

Non-convex matrix sensing: Breaking the quadratic rank barrier in the sample complexity

Anonymous CPAL submission

1 For the problem of reconstructing a low-rank matrix from a few linear measure-
2 ments, two classes of algorithms have been widely studied in the literature: convex
3 approaches based on nuclear norm minimization, and non-convex approaches that
4 use factorized gradient descent. Under certain statistical model assumptions, it is
5 known that nuclear norm minimization recovers the ground truth as soon as the
6 number of samples scales linearly with the number of degrees of freedom of the
7 ground-truth. In contrast, while non-convex approaches are computationally less
8 expensive, existing recovery guarantees assume that the number of samples scales
9 at least quadratically with the rank r of the ground-truth matrix. In this paper, we
10 close this gap by showing that the non-convex approaches can be as efficient as nu-
11 clear norm minimization in terms of sample complexity. Namely, we consider the
12 problem of reconstructing a positive semidefinite matrix from a few Gaussian mea-
13 surements. We show that factorized gradient descent with spectral initialization
14 converges to the ground truth with a linear rate as soon as the number of samples
15 scales with $\Omega(r d \kappa^2)$, where d is the dimension, and κ is the condition number of the
16 ground truth matrix. This improves the previous rank-dependence in the sample
17 complexity of non-convex matrix factorization from quadratic to linear. Our proof
18 relies on a probabilistic decoupling argument, where we show that the gradient
19 descent iterates are only weakly dependent on the individual entries of the mea-
20 surement matrices. We expect that our proof technique is of independent interest
21 for other non-convex problems.

22 1. Introduction

23 Low-rank matrix recovery refers to the problem of reconstructing an unknown matrix $\mathbf{X}_* \in \mathbb{R}^{d_1 \times d_2}$
24 with $\text{rank}(\mathbf{X}_*) =: r \ll \min\{d_1, d_2\}$ from an underdetermined linear set of equations of the form

$$\mathbf{y} = \mathcal{A}(\mathbf{X}_*) \in \mathbb{R}^m,$$

25 where \mathcal{A} represents a known linear measurement operator and $\mathbf{y} \in \mathbb{R}^m$ are the observations. This
26 ill-posed inverse problem has been the topic of intense study over many years, given its relevance
27 to a variety of questions in machine learning, signal processing, and statistics. Notable applications
28 include matrix completion [1], phase retrieval [2], robust PCA [3], blind deconvolution [4] and
29 its extension to blind demixing [5]. A major goal has been to develop methods which are *sample-*
30 *efficient*; that is, they can reconstruct the low-rank matrix \mathbf{X}_* if the number of observations m is
31 roughly of the same order as the number of degrees of freedom of \mathbf{X}_* . In addition, these methods
32 should also be scalable, meaning they remain computationally efficient as the problem dimensions
33 are increasing.

34 Several different algorithmic approaches to solve this problem have been proposed. One line of
35 research revolves around the idea of convex relaxation. Here, the nuclear norm $\|\cdot\|_*$, i.e., the sum
36 of singular values, is considered as a convex proxy for the rank function. For many problem classes,
37 including matrix sensing [6], matrix completion [7, 8], and blind deconvolution and demixing [9], it
38 has been shown that this approach is able to recover the unknown matrix \mathbf{X}_* as soon as the number
39 of samples m scales, up to logarithmic factors, with the information-theoretically optimal sample

40 complexity $r(d_1 + d_2)$. However, a drawback of these convex approaches is that they tend to be
 41 computationally prohibitive.

42 For this reason, many studies have considered non-convex heuristics where one minimizes an ob-
 43 jective of the form

$$f(\mathbf{U}, \mathbf{V}) = \sum_{i=1}^m \ell(y_i, (\mathcal{A}(\mathbf{U}\mathbf{V}^\top))_i), \quad (1)$$

44 with low-rank factors $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{V} \in \mathbb{R}^{d_2 \times r}$ and a loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. To minimize the
 45 objective function, local search methods such as gradient descent or alternating minimization with
 46 a suitable initialization are used. An advantage of these approaches is that they are computationally
 47 less demanding since there are only $r(d_1 + d_2)$ optimization variables instead of at least $d_1 d_2$ opti-
 48 mization variables in the convex approaches. However, due to the non-convexity of the objective
 49 function, it might initially seem unclear that local search methods can find the global minimum of
 50 the objective (1) efficiently.

51 Nevertheless, in recent years a large body of literature has demonstrated that under certain statisti-
 52 cal assumptions, these methods converge to the global minimum and are thus able to recover the
 53 unknown low-rank matrix \mathbf{X}_* . For instance, gradient descent with spectral initialization [10] and
 54 other variants of gradient descent [11–13] have been studied for matrix sensing and related prob-
 55 lems. Similarly, numerous works have established convergence and recovery guarantees for matrix
 56 completion [14–19] and blind deconvolution and demixing [20, 21]. In addition, recent studies
 57 also analyzed overparameterized models, where the exact rank r is either not known or where the
 58 number of parameters exceeds the number of samples [22–28]. Beyond gradient descent, also al-
 59 ternating minimization [29] and other non-convex methods based on matrix factorization such as
 60 GNMR [30] have been proposed and studied. For a more extensive overview of the literature, we
 61 refer the reader to [19].

62 Despite this significant body of literature, the existing theoretical guarantees for non-convex meth-
 63 ods based on matrix factorization in the literature are weaker than the corresponding guarantees
 64 for nuclear norm minimization in terms of sample complexity. Namely, in all these results, it is re-
 65 quired that the number of samples m scales at least quadratically with the rank r and thus the total
 66 number of samples scales at least with $r^2(d_1 + d_2)$. This raises the question of whether this quadratic
 67 rank-dependence is just an artifact of the proof or whether it is inherent to the problem, see, e.g.,
 68 [31, p. 5264].

69 In this paper, we resolve this question in the context of symmetric matrix sensing. Under the as-
 70 sumption that \mathcal{A} is a Gaussian measurement operator and $\mathbf{X}_* \in \mathbb{R}^{d \times d}$ is symmetric and positive
 71 semidefinite, we show that factorized gradient descent with spectral initialization is able to recover
 72 the unknown matrix \mathbf{X}_* if the number of samples scales with rd , which, in particular, is linear in the
 73 rank of \mathbf{X}_* . Our proof is based on a novel probabilistic decoupling argument. Namely, we show that
 74 the trajectory of the gradient descent iterates depends only weakly on any given generalized entry
 75 of the measurement matrices in a suitable sense. This allows us to prove stronger concentration
 76 bounds than what would be possible if one were to rely solely on uniform concentration bounds
 77 (such as the Restricted Isometry Property, for example). To establish this weak dependence, we
 78 construct auxiliary virtual sequences and combine this with an ε -net argument. Our novel proof
 79 approach paves the way to improved sample complexity bounds for other non-convex algorithms
 80 and beyond.

81 Finally, we note that there are also several non-convex algorithms for low-rank matrix recovery that
 82 are not explicitly based on matrix factorization formulation as in equation (1). This includes, for ex-
 83 ample, Singular Value Projection [32, 33], Normalized Iterative Hard Thresholding [34], Iteratively
 84 Reweighted Least Squares (IRLS), see, e.g., [35–38], and Atomic Decomposition for Minimum Rank
 85 Approximation (ADMIRA) [39]. However, since many of these algorithms operate in the full ma-
 86 trix space they are less computationally efficient than algorithms based on matrix factorization. In
 87 the case of IRLS, only local convergence guarantees (with explicit convergence rates) are known.
 88 There have also been algorithms studied that are based on Riemannian optimization, see, e.g., [40–

89 42]. However, these algorithms require that the sample complexity scales quadratically in the rank
 90 r . We believe our work can lead to improved sample size guarantees for these methods as well.

91 **Organization of the paper:** This paper is structured as follows. In the remainder of Section 1,
 92 we will describe the formal setting and the algorithm, and we will state our main theoretical result,
 93 which is Theorem 1.2. In Section 2, we discuss some technical preliminaries regarding the Restricted
 94 Isometry Property and perturbation bounds for eigenspaces. In Section 3, we discuss the proof
 95 strategy, and we introduce the virtual sequences, which are the main ingredient to establish that
 96 the sample complexity depends only linearly on the rank. Section 4 contains the proof of the main
 97 result of this paper, Theorem 1.2. We discuss interesting directions for future research in Section 5.

98 **Notation:** Before we state the problem formulation, we introduce some basic notation. For a matrix
 99 $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$, we denote its transpose by \mathbf{A}^\top and its trace by $\text{trace}(\mathbf{A})$. For matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d_1 \times d_2}$,
 100 we define their inner product via $\langle \mathbf{A}, \mathbf{B} \rangle := \text{trace}(\mathbf{A}\mathbf{B}^\top)$. The Frobenius norm $\|\cdot\|_F$ denotes the
 101 norm induced by this inner product, i.e., $\|\mathbf{A}\|_F := \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$. By $\|\mathbf{A}\|$ we denote the spectral norm
 102 of the matrix \mathbf{A} , i.e., the largest singular value of the matrix \mathbf{A} . By $\|\mathbf{v}\|_2 := \sqrt{\sum_{i=1}^d v_i^2}$ we denote the
 103 Euclidean norm of a vector $\mathbf{v} \in \mathbb{R}^d$. The set $\mathcal{S}^d \subset \mathbb{R}^{d \times d}$ represents the set of all symmetric matrices.
 104 The matrix $\mathbf{Id} \in \mathcal{S}^d$ denotes the identity matrix. Moreover, $\mathcal{I} : \mathcal{S}^d \rightarrow \mathcal{S}^d$ represents the identity
 105 mapping.

106 Furthermore, for a matrix $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ of rank r we denote its singular value decomposition by $\mathbf{A} =$
 107 $\mathbf{V}_\mathbf{A} \Sigma_\mathbf{A} \mathbf{W}_\mathbf{A}^\top$. The matrices $\mathbf{V}_\mathbf{A} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{W}_\mathbf{A} \in \mathbb{R}^{d_2 \times r}$ contain the left-singular and right-singular
 108 vectors of the matrix \mathbf{A} . The matrix $\Sigma_\mathbf{A} \in \mathbb{R}^{r \times r}$ contains the singular values of \mathbf{A} . Moreover,
 109 $\mathbf{V}_{\mathbf{A}, \perp} \in \mathbb{R}^{(d_1-r) \times r}$ represents an orthogonal matrix whose column span is orthogonal to the column
 110 span of $\mathbf{V}_\mathbf{A}$.

111 1.1. Problem formulation

112 In this paper, we focus on symmetric matrix sensing. More precisely, we study the problem of
 113 reconstructing a symmetric, positive semidefinite matrix $\mathbf{X}_\star \in \mathbb{R}^{d \times d}$ with rank r from m linear
 114 observations of the form

$$\mathbf{y}_i = \frac{1}{\sqrt{m}} \langle \mathbf{A}_i, \mathbf{X}_\star \rangle := \frac{1}{\sqrt{m}} \text{trace}(\mathbf{A}_i \mathbf{X}_\star) \quad \text{for } i = 1, 2, \dots, m. \quad (2)$$

115 **Definition 1.1** (Measurement operator). We define the linear measurement operator $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$ by

$$[\mathcal{A}(\mathbf{X})]_i := \frac{1}{\sqrt{m}} \langle \mathbf{A}_i, \mathbf{X} \rangle \quad \text{for } i = 1, 2, \dots, m$$

116 for any matrix $\mathbf{X} \in \mathcal{S}^d$. Recall that $\mathcal{S}^d \subset \mathbb{R}^{d \times d}$ denotes the set of symmetric matrices. The matrices
 117 $\{\mathbf{A}_i\}_{i=1}^m \subset \mathbb{R}^{d \times d}$ represent known, symmetric measurement matrices. We assume that their entries are i.i.d.
 118 with distribution $\mathcal{N}(0, 1)$ on the diagonal and $\mathcal{N}(0, 1/2)$ on the off-diagonal entries. Each \mathbf{A}_i is also known
 119 as a Gaussian orthogonal ensemble [43].

120 This measurement model has been considered before in, e.g., [10, 22]. With this notation in place,
 121 equation (2) can be written more compactly as

$$\mathbf{y} = \mathcal{A}(\mathbf{X}_\star).$$

122 To recover the ground-truth matrix \mathbf{X}_\star , we consider the non-convex objective function

$$\mathcal{L}(\mathbf{U}) := \frac{1}{4} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2 = \frac{1}{4} \|\mathcal{A}(\mathbf{X}_\star - \mathbf{U}\mathbf{U}^\top)\|_2^2, \quad (3)$$

123 where $\mathbf{U} \in \mathbb{R}^{d \times r}$ is a matrix and $\|\cdot\|_2$ denotes the ℓ_2 -norm of a vector. To minimize this objective, we
 124 follow the two-stage approach introduced in [14] for matrix completion, which then subsequently
 125 was studied for matrix sensing in [10]. In the first stage, an initialization \mathbf{U}_0 is constructed via
 126 a so-called spectral initialization. This initialization is subsequently used as a starting point for

127 the gradient descent scheme in the second stage. To precisely define the spectral initialization, we
 128 denote by $\mathcal{A}^* : \mathbb{R}^m \rightarrow \mathcal{S}^d$ the adjoint operator of \mathcal{A} with respect to the trace inner product defined
 129 in equation (2).

130 With this definition in place, we can consider the eigendecomposition of the matrix

$$\mathcal{A}^*(\mathbf{y}) =: \tilde{\mathbf{V}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{V}}^\top,$$

131 where $\tilde{\mathbf{V}} \in \mathbb{R}^{d \times d}$ is an orthogonal matrix and the matrix $\tilde{\mathbf{\Lambda}} \in \mathbb{R}^{d \times d}$ is diagonal matrix which con-
 132 tains the eigenvalues of $\mathcal{A}^*(\mathbf{y})$ sorted by their magnitude, i.e., $|\lambda_1(\mathcal{A}^*(\mathbf{y}))| \geq |\lambda_2(\mathcal{A}^*(\mathbf{y}))| \geq \dots \geq$
 133 $|\lambda_d(\mathcal{A}^*(\mathbf{y}))|$.

134 Since the measurement matrices \mathbf{A}_i are Gaussian we have that

$$\mathbb{E}[\mathcal{A}^*(\mathbf{y})] = \mathbb{E}[(\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)] = \mathbf{X}_*.$$

135 Since \mathbf{X}_* has rank r for a large enough sample size m , one has that the truncated rank- r
 136 eigendecomposition of $\mathcal{A}^*(\mathbf{y})$ fulfills $\tilde{\mathbf{V}}_r \tilde{\mathbf{\Lambda}}_r \tilde{\mathbf{V}}_r^\top \approx \mathbf{X}_*$. Here, by $\tilde{\mathbf{V}}_r \in \mathbb{R}^{d \times r}$ we denote a matrix which
 137 contains the first r columns of $\tilde{\mathbf{V}}$ and by $\tilde{\mathbf{\Lambda}}_r$ we denote a diagonal matrix which contains the largest r
 138 eigenvalues of $\mathcal{A}^*(\mathbf{y})$ in decreasing order. Motivated by this observation, the spectral initialization
 139 \mathbf{U}_0 is defined as

$$\mathbf{U}_0 := \tilde{\mathbf{V}}_r \tilde{\mathbf{\Lambda}}_r^{1/2}.$$

140 Here, the entries of the diagonal matrix $\tilde{\mathbf{\Lambda}}_r^{1/2}$ are given by $\sqrt{|\lambda_i(\mathcal{A}^*(\mathbf{y}))|}$. As we will see, all entries
 141 of $\tilde{\mathbf{\Lambda}}_r$ are positive with high probability.

142 After having computed the initialization \mathbf{U}_0 , we use \mathbf{U}_0 as a starting point of the gradient descent
 143 scheme in the second stage, which is defined as follows

$$\mathbf{U}_{t+1} := \mathbf{U}_t - \mu \nabla \mathcal{L}(\mathbf{U}_t) \quad \text{for } t = 0, 1, \dots,$$

144 where $\mu > 0$ denotes the step size. A direct computation shows that

$$\begin{aligned} \mathbf{U}_{t+1} &= \mathbf{U}_t + \mu [(\mathcal{A}^* \mathcal{A})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_t \\ &= \mathbf{U}_t + \frac{\mu}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \rangle \mathbf{A}_i \mathbf{U}_t. \end{aligned} \quad (4)$$

All steps of the two-stage approach are summarized below in Algorithm 1.1.

Algorithm 1 Two-Stage Approach for Low-Rank Matrix Recovery

Input: Measurement operator $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$, observations $\mathbf{y} \in \mathbb{R}^m$, step size $\mu > 0$

Stage 1 (Spectral Initialization): Compute the truncated eigendecomposition $\tilde{\mathbf{V}}_r \tilde{\mathbf{\Lambda}}_r \tilde{\mathbf{V}}_r^\top$ of the
 data matrix $\mathbf{D} := \mathcal{A}^*(\mathbf{y}) = \frac{1}{\sqrt{m}} \sum_{i=1}^m y_i \mathbf{A}_i$. Here, $\tilde{\mathbf{\Lambda}}_r \in \mathbb{R}^{d \times d}$ is the diagonal matrix which contains
 the r largest eigenvalues of the data matrix \mathbf{D} (in absolute value). The columns of $\tilde{\mathbf{\Lambda}}_r \in \mathbb{R}^{d \times r}$
 contain the corresponding eigenvectors. Define the initialization $\mathbf{U}_0 \in \mathbb{R}^{d \times r}$ via

$$\mathbf{U}_0 := \tilde{\mathbf{V}}_r \tilde{\mathbf{\Lambda}}_r^{1/2}.$$

Stage 2 (Gradient descent):

for $t = 0, 1, 2, \dots$ **do**

$$\mathbf{U}_{t+1} := \mathbf{U}_t - \mu \nabla \mathcal{L}(\mathbf{U}_t)$$

end for

145

146 **1.2. Main result**

147 To formulate our main result, we need to introduce the condition number of \mathbf{X}_* , which is defined
 148 as

$$\kappa := \frac{\|\mathbf{X}_*\|}{\sigma_{\min}(\mathbf{X}_*)}.$$

149 Here, $\sigma_{\min}(\mathbf{X}_*)$ denotes the smallest non-zero singular value of \mathbf{X}_* .

150 Next, let $\mathbf{U}_* \in \mathbb{R}^{d \times r}$ be a matrix such that $\mathbf{X}_* = \mathbf{U}_* \mathbf{U}_*^\top$. The matrix \mathbf{U}_* is uniquely defined only up
 151 to an orthogonal transformation $\mathbf{R} \in \mathbb{R}^{r \times r}$, which is why we can only expect to be able to reconstruct
 152 \mathbf{U}_* up to this ambiguity. To account for this, we will introduce the error metric

$$\text{dist}(\mathbf{U}_t, \mathbf{U}_*) := \min_{\mathbf{R} \in \mathbb{R}^{r \times r}, \mathbf{R}^\top \mathbf{R} = \text{Id}_r} \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_*\|_F. \quad (5)$$

153 With this notation in place, we can state the main result of this paper.

154 **Theorem 1.2.** *Let $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$ be a linear measurement operator as in Definition 1.1 with Gaussian
 155 measurement matrices. Moreover, let $\mathbf{X}_* \in \mathcal{S}^d$ be a positive semidefinite matrix of rank r . Given observations
 156 $\mathbf{y} = \mathcal{A}(\mathbf{X}_*) \in \mathbb{R}^m$, let $\mathbf{U}_0, \mathbf{U}_1, \mathbf{U}_2, \dots$ be the sequence of gradient descent iterates which are obtained via
 157 the two-stage approach described in Algorithm 1. Assume that the number of observations m satisfies*

$$m \geq Crd\kappa^2,$$

158 and that the step size $\mu > 0$ satisfies

$$\frac{32}{6^d \sigma_{\min}(\mathbf{X}_*)} \log(16r) \leq \mu \leq \frac{c_1}{\kappa \|\mathbf{X}_*\|}. \quad (6)$$

159 Then, with probability at least $1 - 7 \exp(-d)$, it holds for all iterations $t \geq 0$ that

$$\text{dist}^2(\mathbf{U}_t, \mathbf{U}_*) \leq c_2 r (1 - c_3 \mu \sigma_{\min}(\mathbf{X}_*))^t \sigma_{\min}(\mathbf{X}_*).$$

160 Here, $C, c_1, c_2, c_3 > 0$ denote absolute constants.

161 **Remark 1.3.** *The lower bound in assumption (6) is rather mild since the left-hand side in this inequality
 162 converges to 0 exponentially as the dimension d increases. If the dimension d is larger than an absolute
 163 constant, then condition (6) can always be satisfied for some step size μ .*

164 Theorem 1.2 shows that factorized gradient descent with spectral initialization converges to the
 165 ground truth with a linear rate as soon as the number of samples scales at least with $rd\kappa^2$. In par-
 166 ticular, the bound on the sample complexity is linear in the rank r . This improves over previous
 167 results in the matrix sensing literature, which have a sample complexity of order at least $r^2 d \kappa^2$, see,
 168 e.g., [10] or [11]. In particular, the sample complexity in Theorem 1.2 is optimal with respect to the
 169 rank r and dimension d . To the best of our knowledge, this is the first result in the literature which
 170 achieves this optimal dependence in the rank for the non-convex low-rank matrix recovery.

171 Compared to approaches based on nuclear norm or trace minimization, which only need $\Omega(rd)$
 172 samples in the matrix sensing scenario, our result is still suboptimal by a factor of κ^2 . However,
 173 all previous results in the literature on non-convex low-rank matrix recovery based on factorized
 174 gradient descent require having at least this quadratic dependence on the condition number. It
 175 remains an interesting open problem whether the dependence of the sample complexity on the
 176 condition number is necessary or an artifact of the proof.

177 Our main result implies that $\text{dist}(\mathbf{U}_t, \mathbf{U}_*) \leq \varepsilon$ after $O\left(\frac{\log(r/(\varepsilon \sigma_{\min}(\mathbf{X}_*)))}{\mu \sigma_{\min}(\mathbf{X}_*)}\right)$ iterations. Thus, if
 178 we choose the largest possible step size $\mu \asymp 1/(\kappa \|\mathbf{X}_*\|)$ we obtain that we reach ε -accuracy af-
 179 ter $O(\kappa^2 \log(r/(\varepsilon \sigma_{\min}(\mathbf{X}_*))))$ iterations. Previous work [10] allows for a larger step size $\mu \lesssim$
 180 $1/(\kappa \|\mathbf{X}_*\|)$ which yields that one can reach ε -accuracy after $O(\kappa \log(r/(\varepsilon \sigma_{\min}(\mathbf{X}_*))))$ iterations,
 181 whereas Theorem 1.2 requires $\mu \lesssim 1/(\kappa \|\mathbf{X}_*\|)$. It remains an open problem whether this additional
 182 condition number in the step size bound can be removed.

183 **Remark 1.4** (Landscape Analysis). Several works [44–47] have shown that if $m \gtrsim rd$, then the loss
 184 landscape of the objective function \mathcal{L} in (3) is benign in the sense that \mathcal{L} has no spurious local minima and all
 185 saddle points have at least one direction of strictly negative curvature. It has been established that in such a
 186 scenario gradient descent starting from random initialization will converge to the ground truth [48]. However,
 187 these results do not imply any guarantees on the convergence rate or on the computational complexity. In fact,
 188 there exist examples [49] where gradient descent may take exponential time to escape saddle points. For this
 189 reason, the results mentioned above are not directly comparable to our results.

190 2. Preliminaries

191 In the following, we will discuss several technical preliminaries, which are needed in our proof.

192 2.1. The Restricted Isometry Property

193 We first recall the Restricted Isometry Property (RIP).

194 **Definition 2.1** (Restricted Isometry Property). The linear measurement operator $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$ satisfies
 195 the Restricted Isometry Property (RIP), of rank r with RIP-constant $\delta_r > 0$, if it holds for all symmetric
 196 matrices $\mathbf{Z} \in \mathbb{R}^{d \times d}$ of rank at most r that

$$(1 - \delta_r) \|\mathbf{Z}\|_F^2 \leq \|\mathcal{A}(\mathbf{Z})\|_2^2 \leq (1 + \delta_r) \|\mathbf{Z}\|_F^2.$$

197 In previous works, it was shown that as soon as the measurement operator \mathcal{A} has the RIP, then
 198 convex approaches based on nuclear norm minimization as well as non-convex approaches are able
 199 to recover the ground truth matrix, see, e.g., [6, 10].

200 It is well known that as soon as the number of samples m satisfies $m \gtrsim rd$ then the measurement
 201 operator \mathcal{A} has the RIP of order r with high probability. This fact is stated in the following lemma.

202 **Lemma 2.2.** Let $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$ be a Gaussian measurement operator as described in Section 1.1. Then the
 203 RIP constant δ_r satisfies $\delta_r \leq \delta \leq 1$ with probability $1 - \varepsilon$ when

$$m \geq C\delta^{-2}(rd + \log(2\varepsilon^{-1})),$$

204 where $C > 0$ is a universal constant. In particular, we have with probability at least $1 - \exp(-d)$, $m \geq$
 205 $C\delta^{-2}rd$.

206 This lemma differs from similar lemmas in the literature (see, e.g., [50]) by specifying how m de-
 207 pends on the RIP-constant δ . A proof of this lemma is provided in Appendix D.1 together with a
 208 more detailed discussion of how this lemma relates to previous work.

209 **Remark 2.3.** The works mentioned in Remark 1.4 have shown that the RIP implies that the optimization
 210 landscape of \mathcal{L} is benign (in the sense of Remark 1.4). Moreover, previous work such as [10] or [11], which
 211 analyzed gradient descent with spectral initialization similar to the paper at hand, relied on their analysis
 212 of gradient descent exclusively on the RIP property of the measurement operator \mathcal{A} . As we will explain in
 213 Section 3, purely relying on the RIP will not suffice to establish Theorem 1.2. For this reason, in addition to
 214 the RIP, we will use the orthogonal invariance of the Gaussian measurement operator \mathcal{A} .

215 The RIP has several important consequences, which we will need throughout our proof. We recall
 216 them in the following lemma.

217 **Lemma 2.4.** Let $\mathcal{A} : \mathcal{S}^d \rightarrow \mathbb{R}^m$ be a linear measurement operator on the set of symmetric matrices as defined
 218 above. Denote by δ_r the RIP constant of the operator \mathcal{A} of order r . Then the following statements hold.

