in the subspace spanned by  $\mathcal{S}_i$  and let  $\widehat{W}(i) \leftarrow \mathcal{P}_k^\epsilon(V(i))$  denote the rank- $k$   $\epsilon$ -approximation to  $V(i)$ , according to (5). Then:

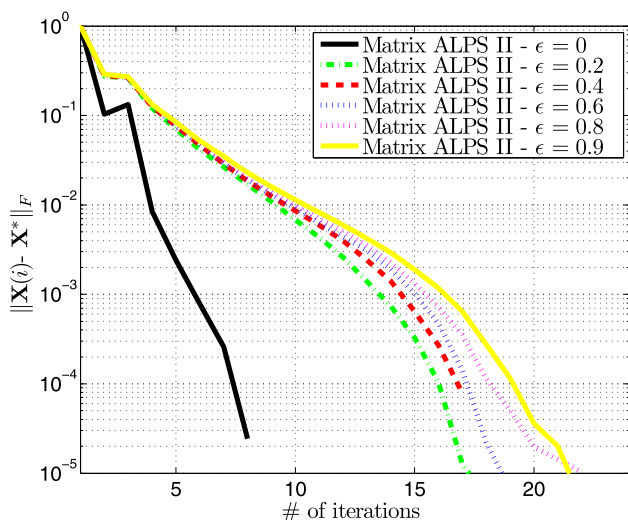
$$\begin{aligned} \|\widehat{W}(i) - V(i)\|_F^2 & \leq (1 + \epsilon) \|W(i) - V(i)\|_F^2 \\ & \leq (1 + \epsilon) \|\mathcal{P}_{\mathcal{S}_i}(V(i) - X^*)\|_F^2 \\ & \leq (1 + \epsilon) \|V(i) - X^*\|_F^2 \end{aligned} \tag{31}$$

where  $W(i) \leftarrow \mathcal{P}_k(V(i))$ .

#### 8.1 MATRIX ALPS I Using $\epsilon$ -Approximate Low-Rank Projection via Column Subset Selection

Using  $\epsilon$ -approximate SVD in MATRIX ALPS I, the following iteration invariant theorem holds:

**Theorem 4** [Iteration invariant with  $\epsilon$ -approximate projections for MATRIX ALPS I] The  $(i + 1)$ -th matrix estimate  $X(i + 1)$  of MATRIX ALPS I with  $\epsilon$ -approximate projections  $\mathcal{D}_i^\epsilon \leftarrow \mathcal{P}_k^\epsilon(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(X(i)))$  and  $\widehat{W}(i) \leftarrow \mathcal{P}_k^\epsilon(V(i))$  in



**Fig. 4** Performance comparison using  $\epsilon$ -approximation SVD [47] in MATRIX ALPS II.  $m = n = 256$ ,  $p = 0.4n^2$ , rank of  $X^*$  equals 2 and  $\mathcal{A}$  constituted by permuted noiselets. The non-smoothness in the error curves is due to the extreme low rankness of  $X^*$  for this setting

Algorithm 1 satisfies the following recursion:

$$\|X(i + 1) - X^*\|_F \leq \rho \|X(i) - X^*\|_F + \gamma \|\epsilon\|_2 + \beta\lambda, \tag{32}$$

where  $\rho := (1 + \frac{3\delta_k}{1-\delta_k})(2 + \epsilon)[(1 + \frac{\delta_{3k}}{1-\delta_{2k}})4\delta_{3k} + \frac{2\delta_{2k}}{1-\delta_{2k}}]$ ,  $\beta := (1 + \frac{3\delta_k}{1-\delta_k})(2 + \epsilon)(1 + \frac{\delta_{3k}}{1-\delta_{2k}})2\sqrt{\epsilon}$ , and  $\gamma := (1 + \frac{3\delta_k}{1-\delta_k})(2 + \epsilon)[(1 + \frac{\delta_{3k}}{1-\delta_{2k}})\sqrt{2(1 + \delta_{2k})} + 2\frac{\sqrt{1+\delta_{2k}}}{1-\delta_{2k}}]$ .

Similar analysis can be conducted for the ADMiRA algorithm. To illustrate the impact of SVD  $\epsilon$ -approximation on the signal reconstruction performance of the proposed methods, we replace the best rank- $k$  projections in steps 1 and 5 of Algorithm 1 by the  $\epsilon$ -approximation SVD algorithm, presented in [47]. In this paper, the column subset selection algorithm satisfies the following theorem:

**Theorem 5** Let  $X \in \mathbb{R}^{m \times n}$  be a signal of interest with arbitrary rank  $< \min\{m, n\}$  and let  $X_k$  represent the best rank- $k$  approximation of  $X$ . After  $2(k + 1)(\log(k + 1) + 1)$  passes over the data, the Linear Time Low-Rank Matrix Approximation algorithm in [47] computes a rank- $k$  approximation  $\mathcal{P}_k^\epsilon(X) \in \mathbb{R}^{m \times n}$  such that Definition 4 is satisfied with probability at least  $3/4$ .

The proof is provided in [47]. In total, Linear Time Low-Rank Matrix Approximation algorithm [47] requires  $O(mn(k/\epsilon + k^2 \log k) + (m + n)(k^2/\epsilon^2 + k^3 \log k/\epsilon + k^4 \log^2 k))$  and  $O(\min\{m, n\}(k/\epsilon + k^2 \log k))$  time and space complexity, respectively. However, while column subset selection methods such as [47] reduce the overall complexity of low-rank projections in theory, in practice this ap-

plies only in very high-dimensional settings. To strengthen this argument, in Fig. 4 we compare SVD-based MATRIX ALPS II with MATRIX ALPS II using the  $\epsilon$ -approximate column subset selection method in [47]. We observe that the total number of iterations for convergence increases due to  $\epsilon$ -approximate low-rank projections, as expected. Nevertheless, we observe that, on average, the column subset selection process [47] is computationally prohibitive compared to regular SVD due to the time overhead in the column selection procedure—fewer passes over the data are desirable in practice to tradeoff the increased number of iterations for convergence. In the next section, we present alternatives based on recent trends in randomized matrix decompositions and how we can use them in low-rank recovery.

### 9 Accelerating MATRIX ALPS: SVD Approximation Using Randomized Matrix Decompositions

Finding low-cost SVD approximations to tackle the above complexity issues is a challenging task. Recent works on probabilistic methods for matrix approximation [26] provide a family of efficient approximate projections on the set of rank-deficient matrices with clear computational advantages over regular SVD computation in practice and attractive theoretical guarantees. In this work, we build on the low-cost, power-iteration subspace tracking scheme, described in Algorithms 4.3 and 4.4 in [26]. Our proposed algorithm is described in Algorithm 4.

The convergence guarantees of Algorithm 4 follow the same motions described in Sect. 8, where  $\epsilon$  is a function of  $m, n, k$  and  $q$ .

## 10 Experiments

### 10.1 List of Algorithms

In the following experiments, we compare the following algorithms: (i) the Singular Value Projection (SVP) algorithm [3], a non-convex first-order projected gradient descent algorithm with constant step size selection (we study the case where  $\mu = 1$ ), (ii) the inexact ALM algorithm [18] based on augmented Lagrange multiplier method, (iii) the OptSpace algorithm [48], a gradient descent algorithm on the Grassmann manifold, (iv) the Grassmannian Rank-One Update Subspace Estimation (GROUSE) and the Grassmannian Robust Adaptive Subspace Tracking methods (GRASTA) [49, 50], two stochastic gradient descent algorithms that operate on the Grassmannian—moreover, to allay the impact of outliers in the subspace selection step, GRASTA incorporates the augmented Lagrangian of  $\ell_1$ -norm loss function into the Grassmannian optimization framework, (v) the Riemannian



**Algorithm 4** Randomized MATRIX ALPS II with QR Factorization

**Input:**  $\mathbf{y}$ ,  $\mathcal{A}$ ,  $k$ ,  $q$ , Tolerance  $\eta$ , MaxIterations

**Initialize:**  $\mathbf{X}(0) \leftarrow 0$ ,  $\mathcal{X}_0 \leftarrow \{\emptyset\}$ ,  $\mathbf{Q}(0) \leftarrow 0$ ,  $\mathcal{Q}_0 \leftarrow \{\emptyset\}$ ,  $\tau_i \forall i, i \leftarrow 0$

**repeat**

- 1:  $\mathcal{D}_i \leftarrow \text{RANDOMIZEDPOWERITERATION}(\mathcal{P}_{\mathcal{Q}_i} \nabla f(\mathbf{Q}(i)), k, q)$  (*Rank- $k$  subspace via Randomized Power Iteration*)
- 2:  $\mathcal{S}_i \leftarrow \mathcal{D}_i \cup \mathcal{Q}_i$  (*Active subspace expansion*)
- 3:  $\mu_i \leftarrow \arg \min_{\mu} \|\mathbf{y} - \mathcal{A}(\mathbf{Q}(i) - \frac{\mu}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i)))\|_2^2 = \frac{\|\mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))\|_F^2}{\|\mathcal{A} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))\|_2^2}$  (*Step size selection*)
- 4:  $\mathbf{V}(i) \leftarrow \mathbf{Q}(i) - \frac{\mu_i}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{Q}(i))$  (*Error norm reduction via gradient descent*)
- 5:  $\mathcal{W} \leftarrow \text{RANDOMIZEDPOWERITERATION}(\mathbf{V}(i), k, q)$  (*Rank- $k$  subspace via Randomized Power Iteration*)
- 6:  $\mathbf{X}(i+1) \leftarrow \mathcal{P}_{\mathcal{W}} \mathbf{V}(i)$  (*Best rank- $k$  subspace selection*)
- 7:  $\mathbf{Q}(i+1) \leftarrow \mathbf{X}(i+1) + \tau_i(\mathbf{X}(i+1) - \mathbf{X}(i))$  (*Momentum update*)
- 8:  $\mathcal{Q}_{i+1} \leftarrow \text{ortho}(\mathcal{X}_i \cup \mathcal{X}_{i+1})$   
 $i \leftarrow i + 1$

**until**  $\|\mathbf{X}(i) - \mathbf{X}(i-1)\|_2 \leq \eta \|\mathbf{X}(i)\|_2$  or MaxIterations.

Trust Region Matrix Completion algorithm (RTRMC) [51], a matrix completion method using first- and second-order Riemannian trust-region approaches, (vi) the Low rank Matrix Fitting algorithm (LMaFit) [52], a nonlinear successive over-relaxation algorithm and (vii) the algorithms MATRIX ALPS I, ADMiRA [21], MATRIX ALPS II and Randomized MATRIX ALPS II with QR Factorization (referred shortly as MATRIX ALPS II with QR) presented in this paper.

10.2 Implementation Details

To properly compare the algorithms in the above list, we present a set of parameters that are common. We denote the ratio between the number of observed samples and the number of variables in  $\mathbf{X}^*$  as  $\text{SR} := p/(m \cdot n)$  (sampling ratio). Furthermore, we reserve FR to represent the degree of freedom in a rank- $k$  matrix to the number of observations—this corresponds to the following definition  $\text{FR} := (k(m+n-k))/p$ . In most of the experiments, we fix the number of observable data  $p = 0.3mn$  and vary the dimensions and the rank  $k$  of the matrix  $\mathbf{X}^*$ . This way, we create a wide range of different problem configurations with variable FR.

Most of the algorithms in comparison as well as the proposed schemes are implemented in MATLAB. We note that the LMaFit software package contains parts implemented in C that reduce the per iteration computational time. This provides insights for further time savings in our schemes; we leave a fully optimized implementation of our algorithms as future work. In this paper, we mostly test cases where  $m \ll n$ . Such settings can be easily found in real-world problems such as recommender systems (e.g. Netflix, Amazon, etc.) where the number of products, movies, etc. is much greater than the number of active users.

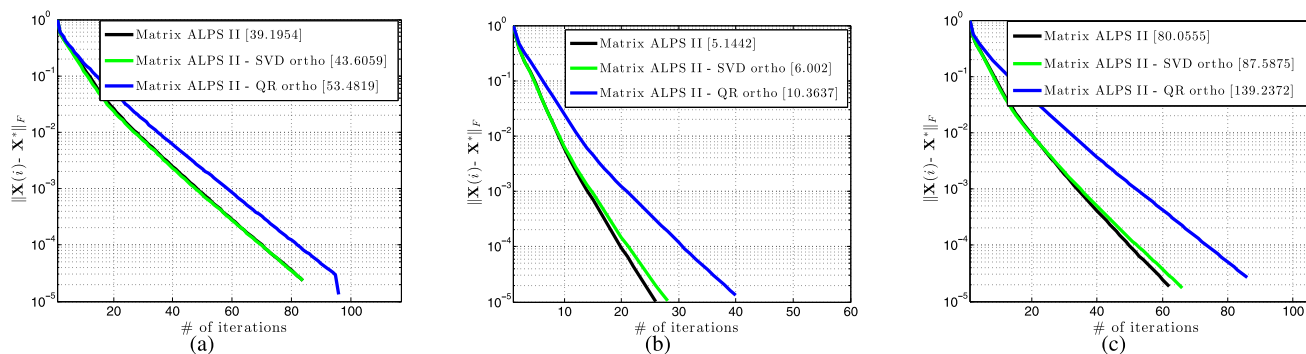
In all algorithms, we fix the maximum number of iterations to 500, unless otherwise stated. To solve a least

squares problem over a restricted low-rank subspace, we use conjugate gradients with maximum number of iterations given by  $\text{cg\_maxiter} := 500$  and tolerance parameter  $\text{cg\_tol} := 10^{-10}$ . We use the same stopping criteria for the majority of algorithms under consideration:

$$\frac{\|\mathbf{X}(i) - \mathbf{X}(i-1)\|_F}{\|\mathbf{X}(i)\|_F} \leq \text{tol}, \tag{33}$$

where  $\mathbf{X}(i)$ ,  $\mathbf{X}(i-1)$  denote the current and the previous estimate of  $\mathbf{X}^*$  and  $\text{tol} := 5 \times 10^{-5}$ . If this is not the case, we tweak the algorithms to minimize the total execution time and achieve similar reconstruction performance as the rest of the algorithms. For SVD calculations, we use the lansvd implementation in PROPACK package [53]—moreover, all the algorithms in comparison use the same linear operators  $\mathcal{A}$  and  $\mathcal{A}^*$  for gradient and SVD calculations and conjugate-gradient least-squares minimizations. For fairness, we modified all the algorithms so that they *exploit the true rank*. Small deviations from the true rank result in relatively small degradation in terms of the reconstruction performance. In case the rank of  $\mathbf{X}^*$  is unknown, one has to predict the dimension of the principal singular space. The authors in [3], based on ideas in [48], propose to compute singular values incrementally until a significant gap between singular values is found. Similar strategies can be found in [18] for the convex case.

In MATRIX ALPS II and MATRIX ALPS II with QR, we perform  $\mathcal{Q}_i \leftarrow \text{ortho}(\mathcal{X}_i \cup \mathcal{X}_{i+1})$  to construct a set of orthonormal rank-1 matrices that span the subspace, spanned by  $\mathcal{X}_i \cup \mathcal{X}_{i+1}$ . While such operation can be implemented using factorization procedures (such as SVD or QR decompositions), in practice this degrades the time complexity of the algorithm substantially as the rank  $k$  and the problem dimensionality increase. In our implementations, we simply *union* the set of orthonormal rank-1 matrices,



**Fig. 5** Median error per iteration for MATRIX ALPS II variants over 10 Monte-Carlo repetitions. In brackets, we present the mean time consumed for convergence in seconds. (a)  $n = 1024, m = 256, p = 0.25n^2$ ,

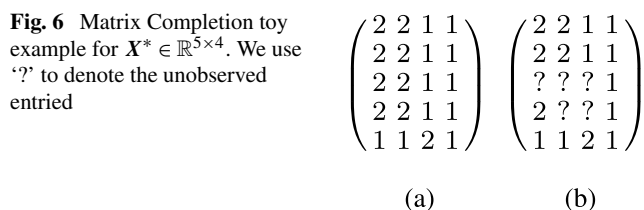
and rank  $k = 20$ . (b)  $n = 2048, m = 512, p = 0.25n^2$ , and rank  $k = 60$ . (c)  $n = 1000, m = 500, p = 0.25n^2$ , and rank  $k = 50$

without further orthogonalization. Thus, we employ *inexact* projections for computational efficiency which results in faster convergence. Figure 5 shows the time overhead due to the additional orthogonalization process. We compare three algorithms: MATRIX ALPS II (no orthogonalization step), MATRIX ALPS II using SVD for orthogonalization and, MATRIX ALPS II using QR for orthogonalization. In Figs. 5(a)–(b), we use subsampled and permuted noiselets for linear map  $\mathcal{A}$  and in Fig. 5(c), we test the MC problem. In all the experimental cases considered in this work, we observed identical performance in terms of reconstruction accuracy for the three variants, as can be also seen in Fig. 5. To this end, for the rest of the paper, we use MATRIX ALPS II where  $\mathcal{Q}_i \leftarrow \mathcal{X}_i \cup \mathcal{X}_{i+1}$ .

