

Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction

Dominik Stöger and Mahdi Soltanolkotabi
Ming Hsieh Department of Electrical and Computer Engineering
University of Southern California

June 28, 2021

Abstract

Recently there has been significant theoretical progress on understanding the convergence and generalization of gradient-based methods on nonconvex losses with overparameterized models. Nevertheless, many aspects of optimization and generalization and in particular the critical role of small random initialization are not fully understood. In this paper, we take a step towards demystifying this role by proving that small random initialization followed by a few iterations of gradient descent behaves akin to popular spectral methods. We also show that this *implicit spectral bias* from small random initialization, which is provably more prominent for overparameterized models, also puts the gradient descent iterations on a particular trajectory towards solutions that are not only globally optimal but also generalize well. Concretely, we focus on the problem of reconstructing a low-rank matrix from a few measurements via a natural nonconvex formulation. In this setting, we show that the trajectory of the gradient descent iterations from small random initialization can be approximately decomposed into three phases: (I) a *spectral or alignment phase* where we show that the iterates have an implicit spectral bias akin to spectral initialization allowing us to show that at the end of this phase the column space of the iterates and the underlying low-rank matrix are sufficiently aligned, (II) a *saddle avoidance/refinement phase* where we show that the trajectory of the gradient iterates moves away from certain degenerate saddle points, and (III) a *local refinement phase* where we show that after avoiding the saddles the iterates converge quickly to the underlying low-rank matrix. Underlying our analysis are insights for the analysis of overparameterized nonconvex optimization schemes that may have implications for computational problems beyond low-rank reconstruction.

1 Introduction

Many contemporary problems in machine learning and signal estimation spanning deep learning to low-rank matrix reconstruction involve fitting nonlinear models to training data. Despite tremendous empirical progress, theoretical understanding of these problems poses two fundamental challenges. First, from an *optimization* perspective, fitting these models often requires solving highly nonconvex optimization problems and except for a few special cases, it is not known how to provably find globally or approximately optimal solutions. Yet simple heuristics such as running (stochastic) gradient descent from (typically) small random initialization is surprisingly effective at finding globally optimal solutions. A second *generalization* challenge is that many modern learning models including neural

network architectures are trained in an overparameterized regime where the parameters of the model exceed the size of the training dataset. It is well understood that in this overparameterized regime, these large models are highly expressive and have the capacity to (over)fit arbitrary training datasets including pure noise. Mysteriously however overparameterized models trained via simple algorithms such as (stochastic) gradient descent when initialized at random continue to predict well or generalize on yet unseen test data. Indeed small random initialization followed by (stochastic) gradient descent iterative updates is arguably the most widely used learning algorithm in modern machine learning and signal estimation.

There has been a large number of exciting results aimed at demystifying both the optimization and generalization aspects over the past few years. We will elaborate on these results in detail in the supplementary, however, we would like to briefly mention the common techniques and their existing limitations. On the optimization front a large body of work has emerged on providing guarantees for nonconvex optimization which can roughly be put into two categories: (I) smart initialization+local convergence and (II) landscape analysis+saddle escaping algorithms. Approaches in (I) focus on showing local convergence of local search techniques from carefully designed spectral initializations [1, 2, 3, 4, 5, 6, 7, 8]. Approaches in (II) focus on showing that in some cases the optimization landscape is benign in the sense that all local minima are global (no spurious local minima) and the saddle points have a direction of strict negative curvature (strict saddle) [9]. Then specialized truncation or saddle escaping algorithms such as trust region, cubic regularization [10, 11], or noisy (stochastic) gradient-based methods [12, 13, 14, 15] are deployed to provably find a global optimum. Both approaches fail to fully explain the typical behavior of local search techniques in practice. Indeed, for many nonconvex problems local search techniques or simple variants, when initialized at random, quickly converge to globally optimal solutions without getting stuck in local optima/saddles without the need for sophisticated initialization or saddle escaping heuristics. We note that while for differentiable losses eventual convergence to local minimizers is known from a random initialization [16] on problems of the form (II), these results cannot rule out exponentially slow cases in the worst-case [17]. Indeed, it has been argued that a more granular analysis of the trajectory of gradient descent beyond the landscape may be necessary [18].

Similarly, there has been a lot of exciting progress on the generalization front, especially for neural networks. Specific to generalization capabilities of gradient-based approaches these results broadly fall into two categories: (1) the first category is based on a linearization principle which characterizes the performance of nonlinear models such as neural networks by comparing it to a linearized kernel problem around the initialization (a.k.a. Neural Tangent Kernels) [19, 20, 21, 22, 23]. (2) the second category is based on a continuous limit analysis in the limit of width going to infinity and learning rate going to zero (mean-field analysis) [24, 25, 26, 27, 28]. However, these existing analyses contain many idealized and non-realistic assumptions (e.g. requiring large widths or large standard deviation at random initialization) and therefore cannot fully explain the success of overparameterized models or serve as a guiding principle for practitioners [29].

Despite the aforementioned exciting recent theoretical progress many aspects of optimization and generalization and in particular the role of random initialization remains mysterious. This leads us to the main challenge of this paper

Why is small random initialization combined with gradient descent updates so effective at finding globally optimal models that generalize well despite the nonconvex nature of the optimization landscape or model overparameterization?

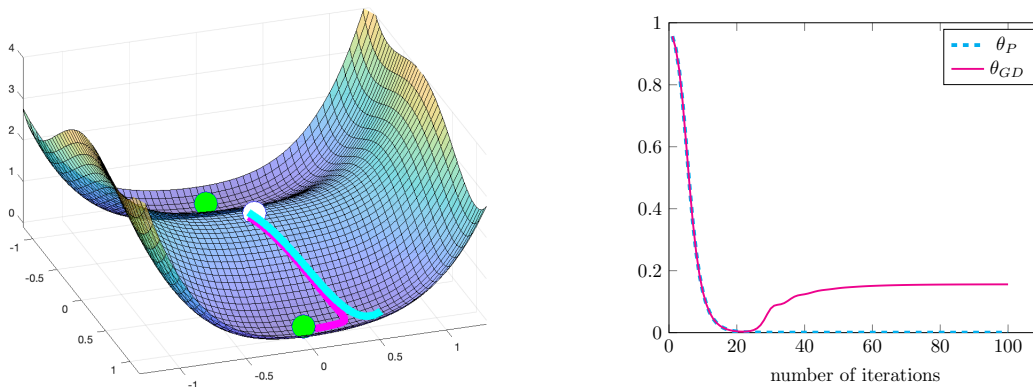


Figure 1: **Gradient descent from small random initialization is akin to spectral initialization.** The left figure depicts the empirical landscape of a low-rank matrix reconstruction problem with the two green circles depicting the two global minima and the white circle the saddle point at the origin. In this figure, we also depict the trajectory of the gradient descent iterations (magenta) together with the power method based on a popular spectral initialization technique (blue). Both gradient descent and power method use the same small initialization near the origin. We see that in the early stage, the two trajectories are almost the same. The figure on the right depicts the angle between the gradient descent (magenta)/power method (blue) iterates and a popular spectral initialization technique, denoted by θ_{GD} and θ_P respectively. This figure clearly demonstrates that for the first iterations these angles are practically the same further confirming that the initial trajectory of gradient descent and power methods are similar. See Section 6 for further detail on the experimental setup.

In this paper we wish to take a step towards addressing the above challenge by demystifying the critical role of small random initialization in gradient-based approaches. Specifically we show that

Small random initialization followed by a few iterations of gradient descent behaves akin to spectral initialization.

By that, we mean more precisely, that if the initialization is chosen small enough, then in the initial stage of the training, gradient descent *implicitly* behaves like spectral initialization techniques such as those commonly used in techniques based on the method of moments. This *implicit spectral bias* of gradient descent from random initialization puts the gradient descent iterations on a particular trajectory towards solutions that are not only globally optimal but also generalize well for overparameterized models. We also show that with small random initialization this implicit spectral bias phenomenon is more prominent for more overparameterized models in the sense that it materializes after fewer iterations. This intriguing phenomenon is depicted in Figure 1 in the context of a low-rank reconstruction problem. This figure clearly demonstrates that the first few iterations of gradient descent starting from a small random initialization are virtually identical to that of running power iterations (a popular algorithm to find the spectral initialization).

Concretely we focus on the problem of low-rank matrix recovery, where our goal is to recover a low-rank matrix of the form XX^T from a few linear measurements. We consider a natural, non-convex approach based on matrix factorization, where we minimize the loss function via gradient

descent. In this paper, we show that, regardless of the amount of overparameterization used, for small random initialization vanilla gradient descent will always converge towards the low-rank solution. This holds as long the measurement operator obeys a popular restricted isometry property [30].¹

Our analysis consists of three phases. The first phase is the aforementioned *spectral* or *alignment* phase where we show gradient descent from small random initialization behaves akin to spectral initialization. Next, we show that after this first *spectral or alignment phase*, gradient descent enters a second phase, which we refer to as *saddle avoidance phase*. In this phase, we show that the trajectory of the gradient iterates moves away from degenerate saddle points, while the iterates maintain almost the same effective rank as XX^T . In the third phase, the *local refinement phase*, we show that the iterates approximately converge towards the underlying low-rank matrix XX^T with a geometric rate up to a certain error floor which depends on the initialization scale. In particular, by decreasing the scale of initialization this error threshold can be made arbitrarily small. While in this paper our main focus is on low-rank matrix reconstruction, we believe that our analysis holds more generally for a variety of contemporary machine learning and signal estimation tasks including neural networks.

2 Low-rank matrix recovery via non-convex optimization

As mentioned earlier in this paper we focus on reconstructing a (possibly overparameterized) Positive Semidefinite (PSD) low rank matrix from a few measurements. In this problem, given m observations of the form

$$y_i = \langle A_i, XX^T \rangle = \text{Tr}(A_i XX^T) \quad i = 1, \dots, m, \quad (1)$$

we wish to reconstruct the unknown matrix XX^T . Here, $X \in \mathbb{R}^{n \times r_\star}$ with $1 \leq r_\star \leq n$ is a factor of the unknown matrix and $\{A_i\}_{i=1}^m$ are known symmetric measurement matrices. A common approach to solving this problem is via minimizing the loss function

$$\min_{\bar{U} \in \mathbb{R}^{n \times r}} f(\bar{U}) := \min_{\bar{U} \in \mathbb{R}^{n \times r}} \frac{1}{4m} \sum_{i=1}^m (y_i - \langle A_i, \bar{U} \bar{U}^T \rangle)^2,$$

with $r \geq r_\star$. More compactly one can rewrite the optimization problem above in the form

$$\min_{\bar{U} \in \mathbb{R}^{n \times r}} f(\bar{U}) := \min_{\bar{U} \in \mathbb{R}^{n \times r}} \frac{1}{4} \|\mathcal{A}(\bar{U} \bar{U}^T - XX^T)\|_{\ell_2}^2, \quad (2)$$

where $\mathcal{A}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ is the measurement operator defined by $[\mathcal{A}(Z)]_i := \frac{1}{\sqrt{m}} \langle A_i, Z \rangle$.

In order to solve the minimization problem (2) we run gradient descent iterations starting from (often small) random initialization. More specifically,

$$U_{t+1} = U_t - \mu \nabla f(U_t) = U_t + \mu [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t.$$

where $U_0 = \alpha U$ is the initialization matrix. Here, $U \in \mathbb{R}^{n \times r}$ is a typically random matrix which represents the form of the initialization and $\alpha > 0$ is a scaling parameter.

There are two challenges associated with analyzing such randomly initialized gradient descent updates. The first is an *optimization* challenge. Since f is *non-convex* it is a priori not clear whether

¹To the best of our knowledge, even for the non-overparameterized scenario, this is the first result that shows that vanilla gradient descent converges to the ground truth from a random initialization, when the measurement operator only satisfies the restricted isometry property.

gradient descent converges to a global optimum or whether it gets stuck in a local minima and/or saddle. The second challenge is that of *generalization*. This is particularly pronounced in the overparameterized scenario where the number of parameters are larger than the number of data points i.e. $rn \geq m$. In this case, there are infinitely many \bar{U} such that $f(\bar{U}) = 0$, but $\|\bar{U}\bar{U}^T - XX^T\|_F$ is arbitrarily large (see, e.g., [31, Proposition 1]). That is, even if gradient descent converges to a global optimum, i.e. $f(\bar{U}) = 0$, it is a priori not clear whether it has found the low-rank solution XX^T (see also Figure 7).

3 Main results

In this section, we present our main results. Stating these results requires a couple of simple definitions. The first definition concerns the measurement operator \mathcal{A} .

Definition 3.1 (Restricted Isometry Property (RIP)). *The measurement operator $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ satisfies RIP of rank r with constant $\delta > 0$, if it holds for all matrices Z of rank at most r*

$$(1 - \delta) \|Z\|_F^2 \leq \|\mathcal{A}(Z)\|_{\ell_2}^2 \leq (1 + \delta) \|Z\|_F^2. \quad (3)$$

We note that for a Gaussian measurement operator \mathcal{A} ², RIP of rank r and constant $\delta > 0$ holds with high probability, if the number of observations satisfies $m \gtrsim nr/\delta^2$ [30, 32].

The second definition concerns the condition number of the factor X .

Definition 3.2 (condition number). *We denote the condition number of $X \in \mathbb{R}^{n \times r_\star}$ by*

$$\kappa := \frac{\|X\|}{\sigma_{r_\star}(X)},$$

where $\sigma_{r_\star}(X)$ denotes r_\star -th largest singular value of X .

With these definitions in place we are now ready to state our main results.

3.1 General case: $r > r_\star$

We begin by stating our first main result.

Theorem 3.3. *Let $X \in \mathbb{R}^{n \times r_\star}$ and assume we have m measurements of the low rank matrix XX^T of the form $y = \mathcal{A}(XX^T)$ with \mathcal{A} the measurement operator. We assume \mathcal{A} satisfies the restricted isometry property for all matrices of rank at most $2r_\star + 1$ with constant $\delta \leq c\kappa^{-4}r_\star^{-1/2}$. To reconstruct XX^T from the measurements we fit a model of the form $\bar{U} \mapsto \mathcal{A}(\bar{U}\bar{U}^T)$ with $\bar{U} \in \mathbb{R}^{n \times r}$ and $r > r_\star$ via running gradient descent iterations of the form $U_{t+1} = U_t - \mu \nabla f(U_t)$ on the objective (2) with a step size obeying $\mu \leq c\kappa^{-4}\|X\|^{-2}$. Here, the initialization is given by $U_0 = \alpha U$, where $U \in \mathbb{R}^{n \times r}$ has i.i.d. entries with distribution $\mathcal{N}(0, 1/\sqrt{r})$. With this setting and assumptions the following two statements hold.*

²By that, we mean that all the entries of the (symmetric) measurement matrices $\{A_i\}_{i=1}^m$ are drawn i.i.d. with distribution $\mathcal{N}(0, 1)$ on the off-diagonal and distribution $\mathcal{N}(0, 1/\sqrt{2})$ on the diagonal.

1. Under the assumption that $r \geq 2r_*$ and that the scale of initialization fulfills

$$\alpha \lesssim \min \left\{ \frac{(\min \{r; n\})^{1/4}}{\kappa^{1/2} n^{3/4}} \left(2\kappa^2 \sqrt{\frac{n}{\min \{r; n\}}} \right)^{-6\kappa^2}; \frac{1}{\kappa^7 n} \right\} \|X\|, \quad (4)$$

after

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\frac{C_1 n \kappa}{\min \{r; n\}} \cdot \max \left\{ 1; \frac{\kappa r_*}{\min \{r; n\} - r_*} \right\} \cdot \frac{\|X\|}{\alpha} \right)$$

iterations we have that

$$\frac{\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F}{\|X\|^2} \lesssim \frac{n^{21/16} \kappa^{81/16} r_*^{1/8}}{(\min \{r; n\})^{15/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}}, \quad (5)$$

holds with probability at least $1 - C e^{-\tilde{c}r}$.

2. Assume that $r_* < r < 2r_*$ and that the scale of initialization fulfills

$$\alpha \lesssim \min \left\{ \frac{\varepsilon^{1/2}}{n^{3/4} \kappa^{1/2}} \left(\frac{2\kappa^2 \sqrt{rn}}{\varepsilon} \right)^{-6\kappa^2}; \frac{\varepsilon}{n \kappa^7} \right\} \|X\|, \quad (6)$$

with $0 < \varepsilon < 1$. Then, after

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\frac{C_2 \kappa n^2}{\varepsilon^2 (r - r_*)} \cdot \frac{\|X\|}{\alpha} \right)$$

iterations we have that

$$\frac{\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F}{\|X\|^2} \lesssim r_*^{1/8} (r - r_*)^{3/8} \kappa^{81/16} \left(\frac{n}{\varepsilon} \cdot \frac{\alpha}{\|X\|} \right)^{21/16} \quad (7)$$

holds with probability at least $1 - (\tilde{C}\varepsilon)^{r-r_*+1} + \exp(-\tilde{c}r)$.

Here, $C_1, C_2, C, \tilde{C}, c > 0$ are fixed numerical constants.

Note that the test error $\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F^2$ can be made arbitrarily small by choosing the scale of initialization α small enough. In particular, the dependence of the test error on α is polynomial and the dependence of the number of iterations on α is logarithmic, which means that reducing the test error by scaling down α introduces only modest additional computational cost. Hence, as long as the rank at most $2r_* + 1$ RIP with constant $\delta \leq c\kappa^{-4} r_*^{-1/2}$ holds, gradient descent converges to a point in the proximity of the low-rank solution, whenever the initialization is chosen small enough regardless of the choice of r . This holds even when the model is overparameterized i.e. $rn \gg m$ and the optimization problem has many global optima many of which do not obey $UU^T \approx XX^T$. This result thus further demonstrates that when initialized with a small random initialization gradient descent has an implicit bias towards solutions of low-rank or small nuclear norm. This is in sharp contrast to Neural Tangent Kernel (NTK)-based theory for low-rank matrix recovery (see [20, Section 4.2]) which will not approximately recover the ground truth matrix XX^T due to the larger scale of initialization required when using that technique.

As discussed in Section 2, the restricted isometry property holds with high probability for a sample complexity $m \gtrsim nr_\star^2 \kappa^8$ for Gaussian measurement matrices. Up to constants, this sample complexity is optimal in n , while it is sub-optimal in r_\star and κ compared to approaches based on nuclear-norm minimization (see, e.g., [30]). While there is numerical evidence that the true scaling of m in r_\star should also be linear in the non-convex case [33], we note that the optimal dependence of the sample complexity on r_\star is a major open problem in the field, as the sample complexities in all theoretical results for non-convex approaches in the literature scale at least quadratically in r_\star .

Interpretation: Recall from Section 1 that our convergence analysis can be divided into three phases: the spectral phase, the saddle avoidance phase, and the local refinement phase. As it will become clear from the proofs in the supplementary when $r \geq 2r_\star$ the bound on the number of iterations can be decomposed as follows

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \left[\underbrace{\ln \left(2\kappa^2 \sqrt{\frac{n}{\min\{r; n\}}} \right)}_{\text{Phase I: spectral/alignment phase}} + \underbrace{\ln \left(\frac{\sigma_{\min}(X)}{\alpha} \right)}_{\text{Phase II: saddle avoidance phase}} + \underbrace{\ln \left(\max \left\{ 1; \frac{\kappa r_\star}{\min\{r; n\} - r_\star} \right\} \frac{\|X\|}{\alpha} \right)}_{\text{Phase III: local refinement phase}} \right]. \quad (8)$$

First, we note that the duration of all three phases scales inversely with $\sigma_{\min}(X)^2$. This is due to the fact that in all three phases the dynamics associated the smallest singular value of X is the slowest one and hence needs the most time to complete.

In the spectral phase, the eigenvectors corresponding to the leading r_\star eigenvalues of $U_t U_t^T$ become aligned with the eigenvectors corresponding to the leading r_\star eigenvalues of $\mathcal{A}^* \mathcal{A}(X X^T)$. We observe in (8) that in the spectral phase increasing r , i.e. the amount of parameters, decreases the number of iterations in this phase. As we will explain in the supplementary, the reason is that increasing r decreases the angle between the column space of the initialization U_0 and the span of the eigenvectors corresponding to the leading r_\star eigenvalues of $\mathcal{A}^* \mathcal{A}(X X^T)$ used in spectral initialization. As a consequence, gradient descent needs fewer iterations to align these two subspaces.

In the saddle avoidance phase (Phase II), $\sigma_{r_\star}(U_t)$, the r_\star th largest singular value of U_t , grows geometrically until it is on the order of $\sigma_{\min}(X)$. Hence, this duration depends on the ratio between the $\sigma_{\min}(X)$ and the the scale of initialization α . This is clearly reflected in the upper bound on the number of needed iterations in equation (8).

In Phase III, the local refinement phase, the matrix $U_t U_t^T$ converges towards $X X^T$. In particular, at iteration \hat{t} the test error obeys (5). We observe that a smaller α allows for a smaller test error in (5) but per (8) this higher accuracy is achieved with a modest increase in the required iterations.

3.2 Special case: $r = r_\star$

The following result deals with the scenario $r = r_\star$, that is, the iterates U_t have as many parameters as the ground truth matrix X .

Theorem 3.4. *Let $X \in \mathbb{R}^{n \times r_\star}$ and assume we have m measurements of the low rank matrix $X X^T$ of the form $y = \mathcal{A}(X X^T)$ with \mathcal{A} the measurement operator. We assume \mathcal{A} satisfies the restricted isometry property for all matrices of rank at most $2r_\star + 1$ with constant $\delta \leq \kappa^{-4} r_\star^{-1/2}$. To reconstruct $X X^T$ from the measurements we fit a model of the form $\bar{U} \mapsto \mathcal{A}(\bar{U} \bar{U}^T)$ with $\bar{U} \in \mathbb{R}^{n \times r_\star}$ via running gradient descent iterations of the form $U_{t+1} = U_t - \mu \nabla f(U_t)$ on the objective (2) with a step size obeying $\mu \leq \kappa^{-4} \|X\|^{-2}$. Here, the initialization is given by $U_0 = \alpha U$, where $U \in \mathbb{R}^{n \times r_\star}$ has i.i.d.*

entries with distribution $\mathcal{N}(0, 1/\sqrt{r_\star})$. Assume that the scale of initialization fulfills

$$\alpha \lesssim \min \left\{ \frac{\varepsilon^{1/2}}{n^{3/4} \kappa^{1/2}} \left(\frac{2\kappa^2 \sqrt{rn}}{\varepsilon} \right)^{-6\kappa^2}; \frac{\varepsilon^2}{n\sqrt{r_\star} \kappa^7} \right\} \|X\|,$$

for some $0 < \varepsilon < 1$. Then with probability at least $1 - C\varepsilon + \exp(-cr_\star)$ after

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\frac{8\kappa^3 n^3}{\varepsilon^2} \cdot \frac{\|X\|}{\alpha} \right)$$

iterations we have that

$$\frac{\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F}{\|X\|^2} \lesssim r_\star^{1/8} \kappa^{81/16} \left(\frac{n}{\varepsilon} \cdot \frac{\alpha}{\|X\|} \right)^{21/16}. \quad (9)$$

Here $C, c > 0$ are fixed numerical constants.

Note that by choosing α small enough we can make the test error in (9) arbitrarily small. In particular, this means that then well-known local convergence results can be applied showing that $U_t U_t^T$ converges linearly to $X X^T$ (see, e.g., [4]).

Thus, this result implies that if the measurement operator fulfills the restricted isometry property, gradient descent with *small, random initialization* will converge to the ground truth matrix X in polynomial time. It is known that under the RIP assumption the loss landscape is benign [34] in the sense that there are no local optima that are not global and all saddles have a direction of negative curvature. However, such results do not imply that vanilla gradient descent converges quickly (i.e. in polynomial time) to a global optimum, as gradient descent may take exponential time to escape from saddle points.

To the best of our knowledge, this is the first result which shows the convergence of vanilla gradient descent for a rank- r_\star matrix recovery problem to the ground truth from a random initialization using only the restricted isometry property in polynomial time. The only other paper in the low-rank matrix recovery literature, which shows fast convergence of vanilla gradient descent to the ground truth from a random initialization, is [35]. In this work, the problem of phase retrieval has been studied, which can be formulated as a low-rank matrix recovery problem with $r_\star = 1$. The paper shows that gradient descent converges from a random initialization to the ground truth with a near-optimal number of iterations. However, the proof in this paper leverages the rotation-invariance of the Gaussian measurements vectors via carefully constructed auxiliary sequences. In contrast, Theorem 3.4 above relies only on the restricted isometry property and no further assumptions on \mathcal{A} are needed.

3.3 Special case: $r = n$ with orthonormal initialization

In the following result, we study the scenario $r = n$, where the initialization matrix $U \in \mathbb{R}^{n \times n}$ is an orthonormal matrix, i.e. $U^T U = \text{Id}$, instead of a Gaussian matrix as in the previous results in this paper. This is the same setting as in [36, Theorem 1.1], and we include this special case so as to explain how our results improves upon prior work in this special case.

Theorem 3.5. *Let $X \in \mathbb{R}^{n \times r_\star}$ and assume we have m measurements of the low rank matrix $X X^T$ of the form $y = \mathcal{A}(X X^T)$ with \mathcal{A} the measurement operator. We assume \mathcal{A} satisfies the restricted*

isometry property for all matrices of rank at most $2r_* + 1$ with constant $\delta \leq c\kappa^{-4}r_*^{-1/2}$. To reconstruct XX^T from the measurements we fit a model of the form $\bar{U} \mapsto \mathcal{A}(\bar{U}\bar{U}^T)$ with $\bar{U} \in \mathbb{R}^{n \times n}$ via running gradient descent iterations of the form $U_{t+1} = U_t - \mu \nabla f(U_t)$ on the objective (2) with a step size obeying $\mu \leq c\kappa^{-4}\|X\|^{-2}$. Here, the initialization is given by $U_0 = \alpha U$, where $U \in \mathbb{R}^{n \times n}$ can be any orthonormal matrix. Assume that the scale of initialization satisfies $\alpha \leq c \frac{\sigma_{\min}(X)}{\kappa^2 n}$. Then, after

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\max \left\{ 1; \frac{\kappa r_*}{n - r_*} \right\} \frac{\|X\|}{\alpha} \right)$$

iterations we have that

$$\frac{\|U_{\hat{t}}U_{\hat{t}}^T - XX^T\|_F}{\|X\|^2} \lesssim \frac{r_*^{1/8} n^{3/8}}{\kappa^{3/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}}.$$

Here $c > 0$ is a fixed numerical constant.