219 1. Let $\mathbf{V} \in \mathbb{R}^{d \times r'}$ be any matrix with orthonormal columns, i.e., $\mathbf{V}^\top \mathbf{V} = \mathbf{Id}$. Then it holds for any
 220 symmetric matrix $\mathbf{Z} \in \mathbb{R}^{d \times d}$ of rank at most r that

$$\|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})\mathbf{V}\|_F \leq \delta_{r+2r'} \|\mathbf{Z}\|_F. \quad (7)$$

221 In particular, it holds that

$$\|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})\| \leq \delta_{r+2} \|\mathbf{Z}\|_F. \quad (8)$$

222 2. Let $\mathbf{w} \in \mathbb{R}^d$ such that $\|\mathbf{w}\|_2 = 1$. Define the orthogonal projection operators

$$\begin{aligned} \mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{Z}) &:= \langle \mathbf{w}\mathbf{w}^\top, \mathbf{Z} \rangle \mathbf{w}\mathbf{w}^\top, \\ \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z}) &:= \mathbf{Z} - \langle \mathbf{w}\mathbf{w}^\top, \mathbf{Z} \rangle \mathbf{w}\mathbf{w}^\top. \end{aligned} \quad (9)$$

223 Then it holds for any symmetric matrix $\mathbf{Z} \in \mathbb{R}^{d \times d}$ of rank at most r that

$$|\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z})) \rangle| \leq \delta_{r+2} \|\mathbf{Z}\|_F. \quad (10)$$

224 Some variants of these inequalities appeared in the literature already before; see, e.g., [23]. For
225 completeness, we decided to include a proof in Appendix D.2.

226 **Remark 2.5.** To keep the notation more concise, we will sometimes drop the subscript and just use the notation
227 δ for the RIP constant. For all results below, the choices of δ satisfy $\delta \leq \delta_{6r}$ due to the monotonicity of the
228 RIP constant with respect to the rank.

229 2.2. Perturbation bounds for eigenspaces

230 The Davis-Kahan $\sin \theta$ -theorem [51] states that the eigenspaces of a symmetric matrix are stable
231 under perturbations of that matrix. Among others, we will need this result in order to show that
232 the spectral initialization recovers the eigenspace of the ground truth matrix sufficiently well. We
233 also will need it in order to show that $\mathbf{U}_{0, \mathbf{w}}$ is sufficiently close to \mathbf{U}_0 .

234 To state this theorem, recall that for a symmetric matrix $\mathbf{Z} \in \mathbb{R}^{n \times n}$ with eigendecomposition $\mathbf{Z} =$
235 $\mathbf{U}_Z \mathbf{\Lambda}_Z \mathbf{U}_Z^\top$ the matrix $\mathbf{U}_{Z, r} \in \mathbb{R}^{n \times r}$ consists of the first r columns of \mathbf{U}_Z and the matrix $\mathbf{U}_{Z, r, \perp} \in$
236 $\mathbb{R}^{n \times (n-r)}$ consists of the remaining $n - r$ columns. Moreover, recall that the eigenvalues of \mathbf{Z} are
237 ordered such that their magnitude is decreasing, i.e., $|\lambda_1(\mathbf{Z})| \geq |\lambda_2(\mathbf{Z})| \geq \dots \geq |\lambda_n(\mathbf{Z})|$.

Lemma 2.6 (Davis-Kahan inequality, Corollary 2.8 in [52]). Set $\|\cdot\| = \|\cdot\|$ or $\|\cdot\| = \|\cdot\|_F$. Let
 $\mathbf{Z}_1 \in \mathbb{R}^{d \times d}$ and $\mathbf{Z}_2 \in \mathbb{R}^{d \times d}$ be two symmetric matrices, such that the eigenvalues of \mathbf{Z}_1 satisfy $|\lambda_r(\mathbf{Z}_1)| >$
 $|\lambda_{r+1}(\mathbf{Z}_1)|$ for an integer $1 \leq r < d$. Let the eigendecompositions of \mathbf{Z}_1 and \mathbf{Z}_2 be given by $\mathbf{Z}_1 = \mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^\top$,
respectively $\mathbf{Z}_2 = \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^\top$. Then, if the assumption

$$\|\mathbf{Z}_1 - \mathbf{Z}_2\| \leq (1 - 1/\sqrt{2}) (|\lambda_r(\mathbf{Z}_1)| - |\lambda_{r+1}(\mathbf{Z}_1)|)$$

238 is fulfilled, it holds that

$$\|\|\mathbf{U}_{2, r, \perp}^\top \mathbf{U}_{1, r}\|\| \leq \frac{\sqrt{2} \|\|\mathbf{Z}_1 - \mathbf{Z}_2\|\| \|\mathbf{U}_{1, r}\|\|}{|\lambda_r(\mathbf{Z}_1)| - |\lambda_{r+1}(\mathbf{Z}_1)|}.$$

239 3. Outline of the proof

240 3.1. A fundamental barrier in previous work

241 Before we give an outline of our proof approach, we want to explain why in previous work the
242 additional r -factor appeared in the sample complexity. As Lemma 4.1 below shows, it holds for the
243 spectral initialization \mathbf{U}_0 with high probability that

$$\|\mathbf{X}_\star - \mathbf{U}_0 \mathbf{U}_0^\top\| \leq C \kappa \sigma_{\min}(\mathbf{X}_\star) \sqrt{\frac{rd}{m}}.$$

244 In particular, for $m \gg \kappa^2 rd$ we have that

$$\|\mathbf{X}_\star - \mathbf{U}_0 \mathbf{U}_0^\top\| \ll \sigma_{\min}(\mathbf{X}_\star).$$

245 Thus, the spectral initialization ensures that the initialization \mathbf{U}_0 is in a neighborhood of the ground
246 truth. We aim to establish that within this neighborhood, gradient descent converges with a linear
247 rate. To show this, we note first that the gradient of our objective function \mathcal{L} depends on the random

248 matrices $(\mathbf{A}_i)_{i=1}^m$. To deal with this, a common technique that has been used in previous works is to
 249 decompose the gradient of the objective function \mathcal{L} into a sum of two terms:

$$\nabla \mathcal{L}(\mathbf{U}) = \mathbb{E}_{(\mathbf{A}_i)_{i=1}^m} [\nabla \mathcal{L}(\mathbf{U})] + [\nabla \mathcal{L}(\mathbf{U}) - \mathbb{E}_{(\mathbf{A}_i)_{i=1}^m} [\nabla \mathcal{L}(\mathbf{U})]].$$

250 The first term is the gradient of the population risk, i.e., the objective function one obtains in the
 251 limit case that the sample size m goes to infinity. The second term can be interpreted as a pertur-
 252 bation term that measures the deviation of the gradient of the empirical risk from the gradient of
 253 the population risk. In particular, this term converges to zero as the sample size m increases. For
 254 this reason, a major task in our proof is to show that the second summand is small with respect to
 255 a suitable norm as soon as the sample size m is sufficiently large. A direct computation shows that

$$\begin{aligned} \nabla \mathcal{L}(\mathbf{U}) - \mathbb{E}_{(\mathbf{A}_i)_{i=1}^m} [\nabla \mathcal{L}(\mathbf{U})] &= [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U} \mathbf{U}^\top - \mathbf{X}_*)] \mathbf{U} \\ &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_* \rangle \mathbf{A}_i - (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*). \end{aligned}$$

256 To deal with this deviation term, in previous works, bounds of the type

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \ll \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \quad (11)$$

257 needed to be established. A major challenge in establishing such bounds is that the gradient descent
 258 iterates $(\mathbf{U}_t)_t$ depend on the measurement matrices $(\mathbf{A}_i)_{i=1}^m$ in an intricate way. For this reason,
 259 standard matrix concentration inequalities are not directly applicable. To circumvent this issue,
 260 previous work establishes *uniform* bounds for the quantity

$$\sup_{\mathbf{Z} \in \mathcal{T}_{2r}} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\|$$

261 where

$$\mathcal{T}_r := \{\mathbf{Z} \in \mathbb{R}^{d \times d} : \mathbf{Z} = \mathbf{Z}^\top, \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\| \leq 1\},$$

262 denotes the collection of matrices with rank at most r and bounded operator norm. Indeed, such a
 263 bound can be directly derived from the Restricted Isometry Property. Namely, when \mathcal{A} has the RIP
 264 of order $2r + 2$ with constant δ_{2r+2} then Lemma 2.4 implies that

$$\sup_{\mathbf{Z} \in \mathcal{T}_{2r}} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\| \leq \delta_{2r+2} \sup_{\mathbf{Z} \in \mathcal{T}_{2r}} \|\mathbf{Z}\|_F \leq \delta_{2r+2} \sqrt{2r},$$

265 where in the second inequality, we used that the matrix \mathbf{Z} has rank at most $2r$ and that $\|\mathbf{Z}\| = 1$.
 266 Thus, it follows from Lemma 2.2 that whenever $m \gg rd$ that with high probability we have that

$$\sup_{\mathbf{Z} \in \mathcal{T}_{2r}} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\| \lesssim \sqrt{\frac{r^2 d}{m}}. \quad (12)$$

267 This shows that if we want to deduce inequality (11) from the uniform bound (12) we must assume
 268 that $m \gg r^2 d$. Indeed, several works, e.g., [22, 23, 53], relied precisely on this bound.

269 This leads to the question of whether the bound (12) can be sharpened. For example, in [53, p. 9],
 270 it was conjectured that using more refined techniques from empirical process theory, one may be
 271 able to refine (12). However, as the following result shows, inequality (12) is tight up to absolute
 272 numerical constants and thus cannot be improved further.

273 **Theorem 3.1.** *Let $(\mathbf{A}_i)_{i \in [m]}$ be independent $d \times d$ symmetric random matrices, where each \mathbf{A}_i has inde-
 274 pendent entries with distribution $\mathcal{N}(0, 1)$ on the diagonal and $\mathcal{N}(0, 1/2)$ on the off-diagonal entries. As-
 275 sume $d \geq 6$, $m \geq C_0$ for some universal constant $C_0 > 0$, and $r \leq \frac{d}{16}$. Then, with probability at least
 276 $1 - 2 \exp(-\frac{m}{32}) - 2 \exp(-\frac{d}{32})$, it holds that*

$$\sup_{\mathbf{Z} \in \mathcal{T}_r} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\| \geq \frac{1}{16} \sqrt{\frac{r^2 d}{m}}.$$

277 Theorem 3.1 shows that we will need to use different proof techniques to establish a bound similar
 278 to (11). In particular, we cannot rely on uniform concentration inequalities. These novel techniques
 279 will be introduced in Section 3.2 below. Before that, we want to prove Theorem 3.1.

280 *Proof.* First, we note that

$$\sup_{\mathbf{Z} \in \mathcal{T}_r} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\| = \sup_{\mathbf{Z} \in \mathcal{T}_r} \left\| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i - \mathbf{Z} \right\| = \sup_{\|\mathbf{u}\|=1} \sup_{\mathbf{Z} \in \mathcal{T}_r} \left| \left\langle \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i - \mathbf{Z}, \mathbf{u} \mathbf{u}^\top \right\rangle \right|.$$

281 Now for any fixed $\mathbf{u} \in \mathbb{R}^d$ with $\|\mathbf{u}\|_2 = 1$, define

$$\mathcal{T}_{\mathbf{u}} := \{ \mathbf{Z} \in \mathbb{R}^{d \times d} : \mathbf{Z} = \mathbf{Z}^\top, \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\| \leq 1, \mathbf{Z} \mathbf{u} = \mathbf{0} \},$$

282 i.e., the set consisting of matrices in \mathcal{T}_r , whose row space is orthogonal to \mathbf{u} . It follows that

$$\begin{aligned} \sup_{\mathbf{Z} \in \mathcal{T}_r} \left\| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i - \mathbf{Z} \right\| &\geq \sup_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}} \left\langle \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i - \mathbf{Z}, \mathbf{u} \mathbf{u}^\top \right\rangle \\ &= \sup_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}} \left\langle \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \right\rangle \\ &= \sup_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}} \frac{1}{m} \sum_{i=1}^m \langle \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle \mathbf{A}_i, \mathbf{Z} \rangle. \end{aligned}$$

283 Now note that $\langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle$ is independent of $(\langle \mathbf{A}_i, \mathbf{Z} \rangle)_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}}$. Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a matrix with the
 284 same distribution as \mathbf{A}_i and which is independent of $(\mathbf{A}_i)_{i=1}^m$. We claim that conditional on
 285 $\{\langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle\}_{i=1}^m$ we have that the following two random variables are equal in distribution:

$$\sup_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}} \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \stackrel{d}{=} \frac{1}{\sqrt{m}} \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle^2} \sup_{\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}} \langle \mathbf{A}, \mathbf{Z} \rangle. \quad (13)$$

286 To show (13), one can check that conditional on $\{\langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle\}_{i=1}^m$, the random variables on both sides
 287 of (13) are the supremum of Gaussian processes indexed by $\mathcal{T}_{\mathbf{u}}$ with the same covariance structure,
 288 so they have the same distribution.

289 In the following, we set

$$\mathbf{u} := (0, \dots, 0, 1)^\top \in \mathbb{R}^d. \quad (14)$$

290 It follows that

$$\sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u} \mathbf{u}^\top \rangle^2 = \sum_{i=1}^m (\mathbf{A}_i)_{d,d}^2.$$

291 By Lipschitz concentration for Gaussian random variables [54, Theorem 5.6], we obtain

$$\mathbb{P} \left(\left| \sqrt{\sum_{i=1}^m (\mathbf{A}_i)_{d,d}^2} - \mathbb{E} \sqrt{\sum_{i=1}^m (\mathbf{A}_i)_{d,d}^2} \right| \geq \sqrt{m}/4 \right) \leq 2 \exp(-m/32).$$

292 This shows that with probability at least $1 - 2 \exp(-m/32)$,

$$\sqrt{\sum_{i=1}^m (\mathbf{A}_i)_{d,d}^2} \geq \mathbb{E} \sqrt{\sum_{i=1}^m (\mathbf{A}_i)_{d,d}^2} - \frac{\sqrt{m}}{4} \geq \sqrt{m}/2 \quad (15)$$

293 for sufficiently large m , where we have used that the expectation of chi-distribution with parameter
 294 m has asymptotic value $\sqrt{m - \frac{1}{2}}$ (see, e.g., [55]). In addition, with \mathbf{u} given in (14), all entries in
 295 the d -th row and d -th column of the matrix $\mathbf{Z} \in \mathcal{T}_{\mathbf{u}}$ are equal to zero. Let $\tilde{\mathbf{A}} \in \mathbb{R}^{(d-1) \times (d-1)}$ be the

296 submatrix \mathbf{A} where the last row and column of \mathbf{A} are removed, and define $\tilde{\mathbf{Z}}$ in the same way. Then
 297 we have

$$\sup_{\mathbf{Z} \in \mathcal{T}_r} \langle \mathbf{A}, \mathbf{Z} \rangle = \sup_{\|\tilde{\mathbf{Z}}\| \leq 1, \tilde{\mathbf{Z}} = \tilde{\mathbf{Z}}^\top, \text{rank}(\tilde{\mathbf{Z}}) \leq r} \langle \tilde{\mathbf{A}}, \tilde{\mathbf{Z}} \rangle = \sum_{i=1}^r \sigma_i(\tilde{\mathbf{A}}).$$

298 Our goal is to bound the sum of singular values on the right-hand side from below. For that, we
 299 define the matrix

$$\hat{\mathbf{A}} := \begin{pmatrix} \mathbf{0}_{\lceil (d-1)/2 \rceil - 1 \times r} & \mathbf{0}_{\lceil (d-1)/2 \rceil \times (d-r)} \\ \tilde{\mathbf{A}}_{\lceil (d-1)/2 \rceil : (d-1), 1:r} & \mathbf{0}_{(d-1 - \lceil (d-1)/2 \rceil) \times (d-r)} \end{pmatrix} \in \mathbb{R}^{(d-1) \times (d-1)}.$$

300 Here, $\tilde{\mathbf{A}}_{\lceil (d-1)/2 \rceil : (d-1), 1:r}$ denotes the submatrix of \mathbf{A} obtained by restricting \mathbf{A} to the $\lceil (d-1)/2 \rceil$ -th
 301 to $(d-1)$ -th rows and the first r columns. By $\mathbf{0}_{a \times b}$ we denote the zero matrix of size a times b . To
 302 relate the singular values of $\tilde{\mathbf{A}}$ with the singular values of $\hat{\mathbf{A}}$, we will use the following lemma.

303 **Lemma 3.2** (Corollary 3.1.3 in [56]). *Let $\mathbf{A} \in \mathbb{R}^{(d-1) \times (d-1)}$ and let $\mathbf{B} \in \mathbb{R}^{(d-1) \times (d-1)}$ be a matrix which
 304 is obtained from the matrix \mathbf{A} by setting the entries of one row or one column to zero. Then it holds that
 305 $\sigma_i(\mathbf{B}) \leq \sigma_i(\mathbf{A})$ for all $i = 1, \dots, d-1$.*

306 By repeatedly applying Lemma 3.2, we find

$$\sum_{i=1}^r \sigma_i(\hat{\mathbf{A}}) \leq \sum_{i=1}^r \sigma_i(\tilde{\mathbf{A}}).$$

307 On the other hand, we can identify the r largest singular values of $\hat{\mathbf{A}}$ with the singular
 308 values of a Gaussian matrix of size $\lfloor \frac{d-1}{2} \rfloor \times r$. By standard concentration inequalities for the singular
 309 values of Gaussian matrices, see, e.g., [57, Corollary 5.35], we find that with probability at least
 310 $1 - 2 \exp(-t^2/2)$,

$$\sigma_r(\hat{\mathbf{A}}) \geq \sqrt{\left\lfloor \frac{d-1}{2} \right\rfloor} - \sqrt{r} - t.$$

311 Taking $t = \frac{\sqrt{d}}{8}$, and using the assumption that $r \leq \frac{d}{16}$, we find for $d \geq 6$,

$$\sum_{i=1}^r \sigma_i(\tilde{\mathbf{A}}) \geq \frac{r\sqrt{d}}{8} \tag{16}$$

312 with probability at least $1 - 2 \exp(-d/32)$. Combining (16) and (15) finishes the proof. \square

313 Note that the key idea in this proof was to fix a vector $\mathbf{u} \in \mathbb{R}^d$ and to pick a matrix $\mathbf{Z} \in \mathcal{T}_r$ based on
 314 eigenvectors corresponding to the largest eigenvalues (of a submatrix) of

$$\mathbf{A} = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top \rangle \mathbf{A}_i.$$

315 By design, this implies that the matrix \mathbf{Z} was chosen in a way which strongly depends on
 316 $(\langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top \rangle)_{i=1}^m$. This observation leads to the key idea in our proof. Namely, we will show that
 317 our gradient descent iterates \mathbf{U}_t depend, in a suitable sense, only weakly $(\langle \mathbf{A}_i, \mathbf{u}\mathbf{u}^\top \rangle)_{i=1}^m$ for fixed
 318 $\mathbf{u} \in \mathbb{R}^d$. This will allow us to prove stronger upper bounds for the term $\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|$
 319 than what can be achieved using uniform concentration inequalities.

320 3.2. Virtual sequences

321 As explained at the end of Section 3.1, we aim to establish that the gradient descent iterates $(\mathbf{U}_t)_t$
 322 depend only weakly on $(\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle)_{i=1}^m$ in a suitable sense. For this aim, we will use so-called *virtual*
 323 *sequences* $(\mathbf{U}_{t,\mathbf{w}})_{t \in \mathbb{N}} \subset \mathcal{S}^d$. The central idea is to introduce for $\mathbf{w} \in \mathcal{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ a
 324 sequence with the following two properties.

- 325 1. The sequence $(\mathbf{U}_{t,\mathbf{w}})_{t \in \mathbb{N}}$ is stochastically independent of $(\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle)_{i=1}^m$.
326 2. The sequence $(\mathbf{U}_{t,\mathbf{w}})_{t \in \mathbb{N}}$ stays sufficiently close to the sequence $(\mathbf{U}_t)_{t \in \mathbb{N}}$. More precisely, we
327 require that $\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F$ stays sufficiently small.

328 The sequences $(\mathbf{U}_{t,\mathbf{w}})_{t \in \mathbb{N}}$ are called *virtual* since they are introduced solely for proof purposes.

329 **Remark 3.3** (Related work). *In the context of non-convex optimization, the use of virtual sequences has
330 been pioneered in the influential works [18] and [33]. In these works, leave-one-out sequences, which can be
331 seen as a special case of virtual sequences, were introduced to show that the gradient descent iterates depend
332 only weakly on the individual samples or measurements. These works lead to a number of follow-up works.
333 For example, several works used virtual sequences to establish convergence from random initialization for
334 gradient descent in phase retrieval [58] or for alternating minimization in rank-one matrix sensing [59]. In
335 [27], leave-one-out sequences were used to establish that in overparameterized matrix completion gradient
336 descent with small random initialization converges to the ground truth. Similar to the paper at hand, the
337 virtual sequence argument was combined with an ε -net argument. However, the technical details are arguably
338 quite different.*

339 Before defining the virtual sequences we recall the notion of an ε -net.

340 **Definition 3.4** (ε -net). *Let $A \subset \mathbb{R}^d$. A subset $B \subset A$ is called ε -net of A if for every $\mathbf{x} \in A$ there is a point
341 $\mathbf{x}_0 \in B$ such that $\|\mathbf{x} - \mathbf{x}_0\|_2 \leq \varepsilon$.*

342 It is well-known that for $S^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$ there exists an ε -net $\mathcal{N}_\varepsilon \subset S^{d-1}$ with cardinal-
343 ity $|\mathcal{N}_\varepsilon| \leq (3/\varepsilon)^d$ [60]. In the remainder of this paper, we will assume that \mathcal{N}_ε is a fixed ε -net of S^{d-1}
344 with $\varepsilon = 1/2$ such that $|\mathcal{N}_\varepsilon| \leq 6^d$. We will define one virtual sequence $(\mathbf{U}_{t,\mathbf{w}})_t$ for each $\mathbf{w} \in \mathcal{N}_\varepsilon$.

345 Recall from equation (9) that for $\mathbf{w} \in \mathcal{N}_\varepsilon$ the orthogonal projection operators $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}$ and $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}$
346 were defined for $\mathbf{Z} \in S^d$ via

$$\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{Z}) = \langle \mathbf{w}\mathbf{w}^\top, \mathbf{Z} \rangle \mathbf{w}\mathbf{w}^\top, \quad \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z}) = \mathbf{Z} - \langle \mathbf{w}\mathbf{w}^\top, \mathbf{Z} \rangle \mathbf{w}\mathbf{w}^\top.$$

347 Next, for $\mathbf{w} \in \mathcal{N}_\varepsilon$ we define the modified measurement matrices via

$$\mathbf{A}_{i,\mathbf{w}} := \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{A}_i) = \mathbf{A}_i - \langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle \mathbf{w}\mathbf{w}^\top.$$

348 Thus, the matrix $\mathbf{A}_{i,\mathbf{w}}$ is obtained from the matrix \mathbf{A}_i by setting the generalized entry $\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle$
349 equal to 0. We observe that by definition the matrices $(\mathbf{A}_{i,\mathbf{w}})_{i=1}^m$ are stochastically independent of
350 $(\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle)_{i=1}^m$. We define the virtual measurement operator $\mathcal{A}_{\mathbf{w}} : S^d \rightarrow \mathbb{R}^{m+1}$ via

$$[\mathcal{A}_{\mathbf{w}}(\mathbf{Z})]_i := \frac{1}{\sqrt{m}} \langle \mathbf{A}_i, \mathbf{Z} \rangle$$

351 for $i \in [m]$ and

$$[\mathcal{A}_{\mathbf{w}}(\mathbf{Z})]_{m+1} := \langle \mathbf{w}\mathbf{w}^\top, \mathbf{Z} \rangle.$$

352 Again, we observe that by construction, the measurement operator $\mathcal{A}_{\mathbf{w}}$ is independent of
353 $(\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle)_{i=1}^m$. As a next step, analogously to the definition of the objective function \mathcal{L} , we can
354 define the modified objective function $\mathcal{L}_{\mathbf{w}} : S^d \rightarrow \mathbb{R}$ via

$$\mathcal{L}_{\mathbf{w}}(\mathbf{U}) := \frac{1}{4} \|\mathcal{A}_{\mathbf{w}}(\mathbf{X}_* - \mathbf{U}\mathbf{U}^\top)\|_2^2.$$

355 With these definitions in place, the virtual sequence $(\mathbf{U}_{t,\mathbf{w}})_t$ can be defined analogously to the orig-
356 inal sequence $(\mathbf{U}_t)_t$. Namely, to define the spectral initialization, we consider the eigendecomposi-
357 tion

$$(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*) =: \tilde{\mathbf{V}}_{\mathbf{w}} \tilde{\mathbf{\Lambda}}_{\mathbf{w}} \tilde{\mathbf{V}}_{\mathbf{w}}^\top. \quad (17)$$

358 Then, analogously as for the original spectral initialization \mathbf{U}_0 , the matrix $\mathbf{U}_{0,\mathbf{w}}$ is defined as

$$\mathbf{U}_{0,\mathbf{w}} =: \tilde{\mathbf{V}}_{r,\mathbf{w}} \tilde{\mathbf{\Lambda}}_{r,\mathbf{w}}^{1/2}. \quad (18)$$

359 Then the virtual sequence $\{\mathbf{U}_{t,\mathbf{w}}\}_{t \in \mathbb{N}}$ via

$$\mathbf{U}_{t+1,\mathbf{w}} := \mathbf{U}_{t,\mathbf{w}} - \mu \nabla \mathcal{L}_{\mathbf{w}}(\mathbf{U}_{t,\mathbf{w}}) = \mathbf{U}_{t,\mathbf{w}} + \mu [(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_{\star} - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^{\top})] \mathbf{U}_{t,\mathbf{w}}.$$

360 It follows directly from the definition of $(\mathbf{U}_{t,\mathbf{w}})_t$ that this sequence is stochastically independent of
 361 $(\langle \mathbf{A}_i, \mathbf{w} \mathbf{w}^{\top} \rangle)_{i=1}^m$. At the end of this section, we state and prove the following lemma, which is a
 362 direct consequence of the definition of $\mathcal{A}_{\mathbf{w}}$. This lemma will be useful in the convergence analysis
 363 where we establish that $\|\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^{\top}\|_F$ stays sufficiently small.