### 10.3 Limitations of $\|\cdot\|_*$ -Based Algorithms: A Toy Example

While nuclear norm heuristic is widely used in solving the low-rank minimization problem, [54] presents simple problem cases where convex, nuclear norm-based, algorithms *fail* in practice. Using the  $\|\cdot\|_*$ -norm in the objective function as the convex surrogate of the  $\text{rank}(\cdot)$  metric might lead to a candidate set with multiple solutions, introducing ambiguity in the selection process. Borrowing the example in [54], we test the list of algorithms above on a toy problem setting that does not satisfy the rank-RIP. To this end, we design the following problem: let  $X^* \in \mathbb{R}^{5 \times 4}$  be the matrix of interest with  $\text{rank}(X^*) = 2$ , as shown in Fig. 6(a). We consider the case where we have access to  $X^*$  only through a subset of its entries, as shown in Fig. 6(b).

In Fig. 7, we present the reconstruction performance of various matrix completion solvers after 300 iterations. Although there are multiple solutions that induce the recovered matrix and have the same rank as  $X^*$ , most of the algorithms in comparison reconstruct  $X^*$  successfully. We note that, in some cases, the inadequacy of an algorithm to reconstruct



$X^*$  is not because of the (relaxed) problem formulation but due to its fast—but inaccurate—implementation (fast convergence versus reconstruction accuracy tradeoff).

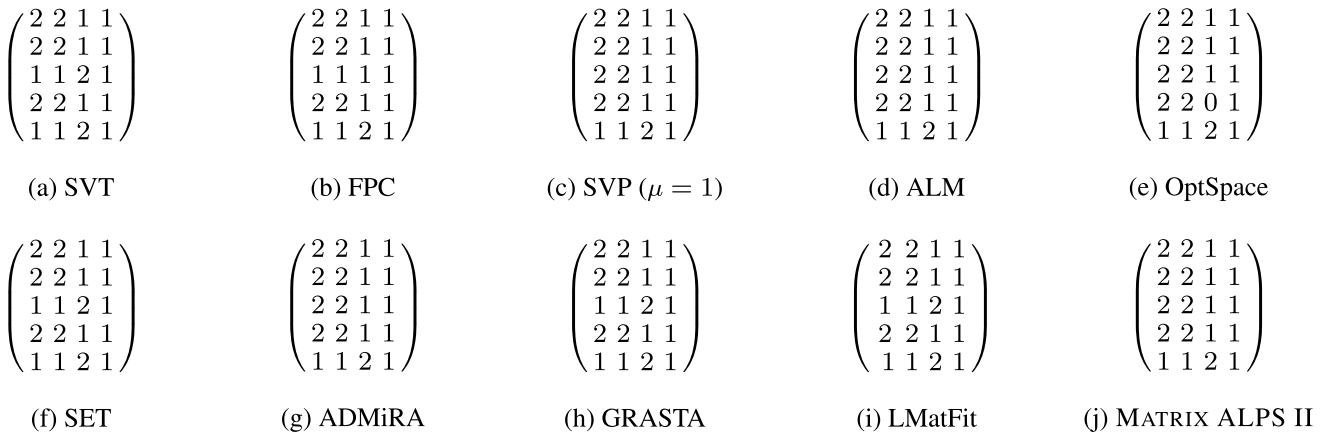
### 10.4 Synthetic Data

**General affine rank minimization using noiselets:** In this experiment, the set of observations  $y \in \mathbb{R}^p$  satisfy:

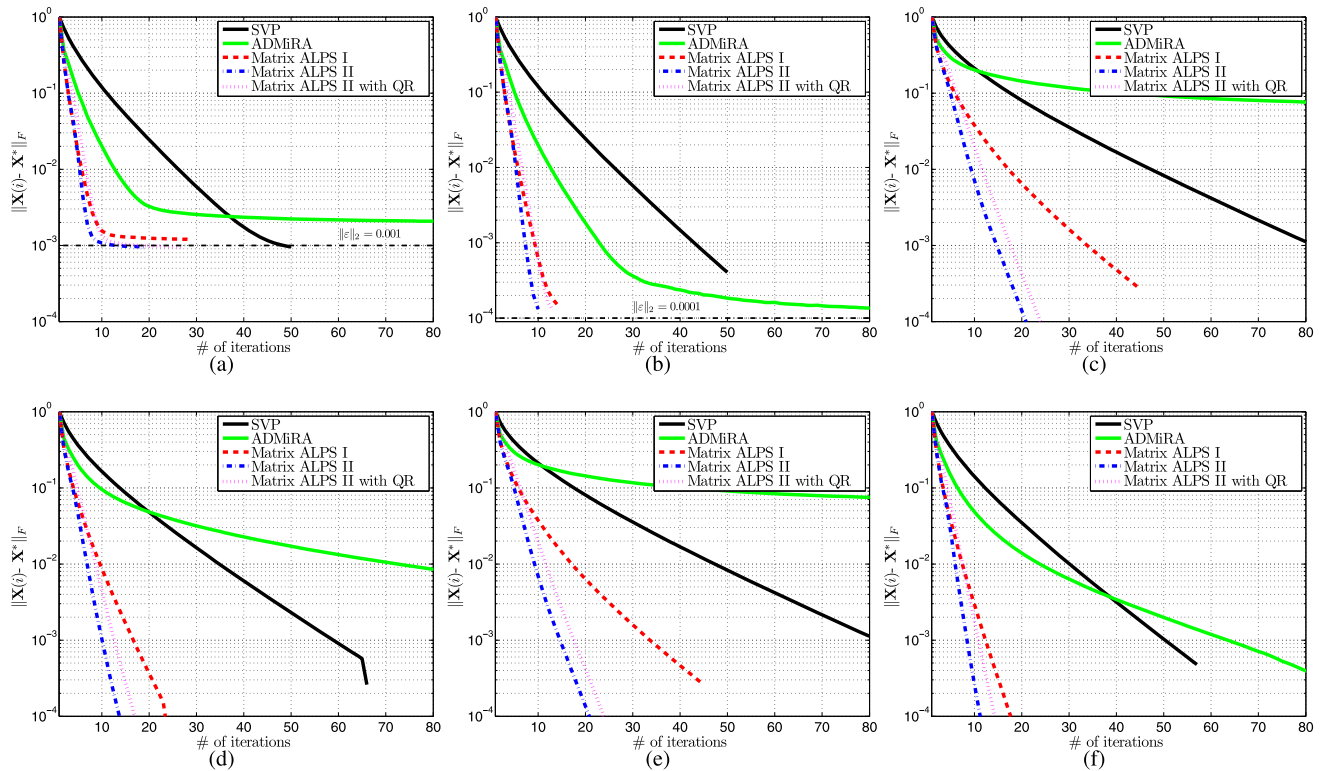
$$y = \mathcal{A}X^* + \epsilon. \tag{34}$$

Here, we use permuted and subsampled noiselets for the linear operator  $\mathcal{A}$  [12]. The signal  $X^*$  is generated as the multiplication of two low-rank matrices,  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{n \times k}$ , such that  $X^* = LR^T$  and  $\|X^*\|_F = 1$ . Both  $L$  and  $R$  have random independent and identically distributed (iid) Gaussian entries with zero mean and unit variance. In the noisy case, the additive noise term  $\epsilon \in \mathbb{R}^p$  contains entries drawn from a zero mean Gaussian distribution with  $\|\epsilon\|_2 \in \{10^{-3}, 10^{-4}\}$ .

We compare the following algorithms: SVP, ADMiRA, MATRIX ALPS I, MATRIX ALPS II and MATRIX ALPS II with QR for various problem configurations, as depicted in Table 1 (there is no available code with arbitrary sensing operators for the rest algorithms). In Table 1, we show the median values of reconstruction error, number of iterations and execution time over 50 Monte Carlo iterations. For all cases, we assume  $\text{SR} = 0.3$  and we set the maximum number of iterations to 500. Bold font denotes the fastest execution time. Furthermore, Fig. 8 illustrates the effectiveness of the algorithms for some representative problem configurations.



**Fig. 7** Toy example reconstruction performance for various algorithms. We observe that  $X^*$  is an integer matrix—since the algorithms under consideration return real matrices as solutions, we round the solution elementwise



**Fig. 8** Low rank signal reconstruction using noiselet linear operator. The error curves are the median values across 50 Monte-Carlo realizations over each iteration. For all cases, we assume  $p = 0.3mn$ . (a)  $m = 256, n = 512, k = 10$  and  $\|\epsilon\|_2 = 10^{-3}$ . (b)  $m = 256, n = 512, k = 10$  and  $\|\epsilon\|_2 = 10^{-4}$ . (c)  $m = 256, n = 512, k = 20$  and  $\|\epsilon\|_2 = 0$ . (d)  $m = 512, n = 1024, k = 30$  and  $\|\epsilon\|_2 = 0$ . (e)  $m = 512, n = 1024, k = 40$  and  $\|\epsilon\|_2 = 0$ . (f)  $m = 1024, n = 2048, k = 50$  and  $\|\epsilon\|_2 = 0$

In Table 1, MATRIX ALPS II and MATRIX ALPS II with QR obtain accurate low-rank solutions much faster than the rest of the algorithms in comparison. In high dimensional settings, MATRIX ALPS II with QR scales better as the problem dimensions increase, leading to faster convergence. Moreover, its execution time is at least a few orders

of magnitude smaller compared to SVP, ADMiRA and MATRIX ALPS I implementations.

**Robust matrix completion:** We design matrix completion problems in the following way. The signal of interest  $X^* \in \mathbb{R}^{m \times n}$  is synthesized as a rank- $k$  matrix, factorized as  $X^* := \mathbf{L}\mathbf{R}^T$  with  $\|X^*\|_F = 1$  where  $\mathbf{L} \in \mathbb{R}^{m \times k}$  and

**Table 1** General ARM using Noiselets

Configuration				FR	SVP			ADMIRA			MATRIX ALPS I		
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
256	512	5	0	0.097	38	$2.2 \times 10^{-4}$	0.78	27	$4.4 \times 10^{-5}$	2.26	13.5	$1 \times 10^{-5}$	0.7
256	512	5	$10^{-3}$	0.097	38	$6 \times 10^{-4}$	0.91	700	$2 \times 10^{-3}$	65.94	16	$7 \times 10^{-4}$	0.92
256	512	5	$10^{-4}$	0.097	38	$2.1 \times 10^{-4}$	0.94	700	$4.1 \times 10^{-4}$	69.03	11.5	$7.9 \times 10^{-5}$	0.72
256	512	10	0	0.193	50	$3.4 \times 10^{-4}$	1.44	38	$5 \times 10^{-5}$	4.42	13	$3.9 \times 10^{-5}$	0.92
256	512	10	$10^{-3}$	0.193	50	$9 \times 10^{-4}$	1.39	700	$1.7 \times 10^{-3}$	56.94	29	$1.2 \times 10^{-3}$	1.78
256	512	10	$10^{-4}$	0.193	50	$3.5 \times 10^{-4}$	1.38	700	$9.3 \times 10^{-5}$	64.69	14	$1.4 \times 10^{-4}$	0.93
256	512	20	0	0.38	86	$7 \times 10^{-4}$	3.32	700	$4.1 \times 10^{-5}$	81.93	45	$2 \times 10^{-4}$	4.09
256	512	20	$10^{-3}$	0.38	86	$1.5 \times 10^{-3}$	3.45	700	$4.2 \times 10^{-2}$	77.35	69	$2.3 \times 10^{-3}$	5.05
256	512	20	$10^{-4}$	0.38	86	$7 \times 10^{-4}$	3.26	700	$4 \times 10^{-2}$	79.47	46	$4 \times 10^{-4}$	4.1
512	1024	30	0	0.287	66	$4.9 \times 10^{-4}$	8.79	295	$5.4 \times 10^{-5}$	143.53	24	$1 \times 10^{-4}$	8.01
512	1024	40	0	0.38	86	$7 \times 10^{-4}$	10.09	700	$4.3 \times 10^{-2}$	251.27	45	$2 \times 10^{-4}$	11.08
1024	2048	50	0	0.24	57	$4.3 \times 10^{-4}$	42.88	103	$5.2 \times 10^{-5}$	312.62	18	$5.7 \times 10^{-5}$	35.86

Configuration				FR	MATRIX ALPS II			MATRIX ALPS II with QR		
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time
256	512	5	0	0.097	8	$7.1 \times 10^{-6}$	0.42	10	$9.1 \times 10^{-6}$	<b>0.39</b>
256	512	5	$10^{-3}$	0.097	9	$7 \times 10^{-4}$	<b>0.56</b>	20	$7 \times 10^{-4}$	0.93
256	512	5	$10^{-4}$	0.097	8	$7 \times 10^{-5}$	0.5	10	$7.8 \times 10^{-5}$	<b>0.46</b>
256	512	10	0	0.193	10	$2.3 \times 10^{-5}$	0.68	13	$2.4 \times 10^{-5}$	<b>0.64</b>
256	512	10	$10^{-3}$	0.193	19	$1 \times 10^{-3}$	<b>1.29</b>	27	$1 \times 10^{-3}$	1.35
256	512	10	$10^{-4}$	0.193	10	$1.1 \times 10^{-4}$	0.68	13	$1.1 \times 10^{-4}$	<b>0.62</b>
256	512	20	0	0.38	21	$1 \times 10^{-4}$	1.92	24	$1 \times 10^{-4}$	<b>1.26</b>
256	512	20	$10^{-3}$	0.38	36	$1.5 \times 10^{-3}$	2.67	39	$1.5 \times 10^{-3}$	<b>1.69</b>
256	512	20	$10^{-4}$	0.38	21	$2 \times 10^{-4}$	1.87	24	$2 \times 10^{-4}$	<b>1.22</b>
512	1024	30	0	0.287	14	$4.5 \times 10^{-5}$	4.7	18	$3.3 \times 10^{-5}$	<b>4.15</b>
512	1024	40	0	0.38	21	$1 \times 10^{-4}$	6.01	24	$1 \times 10^{-4}$	<b>4.53</b>
1024	2048	50	0	0.24	12	$2.5 \times 10^{-5}$	22.76	15	$3.3 \times 10^{-5}$	<b>17.94</b>

$\mathbf{R} \in \mathbb{R}^{n \times k}$  as defined above. In sequence, we subsample  $\mathbf{X}^*$  by observing  $p = 0.3mn$  entries, drawn uniformly at random. We denote the set of ordered pairs that represent the coordinates of the observable entries as  $\Omega = \{(i, j) : [\mathbf{X}^*]_{ij} \text{ is known}\} \subseteq \{1, \dots, m\} \times \{1, \dots, n\}$  and let  $\mathcal{A}_\Omega$  denote the linear operator (mask) that samples a matrix according to  $\Omega$ . Then, the set of observations satisfies:

$$\mathbf{y} = \mathcal{A}_\Omega \mathbf{X}^* + \boldsymbol{\epsilon}, \quad (35)$$

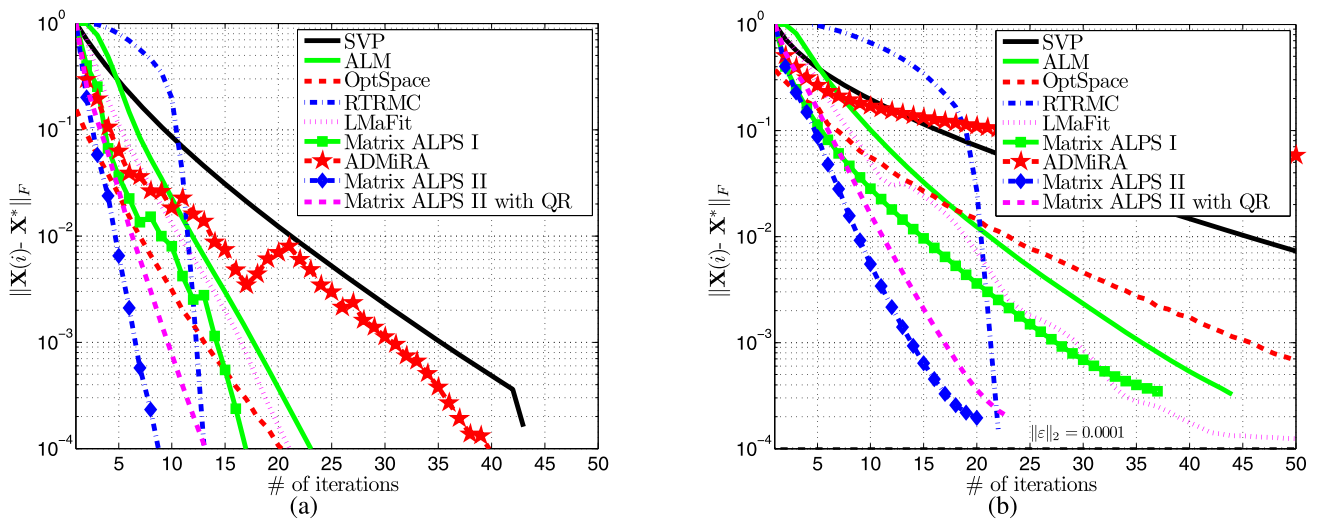
i.e., the known entries of  $\mathbf{X}^*$  are structured as a vector  $\mathbf{y} \in \mathbb{R}^p$ , disturbed by a dense noise vector  $\boldsymbol{\epsilon} \in \mathbb{R}^p$  with fixed-energy, which is populated by iid zero-mean Gaussians.