Note that this result improves over [36, Theorem 4.1] in several aspects. First of all, in [36] it is assumed that the measurement operator \mathcal{A} has the rank- $4r$ restricted isometry property with constant $\delta \lesssim \kappa^{-6}r_*^{-1/2} \log^{-2} \frac{n}{\alpha}$. In particular, this suggests that this result cannot handle the scenario that the scale of initialization α becomes arbitrarily small, as this would also require that the restricted isometry constant δ becomes arbitrarily small as well. This in turn would require an arbitrarily large sample size. Moreover, [36] requires a step size of at most $\mu \lesssim \kappa^{-6}r_*^{-1/2} \log^2 n \|X\|^{-2}$, whereas the above theorem only needs the weaker assumption $\mu \lesssim \kappa^{-2}\|X\|^{-2}$. These improvements aside the main difference between our result and this prior work is that we can handle any r by formalizing an intriguing connection between small random initialization and spectral learning.

4 Related work

Global convergence guarantees for nonconvex low-rank matrix recovery: As mentioned earlier in Section 1, there is a large body of work on developing global convergence guarantees for nonconvex problems. In the context of low-rank matrix recovery, several papers have demonstrated that low-rank reconstruction problems in a variety of domains can be solved via nonconvex gradient descent starting from spectral initialization. More precisely, this has been shown for phase retrieval [1, 2, 3], matrix sensing [37], blind deconvolution [5, 6], and matrix completion [38]. However, in practice often random initialization is used in lieu of specialized spectral initialization techniques. To remedy this issue, more recent literature [39, 40, 41], focusses on studying the loss landscape of such problems. These papers show that despite their non-convexity under certain assumptions these loss landscapes are benign in the sense that there are no *spurious local minima*, (i.e. all minimizers are global minima) and saddles points have a strict direction of negative curvature (a.k.a. strict saddle) [9]. Then specialized truncation or saddle escaping algorithms such as trust region, cubic regularization [10, 11] or noisy (stochastic) gradient-based methods [12, 13, 14, 15] are deployed to provably find a global optimum. These papers however do not directly develop global convergence for gradient descent (without any additional modification) from a random initialization. For differentiable losses eventual convergence to local minimizers is known from random initialization [16] but these results do not provide convergence rates and only guarantee eventual convergence. Indeed, gradient descent may converge exponentially slowly in the worst-case [17]. In contrast to the above literature our result in Theorem 3.4 (in the case of $r = r_*$) shows that gradient descent from a

small random initialization converges rather quickly to the global optima. As mentioned earlier, we are able to establish this result by demonstrating that in the initial phase gradient descent iterates are intimately connected to the spectral initialization techniques discussed above. Furthermore, the above spectral initialization followed by local convergence or landscape analysis techniques cannot be directly applied in the overparameterized case ($r > r_*$) whereas our analysis works regardless of model overparameterization. Finally, we would like to mention that even more recently the papers [35] prove the convergence of gradient descent starting from a random initialization for low-rank recovery problems via an interesting leave-one-out analysis. To the best of our knowledge, this is the only existing result, which provides convergence guarantees for gradient descent from random initialization for low-rank matrix recovery problems when $r = r_*$. However, the leave-one-out analysis heavily relies on the independence and the rotation invariance of the measurements. Also similar to the above this analysis does not seem to easily lend to generalization in the overparameterized regime. In contrast, our proof techniques rely on standard restricted isometry assumptions without requiring the independence of the measurements and does provide generalization guarantees with model overparameterization ($r > r_*$).

Overparameterization in low-rank matrix recovery: In the influential work [42] it has been conjectured and in the special case that the measurement matrices commute proven that gradient descent on overparameterized matrix factorization converges to the solution with the minimal nuclear norm. This phenomenon is now often referred to as *implicit regularization*. In [18], evidence is provided that adding depth even increases the tendency of gradient descent to converge to low-rank solution. In [43] it has been shown that there are certain scenarios where the conjecture in [42] does not hold. In [44] theoretical and empirical evidence has been provided that gradient flow with infinitesimal initialization is equivalent to a certain rank-minimization heuristic.

In this paper, we shed further light on the implicit regularization of gradient descent. In particular, we provide a precise analysis of the initial stage and relate it to the power method and our analysis explains how overparameterization is beneficial in the initial stages. Closest to our work is the paper [36], which studies a special case of the problem analysed in this paper. More precisely, this paper considers the special case $r = n$ with orthonormal initialization. We also applied our theory to this exact same setting, see Section 3.3, where we include a detailed comparison for this special scenario. Most importantly our theory is able to handle the case $\alpha \rightarrow 0$, which the result in [36] seems not to be able to. Moreover, analysing the full range of possible choice of r requires a careful analysis of the spectral phase, which is one key novelty of this paper compared to [36].

In [45, 46] it has been shown that in certain scenarios, where the measurement matrices A_i are positive semidefinite (PSD), the equation $y = \mathcal{A}(UU^T)$ has a unique low-rank solution. This means that in these scenarios the PSD constraint by itself might lead to a low-rank matrix recovery, which makes implicit regularization by gradient descent meaningless in this setting. However, note that these results not apply to the scenario studied in this paper, as we assume the measurement matrices A_i to be Gaussian, which, in particular, means that they are not positive semidefinite. In particular, in our setting it can be shown that there are infinitely many solutions to the equation $y = \mathcal{A}(UU^T)$ with arbitrarily large test error [31].

Gradient-based generalization guarantees for overparameterized tensors and neural networks: A recent line of work is concerned with connecting the analysis of neural network training with the so-called neural tangent kernel (NTK) [19, 20, 21, 22, 23]. The key idea is that for a large enough initialization, it suffices to consider a linearization of the neural network around

the origin. This allows connecting the analysis of neural networks with the well-studied theory of kernel methods. This is also sometimes referred to as *lazy training*, as with such an initialization the parameters of the neural networks stay close to the parameters at initialization. However, there is a line of work, which suggests that NTK-analysis might not be sufficient to completely explain the success of neural networks in practice. The paper [29] provides empirical evidence that by choosing a smaller initialization the test error of the neural network decreases. A similar performance gap between the performance of the NTK and neural networks has been observed in [47], where it has been shown that the performance gap is larger if the covariance matrix is isotropic.

There is also a line of work [24, 25, 26, 27, 28], which is concerned with the mean-field analysis of neural networks. The insight is that for sufficiently large width the training dynamics of the neural network can be coupled with the evolution of a probability distribution described by a PDE. These papers use a smaller initialization than in the NTK-regime and, hence, the parameters can move away from the initialization. However, these results do not provide explicit convergence rates and require an unrealistically large width of the neural network.

For the problem of tensor decomposition it has also been shown that gradient descent with small initialization is able to leverage low-rank structure [48]. This is relevant to neural network analysis, since in [49] a relationship between tensor decomposition and training neural networks has been established. In [50] it has been shown that neural networks with ReLU function and trained by SGD can outperform any kernel method. One crucial element in their analysis is that the early stage of the training is connected with learning the first and second moment of the data.

While in this paper we do not study overparameterized tensor or neural network models we note that the NTK-theory can also be applied to low-rank matrix recovery (see [20, Section 4.2]). This means that if the scale of initialization is chosen large enough and the number of parameters is larger than the number of measurements, i.e. $nr \gtrsim m$, then gradient descent will converge linearly to a global minimizer with zero loss. However, since for this approach the parameters will stay close to the initialization, this approach will not recover the ground truth matrix XX^T . Hence, an NTK analysis will not yield good generalization. In contrast in this paper we have seen that choosing a small initialization is a remedy for low-rank matrix recovery. So in this sense our result can be viewed as going beyond the lazy training in NTK theory. In fact we believe that similar analysis to the one developed in this paper for low-rank recovery can be used to analyze a much broader class of overparameterized models including the analysis of neural networks. We defer this to a future paper.

Linear neural networks: In [51, 52, 53, 54, 55] the convergence of gradient flow and gradient descent is studied for (deep) linear neural networks of the form

$$\min_{W_1, W_2, \dots, W_N} \sum_{i=1}^m \|W_N \dots W_2 W_1 x_i - y_i\|^2.$$

However, note that this model is different from the one studied in this paper. In [56] it is shown that gradient descent for convolutional linear neural networks has a bias towards the ℓ_p -norm, where p depends on the depth of the network.

5 Overview and key ideas of the proof

In this section, we briefly discuss the key ideas and techniques in our proof. We begin by discussing a simple decomposition, which is utilized throughout our proofs. Next, in Sections 5.2 and 5.3 we

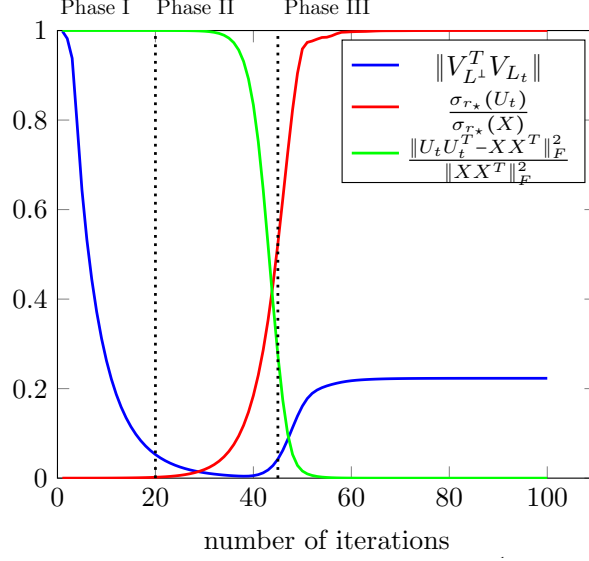


Figure 2: **Depiction of the three phases of convergence.** Let L denote the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $\mathcal{A}^* \mathcal{A}(XX^T)$ and L_t denote the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $U_t U_t^T$. This figure demonstrates that the convergence analysis can be divided into three phases: (I) spectral/alignment phase; (II) saddle avoidance phase and (III) the refinement phase. We see that in the first phase the first r_* eigenvectors of $U_t U_t^T$ rapidly learn the subspace corresponding to the first r_* eigenvectors of $\mathcal{A}^* \mathcal{A}(XX^T)$, i.e. the angle $\|V_{L^\perp}^T V_{L_t}\|$ becomes small. The r_* th largest singular value of U_t is still small in this phase and the (normalized) test error $\|U_t U_t^T - XX^T\|^2 / \|XX^T\|^2$ has not decreased yet. In Phase (II), however, we see that $\sigma_{r_*}(U_t)$ is growing, whereas the loss begins to decrease in this phase and the subspaces stay aligned. In Phase (III) we see that the test error is converging towards 0 rapidly, meaning that $U_t U_t^T$ converges to XX^T . Consequently, $\sigma_{r_*}(U_t) / \sigma_{r_*}(X)$ converges to 1 (red curve). We also see that in this phase the angle $\|V_{L^\perp}^T V_{L_t}\|$ grows again, until it reaches a certain threshold. This is because in this phase the top r_* eigenvalues of $U_t U_t^T$ become aligned with the eigenvectors of XX^T .

show that the trajectory of the gradient descent iterations can be approximately decomposed into three phases: (I) a *spectral or alignment phase* where we show that gradient descent from random initialization behaves akin to spectral initialization allowing us to show that at the end of this phase the column spaces of the iterates U_t and the ground truth matrix X are sufficiently aligned, (II) a *saddle avoidance phase*, where we show that the trajectory of the gradient iterates move away from certain degenerate saddle points, and (III) a *refinement phase*, where the product of the gradient descent iterates $U_t U_t^T$ converges quickly to the underlying low-rank matrix XX^T . The latter result holds up to a small error that is commensurate with the scale of the initialization and tends to zero as the scale of the initialization goes to zero. Figure 2 depicts these three phases.

5.1 Decomposition of U_t into “signal” and “noise” matrices

A key idea in our proof is to decompose the matrix U_t into the sum of two matrices. The first matrix, which is of rank r_* , can be thought of as the “signal” term. We will show that the product

of this matrix with its transpose converges towards the ground truth low-rank matrix XX^T . The second matrix, will have rank at most $r - r_*$ and will have column span orthogonal to the column span of the ground truth matrix X . We will show that the spectral norm of this matrix will remain relatively small depending on the scale of initialization α . Hence, this term can be interpreted as the “noise” term.

We now formally introduce our decomposition. To this aim, consider the matrix $V_X^T U_t \in \mathbb{R}^{r_* \times r}$ and denote its singular value decomposition by $V_X^T U_t = V_t \Sigma_t W_t^T$ with $W_t \in \mathbb{R}^{r \times r_*}$. Similarly, we shall use $W_{t,\perp} \in \mathbb{R}^{r \times (r-r_*)}$ to denote the orthogonal matrix, whose column space is orthogonal to the column space of W_t (i.e. the basis of the subspace orthogonal to the span of W_t). We then can decompose U_t into

$$U_t = \underbrace{U_t W_t W_t^T}_{\text{signal term}} + \underbrace{U_t W_{t,\perp} W_{t,\perp}^T}_{\text{noise term}}.$$

This decomposition has the following two simple properties, which will be useful throughout our proofs.

Lemma 5.1 (Properties of signal-noise decomposition).

1. The column space of the noise term is orthogonal to the column span of X , i.e. $V_X^T U_t W_{t,\perp} = 0$.
2. When $V_X^T U_t$ is full rank, then the signal term has rank r_* and the noise term has rank at most $r - r_*$.

Proof. The first statement follows directly from the observation $V_X^T U_t W_{t,\perp} W_{t,\perp}^T = V_X^T U_t (\text{Id} - W_t W_t^T) = 0$. The second statement is a direct consequence of the definition of W_t . \square

We would like to note that decomposing U_t into two terms has appeared in prior work such as [36] as well as earlier work in the compressive sensing literature. However, [36] uses a different decomposition. A key advantage of our decomposition is that it only depends on U_t and X , whereas the decomposition in [36] depends on all previous iterates U_0, U_1, \dots, U_{t-1} .

5.2 The spectral/alignment phase

In this section we turn our attention to giving an overview of the key ideas and proofs of the spectral/alignment phase. More specifically, we will argue that in the first few iterations gradient descent implicitly performs a form of spectral initialization. By that, we mean that after the first few iterations the column span of the signal term $U_t W_t W_t^T$ is aligned with the column span of X and that $\|U_t W_{t,\perp}\|$ is relatively small compared to $\sigma_{\min}(U_t W_t)$, meaning that the signal term dominates the noise term.

We now provide the main intuition behind the analysis in our spectral/alignment phase. Our starting point is the observation that for the gradient at the initialization $U_0 = \alpha U$ it holds that

$$\begin{aligned} \nabla f(U_0) &= -[\mathcal{A}^* \mathcal{A}(XX^T - U_0 U_0^T)] U_0 \\ &= -\alpha [\mathcal{A}^* \mathcal{A}(XX^T)] U + \alpha^3 [\mathcal{A}^* \mathcal{A}(UU^T)] U. \end{aligned}$$

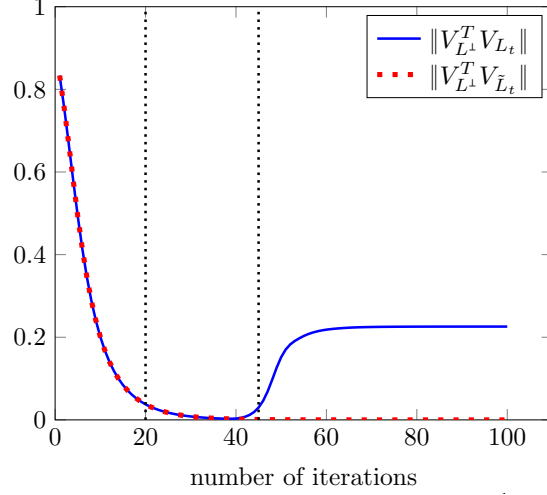


Figure 3: **Depiction of the spectral alignment phase: in the first few iterations, gradient descent with small initialization behaves like a power method.** Here, L denotes the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $\mathcal{A}^* \mathcal{A}(XX^T)$. L_t denotes the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $U_t U_t^T$. Moreover, \tilde{L}_t denotes the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $\tilde{U}_t \tilde{U}_t^T$, where $\tilde{U}_t = (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(XX^T))^t U_0$. We see that in the first iterations U_t and \tilde{U}_t learn the subspace L at the same rate.

In particular, we observe that for $\alpha > 0$ sufficiently small the second term is negligible. Hence, we have that

$$\begin{aligned} U_1 &= U_0 - \mu \nabla f(U_0) \\ &= (\text{Id} + \mu [\mathcal{A}^* \mathcal{A}(XX^T)]) U_0 - \alpha^2 [\mathcal{A}^* \mathcal{A}(UU^T)] U_0 \\ &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(XX^T)) U_0 + O(\alpha^2 \|U_0\|) \end{aligned}$$

In the first few iterations (i.e. small t) we expect the matrix U_t to be small and continue to scale commensurately with α and we expect that a similar approximation holds for the first iterations. Hence, for α sufficiently small we can approximate U_t by

$$U_t \approx \underbrace{(\text{Id} + \mu \mathcal{A}^* \mathcal{A}(XX^T))^t}_{=: Z_t} U_0 := \tilde{U}_t. \quad (10)$$

Figure 3 clearly illustrates that the first few iterations of gradient descent behave essentially identical to (10) confirming our intuition and proofs.

We indeed formally prove that such an approximation holds in Section 8. We note that the matrix $Z_1 = \text{Id} + \mu \mathcal{A}^* \mathcal{A}(XX^T)$ is the basis for the commonly used spectral initialization, where typically a factorization of the rank r_* approximation of this matrix is used as the initialization [4, 3, 38]. Therefore, the approximation (10) suggests that gradient descent iterates modulo the normalization are akin to running power method on Z_1 . Therefore, we expect the column space of the signal term at the end of the spectral phase to be closely aligned with those of the commonly used spectral initialization techniques and in turn the column space of X as we formalize below.

To be more precise about the aforementioned alignment with X , let the singular value decomposition of $\mathcal{A}^* \mathcal{A}(XX^T)$ be given by $\mathcal{A}^* \mathcal{A}(XX^T) = \sum_{i=1}^n \lambda_i v_i v_i^T$. It follows that

$$U_t \approx \left[\sum_{i=1}^n (1 + \mu \lambda_i)^t v_i v_i^T \right] U_0. \quad (11)$$

It is well-known that when the operator \mathcal{A} obeys the restricted isometry property we have

$$\mathcal{A}^* \mathcal{A}(XX^T) \approx XX^T.$$

In particular, we have that

$$\lambda_{r_*+1}(\mathcal{A}^* \mathcal{A}(XX^T)) \ll \lambda_{r_*}(\mathcal{A}^* \mathcal{A}(XX^T)).$$

Hence, it follows from

$$Z_t = \sum_{i=1}^n (1 + \mu \lambda_i)^t v_i v_i^T$$

that $\lambda_{r_*}(Z_t)/\lambda_{r_*+1}(Z_t)$ grows exponentially. In particular, this means that

$$Z_t \approx \sum_{i=1}^{r_*} (1 + \mu \lambda_i)^t v_i v_i^T$$

and, by (10),

$$U_t \approx \left[\sum_{i=1}^{r_*} (1 + \mu \lambda_i)^t v_i v_i^T \right] U_0.$$

Since U_0 is a *random* Gaussian matrix, for an appropriate choice of t , we will be able to show that the matrix U_t has the following two properties with high probability, where $L = \text{span}\{v_1; \dots; v_{r_*}\}$ and L_t is the projection of U_t onto its best rank- r_* approximation:

- There is a sufficiently large gap between $\sigma_{r_*}(U_t)$ and $\sigma_{r_*+1}(U_t)$, i.e., $\frac{\sigma_{r_*}(U_t)}{\sigma_{r_*+1}(U_t)} \geq \Delta > 1$, where Δ is an appropriately chosen constant.
- We have that $\|V_{L_t}^T V_{L_t}\|$ is small. Since the column space of $\mathcal{A}^* \mathcal{A}(XX^T)$ is aligned with the column space of X , this also implies that $\|V_{X^\perp}^T V_{L_t}\|$ is small.

This confirms that in the first few iterations, gradient descent indeed implicitly performs akin to spectral initialization with the column space of U_t aligned with the column space of X . However, this does not yet fully complete our analysis for the spectral/alignment phase, since critical to the analysis of second phase we need certain properties to hold for the signal and noise terms $U_t W_t$ and $U_t W_{t,\perp}$ (see Section 5.1) rather than the singular value decomposition of U_t . However, using the properties of the SVD of U_t , which are listed above, we will establish the following properties of $U_t W_t$ and $U_t W_{t,\perp}$.

- The column space of $U_t W_t$ is aligned with the column space of X : $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-2}$.
- The spectral norm of the noise term is not too large compared to the minimum singular value of the signal term, i.e., $2\sigma_{\min}(U_t W_t) \geq \|U_t W_{t,\perp}\|$.
- The spectral norm of the noise term is bounded from above in the sense that i.e., $\|U_t W_{t,\perp}\| \ll \sigma_{\min}(X)$.
- The spectral norm of U_t is bounded, i.e., $\|U_t\| \leq 3\|X\|$.

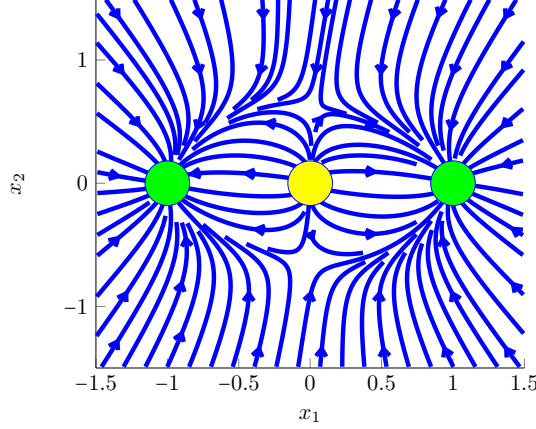


Figure 4: **Depiction of saddle avoidance and local refinement phases.** In this figure we depict the gradient field of the loss function f with $n = 2$, $r_* = r = 1$, $m = 15$, and $X = \begin{pmatrix} 1 & 0 \end{pmatrix}$. The green circles depict the two generalizable global minima of f , namely $(1 \ 0)$ and $(-1 \ 0)$. The red circle depicts the saddle point $(0 \ 0)$. As this figure demonstrates starting from small random initialization after a while the trajectory moves away from the saddle (i.e. avoids it) and then converges to one of the two generalizable global optima (i.e. the local refinement phase).

5.3 The saddle avoidance phase and the refinement phase

In the next two phases, we will show that the signal term $U_t W_t W_t^T U_t^T$ converges towards XX^T , whereas the spectral norm of the noise term, i.e. $\|U_t W_{t,\perp}\|$, stays small. For that, we show that throughout this process the columns of the matrices X and $U_t W_t$ stay approximately aligned, i.e., the angle $\|V_{X^\perp}^T V_{U_t W_t}\|$ stays small. This latter property also ensures that after the spectral phase the iterates are not too close to well known saddle points of the optimization landscape (it is known that this problem may have degenerate saddle points at a point \bar{U} obeying $\text{rank}(\bar{U}) < r_*$ [33]). See Figure 4 for a depiction of the gradient flows of the landscape when $r_* = 1$.

Next we sketch the proofs of Phase II and Phase III in more detail.

Phase II: In this phase, we will show that the minimal singular value of the signal term, $\sigma_{\min}(U_t W_t)$ grows exponentially, until it holds that $\sigma_{\min}(U_t W_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$. To this aim, we show that

$$\sigma_{\min}(U_{t+1} W_{t+1}) \geq \sigma_{\min}(V_X^T U_{t+1}) \geq \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{4} \mu \sigma_{\min}^2(X) - \mu \sigma_{\min}^2(V_X^T U_t)\right)$$

holds under suitable assumptions (see Lemma 9.1). In order to show that the spectral norm of the noise term $\|U_t W_{t,\perp}\|$ grows much slower than $\sigma_{\min}(U_{t+1} W_{t+1})$, we establish the inequality

$$\begin{aligned} & \|U_{t+1} W_{t+1,\perp}\| \\ & \leq \left(1 - \frac{\mu}{2} \|U_t W_{t,\perp}\|^2 + 36\mu \|V_{X^\perp}^T V_{U_t W_t}\|^2 \|X\|^2 + 3\mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|\right) \|U_t W_{t,\perp}\| \end{aligned} \quad (12)$$

(see Lemma 9.2). The next inequality (see Lemma 9.3) shows that $\|V_{X^\perp}^T V_{U_t W_t}\|$ stays sufficiently

small

$$\begin{aligned} & \|V_{X^\perp}^T V_{U_{t+1}W_{t+1}}\| \\ & \leq \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) \|V_{X^\perp}^T V_{U_t W_t}\| + 100\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| + 500\mu^2 \|XX^T - U_t U_t^T\|^2. \end{aligned}$$

As mentioned above, this implies in particular, that U_t stays sufficiently far away from saddle points \bar{U} , which are rank-deficient, e.g., $\text{rank}(\bar{U}) < r_*$.

Phase III: After we have shown that $\sigma_{\min}(U_t W_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$ holds for some t , we enter the *local refinement phase*. We start by observing that the error $\|XX^T - U_t U_t^T\|_F$ can be decomposed into two summands, i.e.

$$\|U_t U_t^T - XX^T\|_F \leq 4\|V_X^T (XX^T - U_t U_t^T)\|_F + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|_F. \quad (13)$$

(see Lemma B.4). We will bound the second summand by using inequality (12), which is also valid for the third phase. We will show that the first summand decreases at a linear rate. For that, we establish the inequality

$$\begin{aligned} & \|V_X^T (XX^T - U_{t+1} U_{t+1}^T)\|_F \\ & \leq \left(1 - \frac{\mu}{200}\sigma_{\min}^2(X)\right) \|V_X^T (XX^T - U_t U_t^T)\|_F + \mu \frac{\sigma_{\min}^2(X)}{100} \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|_F. \end{aligned}$$

Hence, by using inequality (13) we will be able to show that $\|U_t U_t^T - XX^T\|_F$ is decreasing, as long as the spectral norm of the noise term stays sufficiently small.

6 Numerical experiments

In this section, we perform several numerical experiments to corroborate our theoretical results.

Experimental setup. For the experiments we set the ground truth matrix $X \in \mathbb{R}^{n \times r_*}$ to be a random orthogonal matrix with $n = 200$ and $r_* = 5$. Moreover, we use $m = 10nr_* = 50n$ random Gaussian measurements. The initialization U is chosen as in Theorem 3.3 and we use a step size of $\mu = 1/4$ which is consistent with these theorems. We note that while all experimental depictions are based on a single trial, in line with the NeuRIPS guidelines we have drawn these curves multiple times (not depicted) and the behavior of the plots do not change.