364 **Lemma 3.5.** *For any symmetric matrix $\mathbf{Z} \in \mathbb{R}^{d \times d}$ it holds that*

$$\begin{aligned} (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z})) &= \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z}), \\ (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) &= (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) - \langle \mathcal{A}(\mathbf{w} \mathbf{w}^{\top}), \mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) \rangle \mathbf{w} \mathbf{w}^{\top}. \end{aligned}$$

365 *Proof of Lemma 3.5.* To prove the first inequality we note first that it follows directly from the defini-
 366 tion of $\mathbf{A}_{i,\mathbf{w}}$ that $\langle \mathbf{A}_{i,\mathbf{w}}, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z}) \rangle = 0$. It follows that

$$\begin{aligned} (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z})) &= \frac{1}{\sqrt{m}} \sum_{i=1}^m [\mathcal{A}_{\mathbf{w}}(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z}))]_i \mathbf{A}_{i,\mathbf{w}} + (\mathcal{A}_{\mathbf{w}}(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z}))_{m+1} \mathbf{w} \mathbf{w}^{\top} \\ &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_{i,\mathbf{w}}, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}}(\mathbf{Z}) \rangle \mathbf{A}_{i,\mathbf{w}} + \langle \mathbf{w} \mathbf{w}^{\top}, \mathbf{Z} \rangle \mathbf{w} \mathbf{w}^{\top} \\ &= \langle \mathbf{w} \mathbf{w}^{\top}, \mathbf{Z} \rangle \mathbf{w} \mathbf{w}^{\top}. \end{aligned}$$

367 This proves the first equation. In order to prove the second equation, we note that

$$\begin{aligned} (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_{i,\mathbf{w}}, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \mathbf{A}_{i,\mathbf{w}} + \langle \mathbf{w} \mathbf{w}^{\top}, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \mathbf{w} \mathbf{w}^{\top} \\ &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_{i,\mathbf{w}}, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \mathbf{A}_{i,\mathbf{w}} \\ &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \mathbf{A}_{i,\mathbf{w}} \\ &= \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \mathbf{A}_i - \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z}) \rangle \langle \mathbf{w} \mathbf{w}^{\top}, \mathbf{A}_i \rangle \mathbf{w} \mathbf{w}^{\top} \\ &= (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) - \langle \mathcal{A}(\mathbf{w} \mathbf{w}^{\top}), \mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{Z})) \rangle \mathbf{w} \mathbf{w}^{\top}. \end{aligned}$$

368 This proves the second equation. □

369 3.3. Upper bounds for the spectral norm of the deviation term

370 Recall that by construction, it holds for any $\mathbf{w} \in \mathcal{N}_{\varepsilon}$ that the sequence $(\mathbf{U}_{t,\mathbf{w}})_{t=0,1,\dots,T}$ is independent
 371 of $(\langle \mathbf{w} \mathbf{w}^{\top}, \mathbf{A}_i \rangle)_{i=1}^m$. This property allows us to establish the following key lemma which we will use
 372 several times throughout our proof.

373 **Lemma 3.6.** *Let $\mathcal{N}_{\varepsilon}$ be the ε -net with $\varepsilon = 1/2$ introduced in Section 3.2 which we used to construct the
 374 virtual sequences $(\mathbf{U}_{t,\mathbf{w}})_t$. Assume that for the cardinality of $\mathcal{N}_{\varepsilon}$, we have that $|\mathcal{N}_{\varepsilon}| \leq 6^d$. Moreover, let
 375 $T \in \mathbb{N}$ such that $2T \leq 6^d$. Then, with probability at least $1 - 2 \exp(-10d)$, it holds for all $\mathbf{w} \in \mathcal{N}_{\varepsilon}$ and all
 376 $1 \leq t \leq T$ that*

$$|\langle \mathbf{w} \mathbf{w}^{\top}, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{X}_{\star} - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^{\top})) \rangle| \leq 4 \sqrt{\frac{d}{m}} \|\mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^{\top}, \perp}(\mathbf{X}_{\star} - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^{\top}))\|_2.$$

377 *Proof.* We introduce the shorthand

$$\Delta_{t,\mathbf{w}} := \mathbf{X}_{\star} - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^{\top}.$$

378 Due to the definition of $\mathbf{A}_{i,\mathbf{w}}$ and due to the rotation invariance of the Gaussian distribution,
 379 $\{\mathbf{A}_{i,\mathbf{w}}\}_{i=1}^m$ and $\{\langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle\}_{i=1}^m$ are independent. Moreover, note that by construction $\Delta_{t,\mathbf{w}}$
 380 is independent of $\{\langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle\}_{i=1}^m$. Thus, it follows that $\{\langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle\}_{i=1}^m$ is independent of
 381 $\{\langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}) \rangle\}_{i=1}^m$. Moreover, the vector $(\langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle)_{i=1}^m$ has i.i.d. entries with distribution
 382 $\mathcal{N}(0, 1)$. Thus, we have for all $x > 0$ with probability at least $1 - 2 \exp(-x^2/2)$ (see [60, Proposition
 383 2.1.2]) that

$$\begin{aligned} |\langle \mathbf{w}\mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| &= \left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}) \rangle \right| \\ &\leq \frac{x}{m} \sqrt{\sum_{i=1}^m \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}) \rangle^2} \\ &= \frac{x}{\sqrt{m}} \|\mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}))\|_2. \end{aligned} \quad (19)$$

384 Then, by applying inequality (19) with $x = C\sqrt{d}$ and by taking a union bound, it follows that with
 385 probability at least $1 - \xi$ (over the whole probability space), we have for all $\mathbf{w} \in \mathcal{N}_\varepsilon$ and all $t \in [T]$
 386 that

$$|\langle \mathbf{w}\mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| \leq \frac{C\sqrt{d}}{\sqrt{m}} \|\mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}))\|_2,$$

387 where

$$\xi \leq 2T|\mathcal{N}_\varepsilon| \exp(-C^2d) \leq 6^{2d} \exp(-C^2d) = \exp(2d \log(6) - C^2d).$$

388 The claim follows from choosing $C = 4$. □

389 Recall that our goal was to derive an upper bound for the expression $\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|$.
 390 The following lemma provides such a bound for $1 \leq t \leq T$. Here, $T \in \mathbb{N}$ is some fixed number of
 391 iterations, which will be specified later in the proof of our main result.

392 **Proposition 3.7.** *Let \mathcal{N}_ε be the ε -net from above with $\varepsilon = 1/2$ which we used to construct the virtual
 393 sequences $(\mathbf{U}_{t,\mathbf{w}})_{t=0,1,\dots,T}$. Assume that the conclusion of Lemma 3.6 holds. Moreover, assume that the linear
 394 measurement operator \mathcal{A} has the Restricted Isometry Property of order $2r + 2$ with constant $\delta = \delta_{2r+2} \leq 1$.
 395 Then it holds that for all $0 \leq t \leq T$,*

$$\begin{aligned} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| &\leq \left(16\sqrt{\frac{2rd}{m}} + 2\delta \right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\quad + 4 \left(\delta + 4\sqrt{\frac{d}{m}} \right) \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F. \end{aligned}$$

396 As already mentioned, in previous literature, the quantity $\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|$ was con-
 397 trolled via an upper bound of $\sup_{\mathbf{Z} \in \mathcal{T}_{2r}} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z})\|$, where \mathcal{T}_{2r} is a set of all rank- $2r$ matrices
 398 with bounded operator norm. This requires a uniform concentration bound for all matrices of rank
 399 at most $2r$ with bounded spectral norm. As we have seen in Theorem 3.1, this argument necessarily
 400 leads to a multiplicative factor of $\sqrt{r^2 d/m}$.

401 In contrast, Proposition 3.7 bounds $\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|$ by a sum of two terms. The first
 402 term can be controlled with sample complexity $m \gtrsim rd\kappa^2$ since we also have $\delta \lesssim \sqrt{rd/m}$, see
 403 Lemma 2.2. The second term is a uniform bound on the deviation of the “true” sequence from the
 404 “virtual” sequences. This term can be interpreted as a measure of how stable the sequence $(\mathbf{U}_t)_t$
 405 are under perturbation of the generalized entries $(\langle \mathbf{A}_i, \mathbf{w}\mathbf{w}^\top \rangle)_{i=1}^m$ of the symmetric measurement
 406 matrices.

407 *Proof of Proposition 3.7.* We use the shorthand notation

$$\begin{aligned}\Delta_t &:= \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top, \\ \Delta_{t,\mathbf{w}} &:= \mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top.\end{aligned}$$

408 Since \mathcal{N}_ε is an ε -net of S^{d-1} with $\varepsilon = 1/2$ we obtain that

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_t)\| \leq 2 \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_t) \rangle|, \quad (20)$$

409 (see, e.g. [60, Lemma 4.4.1]). Then, for every $\mathbf{w} \in \mathcal{N}_\varepsilon$ using the triangle inequality we obtain that

$$\begin{aligned}|\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_t) \rangle| &\leq |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}}) \rangle| + |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}} - \Delta_t) \rangle| \\ &\leq |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}}) \rangle| + \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}} - \Delta_t)\| \\ &\leq |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}}) \rangle| + \delta \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F.\end{aligned} \quad (21)$$

410 The last line is a consequence of the Restricted Isometry Property and Lemma 2.4, see inequality
411 (8). To estimate the first summand further, we use the triangle inequality again, and we obtain that

$$\begin{aligned}&|\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_{t,\mathbf{w}}) \rangle| \\ &\leq |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top}(\Delta_{t,\mathbf{w}})) \rangle| \\ &\stackrel{(a)}{=} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + \left| \left(\|\mathcal{A}(\mathbf{w} \mathbf{w}^\top)\|_2^2 - 1 \right) \langle \mathbf{w} \mathbf{w}^\top, \Delta_{t,\mathbf{w}} \rangle \right| \\ &\stackrel{(b)}{\leq} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + \delta |\langle \mathbf{w} \mathbf{w}^\top, \Delta_{t,\mathbf{w}} \rangle| \\ &\leq |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + \delta \|\Delta_{t,\mathbf{w}}\|.\end{aligned}$$

412 Equation (a) follows from the definition of $\mathcal{P}_{\mathbf{w} \mathbf{w}^\top}$ and $\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}$ and in inequality (b) we used the Re-
413 stricted Isometry Property; see Definition 2.1. Thus, by combining the last estimate with inequalities
414 (20) and (21) and taking the supremum over all $\mathbf{w} \in \mathcal{N}_\varepsilon$ we obtain that

$$\begin{aligned}&\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_t)\| \\ &\leq 2 \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + 2\delta \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F + 2\delta \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_{t,\mathbf{w}}\| \\ &\leq 2 \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| + 4\delta \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F + 2\delta \|\Delta_t\|.\end{aligned} \quad (22)$$

415 Since we assumed that the conclusion of Lemma 3.6 holds we obtain for the first summand that

$$\begin{aligned}\sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} |\langle \mathbf{w} \mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})) \rangle| &\leq 4\sqrt{\frac{d}{m}} \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}}))\|_2 \\ &\stackrel{(a)}{\leq} 8\sqrt{\frac{d}{m}} \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp}(\Delta_{t,\mathbf{w}})\|_F \\ &\leq 8\sqrt{\frac{d}{m}} \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_{t,\mathbf{w}}\|_F \\ &\leq 8\sqrt{\frac{d}{m}} \|\Delta_t\|_F + 8\sqrt{\frac{d}{m}} \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F \\ &\stackrel{(b)}{\leq} 8\sqrt{\frac{2rd}{m}} \|\Delta_t\| + 8\sqrt{\frac{d}{m}} \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F.\end{aligned}$$

416 Inequality (a) follows from the assumption that the operator \mathcal{A} has the Restricted Isometry Property
417 of order $2r + 2$ with an RIP-constant $\delta \leq 1$. To obtain inequality (b), we have used that the rank of
418 Δ_t is at most $2r$. Inserting the last estimate into (22), we obtain

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_t)\| \leq \left(16\sqrt{\frac{2rd}{m}} + 2\delta \right) \|\Delta_t\| + 4 \left(\delta + 4\sqrt{\frac{d}{m}} \right) \sup_{\mathbf{w} \in \mathcal{N}_\varepsilon} \|\Delta_t - \Delta_{t,\mathbf{w}}\|_F.$$

419 Inserting the definition of Δ_t and $\Delta_{t,\mathbf{w}}$ yields the claim. \square

4. Proof of the main result

4.1. Spectral Initialization

We provide the following lemma to show that both the original sequence and the virtual sequences are close to the ground truth \mathbf{X}_* at the spectral initialization. Moreover, this lemma guarantees that $\|\mathbf{U}_0\mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}}\mathbf{U}_{0,\mathbf{w}}^\top\|_F$ is sufficiently small. The proof of Lemma 4.1 is deferred to Appendix A.

Lemma 4.1. *There exists an absolute constant $C > 0$ such that the following holds:*

1. *With probability at least $1 - \exp(-4d)$, if $m > C^2\kappa^2rd$ is satisfied, it holds that*

$$\|\mathbf{X}_* - \mathbf{U}_0\mathbf{U}_0^\top\| \leq C\kappa\sigma_{\min}(\mathbf{X}_*)\sqrt{\frac{rd}{m}}. \quad (23)$$

2. *With probability at least $1 - \exp(-2d)$, if $m > 4C^2\kappa^2rd$ is satisfied, it holds for every $\mathbf{w} \in \mathcal{N}_\varepsilon$ that*

$$\|\mathbf{X}_* - \mathbf{U}_{0,\mathbf{w}}\mathbf{U}_{0,\mathbf{w}}^\top\| \leq 2C\kappa\sigma_{\min}(\mathbf{X}_*)\sqrt{\frac{rd}{m}}. \quad (24)$$

Consequently, if $m > 4C^2\kappa^2rd$, with probability at least $1 - 2\exp(-2d)$, it holds for every $\mathbf{w} \in \mathcal{N}_\varepsilon$ that

$$\|\mathbf{U}_0\mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}}\mathbf{U}_{0,\mathbf{w}}^\top\| \leq 3C\kappa\sigma_{\min}(\mathbf{X}_*)\sqrt{\frac{rd}{m}}. \quad (25)$$

3. *For any $\alpha \in (0, 1)$, assume $m \geq (51C^2 + \frac{C_1}{\alpha^2})\kappa^2rd$ for an absolute constant $C_1 > 0$. With probability at least $1 - 4\exp(-d)$, for every $\mathbf{w} \in \mathcal{N}_\varepsilon$,*

$$\|\mathbf{U}_0\mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}}\mathbf{U}_{0,\mathbf{w}}^\top\|_F \leq \left(2\alpha + C\kappa\sqrt{\frac{rd}{m}}\right) \left(2\sigma_{\min}(\mathbf{X}_*) + 3\sqrt{2}C\kappa\sqrt{\frac{rd}{m}}\sigma_{\min}(\mathbf{X}_*)\right) \quad (26)$$

4.2. Convergence Analysis

4.2.1. Outline of proof strategy

Before we explain our proof strategy, we want to recall the following convergence lemma which was proven in [10, Theorem 3.2] and [61]. It states that as soon as $\text{dist}(\mathbf{U}_t, \mathbf{U}_*)$ is small enough then $\text{dist}(\mathbf{U}_t, \mathbf{U}_*)$ converges to zero with linear rate. We state it in the version of the overview article [31, Theorem 4].

Lemma 4.2. *Assume that the measurement operator \mathcal{A} satisfies the Restricted Isometry Property for all matrices of rank at most $6r$ with constant $\delta_{6r} < 1/10$. Let $\mathbf{U}_0, \mathbf{U}_1, \mathbf{U}_2, \dots$ be a sequence of gradient descent iterates defined via equation (4). Assume that the step size satisfies $\mu \leq \frac{c_1}{\|\mathbf{X}_*\|}$ and*

$$\text{dist}^2(\mathbf{U}_T, \mathbf{U}_*) \leq \frac{1}{16}\sigma_{\min}(\mathbf{X}_*) \quad (27)$$

for some iteration number T . Then it holds for all $t \geq T$ that

$$\text{dist}^2(\mathbf{U}_t, \mathbf{U}_*) \leq (1 - c_2\mu\sigma_{\min}(\mathbf{X}_*))^{t-T} \text{dist}^2(\mathbf{U}_T, \mathbf{U}_*).$$

Here, $c_1, c_2 > 0$ are absolute numerical constants chosen small enough.

Note that the condition $\delta_{6r} < 1/10$ holds with high probability if the sample size satisfies $m \gtrsim rd$. However, condition (27) cannot be guaranteed for the spectral initialization, i.e., for $T = 0$, when $m \asymp rd\kappa^2$. For this reason, Lemma 4.2 is not directly applicable in our proof. To deal with this, we consider two different phases in our convergence analysis. Namely, we set

$$T := \left\lceil \frac{8}{\mu\sigma_{\min}(\mathbf{X}_*)} \log(16r) \right\rceil.$$

447 We will show that at the end of the first phase, which consists of the iterations $t = 0, 1, \dots, T$, con-
 448 dition (27) holds. The second phase starts at iteration T . For the second phase, we have established
 449 that condition (27) already holds we can directly apply Lemma 4.2 and we obtain linear conver-
 450 gence. Thus, our main focus in this section will be to analyze the first convergence phase.

451 In the following, we will give an outline of the analysis of this first phase. As is typical in the
 452 analysis of non-convex optimization algorithms, we will control several quantities simultaneously
 453 in each iteration via an induction argument. The following list contains an overview of these.

- 454 a) We will show that $\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F$ and $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F$ stay suffi-
 455 ciently small for each $\mathbf{w} \in \mathcal{N}_\varepsilon$. Together with Proposition 3.7, this allows us to control
 456 the deviation term $\|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|$.
- 457 b) We will show that for each iteration $t \in [T]$ it holds that $\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq c \sigma_{\min}(\mathbf{X}_*)$ for
 458 some small constant $c > 0$. This ensures that the gradient descent iterates stay in the basin
 459 of attraction, in which we can establish linear convergence.
- 460 c) We will establish that $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F$ decays linearly in each iteration. Combined
 461 with the result from b) this will allow us to establish linear convergence of $\text{dist}(\mathbf{U}_t, \mathbf{U}_*)$.

462 The remainder of this section is structured as follows. In Section 4.2.2 we will provide the techni-
 463 cal lemmas to control $\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F$ and $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F$ as described in a)
 464 above. In Section 4.2.3, we will provide the technical lemmas which allow us to control the quan-
 465 tities described above in b) and c). In Section 4.2.4, we will combine these ingredients to prove
 466 Proposition 4.10, which is our main result describing the convergence of the iterates $(\mathbf{U}_t)_{0 \leq t \leq T}$
 467 in the first convergence phase.

468 4.2.2. Lemmas for controlling the distance between the virtual sequences and the original 469 sequence

470 The goal of this section is to show that the virtual sequence iterates $(\mathbf{U}_{t,\mathbf{w}})_t$ stay sufficiently close
 471 to the original sequence $(\mathbf{U}_t)_t$. This will be established via induction. In the following, we will
 472 state all key lemmas. To keep the presentation concise, we have moved the proofs, which may be of
 473 independent interest, to Section B.

474 The first lemma in this section provides an a priori estimate. Its proof can be found in Section B.2.

475 **Lemma 4.3.** *For absolute constants $c_1, c_2, c_3 > 0$ chosen small enough the following statement is true. Let*
 476 *$\mathbf{w} \in \mathcal{N}_\varepsilon$ and assume that*

$$\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}, \quad (28)$$

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq c_1 \sigma_{\min}(\mathbf{X}_*), \quad (29)$$

$$\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq \sigma_{\min}(\mathbf{X}_*), \quad (30)$$

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \leq \frac{\sigma_{\min}(\mathbf{X}_*)}{80}, \quad (31)$$

477 and that the step size $\mu > 0$ satisfies $\mu \leq \frac{c_2}{\kappa \|\mathbf{X}_*\|}$. In addition, assume that the conclusions of Lemma 3.6 hold
 478 and that

$$\max \left\{ \delta; 8\sqrt{\frac{rd}{m}} \right\} \leq \frac{c_3}{\kappa}, \quad (32)$$

479 where $\delta = \delta_{4r+1}$ denotes the Restricted Isometry Property of rank $4r + 1$. Then it holds that

$$\|\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,\mathbf{w}} \mathbf{U}_{t+1,\mathbf{w}}^\top\|_F \leq \frac{\sqrt{\sqrt{2}-1}}{40} \sigma_{\min}(\mathbf{X}_*).$$

480 Under the assumption that this a priori estimate holds, the next lemma shows that the quantity
 481 $\|\mathbf{U}_t \mathbf{U}_t - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F$ can be bounded from above by the quantity $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F$.
 482 The proof of this lemma has been deferred to Section B.3.

483 **Lemma 4.4.** Let $\mathbf{w} \in \mathcal{N}_\varepsilon$ and assume that

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*\| \leq \frac{\sigma_{\min}(\mathbf{X}_*)}{1600}, \quad (33)$$

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \leq \frac{\sqrt{3(\sqrt{2}-1)} \cdot \sigma_{\min}(\mathbf{X}_*)}{40}. \quad (34)$$

484 Then it holds that

$$\|\mathbf{V}_{\mathbf{X}_*,\perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*,\perp}\|_F \leq \frac{3\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F}{5}. \quad (35)$$

485 Moreover, it holds that

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \leq 3\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F. \quad (36)$$

486 The following key lemma allows us to control $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F$ iteratively. Its proof
487 can be found in Section B.4.

488 **Lemma 4.5.** For sufficiently small absolute constants $c_1, c_2, c_3, c_4, c_5, c_6 > 0$ the following statement holds.
489 Let $\mathbf{w} \in \mathcal{N}_\varepsilon$ and assume that

$$\|\mathbf{V}_{\mathbf{X}_*,\perp}^\top \mathbf{V}_{\mathbf{U}_t}\| \leq c_1, \quad (37)$$

$$\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}, \quad (38)$$

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*\| \leq c_2 \sigma_{\min}(\mathbf{X}_*), \quad (39)$$

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \leq c_3 \sigma_{\min}(\mathbf{X}_*). \quad (40)$$

490 Moreover, assume that the step size satisfies $\mu \leq \frac{c_4}{\kappa\|\mathbf{X}_*\|}$. Furthermore, assume that the conclusion of Lemma
491 3.6 holds and that

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq c_5 \sigma_{\min}(\mathbf{X}_*), \quad (41)$$

$$\max\left\{\delta; 8\sqrt{\frac{2rd}{m}}\right\} \leq \frac{c_6}{\kappa}, \quad (42)$$

492 where $\delta = \delta_{4r+2}$ denotes the Restricted Isometry Constant of rank $4r+2$. Then, it holds that

$$\begin{aligned} & \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,\mathbf{w}} \mathbf{U}_{t+1,\mathbf{w}}^\top)\|_F \\ & \leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F + \mu \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|. \end{aligned}$$

493 4.2.3. Lemmas controlling the distance between \mathbf{X}_* and $\mathbf{U}_t \mathbf{U}_t^\top$

494 In the following, let $\|\cdot\|$ denote any matrix norm, which satisfies the inequality

$$\|\mathbf{XYZ}\| \leq \|\mathbf{X}\| \|\mathbf{Y}\| \|\mathbf{Z}\| \quad (43)$$

495 for all matrices \mathbf{X} , \mathbf{Y} , and \mathbf{Z} with dimensions such that the matrix product \mathbf{XYZ} is well-defined.
496 Note that all Schatten- p norms have this property. In particular, this includes the spectral norm $\|\cdot\|$
497 and the Frobenius norm $\|\cdot\|_F$.

498 In the following, we are interested in establishing upper bounds for $\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|$, where either
499 $\|\cdot\| = \|\cdot\|_F$ or $\|\cdot\| = \|\cdot\|$. Instead of estimating these quantities directly, we will instead derive
500 upper bounds for the quantity

$$\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|.$$

501 To be able to relate this quantity with $\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|$ one can then use the following lemma.

502 **Lemma 4.6.** Let $\|\cdot\|$ be a norm for which inequality (43) holds. Assume that

$$\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{V}_{\mathbf{U}_t}\| \leq \frac{1}{\sqrt{2}}. \quad (44)$$

503 Then the following inequalities hold:

$$\|\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*, \perp}\|\| \leq 2 \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{V}_{\mathbf{U}_t}\| \|\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_*, \perp}\|\|, \quad (45)$$

$$\|\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*\|\| \leq 2(1 + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{V}_{\mathbf{U}_t}\|) \|\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*)\|\|. \quad (46)$$

504 A comparable lemma was proven in [23] in a more general setting but with less explicit constants.
505 For the sake of completeness, we included in Appendix C.1.

506 The following lemma allows us to control the quantity $\|\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|\|$ iteratively. We note
507 that a similar lemma has already been proven in [23] in a more general setting with less explicit
508 constants. For the sake of completeness, we again included a proof in Appendix C.2.

509 **Lemma 4.7.** Let $\|\cdot\|$ be a norm which is submultiplicative in the sense of inequality (43). Assume that

$$\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{V}_{\mathbf{U}_t}\| \leq \frac{1}{2}, \quad (47)$$

$$\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|},$$

$$\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq \frac{\sigma_{\min}(\mathbf{X}_*)}{48}, \quad (48)$$

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq \frac{1}{48} \sigma_{\min}(\mathbf{X}_*), \quad (49)$$

510 and that the step size satisfies $\mu \leq \frac{1}{1024\kappa\|\mathbf{X}_*\|}$. Then it holds that

$$\begin{aligned} & \|\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{X}_*)\|\| \\ & \leq \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|\| + 5\mu \|\mathbf{X}_*\| \|\|[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_t}\|\|. \end{aligned}$$

511 Given an upper bound for $\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|_F$ we can obtain an estimate for $\text{dist}(\mathbf{U}_t, \mathbf{U}_*)$ by using the
512 following technical lemma.

513 **Lemma 4.8** (Lemma 5.4 in [10]). Let $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$ be two matrices and assume that $\text{rank}(\mathbf{U}) =$
514 $\min\{r, d\}$. Then it holds that

$$\text{dist}^2(\mathbf{U}, \mathbf{V}) \leq \frac{1}{2(\sqrt{2}-1)\sigma_{\min}^2(\mathbf{U})} \|\|\mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top\|_F^2,$$

515 where $\text{dist}(\mathbf{U}, \mathbf{V})$ is defined in (5).