To demonstrate the reconstruction accuracy and the convergence speeds, we generate various problem configurations (both noisy and noiseless settings), according to (35). The energy of the additive noise takes values  $\|\boldsymbol{\epsilon}\|_2 \in$

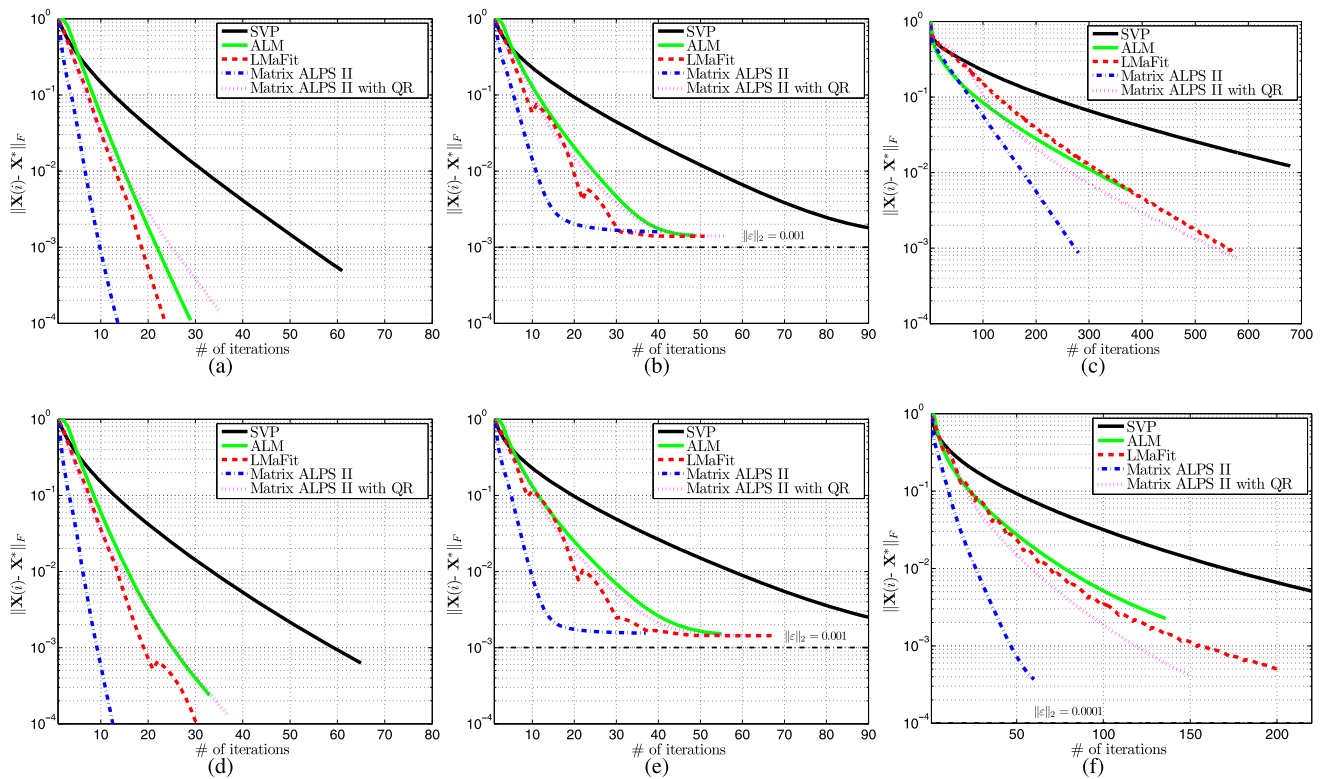
$\{10^{-3}, 10^{-4}\}$ . All the algorithms are tested for the same signal-matrix-noise realizations. A summary of the results can be found in Tables 2, 3 and 4 where we present the median values of reconstruction error, number of iterations and execution time over 50 Monte Carlo iterations. For all cases, we assume  $\text{SR} = 0.3$  and set the maximum number of iterations to 700. Bold font denotes the fastest execution time. Some convergence error curves for specific cases are illustrated in Figs. 9 and 10.

In Table 2, LMaFit [52] implementation has the fastest convergence for small scale problem configuration where  $m = 300$  and  $n = 600$ . We note that part of LMaFit implementation uses C code for acceleration. GROUSE [49] is a competitive low-rank recovery method with attractive execution times for the *extreme low rank* problem settings due to stochastic gradient descent techniques. Nevertheless,





**Fig. 9** Low rank matrix recovery for the matrix completion problem. The error curves are the median values across 50 Monte-Carlo realizations over each iteration. For all cases, we assume  $p = 0.3mn$ . **(a)**  $m = 300, n = 600, k = 5$  and  $\|\epsilon\|_2 = 0$ . **(b)**  $m = 300, n = 600, k = 20$  and  $\|\epsilon\|_2 = 10^{-4}$



**Fig. 10** Low rank matrix recovery for the matrix completion problem. The error curves are the median values across 50 Monte-Carlo realizations over each iteration. For all cases, we assume  $p = 0.3mn$ . **(a)**  $m = 700, n = 1000, k = 30$  and  $\|\epsilon\|_2 = 0$ . **(b)**  $m = 700, n = 1000, k = 50$  and  $\|\epsilon\|_2 = 10^{-3}$ . **(c)**  $m = 700, n = 1000, k = 110$  and  $\|\epsilon\|_2 = 0$ . **(d)**  $m = 500, n = 2000, k = 10$  and  $\|\epsilon\|_2 = 0$ . **(e)**  $m = 500, n = 2000, k = 50$  and  $\|\epsilon\|_2 = 10^{-3}$ . **(f)**  $m = 500, n = 2000, k = 80$  and  $\|\epsilon\|_2 = 10^{-4}$

its execution time performance degrades significantly as we increase the rank of  $X^*$ . Moreover, we observe how randomized low rank projections accelerate the convergence

speed where MATRIX ALPS II with QR converges faster than MATRIX ALPS II. In Tables 3 and 4, we increase the problem dimensions. Here, MATRIX ALPS II with QR has

**Table 2** Matrix Completion problem for  $m = 300$  and  $n = 600$ . “–” depicts no information or not applicable due to time overhead

Configuration				FR	SVP			GROUSE			TFOCS		
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
300	600	5	0	0.083	43	$2.9 \times 10^{-4}$	0.59	–	$1.52 \times 10^{-4}$	0.08	–	$8.69 \times 10^{-5}$	3.36
300	600	5	$10^{-3}$	0.083	42	$6 \times 10^{-4}$	0.65	–	$2 \times 10^{-4}$	0.082	–	$5 \times 10^{-4}$	3.85
300	600	5	$10^{-4}$	0.083	43	$3 \times 10^{-4}$	0.64	–	$2 \times 10^{-4}$	0.079	–	$1 \times 10^{-4}$	3.5
300	600	10	0	0.165	54	$4 \times 10^{-4}$	0.9	–	$4.5 \times 10^{-6}$	0.22	–	$2 \times 10^{-4}$	6.43
300	600	10	$10^{-3}$	0.165	54	$9 \times 10^{-4}$	0.89	–	$2 \times 10^{-4}$	0.16	–	$8 \times 10^{-4}$	7.83
300	600	10	$10^{-4}$	0.165	54	$4 \times 10^{-4}$	0.91	–	$2 \times 10^{-4}$	0.16	–	$1 \times 10^{-4}$	6.75
300	600	20	0	0.326	85	$8 \times 10^{-4}$	2.04	–	$1 \times 10^{-4}$	0.81	–	$2 \times 10^{-4}$	30.04
300	600	40	0	0.637	241	$3.4 \times 10^{-3}$	11.1	–	$3.1 \times 10^{-3}$	13.94	–	–	–
				Inexact ALM			OptSpace			GRASTA			
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
300	600	5	0	0.083	24	$6.7 \times 10^{-5}$	0.47	31	$2.8 \times 10^{-6}$	2.41	–	$2.2 \times 10^{-4}$	2.07
300	600	5	$10^{-3}$	0.083	24	$6 \times 10^{-4}$	0.49	297	$5 \times 10^{-4}$	22.82	–	$1 \times 10^{-4}$	2.07
300	600	5	$10^{-4}$	0.083	24	$1 \times 10^{-4}$	0.49	267	$1 \times 10^{-4}$	21.56	–	$8 \times 10^{-5}$	2.1
300	600	10	0	0.165	26	$1 \times 10^{-4}$	0.6	37	$2.3 \times 10^{-6}$	8.42	–	$8.6 \times 10^{-6}$	4.5
300	600	10	$10^{-3}$	0.165	26	$8 \times 10^{-4}$	0.59	304	$8 \times 10^{-4}$	66.02	–	$5.5 \times 10^{-3}$	3.43
300	600	10	$10^{-4}$	0.165	26	$1 \times 10^{-4}$	0.61	304	$1 \times 10^{-4}$	65.56	–	$5.3 \times 10^{-3}$	3.44
300	600	20	0	0.326	44	$3 \times 10^{-4}$	1.37	–	–	–	–	$5 \times 10^{-4}$	10.51
300	600	40	0	0.637	134	$1.6 \times 10^{-3}$	7.08	–	–	–	–	$5.2 \times 10^{-3}$	251.34
				RTRMC			LMaFit			MATRIX ALPS I			
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
300	600	5	0	0.083	13	$1.2 \times 10^{-4}$	0.59	20	$2.2 \times 10^{-4}$	<b>0.054</b>	22	$1.8 \times 10^{-5}$	0.76
300	600	5	$10^{-3}$	0.083	13	$1 \times 10^{-4}$	0.59	19	$5 \times 10^{-4}$	<b>0.049</b>	37	$7 \times 10^{-4}$	1.34
300	600	5	$10^{-4}$	0.083	13	$2 \times 10^{-4}$	0.59	21	$1 \times 10^{-4}$	<b>0.052</b>	18	$1 \times 10^{-4}$	0.61
300	600	10	0	0.165	16	$1.1 \times 10^{-3}$	1.03	23	$1 \times 10^{-4}$	<b>0.064</b>	16	$1 \times 10^{-4}$	0.65
300	600	10	$10^{-3}$	0.165	17	$1 \times 10^{-4}$	1.09	26	$8 \times 10^{-4}$	<b>0.077</b>	30	$1.1 \times 10^{-3}$	1.16
300	600	10	$10^{-4}$	0.165	17	$2 \times 10^{-4}$	1.09	32	$1 \times 10^{-4}$	<b>0.097</b>	16	$1 \times 10^{-4}$	0.63
300	600	20	0	0.326	22	$4 \times 10^{-4}$	2.99	37	$2 \times 10^{-4}$	<b>0.12</b>	37	$2 \times 10^{-4}$	2.05
300	600	40	0	0.637	35	$3 \times 10^{-5}$	11.83	233	$4.9 \times 10^{-4}$	<b>2.52</b>	500	$6.5 \times 10^{-2}$	45.67
				ADMIRA			MATRIX ALPS II			MATRIX ALPS II with QR			
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
300	600	5	0	0.083	59	$5.2 \times 10^{-5}$	2.86	10	$1.7 \times 10^{-5}$	0.34	14	$3.2 \times 10^{-5}$	0.45
300	600	5	$10^{-3}$	0.083	700	$4 \times 10^{-3}$	30.96	12	$6 \times 10^{-4}$	0.44	24	$6 \times 10^{-4}$	0.81
300	600	5	$10^{-4}$	0.083	700	$4.5 \times 10^{-3}$	31.45	10	$1 \times 10^{-4}$	0.36	14	$1 \times 10^{-4}$	0.47
300	600	10	0	0.165	47	$1 \times 10^{-3}$	2.56	12	$3 \times 10^{-5}$	0.48	16	$3.4 \times 10^{-5}$	0.49
300	600	10	$10^{-3}$	0.165	700	$1.5 \times 10^{-3}$	28.49	19	$9 \times 10^{-4}$	0.74	29	$9 \times 10^{-4}$	0.95
300	600	10	$10^{-4}$	0.165	700	$1 \times 10^{-4}$	31.99	12	$1 \times 10^{-4}$	0.49	16	$1 \times 10^{-4}$	0.54
300	600	20	0	0.326	700	$1.2 \times 10^{-3}$	41.86	20	$1 \times 10^{-4}$	1.16	23	$1 \times 10^{-4}$	0.79
300	600	20	0	0.326	–	–	–	72	$2 \times 10^{-4}$	7.21	68	$2 \times 10^{-4}$	2.6

faster convergence for most of the cases and scales well as the problem size increases. We note that we do not exploit stochastic gradient descent techniques in the recovery process to accelerate convergence which is left for future work.

### 10.5 Real Data

We use real data images to highlight the reconstruction performance of the proposed schemes. To this end, we perform grayscale image denoising from an incomplete set of ob-

**Table 3** Matrix Completion problem for  $m = 700$  and  $n = 1000$ . “–” depicts no information or not applicable due to time overhead

Configuration				FR	SVP			Inexact ALM			GROUSE		
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
700	1000	5	0	0.04	34	$1.9 \times 10^{-4}$	1.77	23	$6.5 \times 10^{-5}$	1.69	–	$3.5 \times 10^{-5}$	<b>0.23</b>
700	1000	5	$10^{-3}$	0.04	34	$4.2 \times 10^{-4}$	1.92	23	$3.7 \times 10^{-4}$	1.87	–	$3.1 \times 10^{-4}$	<b>0.24</b>
700	1000	30	0	0.239	61	$4.6 \times 10^{-4}$	6.39	29	$1.2 \times 10^{-4}$	3.91	–	$3.2 \times 10^{-5}$	3.15
700	1000	30	$10^{-3}$	0.239	61	$1.1 \times 10^{-3}$	6.33	29	$1 \times 10^{-3}$	3.87	–	$8 \times 10^{-4}$	3.14
700	1000	50	0	0.393	95	$8.5 \times 10^{-4}$	14.47	49	$3.2 \times 10^{-4}$	9.02	–	$1.3 \times 10^{-5}$	10.31
700	1000	50	$10^{-3}$	0.393	95	$1.6 \times 10^{-3}$	15.15	49	$1.4 \times 10^{-3}$	9.11	–	$8 \times 10^{-4}$	10.34
700	1000	110	0	0.833	683	$1.2 \times 10^{-2}$	253.1	374	$5.8 \times 10^{-3}$	152.61	–	$1.2 \times 10^{-1}$	110.93
700	1000	110	$10^{-3}$	0.833	682	$1.3 \times 10^{-2}$	256.21	374	$6.8 \times 10^{-3}$	154.34	–	$1.05 \times 10^{-1}$	111.05

				LMaFit			MATRIX ALPS II			MATRIX ALPS II with QR			
$m$	$n$	$k$	$\ \epsilon\ _2$	iter.	err.	time	iter.	err.	time	iter.	err.	time	
700	1000	5	0	0.04	24	$7.2 \times 10^{-6}$	0.67	8	$1.5 \times 10^{-5}$	1.15	15	$8.3 \times 10^{-5}$	1.05
700	1000	5	$10^{-3}$	0.04	17	$3.7 \times 10^{-4}$	0.5	10	$4.5 \times 10^{-4}$	1.38	15	$3.8 \times 10^{-4}$	1.1
700	1000	30	0	0.239	34	$9.2 \times 10^{-6}$	<b>1.95</b>	14	$4.5 \times 10^{-5}$	3.69	35	$1.1 \times 10^{-4}$	2.6
700	1000	30	$10^{-3}$	0.239	30	$1 \times 10^{-3}$	<b>1.71</b>	25	$1.1 \times 10^{-3}$	6.1	35	$1 \times 10^{-3}$	2.61
700	1000	50	0	0.393	53	$2.7 \times 10^{-5}$	4.59	25	$8.6 \times 10^{-5}$	8.87	57	$1.6 \times 10^{-5}$	<b>4.47</b>
700	1000	50	$10^{-3}$	0.393	52	$1.4 \times 10^{-3}$	4.53	40	$1.6 \times 10^{-3}$	14.38	57	$1.4 \times 10^{-3}$	<b>4.49</b>
700	1000	110	0	0.833	584	$9 \times 10^{-4}$	101.95	280	$8 \times 10^{-4}$	214.93	553	$7 \times 10^{-4}$	<b>51.72</b>
700	1000	110	$10^{-3}$	0.833	584	$3.7 \times 10^{-3}$	102.15	336	$4.7 \times 10^{-3}$	261.98	551	$3.7 \times 10^{-3}$	<b>51.62</b>

served pixels—similar experiments can be found in [52]. Based on the matrix completion setting, we observe a limited number of pixels from the original image and perform a low rank approximation based only on the set of measurements. While the true underlying image might not be low-rank, we apply our solvers to obtain low-rank approximations.