Depiction of the three phases and the role of overparameterization. In our first experiment, we want to examine how increasing the number of parameters via increasing the number of the columns r of the matrix $U_t \in \mathbb{R}^{n \times r}$, affects the spectral phase. To this aim we set the scale of initialization to $\alpha = 1/(70n^2)$. Recall from Section 5 that L denotes the subspace spanned by the eigenvectors corresponding to the leading r_* singular values of $\mathcal{A}^* \mathcal{A}(XX^T)$ and L_t denotes the subspace spanned by the left-singular vectors corresponding to the largest r_* singular values of U_t .

Spectral phase and alignment under different levels of overparameterization. First, we examine how the angle between these two subspaces (i.e. $\|V_{L^\perp}^T V_{L_t}\|$) changes in the first few iterations. We depict the results for different r in Figure 5a. We see that in the first few iterations, i.e. in the spectral phase, this angle converges towards zero. This confirms the main conclusion of this paper

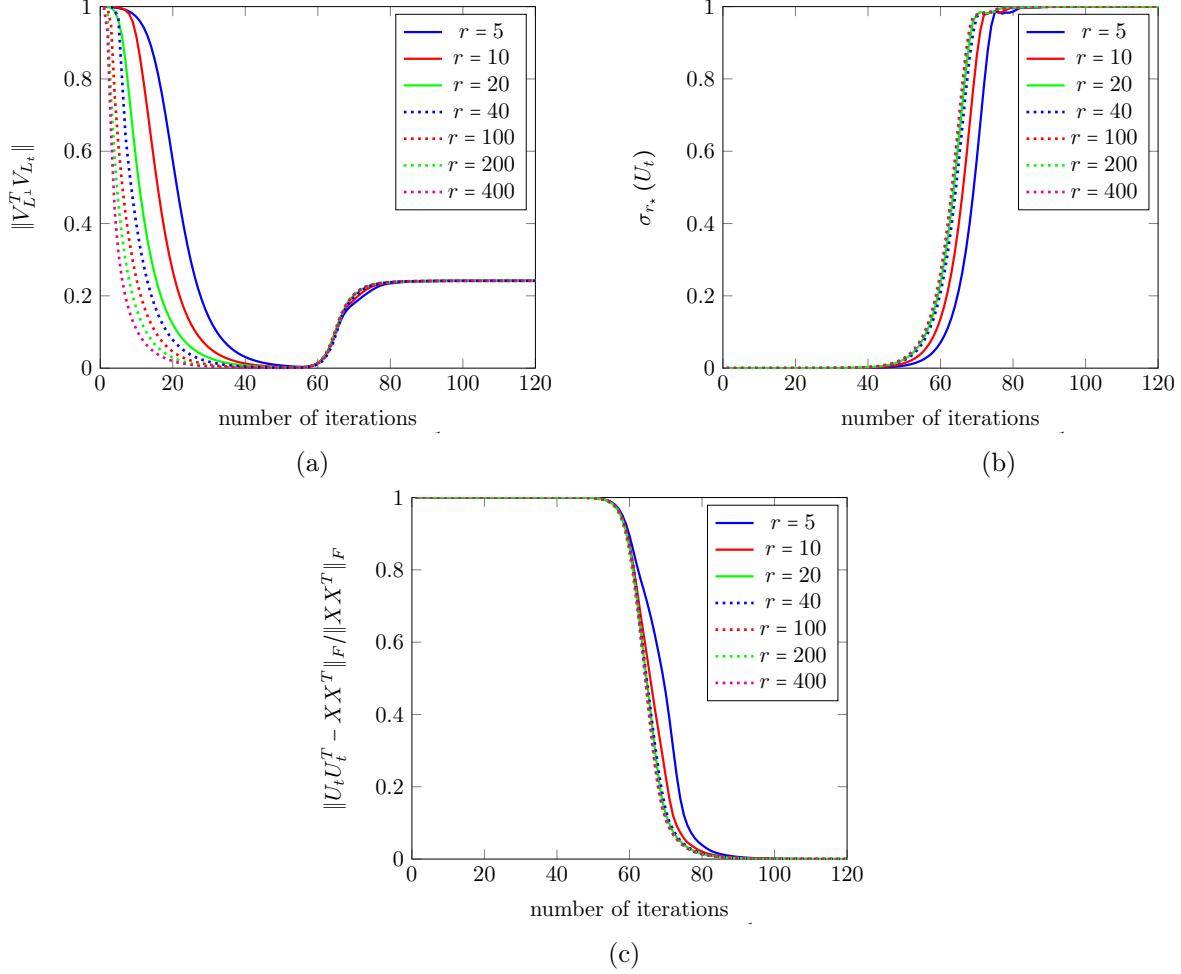


Figure 5: Impact of different levels of overparameterization in the spectral phase on (a) the angle $\|V_{L^\perp}^T V_{L_t}\|$ and (b) the r_* th largest singular value, (c) the trajectory of the (normalized) test error $\|U_t U_t^T - X X^T\|_F / \|X X^T\|_F$.

that the first few iterations of gradient descent from small random initialization indeed behaves akin to running power method for spectral initialization. This experiment also shows that changing the number of columns r of U_t has an interesting effect on the spectral phase. In particular, increasing r allows the gradient descent algorithm to learn the subspace L with fewer iterations, i.e. $\|V_{L^\perp}^T V_{L_t}\|$ becomes small with fewer iterations. This is in accordance with our theory for $r_* \leq r \leq n$ (see, for example, the first summand on the right-hand side of equation (8)), where we show that more overparameterization allows gradient descent to leave the spectral phase earlier. Interestingly, this improvement continues to hold even when increasing r beyond n allowing for even faster convergence of $\|V_{L^\perp}^T V_{L_t}\|$. This holds even though in this case the rank of U_0 is still not larger than n . One potential explanation for this phenomenon might be that for such a choice of r the matrix U_0 is better conditioned.

Growth of $\sigma_{r_}(U_t)$ and saddle avoidance.* In Figure 5b we depict how $\sigma_{r_*}(U_t)$ grows during the training for different choices of r . We see that the curves look similar, although for smaller r the

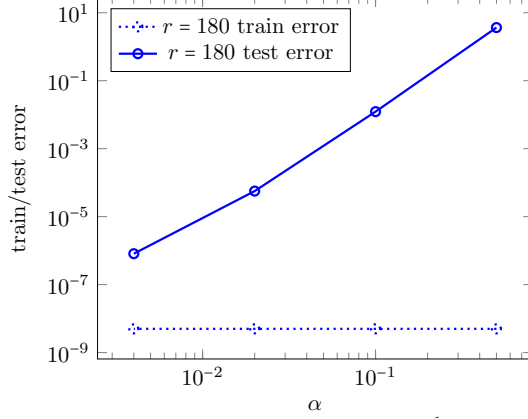


Figure 6: Relative test error $\frac{\|U_t U_t^T - X X^T\|_F}{\|X X^T\|_F}$ for different scales of initialization α .

growth phase sets in at a slightly later time. This is due to the fact that for smaller r , as we have seen in Figure 5a, Phase I, the spectral phase takes longer to complete.

Evolution of the test error and the refinement phase. Similarly, in Figure 5c we depict how the (normalized) test error $\|U_t U_t^T - X X^T\|_F / \|X X^T\|_F$ evolves during the training for different choices of r . We observe that for smaller r the third phase sets in slightly later. Again, this is due to the fact that for smaller r the spectral phase takes slightly longer to complete.

Test error under different scales of initialization. In the next experiment, we focus on understanding how the scale of initialization α affects the generalization error $\|U_t U_t^T - X X^T\|_F^2$. For that, we set $r = 180$ and run gradient descent with for different choices of α . We stop as soon as the training error becomes small ($f(U_t) \leq 0.5 \cdot 10^{-9}$). We depict the results in Figure 6. We see that the test error decreases as α decreases. In particular, this figure indicates that the test error depends polynomially on the scale of initialization α . This is in line with our theory, where we also show that the test error decreases at least with the rate $\alpha^{21/16}$ (see inequality (5) in Theorem 3.3).

Change of test and train error during the training. In the next experiment, we set $r = 180$ and examine how the test error $\|U_t U_t^T - X X^T\|_F^2$ and the train error $f(U_t)$ changes throughout training and, in particular, how this depends on the scale of initialization. To this aim, we run gradient descent with $4 \cdot 10^5$ iterations. We see that for a small scale of initialization, $\alpha = 10^{-3}$, which is the scenario studied in this paper, both test error and train error decrease throughout the training. In particular, the test error is still decreasing even when the train error is already small. For large scale of initialization $\alpha = 0.5$, we observe a very different behaviour. We see that the train error converges with linear rate until machine precision is reached. However, the test error barely changes throughout the training. This scale of initialization corresponds to the *lazy training regime* or *NTK-regime*, where the parameters stay close to the initialization during the training.

Number of iterations until convergence: In the last experiment, we set $\alpha = 10^{-3}$ and examine how many iterations are needed until the test error $\|U_t U_t^T - X X^T\|_F^2$ falls below a certain threshold of 10^{-4} for different values of r obeying $5 \leq r \leq 30$. For each choice of r we run the experiment ten times and then average the number of iterations for each choice of r . The results are depicted in

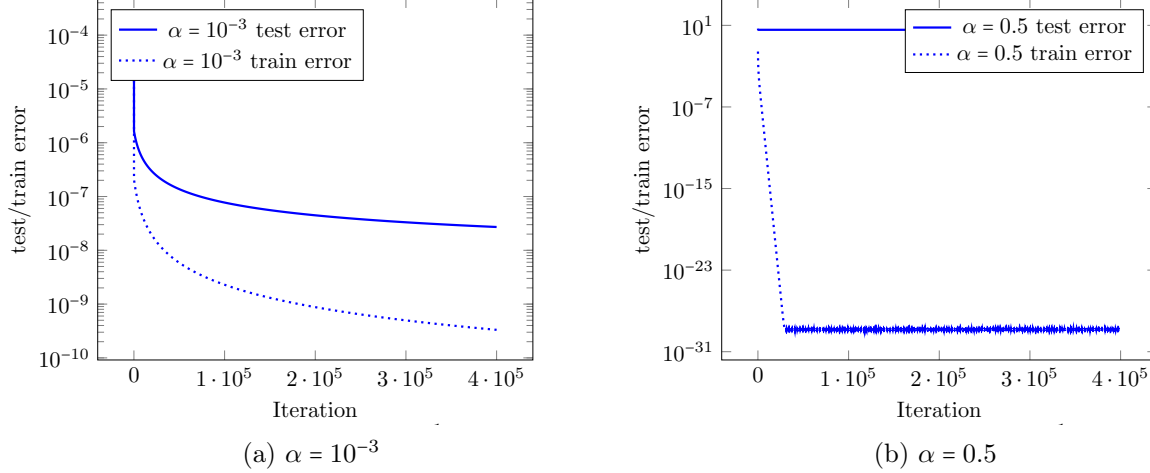


Figure 7: Change of test error $\|U_t U_t^T - X X^T\|_F^2$ and train error $f(U_t)$ for (a) small and (b) large α during training.

Figure 8. We observe that increasing the number of columns r from 5 to 10, i.e., a small amount of overparameterization, decreases the number of iterations needed. After that the number of iterations needed stays roughly constant. This observation is in line with Figure 5, where we have seen that overparameterization leads to fast decrease of the test error in the spectral phase (with diminishing speedup as r becomes larger and larger) without affecting the other two phases.

7 Preliminaries

Before we are going into the details of the proof, we are collecting some useful definitions.

7.1 Notation

For any matrix $A \in \mathbb{R}^{n_1 \times n_2}$ we denote its spectral norm by $\|A\|$ and the Frobenius norm by $\|A\|_F = \sqrt{\text{Tr}(AA^T)}$. By $\|A\|_*$ we denote its nuclear norm, i.e. the sum of the singular values. Moreover, for two symmetric matrices $A, B \in S^d$ we define the Hilbert-Schmidt inner product by $\langle A, B \rangle = \text{Tr}(AB)$. For a positive semidefinite matrix A we denote its square root by $A^{1/2}$, i.e., the unique positive semidefinite matrix B for which it holds that $B^2 = A$. We also set $A^{-1/2} = (A^{1/2})^{-1}$.

For any matrix $A \in \mathbb{R}^{d_1 \times d_2}$ we will denote its singular value decomposition by $A = V_A \Sigma_A W_A^T$ with $V_A \in \mathbb{R}^{d_1 \times \tilde{r}}$, $W_A \in \mathbb{R}^{d_2 \times \tilde{r}}$, $\Sigma_A \in \mathbb{R}^{\tilde{r} \times \tilde{r}}$, where \tilde{r} denotes the rank of A . Moreover, by $V_{A^\perp} \in \mathbb{R}^{(d_1 - \tilde{r}) \times d_1}$ we denote an orthogonal matrix, whose column span is orthogonal to the column span of the matrix V_A . Similarly, if $U \subset \mathbb{R}^n$ is a subspace of dimension \tilde{r} , we will denote by $V_U \in \mathbb{R}^{n \times \tilde{r}}$ a matrix, whose column span is the subspace U . Similarly as before, we will denote by $V_{U^\perp} \in \mathbb{R}^{n \times (n - \tilde{r})}$ a matrix whose column span is orthogonal to the column span of U .

We will measure the angle between two subspaces $U_1, U_2 \subset \mathbb{R}^n$ by $\|V_{U_1}^T V_{U_2}\|$. Moreover, we will also several times rely on the well-known identity (see, e.g., [57, Section 2])

$$\|V_{U_1}^T V_{U_2}\| = \|V_{U_1} V_{U_1}^T - V_{U_2} V_{U_2}^T\|.$$

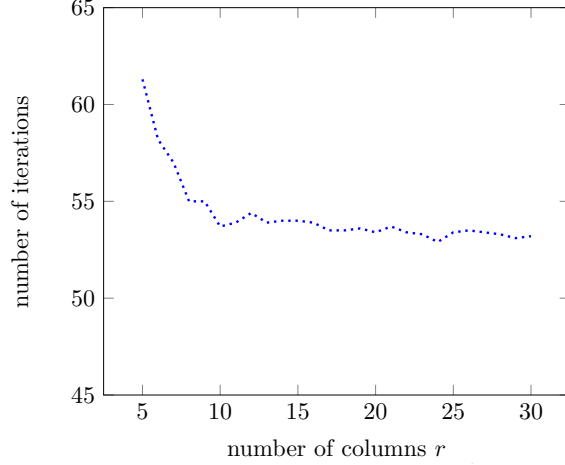


Figure 8: Number of iterations required for the test error to fall below 10^{-4} for different levels of overparameterization.

7.2 Restricted isometry property and related properties

As discussed in Section 2, we are going to assume that the measurement operator \mathcal{A} satisfies the restricted isometry property. However, as it turns out, the following two slightly weaker properties will suffice for our proof.

Definition 7.1. *The measurement operator $\mathcal{A}: S^n \rightarrow \mathbb{R}^m$ satisfies the spectral-to-spectral restricted isometry property of rank r with constant $\delta > 0$, if it holds for all symmetric matrices Z of rank at most r that*

$$\|(Id - \mathcal{A}^* \mathcal{A})(Z)\| \leq \delta \|Z\|.$$

Definition 7.2. *The measurement operator $\mathcal{A}: S^n \rightarrow \mathbb{R}^m$ satisfies the spectral-to-nuclear restricted isometry property with constant $\delta > 0$, if it holds for all symmetric matrices Z that*

$$\|(Id - \mathcal{A}^* \mathcal{A})(Z)\| \leq \delta \|Z\|_*.$$

The following lemma shows that these two properties are induced by the standard restricted isometry property (Definition 3.1).

Lemma 7.3. *Let $\mathcal{A}: S^n \rightarrow \mathbb{R}^m$ be a linear measurement operator. Then the following two statements hold.*

1. *Suppose that \mathcal{A} has the restricted isometry property as in (3) for all matrices of rank $r + 1$ with constant $\delta_1 < 1$. Then \mathcal{A} has the spectral-to-spectral restricted isometry property of rank r with constant $\sqrt{r}\delta_1$.*
2. *Suppose that \mathcal{A} has the restricted isometry property as in (3) for all matrices of rank 2 with constant $\delta_2 < 1$. Then \mathcal{A} has the spectral-to-nuclear restricted isometry property with constant δ_2 .*

Proof. From [30] it follows that if \mathcal{A} has the restricted isometry property of rank $r + r'$ with constant $\delta < 0$, then it holds for all matrices Z with rank at most r and all matrices Y with rank at most r' that

$$\left| \langle (\text{Id} - \mathcal{A}^* \mathcal{A})(Z), Y \rangle \right| = \left| \langle \mathcal{A}^* \mathcal{A}(Z), Y \rangle - \langle Z, Y \rangle \right| \leq \delta_1 \|Z\|_F \|Y\|_F. \quad (14)$$

In order to prove the first statement it suffices to note that there is a vector $v \in \mathbb{R}^n$ with $\|v\|_{\ell_2} = 1$ such that

$$\|(\text{Id} - \mathcal{A}^* \mathcal{A})(Z)\| = \langle (\text{Id} - \mathcal{A}^* \mathcal{A})(Z), vv^T \rangle.$$

The claim follows from (14), $\|v\|_{\ell_2} = 1$, and from $\|Z\|_F \leq \sqrt{r} \|Z\|$.

In order to show the second claim consider the eigenvalue decomposition $Z = \sum_{i=1}^n \lambda_i v_i v_i^T$. We compute that

$$\begin{aligned} \|(\text{Id} - \mathcal{A}^* \mathcal{A})(Z)\| &\leq \sum_{i=1}^n |\lambda_i| \|(\text{Id} - \mathcal{A}^* \mathcal{A})(v_i v_i^T)\| \\ &\leq \delta_2 \sum_{i=1}^n |\lambda_i| \|v_i v_i^T\| \\ &= \|Z\|_*, \end{aligned}$$

where in the second inequality we used the spectral-to-spectral restricted isometry property, which holds due to the first part of this proof. This finishes the proof of the second statement. \square

8 Analysis of the spectral phase

In the following we will provide an analysis of the spectral phase, where the proofs of the technical lemmas are deferred to Appendix A. Our first goal is to show that in the first few iterations U_t can be approximated by

$$\tilde{U}_t := \left(\text{Id} + \underbrace{\mu \mathcal{A}^* \mathcal{A}(XX^T)}_{=: M} \right)^t U_0 =: Z_t U_0,$$

where we have set

$$Z_t = (\text{Id} + \mu M)^t.$$

Next, we define

$$t^* := \min \{i \in \mathbb{N} : \|\tilde{U}_{i-1} - U_{i-1}\| > \|\tilde{U}_{i-1}\|\}.$$

The next lemma shows how well U_t can be approximated by \tilde{U}_t for $t \leq t^*$. To formulate it, we set $E_t = U_t - \tilde{U}_t$.

Lemma 8.1. *Suppose that \mathcal{A} satisfies the rank-1 RIP with constant δ_1 . For all integers t such that $1 \leq t \leq t^*$ it holds that*

$$\|E_t\| = \|U_t - \tilde{U}_t\| \leq \frac{4}{\lambda_1(M)} \alpha^3 \min\{r; n\} (1 + \delta_1) (1 + \mu \lambda_1(M))^{3t} \|U\|^3.$$

The next lemma gives a lower bound for t^* . In particular, this shows how long the approximation in Lemma 8.1 is valid.

Lemma 8.2. Let \tilde{U}_t be as defined before and consider the eigenvalue decomposition $\mathcal{A}^* \mathcal{A} (XX^T) = \sum_{i=1}^n \lambda_i v_i v_i^T$. Then we have that

$$t^* \geq \left\lceil \frac{\ln \left(\frac{\lambda_1(M)}{4\alpha^2(1+\delta_1)\|U\|^3} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r;n\}} \right) \right) \right\rceil}{2 \ln(1 + \mu \lambda_1(M))}.$$

Next, recall the relation

$$U_t = \tilde{U}_t + E_t = Z_t U_0 + E_t$$

and denote by L the subspace spanned by the eigenvectors, which correspond to the largest r_* eigenvalues of the matrix $M = \mathcal{A}^* \mathcal{A} (XX^T)$. Note that L is also the subspace spanned by the eigenvectors corresponding to the largest r_* eigenvalues of the matrix Z_t . Denote by L_t the subspace spanned by the left-singular vectors of $U_t = Z_t U_0 + E_t$, which corresponds to the largest r_* singular values.

Since Z_t is computed via a power method we expect for t large enough that $\lambda_{r_*}(Z_t) \gg \lambda_{r_*+1}(Z_t)$. Moreover, if, in addition, $\|E_t\|$ is sufficiently small, we expect in this case that the subspace L is aligned with the subspace L_t . This is made precise by the following lemma.

Lemma 8.3. Let $Z_t \in \mathbb{R}^{n \times n}$ be a symmetric matrix. Let $U \in \mathbb{R}^{n \times r}$ be a matrix and let $E_t \in \mathbb{R}^{n \times r}$. Set $U_0 = \alpha U$ for some $\alpha > 0$. Moreover, assume that

$$\sigma_{r_*+1}(Z_t) \|U\| + \frac{\|E_t\|}{\alpha} < \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U). \quad (15)$$

Then the following three inequalities hold.

$$\sigma_{r_*}(Z_t U_0 + E_t) \geq \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) - \|E_t\|, \quad (16)$$

$$\sigma_{r_*+1}(Z_t U_0 + E_t) \leq \alpha \sigma_{r_*+1}(Z_t) \|U\| + \|E_t\|, \quad (17)$$

$$\|V_{L_t}^T V_{L_t}\| \leq \frac{\alpha \sigma_{r_*+1}(Z_t) \|U\| + \|E_t\|}{\alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) - \alpha \sigma_{r_*+1}(Z_t) \|U\| - \|E_t\|}. \quad (18)$$

Recall that we are interested in bounds for the quantities $\sigma_{r_*}(U_t W_t)$, $\|V_{X^\perp}^T V_{U_t W_t}\|$, and $\|U_t W_{t,\perp}\|$, i.e., properties of the signal and noise term. However, in the lemma above, we have obtained instead bounds for $\sigma_{r_*}(U_t)$, $\|V_{X^\perp}^T V_{L_t}\|$, and $\sigma_{r_*+1}(U_t)$, i.e. for the singular value decomposition of U_t . However, if $\|V_{X^\perp}^T V_{L_t}\|$ is small, these quantities are closely related to each other, as the next lemma shows.

Lemma 8.4. Assume that $\|V_{X^\perp}^T V_{L_t}\| \leq \frac{1}{8}$ for some $t \geq 1$. Then it holds that

$$\sigma_{r_*}(U_t W_t) \geq \frac{1}{2} \sigma_{r_*}(U_t), \quad (19)$$

$$\|V_{X^\perp}^T V_{U_t W_t}\| \leq 7 \|V_{X^\perp}^T V_{L_t}\|, \quad (20)$$

$$\|U_t W_{t,\perp}\| \leq 2 \sigma_{r_*+1}(U_t). \quad (21)$$

By combining Lemma 8.3 and Lemma 8.4, we obtain the following technical result.

Lemma 8.5. *Let XX^T be a low-rank matrix of rank r_* . Assume that*

$$M := \mathcal{A}^* \mathcal{A}(XX^T) = XX^T + \tilde{E},$$

with $\|\tilde{E}\| \leq \delta \lambda_{r_}(XX^T)$, where $\delta \leq \tilde{c}_1$ where $\tilde{c}_1 > 0$ is a sufficiently small absolute constant. Furthermore, set $E_t = U_t - \tilde{U}_t$. Moreover, assume that*

$$\gamma := \frac{\alpha \sigma_{r_*+1}(Z_t) \|U\| + \|E_t\|}{\alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U)} \leq \tilde{c}_2 \kappa^{-2}, \quad (22)$$

where $\tilde{c}_2 > 0$ is a sufficiently small, absolute constant. Then it holds that

$$\sigma_{\min}(U_t W_t) \geq \frac{\alpha}{4} \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U), \quad (23)$$

$$\|U_t W_{t,\perp}\| \leq \frac{\kappa^{-2}}{8} \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U), \quad (24)$$

$$\|V_{X^\perp}^T V_{U_t W_t}\| \leq 56(\delta + \gamma). \quad (25)$$

In order to utilize Lemma 8.5, we need to insert bounds for the approximation error $\|E_t\|$, which we have derived in the Lemmas 8.1 and 8.2. This yields the following lemma.

Lemma 8.6. *Fix a sufficiently small constant $c > 0$. Let $U \in \mathbb{R}^{n \times r}$. Assume that \mathcal{A} has the spectral-to-nuclear restricted isometry property for some constant $\delta_1 < 1$. Moreover, assume that*

$$M := \mathcal{A}^* \mathcal{A}(XX^T) = XX^T + \tilde{E},$$

with $\|\tilde{E}\| \leq \delta \lambda_{r_}(XX^T)$, where $\delta \leq c_1 \kappa^{-2}$. Denote by L the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $\mathcal{A}^* \mathcal{A}(XX^T)$. Let $U_0 = \alpha U$, where*

$$\alpha^2 \leq \frac{c_2 \|X\|^2}{32 \min\{r; n\} \kappa \|U\|^3} \left(\frac{2\kappa^2 \|U\|}{c_3 \sigma_{\min}(V_L^T U)} \right)^{-12\kappa^2} \min\left(\sigma_{\min}(V_L^T U); \|U^T v_1\|_{\ell_2}\right), \quad (26)$$

where v_1 denotes the eigenvector corresponding to a leading eigenvalue of the matrix $\mathcal{A}^ \mathcal{A}(XX^T)$. Assume that the step size satisfies $\mu \leq c_2 \kappa^{-2} \|X\|^{-2}$. Then after*

$$t_* \asymp \frac{1}{\mu \sigma_{r_*}(X)^2} \cdot \ln \left(\frac{2\kappa^2 \|U\|}{c_3 \sigma_{\min}(V_L^T U)} \right)$$

iterations it holds that

$$\|U_{t_*}\| \leq 3\|X\|, \quad (27)$$

$$\sigma_{\min}(U_{t_*} W_{t_*}) \geq \frac{\alpha \beta}{4}, \quad (28)$$

$$\|U_{t_*} W_{t_*,\perp}\| \leq \frac{\kappa^{-2}}{8} \alpha \beta, \quad (29)$$

$$\|V_{X^\perp}^T V_{U_{t_*} W_{t_*}}\| \leq c \kappa^{-2}, \quad (30)$$

where $\beta > 0$ satisfies

$$\sigma_{\min}(V_L^T U) \lesssim \beta \lesssim \sigma_{\min}(V_L^T U) \left(\frac{\kappa^2 \|U\|}{c_3 \sigma_{\min}(V_L^T U)} \right)^2 = \frac{\kappa^4 \|U\|^2}{c_3^2 \sigma_{\min}(V_L^T U)}. \quad (31)$$

Here $c_1, c_2, c_3 > 0$ are absolute constants only depending on the choice of c .

Note that the result above holds for any initialization U . To complete the proof we are going to utilize the fact that U is a random matrix with Gaussian entries. This yields the following lemma, which is the main result of this section.