516 To check the prerequisite of the Davis-Kahan inequality (Lemma 2.6) in our proof, we will also need
517 the following auxiliary lemma, which provides us with an a priori bound for $\|\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top\|$. Its
518 proof can be found in Appendix C.3.

519 **Lemma 4.9.** There are absolute constants $c_1, c_2, c_3 > 0$ such that the following holds. Assume that $\mu \leq$

520 $\frac{c_1}{\|\mathbf{X}_*\|}$ and

$$\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}, \quad (50)$$

$$\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq c_2 \sigma_{\min}(\mathbf{X}_*), \quad (51)$$

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq c_3 \sigma_{\min}(\mathbf{X}_*). \quad (52)$$

521 Then it holds that

$$\|\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top\| \leq \left(1 - \frac{1}{\sqrt{2}}\right) \sigma_{\min}(\mathbf{X}_*).$$

522 **4.2.4. Statement and proof of the main convergence lemma**

523 We now have all the ingredients in place to prove the main lemma in this section, which is stated
524 below.

525 **Lemma 4.10.** *There are absolute constants $c_1, c_2, c_3, c_4 > 0$ chosen sufficiently small such that the following
526 statement holds. Assume that the spectral initialization \mathbf{U}_0 satisfies*

$$\|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\| \leq c_1 \sigma_{\min}(\mathbf{X}_*) \quad (53)$$

527 and that for every $\mathbf{w} \in \mathcal{N}_\varepsilon$ we have that

$$\|\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top\|_F \leq c_2 \sigma_{\min}(\mathbf{X}_*). \quad (54)$$

528 Moreover, we assume that the conclusion of Lemma 3.6 holds for

$$T = \left\lceil \frac{8}{\mu \sigma_{\min}(\mathbf{X}_*)} \log(16r) \right\rceil.$$

529 Furthermore, we assume that

$$\max \left\{ \delta; 8\sqrt{\frac{2rd}{m}} \right\} \leq \frac{c_3}{\kappa}, \quad (55)$$

530 where $\delta = \delta_{4r+2}$ denotes the Restricted Isometry Property of order $4r + 2$. In addition, assume that $\mu \leq$
531 $\frac{c_4}{\kappa \|\mathbf{X}_*\|}$. Then for every iteration t with $0 \leq t \leq T$ it holds that

$$\text{dist}^2(\mathbf{U}_t, \mathbf{U}_*) \leq r \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16} \right)^{2t} \|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\|. \quad (56)$$

532 In particular, we have that

$$\text{dist}^2(\mathbf{U}_T, \mathbf{U}_*) \leq \frac{1}{16} \sigma_{\min}(\mathbf{X}_*), \quad (57)$$

533 where $\mathbf{U}_* \in \mathbb{R}^{n \times r}$ denotes a matrix which satisfies $\mathbf{U}_* \mathbf{U}_*^\top = \mathbf{X}_*$.

534 *Proof of Lemma 4.10.* We prove by induction that for all iterations t with $0 \leq t \leq T$ the following
535 inequalities hold:

$$\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \leq \left(1 - \frac{\mu}{16} \sigma_{\min}(\mathbf{X}_*) \right)^t \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top)\|_F, \quad (58)$$

$$\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq c_1 \sigma_{\min}(\mathbf{X}_*), \quad (59)$$

$$\|\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t}\| \leq \sqrt{2} c_1, \quad (60)$$

$$\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq 3c_1 \sigma_{\min}(\mathbf{X}_*), \quad (61)$$

536 and, for every $\mathbf{w} \in \mathcal{N}_\varepsilon$,

$$\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \leq c_2 \sigma_{\min}(\mathbf{X}_*), \quad (62)$$

$$\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \leq 3c_2 \sigma_{\min}(\mathbf{X}_*). \quad (63)$$

537 The constants $c_1, c_2 > 0$ are the same as in assumptions (53) and (54) and are thus, in particular,
538 independent of the iteration number t .

539 First, we check that these inequalities hold for $t = 0$. Inequality (58) is immediate. Inequalities (59)
540 and (61) follow from assumption (53). Inequalities (62) and (63) are due to assumption (54). It
541 remains to establish inequality (60) for $t = 0$. Using the Davis-Kahan inequality (see Lemma 2.6)
542 and assumption (53) it follows that

$$\|\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_0}\| \leq \frac{\sqrt{2} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top)\|}{\sigma_{\min}(\mathbf{X}_*)} \leq \sqrt{2} c_1.$$

543 This shows that the above inequalities hold for $t = 0$.

544

545 For the induction step, assume now that these inequalities hold for some t . First, we observe that
 546 it follows from the induction assumption (61) and Weyl's inequalities that $\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}$ for
 547 $c_1 < 1/3$. Moreover, note that since we assumed that the conclusion of Lemma 3.6 holds we obtain
 548 from Proposition 3.7 that

$$\begin{aligned}
 & \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \\
 & \leq \left(16\sqrt{\frac{2rd}{m}} + 2\delta\right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + 4\left(\delta + 4\sqrt{\frac{d}{m}}\right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\
 & \stackrel{(a)}{\leq} \frac{4c_3}{\kappa} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{6c_3}{\kappa} \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\
 & \stackrel{(b)}{\leq} \frac{10c_3}{\kappa} \sigma_{\min}(\mathbf{X}_*), \tag{64}
 \end{aligned}$$

549 where inequality (a) follows from assumption (55). Inequality (b) is due to the induction hypotheses
 550 (61) and (63) with $c_1 \leq 1/3$ and $c_2 \leq 1/3$. Next, we note that from Lemma 4.7 applied with
 551 $\|\cdot\| = \|\cdot\|_F$ it follows that

$$\begin{aligned}
 & \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{X}_*)\|_F \\
 & \leq \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + 5\mu \|\mathbf{X}_*\| \|[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_t}\|_F \\
 & \stackrel{(a)}{\leq} \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + 5\mu\delta \|\mathbf{X}_*\| \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \stackrel{(b)}{\leq} \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + 15\mu\delta \|\mathbf{X}_*\| \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\
 & \stackrel{(c)}{\leq} \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + \frac{15\mu c_3 \|\mathbf{X}_*\|}{\kappa} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\
 & \stackrel{(d)}{\leq} \left(1 - \frac{\mu}{16} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F.
 \end{aligned}$$

552 Inequality (a) follows from the Restricted Isometry Property combined with Lemma 2.4. Inequal-
 553 ity (b) is due to Lemma 4.6 and inequality (60). Inequality (c) follows from assumption (55) and
 554 inequality (d) is due to the fact we can choose $c_3 \leq \frac{1}{240}$. Thus, using the induction assumption, we
 555 see that inequality (58) holds for $t + 1$.

556 Next, our goal is to prove inequality (59) for $t + 1$. For that, we note that it follows from Lemma 4.7
 557 with $\|\cdot\| = \|\cdot\|$ that

$$\begin{aligned}
 & \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{X}_*)\| \\
 & \leq \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + 5\mu \|\mathbf{X}_*\| \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \\
 & \stackrel{(a)}{\leq} \left(1 - \frac{\mu}{8} \sigma_{\min}(\mathbf{X}_*)\right) c_1 \sigma_{\min}(\mathbf{X}_*) + 50c_3 \mu \sigma_{\min}^2(\mathbf{X}_*) \\
 & \stackrel{(b)}{\leq} c_1 \sigma_{\min}(\mathbf{X}_*), \tag{65}
 \end{aligned}$$

558 where inequality (a) follows from the induction hypothesis (59) and inequality (64). Inequality (b)
 559 holds since we can choose c_1 and c_3 in such a way that $c_3 \leq \frac{c_1}{400}$. This proves inequality (59) for
 560 $t + 1$.

561 We observe that Lemma 4.9 yields the a-priori bound

$$\|\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top\| \leq \left(1 - \frac{1}{\sqrt{2}}\right) \sigma_{\min}(\mathbf{X}_*).$$

562 Thus, we can apply the Davis-Kahan inequality (see Lemma 2.6) which together with inequality
 563 (65) yields that

$$\|\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_{t+1}}\| \leq \frac{\sqrt{2} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{X}_*)\|}{\sigma_{\min}(\mathbf{X}_*)} \leq \sqrt{2} c_1.$$

564 This proves inequality (60) for $t + 1$. Next, we apply Lemma 4.6 and (65) to obtain that

$$\begin{aligned} \|\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top\| &\leq 2(1 + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{V}_{\mathbf{U}_{t+1}}\|) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top)\| \\ &\leq 3 \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top)\| \leq 3c_1 \sigma_{\min}(\mathbf{X}_*), \end{aligned}$$

565 which proves inequality (61) for $t + 1$.

566 Next, we can apply Lemma 4.5 since all assumptions are satisfied and it follows that

$$\begin{aligned} &\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1, \mathbf{w}} \mathbf{U}_{t+1, \mathbf{w}}^\top)\|_F \\ &\leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + \mu \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\stackrel{(a)}{\leq} \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right) c_2 \sigma_{\min}(\mathbf{X}_*) + 3c_1 \mu \sigma_{\min}^2(\mathbf{X}_*) \\ &\stackrel{(b)}{\leq} c_2 \sigma_{\min}(\mathbf{X}_*). \end{aligned} \tag{66}$$

567 Inequality (a) is due to inequalities (61) and (62). Inequality (b) holds since we can choose that
 568 $c_1 \leq \frac{c_2}{48}$. This proves inequality (62).

569 Next, we want to prove inequality (63) for $t + 1$. First, we apply Lemma 4.3 and we obtain for all
 570 $\mathbf{w} \in \mathcal{N}_\varepsilon$ the a-priori bound

$$\|\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1, \mathbf{w}} \mathbf{U}_{t+1, \mathbf{w}}^\top\|_F \leq \frac{\sqrt{\sqrt{2}-1}}{40} \cdot \sigma_{\min}(\mathbf{X}_*).$$

571 This allows us to apply Lemma 4.4 and we obtain for all $\mathbf{w} \in \mathcal{N}_\varepsilon$ the sharper bound

$$\|\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1, \mathbf{w}} \mathbf{U}_{t+1, \mathbf{w}}^\top\|_F \leq 3 \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1, \mathbf{w}} \mathbf{U}_{t+1, \mathbf{w}}^\top)\|_F \stackrel{(66)}{\leq} 3c_2 \sigma_{\min}(\mathbf{X}_*),$$

572 which shows inequality (63) for $t + 1$. This completes the induction step.

573

574 To complete the proof of Lemma 4.10 it remains to prove inequalities (56) and (57). For that, we
 575 first observe that

$$\begin{aligned} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|_F &\stackrel{(a)}{\leq} 3 \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\ &\stackrel{(b)}{\leq} 3 \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right)^t \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top)\|_F \\ &\stackrel{(c)}{\leq} 3\sqrt{2}r \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right)^t \|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\|. \end{aligned}$$

576 Inequality (a) follows from Lemma 4.6 with $\|\cdot\| = \|\cdot\|_F$ which is applicable since we have shown
 577 by induction that (60) holds for $0 \leq t \leq T$. Inequality (b) holds since we have proven (58) for all
 578 $0 \leq t \leq T$. Inequality (c) holds since $\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top$ has rank at most $2r$. Thus, we can apply Lemma
 579 4.8 and obtain that

$$\begin{aligned} \text{dist}^2(\mathbf{U}_t, \mathbf{U}_*) &\leq \frac{\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|_F^2}{2(\sqrt{2}-1)\sigma_{\min}(\mathbf{X}_*)} \\ &\leq 18r \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right)^{2t} \cdot \frac{\|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\|^2}{2(\sqrt{2}-1)\sigma_{\min}(\mathbf{X}_*)} \\ &\leq \frac{9c_1 r}{(\sqrt{2}-1)} \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right)^{2t} \|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\|, \end{aligned}$$

580 where in the last inequality, we have used assumption (53). This proves inequality (56) since $c_1 \leq$
 581 $\frac{\sqrt{2}-1}{9}$. Next, we note that for $t = T$, the above inequality yields that

$$\begin{aligned} \text{dist}^2(\mathbf{U}_T, \mathbf{U}_\star) &\stackrel{(a)}{\leq} \frac{9c_1^2 r}{(\sqrt{2}-1)} \left(1 - \frac{\mu\sigma_{\min}(\mathbf{X}_\star)}{16}\right)^{2T} \sigma_{\min}(\mathbf{X}_\star) \\ &\stackrel{(b)}{\leq} \frac{9c_1^2 r}{(\sqrt{2}-1)} \exp\left(\frac{-T\mu\sigma_{\min}(\mathbf{X}_\star)}{8}\right) \sigma_{\min}(\mathbf{X}_\star) \\ &\stackrel{(c)}{\leq} \frac{\sigma_{\min}(\mathbf{X}_\star)}{16}. \end{aligned}$$

582 In inequality (a), we have used again assumption (53). Inequality (b) is due to the elementary
 583 inequality $\ln(1+x) \leq x$ for $-1 < x$ and the assumption $\mu < \frac{c_4}{\kappa\|\mathbf{X}_\star\|}$ for sufficiently small $c_4 > 0$.

584 Inequality (c) follows from $T = \left\lceil \frac{8}{\mu\sigma_{\min}(\mathbf{X}_\star)} \log(16r) \right\rceil$ (and from the fact that we can choose $c_1 \leq$
 585 $\frac{\sqrt{\sqrt{2}-1}}{3}$). This proves inequality (57). Thus, the proof of Lemma 4.10 is complete. \square

586 4.3. Proof of Theorem 1.2

587 Now we have all the ingredients in place to prove the main result of this paper, Theorem 1.2.

588 *Proof of Theorem 1.2.* In the following $c > 0$ denotes a sufficiently small absolute constant. First,
 589 by Lemma 2.2 we know that due to our assumption $m \gtrsim rd\kappa^2$, with probability $1 - \exp(-d)$ the
 590 measurement operator \mathcal{A} satisfies the Restricted Isometry Property of order $6r$ with a constant $\delta =$
 591 $\delta_{6r} \leq \frac{c}{\kappa}$, where $c > 0$ is a sufficiently small absolute constant.

592 Set

$$T := \left\lceil \frac{8}{\mu\sigma_{\min}(\mathbf{X}_\star)} \log(16r) \right\rceil.$$

593 Note that since $r \geq 1$ and the assumption $\mu \leq \frac{c_1}{\sigma_{\min}(\mathbf{X}_\star)}$ for small $c_1 > 0$, we have $T \geq 1$. Let \mathcal{N}_ε be
 594 an ε -net of the unit sphere in \mathbb{R}^d with $\varepsilon = 1/2$ such that $|\mathcal{N}_\varepsilon| \leq 6^d$. Now note that $2T \leq 6^d$, where
 595 we have used the assumption $\mu \geq \frac{32}{\sigma_{\min}(\mathbf{X}_\star)6^d} \log(16r)$. Thus, it follows from Lemma 3.6 that with
 596 probability at least $1 - 2\exp(-10d)$ it holds that

$$|\langle \mathbf{w}\mathbf{w}^\top, (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_\star - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) \rangle| \leq 4\sqrt{\frac{d}{m}} \|\mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_\star - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top))\|_2$$

597 for all $\mathbf{w} \in \mathcal{N}_\varepsilon$ and for all $0 \leq t \leq T$. Next, we know from Lemma 4.1 and due to our assumption
 598 $m \gtrsim rd\kappa^2$ that with probability at least $1 - 5\exp(-d)$, the inequalities

$$\begin{aligned} \|\mathbf{X}_\star - \mathbf{U}_0\mathbf{U}_0^\top\| &\leq c\sigma_{\min}(\mathbf{X}_\star), \\ \|\mathbf{U}_0\mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}}\mathbf{U}_{0,\mathbf{w}}^\top\|_F &\leq c\sigma_{\min}(\mathbf{X}_\star) \end{aligned} \quad (67)$$

599 hold for a sufficiently small constant $c > 0$. Thus, all the assumptions of Lemma 4.10 are fulfilled.
 600 It follows that

$$\text{dist}^2(\mathbf{U}_t, \mathbf{U}_\star) \leq r \left(1 - \frac{\mu\sigma_{\min}(\mathbf{X}_\star)}{16}\right)^{2t} \|\mathbf{X}_\star - \mathbf{U}_0\mathbf{U}_0^\top\| \quad (68)$$

601 for all $0 \leq t \leq T$ and

$$\text{dist}(\mathbf{U}_T, \mathbf{U}_\star) \leq \frac{\sigma_{\min}(\mathbf{X}_\star)}{16}. \quad (69)$$

602 Due to inequality (69) and since $\delta_{6r} < 1/10$ we can apply Lemma 4.2 which yields that for $t \geq T$,

$$\text{dist}^2(\mathbf{U}_t, \mathbf{U}_\star) \leq (1 - c\mu\sigma_{\min}(\mathbf{X}_\star))^{t-T} \text{dist}^2(\mathbf{U}_T, \mathbf{U}_\star). \quad (70)$$

603 Thus, by combining (67), (68), and (70) we obtain the conclusion of Theorem 1.2. \square

5. Discussions

In this paper, we have shown that for symmetric matrix sensing, factorized gradient descent can recover the ground truth matrix as soon as the number of samples satisfies $m \gtrsim rd\kappa^2$. This improves over previous results in the literature with a quadratic rank dependence. The key ingredient in our proof is a combination of a virtual sequence argument with an ε -net argument.

Going forward, our work opens up a number of exciting research directions. In the following, we highlight a few of these.

- *Breaking the quadratic rank barrier in related non-convex matrix sensing problems:* We expect that our novel proof technique will pave the way to break the quadratic rank barrier in the sample complexity in various related non-convex matrix sensing problems. This includes matrix sensing with an asymmetric ground truth matrix or overparameterized matrix sensing with small random initialization [22]. One might also examine whether our new proof technique can be used to remove the additional rank factor in the sample complexity in related algorithms such as *scaled gradient descent* [11] or *GSMR* [30].
- *Removing the condition number dependence in the sample complexity:* Compared to the nuclear norm minimization approach, the sample complexity in Theorem 1.2 is still suboptimal since it depends quadratically on the condition number of the ground truth matrix \mathbf{X}_* . Indeed, all related results in the non-convex low-rank matrix recovery also have such a dependency on the condition number. It would be interesting to examine whether this dependence on the condition number is actually needed.
- *Beyond Gaussian measurement matrices:* It would also be interesting to examine whether the argument in this paper can be adapted to scenarios where the measurement matrices are no longer Gaussian, e.g., the matrix completion problem. Since the proof presented in this paper heavily relies on the orthogonal invariance of the Gaussian distribution, new insights are likely required to handle scenarios where this property is no longer available. We believe that this is an exciting research direction.

References

- [1] Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- [2] Emmanuel J. Candès, Thomas Strohmer, and Vladislav Voroninski. Phaselift: exact and stable signal recovery from magnitude measurements via convex programming. *Commun. Pure Appl. Math.*, 66(8):1241–1274, 2013. ISSN 0010-3640. doi: 10.1002/cpa.21432.
- [3] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):1–37, 2011.
- [4] Ali Ahmed, Benjamin Recht, and Justin Romberg. Blind deconvolution using convex programming. *IEEE Trans. Inf. Theory*, 60(3):1711–1732, 2014. ISSN 0018-9448. doi: 10.1109/TIT.2013.2294644.
- [5] Shuyang Ling and Thomas Strohmer. Blind deconvolution meets blind demixing: algorithms and performance bounds. *IEEE Trans. Inf. Theory*, 63(7):4497–4520, 2017. ISSN 0018-9448. doi: 10.1109/TIT.2017.2701342.
- [6] Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.*, 52(3):471–501, 2010. ISSN 0036-1445. doi: 10.1137/070697835. URL hdl.handle.net/1721.1/60575.
- [7] Emmanuel J. Candès and Terence Tao. The power of convex relaxation: near-optimal matrix completion. *IEEE Trans. Inf. Theory*, 56(5):2053–2080, 2010. ISSN 0018-9448. doi: 10.1109/TIT.2010.2044061.

- 650 [8] David Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inf.*
651 *Theory*, 57(3):1548–1566, 2011. ISSN 0018-9448. doi: 10.1109/TIT.2011.2104999.
- 652 [9] Peter Jung, Felix Krahmer, and Dominik Stöger. Blind demixing and deconvolution at near-
653 optimal rate. *IEEE Trans. Inf. Theory*, 64(2):704–727, 2018. ISSN 0018-9448. doi: 10.1109/TIT.
654 2017.2784481.
- 655 [10] Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank
656 solutions of linear matrix equations via procrustes flow. In *International Conference on Machine*
657 *Learning*, pages 964–973. PMLR, 2016.
- 658 [11] Tian Tong, Cong Ma, and Yuejie Chi. Accelerating ill-conditioned low-rank matrix estimation
659 via scaled gradient descent. *J. Mach. Learn. Res.*, 22:63, 2021. ISSN 1532-4435. URL [jmlr.
660 csail.mit.edu/papers/v22/20-1067.html](http://jmlr.csail.mit.edu/papers/v22/20-1067.html). Id/No 150.
- 661 [12] Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and René Vidal. Nonconvex robust low-rank ma-
662 trix recovery. *SIAM J. Optim.*, 30(1):660–686, 2020. ISSN 1052-6234. doi: 10.1137/18M1224738.
- 663 [13] Vasileios Charisopoulos, Yudong Chen, Damek Davis, Mateo Díaz, Lijun Ding, and Dmitriy
664 Drusvyatskiy. Low-rank matrix recovery with composite optimization: good conditioning
665 and rapid convergence. *Found. Comput. Math.*, 21(6):1505–1593, 2021. ISSN 1615-3375. doi:
666 10.1007/s10208-020-09490-9.
- 667 [14] Raghunandan H. Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a
668 few entries. *IEEE Trans. Inf. Theory*, 56(6):2980–2998, 2010. ISSN 0018-9448. doi: 10.1109/TIT.
669 2010.2046205.
- 670 [15] Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factoriza-
671 tion. *IEEE Trans. Inf. Theory*, 62(11):6535–6579, 2016. ISSN 0018-9448. doi: 10.1109/TIT.2016.
672 2598574.
- 673 [16] Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion
674 using burer-monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.
- 675 [17] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum.
676 *Advances in Neural Information Processing Systems*, 29, 2016.
- 677 [18] Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex
678 statistical estimation: gradient descent converges linearly for phase retrieval, matrix comple-
679 tion, and blind deconvolution. *Found. Comput. Math.*, 20(3):451–632, 2020. ISSN 1615-3375.
680 doi: 10.1007/s10208-019-09429-9.
- 681 [19] Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient
682 descent without ℓ_2, ∞ regularization. *IEEE Trans. Inf. Theory*, 66(9):5806–5841, 2020. ISSN
683 0018-9448. doi: 10.1109/TIT.2020.2992234.
- 684 [20] Shuyang Ling and Thomas Strohmer. Regularized gradient descent: a non-convex recipe for
685 fast joint blind deconvolution and demixing. *Inf. Inference*, 8(1):1–49, 2019. ISSN 2049-8764.
686 doi: 10.1093/imaiai/iax022.
- 687 [21] Jialin Dong and Yuanming Shi. Nonconvex demixing from bilinear measurements. *IEEE Trans.*
688 *Signal Process.*, 66(19):5152–5166, 2018. ISSN 1053-587X. doi: 10.1109/TSP.2018.2864660.
- 689 [22] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-
690 parameterized matrix sensing and neural networks with quadratic activations. In *Conference*
691 *On Learning Theory*, pages 2–47. PMLR, 2018.
- 692 [23] Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral
693 learning: Optimization and generalization guarantees for overparameterized low-rank matrix
694 reconstruction. *Advances in Neural Information Processing Systems*, 34:23831–23843, 2021.

- 695 [24] Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon Shaolei Du, and Jason D Lee. Understanding incre-
696 mental learning of gradient descent: A fine-grained analysis of matrix sensing. In *International*
697 *Conference on Machine Learning*, pages 15200–15238. PMLR, 2023.
- 698 [25] Xingyu Xu, Yandi Shen, Yuejie Chi, and Cong Ma. The power of preconditioning in overpa-
699 rameterized low-rank matrix sensing. In *International Conference on Machine Learning*, pages
700 38611–38654. PMLR, 2023.
- 701 [26] Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regulariza-
702 tion: Generalization and convergence guarantees for overparameterized asymmetric matrix
703 sensing. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 5140–5142. PMLR,
704 2023.
- 705 [27] Jianhao Ma and Salar Fattahi. Convergence of gradient descent with small initialization for
706 unregularized matrix completion. *arXiv preprint arXiv:2402.06756*, 2024.
- 707 [28] Johan S Wind. Asymmetric matrix sensing by gradient descent with small random initializa-
708 tion. *arXiv preprint arXiv:2309.01796*, 2023.
- 709 [29] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using
710 alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of*
711 *computing*, pages 665–674, 2013.
- 712 [30] Pini Zilber and Boaz Nadler. GNMR: a provable one-line algorithm for low rank matrix recov-
713 ery. *SIAM J. Math. Data Sci.*, 4(2):909–934, 2022. ISSN 2577-0187. doi: 10.1137/21M1433812.
- 714 [31] Yuejie Chi, Yue M. Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix fac-
715 torization: an overview. *IEEE Trans. Signal Process.*, 67(20):5239–5269, 2019. ISSN 1053-587X.
716 doi: 10.1109/TSP.2019.2937282.
- 717 [32] Prateek Jain, Raghu Meka, and Inderjit Dhillon. Guaranteed rank minimization via singular
718 value projection. *Advances in Neural Information Processing Systems*, 23, 2010.
- 719 [33] Lijun Ding and Yudong Chen. Leave-one-out approach for matrix completion: primal and
720 dual analysis. *IEEE Trans. Inf. Theory*, 66(11):7274–7301, 2020. ISSN 0018-9448. doi: 10.1109/
721 TIT.2020.2992769.
- 722 [34] Jared Tanner and Ke Wei. Normalized iterative hard thresholding for matrix completion.
723 *SIAM J. Sci. Comput.*, 35(5):s104–s125, 2013. ISSN 1064-8275. doi: 10.1137/120876459. URL
724 semanticsscholar.org/paper/9b785002627fd2066fce004199758ce137a1ce61.
- 725 [35] Karthik Mohan and Maryam Fazel. Iterative reweighted algorithms for matrix rank minimiza-
726 tion. *J. Mach. Learn. Res.*, 13:3441–3473, 2012. ISSN 1532-4435. URL [www.jmlr.org/papers/
727 v13/mohan12a.html](http://www.jmlr.org/papers/v13/mohan12a.html).
- 728 [36] Massimo Fornasier, Holger Rauhut, and Rachel Ward. Low-rank matrix recovery via iteratively
729 reweighted least squares minimization. *SIAM J. Optim.*, 21(4):1614–1640, 2011. ISSN 1052-
730 6234. doi: 10.1137/100811404.
- 731 [37] Christian Kümmerle and Juliane Sigl. Harmonic mean iteratively reweighted least squares for
732 low-rank matrix recovery. *J. Mach. Learn. Res.*, 19:49, 2018. ISSN 1532-4435. URL [jmlr.csail.
733 mit.edu/papers/v19/17-244.html](http://jmlr.csail.mit.edu/papers/v19/17-244.html). Id/No 47.
- 734 [38] Christian Kümmerle and Claudio M Verdun. A scalable second order method for ill-
735 conditioned matrix completion from few samples. In *International Conference on Machine Learn-
736 ing*, pages 5872–5883. PMLR, 2021.
- 737 [39] Kiryung Lee and Yoram Bresler. ADMiRA: atomic decomposition for minimum rank approxi-
738 mation. *IEEE Trans. Inf. Theory*, 56(9):4402–4416, 2010. ISSN 0018-9448. doi: 10.1109/TIT.2010.
739 2054251.