Figures 11 and 12 depict the reconstruction results. In the first test case, we use a  $512 \times 512$  grayscale image as shown in the top left corner of Fig. 11. For this case, we observe only the 35 % of the total number of pixels, randomly selected—a realization is depicted in the top right plot in Fig. 11. In sequel, we fix the desired rank to  $k = 40$ . The best rank-40 approximation using SVD is shown in the top middle of Fig. 11 where the full set of pixels is observed. Given a fixed common tolerance and the same stopping criteria, Fig. 11 shows the recovery performance achieved by a range of algorithms under consideration for 10 Monte-Carlo realizations. We repeat the same experiment for the second image in Fig. 12. Here, the size of the image is  $256 \times 256$ , the desired rank is set to  $k = 30$  and we observe the 33 % of the image pixels. In contrast to the image denoising procedure above, we measure the reconstruction error of the computed solutions with respect to the *best rank-30 approximation* of the true image. In both cases, we note that MATRIX ALPS II has a better phase

transition performance as compared to the rest of the algorithms.

### 11 Discussion

In this paper, we present new strategies and review existing ones for hard thresholding methods to recover low-rank matrices from dimensionality reducing, linear projections. Our discussion revolves around four basic building blocks that exploit the problem structure to reduce computational complexity without sacrificing stability.

In theory, constant  $\mu_i$  selection schemes are accompanied with strong RIP constant conditions but empirical evidence reveal signal reconstruction vulnerabilities. While convergence derivations of adaptive schemes are characterized by weaker bounds, the performance gained by this choice in terms of convergence rate, is quite significant. Memory-based methods lead to convergence speed with (almost) no extra cost on the complexity of hard thresholding methods—we provide theoretical evidence for convergence for simple cases but more theoretical justification is needed to generalize this part as future work. Lastly, further estimate refinement over low rank subspaces using gradient update steps or pseudoinversion optimization techniques provides signal reconstruction efficacy, but more computational power is needed per iteration.

**Table 4** Matrix Completion problem for  $m = 500$  and  $n = 2000$ . “–” depicts no information or not applicable due to time overhead

Configuration				FR	SVP			Inexact ALM			GROUSE		
$m$	$n$	$k$	$\ \epsilon\ _2$		iter.	err.	time	iter.	err.	time	iter.	err.	time
500	2000	30	0	0.083	64	$5.3 \times 10^{-4}$	10.18	32	$1.9 \times 10^{-4}$	6.47	–	$1.6 \times 10^{-4}$	<b>2.46</b>
500	2000	30	$10^{-3}$	0.083	64	$1.1 \times 10^{-3}$	6.69	32	$1 \times 10^{-3}$	4.51	–	$6 \times 10^{-4}$	<b>1.94</b>
500	2000	30	$10^{-4}$	0.083	64	$5.4 \times 10^{-4}$	10.14	32	$2.2 \times 10^{-4}$	6.51	–	$1.6 \times 10^{-4}$	<b>2.46</b>
500	2000	50	0	0.408	103	$1.1 \times 10^{-4}$	15.74	54	$5 \times 10^{-4}$	10.8	–	$8 \times 10^{-5}$	7.32
500	2000	50	$10^{-3}$	0.408	103	$1.8 \times 10^{-3}$	24.97	54	$1.55 \times 10^{-3}$	16.14	–	$9 \times 10^{-4}$	8.6
500	2000	50	$10^{-4}$	0.408	102	$1.1 \times 10^{-3}$	24.85	54	$5 \times 10^{-4}$	16.17	–	$7 \times 10^{-5}$	8.59
500	2000	80	0	0.645	239	$3.5 \times 10^{-3}$	92.91	134	$1.7 \times 10^{-3}$	59.33	–	$1 \times 10^{-4}$	79.64
500	2000	80	$10^{-3}$	0.645	239	$4.2 \times 10^{-3}$	94.86	134	$2.8 \times 10^{-3}$	60.68	–	$1 \times 10^{-4}$	79.98
500	2000	80	$10^{-4}$	0.645	239	$3.6 \times 10^{-3}$	93.95	134	$1.8 \times 10^{-3}$	60.76	–	$1 \times 10^{-4}$	79.48
500	2000	100	0	0.8	523	$1.1 \times 10^{-2}$	259.13	307	$6 \times 10^{-3}$	173.14	–	$4.5 \times 10^{-2}$	143.41
500	2000	100	$10^{-3}$	0.8	525	$1.2 \times 10^{-2}$	262.19	308	$7 \times 10^{-3}$	176.04	–	$5.2 \times 10^{-2}$	142.85
500	2000	100	$10^{-4}$	0.8	523	$1.1 \times 10^{-2}$	262.11	307	$6 \times 10^{-3}$	170.47	–	$5.1 \times 10^{-2}$	144.78

				LMaFit			MATRIX ALPS II			MATRIX ALPS II with QR			
$m$	$n$	$k$	$\ \epsilon\ _2$	iter.	err.	time	iter.	err.	time	iter.	err.	time	
500	2000	30	0	0.083	37	$1.3 \times 10^{-5}$	3.05	13	$3.1 \times 10^{-5}$	4.84	37	$1.2 \times 10^{-5}$	4.04
500	2000	30	$10^{-3}$	0.083	37	$1 \times 10^{-3}$	2.52	22	$1.1 \times 10^{-3}$	5.35	37	$1 \times 10^{-3}$	3.32
500	2000	30	$10^{-4}$	0.083	35	$1 \times 10^{-4}$	2.86	13	$1.3 \times 10^{-4}$	4.85	37	$1.6 \times 10^{-4}$	4.05
500	2000	50	0	0.408	60	$6 \times 10^{-5}$	6.06	22	$1 \times 10^{-4}$	7.6	60	$2 \times 10^{-4}$	<b>5.67</b>
500	2000	50	$10^{-3}$	0.408	60	$1.4 \times 10^{-3}$	7.26	36	$1.6 \times 10^{-3}$	19.64	59	$1.6 \times 10^{-3}$	<b>6.91</b>
500	2000	50	$10^{-4}$	0.408	60	$2 \times 10^{-4}$	7.29	22	$2 \times 10^{-4}$	11.87	59	$2 \times 10^{-4}$	<b>6.75</b>
500	2000	80	0	0.645	183	$3 \times 10^{-4}$	33.65	61	$2 \times 10^{-4}$	49.53	151	$3 \times 10^{-4}$	<b>18.66</b>
500	2000	80	$10^{-3}$	0.645	183	$2.3 \times 10^{-3}$	33.48	92	$2.4 \times 10^{-3}$	75.51	151	$2.3 \times 10^{-3}$	<b>18.87</b>
500	2000	80	$10^{-4}$	0.645	183	$3 \times 10^{-4}$	33.47	61	$4 \times 10^{-4}$	49.52	151	$3 \times 10^{-4}$	<b>18.92</b>
500	2000	100	0	0.8	519	$1.5 \times 10^{-3}$	115.11	148	$4 \times 10^{-4}$	153.74	429	$7 \times 10^{-4}$	<b>55.1</b>
500	2000	100	$10^{-3}$	0.8	529	$3.6 \times 10^{-3}$	117.7	228	$3.7 \times 10^{-3}$	239.92	427	$3.4 \times 10^{-3}$	<b>55.7</b>
500	2000	100	$10^{-4}$	0.8	520	$1.6 \times 10^{-3}$	116.66	148	$6 \times 10^{-4}$	154.46	428	$8 \times 10^{-4}$	<b>55.07</b>

We connect  $\epsilon$ -approximation low-rank revealing schemes with first-order gradient descent algorithms to solve general affine rank minimization problems; to the best of our knowledge, this is the first attempt to theoretically characterize the performance of iterative greedy algorithms with  $\epsilon$ -approximation schemes. In all cases, experimental results illustrate the effectiveness of the proposed schemes on different problem configurations.

**Acknowledgements** This work was supported in part by the European Commission under Grant MIRG-268398, ERC Future Proof, SNF 200021-132548 and DARPA KeCoM program #11-DARPA-1055. VC also would like to acknowledge Rice University for his Faculty Fellowship.

**Appendix**

*Remark 1* Let  $X \in \mathbb{R}^{m \times n}$  with SVD:  $X = U \Sigma V^T$ , and  $Y \in \mathbb{R}^{m \times n}$  with SVD:  $Y = \tilde{U} \tilde{\Sigma} \tilde{V}^T$ . Assume two sets: (i)  $\mathcal{S}_1 =$

$\{u_i u_i^T : i \in \mathcal{I}_1\}$  where  $u_i$  is the  $i$ -th singular vector of  $X$  and  $\mathcal{I}_1 \subseteq \{1, \dots, \text{rank}(X)\}$  and, (ii)  $\mathcal{S}_2 = \{u_i u_i^T, \tilde{u}_j \tilde{u}_j^T : i \in \mathcal{I}_2, j \in \mathcal{I}_3\}$  where  $\tilde{u}_j$  is the  $i$ -th singular vector of  $Y$ ,  $\mathcal{I}_1 \subseteq \mathcal{I}_2 \subseteq \{1, \dots, \text{rank}(X)\}$  and,  $\mathcal{I}_3 \subseteq \{1, \dots, \text{rank}(Y)\}$ . We observe that the subspaces defined by  $u_i u_i^T$  and  $\tilde{u}_j \tilde{u}_j^T$  are not necessarily orthogonal.

To this end, let  $\hat{\mathcal{S}}_2 = \text{ortho}(\mathcal{S}_2)$ ; this operation can be easily computed via SVD. Then, the following commutativity property holds true for any matrix  $W \in \mathbb{R}^{m \times n}$ :

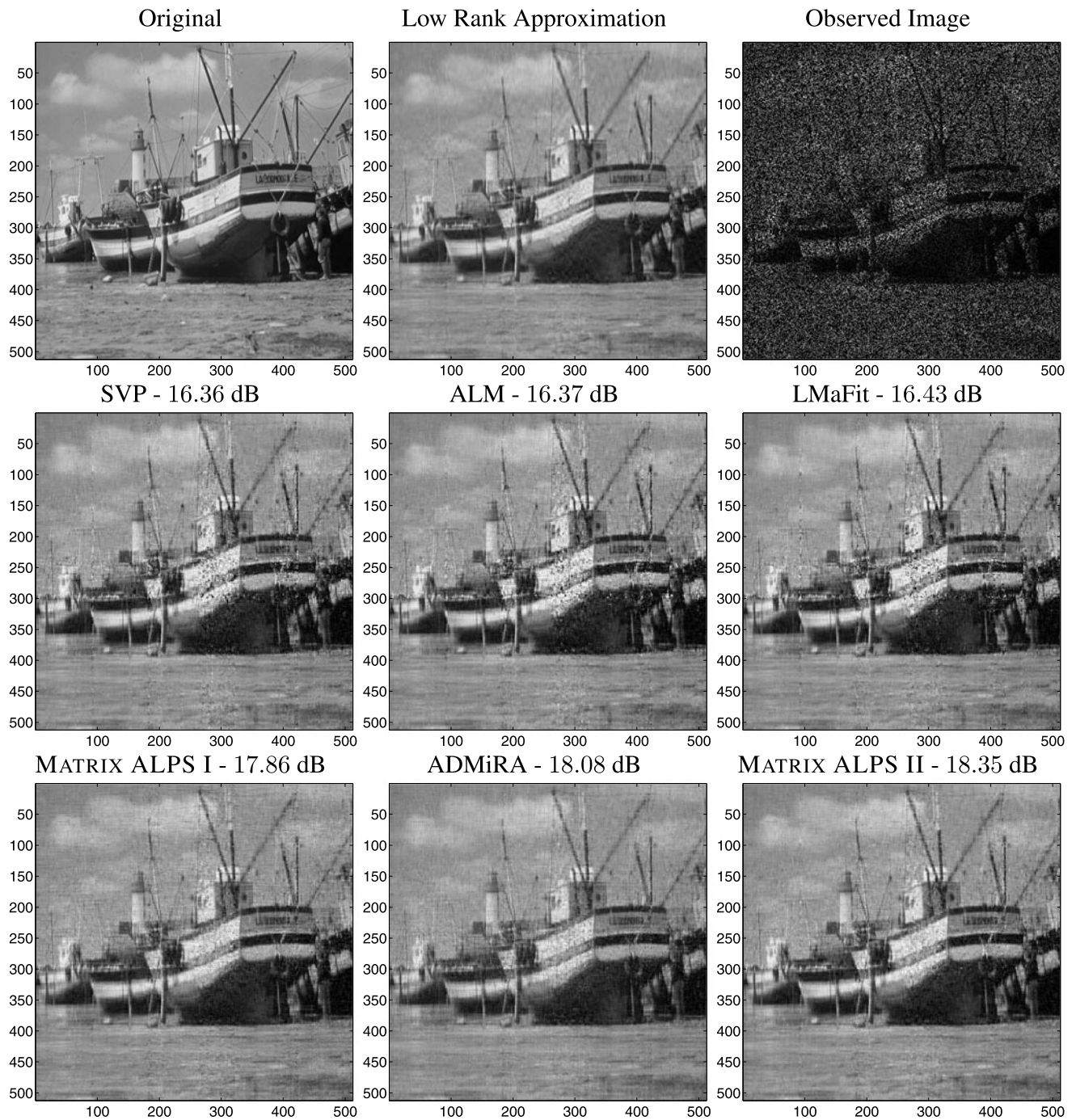
$$\mathcal{P}_{\mathcal{S}_1} \mathcal{P}_{\hat{\mathcal{S}}_2} W = \mathcal{P}_{\hat{\mathcal{S}}_2} \mathcal{P}_{\mathcal{S}_1} W. \tag{36}$$

**A.1 Proof of Lemma 6**

Given  $\mathcal{X}^* \leftarrow \mathcal{P}_k(X^*)$  using SVD factorization, we define the following quantities:  $\mathcal{S}_i \leftarrow \mathcal{X}_i \cup \mathcal{D}_i$ ,  $\mathcal{S}_i^* \leftarrow \text{ortho}(\mathcal{X}_i \cup \mathcal{X}^*)$ . Then, given the structure of the sets  $\mathcal{S}_i$  and  $\mathcal{S}_i^*$

$$\mathcal{P}_{\mathcal{S}_i} \mathcal{P}_{(\mathcal{S}_i^*)^\perp} = \mathcal{P}_{\mathcal{D}_i} \mathcal{P}_{(\mathcal{X}^* \cup \mathcal{X}_i)^\perp}, \tag{37}$$





**Fig. 11** Reconstruction performance in image denoising settings. The image size is  $512 \times 512$  and the desired rank is preset to  $k = 40$ . We observe 35 % of the pixels of the true image. We depict the median re-

construction error with respect to the true image in dB over 10 Monte Carlo realizations