Lemma 8.7. *Fix a sufficiently small constant $c > 0$. Let $U \in \mathbb{R}^{n \times r}$ be a random matrix with i.i.d. entries with distribution $\mathcal{N}(0, 1/\sqrt{r})$ and let $0 < \varepsilon < 1$. Assume that \mathcal{A} has the spectral-to-nuclear restricted isometry property for some constant $\delta_1 < 1$. Moreover, assume that*

$$M := \mathcal{A}^* \mathcal{A}(XX^T) = XX^T + \tilde{E},$$

with $\|\tilde{E}\| \leq \delta \lambda_{r_*}(XX^T)$, where $\delta \leq c_1 \kappa^{-2}$. Let $U_0 = \alpha U$, where

$$\alpha^2 \lesssim \begin{cases} \frac{\sqrt{\min\{r;n\}} \|X\|^2}{\kappa n^{3/2}} \left(2\kappa^2 \sqrt{\frac{n}{\min\{r;n\}}} \right)^{-12\kappa^2} & \text{if } r \geq 2r_* \\ \frac{\|X\|^2}{n^{3/2}\kappa} \left(\frac{2\kappa^2 \sqrt{rn}}{\varepsilon} \right)^{-12\kappa^2} \varepsilon & \text{if } r < 2r_* \end{cases}. \quad (32)$$

Assume that the step size satisfies $\mu \leq c_2 \kappa^{-2} \|X\|^2$. Then with probability at least $1 - p$, where

$$p = \begin{cases} O(\exp(-\tilde{c}r)) & \text{if } r \geq 2r_* \\ (\tilde{C}\varepsilon)^{r-r_*+1} + \exp(-\tilde{c}r) & \text{if } r < 2r_* \end{cases}$$

the following statement holds. After

$$t_* \lesssim \begin{cases} \frac{1}{\mu \sigma_{r_*}(X)^2} \cdot \ln \left(2\kappa^2 \sqrt{\frac{n}{\min\{r;n\}}} \right) & \text{if } r \geq 2r_* \\ \frac{1}{\mu \sigma_{r_*}(X)^2} \cdot \ln \left(\frac{2\kappa^2 \sqrt{rn}}{\varepsilon} \right) & \text{if } r < 2r_* \end{cases}$$

iterations it holds that

$$\|U_{t_*}\| \leq 3\|X\|, \quad (33)$$

$$\sigma_{\min}(U_{t_*} W_{t_*}) \geq \frac{\alpha\beta}{4}, \quad (34)$$

$$\|U_{t_*} W_{t_*, \perp}\| \leq \frac{\kappa^{-2}}{8} \alpha\beta, \quad (35)$$

$$\|V_{X^\perp}^T V_{U_{t_*} W_{t_*}}\| \leq c\kappa^{-2}, \quad (36)$$

where $\beta > 0$ satisfies

$$\beta \lesssim \begin{cases} \frac{n\kappa^4}{c_3^2(\min\{r;n\})} & \text{if } r \geq 2r_* \\ \frac{n\kappa^4}{c_3^2\varepsilon} & \text{if } r < 2r_* \end{cases}$$

as well as

$$\beta \gtrsim \begin{cases} 1 & \text{if } r \geq 2r_\star \\ \frac{\varepsilon}{r} & \text{if } r < 2r_\star \end{cases}.$$

Here $c_1, c_2, c_3 > 0$ are absolute constants only depending on the choice of c . Moreover, $\tilde{C}, \tilde{c} > 0$ are absolute numerical constants.

The proof of Lemma 8.7 requires the following theorem, which gives a non-asymptotic lower bound for the smallest singular value of a Gaussian matrix.

Theorem 8.8. [58] *Let $G \in \mathbb{R}^{r_\star \times r}$ with $r_\star \leq r$ and i.i.d. Gaussian entries with distribution $\mathcal{N}(0, 1/\sqrt{r})$. Then for every $\varepsilon > 0$ we have with probability at least $1 - (C\varepsilon)^{r-r_\star+1} - \exp(-cr)$ that*

$$\sigma_{\min}(G) \geq \varepsilon \frac{\sqrt{r} - \sqrt{r_\star - 1}}{\sqrt{r}}.$$

The constants $C, c > 0$ are universal.

With this theorem in place we can prove Lemma 8.7.

Proof of Lemma 8.7. We will deduce this statement from Lemma 8.6. For that we need to estimate $\|U\|$, $\sigma_{\min}(V_L^T U)$, and $\|U^T v_1\|_{\ell_2}$. It is well-known (see, e.g. [32, Section 4]) that with probability at least $1 - O(\exp(-c \max\{r; n\}))$ it holds that

$$\|U\| \lesssim \sqrt{\max\{r; n\}/r} = \sqrt{\frac{n}{\min\{r; n\}}}. \quad (37)$$

Next, note that again due to rotation invariance of the Gaussian measure the vector $U^T v_1 \in \mathbb{R}^r$ has i.i.d. entries with distribution $\mathcal{N}(0, 1/\sqrt{r})$. Hence, with probability at least $1 - O(\exp(-cr))$ it holds that

$$\|U^T v_1\|_{\ell_2} \asymp 1. \quad (38)$$

Next, we note that due to rotation invariance of the Gaussian distribution the matrix $V_L^T U \in \mathbb{R}^{r_\star \times r}$ has i.i.d. entries with distribution $\mathcal{N}(0, 1/\sqrt{r})$. Moreover, note that using the elementary inequality $\sqrt{1-x} \leq 1 - \frac{1}{2x}$ we obtain that

$$\frac{\sqrt{r} - \sqrt{r_\star - 1}}{\sqrt{r}} \geq \frac{\sqrt{r} - \sqrt{r_\star} \left(1 - \frac{1}{2r_\star}\right)}{\sqrt{r}} \gtrsim \begin{cases} 1 & \text{if } r \geq 2r_\star \\ \frac{1}{r} & \text{else} \end{cases}. \quad (39)$$

In order to proceed we are going to distinguish the following two cases.

Case 1: $r \geq 2r_\star$

Note that by choosing $\varepsilon > 0$ appropriately, we obtain from Theorem 8.8 combined with inequality (39) that with probability at least $1 - O(\exp(-cr))$ it holds that

$$\sigma_{\min}(V_L^T U) \gtrsim 1. \quad (40)$$

By combining the inequalities (37), (38), and (40) with Lemma 8.6 the claim follows in the case that $r \geq 2r_\star$.

Case 2: $r_* \leq r \leq 2r_*$

Similar to the first case, we note that by choosing $\varepsilon > 0$ appropriately, we obtain by applying Theorem 8.8 combined with inequality (39) that with probability at least $1 - (C\varepsilon)^{r-r_*+1} - \exp(-cr)$ it holds that

$$\sigma_{\min}(V_L^T U) \gtrsim \frac{\varepsilon}{r}. \quad (41)$$

By combining the inequalities (37), (38), and (41) with Lemma 8.6 the claim follows. \square

9 Analysis of the saddle avoidance and refinement phases

Before stating and proving the main result of this section, Theorem 9.6, we will first collect some useful lemmas. Their proofs are deferred to Appendix B.

In Phase II we will show that $\sigma_{\min}(U_t W_t)$ grows until it reaches $\sigma_{\min}(U_t W_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$. For that, we note

$$\begin{aligned} \sigma_{\min}(U_t W_t) &\stackrel{(a)}{=} \sigma_{\min}(U_t W_t W_t^T) \\ &\geq \sigma_{\min}(V_X^T U_t W_t W_t^T) \\ &\stackrel{(b)}{=} \sigma_{\min}(V_X^T U_t), \end{aligned}$$

where (a) and (b) follow from the definition of W_t . Hence, in order to show that $\sigma_{\min}(U_t W_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$ it suffices to show that $\sigma_{\min}(V_X^T U_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$. For that, we will use the next lemma, which shows that $\sigma_{\min}(V_X^T U_t)$ grows exponentially.

Lemma 9.1. *Assume that $\mu \leq c\|X\|^{-2}\kappa^{-2}$, $\|U_t\| \leq 3\|X\|$, and that $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$. Moreover, suppose that*

$$\|(\mathcal{A}^* \mathcal{A} - Id)(XX^T - U_t U_t^T)\| \leq c\sigma_{\min}^2(X). \quad (42)$$

Furthermore, assume that $V_X^T U_t$ has full rank. Then it holds that

$$\sigma_{\min}(V_X^T U_{t+1}) \geq \sigma_{\min}(V_X^T U_{t+1} W_t) \geq \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{4}\mu\sigma_{\min}^2(X) - \mu\sigma_{\min}^2(V_X^T U_t)\right).$$

Here $c > 0$ is constant, which is chosen small enough.

The next lemma will allow us to show that the noise term $\|U_t W_{t,\perp}\|$ is growing slower than $\sigma_{\min}(V_X^T U_{t+1})$.

Lemma 9.2. *Assume that $\mu \leq c \min\{\|X\|^{-2}; \|(\mathcal{A}^* \mathcal{A} - Id)(XX^T - U_t U_t^T)\|^{-1}\}$ and that $\|U_t\| \leq 3\|X\|$. Moreover, suppose that $V_X^T U_{t+1} W_t$ has full rank and that $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$. Then it holds that*

$$\begin{aligned} &\|U_{t+1} W_{t+1,\perp}\| \\ &\leq \left(1 - \frac{\mu}{2}\|U_t W_{t,\perp}\|^2 + 36\mu\|V_{X^\perp}^T V_{U_t W_t}\|^2\|X\|^2 + 3\mu\|(\mathcal{A}^* \mathcal{A} - Id)(XX^T - U_t U_t^T)\|\right)\|U_t W_{t,\perp}\|. \end{aligned}$$

Here, $c > 0$ is an absolute constant chosen small enough.

The next lemma shows that the angle between the column space of the signal term $U_t W_t$ and column space of X stays sufficiently small.

Lemma 9.3. Assume that $\|U_t W_{t,\perp}\| \leq 2\sigma_{\min}(U_t W_t)$ and $\|U_t\| \leq 3\|X\|$ holds. Moreover, assume that

$$\|(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \leq c\sigma_{\min}^2(X), \quad (43)$$

$$\|V_{X^\perp}^T V_{U_t W_t}\| \leq c, \quad (44)$$

$$\mu \leq c\kappa^{-2}\|X\|^{-2}, \quad (45)$$

$$\|U_t W_{t,\perp}\| \leq c\kappa^{-2}\|X\|. \quad (46)$$

where $c > 0$ is a small enough constant. Then it holds that

$$\begin{aligned} & \|V_{X^\perp}^T V_{U_{t+1} W_{t+1}}\| \\ & \leq \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) \|V_{X^\perp}^T V_{U_t W_t}\| + 100\mu \|(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| + 500\mu^2 \|XX^T - U_t U_t^T\|^2. \end{aligned}$$

The next lemma will show that we have $\|U_t\| \leq 3\|X\|$ for all t , a technical assumption which is needed in the above lemmas.

Lemma 9.4. Assume that $\|U_t\| \leq 3\|X\|$, $\mu \leq \frac{1}{27\|X\|^2}$, and

$$\|(\mathcal{A}^* \mathcal{A} - Id)(XX^T - U_t U_t^T)\| \leq \|X\|^2.$$

Then it also holds that $\|U_{t+1}\| \leq 3\|X\|$.

With these lemmas in place, we will be able to show that $\sigma_{\min}(U_t W_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}$ holds after sufficiently many iterations. Hence, we can enter Phase III, the *local refinement phase*.

The next lemma is concerned with this third phase. It shows that $U_t W_t W_t^T U_t^T$ converges towards XX^T , when projected onto the column space of X . We are going to provide a somewhat more general version of the lemma than what is needed in the proofs of our main results, since it may be of independent interest. For that, let $\|\cdot\|$ be a matrix norm, which satisfies $\|ABC\| \leq \|A\|\|B\|\|C\|$ for all matrices A, B, C . Furthermore, we assume that $\|A\| = \|A^T\|$ for all matrices A . For example, this property is fulfilled by all Schatten- p norms.

Lemma 9.5. Assume that $\|U_t\| \leq 3\|X\|$ and that $\sigma_{\min}(U_t W_t) \geq \frac{1}{\sqrt{10}}\sigma_{\min}(X)$. Moreover, assume that $\mu \leq c\kappa^{-2}\|X\|^{-2}$, $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-2}$, and

$$\max\{\|(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\|, \|(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\|\} \leq c\kappa^{-2}\|XX^T - U_t U_t^T\|, \quad (47)$$

where the constant $c > 0$ is chosen small enough. Then it holds that

$$\|V_X^T (XX^T - U_{t+1} U_{t+1}^T)\| \leq \left(1 - \frac{\mu}{200}\sigma_{\min}^2(X)\right) \|V_X^T (XX^T - U_t U_t^T)\| + \mu \frac{\sigma_{\min}^2(X)}{100} \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|.$$

When applying this lemma in our proof, we are going to set $\|\cdot\| = \|\cdot\|_F$. However, we believe that this lemma might be of independent interest, as it shows that $U_t U_t^T$ converges linearly towards XX^T with respect to several different norms.

Having collected all the necessary ingredients, we can state and prove the main theorem of this section.

Theorem 9.6. Let $\{U_t\} \subset \mathbb{R}^{d \times r}$ be the sequence created by the gradient descent algorithm. Assume that $\mu \leq c_1 \kappa^{-4} \|X\|^{-2}$ for a sufficiently small constant c_1 . Moreover, assume that \mathcal{A} satisfies the restricted isometry property for rank- $(2r_* + 1)$ matrices with constant $\delta \leq c_1 \kappa^{-4} / \sqrt{r_*}$. Let $\gamma > 0$ and choose the iteration count t_* such that $\sigma_{\min}(U_{t_*} W_{t_*}) \geq \gamma$. Furthermore, assume that the following conditions hold:

$$\|U_{t_*} W_{t_*, \perp}\| \leq 2\gamma, \quad (48)$$

$$\|U_{t_*}\| \leq 3\|X\|, \quad (49)$$

$$\gamma \leq c_2 \frac{\sigma_{\min}(X)}{\min\{r; n\} \kappa^2}, \quad (50)$$

$$\|V_{X^\perp}^T V_{U_{t_*} W_{t_*}}\| \leq c_2 \kappa^{-2}. \quad (51)$$

Then after

$$\hat{t} - t_* \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\max \left\{ 1; \frac{\kappa r_*}{\min\{r; n\} - r_*} \right\} \frac{\|X\|}{\gamma} \right) \quad (52)$$

iterations it holds that

$$\frac{\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F}{\|X\|^2} \lesssim \frac{r_*^{1/8} (\min\{r; n\} - r_*)^{3/8}}{\kappa^{3/16}} \cdot \frac{\gamma^{21/16}}{\|X\|^{21/16}}.$$

Remark 9.7. The proof of Theorem 9.6 shows that the number of iterations needed to complete Phase II is smaller than

$$t_1 - t_* \lesssim \frac{1}{\mu \sigma_{\min}^2(X)} \ln \left(\frac{\sigma_{\min}(X)}{\gamma} \right)$$

and that the number of iterations needed to complete Phase III is smaller than

$$\hat{t} - t_1 \lesssim \frac{1}{\mu \sigma_{\min}^2(X)} \ln \left(\max \left\{ 1; \frac{\kappa r_*}{\min\{r; n\} - r_*} \right\} \frac{\|X\|}{\gamma} \right).$$

Proof of Theorem 9.6. Phase II: In this phase, we will prove that $\sigma_{\min}(V_X^T U_t)$ is growing exponentially until it is at larger than $\frac{\sigma_{\min}(X)}{\sqrt{10}}$, while $\|U_t W_{t, \perp}\|$ stays grows much slower. For that, set

$$t_1 := \min \left\{ t \geq t_* : \sigma_{\min}(V_X^T U_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}} \right\}.$$

We will prove by induction that for $t_* \leq t \leq t_1$ the following inequalities hold:

$$\sigma_{\min}(V_X^T U_t) \geq \frac{1}{2} \left(1 + \frac{1}{8} \mu \sigma_{\min}^2(X) \right)^{t-t_*} \gamma, \quad (53)$$

$$\|U_t W_{t, \perp}\| \leq 2 \left(1 + 80 \mu c_2 \sigma_{\min}^2(X) \right)^{t-t_*} \gamma, \quad (54)$$

$$\|U_t\| \leq 3\|X\|, \quad (55)$$

$$\|V_{X^\perp}^T V_{U_t W_t}\| \leq c_2 \kappa^{-2}. \quad (56)$$

Note that when the inequalities above hold, then from the definition of t_1 above and inequality (53) we can derive that

$$t_1 - t_* \leq \frac{16}{\mu \sigma_{\min}^2(X)} \ln \left(\sqrt{\frac{5}{2}} \cdot \frac{\sigma_{\min}(X)}{\gamma} \right). \quad (57)$$

For $t = t_*$, we first note that inequalities (54), (55), and (56) follow directly from our assumptions. In order to prove inequality (53) we note that

$$\sigma_{\min}(V_X^T U_{t_*}) \geq \sigma_{\min}(V_X^T V_{U_{t_*} W_{t_*}}) \sigma_{\min}(U_{t_*} W_{t_*}) \stackrel{(a)}{\geq} \frac{1}{2} \sigma_{\min}(U_{t_*} W_{t_*}) \stackrel{(b)}{\geq} \frac{\gamma}{2},$$

where inequality (a) is a consequence of assumption (51) and inequality (b) follows from the definition of γ . Assume now that we have shown these four inequalities for t . In order to prove them for $t + 1$ we note first that

$$\begin{aligned} & \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)\| \\ & \stackrel{(a)}{\leq} \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t W_t W_t^T U_t^T)\| + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(U_t W_{t,\perp} W_{t,\perp}^T U_t^T)\| \\ & \stackrel{(b)}{\leq} \delta \sqrt{r_*} \|X X^T - U_t W_t W_t^T U_t^T\| + \delta \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|_* \\ & \leq \delta \sqrt{r_*} (\|X\|^2 + \|U_t W_t\|^2) + \delta \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|_* \\ & \stackrel{(c)}{\leq} 10 \delta \sqrt{r_*} \|X\|^2 + \delta (\min\{r; n\} - r_*) \|U_t W_{t,\perp}\|^2 \\ & \stackrel{(d)}{\leq} 10 c_1 \kappa^{-2} \sigma_{\min}(X)^2 + 4 \delta (\min\{r; n\} - r_*) (1 + 80 \mu c_2 \sigma_{\min}(X)^2)^{2(t-t_*)} \gamma^2 \\ & \stackrel{(e)}{\leq} 10 c_1 \kappa^{-2} \sigma_{\min}(X)^2 + 8 \delta (\min\{r; n\} - r_*) \sigma_{\min}(X)^{1/4} \gamma^{7/4} \\ & \stackrel{(f)}{\leq} 40 c_1 \kappa^{-2} \sigma_{\min}(X)^2. \end{aligned} \tag{58}$$

In inequality (a) we applied the triangle inequality and for inequality (b) we used the restricted isometry property as well as Lemma 7.3. Inequality (c) is due to the induction assumption (55). In inequality (d) we used the assumption $\delta \leq c_1 \kappa^{-4}$ as well as the induction assumption (54). For inequality (e) we used $t \leq t_1$ as well as (57) and for inequality (f) we used (50).

Next, we observe that by Lemma 9.1 we have that

$$\begin{aligned} \sigma_{\min}(V_X^T U_{t+1} W_{t+1}) &= \sigma_{\min}(V_X^T U_{t+1}) \\ &\geq \sigma_{\min}(V_X^T U_{t+1} W_t) \\ &\geq \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{4} \mu \sigma_{\min}^2(X) - \mu \sigma_{\min}^2(V_X^T U)\right) \\ &\stackrel{(a)}{\geq} \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{8} \mu \sigma_{\min}^2(X)\right). \end{aligned}$$

In (a) we have used that $\sigma_{\min}(V_X^T U_t) \leq \frac{\sigma_{\min}(X)}{\sqrt{10}}$, which follows from $t \leq t_1$. Using the induction assumption, this implies inequality (53). Moreover, the inequality chain above shows that $V_X^T U_{t+1} W_{t+1}$ has full rank. This allows us to apply Lemma 9.2, which implies that

$$\begin{aligned} \|U_{t+1} W_{t+1, \perp}\| &\leq \left(1 - \frac{\mu}{2} \|U_t W_{t, \perp}\|^2 + 36 \mu \|V_X^T V_{U_t W_t}\|^2 \|X\|^2 + 2 \mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)\|\right) \|U_t W_{t, \perp}\| \\ &\stackrel{(a)}{\leq} (1 + 80 \mu c_2 \sigma_{\min}^2(X)) \|U_t W_{t, \perp}\| \\ &\leq 2 (1 + 80 \mu c_2 \sigma_{\min}^2(X))^{t+1-t_*} \gamma, \end{aligned}$$

where in inequality (a) we used (56) as well as (58) and that the constant c_1 is chosen sufficiently small. This shows inequality (54). Next, due to inequality (58), our induction assumptions, and Lemma 9.4 we obtain that $\|U_{t+1}\| \leq 3\|X\|$, which shows inequality (55).

Next, we note that by Lemma 9.3 we have that

$$\begin{aligned}
& \|V_{X^\perp}^T V_{U_{t+1}W_{t+1}}\| \\
& \leq \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) \|V_{X^\perp}^T V_{U_tW_t}\| + 100\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_tU_t^T)\| + 500\mu^2 \|XX^T - U_tU_t^T\|^2 \\
& \stackrel{(a)}{\leq} \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) \|V_{X^\perp}^T V_{U_tW_t}\| + 2000c_1\mu\kappa^{-2}\sigma_{\min}(X)^2 + 50000\mu^2\|X\|^4 \\
& \stackrel{(b)}{\leq} \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) \|V_{X^\perp}^T V_{U_tW_t}\| + 2000c_1\mu\kappa^{-2}\sigma_{\min}(X)^2 + 50000c_1\mu\kappa^{-2}\sigma_{\min}^2(X) \\
& \stackrel{(c)}{\leq} \left(1 - \frac{\mu}{4}\sigma_{\min}^2(X)\right) c_2\kappa^{-2} + 2000c_1\mu\kappa^{-2}\sigma_{\min}(X)^2 + 50000c_1\mu\kappa^{-2}\sigma_{\min}^2(X),
\end{aligned}$$

where in inequality (a) we used the induction hypothesis (55) as well as (58). Inequality (b) follows from inequality (58) and our assumption on the step size μ . In inequality (c) we used the induction assumption $\|V_{X^\perp}^T V_{U_tW_t}\| \leq c_2\kappa^{-2}$. By choosing the constant $c_1 > 0$ small enough, this implies inequality (56) and, hence, finishes the induction step.

Note that from the definition of t_1 and from inequality (53) the inequality (57) follows. Hence, we obtain that

$$\begin{aligned}
\|U_{t_1}W_{t_1,\perp}\| & \stackrel{(a)}{\leq} 2(1 + 80\mu c_2\sigma_{\min}^2(X))^{t_1-t_\star} \gamma \\
& \stackrel{(b)}{\leq} 2(\sigma_{\min}(X)/\gamma)^{80c_2} \gamma \\
& \stackrel{(c)}{\leq} 2\left(\frac{\sigma_{\min}(X)}{\gamma}\right)^{1/8} \gamma \\
& = 2\sigma_{\min}(X)^{1/8} \gamma^{7/8},
\end{aligned} \tag{59}$$

where inequality (a) follows from inequality (54) and inequality (b) follows from (57). Inequality (c) follows from choosing $c_2 > 0$ small enough. This finishes the proof of the second phase.

Phase III: In the third phase, we analyse the refinement of the signal U_t . For that, we set

$$\hat{t} := t_1 + \left\lfloor \frac{300}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{5}{8} \kappa^{1/4} \sqrt{\frac{r_\star}{\min\{r; n\} - r_\star}} \frac{\|X\|^{7/4}}{\gamma^{7/4}} \right) \right\rfloor. \tag{60}$$

Similar as in Phase II, we are going to show inductively that the following inequalities are fulfilled for $t_1 \leq t \leq \hat{t}$

$$\sigma_{\min}(U_tW_t) \geq \sigma_{\min}(V_X^T U_t) \geq \frac{\sigma_{\min}(X)}{\sqrt{10}}, \tag{61}$$

$$\|U_tW_{t,\perp}\| \leq (1 + 80\mu c_2\sigma_{\min}^2(X))^{t-t_1} \|U_{t_1}W_{t_1,\perp}\|, \tag{62}$$

$$\|U_t\| \leq 3\|X\|, \tag{63}$$

$$\|V_{X^\perp}^T V_{U_tW_t}\| \leq c_2\kappa^{-2}, \tag{64}$$

$$\|V_X^T (XX^T - U_tU_t^T)\|_F \leq 10\sqrt{r_\star} \left(1 - \frac{\mu}{400}\sigma_{\min}^2(X)\right)^{t-t_1} \|X\|^2. \tag{65}$$

For $t = t_1$ we note the inequalities (61), (63), and (64) follow from the results in Phase 1. The inequality (62) follows directly from setting $t = t_1$. For $t = t_1$, inequality (65) follows from the observation that

$$\begin{aligned}\|V_X^T (XX^T - U_{t_1}U_{t_1}^T)\|_F &= \|V_X^T (XX^T - U_{t_1}W_{t_1}W_{t_1}^T U_{t_1}^T)\|_F \\ &\leq \|XX^T\|_F + \|U_{t_1}W_{t_1}W_{t_1}^T U_{t_1}^T\|_F \\ &\leq \sqrt{r_*} (\|XX^T\| + \|U_{t_1}W_{t_1}W_{t_1}^T U_{t_1}^T\|) \\ &\leq 10\sqrt{r_*}\|X\|^2,\end{aligned}$$

where we have used that $\|U_{t_1}W_{t_1}\| \leq \|U_{t_1}\| \leq 3\|X\|$ by induction assumption (63).