- 740 [40] Ke Wei, Jian-Feng Cai, Tony F Chan, and Shingyu Leung. Guarantees of Riemannian optimization for low rank matrix recovery. *SIAM J. Matrix Anal. Appl.*, 37(3):1198–1222, 2016.
741
- 742 [41] Bart Vandereycken. Low-rank matrix completion by Riemannian optimization. *SIAM J. Optim.*,
743 23(2):1214–1236, 2013. ISSN 1052-6234. doi: 10.1137/110845768. URL [semanticscholar.org/
744 paper/feb9713f4e7614aecdb4778c0bc8c2dced60a325](https://semanticscholar.org/paper/feb9713f4e7614aecdb4778c0bc8c2dced60a325).
- 745 [42] Guillaume Olikier, André Uschmajew, and Bart Vandereycken. Gauss-southwell type descent
746 methods for low-rank matrix optimization. *arXiv preprint arXiv:2306.00897*, 2023.
- 747 [43] Greg W Anderson, Alice Guionnet, and Ofer Zeitouni. *An introduction to random matrices*.
748 Number 118. Cambridge university press, 2010.
- 749 [44] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search
750 for low rank matrix recovery. *Advances in Neural Information Processing Systems*, 29, 2016.
- 751 [45] Dohyung Park, Anastasios Kyrillidis, Constantine Carmanis, and Sujay Sanghavi. Non-square
752 matrix sensing without spurious local minima via the burer-monteiro approach. In *Artificial
753 Intelligence and Statistics*, pages 65–74. PMLR, 2017.
- 754 [46] André Uschmajew and Bart Vandereycken. On critical points of quadratic low-rank matrix
755 optimization problems. *IMA J. Numer. Anal.*, 40(4):2626–2651, 2020. ISSN 0272-4979. doi:
756 10.1093/imanum/drz061.
- 757 [47] Richard Y. Zhang, Somayeh Sojoudi, and Javad Lavaei. Sharp restricted isometry bounds for
758 the inexistence of spurious local minima in nonconvex matrix recovery. *J. Mach. Learn. Res.*, 20:
759 34, 2019. ISSN 1532-4435. URL jmlr.csail.mit.edu/papers/v20/19-020.html. Id/No 114.
- 760 [48] Jason D. Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I. Jordan, and
761 Benjamin Recht. First-order methods almost always avoid strict saddle points. *Math. Program.*,
762 176(1-2 (B)):311–337, 2019. ISSN 0025-5610. doi: 10.1007/s10107-019-01374-3.
- 763 [49] Simon S Du, Chi Jin, Jason D Lee, Michael I Jordan, Aarti Singh, and Barnabas Poczos. Gradient
764 descent can take exponential time to escape saddle points. *Advances in neural information
765 processing systems*, 30, 2017.
- 766 [50] Emmanuel J. Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery
767 from a minimal number of noisy random measurements. *IEEE Trans. Inf. Theory*, 57(4):2342–
768 2359, 2011. ISSN 0018-9448. doi: 10.1109/TIT.2011.2111771.
- 769 [51] Chandler Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM J.*
770 *Numer. Anal.*, 7:1–46, 1970. ISSN 0036-1429. doi: 10.1137/0707001.
- 771 [52] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Spectral methods for data science: a
772 statistical perspective. *Found. Trends Mach. Learn.*, 14(5):1–246, 2021. ISSN 1935-8237. doi:
773 10.1561/22000000079.
- 774 [53] Jiacheng Zhuo, Jeongyeol Kwon, Nhat Ho, and Constantine Caramanis. On the computational
775 and statistical complexity of over-parameterized matrix sensing. *Journal of Machine Learning
776 Research*, 25(169):1–47, 2024.
- 777 [54] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic
778 theory of independence*. Oxford university press, 2013.
- 779 [55] NL Johnson, S Kotz, and N Balakrishnan. Chi-squared distributions including Chi and
780 Rayleigh. *Continuous univariate distributions*, pages 415–493, 1994.
- 781 [56] Roger A Horn and Charles R Johnson. *Topics in matrix analysis*. Cambridge university press,
782 1994.

- 783 [57] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv*
784 *preprint arXiv:1011.3027*, 2010.
- 785 [58] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initial-
786 ization: fast global convergence for nonconvex phase retrieval. *Math. Program.*, 176(1-2 (B)):
787 5–37, 2019. ISSN 0025-5610. doi: 10.1007/s10107-019-01363-6.
- 788 [59] Kiryung Lee and Dominik Stöger. Randomly initialized alternating least squares: Fast conver-
789 gence for matrix sensing. *SIAM Journal on Mathematics of Data Science*, 5(3):774–799, 2023. doi:
790 10.1137/22M1506456. URL <https://doi.org/10.1137/22M1506456>.
- 791 [60] Roman Vershynin. *High-dimensional probability. An introduction with applications in data science*,
792 volume 47 of *Camb. Ser. Stat. Probab. Math.* Cambridge: Cambridge University Press, 2018.
793 ISBN 978-1-108-41519-4; 978-1-108-23159-6. doi: 10.1017/9781108231596.
- 794 [61] Qinqing Zheng and John Lafferty. A convergent gradient descent algorithm for rank minimiza-
795 tion and semidefinite programming from random linear measurements. *Advances in Neural*
796 *Information Processing Systems*, 28, 2015.
- 797 [62] Felix Krahmer, Shahar Mendelson, and Holger Rauhut. Suprema of chaos processes and the
798 restricted isometry property. *Commun. Pure Appl. Math.*, 67(11):1877–1904, 2014. ISSN 0010-
799 3640. doi: 10.1002/cpa.21504.
- 800 [63] Michel Talagrand. *The generic chaining. Upper and lower bounds of stochastic processes*. Springer
801 Monogr. Math. Berlin: Springer, 2005. ISBN 3-540-24518-9; 978-3-642-06386-2; 978-3-540-
802 27499-5. doi: 10.1007/3-540-27499-5.

803 A. Proof for the Spectral Initialization (Proof of Lemma 4.1)

804 *Proof of Lemma 4.1.* (1) We write

$$(\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - \mathbf{X}_* = \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{A}_i - \mathbf{X}_*).$$

805 Let $\widetilde{\mathcal{N}}_\varepsilon$ be any ε -net on S^{d-1} with $\varepsilon = \frac{1}{2}$ of size at most 6^d . Then we have

$$\begin{aligned} \left\| (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - \mathbf{X}_* \right\| &\leq 2 \sup_{\mathbf{x} \in \widetilde{\mathcal{N}}_\varepsilon} \frac{1}{m} \sum_{i=1}^m \mathbf{x}^\top (\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{A}_i - \mathbf{X}_*) \mathbf{x} \\ &= 2 \sup_{\mathbf{x} \in \widetilde{\mathcal{N}}_\varepsilon} \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{x}^\top \mathbf{A}_i \mathbf{x} - \mathbf{x}^\top \mathbf{X}_* \mathbf{x}). \end{aligned}$$

806 For each $i \in [m]$, we have that $\mathbb{E} [\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{x}^\top \mathbf{A}_i \mathbf{x}] = \mathbf{x}^\top \mathbf{X}_* \mathbf{x}$. Moreover, the inner product $\langle \mathbf{A}_i, \mathbf{X}_* \rangle$
807 is a centered Gaussian random variable with variance $\|\mathbf{X}_*\|_F^2$ and $\mathbf{x}^\top \mathbf{A}_i \mathbf{x}$ is a centered Gaussian ran-
808 dom variable with variance 1. Thus, for each fixed \mathbf{x} , $\sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{x}^\top \mathbf{A}_i \mathbf{x} - \mathbf{x}^\top \mathbf{X}_* \mathbf{x})$ is a sum of m
809 independent and centered sub-exponential random variables with subexponential norm bounded
810 by $K \|\mathbf{X}_*\|_F$, where K is an absolute constant (see [60, Lemma 2.7.7]). Therefore, by Bernstein’s
811 inequality (see, for example, [60, Theorem 2.8.1]), it holds that

$$\mathbb{P} \left(\left| \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X}_* \rangle \mathbf{x}^\top \mathbf{A}_i \mathbf{x} - \mathbf{x}^\top \mathbf{X}_* \mathbf{x}) \right| \geq t \right) \leq \exp \left(-C' \min \left\{ \frac{mt^2}{\|\mathbf{X}_*\|_F^2}, \frac{mt}{\|\mathbf{X}_*\|_F} \right\} \right),$$

812 where $C' > 0$ is some absolute constant. Taking $t = \frac{1}{8} C \|\mathbf{X}_*\|_F \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right)$ and a union bound
813 over all points \mathbf{x} on $\widetilde{\mathcal{N}}_\varepsilon$, we obtain

$$\left\| (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - \mathbf{X}_* \right\| \leq \frac{1}{4} C \|\mathbf{X}_*\|_F \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right) \leq \frac{1}{4} C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{r} \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right) \quad (71)$$

814 with probability at least $1 - \exp(d \log(6) - C' C^2 d) \geq 1 - \exp(-4d)$ for some sufficiently large constant
815 $C > 0$.

816 We assume that (71) holds and that $m > C^2 \kappa^2 r d$. Then Weyl's inequalities imply that

$$\lambda_r((\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)) > \frac{1}{2} \sigma_{\min}(\mathbf{X}_*), \quad |\lambda_{r+1}((\mathcal{A}^* \mathcal{A})(\mathbf{X}_*))| < \frac{1}{2} \sigma_{\min}(\mathbf{X}_*).$$

817 Since $\tilde{\Lambda}_r$ is a diagonal matrix with entries $\lambda_1((\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)), \dots, \lambda_r((\mathcal{A}^* \mathcal{A})(\mathbf{X}_*))$, it follows from the
818 definition of $\mathbf{U}_0 = \tilde{\mathbf{V}}_r \tilde{\Lambda}_r^{1/2}$ that $\mathbf{U}_0 \mathbf{U}_0^\top$ is the best rank- r approximation of $(\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)$. Conse-
819 quently, we obtain that

$$\begin{aligned} \|\mathbf{X}_* - \mathbf{U}_0 \mathbf{U}_0^\top\| &\leq \|\mathbf{X}_* - (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)\| + \|(\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - \mathbf{U}_0 \mathbf{U}_0^\top\| \\ &\leq \|\mathbf{X}_* - (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*)\| + \|(\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - \mathbf{X}_*\| \leq C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{\frac{rd}{m}}, \end{aligned}$$

820 where in the second inequality, we used the Eckart-Young-Mirsky theorem.

821 (2) Due to Lemma 3.5 we have

$$(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_*) = (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)) - \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)) \rangle \mathbf{w}\mathbf{w}^\top. \quad (72)$$

822 It follows that

$$\|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_*)\| \leq \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*))\| + |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)) \rangle|. \quad (73)$$

823 For a fixed $\mathbf{w} \in \mathcal{N}_{\varepsilon}$, we obtain with an analogous argument as for (71) that with probability at least
824 $1 - \exp(-4d)$,

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*))\| \leq C \|\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)\|_F \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right) \leq \frac{1}{4} C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{r} \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right).$$

825 The second term in (73) can be rewritten as

$$\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)) \rangle = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*) \rangle.$$

826 Here, $\sum_{i=1}^m \langle \mathbf{w}\mathbf{w}^\top, \mathbf{A}_i \rangle \langle \mathbf{A}_i, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*) \rangle$ is a sum of m independent sub-exponential random vari-
827 ables with mean zero due to the rotation invariance of the Gaussian measure. Moreover, each term
828 has sub-exponential norm $K \|\mathbf{X}_*\|_F$. Applying Bernstein's inequality as in the proof of (71), we
829 obtain that for each fixed \mathbf{w} with probability at least $1 - \exp(-4d)$,

$$\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_*)) \rangle \leq \frac{1}{4} C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{r} \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right). \quad (74)$$

830 Then, by taking a union bound over $\mathbf{w} \in \mathcal{N}_{\varepsilon}$, it follows from (73) that with probability at least
831 $1 - \exp(-2d)$ that for all $\mathbf{w} \in \mathcal{N}_{\varepsilon}$ it holds that

$$\|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_*)\| \leq \frac{1}{2} C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{r} \left(\sqrt{\frac{d}{m}} + \frac{d}{m} \right). \quad (75)$$

832 We now assume that (75) holds and that $m > 4C^2 \kappa^2 r d$. Then it follows from Weyl's inequalities
833 that

$$\lambda_r((\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*)) > \frac{1}{2} \sigma_{\min}(\mathbf{X}_*), \quad |\lambda_{r+1}((\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*))| < \frac{1}{2} \sigma_{\min}(\mathbf{X}_*).$$

834 It follows from the Eckart-Mirsky-Young theorem and the definition of $\mathbf{U}_{0, \mathbf{w}}$ that $\mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top$ is the
835 best rank- r approximation of $(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*)$. Therefore,

$$\begin{aligned} \|\mathbf{X}_* - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top\| &\leq \|\mathbf{X}_* - (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*)\| + \|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*) - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top\| \\ &\leq 2 \|\mathbf{X}_* - (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_*)\| \leq 2C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{\frac{rd}{m}}. \end{aligned}$$

836 This finishes the proof of inequality (24). Finally, (25) follows from (23) and (24) via the triangle
837 inequality.

838 (3) From (72), we have

$$\begin{aligned} (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) - (\mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_*) &= (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_*) - (\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I})(\mathbf{X}_*) \\ &= \langle \mathbf{w} \mathbf{w}^\top, \mathbf{X}_* \rangle (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{w} \mathbf{w}^\top) + \langle \mathcal{A}(\mathbf{w} \mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp})(\mathbf{X}_*) \rangle \mathbf{w} \mathbf{w}^\top. \end{aligned}$$

839 It follows from Lemma 2.2 that there exists an absolute constant $C_1 > 0$ such that for any $\alpha \in (0, 1)$
840 and $m \geq \frac{C_1}{\alpha^2} \kappa^2 r d$, with probability at least $1 - \exp(-d)$, the measurement operator \mathcal{A} satisfies the
841 Restricted Isometry Property of order $6r$ with constant

$$\delta := \delta_{6r} \leq \frac{\alpha}{\kappa}. \quad (76)$$

842 Then for any $\mathbf{V} \in \mathbb{R}^{d \times r}$ with orthonormal columns and for all $\mathbf{w} \in \mathcal{N}_\varepsilon$, when $m \geq \frac{C_1}{\alpha^2} \kappa^2 r d$, with
843 probability at least $1 - 2 \exp(-d)$,

$$\begin{aligned} & \|(\mathcal{A}^* \mathcal{A} - \mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_*) \mathbf{V}\|_F \\ & \leq |\langle \mathbf{w} \mathbf{w}^\top, \mathbf{X}_* \rangle| \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{w} \mathbf{w}^\top) \mathbf{V}\|_F + |\langle \mathcal{A}(\mathbf{w} \mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp})(\mathbf{X}_*) \rangle| \|\mathbf{w} \mathbf{w}^\top \mathbf{V}\|_F \\ & \stackrel{(a)}{\leq} \delta \|\mathbf{X}_*\| \|\mathbf{w} \mathbf{w}^\top\|_F + |\langle \mathcal{A}(\mathbf{w} \mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w} \mathbf{w}^\top, \perp})(\mathbf{X}_*) \rangle| \\ & \stackrel{(b)}{\leq} \alpha \sigma_{\min}(\mathbf{X}_*) + \frac{1}{2} C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{\frac{r d}{m}}. \end{aligned} \quad (77)$$

844 Here in (a) we use property (7) in Lemma 2.4 and the fact that $\mathbf{w} \mathbf{w}^\top \mathbf{V}$ is of rank 1, and in (b) we
845 use (76) and, moreover, (74) with a union bound over $\mathbf{w} \in \mathcal{N}_\varepsilon$.

846 We now proceed under the assumption that the inequalities in parts (1) and (2) hold. We use the
847 following notations for spectral initialization:

$$\begin{aligned} (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*) &= \tilde{\mathbf{V}} \tilde{\Lambda} \tilde{\mathbf{V}}^\top, \quad \mathbf{U}_0 = \tilde{\mathbf{V}}_r \tilde{\Lambda}_r^{1/2}, \\ (\mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_*) &= \tilde{\mathbf{V}}_w \tilde{\Lambda}_w \tilde{\mathbf{V}}_w^\top, \quad \mathbf{U}_{0, \mathbf{w}} = \tilde{\mathbf{V}}_{r, \mathbf{w}} \tilde{\Lambda}_{r, \mathbf{w}}^{1/2}. \end{aligned} \quad (78)$$

848 Denote

$$\mathbf{Z}_1 := (\mathcal{A}^* \mathcal{A})(\mathbf{X}_*), \quad \mathbf{Z}_2 := (\mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_*),$$

849 and

$$\mathbf{Z}_{1, r} := \mathbf{U}_0 \mathbf{U}_0^\top, \quad \mathbf{Z}_{2, r} := \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top.$$

850 Recall the definition of $\tilde{\mathbf{V}}_r$ and $\tilde{\mathbf{V}}_{r, \mathbf{w}}$ in (78) and (17). We have

$$\begin{aligned} \|\mathbf{Z}_{1, r} - \mathbf{Z}_{2, r}\|_F &= \|\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top\|_F \\ &\leq \|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top) \tilde{\mathbf{V}}_r\|_F + \|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top) \tilde{\mathbf{V}}_{r, \perp}\|_F. \end{aligned} \quad (79)$$

851 For the first term in (79), we have

$$\begin{aligned} & \|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0, \mathbf{w}} \mathbf{U}_{0, \mathbf{w}}^\top) \tilde{\mathbf{V}}_r\|_F \\ &= \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F \\ &\leq \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F + \|(\mathbf{Z}_2 - \mathbf{Z}_{2, r}) \tilde{\mathbf{V}}_r\|_F \\ &= \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F + \|(\tilde{\mathbf{V}}_{r, \mathbf{w}, \perp} \tilde{\Lambda}_{r, \mathbf{w}, \perp} \tilde{\mathbf{V}}_{r, \mathbf{w}, \perp}^\top) \tilde{\mathbf{V}}_r\|_F \\ &\leq \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F + \sigma_{r+1}(\mathbf{Z}_2) \|\tilde{\mathbf{V}}_{r, \mathbf{w}, \perp}^\top \tilde{\mathbf{V}}_r\|_F \\ &\leq \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F + C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{\frac{r d}{m}} \|\tilde{\mathbf{V}}_{r, \mathbf{w}, \perp}^\top \tilde{\mathbf{V}}_r\|_F, \end{aligned} \quad (80)$$

852 where in the last inequality we used Weyl's inequality and (75), which implies

$$\sigma_{r+1}(\mathbf{Z}_2) = |\sigma_{r+1}(\mathbf{Z}_2) - \sigma_{r+1}(\mathbf{X}_*)| \leq \|\mathbf{Z}_2 - \mathbf{X}_*\| \leq C \kappa \sigma_{\min}(\mathbf{X}_*) \sqrt{\frac{r d}{m}} \|\tilde{\mathbf{V}}_{r, \mathbf{w}, \perp}^\top \tilde{\mathbf{V}}_r\|_F. \quad (81)$$

853 From (75) and (71), it follows that when $m \geq C^2 \kappa^2 r d$,

$$\|\mathbf{Z}_1 - \mathbf{Z}_2\| \leq \frac{3C}{2} \kappa \sigma_{\min}(\mathbf{X}_\star) \sqrt{\frac{rd}{m}}. \quad (82)$$

854 Similar to (81), using (75) and Weyl's inequalities we obtain that

$$\begin{aligned} |\sigma_r(\mathbf{Z}_1) - \sigma_{\min}(\mathbf{X}_\star)| &\leq C \kappa \sigma_{\min}(\mathbf{X}_\star) \sqrt{\frac{rd}{m}}, \\ \sigma_{r+1}(\mathbf{Z}_1) &\leq C \kappa \sigma_{\min}(\mathbf{X}_\star) \sqrt{\frac{rd}{m}}. \end{aligned}$$

855 Therefore, if $m > 16C^2 \kappa^2 r d$, the spectral gap between $\sigma_r(\mathbf{Z}_1)$ and $\sigma_{r+1}(\mathbf{Z}_2)$ can be bounded from
856 below by

$$\sigma_r(\mathbf{Z}_1) - \sigma_{r+1}(\mathbf{Z}_1) \geq \left(1 - 2C \kappa \sqrt{\frac{rd}{m}}\right) \sigma_{\min}(\mathbf{X}_\star) \geq \frac{1}{2} \sigma_{\min}(\mathbf{X}_\star). \quad (83)$$

857 When $m \geq 51C^2 \kappa^2 r d$, we have from (82) and (83),

$$\begin{aligned} \|\mathbf{Z}_1 - \mathbf{Z}_2\| &\leq \frac{3C}{2} \kappa \sqrt{\frac{rd}{m}} \sigma_{\min}(\mathbf{X}_\star) \\ &\leq \left(1 - \frac{1}{\sqrt{2}}\right) \left(1 - 2C \kappa \sqrt{\frac{rd}{m}}\right) \sigma_{\min}(\mathbf{X}_\star) \\ &\leq \left(1 - \frac{1}{\sqrt{2}}\right) (\sigma_r(\mathbf{Z}_1) - \sigma_{r+1}(\mathbf{Z}_1)). \end{aligned}$$

858 Thus, the prerequisites of Lemma 2.6 (Davis-Kahan inequality) are satisfied. It follows that when
859 $m \geq 51C^2 \kappa^2 r d$,

$$\|\tilde{\mathbf{V}}_{r,\mathbf{w},\perp}^\top \tilde{\mathbf{V}}_r\|_F \leq \frac{2\sqrt{2} \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F}{\sigma_{\min}(\mathbf{X}_\star)}. \quad (84)$$

860 Hence, when $m \geq (51C^2 + \frac{C_1}{\alpha^2}) \kappa^2 r d$, we obtain from (80) and (77) that

$$\begin{aligned} \|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top) \tilde{\mathbf{V}}_r\|_F &\leq \left(1 + 2\sqrt{2} C \kappa \sqrt{\frac{rd}{m}}\right) \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F \\ &\leq 2 \|(\mathbf{Z}_1 - \mathbf{Z}_2) \tilde{\mathbf{V}}_r\|_F \leq \left(2\alpha + C \kappa \sqrt{\frac{rd}{m}}\right) \sigma_{\min}(\mathbf{X}_\star). \end{aligned} \quad (85)$$

861 For the second term in (79), we have when $m \geq (51C^2 + \frac{C_1}{\alpha^2}) \kappa^2 r d$,

$$\begin{aligned} &\|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top) \tilde{\mathbf{V}}_{r,\perp}\|_F \\ &\leq \|\tilde{\mathbf{V}}_r^\top (\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top) \tilde{\mathbf{V}}_{r,\perp}\|_F + \|\tilde{\mathbf{V}}_{r,\perp}^\top (\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top) \tilde{\mathbf{V}}_{r,\perp}\|_F \\ &\leq \|\tilde{\mathbf{V}}_r^\top (\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top)\|_F + \|\tilde{\mathbf{V}}_{r,\perp}^\top \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\ &\leq \left(2\alpha + C \kappa \sqrt{\frac{rd}{m}}\right) \sigma_{\min}(\mathbf{X}_\star) + \|\tilde{\mathbf{V}}_{r,\perp}^\top \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F, \end{aligned} \quad (86)$$

862 where the last inequality is due to (85).