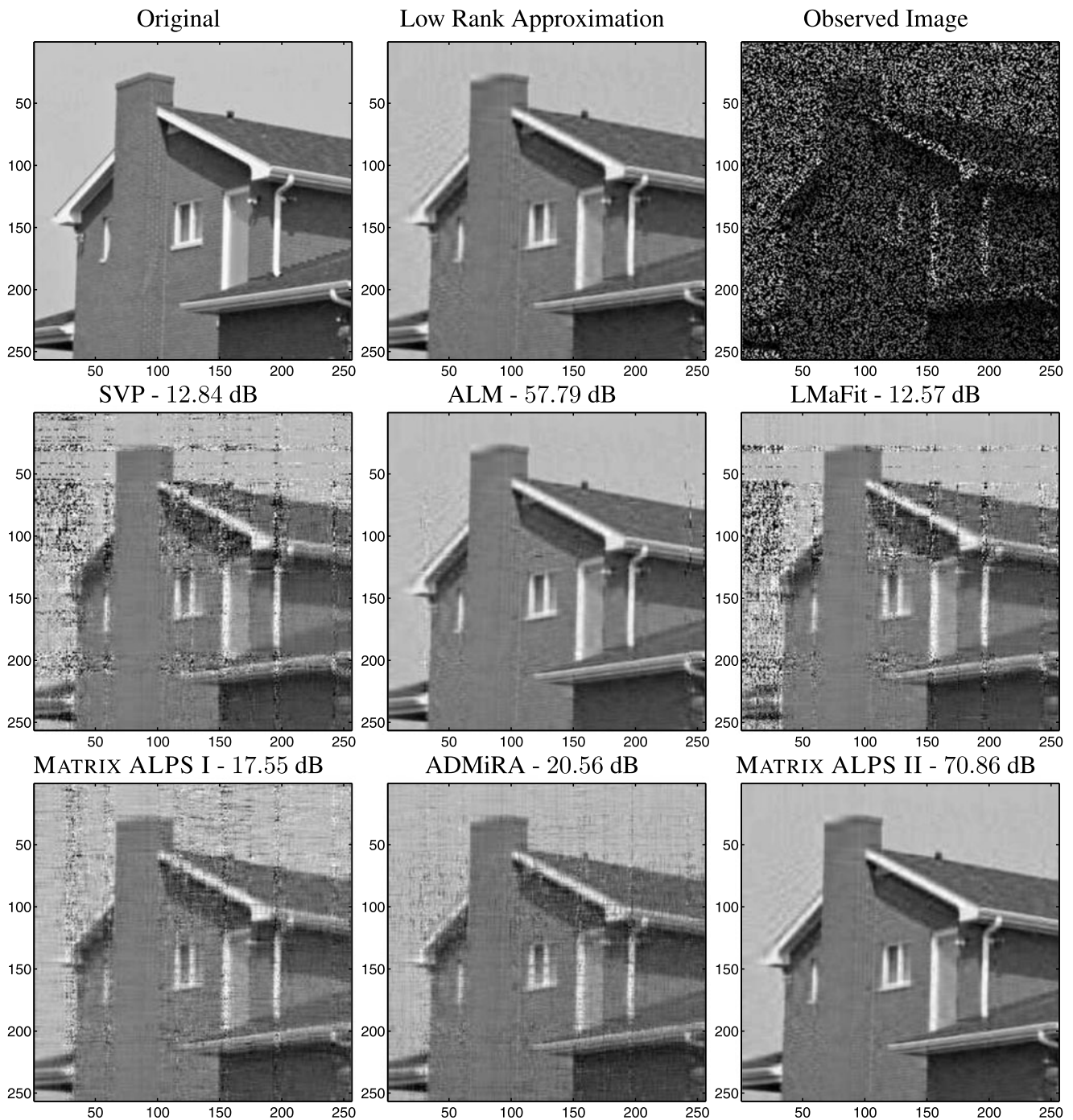
and

$$\mathcal{P}_{S_i^*} \mathcal{P}_{S_i^\perp} = \mathcal{P}_{\mathcal{X}^*} \mathcal{P}_{(\mathcal{D}_i \cup \mathcal{X}_i)^\perp}. \tag{38}$$

Since the subspace defined in  $\mathcal{D}_i$  is the best rank- $k$  subspace, orthogonal to the subspace spanned by  $\mathcal{X}_i$ , the following

holds true:

$$\begin{aligned} \|\mathcal{P}_{\mathcal{D}_i} \mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i))\|_F^2 &\geq \|\mathcal{P}_{\mathcal{X}^*} \mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i))\|_F^2 \\ \Rightarrow \|\mathcal{P}_{S_i} \nabla f(\mathbf{X}(i))\|_F^2 &\geq \|\mathcal{P}_{S_i^*} \nabla f(\mathbf{X}(i))\|_F^2. \end{aligned}$$



**Fig. 12** Reconstruction performance in image denoising settings. The image size is  $256 \times 256$  and the desired rank is preset to  $k = 30$ . We observe 33 % of the pixels of the best rank-30 approximation of the

image. We depict the median reconstruction with respect to the best rank-30 approximation in dB over 10 Monte Carlo realizations

Removing the common subspaces in  $\mathcal{S}_i$  and  $\mathcal{S}_i^*$  by the commutativity property of the projection operation and using the shortcut  $\mathcal{P}_{\mathcal{A} \setminus \mathcal{B}} \equiv \mathcal{P}_{\mathcal{A}} \mathcal{P}_{\mathcal{B}^\perp}$  for sets  $\mathcal{A}, \mathcal{B}$ , we get:

$$\| \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \nabla f(\mathbf{X}(i)) \|_F^2 \geq \| \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \nabla f(\mathbf{X}(i)) \|_F^2$$

$$\begin{aligned} &\Rightarrow \| \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathbf{A}^* \mathbf{A} (\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathbf{A}^* \boldsymbol{\varepsilon} \|_F \\ &\geq \| \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathbf{A}^* \mathbf{A} (\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathbf{A}^* \boldsymbol{\varepsilon} \|_F. \end{aligned} \quad (39)$$

Next, we assume that  $\mathcal{P}_{(\mathcal{A} \setminus \mathcal{B})^\perp}$  denotes the orthogonal projection onto the subspace spanned by  $\mathcal{P}_{\mathcal{A}} \mathcal{P}_{\mathcal{B}^\perp}$ . Then, on the

left hand side of (39), we have:

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \stackrel{(i)}{\leq} \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i))\|_F + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \stackrel{(ii)}{=} \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*}(\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \stackrel{(iii)}{=} \|(\mathbf{I} - \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A} \mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp})(\mathbf{X}^* - \mathbf{X}(i)) \\
 & \quad + \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A} \mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \leq \|(\mathbf{I} - \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*})(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A} \mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \stackrel{(iv)}{\leq} \delta_{3k} \|\mathbf{X}^* - \mathbf{X}(i)\|_F + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A} \mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \stackrel{(v)}{\leq} \delta_{3k} \|\mathbf{X}^* - \mathbf{X}(i)\|_F + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \quad + \delta_{3k} \|\mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & \stackrel{(vi)}{\leq} 2\delta_{3k} \|\mathbf{X}^* - \mathbf{X}(i)\|_F + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F, \tag{40}
 \end{aligned}$$

where (i) due to triangle inequality over Frobenius metric norm, (ii) since  $\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*}(\mathbf{X}(i) - \mathbf{X}^*) = \mathbf{0}$ , (iii) by using the fact that  $\mathbf{X}(i) - \mathbf{X}^* := \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*}(\mathbf{X}(i) - \mathbf{X}^*) + \mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}(i) - \mathbf{X}^*)$ , (iv) due to Lemma 4, (v) due to Lemma 5 and (vi) since  $\|\mathcal{P}_{(\mathcal{S}_i \setminus \mathcal{S}_i^*)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \leq \|\mathbf{X}(i) - \mathbf{X}^*\|_F$ .

For the right hand side of (39), we calculate:

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i}(\mathbf{X}^* - \mathbf{X}(i))\|_F
 \end{aligned}$$

---


$$\begin{aligned}
 & \|\mathbf{W}(i) - \mathbf{X}^*\|_F^2 \leq 2\langle \mathcal{P}_{\mathcal{E}}(\mathbf{W}(i) - \mathbf{X}^*), \mathcal{P}_{\mathcal{E}}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \Rightarrow \\
 & = 2\langle \underbrace{\mathcal{P}_{\mathcal{E}}(\mathbf{W}(i) - \mathbf{X}^*)}_{\doteq A}, \underbrace{\mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^* - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A}(\mathbf{X}(i) - \mathbf{X}^*))}_{\doteq B} \rangle + 2\mu_i \langle \mathcal{P}_{\mathcal{E}}(\mathbf{W}(i) - \mathbf{X}^*), \mathcal{P}_{\mathcal{E}} \mathcal{P}_{\mathcal{S}_i}(\mathcal{A}^* \boldsymbol{\varepsilon}) \rangle. \tag{43}
 \end{aligned}$$


---

In B, we observe:

$$\begin{aligned}
 B & := 2\mu_i \langle \mathcal{P}_{\mathcal{E}}(\mathbf{W}(i) - \mathbf{X}^*), \mathcal{P}_{\mathcal{E}} \mathcal{P}_{\mathcal{S}_i}(\mathcal{A}^* \boldsymbol{\varepsilon}) \rangle \\
 & \stackrel{(i)}{=} 2\mu_i \langle \mathbf{W}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathcal{A}^* \boldsymbol{\varepsilon}) \rangle \\
 & \stackrel{(ii)}{\leq} 2\mu_i \|\mathbf{W}(i) - \mathbf{X}^*\|_F \|\mathcal{P}_{\mathcal{S}_i}(\mathcal{A}^* \boldsymbol{\varepsilon})\|_F
 \end{aligned}$$

$$\begin{aligned}
 & - \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{(\mathcal{S}_i^* \setminus \mathcal{S}_i)^\perp}(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & - \|(\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} - \mathbf{I})(\mathbf{X}^* - \mathbf{X}(i))\|_F \\
 & - \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \\
 & \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i}(\mathbf{X}^* - \mathbf{X}(i))\|_F - 2\delta_{2k} \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 & - \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \boldsymbol{\varepsilon}\|_F \tag{41}
 \end{aligned}$$

by using Lemmas 4 and 5. Combining (40) and (41) in (39), we get:

$$\begin{aligned}
 \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{S}_i} \mathbf{X}^*\|_F & \leq (2\delta_{2k} + 2\delta_{3k}) \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 & \quad + \sqrt{2(1 + \delta_{2k})} \|\boldsymbol{\varepsilon}\|_2.
 \end{aligned}$$

### A.2 Proof of Theorem 1

Let  $\mathcal{X}^* \leftarrow \mathcal{P}_k(\mathbf{X}^*)$  be a set of orthonormal, rank-1 matrices that span the range of  $\mathbf{X}^*$ . In Algorithm 1,  $\mathbf{W}(i) \leftarrow \mathcal{P}_k(\mathbf{V}(i))$ . Thus:

$$\begin{aligned}
 \|\mathbf{W}(i) - \mathbf{V}(i)\|_F^2 & \leq \|\mathbf{X}^* - \mathbf{V}(i)\|_F^2 \\
 \Rightarrow \|\mathbf{W}(i) - \mathbf{X}^* + \mathbf{X}^* - \mathbf{V}(i)\|_F^2 & \leq \|\mathbf{X}^* - \mathbf{V}(i)\|_F^2 \\
 \Rightarrow \|\mathbf{W}(i) - \mathbf{X}^*\|_F^2 & \leq 2\langle \mathbf{W}(i) - \mathbf{X}^*, \mathbf{V}(i) - \mathbf{X}^* \rangle. \tag{42}
 \end{aligned}$$

From Algorithm 1, (i)  $\mathbf{V}(i) \in \text{span}(\mathcal{S}_i)$ , (ii)  $\mathbf{X}(i) \in \text{span}(\mathcal{S}_i)$  and (iii)  $\mathbf{W}(i) \in \text{span}(\mathcal{S}_i)$ . We define  $\mathcal{E} \leftarrow \text{ortho}(\mathcal{S}_i \cup \mathcal{X}^*)$  where  $\text{rank}(\text{span}(\mathcal{E})) \leq 3k$  and let  $\mathcal{P}_{\mathcal{E}}$  be the orthogonal projection onto the subspace defined by  $\mathcal{E}$ .

Since  $\mathbf{W}(i) - \mathbf{X}^* \in \text{span}(\mathcal{E})$  and  $\mathbf{V}(i) - \mathbf{X}^* \in \text{span}(\mathcal{E})$ , the following hold true:

$$\begin{aligned}
 \mathbf{W}(i) - \mathbf{X}^* & = \mathcal{P}_{\mathcal{E}}(\mathbf{W}(i) - \mathbf{X}^*) \quad \text{and} \\
 \mathbf{V}(i) - \mathbf{X}^* & = \mathcal{P}_{\mathcal{E}}(\mathbf{V}(i) - \mathbf{X}^*).
 \end{aligned}$$

Then, (42) can be written as:

$$\stackrel{(iii)}{\leq} 2\mu_i \sqrt{1 + \delta_{2k}} \|\mathbf{W}(i) - \mathbf{X}^*\|_F \|\boldsymbol{\varepsilon}\|_2, \tag{44}$$

where (i) holds since  $\mathcal{P}_{\mathcal{S}_i} \mathcal{P}_{\mathcal{E}} = \mathcal{P}_{\mathcal{E}} \mathcal{P}_{\mathcal{S}_i} = \mathcal{P}_{\mathcal{S}_i}$  for  $\text{span}(\mathcal{S}_i) \in \text{span}(\mathcal{E})$ , (ii) is due to Cauchy-Schwarz inequality and, (iii) is easily derived using Lemma 2.



In A, we perform the following motions:

$$\begin{aligned}
 A &:= 2\langle \mathbf{W}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \\
 &\quad - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \rangle \\
 &\stackrel{(i)}{=} 2\langle \mathbf{W}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \\
 &\quad - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} [\mathcal{P}_{\mathcal{S}_i} + \mathcal{P}_{\mathcal{S}_i^\perp}] \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \rangle \\
 &= 2\langle \mathbf{W}(i) - \mathbf{X}^*, (\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \rangle \\
 &\quad - 2\mu_i \langle \mathbf{W}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) \rangle \\
 &\stackrel{(ii)}{\leq} 2\|\mathbf{W}(i) - \mathbf{X}^*\|_F \\
 &\quad \times \|(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 &\quad + 2\mu_i \|\mathbf{W}(i) - \mathbf{X}^*\|_F \\
 &\quad \times \|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F, \tag{45}
 \end{aligned}$$

where (i) is due to  $\mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) := \mathcal{P}_{\mathcal{S}_i} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) + \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)$  and (ii) follows from Cauchy-Schwarz inequality. Since  $\frac{1}{1+\delta_{2k}} \leq \mu_i \leq \frac{1}{1-\delta_{2k}}$ , Lemma 4 implies:

$$\begin{aligned}
 \lambda(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) &\in \left[ 1 - \frac{1 - \delta_{2k}}{1 + \delta_{2k}}, \frac{1 + \delta_{2k}}{1 - \delta_{2k}} - 1 \right] \\
 &\leq \frac{2\delta_{2k}}{1 - \delta_{2k}},
 \end{aligned}$$

and thus:

$$\begin{aligned}
 &\|(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 &\leq \frac{2\delta_{2k}}{1 - \delta_{2k}} \|\mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F.
 \end{aligned}$$

Furthermore, according to Lemma 5:

$$\begin{aligned}
 &\|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 &\leq \delta_{3k} \|\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F
 \end{aligned}$$

since  $\text{rank}(\mathcal{P}_{\mathcal{K}} \mathbf{X}) \leq 3k, \forall \mathbf{X} \in \mathbb{R}^{m \times n}$  for  $\mathcal{K} \leftarrow \text{ortho}(\mathcal{E} \cup \mathcal{S}_i)$ . Since  $\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*) = \mathcal{P}_{\mathcal{X}^* \setminus (\mathcal{D}_i \cup \mathcal{X}_i)} \mathbf{X}^*$  where

$$\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i))),$$

then:

$$\begin{aligned}
 &\|\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 &= \|\mathcal{P}_{\mathcal{X}^* \setminus (\mathcal{D}_i \cup \mathcal{X}_i)} \mathbf{X}^*\|_F \\
 &\leq (2\delta_{2k} + 2\delta_{3k}) \|\mathbf{X}(i) - \mathbf{X}^*\|_F + \sqrt{2(1 + \delta_{2k})} \|\boldsymbol{\epsilon}\|_2,
 \end{aligned}$$

using Lemma 6. Combining the above in (45), we compute:

$$\begin{aligned}
 A &\leq \left( \frac{4\delta_{2k}}{1 - \delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1 - \delta_{2k}} \right) \|\mathbf{W}(i) - \mathbf{X}^*\|_F \\
 &\quad \times \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 &\quad + \frac{2\delta_{3k}}{1 - \delta_{2k}} \|\mathbf{W}(i) - \mathbf{X}^*\|_F \sqrt{2(1 + \delta_{2k})} \|\boldsymbol{\epsilon}\|_2. \tag{46}
 \end{aligned}$$