For the induction step from t to $t+1$, we note first that with similar arguments as in Phase 1 we can show that

$$\begin{aligned}&\|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \\ &\leq 10\delta\sqrt{r_*}\|X\|^2 + \delta(\min\{r; n\} - r_*)\|U_t W_{t,\perp}\|^2 \\ &\stackrel{(a)}{\leq} 10c_1\kappa^{-2}\sigma_{\min}(X)^2 + 4\delta(\min\{r; n\} - r_*)(1 + 80\mu c_2\sigma_{\min}(X)^2)^{2(t-t_1)}\sigma_{\min}(X)^{1/4}\gamma^{7/4} \\ &\stackrel{(b)}{\leq} 10c_1\kappa^{-2}\sigma_{\min}(X)^2 + 4\delta(\min\{r; n\} - r_*)\left(\frac{5}{8}\kappa^{1/4}\sqrt{\frac{r_*}{\min\{r; n\} - r_*}}\frac{\|X\|^{7/4}}{\gamma^{7/4}}\right)^{O(c_2)}\sigma_{\min}(X)^{1/4}\gamma^{7/4} \\ &\stackrel{(c)}{\leq} 40c_1\kappa^{-2}\sigma_{\min}(X)^2,\end{aligned}$$

where in inequality (a) follows from (62). Inequality (b) follows from (60) as well as the elementary inequality $\ln(1+x) \leq x$. Inequality (c) follows from the assumption $\gamma \leq c_2 \frac{\sigma_{\min}(X)}{\min\{r; n\}\kappa}$. This puts us in a position to apply our technical lemmas. We note that by Lemma 9.1 we have that

$$\begin{aligned}\sigma_{\min}(U_{t+1}W_{t+1}) &\geq \sigma_{\min}(V_X^T U_{t+1}) \geq \sigma_{\min}(V_X^T U_{t+1}W_t) \\ &\geq \underbrace{\sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{4}\mu\sigma_{\min}(X)^2 - \mu\sigma_{\min}(V_X^T U_t)^2\right)}_{= (*)}.\end{aligned}$$

Note that for $\sigma_{\min}(V_X^T U_t) \leq \frac{1}{2}\sigma_{\min}(X)$ it holds that $(*) \geq 1$ and thus it follows that (61) holds for $t+1$ in this case. In the case of $\frac{1}{2}\sigma_{\min}(X) \leq \sigma_{\min}(V_X^T U_t)$ we obtain that

$$(*) \geq 1 - \mu\sigma_{\min}(V_X^T U_t)^2 \stackrel{(a)}{\geq} 1 - 9\mu\|X\|^2 \stackrel{(b)}{\geq} 4/5,$$

where in inequality (a) we used the induction hypothesis (55) and in inequality (b) we used the assumption $\mu \leq c_1\kappa^{-2}\|X\|^{-2}$. Hence, we have shown that also in this case the inequality (61) holds for $t+1$.

Note that the previous inequality chain also implies that $V_X^T U_{t+1}W_t$ is invertible. Hence, in a similar way as in Phase II for inequality (54) we can verify that (62) holds for $t+1$.

Note that from Lemma 9.4, induction assumption (55), and the assumption on the step size μ it follows that $\|U_{t+1}\| \leq 3\|X\|$. Moreover, inequality (64) can be shown analogously as in Phase 1. Next, we note that due to the restricted isometry property we have that

$$\|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\|_F \leq c_1\kappa^{-2}\|XX^T - U_t U_t^T\|_F,$$

which shows that inequality (47) is fulfilled (with $\|\cdot\|$ being the Frobenius norm $\|\cdot\|_F$). Hence, we obtain from Lemma 9.5 that

$$\begin{aligned} & \|V_X^T (XX^T - U_{t+1}U_{t+1}^T)\|_F \\ & \leq \left(1 - \frac{\mu}{200}\sigma_{\min}(X)^2\right) \|V_X^T (XX^T - U_tU_t^T)\|_F + \mu \frac{\sigma_{\min}(X)^2}{100} \|U_tW_{t,\perp}W_{t,\perp}^T U_t^T\|_F \\ & \leq 10\sqrt{r_\star} \left(1 - \frac{\mu}{200}\sigma_{\min}(X)^2\right) \left(1 - \frac{\mu}{400}\sigma_{\min}(X)^2\right)^{t-t_1} \|X\|^2 + \mu \frac{\sigma_{\min}(X)^2}{100} \|U_tW_{t,\perp}W_{t,\perp}^T U_t^T\|_F, \end{aligned}$$

where in the last inequality we used the induction assumption (65). We note that this shows (65) holds for $t+1$, if we can show that

$$\|U_tW_{t,\perp}W_{t,\perp}^T U_t^T\|_F \leq \frac{5}{2}\sqrt{r_\star} \left(1 - \frac{\mu}{400}\sigma_{\min}(X)^2\right)^{t-t_1} \|X\|^2. \quad (66)$$

For that, we note that

$$\begin{aligned} \|U_tW_{t,\perp}W_{t,\perp}^T U_t^T\|_F & \leq \sqrt{\min\{r; n\} - r_\star} \|U_tW_{t,\perp}\|^2 \\ & \leq 4\sqrt{\min\{r; n\} - r_\star} (1 + 80\mu c_2 \sigma_{\min}(X)^2)^{2(t-t_1)} \sigma_{\min}(X)^{1/4} \gamma^{7/4}, \end{aligned}$$

where in the last inequality we used (59) and (62). Hence, for $c_2 > 0$ small enough, inequality (66) is implied by

$$\frac{8}{5} \sqrt{\frac{\min\{r; n\} - r_\star}{r_\star}} \sigma_{\min}(X)^{1/4} \gamma^{7/4} \leq \left(1 - \frac{\mu}{350}\sigma_{\min}(X)^2\right)^{t-t_1} \|X\|^2.$$

By rearranging terms and using the elementary inequality $\ln(1+x) \geq \frac{x}{1+x}$, we see that this in turn is implied by

$$t - t_1 \leq \frac{300}{\mu \sigma_{\min}(X)^2} \ln \left(\frac{5}{8} \sqrt{\frac{r_\star}{\min\{r; n\} - r_\star}} \cdot \frac{\|X\|^2}{\gamma^{7/4} \sigma_{\min}(X)^{1/4}} \right).$$

Hence, (60) shows (66), which shows inequality (65) for $t+1$. This finishes the induction step.

Conclusion: In order to finish the proof we note that

$$\begin{aligned} \|U_{\hat{t}}U_{\hat{t}}^T - XX^T\|_F & \stackrel{(a)}{\leq} 4\|V_X^T (XX^T - U_{\hat{t}}U_{\hat{t}}^T)\|_F + \|U_{\hat{t}}W_{\hat{t},\perp}W_{\hat{t},\perp}^T U_{\hat{t}}^T\|_F \\ & \stackrel{(b)}{\lesssim} \sqrt{r_\star} \left(1 - \frac{\mu}{400}\sigma_{\min}(X)^2\right)^{\hat{t}-t_1} \|X\|^2 \\ & \stackrel{(c)}{\lesssim} \sqrt{r_\star} \left(\frac{5}{8} \kappa^{1/4} \sqrt{\frac{r_\star}{\min\{r; n\} - r_\star}} \frac{\|X\|^{7/4}}{\gamma^{7/4}} \right)^{-3/4} \|X\|^2 \\ & \lesssim r_\star^{1/8} \kappa^{-3/16} (\min\{r; n\} - r_\star)^{3/8} \gamma^{21/16} \|X\|^{11/16}, \end{aligned}$$

where inequality (a) follows from the triangle inequality and the definition of $W_{\hat{t}}$. Inequality (b) follows from (65) and inequality (66). In (c) we used the definition of \hat{t} .

In order to finish the proof we need to show (52). For that we note that

$$\begin{aligned}
t - t_1 &\leq \frac{300}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{5}{8} \sqrt{\frac{r_\star}{\min\{r; n\} - r_\star}} \cdot \frac{\|X\|^2}{\gamma^{7/4} \sigma_{\min}(X)^{1/4}} \right) \\
&= \frac{300}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{5}{8} \kappa^{1/4} \sqrt{\frac{r_\star}{\min\{r; n\} - r_\star}} \frac{\|X\|^{7/4}}{\gamma^{7/4}} \right) \\
&\leq \frac{300}{\mu\sigma_{\min}(X)^2} \ln \left(\min \left\{ 1; \frac{\kappa r_\star}{\min\{r; n\} - r_\star} \right\} \frac{\|X\|^{7/4}}{\gamma^{7/4}} \right) \\
&\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\min \left\{ 1; \frac{\kappa r_\star}{\min\{r; n\} - r_\star} \right\} \frac{\|X\|}{\gamma} \right).
\end{aligned}$$

Combining this with inequality (57) shows (52). \square

10 Proof of the main results

10.1 Proof of Theorem 3.3

Proof of Theorem 3.3. Set $\tilde{E} = \mathcal{A}^* \mathcal{A}(XX^T) - XX^T$. From the spectral-to-spectral restricted isometry property, which follows from Lemma 7.3 as well as from our assumption on the restricted isometry property, it follows that

$$\|\tilde{E}\| = \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T)\| \leq c\kappa^{-4} \|X\|^2 = c\kappa^{-2} \sigma_{\min}(X)^2. \quad (67)$$

In order to finish the proof we will distinguish two cases:

Case $r \geq 2r_\star$: Due to (4) and (67) we can apply Lemma 8.7. Hence, with probability at least $1 - O(\exp(-cr))$ after

$$t_\star \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \cdot \ln \left(2\kappa^2 \sqrt{\frac{n}{\min\{r; n\}}} \right)$$

iterations we have that

$$\|U_{t_\star}\| \leq 3\|X\|, \quad (68)$$

$$\sigma_{\min}(U_{t_\star} W_{t_\star}) \geq \frac{\alpha\beta}{4}, \quad (69)$$

$$\|U_{t_\star} W_{t_\star, \perp}\| \leq \frac{\kappa^{-2}}{8} \alpha\beta, \quad (70)$$

$$\|V_{X^\perp}^T V_{U_{t_\star} W_{t_\star}}\| \leq c\kappa^{-2} \quad (71)$$

with $1 \lesssim \beta \lesssim \frac{n\kappa^4}{\min\{r; n\}}$. Our goal is to apply Theorem 9.6 with $\gamma = \frac{\alpha\beta}{4}$. For that we need to check that

$$\frac{c_2 \sigma_{\min}(X)}{\min\{r; n\} \kappa^2} \geq \gamma = \frac{\alpha\beta}{4} \quad (72)$$

holds. Note that since

$$\frac{4c_2\|X\|}{\min\{r;n\}\kappa^3\beta} \gtrsim \frac{\|X\|}{\kappa^7 n} \gtrsim \alpha$$

condition (72) is fulfilled, when the constant in (4) is chosen sufficiently small. Hence, by Theorem 9.6 after

$$\hat{t} - t_\star \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\max \left\{ 1; \frac{\kappa r_\star}{\min\{r;n\} - r_\star} \right\} \frac{4\|X\|}{\alpha\beta} \right)$$

iterations it holds that

$$\begin{aligned} \frac{\|U_{\hat{t}}U_{\hat{t}}^T - XX^T\|_F}{\|X\|^2} &\lesssim \frac{r_\star^{1/8} (\min\{r;n\} - r_\star)^{3/8}}{\kappa^{3/16}} \cdot \frac{\gamma^{21/16}}{\|X\|^{21/16}} \\ &\lesssim \frac{r_\star^{1/8} (\min\{r;n\} - r_\star)^{3/8} (\alpha\beta)^{21/16}}{\kappa^{3/16} \|X\|^{21/16}} \\ &\lesssim \frac{\kappa^{81/16} n^{21/16} r_\star^{1/8} (\min\{r;n\} - r_\star)^{3/8}}{(\min\{r;n\})^{21/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}} \\ &\leq \frac{n^{21/16} \kappa^{81/16} r_\star^{1/8}}{(\min\{r;n\})^{15/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}}. \end{aligned}$$

Note that for the total amount of iterations we have that

$$\begin{aligned} \hat{t} &\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \left(\ln \left(2\kappa^2 \sqrt{\frac{n}{\min\{r;n\}}} \right) + \ln \left(\max \left\{ 1; \frac{\kappa r_\star}{\min\{r;n\} - r_\star} \right\} \frac{4\|X\|}{\alpha\beta} \right) \right) \\ &= \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(8\kappa^3 \sqrt{\frac{n}{\min\{r;n\}}} \cdot \max \left\{ 1; \frac{\kappa r_\star}{\min\{r;n\} - r_\star} \right\} \cdot \frac{\|X\|}{\alpha\beta} \right) \\ &\stackrel{(b)}{\leq} \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(C_1 \kappa^3 \sqrt{\frac{n}{\min\{r;n\}}} \cdot \max \left\{ 1; \frac{\kappa r_\star}{\min\{r;n\} - r_\star} \right\} \cdot \frac{\|X\|}{\alpha} \right) \\ &\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{C_1 \kappa n}{\min\{r;n\}} \cdot \max \left\{ 1; \frac{\kappa r_\star}{\min\{r;n\} - r_\star} \right\} \cdot \frac{\|X\|}{\alpha} \right), \end{aligned}$$

where in inequality (b) we have used $\beta \gtrsim 1$ and chosen the constant $C_1 > 0$ large enough. This finishes the proof of the first part.

Case $r_\star < r < 2r_\star$: As in the first case, we can apply Lemma 8.7. Hence, with probability at least $1 - (C\varepsilon)^{r-r_\star+1} + O(\exp(-cr))$ after

$$t_\star \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \cdot \ln \left(\frac{2\kappa^2 \sqrt{rn}}{\varepsilon} \right)$$

iterations the inequalities (33), (34), (35), and (36) hold with $\frac{\varepsilon}{r} \lesssim \beta \lesssim \frac{\kappa^4 n}{\varepsilon}$. Again, we want to apply Theorem 9.6 with $\gamma = \frac{\alpha\beta}{4}$. For that we need to check that

$$\frac{c_2\sigma_{\min}(X)}{\min\{r;n\}\kappa^2} \geq \gamma = \frac{\alpha\beta}{4} \quad (73)$$

holds. Note that since

$$\frac{4c_2\|X\|}{\min\{r;n\}\kappa^3\beta} \gtrsim \frac{\varepsilon r\|X\|}{\min\{r;n\}n\kappa^7} = \frac{\varepsilon\|X\|}{n\kappa^7} \gtrsim \alpha$$

condition (73) holds true, when the constant in (6) is chosen sufficiently small. Hence, by Theorem 9.6 after

$$\hat{t} - t_* \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\max \left\{ 1; \frac{\kappa r_*}{r - r_*} \right\} \frac{4\|X\|}{\alpha\beta} \right)$$

iterations it holds that

$$\begin{aligned} \frac{\|U_{\hat{t}}U_{\hat{t}}^T - XX^T\|_F}{\|X\|^2} &\lesssim \frac{r_*^{1/8}(r-r_*)^{3/8}}{\kappa^{3/16}} \cdot \frac{\gamma^{21/16}}{\|X\|^{21/16}} \\ &\lesssim \frac{r_*^{1/8}(r-r_*)^{3/8}}{\kappa^{3/16}} \cdot \frac{(\alpha\beta)^{21/16}}{\|X\|^{21/16}} \\ &\leq r_*^{1/8}(r-r_*)^{3/8} \kappa^{81/16} \left(\frac{n}{\varepsilon} \cdot \frac{\alpha}{\|X\|} \right)^{21/16}. \end{aligned}$$

Note that for the total amount of iterations we have that

$$\begin{aligned} \hat{t} &\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \left(\ln \left(\frac{2\kappa^2\sqrt{rn}}{\varepsilon} \right) + \ln \left(\max \left\{ 1; \frac{\kappa r_*}{r - r_*} \right\} \frac{4\|X\|}{\alpha\beta} \right) \right) \\ &\stackrel{(a)}{=} \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{8\kappa^3 r_* \sqrt{rn}}{\varepsilon(r-r_*)} \cdot \frac{\|X\|}{\alpha\beta} \right) \\ &\stackrel{(b)}{\leq} \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{C_2\kappa^3 r_* r \sqrt{rn}}{\varepsilon^2(r-r_*)} \cdot \frac{\|X\|}{\alpha} \right) \\ &\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{C_2\kappa n^2}{\varepsilon^2(r-r_*)} \cdot \frac{\|X\|}{\alpha} \right), \end{aligned}$$

where in equality (a) we have used $\frac{r_*}{r-r_*} \geq 1$, which follows from $r_* < r \leq 2r_*$. Inequality (b) follows from $\frac{\varepsilon}{r} \lesssim \beta$ as well as from choosing $C_2 > 0$ large enough. This finishes the proof. \square

10.2 Proof of Theorem 3.4

Proof of Theorem 3.4. As in the proof of the second part of Theorem 3.3 we can show that with probability at least $1 - C\varepsilon + O(\exp(-cr_*))$ after

$$t_* \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \cdot \ln \left(\frac{2\kappa^2\sqrt{n}}{\varepsilon} \right)$$

iterations the inequalities (33), (34), (35), and (36) hold with $\frac{\varepsilon}{r_*} \lesssim \beta \lesssim \frac{\kappa^4 n}{\varepsilon}$. Now define the matrix \widehat{U}_{t_*} by adding a zero column to U_{t_*} , i.e.,

$$\widehat{U}_{t_*} = (U_{t_*} \quad 0) \in \mathbb{R}^{n \times (r_*+1)}.$$

Clearly, we can run gradient descent on \widehat{U}_{t_\star} instead of U_{t_\star} with the same step size, which gives us a sequence $\widehat{U}_{t_\star}, \widehat{U}_{t_\star+1}, \widehat{U}_{t_\star+2}, \dots$ to which we can apply Theorem 9.6 with $\gamma = \frac{\alpha\beta}{4}$ and $r = r_\star + 1$. However, note that the last column always stays zero, which means that the results of this theorem also apply to $U_{t_\star}, U_{t_\star+1}, U_{t_\star+2}, \dots$. Hence, after

$$\hat{t} - t_\star \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(4\kappa r_\star \frac{\|X\|}{\alpha\beta} \right) \lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(4\kappa r_\star^2 \frac{\|X\|}{\alpha\varepsilon} \right)$$

iterations it holds that

$$\begin{aligned} \frac{\|U_{\hat{t}}U_{\hat{t}}^T - XX^T\|_F}{\|X\|^2} &\lesssim \frac{r_\star^{1/8}}{\kappa^{3/16}} \cdot \frac{(\alpha\beta)^{21/16}}{\|X\|^{21/16}} \\ &\lesssim \frac{r_\star^{1/8}}{\kappa^{3/16}} \left(\frac{\kappa^4 n}{\varepsilon} \cdot \frac{\alpha}{\|X\|} \right)^{21/16} \\ &= r_\star^{1/8} \kappa^{81/16} \left(\frac{n}{\varepsilon} \cdot \frac{\alpha}{\|X\|} \right)^{21/16}. \end{aligned}$$

Note that for the total amount of iterations we have that

$$\begin{aligned} \hat{t} &\lesssim \frac{1}{\mu\sigma_{\min}(X)^2} \left(\ln \left(\frac{2\kappa^2 \sqrt{n}}{\varepsilon} \right) + \ln \left(4\kappa r_\star^2 \frac{\|X\|}{\alpha\varepsilon} \right) \right) \\ &= \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{8\kappa^3 r_\star^2 \sqrt{n} \|X\|}{\alpha\varepsilon^2} \right) \\ &\leq \frac{1}{\mu\sigma_{\min}(X)^2} \ln \left(\frac{8\kappa^3 n^3}{\varepsilon^2} \cdot \frac{\|X\|}{\alpha} \right). \end{aligned}$$

This finishes the proof. □

10.3 Proof of Theorem 3.5

We start by noting that in the special case $r = n$, the required assumptions for Theorem 9.6 are already fulfilled at the initialization $t_0 = 0$. This means that in this special case we do not need to analyze the spectral phase. This is shown by the following lemma.

Lemma 10.1. *Assume that $r = n$ and let $U_0 = \alpha U$, where $U \in \mathbb{R}^{n \times n}$ is an orthonormal matrix. Then it holds that*

$$\begin{aligned} \|V_{X^\perp}^T V_{U_0 W_0}\| &= 0, \\ \sigma_{\min}(U_0 W_0) &= \alpha, \\ \|U_0\| &= \alpha. \end{aligned}$$

Proof. Note that $V_X^T U \in \mathbb{R}^{r_\star \times n}$ is an isometric embedding. Hence, a feasible choice for W_0 is given by $W_0 = U^T V_X$, which implies that

$$U_0 W_0 = \alpha U U^T V_X = \alpha V_X.$$

It follows that $\|V_{X^\perp}^T V_{U_0 W_0}\| = 0$, which verifies the first equality. In order to see that the second equality holds we note that

$$\sigma_{\min}(U_0 W_0) = \sigma_{\min}(\alpha V_X) = \alpha.$$

The third equality follows directly from the definition of U_0 . This finishes the proof. \square

Now we are in a position to give a proof of Theorem 3.5.

Proof of Theorem 3.5. By Lemma 10.1 we have that $\|V_{X^\perp}^T V_{U_0 W_0}\| = 0$, $\sigma_{\min}(U_0 W_0) = \alpha$, and $\|U_0\| = \alpha$. This allows us to apply Theorem 9.6 with $t_0 = 0$ and $\gamma = \alpha$, which yields that after

$$\hat{t} \lesssim \frac{1}{\mu \sigma_{\min}(X)^2} \ln \left(\max \left\{ 1; \frac{\kappa r_\star}{n - r_\star} \right\} \frac{\|X\|}{\alpha} \right)$$

iterations we have that

$$\begin{aligned} \frac{\|U_{\hat{t}} U_{\hat{t}}^T - X X^T\|_F}{\|X\|^2} &\lesssim \frac{r_\star^{1/8} (n - r_\star)^{3/8}}{\kappa^{3/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}} \\ &\leq \frac{r_\star^{1/8} n^{3/8}}{\kappa^{3/16}} \cdot \frac{\alpha^{21/16}}{\|X\|^{21/16}}. \end{aligned}$$

This finishes the proof. \square

11 Conclusion

In this paper we focused on demystifying the role of initialization when training overparameterized models by showing that small random initialization followed by a few iterations of gradient descent behaves akin to popular spectral methods. We also show that this *implicit spectral bias* from small random initialization, which is provably more prominent for overparameterized models, also puts the gradient descent iterations on a particular trajectory towards solutions that are not only globally optimal but also generalize well.

We think that our results give rise to a number of interesting future research directions. For example, one could extend our results to scenarios where the measurement matrices are more structured such as in matrix completion [59] or in blind deconvolution [60]. Moreover, while our main results, e.g. Theorem 3.3 do require early stopping, our simulations (e.g. Figure 7a) indicate that early stopping is not needed. It would be interesting to examine whether we can remove the early stopping requirement. It is also an interesting future avenue to examine whether the quadratic dependence of the sample complexity m on r_\star in our results is really needed.

Moreover, while in this paper our main focus was on low-rank matrix reconstruction, we believe that our analysis holds more generally for a variety of contemporary overparameterized machine learning and signal estimation tasks including neural network training. This is a tantalizing future research direction.

Acknowledgements

M.S. is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, an NSF-CAREER under award #1846369, the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award #FA9550-18-1-0078, DARPA Learning with Less Labels (LwLL) and FastNICS programs, and NSF-CIF awards #1813877 and #2008443.

References

- [1] Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow: theory and algorithms. *IEEE Trans. Inf. Theory*, 61(4):1985–2007, 2015.
- [2] Yuxin Chen and Emmanuel J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Commun. Pure Appl. Math.*, 70(5):822–883, 2017.
- [3] Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Found. Comput. Math.*, 20(3):451–632, 2020.
- [4] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning (ICML)*, volume 48, pages 964–973. Journal of Machine Learning Research, 2016.
- [5] Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Appl. Comput. Harmon. Anal.*, 47(3):893–934, 2019.
- [6] Shuyang Ling and Thomas Strohmer. Regularized gradient descent: a non-convex recipe for fast joint blind deconvolution and demixing. *Inf. Inference*, 8(1):1–49, 2019.
- [7] Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. Phase retrieval using alternating minimization. *IEEE Trans. Signal Process.*, 63(18):4814–4826, 2015.
- [8] Irène Waldspurger. Phase retrieval with random Gaussian sensing vectors by alternating projections. *IEEE Trans. Inf. Theory*, 64(5):3301–3312, 2018.
- [9] Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.
- [10] Yurii Nesterov and B. T. Polyak. Cubic regularization of Newton method and its global performance. *Math. Program.*, 108(1 (A)):177–205, 2006.
- [11] Jorge Nocedal and Stephen J. Wright. Trust-region methods. *Numerical Optimization*, pages 66–100, 2006.
- [12] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape saddle points efficiently. page 1724–1732, 2017.
- [13] Rong Ge, Furong Huang, Chi Jin, and Yang. Yuan. Escaping from saddle points: online stochastic gradient for tensor decomposition. In *Proceedings of The 28th Conference on Learning Theory*, pages 797–842, 2015.
- [14] Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. pages 1674–1703, 2017.
- [15] Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient langevin dynamics. In Satyen Kale and Ohad Shamir, editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings of Machine Learning Research*, pages 1980–2022. PMLR, 07–10 Jul 2017.

- [16] Jason D. Lee, Max Simchowitz, Michael I. Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 1246–1257, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR.
- [17] Simon S Du, Chi Jin, Jason D Lee, Michael I Jordan, Aarti Singh, and Barnabas Poczos. Gradient descent can take exponential time to escape saddle points. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [18] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, pages 7413–7424, 2019.
- [19] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [20] Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *International Conference on Machine Learning*, pages 4951–4960. PMLR, 2019.
- [21] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- [22] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.
- [23] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *36th International Conference on Machine Learning, ICML 2019*, pages 477–502. International Machine Learning Society (IMLS), 2019.
- [24] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [25] Lénaïc Chizat and Francis Bach. On the global convergence of gradient descent for overparameterized models using optimal transport. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 3036–3046. Curran Associates, Inc., 2018.
- [26] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- [27] Adel Javanmard, Marco Mondelli, Andrea Montanari, et al. Analysis of a two-layer neural network via displacement convexity. *Ann. Statist.*, 48(6):3619–3642, 2020.

- [28] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: a central limit theorem. *Stochastic Processes Appl.*, 130(3):1820–1852, 2020.
- [29] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32:2937–2947, 2019.
- [30] Emmanuel J. Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Trans. Inf. Theory*, 57(4):2342–2359, 2011.
- [31] Martin Slawski, Ping Li, and Matthias Hein. Regularization-free estimation in trace regression with symmetric positive semidefinite matrices. *Advances in Neural Information Processing Systems*, 28:2782–2790, 2015.
- [32] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [33] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank matrix recovery. 29, 2016.
- [34] Richard Y. Zhang, Somayeh Sojoudi, and Javad Lavaei. Sharp restricted isometry bounds for the inexistence of spurious local minima in nonconvex matrix recovery. *J. Mach. Learn. Res.*, 20(114):1–34, 2019.
- [35] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initialization: fast global convergence for nonconvex phase retrieval. *Math. Program.*, 176(1-2 (B)):5–37, 2019.
- [36] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. volume 75 of *Proceedings of Machine Learning Research*, pages 2–47. PMLR, 06–09 Jul 2018.
- [37] Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- [38] Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient descent without ℓ_2, ∞ regularization. *IEEE Trans. Inf. Theory*, 66(9):5806–5841, 2020.
- [39] Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Found. Comput. Math.*, 18(5):1131–1198, 2018.
- [40] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in Neural Information Processing Systems*, 29:2973–2981, 2016.
- [41] Richard Y. Zhang, Somayeh Sojoudi, and Javad Lavaei. Sharp restricted isometry bounds for the inexistence of spurious local minima in nonconvex matrix recovery. *J. Mach. Learn. Res.*, 20(114):1–34, 2019.
- [42] Suriya Gunasekar, Blake Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nathan Srebro. Implicit regularization in matrix factorization. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6152–6160, 2017.