863 We now consider the second term in (86). Recall the definition of $\mathbf{U}_{0,\mathbf{w}}$ in (18). We have for $m \geq$
 864 $(51C^2 + \frac{C_1}{\alpha^2}) \kappa^2 rd$,

$$\begin{aligned}
 \|\tilde{\mathbf{V}}_{r,\perp}^\top \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F &= \|\tilde{\mathbf{V}}_{r,\perp}^\top \tilde{\mathbf{V}}_{r,\mathbf{w}} \mathbf{\Lambda}_{r,\mathbf{w}} \tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &\leq \|\tilde{\mathbf{V}}_{r,\perp}^\top \tilde{\mathbf{V}}_{r,\mathbf{w}} \mathbf{\Lambda}_{r,\mathbf{w}}\| \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &= \sqrt{\|\tilde{\mathbf{V}}_{r,\perp}^\top \tilde{\mathbf{V}}_{r,\mathbf{w}} \mathbf{\Lambda}_{r,\mathbf{w}}^2 \tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|} \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &= \sqrt{\|\tilde{\mathbf{V}}_{r,\perp}^\top (\mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top)^2 \tilde{\mathbf{V}}_{r,\perp}\|} \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &= \|\tilde{\mathbf{V}}_{r,\perp}^\top \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top\| \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &= \|\tilde{\mathbf{V}}_{r,\perp}^\top (\mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top - \mathbf{U}_0 \mathbf{U}_0^\top)\| \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &\leq \|\mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top - \mathbf{U}_0 \mathbf{U}_0^\top\| \|\tilde{\mathbf{V}}_{r,\mathbf{w}}^\top \tilde{\mathbf{V}}_{r,\perp}\|_F \\
 &\stackrel{(a)}{\leq} 3C\kappa\sigma_{\min}(\mathbf{X}_\star) \sqrt{\frac{rd}{m}} \cdot \frac{2\sqrt{2}\|(\mathbf{Z}_1 - \mathbf{Z}_2)\tilde{\mathbf{V}}_r\|_F}{\sigma_{\min}(\mathbf{X}_\star)} \\
 &\stackrel{(b)}{\leq} 6\sqrt{2}C\kappa \left(\alpha + \frac{1}{2}C\kappa\sqrt{\frac{rd}{m}} \right) \sqrt{\frac{rd}{m}} \sigma_{\min}(\mathbf{X}_\star), \tag{87}
 \end{aligned}$$

865 where (a) is due to (25) and (84), and (b) is due to (77). Therefore from (86) and (87), we obtain
 866 for $m \geq (51C^2 + \frac{C_1}{\alpha^2}) \kappa^2 rd$,

$$\begin{aligned}
 &\|(\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top) \tilde{\mathbf{V}}_{r,\perp}\|_F \tag{88} \\
 &\leq \left(2\alpha + C\kappa\sqrt{\frac{rd}{m}} \right) \sigma_{\min}(\mathbf{X}_\star) + 6\sqrt{2}C\kappa \left(\alpha + \frac{1}{2}C\kappa\sqrt{\frac{rd}{m}} \right) \sqrt{\frac{rd}{m}} \sigma_{\min}(\mathbf{X}_\star).
 \end{aligned}$$

867 From (85), (88), and (79), we conclude that if $m \geq (51C^2 + \frac{C_1}{\alpha^2}) \kappa^2 rd$,

$$\|\mathbf{U}_0 \mathbf{U}_0^\top - \mathbf{U}_{0,\mathbf{w}} \mathbf{U}_{0,\mathbf{w}}^\top\|_F \leq \left(2\alpha + C\kappa\sqrt{\frac{rd}{m}} \right) \left(2\sigma_{\min}(\mathbf{X}_\star) + 3\sqrt{2}C\kappa\sqrt{\frac{rd}{m}} \sigma_{\min}(\mathbf{X}_\star) \right).$$

868 This finishes the proof of (26). \square

869 B. Proofs of lemmas concerning the distance between the virtual 870 sequences and the original sequence

871 B.1. Some auxiliary estimates

872 In order to prove Lemma 4.3 and Lemma 4.5 we will need several auxiliary estimates. These are
 873 summarized in the following lemma.

874 **Lemma B.1.** *Assume that the measurement operator \mathcal{A} has the Restricted Isometry Property with constant*
 875 *$\delta = \delta_{4r+1} \leq 1$. Moreover, assume that the conclusion of Lemma 3.6 holds. Then, the following inequalities*
 876 *hold.*

1.

$$\begin{aligned}
 &\|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}}) (\mathbf{X}_\star - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\
 &\leq \left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}} \right) \|\mathbf{X}_\star - \mathbf{U}_t \mathbf{U}_t^\top\| + \left(\delta + \frac{4\sqrt{2d}}{\sqrt{m}} \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F, \tag{89}
 \end{aligned}$$

2.

$$\|[(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I}) (\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \leq 2\delta \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F, \tag{90}$$

3.

$$\|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \leq \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \|\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \quad (91)$$

877 4. and

$$\begin{aligned} \|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\| &\leq \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\quad + \left(2\delta + 4\sqrt{\frac{2d}{m}} \right) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F. \end{aligned} \quad (92)$$

878 *Proof of Lemma B.1.* To prove inequality (89), we compute that

$$\begin{aligned} (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) &= (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) + (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \\ &\stackrel{(a)}{=} (\mathcal{A}^* \mathcal{A})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) + \mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \\ &\quad - \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle \mathbf{w}\mathbf{w}^\top, \end{aligned}$$

879 where in equation (a) we used Lemma 3.5. It follows that

$$\begin{aligned} (\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) &= (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \\ &\quad + \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle \mathbf{w}\mathbf{w}^\top. \end{aligned}$$

880 By using the triangle inequality, we obtain the estimate

$$\begin{aligned} &\|(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\leq \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F + |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle \mathbf{w}\mathbf{w}^\top|_F \\ &\stackrel{(a)}{\leq} \delta \|\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\|_F + |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle| \\ &\stackrel{(b)}{\leq} \delta \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)) \rangle| \\ &\quad + |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)) \rangle| \\ &\stackrel{(c)}{\leq} \delta \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{4\sqrt{d}}{\sqrt{m}} \|\mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top))\|_2 + \delta \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(d)}{\leq} \delta \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{4\sqrt{2d}}{\sqrt{m}} \|\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F + \delta \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(e)}{\leq} \delta \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{4\sqrt{2d}}{\sqrt{m}} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|_F + \left(\delta + \frac{4\sqrt{2d}}{\sqrt{m}} \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(f)}{\leq} \left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}} \right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \left(\delta + \frac{4\sqrt{2d}}{\sqrt{m}} \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F. \end{aligned}$$

881 Inequality (a) follows from the RIP-assumption combined with Lemma 2.4 and from the fact that
 882 $\|\mathbf{w}\|_2 = 1$. Inequality (b) is a consequence of the fact that $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top}$ is a rank-one projection and of the
 883 triangle inequality. In inequality (c), we used that the conclusion of Lemma 3.6 holds and Lemma
 884 2.4. In inequality (d), we used the RIP of rank $2r+1$. Inequality (e) is due to the fact that $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}$ is an
 885 orthogonal projection and due to the triangle inequality. In inequality (f), we used that $\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top$
 886 has rank at most $2r$. This proves inequality (89).

887 To prove inequality (90) we compute first that

$$\begin{aligned} &(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \\ &= (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) - \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle \mathbf{w}\mathbf{w}^\top. \end{aligned}$$

888 It follows that

$$\begin{aligned}
& \left\| [(\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I})(\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,w}} \right\|_F \\
& \stackrel{(a)}{\leq} \delta \left\| \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \right\|_F + \left| \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top)) \rangle \right| \\
& \stackrel{(b)}{\leq} 2\delta \left\| \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \right\|_F \\
& \leq 2\delta \left\| \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F.
\end{aligned}$$

889 In inequalities (a) and (b) we used Lemma 2.4. This proves inequality (90).

890 Next, we prove the third inequality. For that, we observe that using Lemma 3.5 it holds that

$$\begin{aligned}
(\mathcal{A}^* \mathcal{A} - \mathcal{A}_w^* \mathcal{A}_w) (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) &= (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathcal{P}_{\mathbf{w}\mathbf{w}^\top} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \\
&\quad + \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \rangle \mathbf{w}\mathbf{w}^\top.
\end{aligned}$$

891 Then it follows that

$$\begin{aligned}
& \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{A}_w^* \mathcal{A}_w) (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{U}_{t,w}} \right\|_F \\
& \leq \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathcal{P}_{\mathbf{w}\mathbf{w}^\top} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \mathbf{V}_{\mathbf{U}_{t,w}} \right\|_F + \left| \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \rangle \right| \\
& \stackrel{(a)}{\leq} \delta \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\| + 4\sqrt{\frac{d}{m}} \left\| \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \right\|_2 \\
& \stackrel{(b)}{\leq} \delta \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\| + 4\sqrt{\frac{2d}{m}} \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\|_F \\
& \leq \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\|,
\end{aligned}$$

892 where inequality (a) holds due to Lemma 2.4, since $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}$ is a rank-one projection, and since we
893 assumed that the conclusion of Lemma 3.6 holds. Inequality (b) is again due to Lemma 2.4 and since
894 $\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}$ is an orthogonal projection. This proves inequality (91).

895 It remains to prove inequality (92). We note that it holds that

$$\begin{aligned}
& (\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \\
&= (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) - \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \rangle \mathbf{w}\mathbf{w}^\top,
\end{aligned}$$

896 where in the last line we applied Lemma 3.5. It follows from the triangle inequality that

$$\begin{aligned}
& \left\| (\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \right\| \\
& \leq \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathcal{P}_{\mathbf{w}\mathbf{w}^\top} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \right\| \\
& \quad + \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \right\| + \left| \langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)) \rangle \right| \\
& \stackrel{(a)}{\leq} \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + \delta \left\| \mathcal{P}_{\mathbf{w}\mathbf{w}^\top} (\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \right\|_F + \delta \left\| \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F \\
& \quad + 4\sqrt{\frac{2d}{m}} \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\|_F \\
& \leq \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + \delta \left\| \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \right\| \\
& \quad + 2\delta \left\| \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F + 4\sqrt{\frac{2d}{m}} \left\| \mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top \right\|_F \\
& \leq \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \left\| \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \right\| \\
& \quad + \left(2\delta + 4\sqrt{\frac{2d}{m}} \right) \left\| \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F.
\end{aligned}$$

897 In inequality (a) we applied Lemma 2.4 and that the conclusion of Lemma 3.6 holds. This proves
898 inequality (92). Thus, the proof of Lemma B.1 is complete. \square

899 **B.2. Proof of Lemma 4.3**

900 *Proof of Lemma 4.3.* We define the shorthand notation

$$\begin{aligned}\mathbf{M}_t &:= (\mathcal{A}^* \mathcal{A}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top), \\ \mathbf{M}_{t,\mathbf{w}} &:= (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}}) (\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top).\end{aligned}$$

901 It follows that

$$\begin{aligned}\mathbf{U}_{t+1} &= (\mathbf{Id} + \mu \mathbf{M}_t) \mathbf{U}_t, \\ \mathbf{U}_{t+1,\mathbf{w}} &= (\mathbf{Id} + \mu \mathbf{M}_{t,\mathbf{w}}) \mathbf{U}_{t,\mathbf{w}}.\end{aligned}$$

902 We compute that

$$\begin{aligned}\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t+1,\mathbf{w}}^\top &= (\mathbf{Id} + \mu \mathbf{M}_t) \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{Id} + \mu \mathbf{M}_t) - (\mathbf{Id} + \mu \mathbf{M}_{t,\mathbf{w}}) \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top (\mathbf{Id} + \mu \mathbf{M}_{t,\mathbf{w}}) \\ &= \underbrace{\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top}_{=: (i)} + \underbrace{\mu \mathbf{M}_t (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)}_{=: (ii)} + \underbrace{\mu (\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}) \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top}_{=: (ii)} \\ &\quad + \underbrace{\mu (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{M}_t}_{=: (iii)} + \underbrace{\mu \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top (\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})}_{=: (iv)} \\ &\quad + \underbrace{\mu^2 (\mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top \mathbf{M}_{t,\mathbf{w}})}_{=: (v)}.\end{aligned}$$

903 We want to estimate the spectral norm of these terms individually. Before that, we note that

$$\begin{aligned}\|\mathbf{M}_t\| &\stackrel{(a)}{\leq} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|(\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ &\stackrel{(b)}{\leq} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + c_1 \sigma_{\min}(\mathbf{X}_*)\end{aligned}\tag{93}$$

$$\stackrel{(c)}{\leq} 2\sigma_{\min}(\mathbf{X}_*).\tag{94}$$

904 Inequality (a) follows from the triangle inequality and inequality (b) follows from assumption (29).

905 Inequality (c) is a consequence of assumption (30). Moreover, we note that

$$\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}} = (\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) - (\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}}) (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top).$$

906 It follows that

$$\begin{aligned}&\|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\leq \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F + \|[(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I}) (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\quad + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(a)}{\leq} \left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}} \right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \left(3\delta + \frac{4\sqrt{2d}}{\sqrt{m}} + 1 \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(b)}{\leq} \frac{2c_3}{\kappa} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \left(\frac{4c_3}{\kappa} + 1 \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F,\end{aligned}\tag{95}$$

907 where in inequality (a) we used inequalities (89) and (90) from Lemma B.1. Inequality (b) is due
908 to assumption (32). Note that it also follows from these estimates that

$$\begin{aligned}\|\mathbf{M}_{t,\mathbf{w}} \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\| &\leq \|\mathbf{M}_t\| + \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(a)}{\leq} 2\sigma_{\min}(\mathbf{X}_*) + \frac{2c_3}{\kappa} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \left(\frac{4c_3}{\kappa} + 1 \right) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(b)}{\leq} 3\sigma_{\min}(\mathbf{X}_*),\end{aligned}\tag{96}$$

909 where inequality (a) follows from (95). Inequality (b) is a consequence of the assumptions (30) and
910 (31) (and by choosing the absolute constant $c_3 > 0$ small enough).

911 Now we are in a position to estimate the spectral norms of the terms (i)-(v).

912 * Estimating term (i): We compute that that

$$\begin{aligned} \|\mathbf{M}_t(\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)\|_F &\leq \|\mathbf{M}_t\| \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\stackrel{(93)}{\leq} (\|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + c_1\sigma_{\min}(\mathbf{X}_\star)) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F. \end{aligned}$$

913 * Estimating term (ii): We compute that

$$\begin{aligned} \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F &\leq \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \|\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\| \\ &\leq \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F (\|\mathbf{U}_t\mathbf{U}_t^\top\| + \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|) \\ &\leq 3\|\mathbf{X}_\star\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F, \end{aligned}$$

914 where in the last inequality we used assumptions (28) and (30).

915 * Estimating term (iii): With the same argument as for term (i) we observe that

$$\|(\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)\mathbf{M}_t\|_F \leq (\|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + c_1\sigma_{\min}(\mathbf{X}_\star)) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F.$$

916 * Estimating term (iv): With the same argument as for term (ii) we compute that

$$\|\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\|_F \leq 3\|\mathbf{X}_\star\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F.$$

917 * Estimating term (v): First, we compute that

$$\begin{aligned} \mathbf{M}_t\mathbf{U}_t\mathbf{U}_t^\top\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\mathbf{M}_{t,\mathbf{w}} &= \mathbf{M}_t(\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)\mathbf{M}_t + (\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\mathbf{M}_t \\ &\quad + \mathbf{M}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}). \end{aligned}$$

918 It follows that

$$\begin{aligned} &\|\mathbf{M}_t\mathbf{U}_t\mathbf{U}_t^\top\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\mathbf{M}_{t,\mathbf{w}}\|_F \\ &\leq \|\mathbf{M}_t\|^2 \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + \left(\|\mathbf{U}_t\|^2 + \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|\right) \|\mathbf{M}_t\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\quad + \|\mathbf{M}_{t,\mathbf{w}}\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\| \left(\|\mathbf{U}_t\|^2 + \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|\right) \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(a)}{\leq} \|\mathbf{M}_t\|^2 \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 3\|\mathbf{X}_\star\| \|\mathbf{M}_t\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\quad + 3\|\mathbf{X}_\star\| \|\mathbf{M}_{t,\mathbf{w}}\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(b)}{\leq} 4\sigma_{\min}^2(\mathbf{X}_\star) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 15\sigma_{\min}(\mathbf{X}_\star) \|\mathbf{X}_\star\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F. \end{aligned}$$

919 For inequality (a) we used the assumptions (28) and (31). Inequality (b) is a consequence of in-
920 equalities (94) and (96).

921 * Conclusion: By summing up all terms we obtain that

$$\begin{aligned} &\|\mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,\mathbf{w}}\mathbf{U}_{t+1,\mathbf{w}}^\top\|_F \\ &\leq \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 2\mu (\|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + c_1\sigma_{\min}(\mathbf{X}_\star)) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\quad + 6\mu \|\mathbf{X}_\star\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\quad + \mu^2 (4\sigma_{\min}^2(\mathbf{X}_\star) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 15\sigma_{\min}(\mathbf{X}_\star) \|\mathbf{X}_\star\| \|(\mathbf{M}_t - \mathbf{M}_{t,\mathbf{w}})\mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F) \\ &\stackrel{(a)}{\leq} (1 + 2\mu \|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + 2c_1\sigma_{\min}(\mathbf{X}_\star)) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\quad + 12\mu\sigma_{\min}(\mathbf{X}_\star)c_3 \|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + 6\mu \|\mathbf{X}_\star\| \left(\frac{4c_3}{\kappa} + 1\right) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &\quad + 4\mu^2\sigma_{\min}^2(\mathbf{X}_\star) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 30c_3\mu^2\sigma_{\min}^2(\mathbf{X}_\star) \|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| \\ &\quad + 60c_3\mu^2\sigma_{\min}^2(\mathbf{X}_\star) \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + 15\mu^2\sigma_{\min}(\mathbf{X}_\star) \|\mathbf{X}_\star\| \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F \\ &= (1 + 2\mu \|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| + (2c_1 + 24c_3)\mu\sigma_{\min}(\mathbf{X}_\star) + 6\mu \|\mathbf{X}_\star\| + 4\mu^2\sigma_{\min}^2(\mathbf{X}_\star) + 60c_3\mu^2\sigma_{\min}^2(\mathbf{X}_\star)) \\ &\quad \cdot \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\|_F + (12c_3\mu\sigma_{\min}(\mathbf{X}_\star) + 30c_3\mu^2\sigma_{\min}^2(\mathbf{X}_\star)) \|\mathbf{X}_\star - \mathbf{U}_t\mathbf{U}_t^\top\| \\ &\stackrel{(b)}{\leq} \frac{\sqrt{\sqrt{2}-1}}{40} \sigma_{\min}(\mathbf{X}_\star). \end{aligned}$$

922 Inequality (a) follows from inequality (95). Inequality (b) is due to assumptions (30), (31), and the
 923 assumption $\mu \leq \frac{c_2}{\kappa \|\mathbf{X}_*\|}$ for a sufficiently small absolute constant $c_2 > 0$. This completes the proof
 924 of Lemma 4.3. \square

925 B.3. Proof of Lemma 4.4

926 *Proof of Lemma 4.4.* Let $\mathbf{R} \in \mathbb{R}^{r \times r}$ be an orthogonal matrix. We compute that

$$\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top = \mathbf{U}_t \mathbf{R} (\mathbf{U}_t \mathbf{R})^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top = \mathbf{U}_t \mathbf{R} (\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w})^\top - (\mathbf{U}_{t,w} - \mathbf{U}_t \mathbf{R}) \mathbf{U}_{t,w}^\top.$$

927 It follows that

$$\begin{aligned} & \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{X}_*, \perp}\|_F \\ & \leq \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t \mathbf{R}\| \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F + \|\mathbf{U}_{t,w} - \mathbf{U}_t \mathbf{R}\|_F \|\mathbf{U}_{t,w}^\top \mathbf{V}_{\mathbf{X}_*, \perp}\| \\ & \leq (\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t \mathbf{R}\| + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_{t,w}\|) \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F \\ & \leq (2\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t\| + \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|) \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F \\ & = \left(2\sqrt{\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*, \perp}\|} + \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|\right) \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F \\ & = \left(2\sqrt{\|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_*, \perp}\|} + \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|\right) \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F \\ & \stackrel{(a)}{\leq} \left(\frac{1}{20} \sqrt{\sigma_{\min}(\mathbf{X}_*)} + \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F\right) \|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F. \end{aligned} \quad (97)$$

928 In inequality (a) we used Assumption (33). By choosing the orthogonal matrix \mathbf{R} as the minimizer
 929 of Procruste's problem, i.e., such that $\|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F$ is minimal, we obtain by Lemma 4.8 that

$$\|\mathbf{U}_t \mathbf{R} - \mathbf{U}_{t,w}\|_F \leq \frac{\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F}{\sqrt{2(\sqrt{2}-1)\sigma_{\min}^2(\mathbf{U}_t)}} \stackrel{(a)}{\leq} \frac{\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F}{\sqrt{(\sqrt{2}-1)\frac{3}{2}\sigma_{\min}(\mathbf{X}_*)}} \stackrel{(b)}{\leq} \frac{\sqrt{\sigma_{\min}(\mathbf{X}_*)}}{20}.$$

930 Inequality (a) follows from Assumption (33) and Weyl's inequalities for singular values. For in-
 931 equality (b) we used Assumption (34). Inequality (97) combined with this inequality chain yields
 932 that

$$\begin{aligned} \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{X}_*, \perp}\|_F & \leq \frac{\sqrt{\sigma_{\min}(\mathbf{X}_*)}}{10} \cdot \frac{\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F}{\sqrt{(\sqrt{2}-1)\frac{3}{2}\sigma_{\min}(\mathbf{X}_*)}} \\ & \leq \frac{\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F}{5}. \end{aligned} \quad (98)$$

933 In order to proceed we note that

$$\begin{aligned} \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F & \leq \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{X}_*}\|_F \\ & \quad + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{X}_*, \perp}\|_F \\ & \leq 2\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F + \|\mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) \mathbf{V}_{\mathbf{X}_*, \perp}\|_F \\ & \stackrel{(a)}{\leq} 2\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F + \frac{1}{5}\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F. \end{aligned}$$

934 In inequality (a) we have used inequality (98). By rearranging terms we obtain that

$$\begin{aligned} \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F & \leq \frac{2}{1-\frac{1}{5}} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F \\ & \leq 3\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F. \end{aligned}$$

935 This shows inequality (36). Then (35) follows directly from inserting the above inequality into
 936 (98). \square

937 **B.4. Proof of Lemma 4.5**

938 The key idea in the proof of Lemma 4.5 is to decompose $\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,\mathbf{w}}\mathbf{U}_{t+1,\mathbf{w}}^\top)$ into a
939 sum of the form

$$\begin{aligned} & \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,\mathbf{w}}\mathbf{U}_{t+1,\mathbf{w}}^\top) \\ &= \mathbf{V}_{\mathbf{X}_*}^\top (1 + \mu (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) (\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top) (1 + \mu (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) \\ & \quad + \mathbf{V}_{\mathbf{X}_*}^\top \Delta. \end{aligned} \tag{99}$$

940 The first summand can be interpreted as a contraction mapping applied to the matrix
941 $\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top$ and thus can be expected to have a smaller Frobenius norm than
942 $\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)\|_F$. In contrast, the term Δ , which will be determined explicitly in the
943 proof of Lemma 4.5, can be interpreted as an additive error term which, as we will show, has rela-
944 tively small Frobenius norm.

945 To deal with the first summand we need the following auxiliary lemma.

946 **Lemma B.2.** Denote by $\lambda_{\max}(\mathbf{A})$ the largest eigenvalue of a symmetric matrix \mathbf{A} and by $\lambda_{\min}(\mathbf{A})$ the smallest
947 eigenvalue of \mathbf{A} . Assume that the assumptions of Lemma 4.5 are satisfied. Then it holds that

$$\lambda_{\min}(\mathbf{Id} + \mu (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) \geq 0, \tag{100}$$

$$\lambda_{\max}(\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*}) \leq -\frac{\sigma_{\min}(\mathbf{X}_*)}{2}, \tag{101}$$

$$\|\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)\| \leq 1 + \frac{\mu\sigma_{\min}(\mathbf{X}_*)}{128}. \tag{102}$$

948 *Proof of Lemma B.2.* Note that the assumptions $\mu \leq \frac{c_4}{\kappa\|\mathbf{X}_*\|}$, (38), and (40) together with Weyl's
949 inequalities imply

$$\begin{aligned} & \lambda_{\min}(\mathbf{Id} + \mu (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) \\ &= \lambda_{\min}(\mathbf{Id} + \mu ((\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) - \mathbf{U}_t\mathbf{U}_t^\top + \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top)) \\ &\geq 1 - \mu\|\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top\| - \mu\|\mathbf{U}_t\mathbf{U}_t^\top\| - \mu\|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\| \\ &\geq 0. \end{aligned}$$

950 for sufficiently small $c_2, c_3, c_4 > 0$. This shows inequality (100).

951 We observe that

$$\begin{aligned} & \lambda_{\max}(\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*}) \\ &\stackrel{(a)}{\leq} \lambda_{\max}(-\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*}) + \|\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top\| + \|\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top\| \\ &\stackrel{(b)}{\leq} \lambda_{\max}(-\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*}) + (c_2 + c_3)\sigma_{\min}(\mathbf{X}_*) \\ &= -\lambda_{\min}(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t} \mathbf{V}_{\mathbf{U}_t}^\top \mathbf{U}_t\mathbf{U}_t^\top \mathbf{V}_{\mathbf{U}_t} \mathbf{V}_{\mathbf{U}_t}^\top \mathbf{V}_{\mathbf{X}_*}) + (c_2 + c_3)\sigma_{\min}(\mathbf{X}_*) \\ &\leq -\sigma_{\min}(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})^2 \lambda_{\min}(\mathbf{U}_t\mathbf{U}_t^\top) + (c_2 + c_3)\sigma_{\min}(\mathbf{X}_*) \\ &\stackrel{(c)}{\leq} -\frac{\sigma_{\min}(\mathbf{X}_*)}{2}. \end{aligned}$$

952 Inequality (a) follows from Weyl's inequalities. Inequality (b) follows from assumption (39) and
953 (40). For inequality (c) we used assumptions (37), (39) for sufficiently small c_1, c_2, c_3 , and Weyl's
954 inequalities. This proves inequality (101).

955 To prove inequality (102), we first establish an upper bound for the largest eigenvalue of $\mathbf{X}_* -$
956 $\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}}\mathbf{U}_{t,\mathbf{w}}^\top$. For that let $\mathbf{x} \in \mathbb{R}^d$ be arbitrary. We use the orthogonal decomposition $\mathbf{x} =$

957 $\mathbf{x}_{\parallel} + \mathbf{x}_{\perp}$, where \mathbf{x}_{\parallel} is the orthogonal projection of \mathbf{x} onto the column span of \mathbf{X}_{\star} . We compute that

$$\begin{aligned}
& \mathbf{x}^{\top} (\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x} \\
&= \mathbf{x}_{\parallel}^{\top} (\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x}_{\parallel} - \mathbf{x}_{\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x}_{\perp} - 2\mathbf{x}_{\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x}_{\parallel} \\
&\stackrel{(101)}{\leq} -\frac{\sigma_{\min}(\mathbf{X}_{\star})}{2} \|\mathbf{x}_{\parallel}\|_2^2 - 2\mathbf{x}_{\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x}_{\parallel}. \tag{103}
\end{aligned}$$

958 Next, we observe that

$$\begin{aligned}
-\mathbf{x}_{\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x}_{\parallel} &\leq \|\mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{V}_{\mathbf{X}_{\star}}\| \|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2 \\
&\leq (2\|\mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top} \mathbf{U}_t \mathbf{U}_t^{\top} \mathbf{V}_{\mathbf{X}_{\star}}\| + \|\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}\|) \|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2 \\
&= (2\|\mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{X}_{\star}) \mathbf{V}_{\mathbf{X}_{\star}}\| + \|\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}\|) \|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2 \\
&\leq (2\|\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top}\| + \|\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}\|) \|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2 \\
&\leq \frac{\sigma_{\min}(\mathbf{X}_{\star}) \|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2}{16}.
\end{aligned}$$

959 In the last inequality we have used the assumptions (39) and (40) for sufficiently small $c_2, c_3 > 0$.

960 Combining this estimate with (103) we obtain that

$$\begin{aligned}
\mathbf{x}^{\top} (\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{x} &\leq \sigma_{\min}(\mathbf{X}_{\star}) \left(\frac{\|\mathbf{x}_{\parallel}\|_2 \|\mathbf{x}_{\perp}\|_2}{8} - \frac{\|\mathbf{x}_{\parallel}\|_2^2}{2} \right) \\
&\leq \frac{\sigma_{\min}(\mathbf{X}_{\star}) \|\mathbf{x}_{\perp}\|_2^2}{128} \leq \frac{\sigma_{\min}(\mathbf{X}_{\star}) \|\mathbf{x}\|_2^2}{128}.
\end{aligned}$$

961 This implies that

$$\lambda_{\max}(\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) \leq 1 + \frac{\mu\sigma_{\min}(\mathbf{X}_{\star})}{128}.$$

962 This inequality, together with inequality (100), yields inequality (102). Thus, the proof of Lemma
963 B.2 is complete. \square

964 With Lemma B.2 in place, we can show that the first term in the decomposition (99) indeed has a
965 smaller Frobenius norm than the term $\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})$.