Combining (44) and (46) in (43), we get:

$$\begin{aligned}
 &\|\mathbf{W}(i) - \mathbf{X}^*\|_F \\
 &\leq \left( \frac{4\delta_{2k}}{1 - \delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1 - \delta_{2k}} \right) \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 &\quad + \left( \frac{2\sqrt{1 + \delta_{2k}}}{1 - \delta_{2k}} + \frac{2\delta_{3k}}{1 - \delta_{2k}} \sqrt{2(1 + \delta_{2k})} \right) \|\boldsymbol{\epsilon}\|_2. \tag{47}
 \end{aligned}$$

Focusing on steps 5 and 6 of Algorithm 1, we perform similar motions to obtain:

$$\begin{aligned}
 \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F &\leq \left( \frac{1 + 2\delta_{2k}}{1 - \delta_{2k}} \right) \|\mathbf{W}(i) - \mathbf{X}^*\|_F \\
 &\quad + \frac{\sqrt{1 + \delta_k}}{1 - \delta_k} \|\boldsymbol{\epsilon}\|_2. \tag{48}
 \end{aligned}$$

Combining the recursions in (47) and (48), we finally compute:

$$\|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F \leq \rho \|\mathbf{X}(i) - \mathbf{X}^*\|_F + \gamma \|\boldsymbol{\epsilon}\|_2,$$

for  $\rho := \left( \frac{1 + 2\delta_{2k}}{1 - \delta_{2k}} \right) \left( \frac{4\delta_{2k}}{1 - \delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1 - \delta_{2k}} \right)$  and

$$\begin{aligned}
 \gamma &:= \left( \left( \frac{1 + 2\delta_{2k}}{1 - \delta_{2k}} \right) \left( \frac{2\sqrt{1 + \delta_{2k}}}{1 - \delta_{2k}} + \frac{2\delta_{3k}}{1 - \delta_{2k}} \sqrt{2(1 + \delta_{2k})} \right) \right. \\
 &\quad \left. + \frac{\sqrt{1 + \delta_k}}{1 - \delta_k} \right).
 \end{aligned}$$

For the convergence parameter  $\rho$ , further compute:

$$\begin{aligned}
 &\left( \frac{1 + 2\delta_{2k}}{1 - \delta_{2k}} \right) \left( \frac{4\delta_{2k}}{1 - \delta_{2k}} + (2\delta_{2k} + 2\delta_{3k}) \frac{2\delta_{3k}}{1 - \delta_{2k}} \right) \\
 &\leq \frac{1 + 2\delta_{3k}}{(1 - \delta_{3k})^2} (4\delta_{3k} + 8\delta_{3k}^2) =: \hat{\rho} \tag{49}
 \end{aligned}$$

for  $\delta_k \leq \delta_{2k} \leq \delta_{3k}$ . Calculating the roots of this expression, we easily observe that  $\rho < \hat{\rho} < 1$  for  $\delta_{3k} < 0.1235$ .

### A.3 Proof of Theorem 2

Before we present the proof of Theorem 2, we list a series of lemmas that correspond to the motions Algorithm 2 performs.

**Lemma 9** [Error norm reduction via least-squares optimization] *Let  $\mathcal{S}_i$  be a set of orthonormal, rank-1 matrices that span a rank- $2k$  subspace in  $\mathbb{R}^{m \times n}$ . Then, the least squares solution  $\mathbf{V}(i)$  given by:*

$$\mathbf{V}(i) \leftarrow \underset{\mathbf{V}: \mathbf{V} \in \text{span}(\mathcal{S}_i)}{\text{arg min}} \|\mathbf{y} - \mathcal{A}\mathbf{V}\|_2^2, \tag{50}$$

satisfies:

$$\begin{aligned} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \leq & \frac{1}{\sqrt{1 - \delta_{3k}^2(\mathcal{A})}} \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ & + \frac{\sqrt{1 + \delta_{2k}}}{1 - \delta_{3k}} \|\boldsymbol{\varepsilon}\|_2. \end{aligned} \tag{51}$$

*Proof* We observe that  $\|\mathbf{V}(i) - \mathbf{X}^*\|_F^2$  is decomposed as follows:

$$\begin{aligned} \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 &= \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 + \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2. \end{aligned} \tag{52}$$

In (50),  $\mathbf{V}(i)$  is the minimizer over the low-rank subspace spanned by  $\mathcal{S}_i$  with  $\text{rank}(\text{span}(\mathcal{S}_i)) \leq 2k$ . Using the optimality condition (Lemma 1) over the convex set  $\Theta = \{\mathbf{X} : \text{span}(\mathbf{X}) \in \mathcal{S}_i\}$ , we have:

$$\begin{aligned} \langle \nabla f(\mathbf{V}(i)), \mathcal{P}_{\mathcal{S}_i}(\mathbf{X}^* - \mathbf{V}(i)) \rangle &\geq 0 \Rightarrow \\ \langle \mathcal{A}\mathbf{V}(i) - \mathbf{y}, \mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle &\leq 0 \end{aligned} \tag{53}$$

for  $\mathcal{P}_{\mathcal{S}_i}\mathbf{X}^* \in \text{span}(\mathcal{S}_i)$ . Given condition (53), the first term on the right hand side of (52) becomes:

$$\begin{aligned} \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 &= \langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &\stackrel{(53)}{\leq} \langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &\quad - \langle \mathcal{A}\mathbf{V}(i) - \mathbf{y}, \mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &\leq \|\langle \mathbf{V}(i) - \mathbf{X}^*, (\mathbf{I} - \mathcal{A}^*\mathcal{A})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle\| \\ &\quad + \langle \boldsymbol{\varepsilon}, \mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle. \end{aligned} \tag{54}$$

Focusing on the term  $|\langle \mathbf{V}(i) - \mathbf{X}^*, (\mathbf{I} - \mathcal{A}^*\mathcal{A})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle|$ , we derive the following:

$$\begin{aligned} &|\langle \mathbf{V}(i) - \mathbf{X}^*, (\mathbf{I} - \mathcal{A}^*\mathcal{A})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle| \\ &= |\langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle| \\ &\quad - \langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{A}^*\mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &\stackrel{(i)}{=} |\langle \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}(\mathbf{V}(i) - \mathbf{X}^*), \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle| \\ &\quad - \langle \mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}(\mathbf{V}(i) - \mathbf{X}^*), \mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \end{aligned}$$

$$\begin{aligned} &\stackrel{(ii)}{=} |\langle \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}(\mathbf{V}(i) - \mathbf{X}^*), \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle| \\ &\quad - \langle \mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}(\mathbf{V}(i) - \mathbf{X}^*), \\ &\quad \quad \mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle| \\ &= |\langle \mathbf{V}(i) - \mathbf{X}^*, \\ &\quad (\mathbf{I} - \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{A}^*\mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle|, \end{aligned}$$

where (i) follows from the facts that  $\mathbf{V}(i) - \mathbf{X}^* \in \text{span}(\text{ortho}(\mathcal{S}_i \cup \mathcal{X}^*))$  and thus  $\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}(\mathbf{V}(i) - \mathbf{X}^*) = \mathbf{V}(i) - \mathbf{X}^*$  and (ii) is due to  $\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{P}_{\mathcal{S}_i} = \mathcal{P}_{\mathcal{S}_i}$  since  $\text{span}(\mathcal{S}_i) \subseteq \text{span}(\text{ortho}(\mathcal{S}_i \cup \mathcal{X}^*))$ . Then, (54) becomes:

$$\begin{aligned} &\|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 \\ &\leq \|\langle \mathbf{V}(i) - \mathbf{X}^*, \\ &\quad (\mathbf{I} - \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{A}^*\mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle\| \\ &\quad + \langle \boldsymbol{\varepsilon}, \mathcal{A}\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &\stackrel{(i)}{\leq} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\quad \times \|\langle (\mathbf{I} - \mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*}\mathcal{A}^*\mathcal{A}\mathcal{P}_{\mathcal{S}_i \cup \mathcal{X}^*})\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*) \rangle\|_F \\ &\quad + \|\mathcal{P}_{\mathcal{S}_i}\mathcal{A}^*\boldsymbol{\varepsilon}\|_F \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ &\stackrel{(ii)}{\leq} \delta_{3k} \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\quad + \sqrt{1 + \delta_{2k}} \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \|\boldsymbol{\varepsilon}\|_2, \end{aligned} \tag{55}$$

where (i) comes from Cauchy-Swartz inequality and (ii) is due to Lemmas 2 and 4. Simplifying the above quadratic expression, we obtain:

$$\|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \leq \delta_{3k} \|\mathbf{V}(i) - \mathbf{X}^*\|_F + \sqrt{1 + \delta_{2k}} \|\boldsymbol{\varepsilon}\|_2. \tag{56}$$

As a consequence, (52) can be upper bounded by:

$$\begin{aligned} \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 &\leq (\delta_{3k} \|\mathbf{V}(i) - \mathbf{X}^*\|_F + \sqrt{1 + \delta_{2k}} \|\boldsymbol{\varepsilon}\|_2)^2 \\ &\quad + \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2. \end{aligned} \tag{57}$$

We form the quadratic polynomial for this inequality assuming as unknown variable the quantity  $\|\mathbf{V}(i) - \mathbf{X}^*\|_F$ . Bounding by the largest root of the resulting polynomial, we get:

$$\begin{aligned} \|\mathbf{V}(i) - \mathbf{X}^*\|_F &\leq \frac{1}{\sqrt{1 - \delta_{3k}^2(\mathcal{A})}} \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ &\quad + \frac{\sqrt{1 + \delta_{2k}}}{1 - \delta_{3k}} \|\boldsymbol{\varepsilon}\|_2. \end{aligned} \tag{58}$$

□

The following lemma characterizes how subspace pruning affects the recovered energy:



**Lemma 10** [Best rank- $k$  subspace selection] *Let  $\mathbf{V}(i) \in \mathbb{R}^{m \times n}$  be a rank- $2k$  proxy matrix in the subspace spanned by  $\mathcal{S}_i$  and let  $\mathbf{X}(i+1) \leftarrow \mathcal{P}_k(\mathbf{V}(i))$  denote the best rank- $k$  approximation to  $\mathbf{V}(i)$ , according to (5). Then:*

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{V}(i)\|_F &\leq \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ &\leq \|\mathbf{V}(i) - \mathbf{X}^*\|_F. \end{aligned} \tag{59}$$

*Proof* Since  $\mathbf{X}(i+1)$  denotes the best rank- $k$  approximation to  $\mathbf{V}(i)$ , the following inequality holds for any rank- $k$  matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  in the subspace spanned by  $\mathcal{S}_i$ , i.e.  $\forall \mathbf{X} \in \text{span}(\mathcal{S}_i)$ :

$$\|\mathbf{X}(i+1) - \mathbf{V}(i)\|_F \leq \|\mathbf{X} - \mathbf{V}(i)\|_F. \tag{60}$$

Since  $\mathcal{P}_{\mathcal{S}_i} \mathbf{V}(i) = \mathbf{V}(i)$ , the left inequality in (59) is satisfied for  $\mathbf{X} := \mathcal{P}_{\mathcal{S}_i} \mathbf{X}^*$  in (60).  $\square$

**Lemma 11** *Let  $\mathbf{V}(i)$  be the least squares solution in Step 2 of the ADMiRA algorithm and let  $\mathbf{X}(i+1)$  be a proxy, rank- $k$  matrix to  $\mathbf{V}(i)$  according to:  $\mathbf{X}(i+1) \leftarrow \mathcal{P}_k(\mathbf{V}(i))$ . Then,  $\|\mathbf{X}(i+1) - \mathbf{X}^*\|_F$  can be expressed in terms of the distance from  $\mathbf{V}(i)$  to  $\mathbf{X}^*$  as follows:*

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F &\leq \sqrt{1 + 3\delta_{3k}^2} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\quad + \sqrt{1 + 3\delta_{3k}^2} \sqrt{\frac{3(1 + \delta_{2k})}{1 + 3\delta_{3k}^2}} \|\boldsymbol{\epsilon}\|_2. \end{aligned} \tag{61}$$

*Proof* We observe the following

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F^2 &= \|\mathbf{X}(i+1) - \mathbf{V}(i) + \mathbf{V}(i) - \mathbf{X}^*\|_F^2 \\ &= \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 + \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F^2 \\ &\quad - 2\langle \mathbf{V}(i) - \mathbf{X}^*, \mathbf{V}(i) - \mathbf{X}(i+1) \rangle. \end{aligned} \tag{62}$$

Focusing on the right hand side of expression (62),  $\langle \mathbf{V}(i) - \mathbf{X}^*, \mathbf{V}(i) - \mathbf{X}(i+1) \rangle = \langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}(i+1)) \rangle$  can be similarly analysed as in Lemma 10 where we obtain the following expression:

$$\begin{aligned} &\|\langle \mathbf{V}(i) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}(i+1)) \rangle\| \\ &\leq \delta_{3k} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F \\ &\quad + \sqrt{1 + \delta_{2k}} \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F \|\boldsymbol{\epsilon}\|_2. \end{aligned} \tag{63}$$

Now, expression (62) can be further transformed as:

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F^2 &\stackrel{(i)}{\leq} \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 + \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F^2 \end{aligned}$$

$$\begin{aligned} &+ 2(\delta_{3k} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F \\ &+ \sqrt{1 + \delta_{2k}} \|\mathbf{V}(i) - \mathbf{X}(i+1)\|_F \|\boldsymbol{\epsilon}\|_2), \end{aligned} \tag{64}$$

where (i) is due to (63). Using Lemma 10, we further have:

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F^2 &\leq \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 + \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 \\ &\quad + 2(\delta_{3k} \|\mathbf{V}(i) - \mathbf{X}^*\|_F \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ &\quad + \sqrt{1 + \delta_{2k}} \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \|\boldsymbol{\epsilon}\|_2). \end{aligned} \tag{65}$$

Furthermore, replacing  $\|\mathcal{P}_{\mathcal{S}_i}(\mathbf{X}^* - \mathbf{V}(i))\|_F$  with its upper bound defined in (56), we get:

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F^2 &\stackrel{(i)}{\leq} (1 + 3\delta_{3k}^2) \left( \|\mathbf{V}(i) - \mathbf{X}^*\|_2 + \sqrt{\frac{3(1 + \delta_{2k})}{1 + 3\delta_{3k}^2}} \|\boldsymbol{\epsilon}\| \right)^2, \end{aligned} \tag{66}$$

where (i) is obtained by completing the squares and eliminating negative terms.  $\square$

Applying basic algebra tools in (61) and (51), we get:

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F &\leq \sqrt{\frac{1 + 3\delta_{3k}^2}{1 - \delta_{3k}^2}} \|\mathcal{P}_{\mathcal{S}_i^{\perp}}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ &\quad + \left( \frac{\sqrt{1 + 3\delta_{3k}^2}}{1 - \delta_{3k}} + \sqrt{3} \right) \sqrt{1 + \delta_{2k}} \|\boldsymbol{\epsilon}\|_2. \end{aligned}$$

Since  $\mathbf{V}(i) \in \text{span}(\mathcal{S}_i)$ , we observe  $\mathcal{P}_{\mathcal{S}_i^{\perp}}(\mathbf{V}(i) - \mathbf{X}^*) = -\mathcal{P}_{\mathcal{S}_i^{\perp}} \mathbf{X}^* = -\mathcal{P}_{\mathcal{X}^* \setminus (\mathcal{D}_i \cup \mathcal{X}_i)} \mathbf{X}^*$ . Then, using Lemma 6, we obtain:

$$\begin{aligned} \|\mathbf{X}(i+1) - \mathbf{X}^*\|_F &\leq (2\delta_{2k} + 2\delta_{3k}) \sqrt{\frac{1 + 3\delta_{3k}^2}{1 - \delta_{3k}^2}} \|\mathbf{X}^* - \mathbf{X}(i)\|_F \\ &\quad + \left[ \sqrt{\frac{1 + 3\delta_{3k}^2}{1 - \delta_{3k}^2}} \sqrt{2(1 + \delta_{3k})} \right. \\ &\quad \left. + \left( \frac{\sqrt{1 + 3\delta_{3k}^2}}{1 - \delta_{3k}} + \sqrt{3} \right) \sqrt{1 + \delta_{2k}} \right] \|\boldsymbol{\epsilon}\|_2. \end{aligned} \tag{67}$$

Given  $\delta_{2k} \leq \delta_{3k}$ ,  $\rho$  is upper bounded by  $\rho < 4\delta_{3k} \sqrt{\frac{1 + 3\delta_{3k}^2}{1 - \delta_{3k}^2}}$ .