- [43] Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by norms. *arXiv preprint arXiv:2005.06398*, 2020.
- [44] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. In *International Conference on Learning Representations*, 2021.
- [45] Kelly Geyer, Anastasios Kyrillidis, and Amir Kalev. Low-rank regularization and solution uniqueness in over-parameterized matrix sensing. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 930–940. PMLR, 26–28 Aug 2020.
- [46] Maryia Kabanava, Richard Kueng, Holger Rauhut, and Ulrich Terstiege. Stable low-rank matrix recovery via null space properties. *Inf. Inference*, 5(4):405–441, 2016.
- [47] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? *arXiv preprint arXiv:2006.13409*, 2020.
- [48] Xiang Wang, Chenwei Wu, Jason D Lee, Tengyu Ma, and Rong Ge. Beyond lazy training for over-parameterized tensor decomposition. *Advances in Neural Information Processing Systems*, 33, 2020.
- [49] Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *International Conference on Learning Representations*, 2018.
- [50] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer relu neural networks beyond ntk. *arXiv preprint arXiv:2007.04596*, 2020.
- [51] Peter Bartlett, Dave Helmbold, and Philip Long. Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks. 80:521–530, 10–15 Jul 2018.
- [52] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International Conference on Machine Learning*, pages 244–253. PMLR, 2018.
- [53] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- [54] Bubacarr Bah, Holger Rauhut, Ulrich Terstiege, and Michael Westdickenberg. Learning deep linear neural networks: Riemannian gradient flows and convergence to global minimizers. *arXiv preprint arXiv:1910.05505*, 2019.
- [55] Hung-Hsu Chou, Carsten Gieshoff, Johannes Maly, and Holger Rauhut. Gradient descent for deep matrix factorization: Dynamics and implicit bias towards low rank. *arXiv preprint arXiv:2011.13772*, 2020.
- [56] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. 31, 2018.

- [57] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Spectral methods for data science: A statistical perspective. *arXiv preprint arXiv:2012.08496*, 2020.
- [58] Mark Rudelson and Roman Vershynin. Smallest singular value of a random rectangular matrix. *Commun. Pure Appl. Math.*, 62(12):1707–1739, 2009.
- [59] Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Found. Comput. Math.*, 9(6):717–772, 2009.
- [60] Ali Ahmed, Benjamin Recht, and Justin Romberg. Blind deconvolution using convex programming. *IEEE Trans. Inf. Theory*, 60(3):1711–1732, 2014.
- [61] Simon Foucart and Holger Rauhut. *A mathematical introduction to compressive sensing*. New York, NY: Birkhäuser/Springer, 2013.
- [62] Per-Ake Wedin. Perturbation bounds in connection with singular value decomposition. *BIT, Nord. Tidskr. Inf.-behandl.*, 12:99–111, 1972.
- [63] Chandler Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.*, 7:1–46, 1970.

A Proofs for the spectral phase

A.1 Proof of Lemma 8.1

Proof of Lemma 8.1. We are first going to derive a formula for $\tilde{U}_t - U_t$.

Claim: Set $\hat{E}_i := \mu \mathcal{A}^* \mathcal{A}(U_{i-1} U_{i-1}^T) U_{i-1}$. Then, for $t \geq 1$ it holds that

$$\tilde{U}_t - U_t = \sum_{i=1}^t (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T))^{t-i} \hat{E}_i. \quad (74)$$

Proof of the claim: We will prove the claim by induction. For $t = 1$ we note that

$$\begin{aligned} U_1 &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T - U_0 U_0^T)) U_0 \\ &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T)) U_0 - \mu \mathcal{A}^* \mathcal{A}(U_0 U_0^T) U_0 \\ &= \tilde{U}_1 - \hat{E}_1, \end{aligned}$$

which proves the claim for $t = 1$. Now suppose that the claim holds for some t . We obtain that

$$\begin{aligned} U_{t+1} &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)) U_t \\ &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T)) U_t - \mu \mathcal{A}^* \mathcal{A}(U_t U_t^T) U_t \\ &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T)) U_t - \hat{E}_{t+1}, \end{aligned}$$

where the last line follows from the definition of \hat{E}_{t+1} . By using the induction hypothesis we obtain that

$$\begin{aligned} U_{t+1} &= (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T)) \left(\tilde{U}_t - \sum_{i=1}^t (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T))^{t-i} \hat{E}_i \right) - \hat{E}_{t+1} \\ &= \tilde{U}_{t+1} - \sum_{i=1}^t (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T))^{t+1-i} \hat{E}_i - \hat{E}_{t+1} \\ &= \tilde{U}_{t+1} - \sum_{i=1}^{t+1} (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T))^{t+1-i} \hat{E}_i, \end{aligned}$$

which shows the claimed equation (74).

In order to estimate $\|U_t - \tilde{U}_t\|$ we note first that

$$\begin{aligned} \|\hat{E}_i\| &= \mu \|\mathcal{A}^* \mathcal{A}(U_{i-1} U_{i-1}^T) U_{i-1}\| \\ &\leq \mu \|\mathcal{A}^* \mathcal{A}(U_{i-1} U_{i-1}^T)\| \|U_{i-1}\| \\ &\leq (1 + \delta_1) \mu \|U_{i-1} U_{i-1}^T\|_* \|U_{i-1}\| \\ &= (1 + \delta_1) \mu \|U_{i-1}\|_F^2 \|U_{i-1}\|. \end{aligned}$$

Moreover, we observe that

$$\begin{aligned} \|(\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T))^{t-i} \hat{E}_i\| &\leq \|\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T)\|^{t-i} \|\hat{E}_i\| \\ &\leq (1 + \mu \|\mathcal{A}^* \mathcal{A}(X X^T)\|)^{t-i} \|\hat{E}_i\| \\ &\leq (1 + \mu \lambda_1(M))^{t-i} \|\hat{E}_i\|, \end{aligned}$$

where in the first line we used the submultiplicativity of the operator norm and in the second line we used the triangle inequality. In the third line we used that $\|\mathcal{A}^* \mathcal{A}(XX^T)\| = \lambda_1(M)$. Hence, we have shown that

$$\|U_t - \tilde{U}_t\| \leq \sum_{i=1}^t (1 + \mu\lambda_1(M))^{t-i} (1 + \delta_1) \mu \|U_{i-1}\|_F^2 \|U_{i-1}\|. \quad (75)$$

Note that for all $1 \leq i \leq t^*$ we have that $\|\tilde{U}_{i-1} - U_{i-1}\| \leq \|\tilde{U}_{i-1}\|$, which implies that

$$\begin{aligned} \|U_{i-1}\|_F^2 \|U_{i-1}\| &\leq \min\{r; n\} \|U_{i-1}\|^3 \\ &\leq \min\{r; n\} (\|\tilde{U}_{i-1}\| + \|\tilde{U}_{i-1} - U_{i-1}\|)^3 \\ &\leq 8 \min\{r; n\} \|\tilde{U}_{i-1}\|^3 \\ &\leq 8 \min\{r; n\} \|\text{Id} + \mu\mathcal{A}^* \mathcal{A}(XX^T)\|^{3(i-1)} \|U_0\|^3 \\ &\leq 8 \min\{r; n\} (1 + \mu\lambda_1(M))^{3i-3} \alpha^3 \|U\|^3. \end{aligned}$$

In order to proceed assume now $t \leq t^*$. Then by inequality (75) and the previous inequality we obtain that

$$\begin{aligned} \|U_t - \tilde{U}_t\| &\leq \sum_{i=1}^t (1 + \mu\lambda_1(M))^{t-i} (1 + \delta_1) \mu \|U_{i-1}\|_F^2 \|U_{i-1}\| \\ &\leq 8 \sum_{i=1}^t (1 + \mu\lambda_1(M))^{t-i} (1 + \delta_1) \mu r (1 + \mu\lambda_1(M))^{3i-3} \alpha^3 \|U\|^3 \\ &= 8\alpha^3 \mu \min\{r; n\} (1 + \delta_1) (1 + \mu\lambda_1(M))^{t-1} \sum_{i=1}^t (1 + \mu\lambda_1(M))^{2(i-1)} \|U\|^3 \\ &= 8\alpha^3 \mu \min\{r; n\} (1 + \delta_1) (1 + \mu\lambda_1(M))^{t-1} \frac{(1 + \mu\lambda_1(M))^{2t} - 1}{(1 + \mu\lambda_1(M))^2 - 1} \|U\|^3 \\ &\leq \frac{4}{\lambda_1(M)} \alpha^3 \min\{r; n\} (1 + \delta_1) (1 + \mu\lambda_1(M))^{3t} \|U\|^3. \end{aligned} \quad (76)$$

This shows the claim. \square

A.2 Proof of Lemma 8.2

Proof of Lemma 8.2. First, we note that $\|\tilde{U}_t\| \geq \|\tilde{U}_t^T v_1\|_{\ell_2}$. Then, we observe that

$$\begin{aligned} \tilde{U}_t^T v_1 &= U_0^T (\text{Id} + \mu\mathcal{A}^* \mathcal{A}(XX^T))^t v_1 \\ &= U_0^T \left(\sum_{i=1}^n (1 + \mu\lambda_i) v_i v_i^T \right)^t v_1 \\ &= U_0^T \sum_{i=1}^n (1 + \mu\lambda_i)^t v_i v_i^T v_1 \\ &= (1 + \mu\lambda_1)^t U_0^T v_1. \end{aligned}$$

This proves that $\|\tilde{U}_t\| \geq (1 + \mu\lambda_1(M))^t \|U_0^T v_1\|_{\ell_2}$. From this observation together with Lemma 8.1 it follows for all $t < t^*$ that

$$\frac{\|U_t - \tilde{U}_t\|}{\|\tilde{U}_t\|} \leq \frac{4}{\lambda_1(M)} \alpha^2 \left(\frac{\alpha \min\{r; n\}}{\|U_0^T v_1\|_{\ell_2}} \right) (1 + \delta_1) (1 + \mu\lambda_1(M))^{2t} \|U\|^3.$$

In order to finish the proof, we are going to derive a lower bound for t^* . First we note that by the definition of t^* followed by elementary algebraic manipulations for $t < t^*$ we have

$$\begin{aligned}
& \frac{4}{\lambda_1(M)} \alpha^2 \left(\frac{\alpha \min\{r; n\}}{\|U_0^T v_1\|_{\ell_2}} \right) (1 + \delta_1) (1 + \mu \lambda_1(M))^{2t} \|U\|^3 < 1 \\
& \iff (1 + \mu \lambda_1(M))^{2t} \|U\|^3 < \frac{\lambda_1(M)}{4\alpha^2(1 + \delta_1)} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right) \\
& \iff t < \frac{\ln \left(\frac{\lambda_1(M)}{4\alpha^2(1 + \delta_1)} \frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right)}{2 \ln(1 + \mu \lambda_1(M))}.
\end{aligned}$$

Therefore, we must have

$$t^* \geq \left\lceil \frac{\ln \left(\frac{\lambda_1(M)}{4\alpha^2(1 + \delta_1)} \frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right)}{2 \ln(1 + \mu \lambda_1(M))} \right\rceil.$$

□

A.3 Proof of Lemma 8.3

Proof of Lemma 8.3. Proof of inequality (16): Due to Weyl's inequality we have that

$$\sigma_{r_*}(Z_t U_0 + E_t) \geq \sigma_{r_*}(Z_t U_0) - \|E_t\| \geq \sigma_{r_*}(V_L^T Z_t U_0) - \|E_t\|,$$

where the second inequality follows from the Courant-Fisher minimax theorem (see, e.g., [61, Appendix A]). Now we note that

$$\begin{aligned}
\sigma_{r_*}(V_L^T Z_t U_0) &= \sigma_{\min}(V_L^T Z_t V_L V_L^T U_0) \\
&\geq \sigma_{\min}(V_L^T Z_t V_L) \sigma_{\min}(V_L^T U_0) \\
&= \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U_0) \\
&= \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U).
\end{aligned}$$

This shows the second statement.

Proof of inequality (17): From Weyl's inequality it follows that

$$\sigma_{r_*+1}(Z_t U_0 + E_t) \leq \sigma_{r_*+1}(Z_t U_0) + \|E_t\|. \quad (77)$$

Denote by $U = V_U \Sigma_U W_U^T$ the singular value decomposition of U . Then we can compute that

$$\begin{aligned}
\sigma_{r_*+1}(Z_t U_0) &= \alpha \max_{\mathcal{V}, \dim \mathcal{V} = r_*+1} \min_{x \in \mathcal{V}, \|x\|=1} \|Z_t U x\| \\
&= \alpha \max_{\mathcal{V}, \dim \mathcal{V} = r_*+1} \min_{x \in \mathcal{V}, \|x\|=1} \|Z_t V_U V_U^T U x\| \\
&= \alpha \max_{\mathcal{V}, \dim \mathcal{V} = r_*+1} \min_{x \in \mathcal{V}, \|x\|=1} \|Z_t V_U x\| \|U\| \\
&\leq \alpha \max_{\mathcal{V}, \dim \mathcal{V} = r_*+1} \min_{x \in \mathcal{V}, \|x\|=1} \|Z_t x\| \|U\| \\
&= \alpha \sigma_{r_*+1}(Z_t) \|U\|.
\end{aligned}$$

The first line is due to the Courant-Fisher minimax theorem and $U_0 = \alpha U$. The last line follows again from the Courant-Fisher minimax theorem. Together with inequality (77) this implies the third claim.

Proof of inequality (18): First, we note that

$$Z_t U_0 + E_t = Z_t V_L V_L^T U_0 + \underbrace{Z_t V_{L^\perp} V_{L^\perp}^T U_0 + E_t}_{=: H}.$$

Note that since $V_L^T V_U$ has rank r_* , the matrix $Z_t V_L V_L^T U$ must have rank r_* as well. In particular, since $Z_t V_L V_L^T U = V_L V_L^T Z_t V_L V_L^T U$ this means that L is the subspace spanned by the left-singular vectors of $Z_t V_L V_L^T U$ corresponding to the largest r_* singular values. Due to Wedin's $\sin \theta$ theorem [62] we obtain that

$$\begin{aligned} \|V_{L^\perp}^T V_{L_t}\| &\leq \frac{\|H\|}{\sigma_{r_*}(Z_t V_L V_L^T U_0) - \sigma_{r_*+1}(Z_t U_0 + E_t)} \\ &\leq \frac{\|H\|}{\alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U_0) - \sigma_{r_*+1}(Z_t U_0 + E_t)} \\ &\leq \frac{\|H\|}{\alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) - \alpha \sigma_{r_*+1}(Z_t) \|U\| - \|E_t\|}, \end{aligned}$$

where in the last line we also used (17). (Note that the assumption (15) guarantees that the denominator is positive, which is a necessary condition for an application of Wedin's $\sin \theta$ theorem.) Now we observe that

$$\|H\| \leq \|Z_t V_{L^\perp} V_{L^\perp}^T U\| + \|E_t\| \leq \alpha \|Z_t V_{L^\perp}\| \|U\| + \|E_t\| = \alpha \sigma_{r_*+1}(Z_t) \|U\| + \|E_t\|.$$

Together with the previous inequality chain, this shows inequality (18). \square

A.4 Proof of Lemma 8.4

Before proving Lemma 8.4, we are going to introduce some notation. Let $U_t = \sum_{i=1}^r \sigma_i u_i v_i^T$ be the singular value decomposition of U_t . Define $L_t := \sum_{i=1}^{r_*} \sigma_i u_i v_i^T$ and $N_t := \sum_{i=r_*+1}^r \sigma_i u_i v_i^T$. Denote by $L_t = V_{L_t} \Sigma_{L_t} W_{L_t}^T$ and $N_t = V_{N_t} \Sigma_{N_t} W_{N_t}^T$ the singular value decomposition of those two matrices.

We start by proving the following technical lemma. It says that if the subspace spanned by the columns of X and L_t are aligned, then also the subspaces given by W_t and $W_{L_t^\perp}$ will be closely aligned.

Lemma A.1. *Assume that $\|V_{X^\perp}^T V_{L_t}\| \leq 1/2$. Then it holds that*

$$\|W_{L_t^\perp}^T W_t\| \leq \frac{2\sigma_{r_*+1}(U_t) \|V_{X^\perp}^T V_{L_t}\|}{\sigma_{r_*}(U_t)}.$$

Proof. We note that

$$\begin{aligned}
\|W_{L_t^\perp}^T W_t\| &= \sqrt{\|W_{L_t^\perp}^T W_t W_t^T W_{L_t^\perp}\|} \\
&= \sqrt{\|W_{L_t^\perp}^T U_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U_t W_{L_t^\perp}\|} \\
&= \sqrt{\|W_{L_t^\perp}^T U_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U_t W_{L_t^\perp}\|} \\
&= \sqrt{\|W_{L_t^\perp}^T N_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T N_t W_{L_t^\perp}\|} \\
&= \sqrt{\|W_{L_t^\perp}^T W_{N_t} \Sigma_{N_t} V_{N_t}^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T V_{N_t} \Sigma_{N_t} W_{N_t}^T W_{L_t^\perp}\|} \\
&= \sqrt{\|\Sigma_{N_t} V_{N_t}^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T V_{N_t} \Sigma_{N_t}\|} \\
&\leq \frac{\|\Sigma_{N_t}\| \|V_{N_t}^T V_X\|}{\sigma_{\min}(V_X^T U_t)}.
\end{aligned}$$

In order to control the denominator we note that

$$\begin{aligned}
\sigma_{\min}(V_X^T U_t) &= \sqrt{\sigma_{\min}(V_X^T U_t U_t^T V_X)} \\
&= \sqrt{\sigma_{\min}(V_X^T (L_t L_t^T + N_t N_t^T) V_X)} \\
&\geq \sqrt{\sigma_{\min}(V_X^T L_t L_t^T V_X)} \\
&= \sigma_{\min}(V_X^T L_t) \\
&\geq \sigma_{\min}(V_X^T V_{L_t}) \sigma_{\min}(L_t) \\
&\geq \frac{\sigma_{\min}(L_t)}{2}.
\end{aligned}$$

In the last line we have used the assumption $\|V_{X^\perp}^T V_{L_t}\| \leq 1/2$. Hence, we have shown that

$$\begin{aligned}
\|W_{L_t^\perp}^T W_t\| &\leq \frac{2\|\Sigma_{N_t}\| \|V_{N_t}^T V_X\|}{\sigma_{\min}(\tilde{L})} \\
&= \frac{2\sigma_{r_*+1}(U_t) \|V_{N_t}^T V_X\|}{\sigma_{r_*}(U_t)} \\
&\leq \frac{2\sigma_{r_*+1}(U_t) \|V_{L_t^\perp}^T V_X\|}{\sigma_{r_*}(U_t)} \\
&= \frac{2\sigma_{r_*+1}(U_t) \|V_{X^\perp}^T V_{L_t}\|}{\sigma_{r_*}(U_t)},
\end{aligned}$$

which finishes the proof. □

Now we are in a position to prove Lemma 8.4.

Proof of Lemma 8.4. Proof of inequality (19): First, we observe that due to Lemma A.1 and the assumption $\|V_{X^\perp}^T V_{L_t}\| \leq \frac{1}{8}$ we have that

$$\|W_{L_t^\perp}^T W_t\| \leq \frac{2\sigma_{r_*+1}(U_t) \|V_{X^\perp}^T V_{L_t}\|}{\sigma_{r_*}(U_t)} \leq 1/4. \quad (78)$$

Then, we note that

$$\begin{aligned} \sigma_{r_*}(U_t W_t)^2 &= \sigma_{r_*}(W_t^T U_t^T U_t W_t) \\ &= \sigma_{r_*}(W_t^T (L_t^T L_t + N_t^T N_t) W_t) \\ &\geq \sigma_{r_*}(W_t^T L_t^T L_t W_t) \\ &\geq \sigma_{r_*}(W_t^T W_{L_t})^2 \sigma_{r_*}(L_t)^2 \\ &= (1 - \|W_{L_t^\perp}^T W_t\|^2) \sigma_{r_*}(U_t)^2. \end{aligned}$$

Using inequality (78) we obtain inequality (19).

Proof of inequality (20): Note that

$$\begin{aligned} V_{X^\perp}^T V_{U_t W_t} &= V_{X^\perp}^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t (V_{U_t W_t}^T U_t W_t)^{-1} \\ &= V_{X^\perp}^T U_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}. \end{aligned}$$

By the triangle inequality it follows that

$$\|V_{X^\perp}^T V_{U_t W_t}\| \leq \|V_{X^\perp}^T L_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| + \|V_{X^\perp}^T N_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\|.$$

The second term can be bounded as follows.

$$\begin{aligned} \|V_{X^\perp}^T N_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| &\leq \frac{\|N_t W_t\|}{\sigma_{r_*}(U_t W_t)} \\ &\leq \frac{\|N_t W_{N_t}\| \|W_{N_t}^T W_t\|}{\sigma_{r_*}(U_t W_t)} \\ &= \frac{\sigma_{r_*+1}(U_t) \|W_{N_t}^T W_t\|}{\sigma_{r_*}(U_t W_t)} \\ &\leq \frac{\sigma_{r_*+1}(U_t) \|W_{L_t^\perp}^T W_t\|}{\sigma_{r_*}(U_t W_t)} \\ &\leq 2 \frac{\sigma_{r_*+1}(U_t) \|W_{L_t^\perp}^T W_t\|}{\sigma_{r_*}(U_t)}. \end{aligned} \quad (79)$$

In order to bound the first term, we note that

$$\begin{aligned}
\|V_{X^\perp}^T L_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| &\leq \|V_{X^\perp}^T V_{L_t}\| \|L_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| \\
&\stackrel{(a)}{\leq} \|V_{X^\perp}^T V_{L_t}\| \left(\|U_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| + \|N_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| \right) \\
&= \|V_{X^\perp}^T V_{L_t}\| \left(1 + \|N_t W_t (V_{U_t W_t}^T U_t W_t)^{-1}\| \right) \\
&\stackrel{(c)}{\leq} \|V_{X^\perp}^T V_{L_t}\| \left(1 + \frac{\|N_t W_t\|}{\sigma_{r_*}(U_t W_t)} \right) \\
&\stackrel{(c)}{\leq} \|V_{X^\perp}^T V_{L_t}\| \left(1 + \frac{\sigma_{r_*+1}(U_t) \|W_{L_t^\perp}^T W_t\|}{\sigma_{r_*}(U_t W_t)} \right) \\
&\leq \|V_{X^\perp}^T V_{L_t}\| \left(1 + 2 \frac{\sigma_{r_*+1}(U_t) \|W_{L_t^\perp}^T W_t\|}{\sigma_{r_*}(U_t)} \right) \\
&\stackrel{(d)}{\leq} 3 \|V_{X^\perp}^T V_{L_t}\|.
\end{aligned}$$

In (a) we have used the triangle inequality and inequality (b) follows from inspecting the inequality chain (79). In (c) we used inequality (19) and (d) follows from (78). Combining our results we obtain that

$$\begin{aligned}
\|V_{X^\perp}^T V_{U_t W_t}\| &\leq 3 \|V_{X^\perp}^T V_{L_t}\| + 2 \frac{\sigma_{r_*+1}(U_t) \|W_{L_t^\perp}^T W_t\|}{\sigma_{r_*}(U_t)} \\
&\leq 3 \|V_{X^\perp}^T V_{L_t}\| + 4 \frac{\sigma_{r_*+1}^2(U_t) \|V_{X^\perp}^T V_{L_t}\|}{\sigma_{r_*}^2(U_t)} \\
&\leq 7 \|V_{X^\perp}^T V_{L_t}\|,
\end{aligned}$$

where in the second line we used Lemma A.1. This shows (20).

Proof of inequality (21): We note that

$$\begin{aligned}
\|U_t W_{t,\perp}\| &\leq \|L_t W_{t,\perp}\| + \|N_t W_{t,\perp}\| \\
&\leq \|L_t W_{t,\perp}\| + \|N_t\| \\
&= \|L_t W_{t,\perp}\| + \sigma_{r_*+1}(U_t).
\end{aligned} \tag{80}$$

Observe that $\|L_t W_{t,\perp}\| = \|L_t W_{t,\perp} W_{t,\perp}^T\|$. Then we compute that

$$\begin{aligned}
L_t W_{t,\perp} W_{t,\perp}^T &= L_t \left(\text{Id} - U^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U \right) \\
&= L_t \left(\text{Id} - L_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U \right) \\
&= L_t \left(W_{L_t} W_{L_t}^T - L_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U_t \right).
\end{aligned}$$

Next, we note that

$$\begin{aligned}
V_X^T U U^T V_X &= V_X^T L_t L_t^T V_X + V_X^T N_t N_t^T V_X \\
&= V_X^T L_t W_{L_t} W_{L_t}^T L_t^T V_X + V_X^T N_t N_t^T V_X \\
&= V_X^T L_t W_{L_t} \left(\text{Id} + \underbrace{(V_X^T L_t W_{L_t})^{-1} V_X^T N_t N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1}}_{=:A} \right) W_{L_t}^T L_t^T V_X.
\end{aligned}$$

Now observe that

$$\begin{aligned}
\|A\| &\leq \frac{\|V_X^T N_t N_t^T V_X\|}{\sigma_{\min}(V_X^T L_t W_{L_t})^2} \\
&\leq \frac{\|V_X^T V_{N_t}\|^2 \|N_t\|^2}{\sigma_{\min}(V_X^T V_{L_t})^2 \sigma_{\min}(L_t)^2} \\
&\leq \frac{\|V_X^T V_{L_t^\perp}\|^2 \sigma_{r_*+1}(U_t)^2}{\sigma_{\min}(V_X^T V_{L_t})^2 \sigma_{r_*}(U_t)^2} \\
&= \frac{\|V_{X^\perp}^T V_{L_t}\|^2 \sigma_{r_*+1}(U_t)^2}{\sigma_{\min}(V_X^T V_{L_t})^2 \sigma_{r_*}(U_t)^2} \\
&\leq \frac{\|V_{X^\perp}^T V_{L_t}\|^2 \sigma_{r_*+1}(U_t)^2}{(1 - \|V_{X^\perp}^T V_{L_t}\|^2) \sigma_{r_*}(U_t)^2} \\
&\leq 1/2.
\end{aligned}$$

In the last line we have used the assumption $\|V_{X^\perp}^T V_{L_t}\| \leq \frac{1}{8}$. To continue note that we have

$$\begin{aligned}
&L_t^T V_X (V_X^T U_t U_t^T V_X)^{-1} V_X^T U_t \\
&= L_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T U_t \\
&= W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T U_t \\
&= W_{L_t} (\text{Id} + A)^{-1} W_{L_t}^T + W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \\
&= W_{L_t} W_{L_t}^T - W_{L_t} A (\text{Id} + A)^{-1} W_{L_t}^T + W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t.
\end{aligned}$$

Note that in the last line we used that $\|A\| \leq 1/2$, which we have shown above. It follows that

$$L_t W_{t,\perp} W_{t,\perp}^T = L_t W_{L_t} A (\text{Id} + A)^{-1} W_{L_t}^T - L_t W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t.$$

In particular, by the triangle inequality it follows that

$$\|L_t W_{t,\perp}\| \leq \underbrace{\|L_t W_{L_t} A (\text{Id} + A)^{-1} W_{L_t}^T\|}_{=: (I)} + \underbrace{\|L_t W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t\|}_{=: (II)}. \quad (81)$$

Bounding (I): In order to bound the first term, we note that

$$\begin{aligned}
& L_t W_{L_t} A (\text{Id} + A)^{-1} W_{L_t}^T \\
&= L_t W_{L_t} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} W_{L_t}^T \\
&= L_t W_{L_t} (V_{L_t}^T L_t W_{L_t})^{-1} (V_X^T V_{L_t})^{-1} V_X^T N_t N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} W_{L_t}^T \\
&= V_{L_t} (V_X^T V_{L_t})^{-1} V_X^T N_t N_t^T V_X (V_{L_t}^T V_X)^{-1} (W_{L_t}^T L_t^T V_{L_t})^{-1} (\text{Id} + A)^{-1} W_{L_t}^T.
\end{aligned}$$

It follows that

$$\begin{aligned}
\|L_t W_{L_t} A (\text{Id} + A)^{-1} W_{L_t}^T\| &\leq \frac{\|V_X^T N_t N_t^T V_X\|}{\sigma_{\min}(V_X^T V_{L_t})^2 \sigma_{\min}(\text{Id} + A) \sigma_{\min}(W_{L_t}^T L_t^T V_{L_t})} \\
&\leq \frac{\|V_X^T V_{N_t}\|^2 \|N_t\|^2}{\sigma_{\min}(V_X^T V_{L_t})^2 (1 - \|A\|) \sigma_{r_*}(U_t)} \\
&\leq \frac{\|V_X^T V_{L_t}^\perp\|^2 \sigma_{r_*+1}(U_t)^2}{\sigma_{\min}(V_X^T V_{L_t})^2 (1 - \|A\|) \sigma_{r_*}(U_t)} \\
&= \frac{\|V_{X^\perp}^T V_{L_t}\|^2 \sigma_{r_*+1}(U_t)^2}{\sigma_{\min}(V_X^T V_{L_t})^2 (1 - \|A\|) \sigma_{r_*}(U_t)} \\
&\leq \frac{\sigma_{r_*+1}(U_t)}{2}
\end{aligned}$$

In the last line we have used the assumption $\|V_{X^\perp}^T V_{L_t}\| \leq \frac{1}{8}$ as well as $\|A\| \leq 1/2$.