966 **Lemma B.3.** *Assume that the assumptions of Lemma 4.5 are satisfied. Then, it holds that*

$$\begin{aligned}
& \|\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}))\|_F \\
&\leq \left(1 - \frac{\mu\sigma_{\min}(\mathbf{X}_{\star})}{8}\right) \|\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})\|_F.
\end{aligned}$$

967 *Proof of Lemma B.3.* We first compute that

$$\begin{aligned}
& \|\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}))\|_F \\
&\leq \|\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})\|_F \|\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})\|_F \\
&\leq \left(1 + \frac{\mu\sigma_{\min}(\mathbf{X}_{\star})}{128}\right) \|\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})\|_F, \tag{104}
\end{aligned}$$

968 where in the last line we used inequality (102) from Lemma B.2. In order to proceed, we consider
969 the decomposition

$$\begin{aligned}
& \mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \\
& \underbrace{\mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{Id} + \mu(\mathbf{X}_{\star} - \mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})) \mathbf{V}_{\mathbf{X}_{\star}} \mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top})}_{=:\mathbf{N}_1} \\
& - \underbrace{\mu \mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{V}_{\mathbf{X}_{\star},\perp} \mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{V}_{\mathbf{X}_{\star}} \mathbf{V}_{\mathbf{X}_{\star}}^{\top}}_{=:\mathbf{N}_2} \\
& - \underbrace{\mu \mathbf{V}_{\mathbf{X}_{\star}}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} + \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{V}_{\mathbf{X}_{\star},\perp} \mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top} (\mathbf{U}_t \mathbf{U}_t^{\top} - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^{\top}) \mathbf{V}_{\mathbf{X}_{\star},\perp} \mathbf{V}_{\mathbf{X}_{\star},\perp}^{\top}}_{=:\mathbf{N}_3}.
\end{aligned}$$

970 We estimate the Frobenius norm of the three terms individually. For the first term we obtain that

$$\begin{aligned}
\|\mathbf{N}_1\|_F &\leq \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)) \mathbf{V}_{\mathbf{X}_*}\| \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&= \|\mathbf{Id} + \mu \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*}\| \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\stackrel{(a)}{\leq} (1 + \mu \lambda_{\max}(\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*})) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\stackrel{(b)}{\leq} \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{2}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F,
\end{aligned}$$

971 where in inequality (a) we have used (100) and in (b) we have used inequality (101) from Lemma
972 B.2. The Frobenius norm of the term \mathbf{N}_2 can be estimated by

$$\begin{aligned}
\|\mathbf{N}_2\|_F &\leq \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top + \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_{*,\perp}}\| \|\mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_*}\|_F \\
&= (\|\mathbf{V}_{\mathbf{X}_{*,\perp}}^\top [2(\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) + (\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)]\|) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\leq (2\|\mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*)\| + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\leq (2c_2 \sigma_{\min}(\mathbf{X}_*) + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\leq (2c_2 + c_3) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F,
\end{aligned}$$

973 where we have used Assumptions (39) and (40). With similar arguments, we can estimate the
974 Frobenius norm of the term \mathbf{N}_3 by

$$\|\mathbf{N}_3\|_F \leq (2c_2 + c_3) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_{*,\perp}}\|_F.$$

975 By using Lemma 4.4 we obtain that

$$\|\mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{X}_{*,\perp}}\|_F \leq \frac{3\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F}{5}.$$

976 It follows that

$$\|\mathbf{N}_3\|_F \leq \frac{3(2c_2 + c_3) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F}{5}.$$

977 By summing up our estimates for $\|\mathbf{N}_1\|_F$, $\|\mathbf{N}_2\|_F$, and $\|\mathbf{N}_3\|_F$ and choosing the constants $c_1, c_2 > 0$
978 small enough we obtain that

$$\begin{aligned}
&\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)) (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{4}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F.
\end{aligned}$$

979 Inserting this estimate into (104) yields that

$$\begin{aligned}
&\|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)) (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top))\|_F \\
&\leq \left(1 + \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{128}\right) \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{4}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F \\
&\leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{8}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\|_F,
\end{aligned}$$

980 where in the last line, we used our assumption on the step size μ . This completes the proof of Lemma
981 B.3. \square

982 With the auxiliary estimates in Lemma B.3 we can give a proof of Lemma 4.5.

983 *Proof of Lemma 4.5.* First, we compute that

$$\begin{aligned}
\mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top &= (\mathbf{Id} + \mu [(\mathcal{A}^*\mathcal{A})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)]) \mathbf{U}_t\mathbf{U}_t^\top (\mathbf{Id} + \mu [(\mathcal{A}^*\mathcal{A})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)]) \\
&= (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) \mathbf{U}_t\mathbf{U}_t^\top (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) \\
&\quad + \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top + \mu\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top \\
&\quad + \mu^2\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) + \mu^2(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top \\
&\quad - \mu^2\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top \\
&\quad + \mu [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \mathbf{U}_t\mathbf{U}_t^\top (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_t\mathbf{U}_t^\top) \\
&\quad + \mu (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_t\mathbf{U}_t^\top) \mathbf{U}_t\mathbf{U}_t^\top [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \\
&\quad + \mu^2 [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \mathbf{U}_t\mathbf{U}_t^\top [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)].
\end{aligned}$$

984 Analogously, we can compute that

$$\begin{aligned}
&\mathbf{U}_{t+1,w}\mathbf{U}_{t+1,w}^\top \\
&= (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) \\
&\quad + \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top + \mu\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top \\
&\quad + \mu^2\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) + \mu^2(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top \\
&\quad - \mu^2\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top \\
&\quad + \mu [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) \\
&\quad + \mu (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)] \\
&\quad + \mu^2 [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)].
\end{aligned}$$

985 Thus, we obtain that

$$\begin{aligned}
&\mathbf{U}_{t+1}\mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1,w}\mathbf{U}_{t+1,w}^\top \\
&= \mathbf{M}_1 + \mu^2\mathbf{M}_2 + \mu^2\mathbf{M}_3 + \mu^2\mathbf{M}_4 + \mu^2\mathbf{M}_4 + \mu\mathbf{M}_5 + \mu\mathbf{M}_6 + \mu^2\mathbf{M}_7,
\end{aligned} \tag{105}$$

986 where

$$\begin{aligned}
\mathbf{M}_1 &:= (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) (\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) (\mathbf{Id} + \mu(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)) \\
\mathbf{M}_2 &:= \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) - \mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top), \\
\mathbf{M}_3 &:= (\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) \mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top - (\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top, \\
\mathbf{M}_4 &:= \mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top, \\
\mathbf{M}_5 &:= [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \mathbf{U}_t\mathbf{U}_t^\top (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_t\mathbf{U}_t^\top) \\
&\quad - [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top), \\
\mathbf{M}_6 &:= (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_t\mathbf{U}_t^\top) \mathbf{U}_t\mathbf{U}_t^\top [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \\
&\quad - (\mathbf{Id} + \mu\mathbf{X}_* - \mu\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)], \\
\mathbf{M}_7 &:= [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \mathbf{U}_t\mathbf{U}_t^\top [(\mathcal{A}^*\mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top)] \\
&\quad - [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top [(\mathcal{A}_w^*\mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)].
\end{aligned}$$

987 Recall that Lemma B.3 shows that

$$\|\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{M}_1\|_F \leq \left(1 - \frac{\mu\sigma_{\min}(\mathbf{X}_*)}{8}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)\|_F.$$

988 To complete the proof, we need to derive upper bounds for $\|\mathbf{M}_i\|_F$, where $i = 2, 3, \dots, 7$.

989

990 **Estimating $\|\mathbf{M}_2\|_F$:** We compute that

$$\begin{aligned}
\mathbf{M}_2 &= \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top\mathbf{U}_t\mathbf{U}_t^\top(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) - \mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top(\mathbf{X}_* - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top) \\
&= (\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top - \mathbf{U}_t\mathbf{U}_t^\top) \mathbf{U}_t\mathbf{U}_t^\top(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) + \mathbf{U}_t\mathbf{U}_t^\top(\mathbf{U}_t\mathbf{U}_t^\top - \mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top)(\mathbf{X}_* - \mathbf{U}_t\mathbf{U}_t^\top) \\
&\quad + \mathbf{U}_t\mathbf{U}_t^\top\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top(\mathbf{U}_{t,w}\mathbf{U}_{t,w}^\top - \mathbf{U}_t\mathbf{U}_t^\top).
\end{aligned}$$

991 Thus, we obtain that

$$\begin{aligned}
& \|M_2\|_F \\
& \leq 2\|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F \|U_tU_t^\top\| \|X_\star - U_tU_t^\top\| + \|U_tU_t^\top\| \|U_{t,w}U_{t,w}^\top\| \|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F \\
& \leq 2\|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F \|U_tU_t^\top\| \|X_\star - U_tU_t^\top\| \\
& \quad + \|U_tU_t^\top\| (\|U_tU_t^\top\| + \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\|) \|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F \\
& \leq 5\|X_\star\|^2 \|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F.
\end{aligned}$$

992 In the last inequality we used assumptions (38), (39), and (40) for sufficiently small $c_2, c_3 > 0$.

993

994 **Estimating $\|M_3\|_F$:** Since $M_3 = M_2^\top$ it follows that

$$\|M_3\|_F \leq 5\|X_\star\|^2 \|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F.$$

995 **Estimating $\|M_4\|_F$:** We compute that

$$\begin{aligned}
M_4 &= (U_tU_t^\top - U_{t,w}U_{t,w}^\top) U_{t,w}U_{t,w}^\top U_tU_t^\top + U_{t,w}U_{t,w}^\top (U_{t,w}U_{t,w}^\top - U_tU_t^\top) U_tU_t^\top \\
& \quad + U_{t,w}U_{t,w}^\top U_tU_t^\top (U_tU_t^\top - U_{t,w}U_{t,w}^\top).
\end{aligned}$$

996 Again, using the assumptions (38) and (40), and the triangle inequality we obtain that

$$\|M_4\|_F \leq 20\|X_\star\|^2 \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\|_F.$$

997 **Estimating $\|M_5\|_F$:** We compute

$$\begin{aligned}
M_5 &= \underbrace{[(A^*A - \mathcal{I})(X_\star - U_tU_t^\top)] (U_tU_t^\top - U_{t,w}U_{t,w}^\top) (\text{Id} + \mu X_\star - \mu U_tU_t^\top)}_{=:O_1} \\
& \quad + \underbrace{\mu[(A^*A - \mathcal{I})(X_\star - U_tU_t^\top)] U_{t,w}U_{t,w}^\top (U_{t,w}U_{t,w}^\top - U_tU_t^\top)}_{=:O_2} \\
& \quad + \underbrace{[(A^*A - A_w^*A_w)(X_\star - U_tU_t^\top)] U_{t,w}U_{t,w}^\top (\text{Id} + \mu X_\star - \mu U_{t,w}U_{t,w}^\top)}_{=:O_3} \\
& \quad + \underbrace{[(A_w^*A_w - \mathcal{I})(U_{t,w}U_{t,w}^\top - U_tU_t^\top)] U_{t,w}U_{t,w}^\top (\text{Id} + \mu X_\star - \mu U_{t,w}U_{t,w}^\top)}_{=:O_4}.
\end{aligned}$$

998 We estimate the Frobenius norm of these summands individually. For the first term we observe that

$$\begin{aligned}
\|O_1\|_F &\leq \|(A^*A - \mathcal{I})(X_\star - U_tU_t^\top)\| \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\|_F (1 + \mu\|X_\star\| + \mu\|U_{t,w}U_{t,w}^\top\|) \\
&\stackrel{(a)}{\leq} 2\|(A^*A - \mathcal{I})(X_\star - U_tU_t^\top)\| \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\|_F \\
&\stackrel{(b)}{\leq} 2c_5\sigma_{\min}(X_\star) \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\|_F,
\end{aligned}$$

999 where in inequality (a) we have used assumptions (38), (40), and the assumption on the step size

1000 μ . In inequality (b) we have used assumption (41).

1001 Using again assumptions (38), (40), and (41) we obtain that

$$\|O_2\|_F \leq 3c_5\sigma_{\min}(X_\star) \|X_\star\| \|U_{t,w}U_{t,w}^\top - U_tU_t^\top\|_F.$$

1002 For the term $\|O_3\|_F$ we obtain that

$$\begin{aligned}
\|O_3\|_F &\leq \|(A^*A - A_w^*A_w)(X_\star - U_tU_t^\top)\| V_{U_{t,w}} \|U_{t,w}U_{t,w}^\top\| (1 + \mu\|X_\star\| + \mu\|U_{t,w}U_{t,w}^\top\|) \\
&\leq \|(A^*A - A_w^*A_w)(X_\star - U_tU_t^\top)\| V_{U_{t,w}} (\|U_tU_t^\top - U_{t,w}U_{t,w}^\top\| + \|U_tU_t^\top\|) \\
&\quad (1 + \mu\|X_\star\| + \mu\|U_{t,w}U_{t,w}^\top - U_tU_t^\top\| + \mu\|U_tU_t^\top\|) \\
&\stackrel{(a)}{\leq} 4\|(A^*A - A_w^*A_w)(X_\star - U_tU_t^\top)\| V_{U_{t,w}} \|X_\star\| \\
&\stackrel{(b)}{\leq} 4\left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}}\right) \|X_\star - U_tU_t^\top\| \|X_\star\| + 4\left(\delta + \frac{8\sqrt{2d}}{\sqrt{m}}\right) \|U_tU_t^\top - U_{t,w}U_{t,w}^\top\| \|X_\star\|.
\end{aligned}$$

1003 Inequality (a) follows from the assumptions (38) and (40), and the assumption on the step size μ .
 1004 In inequality (b) we used the estimate (89) from Lemma B.1.

1005 For the term $\|\mathbf{O}_4\|_F$ we obtain that

$$\begin{aligned} \|\mathbf{O}_4\|_F &\leq \|(\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I})(\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{U}_{t,w}}\|_F (\|\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_t \mathbf{U}_t^\top\|) \\ &\quad \cdot (1 + \mu \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \mu \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|) \\ &\stackrel{(a)}{\leq} 3 \|\mathbf{X}_*\| \|[(\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I})(\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,w}}\|_F \\ &\stackrel{(b)}{\leq} 6\delta \|\mathbf{X}_*\| \|\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F. \end{aligned}$$

1006 Inequality (a) follows from assumptions (39) and (40), and the assumption on the step size μ . In-
 1007 equality (b) is due to inequality (90) in Lemma B.1. By summing up all terms we obtain that

$$\begin{aligned} \|\mathbf{M}_5\|_F &\leq \|\mathbf{O}_1\|_F + \mu \|\mathbf{O}_2\|_F + \|\mathbf{O}_3\|_F + \|\mathbf{O}_4\|_F \\ &\leq 2c_5 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F + 3\mu c_5 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\quad + 4 \left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}} \right) \|\mathbf{X}_*\| \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + 4 \left(\delta + \frac{4\sqrt{2d}}{\sqrt{m}} \right) \|\mathbf{X}_*\| \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F \\ &\quad + 6\delta \|\mathbf{X}_*\| \|\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &= \left[((2 + 3\mu) c_5 + 6\kappa\delta) \sigma_{\min}(\mathbf{X}_*) + 4 \left(\delta + \frac{4\sqrt{2d}}{\sqrt{m}} \right) \|\mathbf{X}_*\| \right] \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F \\ &\quad + 4 \left(\delta + \frac{8\sqrt{rd}}{\sqrt{m}} \right) \|\mathbf{X}_*\| \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\stackrel{(a)}{\leq} (((2 + 3\mu) c_5 + 6c_6) \sigma_{\min}(\mathbf{X}_*) + 8c_6 \sigma_{\min}(\mathbf{X}_*)) \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F + 8c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\stackrel{(b)}{\leq} \frac{\sigma_{\min}(\mathbf{X}_*)}{100} \cdot \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top\|_F + 8c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\stackrel{(c)}{\leq} \frac{3\sigma_{\min}(\mathbf{X}_*)}{100} \cdot \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F + 8c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|, \end{aligned}$$

1008 where in inequality (a) we used the assumption (42). Inequality (b) follows from choosing the
 1009 constants c_5 and c_6 small enough. To obtain inequality (c) we applied Lemma 4.4.

1010

1011 **Estimating $\|\mathbf{M}_6\|_F$:**

1012 Since $\mathbf{M}_6 = \mathbf{M}_5^\top$ we obtain that

$$\|\mathbf{M}_6\|_F \leq \frac{3\sigma_{\min}(\mathbf{X}_*)}{100} \cdot \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)\|_F + 8c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|.$$

1013 **Estimating $\|\mathbf{M}_7\|_F$:** To deal with the term \mathbf{M}_7 we first compute that

$$\begin{aligned} \mathbf{M}_7 &= \underbrace{[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top) [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)]}_{=: \mathbf{L}_1} \\ &\quad + \underbrace{[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)]}_{=: \mathbf{L}_2} \\ &\quad + \underbrace{[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top - \mathbf{U}_t \mathbf{U}_t^\top)]}_{=: \mathbf{L}_3} \\ &\quad + \underbrace{[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)]}_{=: \mathbf{L}_4} \\ &\quad + \underbrace{[(\mathcal{A}_w^* \mathcal{A}_w - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)] \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top [(\mathcal{A}^* \mathcal{A} - \mathcal{A}_w^* \mathcal{A}_w)(\mathbf{X}_* - \mathbf{U}_{t,w} \mathbf{U}_{t,w}^\top)]}_{=: \mathbf{L}_5}. \end{aligned}$$

1014 We estimate the Frobenius norm of the summands individually. For $\|\mathbf{L}_1\|_F$ we obtain that

$$\begin{aligned} \|\mathbf{L}_1\|_F &\leq \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \| \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ &\leq c_5^2 \sigma_{\min}(\mathbf{X}_*)^2 \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F, \end{aligned}$$

1015 where we have used assumption (41). Next, we note that

$$\begin{aligned} \|\mathbf{L}_2\|_F &\leq \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F (\|\mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|) \\ &\quad \cdot \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \\ &\stackrel{(a)}{\leq} 3c_5 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(b)}{\leq} 3c_5 \delta \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\stackrel{(c)}{\leq} 3c_5 c_6 \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F. \end{aligned}$$

1016 Inequality (a) follows from assumptions (38), (40), and (41). Inequality (b) is due to Lemma 2.4
1017 and inequality (c) is due to assumption (42). In order to estimate $\|\mathbf{L}_3\|_F$ we note that

$$\begin{aligned} \|\mathbf{L}_3\|_F &(\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + \|[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top) \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F] \\ &\quad \cdot (\|\mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F) \|[(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(a)}{\leq} (c_5 \sigma_{\min}(\mathbf{X}_*) + \delta \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|_F) (2\|\mathbf{X}_*\| + c_3 \sigma_{\min}(\mathbf{X}_*)) \delta \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\stackrel{(b)}{\leq} 3(c_5 + \delta c_3) \delta \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\ &\stackrel{(c)}{\leq} 3c_6 (c_5 + \delta c_3) \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F. \end{aligned}$$

1018 In inequality (a) we used the assumptions (38), (40), (41), and Lemma 2.4. Inequality (b) follows
1019 from assumption (40) and since the constant $c_3 > 0$ is chosen small enough. Inequality (c) is due
1020 to assumption (42).

1021 Next, we can estimate $\|\mathbf{L}_4\|_F$ by

$$\begin{aligned} \|\mathbf{L}_4\|_F &\leq \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F (\|\mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\|) \\ &\quad \cdot (\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|) \\ &\stackrel{(a)}{\leq} \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F (2\|\mathbf{X}_*\| + c_3 \sigma_{\min}(\mathbf{X}_*)) \\ &\quad \cdot (c_5 \sigma_{\min}(\mathbf{X}_*) + \delta \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F) \\ &\stackrel{(b)}{\leq} 3(c_5 + c_3 \delta) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F \\ &\stackrel{(c)}{\leq} 3(c_5 + c_3 \delta) \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\| \\ &\leq 3(c_5 + c_3 \delta) \left(\delta + 8\sqrt{\frac{rd}{m}} \right) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|) \\ &\stackrel{(d)}{\leq} 6c_6 (c_5 + c_3 \delta) \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|). \end{aligned}$$

1022 In inequality (a) we used assumptions (38), (40), and (41) as well as Lemma 2.4. Inequality (b) uses
1023 assumption (40). Inequality (c) follows from inequality (91) in Lemma B.1. Inequality (d) is due to
1024 assumption (42).

1025 The norm $\|\mathbf{L}_5\|_F$ can be estimated by

$$\begin{aligned} \|\mathbf{L}_5\|_F &\leq \|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\| \| \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\| \| \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}^\top [(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \|_F \\ &\stackrel{(a)}{\leq} 3\|\mathbf{X}_*\| \|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\| \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F. \end{aligned} \tag{106}$$

1026 In inequality (a) we used the triangle inequality and the assumptions (38), (40). In order to proceed,
 1027 we note first that

$$\begin{aligned}
 & \|(\mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)\| \\
 & \stackrel{(a)}{\leq} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + \left(\delta + 8\sqrt{\frac{rd}{m}}\right) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\
 & \quad + \left(2\delta + 4\sqrt{\frac{2d}{m}}\right) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \stackrel{(b)}{\leq} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| + \frac{2c_6}{\kappa} \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{3c_6}{\kappa} \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \stackrel{(c)}{\leq} \left(c_5 + \frac{2c_2c_6}{\kappa} + \frac{3c_3c_6}{\kappa}\right) \sigma_{\min}(\mathbf{X}_*),
 \end{aligned}$$

1028 where in inequality (a) we used Lemma B.1. Inequality (b) follows from the assumptions (42).
 1029 Inequality (c) is due to assumption (39), (40), and (41). Moreover, it holds that

$$\begin{aligned}
 \|[(\mathcal{A}^* \mathcal{A} - \mathcal{A}_{\mathbf{w}}^* \mathcal{A}_{\mathbf{w}})(\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top)] \mathbf{V}_{\mathbf{U}_{t,\mathbf{w}}}\|_F & \stackrel{(a)}{\leq} \left(\delta + 8\sqrt{\frac{rd}{m}}\right) \|\mathbf{X}_* - \mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top\| \\
 & \stackrel{(b)}{\leq} \frac{2c_6}{\kappa} (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|).
 \end{aligned}$$

1030 Inequality (a) follows from inequality (91) in Lemma B.1. Inequality (b) is due to assumption (42).
 1031 Inserting the last two inequality chains into inequality (106) we obtain that

$$\|\mathbf{L}_5\|_F \leq 6c_6 \left(c_5 + \frac{2c_2c_6}{\kappa} + \frac{3c_3c_6}{\kappa}\right) \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|)$$

1032 By summing up all terms $\|\mathbf{L}_i\|_F$ for $i = 1, \dots, 5$ it follows that

$$\begin{aligned}
 \|\mathbf{M}_7\|_F & \leq c_5^2 \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \quad + 3c_5c_6 \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \quad + 3c_6 (c_5 + c_3\delta) \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F \\
 & \quad + 6c_6 (c_5 + c_3\delta) \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|) \\
 & \quad + 6c_6 \left(c_5 + \frac{2c_2c_6}{\kappa} + \frac{3c_3c_6}{\kappa}\right) \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|) \\
 & \leq \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t,\mathbf{w}} \mathbf{U}_{t,\mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F),
 \end{aligned}$$

1033 where the last inequality holds since the absolute constants $c_3, c_5, c_6 > 0$ are chosen small enough.