Then,  $4\delta_{3k} \sqrt{\frac{1 + 3\delta_{3k}^2}{1 - \delta_{3k}^2}} < 1 \Leftrightarrow \delta_{3k} < 0.2267$ .

A.4 Proof of Theorem 3

Let  $\mathcal{X}^* \leftarrow \mathcal{P}_k(\mathbf{X}^*)$  be a set of orthonormal, rank-1 matrices that span the range of  $\mathbf{X}^*$ . In Algorithm 3,  $\mathbf{X}(i + 1)$  is the best rank- $k$  approximation of  $\mathbf{V}(i)$ . Thus:

$$\begin{aligned} \|\mathbf{X}(i + 1) - \mathbf{V}(i)\|_F^2 &\leq \|\mathbf{X}^* - \mathbf{V}(i)\|_F^2 \\ \Rightarrow \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F^2 &\leq 2\langle \mathbf{X}(i + 1) - \mathbf{X}^*, \mathbf{V}(i) - \mathbf{X}^* \rangle. \end{aligned} \tag{68}$$

From Algorithm 3, (i)  $\mathbf{V}(i) \in \text{span}(\mathcal{S}_i)$ , (ii)  $\mathbf{Q}_i \in \text{span}(\mathcal{S}_i)$  and (iii)  $\mathbf{W}(i) \in \text{span}(\mathcal{S}_i)$ . We define  $\mathcal{E} \leftarrow \text{ortho}(\mathcal{S}_i \cup \mathcal{X}^*)$  where we observe  $\text{rank}(\text{span}(\mathcal{E})) \leq 4k$  and let  $\mathcal{P}_{\mathcal{E}}$  be the orthogonal projection onto the subspace defined by  $\mathcal{E}$ .

Since  $\mathbf{X}(i + 1) - \mathbf{X}^* \in \text{span}(\mathcal{E})$  and  $\mathbf{V}(i) - \mathbf{X}^* \in \text{span}(\mathcal{E})$ , the following hold true:

$$\mathbf{X}(i + 1) - \mathbf{X}^* = \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i + 1) - \mathbf{X}^*),$$

and,

$$\mathbf{V}(i) - \mathbf{X}^* = \mathcal{P}_{\mathcal{E}}(\mathbf{V}(i) - \mathbf{X}^*).$$

$$\begin{aligned} g(i + 1) &\leq \left[ b_1 \left( \frac{\alpha(1 + \tau_i) + \sqrt{\Delta}}{2} \right)^{i+1} + b_2 \left( \frac{\alpha(1 + \tau_i) - \sqrt{\Delta}}{2} \right)^{i+1} \right] \|\mathbf{X}(0) - \mathbf{X}^*\|_F \\ &\leq \left[ (b_1 + b_2) \left( \frac{\alpha(1 + \tau_i) + \sqrt{\Delta}}{2} \right)^{i+1} \right] \\ &\quad \times \|\mathbf{X}(0) - \mathbf{X}^*\|_F. \end{aligned} \tag{69}$$

Then, (68) can be written as:

$$\begin{aligned} \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F^2 &\leq 2\langle \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i + 1) - \mathbf{X}^*), \mathcal{P}_{\mathcal{E}}(\mathbf{V}(i) - \mathbf{X}^*) \rangle \\ &= 2\langle \mathcal{P}_{\mathcal{E}}(\mathbf{X}(i + 1) - \mathbf{X}^*), \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i + \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{Q}_i) - \mathbf{X}^*) \rangle \\ &\stackrel{(i)}{=} 2\langle \mathbf{X}(i + 1) - \mathbf{X}^*, \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A}[\mathcal{P}_{\mathcal{S}_i} + \mathcal{P}_{\mathcal{S}_i^\perp}] \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) \rangle \\ &= 2\langle \mathbf{X}(i + 1) - \mathbf{X}^*, (\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) - 2\mu_i \langle \mathbf{X}(i + 1) - \mathbf{X}^*, \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) \rangle \rangle \\ &\stackrel{(ii)}{\leq} 2\|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F \\ &\quad \times \|(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F \end{aligned}$$

$$\begin{aligned} &+ 2\mu_i \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F \\ &\quad \times \|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F, \end{aligned} \tag{70}$$

where (i) is due to  $\mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) := \mathcal{P}_{\mathcal{S}_i} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) + \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)$  and (ii) follows from Cauchy-Schwarz inequality. Since  $\frac{1}{1 + \delta_{3k}} \leq \mu_i \leq \frac{1}{1 - \delta_{3k}}$ , Lemma 4 implies:

$$\begin{aligned} \lambda(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) &\in \left[ 1 - \frac{1 - \delta_{3k}}{1 + \delta_{3k}}, \frac{1 + \delta_{3k}}{1 - \delta_{3k}} - 1 \right] \\ &\leq \frac{2\delta_{3k}}{1 - \delta_{3k}}, \end{aligned}$$

and thus:

$$\begin{aligned} \|(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i}) \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F &\leq \frac{2\delta_{3k}}{1 - \delta_{3k}} \|\mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F. \end{aligned}$$

Furthermore, according to Lemma 5:

$$\begin{aligned} \|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F &\leq \delta_{4k} \|\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F \end{aligned}$$

since  $\text{rank}(\mathcal{P}_{\mathcal{K}} \mathbf{Q}) \leq 4k, \forall \mathbf{Q} \in \mathbb{R}^{m \times n}$  where  $\mathcal{K} \leftarrow \text{ortho}(\mathcal{E} \cup \mathcal{S}_i)$ . Since  $\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*) = \mathcal{P}_{\mathcal{X}^* \setminus (\mathcal{D}_i \cup \mathcal{X}_i)} \mathbf{X}^*$  where

$$\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{Q}_i^\perp} \nabla f(\mathbf{Q}_i)),$$

then:

$$\begin{aligned} \|\mathcal{P}_{\mathcal{S}_i^\perp} \mathcal{P}_{\mathcal{E}}(\mathbf{Q}_i - \mathbf{X}^*)\|_F &= \|\mathcal{P}_{\mathcal{X}^* \setminus (\mathcal{D}_i \cup \mathcal{X}_i)} \mathbf{X}^*\|_F \leq (2\delta_{3k} + 2\delta_{4k}) \|\mathbf{Q}_i - \mathbf{X}^*\|_F, \end{aligned} \tag{71}$$

using Lemma 6. Using the above in (70), we compute:

$$\begin{aligned} \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F &\leq \left( \frac{4\delta_{3k}}{1 - \delta_{3k}} + (2\delta_{3k} + 2\delta_{4k}) \frac{2\delta_{3k}}{1 - \delta_{3k}} \right) \|\mathbf{Q}_i - \mathbf{X}^*\|_F. \end{aligned} \tag{72}$$

Furthermore:

$$\begin{aligned} \|\mathbf{Q}_i - \mathbf{X}^*\|_F &= \|\mathbf{X}(i) + \tau_i(\mathbf{X}(i) - \mathbf{X}(i - 1))\|_F \\ &= \|(1 + \tau_i)(\mathbf{X}(i) - \mathbf{X}^*) + \tau_i(\mathbf{X}^* - \mathbf{X}(i - 1))\|_F \\ &\leq (1 + \tau_i) \|\mathbf{X}(i) - \mathbf{X}^*\|_F + \tau_i \|\mathbf{X}(i - 1) - \mathbf{X}^*\|_F. \end{aligned} \tag{73}$$

Combining (72) and (73), we get:

$$\begin{aligned} & \| \mathbf{X}(i + 1) - \mathbf{X}^* \|_F \\ & \leq (1 + \tau_i) \left( \frac{4\delta_{3k}}{1 - \delta_{3k}} + (2\delta_{3k} + 2\delta_{4k}) \frac{2\delta_{3k}}{1 - \delta_{3k}} \right) \\ & \quad \times \| \mathbf{X}(i) - \mathbf{X}^* \|_F \\ & \quad + \tau_i \left( \frac{4\delta_{3k}}{1 - \delta_{3k}} + (2\delta_{3k} + 2\delta_{4k}) \frac{2\delta_{3k}}{1 - \delta_{3k}} \right) \\ & \quad \times \| \mathbf{X}(i - 1) - \mathbf{X}^* \|_F. \end{aligned} \tag{74}$$

Let  $\alpha := \frac{4\delta_{3k}}{1 - \delta_{3k}} + (2\delta_{3k} + 2\delta_{4k}) \frac{2\delta_{3k}}{1 - \delta_{3k}}$  and  $g(i) := \| \mathbf{X}(i + 1) - \mathbf{X}^* \|_F$ . Then, (74) defines the following homogeneous recurrence:

$$g(i + 1) - \alpha(1 + \tau_i)g(i) + \alpha\tau_i g(i - 1) \leq 0. \tag{75}$$

Using the *method of characteristic roots* to solve the above recurrence, we assume that the homogeneous linear recursion has solution of the form  $g(i) = r^i$  for  $r \in \mathbb{R}$ . Thus, replacing  $g(i) = r^i$  in (75) and factoring out  $r^{(i-2)}$ , we form the following characteristic polynomial:

$$r^2 - \alpha(1 + \tau_i)r - \alpha\tau_i \leq 0. \tag{76}$$

Focusing on the worst case where (76) is satisfied with equality, we compute the roots  $r_{1,2}$  of the quadratic characteristic polynomial as:

$$r_{1,2} = \frac{\alpha(1 + \tau_i) \pm \sqrt{\Delta}}{2}, \quad \text{where } \Delta := \alpha^2(1 + \tau_i)^2 + 4\alpha\tau_i.$$

Then, as a general solution, we combine the above roots with unknown coefficients  $b_1, b_2$  to obtain (69). Using the initial condition  $g(0) := \| \mathbf{X}(0) - \mathbf{X}^* \|_F \stackrel{\mathbf{X}(0)=\mathbf{0}}{=} \| \mathbf{X}^* \|_F = 1$ , we get  $b_1 + b_2 = 1$ . Thus, we conclude to the following recurrence:

$$\| \mathbf{X}(i + 1) - \mathbf{X}^* \|_F \leq \left( \frac{\alpha(1 + \tau_i) + \sqrt{\Delta}}{2} \right)^{i+1}.$$

### A.5 Proof of Lemma 7

Let  $\mathcal{D}_i^\epsilon \leftarrow \mathcal{P}_k^\epsilon(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i)))$  and  $\mathcal{D}_i \leftarrow \mathcal{P}_k(\mathcal{P}_{\mathcal{X}_i^\perp} \nabla f(\mathbf{X}(i)))$ . Using Definition 4, the following holds

true:

$$\begin{aligned} & \| \mathcal{P}_{\mathcal{D}_i^\epsilon} \nabla f(\mathbf{X}(i)) - \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \leq (1 + \epsilon) \| \mathcal{P}_{\mathcal{D}_i} \nabla f(\mathbf{X}(i)) - \nabla f(\mathbf{X}(i)) \|_F^2. \end{aligned} \tag{77}$$

Furthermore, we observe:

$$\begin{aligned} & \| \nabla f(\mathbf{X}(i)) \|_F^2 \\ & = \| \nabla f(\mathbf{X}(i)) \|_F^2 \Leftrightarrow \| \mathcal{P}_{\mathcal{D}_i^\epsilon} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \quad + \| \mathcal{P}_{(\mathcal{D}_i^\epsilon)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & = \| \mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i)) \|_F^2 + \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2. \end{aligned} \tag{78}$$

Here, we use the notation defined in the proof of Lemma 6. Since  $\mathcal{P}_{\mathcal{D}_i} \nabla f(\mathbf{X}(i))$  is the best rank- $k$  approximation to  $\nabla f(\mathbf{X}(i))$ , we have:

$$\begin{aligned} & \| \mathcal{P}_{\mathcal{D}_i} \nabla f(\mathbf{X}(i)) - \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \leq \| \mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i)) - \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \Leftrightarrow \| \mathcal{P}_{\mathcal{D}_i^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \leq \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \Leftrightarrow (1 + \epsilon) \| \mathcal{P}_{\mathcal{D}_i^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \leq (1 + \epsilon) \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2, \end{aligned} \tag{79}$$

where  $\text{rank}(\text{span}(\text{ortho}(\mathcal{X}^* \setminus \mathcal{X}_i))) \leq k$ . Using (77) in (79), the following series of inequalities are observed:

$$\begin{aligned} & \| \mathcal{P}_{(\mathcal{D}_i^\epsilon)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \leq (1 + \epsilon) \| \mathcal{P}_{\mathcal{D}_i^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \leq (1 + \epsilon) \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2. \end{aligned} \tag{80}$$

Now, in (78), we compute the series of inequalities in (81)-(82).

$$\begin{aligned} & \| \mathcal{P}_{\mathcal{D}_i^\epsilon} \nabla f(\mathbf{X}(i)) \|_F^2 + \| \mathcal{P}_{(\mathcal{D}_i^\epsilon)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 = \| \mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \quad + \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \stackrel{(79)}{\Leftrightarrow} \tag{81} \\ & \| \mathcal{P}_{\mathcal{D}_i^\epsilon} \nabla f(\mathbf{X}(i)) \|_F^2 + (1 + \epsilon) \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \geq \| \mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i)) \|_F^2 \\ & \quad + \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \Leftrightarrow \\ & \| \mathcal{P}_{\mathcal{D}_i^\epsilon} \nabla f(\mathbf{X}(i)) \|_F^2 + \epsilon \| \mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i)) \|_F^2 \geq \| \mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i)) \|_F^2 \Leftrightarrow \end{aligned}$$

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{D}_i^c} \nabla f(\mathbf{X}(i))\|_F^2 + \|\mathcal{P}_{\mathcal{X}_i} \nabla f(\mathbf{X}(i))\|_F^2 + \epsilon \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i))\|_F^2 \geq \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \nabla f(\mathbf{X}(i))\|_F^2 + \|\mathcal{P}_{\mathcal{X}_i} \nabla f(\mathbf{X}(i))\|_F^2 \Leftrightarrow \\
 & \quad \|\mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i))\|_F^2 + \epsilon \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i))\|_F^2 \geq \|\mathcal{P}_{\mathcal{S}_i^*} \nabla f(\mathbf{X}(i))\|_F^2 \Leftrightarrow \\
 & \quad \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \nabla f(\mathbf{X}(i))\|_F^2 + \epsilon \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \nabla f(\mathbf{X}(i))\|_F^2 \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \nabla f(\mathbf{X}(i))\|_F^2 \Leftrightarrow \\
 & \quad \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F^2 + \epsilon \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F^2 \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F^2 \Leftrightarrow \\
 & \quad \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F + \sqrt{\epsilon} \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F. \tag{82}
 \end{aligned}$$

Focusing on  $\|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i}^\perp \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F$ , we observe:

$$-\|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \boldsymbol{\epsilon}\|_F. \tag{85}$$

$$\begin{aligned}
 & \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \mathcal{A}^*(\mathbf{y} - \mathcal{A}\mathbf{X}(i))\|_F \\
 & = \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \mathcal{A}^*(\mathcal{A}\mathbf{X}^* + \boldsymbol{\epsilon} - \mathcal{A}\mathbf{X}(i))\|_F \\
 & \leq \|\mathcal{P}_{(\mathcal{X}^* \setminus \mathcal{X}_i)^\perp} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i))\|_F + \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i}^\perp \mathcal{A}^* \boldsymbol{\epsilon}\|_F \\
 & \leq \|\mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i))\|_F + \|\mathcal{A}^* \boldsymbol{\epsilon}\|_F \leq 2\lambda. \tag{83}
 \end{aligned}$$

Moreover, we know the following hold true from Lemma 6:

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\epsilon}\|_F \\
 & \leq 2\delta_{3k} \|\mathbf{X}^* - \mathbf{X}(i)\|_F + \|\mathcal{P}_{\mathcal{S}_i \setminus \mathcal{S}_i^*} \mathcal{A}^* \boldsymbol{\epsilon}\|_F \tag{84}
 \end{aligned}$$

and

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \mathcal{A}(\mathbf{X}^* - \mathbf{X}(i)) + \mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathcal{A}^* \boldsymbol{\epsilon}\|_F \\
 & \geq \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} (\mathbf{X}^* - \mathbf{X}(i))\|_F - 2\delta_{2k} \|\mathbf{X}(i) - \mathbf{X}^*\|_F
 \end{aligned}$$

Combining (83)–(85) in (82), we obtain:

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i^* \setminus \mathcal{S}_i} \mathbf{X}^*\|_F = \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{S}_i} \mathbf{X}^*\|_F \\
 & \leq (2\delta_{2k} + 2\delta_{3k}) \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 & \quad + \sqrt{2(1 + \delta_{2k})} \|\boldsymbol{\epsilon}\|_2 + 2\lambda\sqrt{\epsilon}.
 \end{aligned}$$

#### A.6 Proof of Theorem 4

To prove Theorem 4, we combine the following series of lemmas for each step of Algorithm 1.