Bounding (II): We observe that

$$\begin{aligned}
& L_t W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \\
&= L_t W_{L_t} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t - L_t W_{L_t} A (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \\
&= V_{L_t} (V_X^T V_{L_t})^{-1} V_X^T N_t \\
&\quad - L_t W_{L_t} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \\
&= V_{L_t} (V_X^T V_{L_t})^{-1} V_X^T N_t \\
&\quad - V_{L_t} (V_X^T V_{L_t})^{-1} V_X^T N_t N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \\
&= V_{L_t} (V_X^T V_{L_t})^{-1} V_X^T N_t \left(\text{Id} - N_t^T V_X (W_{L_t}^T L_t^T V_X)^{-1} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t \right).
\end{aligned}$$

It follows that

$$\begin{aligned}
& \|L_t W_{L_t} (\text{Id} + A)^{-1} (V_X^T L_t W_{L_t})^{-1} V_X^T N_t\| \\
& \leq \frac{\|V_X^T V_{N_t}\| \sigma_{r_*+1}(U_t)}{\sigma_{\min}(V_X^T V_{L_t})} \left(1 + \frac{\sigma_{r_*+1}(U_t)^2 \|V_X^T V_{N_t}\|^2}{(1 - \|A\|) \sigma_{r_*}(U_t)^2 \sigma_{\min}(V_X^T V_{L_t})^2} \right) \\
& \leq \frac{\|V_X^T V_{L_t^\perp}\| \sigma_{r_*+1}(U_t)}{\sigma_{\min}(V_X^T V_{L_t})} \left(1 + \frac{\sigma_{r_*+1}(U_t)^2 \|V_X^T V_{L_t}\|^2}{(1 - \|A\|) \sigma_{r_*}(U_t)^2 \sigma_{\min}(V_X^T V_{L_t})^2} \right) \\
& \leq \frac{\sigma_{r_*+1}(U_t)}{2}.
\end{aligned}$$

In the last line we have used the assumption $\|V_{X^\perp}^T V_{L_t}\| \leq \frac{1}{8}$ as well as $\|A\| \leq 1/2$. Hence, from inequality (81) it follows that $\|L_t W_{L_t}\| \leq \sigma_{r_*+1}(U_t)$. Inserting this result into inequality (80) we obtain inequality (21), which finishes the proof. \square

A.5 Proof of Lemma 8.5

Before we can prove Lemma 8.5 we will need a technical lemma. In order to state it, recall that L denotes the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of the matrix $M := \mathcal{A}^* \mathcal{A}(XX^T)$ and that $V_L \in \mathbb{R}^{n \times r_*}$ is an orthogonal matrix, whose column span is the subspace L . The following lemma, which follows from standard matrix perturbation theory arguments, shows that for if $\mathcal{A}^* \mathcal{A}(XX^T)$ is sufficiently close to XX^T in spectral norm, then L is aligned with the column space of X . Moreover, it says that the eigenvalues of XX^T are close to the ones of M .

Lemma A.2. *Suppose that $M := \mathcal{A}^* \mathcal{A}(XX^T) = XX^T + \tilde{E}$ with $\|\tilde{E}\| \leq \delta \lambda_{r_*}(XX^T)$ and $\delta < 1/2$. Then it holds that*

$$\begin{aligned}
(1 - \delta) \lambda_1(XX^T) & \leq \lambda_1(M) \leq (1 + \delta) \lambda_1(XX^T), \\
\lambda_{r_*+1}(M) & \leq \delta \lambda_{r_*}(XX^T), \\
\lambda_{r_*}(M) & \geq (1 - \delta) \lambda_{r_*}(XX^T), \\
\|V_{X^\perp}^T V_L\| & \leq 2\delta.
\end{aligned}$$

Proof. The first three inequalities are a direct consequence of Weyl's inequality. In order to prove the fourth inequality, we denote by L the subspace spanned by the eigenvectors corresponding to the r_* largest eigenvalues of M . From the Davis-Kahan sin Θ theorem [63] it follows that

$$\|V_{X^\perp}^T V_L\| \leq \frac{\|\tilde{E}\|}{\lambda_{r_*}(XX^T) - \|\tilde{E}\|} \stackrel{(a)}{\leq} \frac{\delta}{1 - \delta} \stackrel{(b)}{\leq} 2\delta.$$

Inequality (a) follows from the assumption $\|\tilde{E}\| \leq \delta \lambda_{r_*}(XX^T)$. In (b) we used that $\delta \leq \frac{1}{2}$. \square

This allows us to prove Lemma 8.5.

Proof of Lemma 8.5. Due to the assumption (22) we have that $\gamma < 1/2$, if \tilde{c}_2 is chosen small enough, and hence we can apply Lemma 8.3. Hence, we obtain that

$$\sigma_{r_*}(U_t) \geq \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) - \|E_t\| \stackrel{(a)}{\geq} \frac{\alpha}{2} \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U), \quad (82)$$

$$\sigma_{r_*+1}(U_t) \leq \gamma \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U), \quad (83)$$

where in (a) we used that $\gamma \leq 1/2$. Moreover, we also have that

$$\|V_{L^\perp}^T V_{L_t}\| \leq \frac{\alpha \sigma_{r_*+1}(Z_t) \|U\| + \|E_t\|}{\alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) - \alpha \sigma_{r_*+1}(Z_t) \|U\| - \|E_t\|} \leq \frac{\gamma}{1-\gamma}.$$

Now note that

$$\begin{aligned} \|V_{X^\perp}^T V_{L_t}\| &= \|V_X^T V_X^T - V_{L_t} V_{L_t}^T\| \\ &\leq \|V_X^T V_X^T - V_L V_L^T\| + \|V_L V_L^T - V_{L_t} V_{L_t}^T\| \\ &= \|V_{X^\perp}^T V_L\| + \|V_{L^\perp}^T V_{L_t}\| \\ &\leq 2\delta + \frac{\gamma}{1-\gamma}, \end{aligned}$$

where in the last inequality we applied Lemma 8.3 and Lemma A.2. Hence, by our assumptions on δ and γ we can apply Lemma 8.4. Together with the inequality (82) we obtain

$$\sigma_{\min}(U_t W_t) \geq \frac{1}{2} \sigma_{r_*}(U_t) \geq \frac{\alpha}{4} \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U)$$

as well as

$$\begin{aligned} \|V_{X^\perp}^T V_{U_t W_t}\| &\leq 7 \|V_{X^\perp}^T V_{L_t}\| \\ &\leq 7 \left(8\delta + \frac{\gamma}{1-\gamma} \right) \\ &\leq 56(\delta + \gamma). \end{aligned}$$

Moreover, it also follows from Lemma 8.4, inequality (83) and our assumption on γ that

$$\begin{aligned} \|U_t W_{t,\perp}\| &\leq 2 \sigma_{r_*+1}(U_t) \\ &\leq 2 \gamma \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U) \\ &\leq \frac{\kappa^{-2}}{8} \alpha \sigma_{r_*}(Z_t) \sigma_{\min}(V_L^T U). \end{aligned}$$

This finishes the proof. □

A.6 Proof of Lemma 8.6

Proof of Lemma 8.6. In order to apply Lemma 8.5, we want to show that $\gamma \leq \tilde{c}_2 \kappa^{-2}$ for an appropriately chosen $t_* = t$. We are going to show the stronger statement $\gamma \leq c_3 \kappa^{-2}$, where c_3 is a sufficiently

small constant depending only on c , which will be specified later. Note that by the definition of γ it suffices to check the following two conditions.

$$\sigma_{r_*+1}(Z_t)\|U\| \leq \frac{c_3}{2}\sigma_{r_*}(Z_t)\sigma_{\min}(V_L^T U)\kappa^{-2}, \quad (84)$$

$$\|E_t\| \leq \frac{c_3}{2}\alpha\sigma_{r_*}(Z_t)\sigma_{\min}(V_L^T U)\kappa^{-2}. \quad (85)$$

By using the identity $Z_t = (\text{Id} + \mu M)^t$ and by rearranging terms we see that the first inequality is equivalent to the inequality

$$\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \leq \left(\frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)} \right)^t.$$

Hence, if we set

$$t_* = \underbrace{\left\lceil \ln \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right) \left(\ln \left(\frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)} \right) \right)^{-1} \right\rceil}_{=: \sigma},$$

we see that condition (84) is satisfied. Let us check that this choice is feasible, i.e. $t_* \leq t^*$. By Lemma 8.2 and the definition of t_* it suffices to show that

$$\ln \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right) \left(\ln \left(\frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)} \right) \right)^{-1} \leq \frac{\ln \left(\frac{\lambda_1(M)}{4\alpha^2(1+\delta_1)\|U\|^3} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha r} \right) \right)}{8 \ln(1 + \mu\lambda_1(M))}.$$

Next, we note that

$$\begin{aligned} \frac{\ln(1 + \mu\lambda_1(M))}{\ln \left(\frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)} \right)} &= \frac{\ln(1 + \mu\lambda_1(M))}{\ln \left(1 + \frac{\mu(\lambda_{r_*}(M) - \lambda_{r_*+1}(M))}{1 + \mu\lambda_{r_*+1}(M)} \right)} \\ &\leq \frac{\lambda_1(M) \cdot \frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)}}{\frac{\lambda_{r_*}(M) - \lambda_{r_*+1}(M)}{1 + \mu\lambda_{r_*+1}(M)}} \\ &= \frac{\lambda_1(M)(1 + \mu\lambda_{r_*}(M))}{\lambda_{r_*}(M) - \lambda_{r_*+1}(M)} \leq 2\kappa^2, \end{aligned} \quad (86)$$

where in the first inequality we have used the elementary inequality $\frac{x}{1+x} \leq \ln(1+x) \leq x$. in the last inequality we used our assumption on the step size μ , Lemma A.2 as well as our assumption on $\delta > 0$ with a sufficiently small constant c_1 . Hence, $t_* \leq t^*$ is implied by

$$\ln \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right) \leq \frac{1}{9\kappa^2} \ln \left(\frac{\lambda_1(M)}{4\alpha^2(1+\delta_1)\|U\|^3} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right) \right).$$

By rearranging terms we see that this inequality is equivalent to

$$\alpha^2 \leq \frac{\lambda_1(M)}{4(1+\delta_1)\|U\|^3} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right) \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right)^{-9\kappa^2}.$$

Since by assumption $\delta_1 < 1$ and since by Lemma A.2 we have $\lambda_1(M) \geq \frac{1}{2}\|X\|^2$, we observe that this inequality is implied by

$$\alpha^2 \leq \frac{\|X\|^2}{16\|U\|^3} \left(\frac{\|U_0^T v_1\|_{\ell_2}}{\alpha \min\{r; n\}} \right) \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right)^{-9\kappa^2} = \frac{\|X\|^2}{16\|U\|^3} \left(\frac{\|U^T v_1\|_{\ell_2}}{\min\{r; n\}} \right) \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right)^{-9\kappa^2},$$

which follows from assumption (26), which shows $t_* \leq t^*$.

In order to show condition (85), we recall that by Lemma 8.1 (which we can apply since we just showed $t_* \leq t^*$)

$$\|E_{t_*}\| \leq \frac{4}{\lambda_1(M)} \alpha^3 \min\{r; n\} (1 + \delta_1) (1 + \mu\lambda_1(M))^{3t_*} \|U\|^3.$$

Hence, inequality (85) is implied by the inequality

$$\frac{8}{\lambda_1(M)} \alpha^2 \min\{r; n\} (1 + \delta_1) (1 + \mu\lambda_1(M))^{3t_*} \|U\|^3 \leq c_3 (1 + \mu\lambda_{r_*}(M))^t \sigma_{\min}(V_L^T U) \kappa^{-2}.$$

This, in turn, is equivalent to

$$\alpha^2 \leq \frac{c_3 \lambda_1(M) \sigma_{\min}(V_L^T U)}{8 \min\{r; n\} (1 + \delta_1) \kappa^2 \|U\|^3} \left[\frac{1 + \mu\lambda_{r_*}(M)}{(1 + \mu\lambda_1(M))^3} \right]^{t_*}. \quad (87)$$

In order to proceed, we note that

$$\begin{aligned} \left[\frac{1 + \mu\lambda_{r_*}(M)}{(1 + \mu\lambda_1(M))^3} \right]^{t_*} &\geq \exp(-3t_* \ln(1 + \mu\lambda_1(M))) \\ &\geq \exp\left(-\sigma \frac{6 \ln(1 + \mu\lambda_1(M))}{\ln\left(\frac{1 + \mu\lambda_{r_*}(M)}{1 + \mu\lambda_{r_*+1}(M)}\right)}\right). \end{aligned}$$

Hence, using (86), we have shown that

$$\left[\frac{1 + \mu\lambda_{r_*}(M)}{(1 + \mu\lambda_1(M))^3} \right]^{t_*} \geq \exp(-12\sigma\kappa^2).$$

Inserting this into (87) and using the definition of t_* , we have shown that inequality (85) holds, if

$$\alpha^2 \leq \frac{c_2 \|X\|^2 \sigma_{\min}(V_L^T U)}{32r\kappa \|U\|^3} \left(\frac{2\kappa^2\|U\|}{c_3\sigma_{\min}(V_L^T U)} \right)^{-12\kappa^2}$$

holds, which is precisely our assumption on α . In particular, we have shown that $\gamma \leq c_3\kappa^{-2}$, which allows us to apply Lemma 8.5. We obtain that

$$\begin{aligned} \sigma_{\min}(U_{t_*} W_{t_*}) &\geq \frac{\alpha}{4} \sigma_{r_*}(Z_{t_*}) \sigma_{\min}(V_L^T U), \\ \|U_t W_{t_*, \perp}\| &\leq \frac{\kappa^{-2}}{8} \alpha \sigma_{r_*}(Z_{t_*}) \sigma_{\min}(V_L^T U), \\ \|V_{X^\perp}^T V_{U_{t_*} W_{t_*}}\| &\leq 56(\delta + \gamma) \stackrel{(a)}{\leq} c\kappa^{-2}, \end{aligned}$$

where inequality (a) follows from choosing c_1 and c_3 small enough. Setting $\beta := \sigma_{r_*}(Z_{t_*}) \sigma_{\min}(V_L^T U)$ shows inequalities (28), (29), and (30). It remains to verify that $\|U_{t_*}\|$, t_* , and β have the desired properties. We start with t_* . Note that

$$\ln \left(\frac{1 + \mu \lambda_{r_*}(M)}{1 + \mu \lambda_{r_*+1}(M)} \right) \leq \ln(1 + \mu \lambda_{r_*}(M)) \leq \mu \lambda_{r_*}(M) \leq \mu(1 + \delta) \sigma_{r_*}(X)^2$$

as well as

$$\ln \left(\frac{1 + \mu \lambda_{r_*}(M)}{1 + \mu \lambda_{r_*+1}(M)} \right) \geq \frac{\mu \lambda_{r_*}(M)}{1 + \mu \lambda_{r_*}(M)} - \mu \lambda_{r_*+1}(M) \geq \frac{1}{2} \mu \sigma_{r_*}(X)^2.$$

Here we have used the inequalities $\frac{x}{1+x} \leq \ln(1+x) \leq x$, $\lambda_{r_*}(M) \leq \delta \sigma_{\min}(X)^2$, and $(1-\delta) \sigma_{\min}(X)^2 \leq \lambda_{r_*}(M) \leq (1+\delta) \sigma_{r_*}(X)^2$ from Lemma A.2. Hence, these estimates show that t_* has the desired property.

Next, we are going to prove the desired bound for $\|U_{t_*}\|$. We obtain that

$$\begin{aligned} \|U_{t_*}\| &\leq \alpha \|Z_{t_*}\| \|U\| + \|E_{t_*}\| \\ &= \alpha (1 + \mu \lambda_1(M))^{t_*} \|U\| + \|E_{t_*}\| \\ &\stackrel{(a)}{\leq} 2\alpha (1 + \mu \lambda_1(M))^{t_*} \|U\| \\ &\leq 2\alpha \exp \left(2\sigma \frac{\ln(1 + \mu \lambda_1(M))}{\ln \left(\frac{1 + \mu \lambda_{r_*}(M)}{1 + \mu \lambda_{r_*+1}(M)} \right)} \right) \|U\| \\ &\leq 2\alpha \exp(4\sigma \kappa^2) \|U\|, \end{aligned}$$

where (a) follows from (84) and in the last line we used inequality (86). Hence, by inserting the definition of σ we have shown that

$$\begin{aligned} \|U_{t_*}\| &\leq 2\alpha \left(\frac{2\kappa^2 \|U\|}{c_3 \sigma_{\min}(V_L^T U)} \right)^{4\kappa^2} \|U\| \\ &\stackrel{(a)}{\leq} 2\sqrt{\frac{c_2 \|X\|^2 \sigma_{\min}(V_L^T U)}{32 \min\{r; n\} \kappa \|U\|}} \left(\frac{2\kappa^2 \|U\|}{c_3 \sigma_{\min}(V_L^T U)} \right)^{-2\kappa^2} \\ &\leq 3\|X\| \end{aligned}$$

where in inequality (a) we have used the assumption on α . This shows inequality (27). Now let us check that β has the desired property. For that, note that

$$\beta = (1 + \mu \lambda_{r_*}(M))^{t_*} \sigma_{\min}(V_L^T U) = \sigma_{\min}(V_L^T U) \exp(t_* \ln(1 + \mu \lambda_{r_*}(M))).$$

By inserting the definition of t_* and using inequality (86) we can show the upper bound for β in inequality (31). The lower bound follows immediately from the definition of β . This finishes the proof. \square

B Proofs for the saddle avoidance phase and the refinement phase

B.1 Proof of Lemma 9.1

Proof of Lemma 9.1. Let W_t and $W_{t,\perp}$ be defined as before. We note that

$$\begin{aligned}
V_X^T U_{t+1} W_t &= V_X^T (\text{Id} + \mu \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T)) U_t W_t \\
&= V_X^T (\text{Id} + \mu (X X^T - U_t U_t^T) + \mu [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)]) U_t W_t \\
&= V_X^T U_t W_t + \mu \Sigma_X^2 V_X^T U_t W_t - \mu V_X^T U_t U_t^T U_t W_t + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t \\
&= (\text{Id} + \mu \Sigma_X^2) V_X^T U_t W_t - \mu V_X^T U_t W_t W_t^T U_t^T U_t W_t + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t \\
&= (\text{Id} + \mu \Sigma_X^2) V_X^T U_t W_t - \mu V_X^T U_t W_t W_t^T U_t^T V_X V_X^T U_t W_t - \mu V_X^T U_t W_t W_t^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_t \\
&\quad + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t \\
&= (\text{Id} + \mu \Sigma_X^2) V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) - \mu V_X^T U_t W_t W_t^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_t \\
&\quad + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t + \mu^2 \Sigma_X^2 V_X^T U_t W_t W_t^T U_t^T V_X V_X^T U_t W_t \\
&= (\text{Id} + \mu \Sigma_X^2) V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) - \underbrace{\mu V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_t}_{=:A_1} \\
&\quad + \underbrace{\mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t}_{=:A_2} + \underbrace{\mu^2 \Sigma_X^2 V_X^T U_t U_t^T V_X V_X^T U_t W_t}_{=:A_3}.
\end{aligned}$$

First, we want to bring all A_i into the form $P_i V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)$ for $i \in \{1; 2; 3\}$.

Rewriting A_1 : Now let the singular value decomposition of $U_{t+1} W_t \in \mathbb{R}^{n \times r^*}$ be given by $V_{U_{t+1} W_t} \Sigma_{U_{t+1} W_t} W_{U_{t+1} W_t}^T$ with $V_{U_{t+1} W_t} \in \mathbb{R}^{n \times r^*}$. This allows us to compute

$$\begin{aligned}
V_{X^\perp}^T U_t W_t &= V_{X^\perp}^T U_t W_t (V_X^T U_t W_t)^{-1} V_X^T U_t W_t \\
&= V_{X^\perp}^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t (V_X^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t)^{-1} V_X^T U_t W_t \\
&= V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t.
\end{aligned}$$

We compute that

$$\begin{aligned}
&V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_t \\
&= V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t \\
&= V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)^{-1} (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) \\
&= V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} (\text{Id} - \mu V_X^T U_t W_t W_t^T U_t^T V_X)^{-1} V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) \\
&= \underbrace{V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} (\text{Id} - \mu V_X^T U_t W_t U_t^T V_X)^{-1} V_X^T U_t W_t}_{=:P_1} (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t).
\end{aligned} \tag{88}$$

Rewriting A_2 : We observe that

$$\begin{aligned}
U_t W_t &= V_{U_t W_t} V_{U_t W_t}^T U_t W_t (V_X^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t)^{-1} V_X^T U_t W_t \\
&= V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t.
\end{aligned}$$

Hence, we can write

$$\begin{aligned}
& V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_t \\
&= V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t \\
&= V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)^{-1} \\
&\quad \cdot (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) \\
&= \underbrace{V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} (\text{Id} - \mu V_X^T U_t W_t W_t^T U_t^T V_X)^{-1} V_X^T U_t W_t}_{=: P_2} \\
&\quad \cdot (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t).
\end{aligned}$$

Rewriting A_3 : Note that

$$\begin{aligned}
& \Sigma_X^2 V_X^T U_t U_t^T V_X V_X^T U_t W_t \\
&= \Sigma_X^2 V_X^T U_t W_t W_t^T U_t^T V_X V_X^T U_t W_t \\
&= \Sigma_X^2 V_X^T U_t W_t W_t^T U_t^T V_X V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)^{-1} (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t) \\
&= \underbrace{\Sigma_X^2 V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)^{-1} W_t^T U_t^T V_X V_X^T U_t W_t}_{=: P_3} (\text{Id} - \mu W_t^T U_t V_X V_X^T U_t W_t).
\end{aligned}$$

Hence, we have computed that

$$V_X^T U_{t+1} W_t = (\text{Id} + \mu \Sigma_X^2 - \mu P_1 + \mu P_2 + \mu^2 P_3) V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t). \quad (89)$$

It follows that

$$\begin{aligned}
& \sigma_{\min} (V_X^T U_{t+1} W_t) \\
& \geq \sigma_{\min} (\text{Id} + \mu \Sigma_X^2 - \mu P_1 + \mu P_2 + \mu^2 P_3) \sigma_{\min} (V_X^T U_t W_t (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)) \\
& \stackrel{(a)}{=} \sigma_{\min} (\text{Id} + \mu \Sigma_X^2 - \mu P_1 + \mu P_2 + \mu^2 P_3) \sigma_{\min} (V_X^T U_t W_t) (1 - \mu \sigma_{\min}^2 (V_X^T U_t W_t)) \\
& = \sigma_{\min} (\text{Id} + \mu \Sigma_X^2 - \mu P_1 + \mu P_2 + \mu^2 P_3) \sigma_{\min} (V_X^T U_t) (1 - \mu \sigma_{\min}^2 (V_X^T U_t)) \\
& \stackrel{(b)}{\geq} (\sigma_{\min} (\text{Id} + \mu \Sigma_X^2) - \mu \|P_1\| - \mu \|P_2\| - \mu^2 \|P_3\|) \sigma_{\min} (V_X^T U_t) (1 - \mu \sigma_{\min}^2 (V_X^T U_t)) \\
& = (1 + \mu \sigma_{\min}^2 (X) - \mu \|P_1\| - \mu \|P_2\| - \mu^2 \|P_3\|) \sigma_{\min} (V_X^T U_t) (1 - \mu \sigma_{\min}^2 (V_X^T U_t)).
\end{aligned} \quad (90)$$

Equality (a) can be obtained by using the singular value decomposition of $V_X^T U_t W_t$ and the fact that $\mu \leq 1/(\sqrt{3}\|V_X^T U_t\|^2)$, which follows from our assumption on μ . For inequality (b) we used Weyl's

inequality. In order to proceed, we are going to estimate $\|P_1\|$, $\|P_2\|$, and $\|P_3\|$. First, we note that

$$\begin{aligned}
\|P_1\| &\stackrel{(a)}{\leq} \|V_X^T U_t W_t W_t^T U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t}\| \| (V_X^T V_{U_t W_t})^{-1} \| \| (\text{Id} - \mu V_X^T U_t U_t^T V_X)^{-1} \| \\
&\leq \|U_t W_t\| \|V_{X^\perp}^T U_t W_t\| \|V_{X^\perp}^T V_{U_t W_t}\| \| (V_X^T V_{U_t W_t})^{-1} \| \| (\text{Id} - \mu V_X^T U_t U_t^T V_X)^{-1} \| \\
&\leq \|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|^2 \| (V_X^T V_{U_t W_t})^{-1} \| \| (\text{Id} - \mu V_X^T U_t U_t^T V_X)^{-1} \| \\
&= \frac{\|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|^2}{\sigma_{\min}(V_X^T V_{U_t W_t}) \sigma_{\min}(\text{Id} - \mu V_X^T U_t U_t^T V_X)} \\
&= \frac{\|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|^2}{\sigma_{\min}(V_X^T V_{U_t W_t}) (1 - \mu \|V_X^T U_t\|^2)} \\
&\stackrel{(b)}{\leq} 4 \|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|^2 \\
&\stackrel{(c)}{\leq} 36 \|X\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|^2 \\
&\stackrel{(d)}{\leq} \frac{1}{4} \sigma_{\min}(X)^2.
\end{aligned}$$

For inequality (a) we used the submultiplicativity of the spectral norm and the fact that $V_X^T U_t = V_X^T U_t W_t W_t^T$. In (b) we used the assumption $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$ and $\mu \leq c_1 \|X\|^{-2} \kappa^{-2} \leq c \|V_X^T U_t\|^{-2}/9$. In inequality (c) we used the assumption $\|U_t\| \leq 3\|X\|$. Inequality (d) follows from the assumption $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$, where the constant c is chosen sufficiently small.