1034 Using the decomposition (105), the triangle inequality, combined with our estimates for
 1035 $\|\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{M}_1\|_F$ and for $\|\mathbf{M}_i\|_F$, where $2 \leq i \leq 7$, we obtain that

$$\begin{aligned}
& \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top - \mathbf{U}_{t+1, \mathbf{w}} \mathbf{U}_{t+1, \mathbf{w}}^\top)\|_F \\
& \leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{8}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + 30\mu^2 \|\mathbf{X}_*\|^2 \|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top\|_F \\
& \quad + \frac{3\mu \sigma_{\min}(\mathbf{X}_*)}{50} \cdot \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + 16\mu c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\
& \quad + \mu^2 \sigma_{\min}^2(\mathbf{X}_*) (\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|\mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top\|_F) \\
& \stackrel{(a)}{\leq} \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{8}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + 90\mu c_4 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F \\
& \quad + \frac{3\mu \sigma_{\min}(\mathbf{X}_*)}{50} \cdot \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + 16\mu c_6 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\
& \quad + \mu^2 \sigma_{\min}^2(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \frac{3\mu c_4 \sigma_{\min}(\mathbf{X}_*)}{\kappa} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top - \mathbf{U}_t \mathbf{U}_t^\top)\|_F \\
& \stackrel{(b)}{\leq} \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + \mu (16c_6 + \mu \sigma_{\min}(\mathbf{X}_*)) \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\
& \leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16}\right) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_{t, \mathbf{w}} \mathbf{U}_{t, \mathbf{w}}^\top)\|_F + \mu \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\|,
\end{aligned}$$

1036 where inequality (a) is due to Lemma 4.4 and the assumption on the step size μ . Inequality (b) is
 1037 obtained by choosing $c_4 < 1/2$, and the last inequality is obtained by choosing $c_6 < \frac{1}{32}$. \square

1038 C. Proof of the lemmas controlling the distance between \mathbf{X}_* and 1039 $\mathbf{U}_t \mathbf{U}_t^\top$ (Lemma 4.6, Lemma 4.7, and Lemma 4.9)

1040 C.1. Proof of Lemma 4.6

1041 *Proof of Lemma 4.6.* We first note that

$$\begin{aligned}
\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}} &= \mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t} \mathbf{V}_{\mathbf{U}_t}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}} \\
&= \mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t} (\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})^{-1} \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t} \mathbf{V}_{\mathbf{U}_t}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}} \\
&= \mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t} (\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})^{-1} \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}} \\
&= \mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t} (\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})^{-1} \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}.
\end{aligned}$$

1042 Using the submultiplicativity property of the $\|\cdot\|$ -norm it follows that

$$\begin{aligned}
\|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}}\| &\leq \|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t}\| \|(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})^{-1}\| \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}\| \\
&= \frac{\|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t}\|}{\sigma_{\min}(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t})} \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}\|.
\end{aligned}$$

1043 Recall that

$$\sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t}) = 1 - \|\mathbf{V}_{\mathbf{U}_t}^\top \mathbf{V}_{\mathbf{X}_{*, \perp}} \mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t}\| = 1 - \|\mathbf{V}_{\mathbf{U}_t}^\top \mathbf{V}_{\mathbf{X}_{*, \perp}}\|^2 \geq \frac{1}{4},$$

1044 where in the last inequality, we used assumption (44). It follows that

$$\|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}}\| \leq 2 \|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t}\| \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}\|.$$

1045 This proves inequality (45). To prove inequality (46) we note that

$$\begin{aligned}
\|\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*\| &\leq \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*)\| + \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}\| \\
&\quad + \|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*, \perp}}\| \\
&\leq 2 \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*)\| + \|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*, \perp}}\| \\
&\leq 2(1 + \|\mathbf{V}_{\mathbf{X}_{*, \perp}}^\top \mathbf{V}_{\mathbf{U}_t}\|) \|\mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*)\|,
\end{aligned}$$

1046 where in the last inequality we used (45). This completes the proof of Lemma 4.6. \square

1047 **C.2. Proof of Lemma 4.7**

1048 *Proof of Lemma 4.7.* We define the shorthand notation

$$\mathbf{M}_t := (\mathcal{A}^* \mathcal{A}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) = \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top + \underbrace{(\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)}_{=:\mathbf{E}_t}.$$

1049 Thus, we have that

$$\mathbf{U}_{t+1} = (\mathbf{Id} + \mu \mathbf{M}_t) \mathbf{U}_t.$$

1050 We compute that

$$\begin{aligned} & \mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top \\ &= \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t - \mu^2 \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \\ &= \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \mu (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) - \mu \mathbf{E}_t \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_t \mathbf{U}_t^\top \mathbf{E}_t \\ & \quad - \mu^2 \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \\ &= (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) - \mu^2 \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top \\ & \quad - \mu \mathbf{E}_t \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top \mathbf{E}_t - \mu^2 \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t. \end{aligned}$$

1051 It follows that

$$\begin{aligned} & \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top) \\ &= \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_*} \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \\ & \quad + \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{V}_{\mathbf{X}_*, \perp} \mathbf{V}_{\mathbf{X}_*, \perp}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \\ & \quad - \mu^2 \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{E}_t \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{E}_t - \mu^2 \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t \\ &= \underbrace{(\mathbf{Id} - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*}) \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top)}_{=:(I)} \\ & \quad + \underbrace{\mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*, \perp} \mathbf{V}_{\mathbf{X}_*, \perp}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top)}_{=:(II)} \\ & \quad - \underbrace{(\mu^2 \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top + \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{E}_t \mathbf{U}_t \mathbf{U}_t^\top + \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{E}_t + \mu^2 \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t)}_{=:(III)}. \end{aligned}$$

1052 We estimate the spectral norm of these terms individually.

1053 * Estimating term (I): We obtain that

$$\begin{aligned} & \left| \left| (\mathbf{Id} - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*}) \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \\ & \stackrel{(a)}{\leq} \left| \left| \mathbf{Id} - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*} \right| \right| \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \left| \left| \mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top \right| \right| \\ & \stackrel{(b)}{\leq} \left| \left| \mathbf{Id} - \mu \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_*} \right| \right| \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \\ & \stackrel{(c)}{=} (1 - \mu \sigma_{\min}^2(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t)) \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \\ & \leq \left(1 - \mu (\sigma_{\min}(\mathbf{V}_{\mathbf{X}_*}^\top \mathbf{V}_{\mathbf{U}_t}) \sigma_{\min}(\mathbf{U}_t))^2 \right) \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \\ & \stackrel{(d)}{\leq} \left(1 - \frac{\mu}{2} \sigma_{\min}^2(\mathbf{U}_t) \right) \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right| \\ & \stackrel{(e)}{\leq} \left(1 - \frac{\mu}{4} \sigma_{\min}(\mathbf{X}_*) \right) \left| \left| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right| \right|. \end{aligned}$$

1054 Inequality (a) is due to the submultiplicativity of the $\|\cdot\|$ -norm. In inequality (b) and equality (c)

1055 we used the assumptions $\|\mathbf{U}_t\| \leq \sqrt{2} \|\mathbf{X}_*\|$ and $\mu \leq \frac{1}{1024\kappa \|\mathbf{X}_*\|}$. In inequality (d) we used assump-

1056 tion (47). Inequality (e) follows from assumption (48), which, due to Weyl's inequality, implies

1057 $\sigma_{\min}^2(\mathbf{U}_t) \geq \frac{1}{2} \sigma_{\min}(\mathbf{X}_*)$.

1058 * Estimating term (II): We note that

$$\begin{aligned}
& \left\| \mathbf{V}_{\mathbf{X}_*}^\top \mathbf{U}_t \mathbf{U}_t^\top \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \right\| \\
&= \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*,\perp}} \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) \right\| \\
&\stackrel{(a)}{\leq} \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \left\| \mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top \right\| \\
&\stackrel{(b)}{\leq} \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \\
&\leq \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \left(\left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_*} \right\| + \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*,\perp}} \right\| \right) \\
&\leq \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \left(\left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| + \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \mathbf{V}_{\mathbf{X}_{*,\perp}} \right\| \right) \\
&\stackrel{(c)}{\leq} 2 \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| (1 + \left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top \mathbf{V}_{\mathbf{U}_t} \right\|) \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\| \\
&\stackrel{(d)}{\leq} 3 \left\| \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_* \right\| \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\|.
\end{aligned}$$

1059 In inequality (a) we used the submultiplicativity of the $\|\cdot\|$ -norm. Inequality (b) follows from
1060 the assumption $\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}$ and $\mu \leq \frac{1}{1024\kappa\|\mathbf{X}_*\|}$. In inequality (c), we used Lemma 4.6.

1061 In inequality (d) we used the assumption $\left\| \mathbf{V}_{\mathbf{X}_{*,\perp}}^\top \mathbf{V}_{\mathbf{U}_t} \right\| \leq \frac{1}{2}$. Thus, by using the assumption
1062 $\|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \leq \frac{\sigma_{\min}(\mathbf{X}_*)}{48}$ it follows that

$$\|(II)\| \leq \frac{\sigma_{\min}(\mathbf{X}_*)}{16} \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{U}_t \mathbf{U}_t^\top - \mathbf{X}_*) \right\|.$$

1063 * Estimating term (III): We first note that

$$\begin{aligned}
\|\mathbf{M}_t \mathbf{V}_{\mathbf{U}_t}\| &\stackrel{(a)}{\leq} \left\| \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \right\| + \left\| [(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)] \mathbf{V}_{\mathbf{U}_t} \right\| \\
&\stackrel{(b)}{\leq} 4 \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\|, \tag{107}
\end{aligned}$$

1064 where (a) follows from the triangle inequality and (b) follows from Lemma 4.6. Moreover, we have
1065 that

$$\|\mathbf{M}_t\| \leq \left\| \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \right\| + \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| \stackrel{(a)}{\leq} \sigma_{\min}(\mathbf{X}_*). \tag{108}$$

1066 Inequality (a) follows from assumptions (48) and (49). Thus, we obtain for term (III) that

$$\begin{aligned}
\|(III)\| &\leq \mu^2 \|\mathbf{U}_t\|^4 \left\| \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top \right\| + 2\mu \|\mathbf{U}_t\|^2 \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\| + \mu^2 \|\mathbf{U}_t\|^2 \|\mathbf{M}_t \mathbf{V}_{\mathbf{U}_t}\| \|\mathbf{M}_t\| \\
&\stackrel{(a)}{\leq} 16\mu^2 \|\mathbf{X}_*\|^2 \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + 4\mu \|\mathbf{X}_*\| \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\| + 2\mu^2 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \|\mathbf{M}_t \mathbf{V}_{\mathbf{U}_t}\| \\
&\stackrel{(b)}{\leq} \left(16\mu^2 \|\mathbf{X}_*\|^2 + 8\mu^2 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\| \right) \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| \\
&\quad + (4\mu \|\mathbf{X}_*\| + 2\mu^2 \sigma_{\min}(\mathbf{X}_*) \|\mathbf{X}_*\|) \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\| \\
&\stackrel{(c)}{\leq} \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{16} \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + 5\mu \|\mathbf{X}_*\| \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\|.
\end{aligned}$$

1067 In inequality (a) we used the assumption $\|\mathbf{U}_t\| \leq \sqrt{2\|\mathbf{X}_*\|}$, Lemma 4.6, and inequality (108).
1068 Inequality (b) is due to inequalities (107). In inequality (c) we used the assumption that $\mu \leq$
1069 $\frac{1}{1024\kappa\|\mathbf{X}_*\|}$.

1070 * Conclusion: By adding up all terms, we obtain that

$$\begin{aligned}
\left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top) \right\| &\leq \|(I)\| + \mu \|(II)\| + \|(III)\| \\
&\leq \left(1 - \frac{\mu \sigma_{\min}(\mathbf{X}_*)}{8} \right) \left\| \mathbf{V}_{\mathbf{X}_*}^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \right\| + 5\mu \|\mathbf{X}_*\| \|\mathbf{E}_t \mathbf{V}_{\mathbf{U}_t}\|.
\end{aligned}$$

1071 This completes the proof. \square

1072 **C.3. Proof of Lemma 4.9**

1073 *Proof of Lemma 4.9.* Analogously, as in the proof of Lemma 4.7 we define the shorthand notation

$$\mathbf{M}_t := (\mathcal{A}^* \mathcal{A}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) = \mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top - \underbrace{(\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)}_{=: \mathbf{E}_t}.$$

1074 We note that

$$\|\mathbf{M}_t\| \leq \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + \|(\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top)\| \leq (c_2 + c_3) \sigma_{\min}(\mathbf{X}_*).$$

1075 With an analogous computation as in the proof of Lemma 4.7, it follows that

$$\begin{aligned} \mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top &= (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) (\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top) - \mu^2 \mathbf{U}_t \mathbf{U}_t^\top (\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top) \mathbf{U}_t \mathbf{U}_t^\top \\ &\quad - \mu \mathbf{E}_t \mathbf{U}_t \mathbf{U}_t^\top - \mu \mathbf{U}_t \mathbf{U}_t^\top \mathbf{E}_t - \mu^2 \mathbf{M}_t \mathbf{U}_t \mathbf{U}_t^\top \mathbf{M}_t. \end{aligned}$$

1076 When $c_1 \leq 1/2$, we have $\|\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top\| \leq 1$ by assumption (50). It follows from the assumptions
1077 $\mu \leq \frac{c_1}{\|\mathbf{X}_*\|}$, (51), and (52) that for sufficiently small $c_1, c_2, c_3 > 0$

$$\begin{aligned} \|\mathbf{X}_* - \mathbf{U}_{t+1} \mathbf{U}_{t+1}^\top\| &\leq \|\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top\| \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \|\mathbf{Id} - \mu \mathbf{U}_t \mathbf{U}_t^\top\| + \mu^2 \|\mathbf{U}_t\|^4 \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| \\ &\quad + 2\mu \|\mathbf{E}_t\| \|\mathbf{U}_t\|^2 + \mu^2 \|\mathbf{M}_t\|^2 \|\mathbf{U}_t\|^2 \\ &\leq \|\mathbf{X}_* - \mathbf{U}_t \mathbf{U}_t^\top\| + 4\mu^2 c_2 \|\mathbf{X}_*\|^2 \sigma_{\min}(\mathbf{X}_*) + 4\mu c_3 \|\mathbf{X}_*\| \sigma_{\min}(\mathbf{X}_*) \\ &\quad + 2(c_2 + c_3)^2 \mu^2 \|\mathbf{X}_*\| \sigma_{\min}^2(\mathbf{X}_*) \\ &\leq (c_2 + 4c_1^2 c_2 + 4c_1 c_3 + 2(c_2 + c_3)^2 c_1^2) \sigma_{\min}(\mathbf{X}_*) \\ &\leq \left(1 - \frac{1}{\sqrt{2}}\right) \sigma_{\min}(\mathbf{X}_*). \end{aligned}$$

1078 This completes the proof. □

1079 **D. Proofs regarding the Restricted Isometry Property and its**
1080 **consequences**

1081 **D.1. Proof of Lemma 2.2**

1082 As already mentioned in Section 2.1, there exist similar versions of Lemma 2.1 in the literature
1083 (see, e.g., [50]), which, however, do not specify the dependence of the number of samples m on the
1084 constant $\delta > 0$. It would be possible to trace the steps of the ε -net argument in [50] and work out the
1085 δ -dependence explicitly. However, this would lead to an extra $\log(1/\delta)$ -factor, which is unnecessary.
1086 The reason is that as δ is decreased, a covering with smaller balls is required, leading to a larger ε -
1087 net. This observation suggests a proof strategy based on generic chaining. Indeed, we will use the
1088 following general theorem from [62], which is proven via the generic chaining technique. To state
1089 it, we define the diameter of a set of matrices \mathcal{B} with respect to some norm $\|\cdot\|$ as

$$d_{\|\cdot\|}(\mathcal{B}) := \sup_{\mathbf{B} \in \mathcal{B}} \|\mathbf{B}\|.$$

1090 Moreover, we will also need Talagrand's functional $\gamma_2(\mathcal{B}, \|\cdot\|)$ [63], where for a precise definition,
1091 we refer to [62].

1092 **Theorem D.1** (Theorem 3.1 in [62]). *Let \mathcal{B} be a set of matrices, and $\boldsymbol{\xi}$ be a random Gaussian vector, i.e.,*
1093 *$\boldsymbol{\xi}$ has i.i.d. entries with distribution $\mathcal{N}(0, 1)$. Set*

$$\begin{aligned} E &:= \gamma_2(\mathcal{B}, \|\cdot\|) (\gamma_2(\mathcal{B}, \|\cdot\|) + d_{\|\cdot\|_F}(\mathcal{B})) + d_{\|\cdot\|_F}(\mathcal{B}) d_{\|\cdot\|}(\mathcal{B}), \\ V &:= d_{\|\cdot\|}(\mathcal{B}) (\gamma_2(\mathcal{B}, \|\cdot\|) + d_{\|\cdot\|_F}(\mathcal{B})), \quad U := d_{\|\cdot\|}^2(\mathcal{B}). \end{aligned}$$

1094 Then, for any $t > 0$,

$$\mathbb{P} \left(\sup_{\mathbf{B} \in \mathcal{B}} \left| \|\mathbf{B}\boldsymbol{\xi}\|_2^2 - \mathbb{E} \|\mathbf{B}\boldsymbol{\xi}\|_2^2 \right| > c_1 E + t \right) \leq 2 \exp \left(-c_2 \min \left\{ \frac{t^2}{V^2}, \frac{t}{U} \right\} \right),$$

1095 where $c_1, c_2 > 0$ denote absolute constants.

1096 With this result in place, we can give a proof of Lemma 2.2. This proof strategy has been used in
 1097 [62, Section A.3].

1098 *Proof of Lemma 2.2.* Since \mathcal{A} is a linear operator we can write $\mathcal{A}(\mathbf{X}) = \mathbf{V}_{\mathbf{X}}\boldsymbol{\xi}$, where $\boldsymbol{\xi}$ is a Gaussian
 1099 random vector with independent entries of length $m \binom{d+1}{2}$ and

$$\mathbf{V}_{\mathbf{X}} := \frac{1}{\sqrt{m}} \begin{bmatrix} \text{vec}(\mathbf{X})^\top & & & \\ & \text{vec}(\mathbf{X})^\top & & \\ & & \ddots & \\ & & & \text{vec}(\mathbf{X})^\top \end{bmatrix}$$

1100 is an $m \times (m \binom{d+1}{2})$ block-diagonal matrix. Here, $\text{vec}(\mathbf{X}) \in \mathbb{R}^{\binom{d+1}{2}}$ is a vector indexed by $\{(i, j) \in$
 1101 $[d] \times [d] : i \leq j\}$ such that

$$\text{vec}(\mathbf{X})(i, j) = \begin{cases} \sqrt{2}\mathbf{X}_{ij} & i \neq j \\ \mathbf{X}_{ii} & i = j. \end{cases}$$

1102 Let

$$D_r := \{\mathbf{X} \in \mathcal{S}^d : \|\mathbf{X}\|_F = 1, \text{rank}(\mathbf{X}) \leq r\}.$$

1103 Then it follows from the identity $\mathcal{A}(\mathbf{X}) = \mathbf{V}_{\mathbf{X}}\boldsymbol{\xi}$ that

$$\delta_r := \sup_{\mathbf{X} \in D_r} \left| \|\mathcal{A}(\mathbf{X})\|_2^2 - \|\mathbf{X}\|_F^2 \right| = \sup_{\mathbf{X} \in D_r} \left| \|\mathbf{V}_{\mathbf{X}}\boldsymbol{\xi}\|_2^2 - \mathbb{E}\|\mathbf{V}_{\mathbf{X}}\boldsymbol{\xi}\|_2^2 \right|.$$

1104 Denote $\mathcal{B} := \{\mathbf{V}_{\mathbf{X}} : \mathbf{X} \in D_r\}$. We now estimate the parameters in Theorem D.1. Note that it follows
 1105 directly from the definition of $\text{vec}(\mathbf{X})$ that $\|\text{vec}(\mathbf{X})\|_2 = \|\mathbf{X}\|_F = 1$ and hence $\|\mathbf{V}_{\mathbf{X}}\|_F = \|\mathbf{X}\|_F$
 1106 for all $\mathbf{X} \in \mathcal{S}^d$. Thus, we have $d_F(\mathcal{B}) = 1$ since $\|\mathbf{V}_{\mathbf{X}}\|_F = \|\mathbf{X}\|_F$ for all $\mathbf{X} \in D_r$. On the other hand, for
 1107 $\mathbf{X} \in D_r$,

$$m\mathbf{V}_{\mathbf{X}}\mathbf{V}_{\mathbf{X}}^T = \mathbf{Id}_m,$$

1108 which implies that

$$\|\mathbf{V}_{\mathbf{X}}\| = \frac{1}{\sqrt{m}} \|\text{vec}(\mathbf{X})\|_2 = \frac{1}{\sqrt{m}} \|\mathbf{X}\|_F \quad (109)$$

1109 and $d_{\|\cdot\|}(\mathcal{B}) = \frac{1}{\sqrt{m}}$. From [50, Lemma 3.1], it follows that the covering number for $d \times d$ symmetric
 1110 matrices with Frobenius norm 1 and rank at most r satisfies

$$\mathcal{N}(D_r, \|\cdot\|_F, \varepsilon) \leq (1 + 6/\varepsilon)^{(2d+1)r}. \quad (110)$$

1111 Using Dudley's integral estimate (see, e.g., [63]), combined with (109) and (110), we obtain that

$$\gamma_2(\mathcal{B}, \|\cdot\|) = \gamma_2(D_r, \|\cdot\|_F) \leq C \frac{1}{\sqrt{m}} \int_0^1 \sqrt{\log(\mathcal{N}(D_r, \|\cdot\|_F, u))} du \leq C' \sqrt{\frac{dr}{m}}.$$

1112 With the notations in Theorem D.1, we have

$$E = C' \sqrt{\frac{dr}{m}} \left(C' \sqrt{\frac{dr}{m}} + 1 \right) + \frac{1}{\sqrt{m}}, \quad V = \frac{1}{\sqrt{m}} \left(C' \sqrt{\frac{dr}{m}} + 1 \right), \quad U = \frac{1}{m}.$$

Therefore, applying Theorem D.1, we have $\delta_r \leq \delta$ with probability at least $1 - \varepsilon$ when

$$m \geq C\delta^{-2}(rd + \log(2\varepsilon^{-1})).$$

1113 Here, $C > 0$ denotes some universal constant. This completes the proof of Lemma 2.2. \square

1114 **D.2. Proof of Lemma 2.4**

1115 *Proof of Lemma 2.4.* We will establish first that for all symmetric matrices $\mathbf{Z}_1, \mathbf{Z}_2 \in \mathbb{R}^{d \times d}$ with rank
1116 $\text{rank}(\mathbf{Z}_1) = r$ and $\text{rank}(\mathbf{Z}_2) = r'$ it holds that

$$|\langle (\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z}_1), \mathbf{Z}_2 \rangle| \leq \delta_{r+r'} \|\mathbf{Z}_1\|_F \|\mathbf{Z}_2\|_F. \quad (111)$$

1117 Let us remark that in the case of $\langle \mathbf{Z}_1, \mathbf{Z}_2 \rangle = 0$, this inequality has been proven in [50, Lemma 3.3].
1118 The following proof of this slightly more general statement is analogous.

1119 To prove inequality (111) we assume without loss of generality that $\|\mathbf{Z}_1\|_F = \|\mathbf{Z}_2\|_F = 1$. We note
1120 first that from the parallelogram identity, it follows that

$$\begin{aligned} \langle \mathcal{A}(\mathbf{Z}_1), \mathcal{A}(\mathbf{Z}_2) \rangle &= \frac{1}{4} \|\mathcal{A}(\mathbf{Z}_1 + \mathbf{Z}_2)\|_2^2 - \frac{1}{4} \|\mathcal{A}(\mathbf{Z}_1 - \mathbf{Z}_2)\|_2^2 \\ &\leq \frac{1 + \delta_{r+r'}}{4} \|\mathbf{Z}_1 + \mathbf{Z}_2\|_F^2 - \frac{1 - \delta_{r+r'}}{4} \|\mathbf{Z}_1 - \mathbf{Z}_2\|_F^2 \\ &= \frac{\delta_{r+r'}}{2} \left(\|\mathbf{Z}_1\|_F^2 + \|\mathbf{Z}_2\|_F^2 \right) + \langle \mathbf{Z}_1, \mathbf{Z}_2 \rangle. \end{aligned}$$

1121 By rearranging terms and using the assumption $\|\mathbf{Z}_1\|_F = \|\mathbf{Z}_2\|_F = 1$ we obtain that

$$\langle (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z}_1), \mathbf{Z}_2 \rangle = \langle \mathcal{A}(\mathbf{Z}_1), \mathcal{A}(\mathbf{Z}_2) \rangle - \langle \mathbf{Z}_1, \mathbf{Z}_2 \rangle \leq \delta_{r+r'}.$$

1122 Since the reverse bound

$$\langle (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{Z}_1), \mathbf{Z}_2 \rangle \geq -\delta_{r+r'}$$

1123 can be shown analogously, inequality (111) follows.

1124 Next, we prove inequality (7). For that, we note that there exists a matrix $\mathbf{M} \in \mathbb{R}^{d \times r'}$ with $\|\mathbf{M}\|_F = 1$
1125 such that

$$\begin{aligned} \|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})\mathbf{V}\|_F &= \langle [(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})]\mathbf{V}, \mathbf{M} \rangle = \langle [(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})], \mathbf{V}\mathbf{M}^\top \rangle \\ &= \langle [(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})], \frac{1}{2}\mathbf{V}\mathbf{M}^\top + \frac{1}{2}\mathbf{M}\mathbf{V}^\top \rangle. \end{aligned}$$

1126 holds. Using inequality (111) we obtain that

$$\|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})\mathbf{V}\|_F \leq \delta_{r+2r'} \|\mathbf{Z}\|_F \left\| \frac{1}{2}\mathbf{V}\mathbf{M}^\top + \frac{1}{2}\mathbf{M}\mathbf{V}^\top \right\|_F \leq \delta_{r+2r'} \|\mathbf{Z}\|_F \|\mathbf{V}\| \|\mathbf{M}\|_F = \delta_{r+2r'} \|\mathbf{Z}\|_F.$$

1127 This proves inequality (7).

1128 Inequality (8) is a direct consequence of (7). Indeed, let $\mathbf{v} \in \mathbb{R}^d$ with $\|\mathbf{v}\|_2 = 1$ be an eigenvector
1129 of $(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})$ corresponding to the largest eigenvalue in absolute value. It then follows from
1130 inequality (7) that

$$\|(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})\| = \|[(\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{Z})] \mathbf{v}\|_2 \leq \delta_{r+2} \|\mathbf{Z}\|_F.$$

1131 It remains to prove inequality (10). Note that using the fact $\langle \mathbf{w}\mathbf{w}^\top, \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z}) \rangle = 0$, we have

$$\begin{aligned} |\langle \mathcal{A}(\mathbf{w}\mathbf{w}^\top), \mathcal{A}(\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z})) \rangle| &= |\langle (\mathcal{A}^* \mathcal{A})(\mathbf{w}\mathbf{w}^\top), \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z}) \rangle| \\ &= |\langle (\mathcal{I} - \mathcal{A}^* \mathcal{A})(\mathbf{w}\mathbf{w}^\top), \mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z}) \rangle| \\ &\stackrel{(a)}{\leq} \delta_{(r+1)+1} \|\mathbf{w}\mathbf{w}^\top\|_F \|\mathcal{P}_{\mathbf{w}\mathbf{w}^\top, \perp}(\mathbf{Z})\|_F \\ &\leq \delta_{r+2} \|\mathbf{Z}\|_F, \end{aligned}$$

1132 where in inequality (a) we used (111). This completes the proof of Lemma 2.4. \square