**Lemma 12** [Error norm reduction via gradient descent] *Let  $\mathcal{S}_i \leftarrow \text{ortho}(\mathcal{X}_i \cup \mathcal{D}_i^c)$  be a set of orthonormal, rank-1 matrices that span a rank-2k subspace in  $\mathbb{R}^{m \times n}$ . Then (86) holds.*

$$\begin{aligned}
 \|\mathbf{V}(i) - \mathbf{X}^*\|_F & \leq \left[ \left(1 + \frac{\delta_{3k}}{1 - \delta_{2k}}\right) (2\delta_{2k} + 2\delta_{3k} + \delta_k) + \frac{2\delta_{2k}}{1 - \delta_{2k}} \right] \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\
 & \quad + \left[ \left(1 + \frac{\delta_{3k}}{1 - \delta_{2k}}\right) \sqrt{2(1 + \delta_{2k})} + \frac{\sqrt{1 + \delta_{2k}}}{1 - \delta_{2k}} \right] \|\boldsymbol{\epsilon}\|_2 + \left(1 + \frac{\delta_{3k}}{1 - \delta_{2k}}\right) 2\lambda\sqrt{\epsilon}. \tag{86}
 \end{aligned}$$

*Proof* We observe the following:

$$\begin{aligned}
 \|\mathbf{V}(i) - \mathbf{X}^*\|_F^2 & = \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 \\
 & \quad + \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2. \tag{87}
 \end{aligned}$$

The following equations hold true:

$$\|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F^2 = \|\mathcal{P}_{\mathcal{S}_i^\perp} \mathbf{X}^*\|_F^2 = \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{S}_i} \mathbf{X}^*\|_F^2.$$

Furthermore, we compute:

$$\begin{aligned}
 & \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\
 & = \left\| \mathcal{P}_{\mathcal{S}_i} \left( \mathbf{X}(i) - \frac{\mu_i}{2} \mathcal{P}_{\mathcal{S}_i} \nabla f(\mathbf{X}(i)) - \mathbf{X}^* \right) \right\|_F
 \end{aligned}$$

$$\begin{aligned}
 & = \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{X}(i) - \mathbf{X}^*) - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 & \quad + \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \boldsymbol{\epsilon}\|_F \\
 & \leq \|(\mathbf{I} - \mu_i \mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i} \mathcal{P}_{\mathcal{S}_i})(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 & \quad + \mu_i \|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \mathcal{A} \mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 & \quad + \mu_i \|\mathcal{P}_{\mathcal{S}_i} \mathcal{A}^* \boldsymbol{\epsilon}\|_F \\
 & \stackrel{(i)}{\leq} \frac{2\delta_{2k}}{1 - \delta_{2k}} \|\mathcal{P}_{\mathcal{S}_i}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \\
 & \quad + \frac{\delta_{3k}}{1 - \delta_{2k}} \|\mathcal{P}_{\mathcal{S}_i^\perp}(\mathbf{X}(i) - \mathbf{X}^*)\|_F
 \end{aligned}$$

$$+ \frac{\sqrt{1 + \delta_{2k}}}{1 - \delta_{2k}} \|\mathbf{e}\|_2, \tag{88}$$

where (i) is due to Lemmas 2, 4, 5 and  $\frac{1}{1+\delta_{2k}} \leq \mu_i \leq \frac{1}{1-\delta_{2k}}$ . Using the subadditivity property of the square root in (87), (88), Lemma 7 and the fact that  $\|\mathcal{P}_{S_i}(\mathbf{X}(i) - \mathbf{X}^*)\|_F \leq \|\mathbf{X}(i) - \mathbf{X}^*\|_F$ , we obtain:

$$\begin{aligned} & \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ & \leq \|\mathcal{P}_{S_i}(\mathbf{V}(i) - \mathbf{X}^*)\|_F + \|\mathcal{P}_{S_i^\perp}(\mathbf{V}(i) - \mathbf{X}^*)\|_F \\ & \leq \hat{\rho} \|\mathbf{X}(i) - \mathbf{X}^*\|_F \\ & \quad + \left(1 + \frac{\delta_{3k}}{1 - \delta_{2k}}\right) \sqrt{\epsilon} \|\mathcal{P}_{\mathcal{X}^* \setminus \mathcal{X}_i} \mathcal{A}^* \mathbf{e}\|_F \\ & \quad + \left[\left(1 + \frac{\delta_{3k}}{1 - \delta_{2k}}\right) \sqrt{2(1 + \delta_{2k})} + \frac{\sqrt{1 + \delta_{2k}}}{1 - \delta_{2k}}\right] \|\mathbf{e}\|_2, \end{aligned} \tag{89}$$

where  $\hat{\rho} := (1 + \frac{\delta_{3k}}{1-\delta_{2k}})(2\delta_{2k} + 2\delta_{3k}) + \frac{2\delta_{2k}}{1-\delta_{2k}}$ .  $\square$

We exploit Lemma 8 to obtain the following inequalities:

$$\begin{aligned} \|\widehat{\mathbf{W}}_i - \mathbf{X}^*\|_F &= \|\widehat{\mathbf{W}}_i - \mathbf{V}(i) + \mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\leq \|\widehat{\mathbf{W}}_i - \mathbf{V}(i)\|_F + \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\leq (1 + \epsilon) \|\mathbf{W}(i) - \mathbf{V}(i)\|_F + \|\mathbf{V}(i) - \mathbf{X}^*\|_F \\ &\leq (2 + \epsilon) \|\mathbf{V}(i) - \mathbf{X}^*\|_F, \end{aligned} \tag{90}$$

where the last inequality holds since  $\mathbf{W}(i)$  is the best rank- $k$  matrix estimate of  $\mathbf{V}(i)$  and, thus,  $\|\mathbf{W}(i) - \mathbf{V}(i)\|_F \leq \|\mathbf{V}(i) - \mathbf{X}^*\|_F$ .

Following similar motions for steps 6 and 7 in Matrix ALPS I, we obtain:

$$\begin{aligned} & \|\mathbf{X}(i + 1) - \mathbf{X}^*\|_F \\ & \leq \left(1 + \frac{2\delta_k}{1 - \delta_k} + \frac{\delta_{2k}}{1 - \delta_k}\right) \|\widehat{\mathbf{W}}_i - \mathbf{X}^*\|_F \\ & \quad + \frac{\sqrt{1 + \delta_k}}{1 - \delta_k} \|\mathbf{e}\|_2. \end{aligned} \tag{91}$$

Combining (91), (90) and (89), we obtain the desired inequality.

**References**

1. Baraniuk, R.G., Cevher, V., Wakin, M.B.: Low-dimensional models for dimensionality reduction and signal recovery: a geometric perspective. *Proc. IEEE* **98**(6), 959–971 (2010)
2. Candès, E.J., Recht, B.: Exact matrix completion via convex optimization. *Found. Comput. Math.* **9**(6), 717–772 (2009)
3. Meka, R., Jain, P., Dhillon, I.S.: Guaranteed rank minimization via singular value projection. In: *NIPS Workshop on Discrete Optimization in Machine Learning* (2010)

4. Tyagi, H., Cevher, V.: Learning ridge functions with randomized sampling in high dimensions. In: *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2025–2028. IEEE Press, New York (2012)
5. Tyagi, H., Cevher, V.: Learning non-parametric basis independent models from point queries via low-rank methods. Technical report, EPFL (2012)
6. Liu, Y.K.: Universal low-rank matrix recovery from Pauli measurements (2011)
7. Tyagi, H., Cevher, V.: Active learning of multi-index function models. In: *Advances in Neural Information Processing Systems*, vol. 25, pp. 1475–1483 (2012)
8. Candès, E.J., Li, X.: Solving quadratic equations via phaselift when there are about as many equations as unknowns (2012). Preprint arXiv:1208.6247
9. Bennett, J., Lanning, S.: The netflix prize. In: *KDD Cup and Workshop in Conjunction with KDD* (2007)
10. Candès, E.J., Li, X., Ma, Y., Wright, J.: Robust principal component analysis? *J. ACM* **58**(3) (2011)
11. Kyrillidis, A., Cevher, V.: Matrix alps: Accelerated low rank and sparse matrix reconstruction. Technical report, EPFL (2012)
12. Waters, A.E., Sankaranarayanan, A.C., Baraniuk, R.G.: Sparcs: recovering low-rank and sparse matrices from compressive measurements. In: *NIPS* (2011)
13. Fazel, M., Recht, B., Parrilo, P.A.: Guaranteed minimum rank solutions to linear matrix equations via nuclear norm minimization. *SIAM Rev.* **52**(3), 471–501 (2010)
14. Liu, Z., Vandenberghe, L.: Interior-point method for nuclear norm approximation with application to system identification. *SIAM J. Matrix Anal. Appl.* **31**, 1235–1256 (2009)
15. Mohan, K., Fazel, M.: Reweighted nuclear norm minimization with application to system identification. In: *American Control Conference (ACC)*. IEEE Press, New York (2010)
16. Cai, J.-F., Candès, E.J., Shen, Z.: A singular value thresholding algorithm for matrix completion. *SIAM J. Optim.* **20**, 1956–1982 (2010)
17. Recht, B., Re, C.: Parallel stochastic gradient algorithms for large-scale matrix completion. Preprint (2011)
18. Lin, Z., Chen, M., Ma, Y.: The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices (2010). preprint arXiv:1009.5055
19. Wright, J., Wu, L., Chen, M., Lin, Z., Ganesh, A., Ma, Y.: Fast convex optimization algorithms for exact recovery of a corrupted low-rank matrix. UIUC Technical Report UILU-ENG-09-2214
20. Natarajan, B.K.: Sparse approximate solutions to linear systems. *SIAM J. Comput.* **24**(2), 227–234 (1995)
21. Lee, K., Bresler, Y.: Admira: atomic decomposition for minimum rank approximation. *IEEE Trans. Inf. Theory* **56**(9), 4402–4416 (2010)
22. Goldfarb, D., Ma, S.: Convergence of fixed-point continuation algorithms for matrix rank minimization. *Found. Comput. Math.* **11**, 183–210 (2011)
23. Beck, A., Teboulle, M.: A linearly convergent algorithm for solving a class of nonconvex/affine feasibility problems. In: *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, pp. 33–48 (2011)
24. Kyrillidis, A., Cevher, V.: Recipes on hard thresholding methods. In: *Computational Advances in Multi-Sensor Adaptive Processing*, Dec. 2011
25. Kyrillidis, A., Cevher, V.: Combinatorial selection and least absolute shrinkage via the Clash algorithm. In: *IEEE International Symposium on Information Theory*, July 2012
26. Halko, N., Martinsson, P.G., Tropp, J.A.: Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.* **53**, 217–288 (2011)
27. Bertsekas, D.: *Nonlinear Programming*. Athena Scientific, Nashua (1995)



28. Horn, R.A., Johnson, C.R.: *Matrix Analysis*. Cambridge Univ. Press, Cambridge (1990)
29. Cohen, A., Dahmen, W., DeVore, R.: Compressed sensing and best k-term approximation. *J. Am. Math. Soc.* **22**(1), 211–231 (2009)
30. Tropp, J.A.: Greed is good: algorithmic results for sparse approximation. *IEEE Trans. Inf. Theory* **50**(10), 2231–2242 (2004)
31. Cevher, V.: An alps view of sparse recovery. In: *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5808–5811. IEEE Press, New York (2011)
32. Needell, D., Tropp, J.A.: Cosamp: iterative signal recovery from incomplete and inaccurate samples. *Appl. Comput. Harmon. Anal.* **26**(3), 301–321 (2009)
33. Dai, W., Milenkovic, O.: Subspace pursuit for compressive sensing signal reconstruction. *IEEE Trans. Inf. Theory* **55**, 2230–2249 (2009)
34. Foucart, S.: Hard thresholding pursuit: an algorithm for compressed sensing. *SIAM J. Numer. Anal.* **49**(6), 2543–2563 (2011)
35. Blumensath, T., Davies, M.E.: Iterative hard thresholding for compressed sensing. *Appl. Comput. Harmon. Anal.* **27**(3), 265–274 (2009)
36. Garg, R., Khandekar, R.: Gradient descent with sparsification: an iterative algorithm for sparse recovery with restricted isometry property. In: *ICML*. ACM Press, New York (2009)
37. Blumensath, T., Davies, M.E.: Normalized iterative hard thresholding: guaranteed stability and performance. *IEEE J. Sel. Top. Signal Process.* **4**(2), 298–309 (2010)
38. Blumensath, T.: Accelerated iterative hard thresholding. *Signal Process.* **92**, 752–756 (2012)
39. Tanner, J., Wei, K.: Normalized iterative hard thresholding for matrix completion. Preprint (2012)
40. Coifman, R., Geshwind, F., Meyer, Y.: Noiselets. *Appl. Comput. Harmon. Anal.* **10**(1), 27–44 (2001)
41. Foucart, S.: Sparse recovery algorithms: sufficient conditions in terms of restricted isometry constants. In: *Proceedings of the 13th International Conference on Approximation Theory* (2010)
42. Nesterov, Y.: Gradient methods for minimizing composite objective function. core discussion papers 2007076, universit  catholique de louvain. Center for Operations Research and Econometrics (CORE) (2007)
43. Nesterov, Y.: *Introductory Lectures on Convex Optimization*. Kluwer Academic, Dordrecht (1996)
44. Drineas, P., Frieze, A., Kannan, R., Vempala, S., Vinay, V.: Clustering large graphs via the singular value decomposition. *Mach. Learn.* **56**(1), 9–33 (2004)
45. Drineas, P., Kannan, R., Mahoney, M.W.: Fast Monte Carlo algorithms for matrices ii: computing a low-rank approximation to a matrix. *SIAM J. Comput.* **36**, 158–183 (2006)
46. Deshpande, A., Rademacher, L., Vempala, S., Wang, G.: Matrix approximation and projective clustering via volume sampling. In: *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm, SODA '06*, New York, NY, USA, pp. 1117–1126. ACM Press, New York (2006)
47. Deshpande, A., Vempala, S.: Adaptive sampling and fast low-rank matrix approximation. *Electron. Colloq. Comput. Complex.* **13**, 042 (2006)
48. Keshavan, R.H., Montanari, A., Oh, S.: Matrix completion from a few entries. *IEEE Trans. Inf. Theory* **56**(6), 2980–2998 (2010)
49. Balzano, L., Nowak, R., Recht, B.: Online identification and tracking of subspaces from highly incomplete information. In: *48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pp. 704–711. IEEE Press, New York (2010)
50. He, J., Balzano, L., Lui, J.C.S.: Online robust subspace tracking from partial information (2011). [arXiv:1109.3827](https://arxiv.org/abs/1109.3827)
51. Boumal, N., Absil, P.A.: Rtrmc: a Riemannian trust-region method for low-rank matrix completion. In: *NIPS* (2011)
52. Wen, Z., Yin, W., Zhang, Y.: Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. Rice University CAAM Technical Report TR10-07 (2010) Submitted
53. Larsen, R.M.: Propack: Software for large and sparse svd calculations. <http://soi.stanford.edu/~rmunk/PROPACK>
54. Shi, X., Yu, P.S.: Limitations of matrix completion via trace norm minimization. *ACM SIGKDD Explor. Newsl.* **12**(2), 16–20 (2011)