In order to estimate $\|P_2\|$ we note that

$$\begin{aligned}
\|P_2\| &\stackrel{(a)}{\leq} \|[(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)]\| \| (V_X^T V_{U_t W_t})^{-1} \| \| (\text{Id} - \mu V_X^T U_t U_t^T V_X)^{-1} \| \\
&= \frac{\|[(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)]\|}{\sigma_{\min}(V_X^T V_{U_t W_t}) (1 - \mu \|V_X^T U_t\|^2)} \\
&\stackrel{(b)}{\leq} 4 \|[(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)]\|.
\end{aligned}$$

In (a) we used the submultiplicativity of the spectral norm. In inequality (b) we used the assumption $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$ and $\mu \leq c \|X\|^{-2} \kappa^{-2} \leq c \|V_X^T U_t\|^{-2}/9$. Next, we are going to estimate $\|P_3\|$ by

$$\begin{aligned}
\|P_3\| &\leq \|\Sigma_X^2\| \|V_X^T U_t W_t\| \| (\text{Id} - \mu W_t^T U_t^T V_X V_X^T U_t W_t)^{-1} \| \|W_t^T U_t^T V_X\| \\
&= \frac{\|X\|^2 \|V_X^T U_t W_t\|^2}{1 - \mu \|V_X^T U_t W_t\|^2} \\
&\leq 2 \|X\|^2 \|V_X^T U_t W_t\|^2 \\
&\leq 2 \|X\|^2 \|U_t W_t\|^2 \\
&\leq 18 \|X\|^4.
\end{aligned}$$

In the last line we used the assumption $\|U_t\| \leq 3\|X\|$. Inserting our estimates for $\|P_1\|$, $\|P_2\|$, and

$\|P_3\|$ into (90) we obtain that

$$\begin{aligned}
\sigma_{\min}(V_X^T U_{t+1} W_t) &\geq \left(1 + \frac{3}{4} \mu \sigma_{\min}(X)^2 - 4\mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)\| - 18\mu^2 \|X\|^4\right) \\
&\quad \sigma_{\min}(V_X^T U_t) (1 - \mu \sigma_{\min}^2(V_X^T U_t)) \\
&\stackrel{(a)}{\geq} \left(1 + \frac{1}{2} \mu \sigma_{\min}^2(X)\right) \sigma_{\min}(V_X^T U_t) (1 - \mu \sigma_{\min}^2(V_X^T U_t)) \\
&= \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{2} \mu \sigma_{\min}(X)^2 (1 - \mu \sigma_{\min}^2(V_X^T U_t)) - \mu \sigma_{\min}^2(V_X^T U_t)\right) \\
&\stackrel{(b)}{\geq} \sigma_{\min}(V_X^T U_t) \left(1 + \frac{1}{4} \mu \sigma_{\min}^2(X) - \mu \sigma_{\min}^2(V_X^T U_t)\right).
\end{aligned}$$

Inequality (a) follows from assumption (42) and the assumption $\mu \leq c\kappa^{-2} \|X\|^{-2}$. Inequality (b) is a consequence of our assumption on the step size μ and the assumption $\|U_t\| \leq 3\|X\|$. The final claim follows from the observation that $\sigma_{\min}(V_X^T U_{t+1}) \geq \sigma_{\min}(V_X^T U_{t+1} W_t)$. \square

B.2 Proof of Lemma 9.2

Before we can prove Lemma 9.2, we first need the following technical lemma.

Lemma B.1. *Suppose that the assumptions of Lemma 9.2 are fulfilled with a small enough constant $c > 0$. Then we have that*

$$\|V_{X^\perp}^T V_{U_{t+1} W_t}\| \leq 2\|V_{X^\perp}^T V_{U_t W_t}\| + 2\mu \|(\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T))\|. \quad (91)$$

In particular, it holds that $\|V_{X^\perp}^T V_{U_{t+1} W_t}\| \leq 1/50$.

Proof. We note that

$$U_{t+1} W_t = (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)) U_t W_t.$$

Let $V_{U_t W_t} \Sigma_{U_t W_t} W_{U_t W_t}^T = U_t W_t$ be the singular value decomposition of $U_t W_t$. Set

$$Z := (\text{Id} + \mu \mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)) V_{U_t W_t}.$$

Since $\Sigma_{U_t W_t} W_{U_t W_t}^T$ has full rank by assumption, the matrix $Z = V_Z \Sigma_Z W_Z^T$ has the same column space as the matrix $U_{t+1} W_t$. In particular, it follows that

$$\begin{aligned}
\|V_{X^\perp}^T V_{U_{t+1} W_t}\| &= \|V_{X^\perp}^T V_Z\| \\
&\leq \|V_{X^\perp}^T V_Z \Sigma_Z W_Z^T\| \|(\Sigma_Z W_Z^T)^{-1}\| \\
&= \|V_{X^\perp}^T Z\| \|Z^{-1}\| \\
&= \frac{\|V_{X^\perp}^T Z\|}{\sigma_{\min}(Z)}.
\end{aligned}$$

By Weyl's inequality it holds that

$$\begin{aligned}
\sigma_{\min}(Z) &\geq \sigma_{\min}(V_{U_t W_t}) - \mu \|(\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)) V_{U_t W_t}\| \\
&= 1 - \mu \|(\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)) V_{U_t W_t}\| \\
&= 1 - \mu \|(\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T))\| \\
&\geq 1/2,
\end{aligned}$$

where in the last inequality we used the assumption on the step size μ . Moreover, note that

$$\|V_{X^\perp}^T Z\| \leq \|V_{X^\perp}^T V_{U_t W_t}\| + \mu \|(\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T))\|.$$

This implies inequality (91). Using the assumptions on $\|V_{X^\perp}^T V_{U_t W_t}\|$ and μ , where the constant c is chosen small enough, it follows that $\|V_{X^\perp}^T V_{U_{t+1} W_t}\| \leq 1/50$, which finishes the proof. \square

With all ingredients in place, we can give a proof of Lemma 9.2.

Proof of Lemma 9.2. As a first step we are going to establish a formula for $W_t^T W_{t+1,\perp}$. Recall that $V_X^T U_{t+1} W_{t+1,\perp} = 0$ due to the definition of $W_{t+1,\perp}$. Since $W_t W_t^T + W_{t,\perp} W_{t,\perp}^T = \text{Id}$ we obtain that

$$V_X^T U_{t+1} W_t W_t^T W_{t+1,\perp} = -V_X^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp},$$

or, equivalently,

$$W_t^T W_{t+1,\perp} = -(V_X^T U_{t+1} W_t)^{-1} V_X^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}. \quad (92)$$

Now recall that we want to bound $\|U_{t+1} W_{t+1,\perp}\|$ from above. Note that using $V_X^T U_{t+1} W_{t+1,\perp} = 0$ we have

$$U_{t+1} W_{t+1,\perp} = V_X V_X^T U_{t+1} W_{t+1,\perp} + V_{X^\perp} V_{X^\perp}^T U_{t+1} W_{t+1,\perp} = V_{X^\perp} V_{X^\perp}^T U_{t+1} W_{t+1,\perp},$$

which implies that $\|U_{t+1} W_{t+1,\perp}\| = \|V_{X^\perp}^T U_{t+1} W_{t+1,\perp}\|$. Due to $W_t W_t^T + W_{t,\perp} W_{t,\perp}^T = \text{Id}$ we have that

$$V_{X^\perp}^T U_{t+1} W_{t+1,\perp} = \underbrace{V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp}}_{=(a)} + \underbrace{V_{X^\perp}^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}}_{=(b)}. \quad (93)$$

We are going to consider the two summands individually.

Summand (a): We note that from (92) it follows that

$$V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp} = -V_{X^\perp}^T U_{t+1} W_t (V_X^T U_{t+1} W_t)^{-1} V_X^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}.$$

Let the singular value decomposition of $U_{t+1} W_t \in \mathbb{R}^{n \times r^*}$ be given by $V_{U_{t+1} W_t} \Sigma_{U_{t+1} W_t} W_{U_{t+1} W_t}^T$ with $V_{U_{t+1} W_t} \in \mathbb{R}^{n \times r^*}$. By assumption we have that $V_X^T U_{t+1} W_t$ is invertible, which also implies that $U_{t+1} W_t$ has full-rank. Hence, we can compute that

$$\begin{aligned} V_{X^\perp}^T U_{t+1} W_t (V_X^T U_{t+1} W_t)^{-1} &= V_{X^\perp}^T V_{U_{t+1} W_t} V_{U_{t+1} W_t}^T U_{t+1} W_t (V_X^T V_{U_{t+1} W_t} V_{U_{t+1} W_t}^T U_{t+1} W_t)^{-1} \\ &= V_{X^\perp}^T V_{U_{t+1} W_t} V_{U_{t+1} W_t}^T U_{t+1} W_t (V_{U_{t+1} W_t}^T U_{t+1} W_t)^{-1} (V_X^T V_{U_{t+1} W_t})^{-1} \\ &= V_{X^\perp}^T V_{U_{t+1} W_t} (V_X^T V_{U_{t+1} W_t})^{-1}, \end{aligned}$$

which shows that

$$V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp} = -V_{X^\perp}^T V_{U_{t+1} W_t} (V_X^T V_{U_{t+1} W_t})^{-1} V_X^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}.$$

Moreover, we note that

$$\begin{aligned}
V_X^T U_{t+1} W_{t,\perp} &= V_X^T U_t W_{t,\perp} + \mu V_X^T [\mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T)] U_t W_{t,\perp} \\
&= V_X^T U_t W_{t,\perp} + \mu V_X^T (X X^T - U_t U_t^T) U_t W_{t,\perp} + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_{t,\perp} \\
&\stackrel{(a)}{=} -\mu V_X^T U_t U_t^T U_t W_{t,\perp} + \mu V_X^T [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] U_t W_{t,\perp} \\
&= \mu V_X^T [-U_t U_t^T + [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)]] U_t W_{t,\perp}.
\end{aligned}$$

In equality (a) we used that $V_X^T U_t W_{t,\perp} = 0$ and $X^T U W_{t,\perp} = 0$, which follows from the definition of $W_{t,\perp}$. Hence, we have shown that

$$\begin{aligned}
&V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp} \\
&= \mu V_{X^\perp}^T V_{U_{t+1} W_t} (V_X^T V_{U_{t+1} W_t})^{-1} V_X^T [U_t U_t^T - [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)]] U_t W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp} \\
&\stackrel{(b)}{=} \mu V_{X^\perp}^T V_{U_{t+1} W_t} (V_X^T V_{U_{t+1} W_t})^{-1} \underbrace{V_X^T [U_t U_t^T V_{X^\perp} - [(\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)] V_{X^\perp}]}_{=: M_1} V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp} \\
&= \mu V_{X^\perp}^T V_{U_{t+1} W_t} (V_X^T V_{U_{t+1} W_t})^{-1} M_1 V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}
\end{aligned}$$

In equality (b) above we used that $V_X V_{X^\perp}^T U_t W_{t,\perp} = 0$, which is a consequence of $V_X^T U_t W_{t,\perp} = 0$. It follows that

$$\begin{aligned}
\|V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp}\| &\leq \mu \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \| (V_X^T V_{U_{t+1} W_t})^{-1} \| \|M_1\| \|V_{X^\perp}^T U_t W_{t,\perp}\| \|W_{t,\perp}^T W_{t+1,\perp}\| \\
&\leq \mu \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \| (V_X^T V_{U_{t+1} W_t})^{-1} \| \|M_1\| \|V_{X^\perp}^T U_t W_{t,\perp}\| \\
&= \mu \frac{\|V_{X^\perp}^T V_{U_{t+1} W_t}\| \|M_1\| \|V_{X^\perp}^T U_t W_{t,\perp}\|}{\sigma_{\min}(V_X^T V_{U_{t+1} W_t})}.
\end{aligned}$$

In order to proceed we note that by Lemma B.1 it holds that $\|V_{X^\perp}^T V_{U_{t+1} W_t}\| \leq 1/50$. This implies that

$$\begin{aligned}
\sigma_{\min}(V_X^T V_{U_{t+1} W_t}) &= \sqrt{\sigma_{\min}(V_{U_{t+1} W_t}^T V_X V_X^T V_{U_{t+1} W_t})} \\
&= \sqrt{\sigma_{\min}(V_{U_{t+1} W_t}^T (\text{Id} - V_{X^\perp} V_{X^\perp}^T) V_{U_{t+1} W_t})} \\
&= \sqrt{1 - \|V_{U_{t+1} W_t}^T V_{X^\perp} V_{X^\perp}^T V_{U_{t+1} W_t}\|} \\
&= \sqrt{1 - \|V_{X^\perp}^T V_{U_{t+1} W_t}\|^2} \\
&\geq 1/2.
\end{aligned}$$

Hence, we have shown that

$$\|V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp}\| \leq 2\mu \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \|M_1\| \|V_{X^\perp}^T U_t W_{t,\perp}\|.$$

We can estimate $\|M_1\|$ by

$$\begin{aligned}
\|M_1\| &\stackrel{(a)}{\leq} \|V_X^T U_t U_t^T V_{X^\perp}\| + \|[(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}\| \\
&\stackrel{(b)}{=} \|V_X^T U_t W_t W_t^T U_t^T V_{X^\perp}\| + \|[(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}\| \\
&\leq \|V_X^T U_t W_t\| \|V_{X^\perp}^T U_t W_t\| + \|[(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}\| \\
&\leq \|V_X^T U_t W_t\| \|V_{X^\perp}^T U_t W_t\| + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \\
&\leq \|V_X^T U_t W_t\| \|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_t\| + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \\
&\leq \|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|.
\end{aligned}$$

In inequality (a) we used the triangle inequality and in equality (b) we used that $V_X^T U_t = V_X^T U_t W_t W_t^T$. Hence, we have shown that

$$\begin{aligned}
&\|V_{X^\perp}^T U_{t+1} W_t W_t^T W_{t+1,\perp}\| \\
&\leq 2\mu (\|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \|V_{X^\perp}^T U_t W_{t,\perp}\|) \\
&\leq 2\mu (\|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \|U_t W_{t,\perp}\|) \\
&\leq 2\mu (9\|V_{X^\perp}^T V_{U_t W_t}\| \|X\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \|V_{X^\perp}^T V_{U_{t+1} W_t}\| \|U_t W_{t,\perp}\|) \\
&\stackrel{(a)}{\leq} \mu (36\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|X\|^2 + 4\|V_{X^\perp}^T V_{U_t W_t}\| \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \\
&\quad + \mu^2 \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|^2) \|U_t W_{t,\perp}\| \\
&\stackrel{(b)}{\leq} \mu (36\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|X\|^2 + 2\|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \|U_t W_{t,\perp}\|),
\end{aligned}$$

where in inequality (a) we used Lemma B.1. In the inequality (b) we used assumption $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-1}$ and the assumption on the step size μ .

Summand (b): First, we compute that

$$\begin{aligned}
&V_{X^\perp}^T U_{t+1} W_{t,\perp} \\
&= V_{X^\perp}^T U_t W_{t,\perp} + \mu V_{X^\perp}^T (XX^T - U_t U_t^T) U_t W_{t,\perp} + \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\
&= V_{X^\perp}^T U_t W_{t,\perp} - \mu V_{X^\perp}^T U_t U_t^T U_t W_{t,\perp} + \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\
&= V_{X^\perp}^T U_t W_{t,\perp} - \mu V_{X^\perp}^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} + \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\
&= (\text{Id} - \mu V_{X^\perp}^T U_t U_t^T V_{X^\perp} - \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}) V_{X^\perp}^T U_t W_{t,\perp} \\
&= V_{X^\perp}^T U_t W_{t,\perp} - \mu V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} - \mu V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} \\
&\quad + \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} \\
&= (\text{Id} - \mu V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp} + \mu V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T U_t W_{t,\perp}) \\
&\quad - \mu^2 (V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp} - V_{X^\perp}^T [(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)] V_{X^\perp}) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp}.
\end{aligned}$$

Set for brevity of notation $M_2 := V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp}$ and $M_3 := V_{X^\perp}^T (\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T) V_{X^\perp}$.

Hence, we have computed that

$$\begin{aligned}
& V_{X^\perp}^T U_{t+1} W_{t,\perp} \\
&= (\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T U_t W_{t,\perp}) - \mu^2 (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp} \\
&= (\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp}) - \mu^2 (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp}.
\end{aligned}$$

Hence, we obtain that

$$\begin{aligned}
& \|V_{X^\perp}^T U_{t+1} W_{t,\perp} W_{t,\perp}^T W_{t+1,\perp}\| \\
&\leq \|(\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp}) \\
&\quad - \mu^2 (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp}\| \|W_{t,\perp}^T W_{t+1,\perp}\| \\
&\leq \|(\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp}) \\
&\quad - \mu^2 (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp}\| \\
&\leq \|(\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp})\| \\
&\quad + \mu^2 \|(M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp}\|.
\end{aligned}$$

In order to proceed, we compute that

$$\begin{aligned}
& \|(\text{Id} - \mu M_2 + \mu M_3) V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp})\| \\
&\stackrel{(a)}{\leq} \|\text{Id} - \mu M_2 + \mu M_3\| \|V_{X^\perp}^T U_t W_{t,\perp} (\text{Id} - \mu W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp})\| \\
&\stackrel{(b)}{=} \|\text{Id} - \mu M_2 + \mu M_3\| \|V_{X^\perp}^T U_t W_{t,\perp}\| (1 - \mu \|V_{X^\perp}^T U_t W_{t,\perp}\|^2) \\
&\stackrel{(c)}{\leq} (\|\text{Id} - \mu M_2\| + \mu \|M_3\|) \|V_{X^\perp}^T U_t W_{t,\perp}\| (1 - \mu \|V_{X^\perp}^T U_t W_{t,\perp}\|^2) \\
&\stackrel{(d)}{\leq} (1 + \mu \|M_3\|) \|V_{X^\perp}^T U_t W_{t,\perp}\| (1 - \mu \|V_{X^\perp}^T U_t W_{t,\perp}\|^2) \\
&= (1 + \mu \|M_3\|) \|U_t W_{t,\perp}\| (1 - \mu \|U_t W_{t,\perp}\|^2) \\
&\leq \|U_t W_{t,\perp}\| (1 - \mu \|U_t W_{t,\perp}\|^2 + \mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T) V_{X^\perp}\| - \mu^2 \|M_3\| \|U_t W_{t,\perp}\|^2) \\
&\leq \|U_t W_{t,\perp}\| (1 - \mu \|U_t W_{t,\perp}\|^2 + \mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T) V_{X^\perp}\|) \\
&\leq \|U_t W_{t,\perp}\| (1 - \mu \|U_t W_{t,\perp}\|^2 + \mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(X X^T - U_t U_t^T)\|).
\end{aligned}$$

The inequality (a) follows from the submultiplicativity of the spectral norm. Equality (b) can be seen by using the singular value decomposition of $V_{X^\perp}^T U_t W_{t,\perp}$ and the assumption $\mu \leq c_1 \|X\|^{-2} \leq 1/(\sqrt{3} \|V_{X^\perp}^T U_t W_{t,\perp}\|^2)$. Inequality (c) follows from the triangle inequality. For inequality (d) we used that $0 \leq \text{Id} - \mu V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{X^\perp} \leq \text{Id}$, which again is a consequence of our assumptions on μ and

$\|U_t\| \leq 3\|X\|$. For the $O(\mu^2)$ -term we note that

$$\begin{aligned}
& \| (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp} \| \\
&= \| (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} \| \\
&\leq \|M_2 - M_3\| \|V_{X^\perp}^T U_t W_{t,\perp}\|^3 \\
&\leq (\|M_2\| + \|M_3\|) \|V_{X^\perp}^T U_t W_{t,\perp}\|^3 \\
&\leq (\|V_{X^\perp}^T U_t W_t\|^2 + \|M_3\|) \|V_{X^\perp}^T U_t W_{t,\perp}\|^3 \\
&\leq (\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|U_t W_t\|^2 + \|M_3\|) \|V_{X^\perp}^T U_t W_{t,\perp}\|^3 \\
&= (\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|U_t W_t\|^2 + \|M_3\|) \|U_t W_{t,\perp}\|^3 \\
&= (\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T) V_{X^\perp}\|) \|U_t W_{t,\perp}\|^3 \\
&\leq (\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|) \|U_t W_{t,\perp}\|^3.
\end{aligned}$$

It follows that

$$\begin{aligned}
& \mu^2 \| (M_2 - M_3) V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T U_t W_{t,\perp} \| \\
&\leq \mu^2 (\|V_{X^\perp}^T V_{U_t W_t}\|^2 \|U_t W_t\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|) \|U_t W_{t,\perp}\|^3 \\
&\stackrel{(a)}{\leq} \mu^2 (9\|X\|^2 + \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|) \|U_t W_{t,\perp}\|^3 \\
&\stackrel{(b)}{\leq} \frac{\mu}{2} \|U_t W_{t,\perp}\|^3,
\end{aligned}$$

In (a) we used that $\|U_t\| \leq 3\|X\|$. In (b) we used our assumption on the step size μ . This implies that

$$\|(b)\| \leq \|U_t W_{t,\perp}\| \left(1 - \frac{\mu}{2} \|U_t W_{t,\perp}\|^2 + \mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \right).$$

Conclusion: Putting things together it follows

$$\begin{aligned}
& \|V_{X^\perp}^T U_{t+1} W_{t+1,\perp}\| \\
&\leq \|(a)\| + \|(b)\| \\
&\leq \left(1 - \frac{\mu}{2} \|U_t W_{t,\perp}\|^2 + 36\mu \|V_{X^\perp}^T V_{U_t W_t}\|^2 \|X\|^2 + 3\mu \|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\| \right) \|U_t W_{t,\perp}\|.
\end{aligned}$$

This finishes the proof. \square

B.3 Proof of Lemma 9.3

We define the inverse of the square root of a symmetric, positiv definite matrix $A = V_A \Sigma_A V_A^T$ by $A^{-1/2} := V_A \Sigma_A^{-1/2} V_A^T$, where $(\Sigma_A^{-1/2})_{ii} := \frac{1}{\sqrt{A_{ii}}}$. We will need the following technical lemma, which gives a bound on the first order Taylor-approximation of the matrix inverse square root.

Lemma B.2. *Let A be a symmetric, positiv definite matrix such that $\|A\| \leq 1/2$. Then there holds*

$$\|(Id + A)^{-1/2} - Id + \frac{1}{2}A\| \leq 3\|A\|^2.$$

Proof. Since A is symmetric, this can be readily deduced from the (one-dimensional) Taylor's theorem. Indeed, we have that for $|x| \leq 1/2$ that

$$\left| \frac{1}{\sqrt{1+x}} - 1 + \frac{x}{2} \right| \leq \sup_{|z| \leq 1/2} \left| \frac{3}{8} \cdot (1+z)^{-5/2} \cdot x^2 \right| \leq 3x^2.$$

□

The next technical lemma shows that the orthogonal matrices W_t and W_{t+1} span approximately the same column space.

Lemma B.3. *Assume that the assumptions of Lemma 9.3 are fulfilled. Then it holds that*

$$\|W_{t,\perp}^T W_{t+1}\| \leq \mu \left(\frac{1}{6400} \sigma_{\min}(X)^2 + \|U_t W_t\| \|U_t W_{t,\perp}\| \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 4\mu \|[(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \quad (94)$$

and

$$\sigma_{\min}(W_t^T W_{t+1}) \geq 1/2.$$

Proof. Due to $V_X^T U_{t+1} = V_X^T U_{t+1} W_{t+1} W_{t+1}^T$ we observe that

$$\|W_{t,\perp}^T W_{t+1}\| = \|W_{t,\perp}^T U_{t+1}^T V_X (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2}\|.$$

We note that

$$\begin{aligned} & V_X^T U_{t+1} W_{t,\perp} \\ &= V_X^T (Id + \mu (XX^T - U_t U_t^T)) U_t W_{t,\perp} - \mu V_X^T [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\ &= -\mu V_X^T U_t U_t^T U_t W_{t,\perp} - \mu V_X^T [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\ &= -\mu V_X^T U_t W_t W_t^T U_t^T U_t W_{t,\perp} - \mu V_X^T [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\ &= -\mu V_X^T U_t W_t W_t^T U_t^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} - \mu V_X^T [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t W_{t,\perp} \\ &= -\mu V_X^T U_t W_t W_t^T U_t^T V_{U_t W_t} V_{U_t W_t}^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} - \mu V_X^T [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t W_{t,\perp}. \end{aligned}$$

It follows that

$$\begin{aligned} \|W_{t,\perp}^T W_{t+1}\| &\leq \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_t W_t W_t^T U_t^T V_{U_t W_t} V_{U_t W_t}^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} \| \\ &\quad + \mu \| [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \| \|U_t W_{t,\perp}\| \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} \| \\ &= \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_t W_t W_t^T U_t^T V_{U_t W_t} V_{U_t W_t}^T V_{X^\perp} V_{X^\perp}^T U_t W_{t,\perp} \| \\ &\quad + \mu \frac{\| [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \| \|U_t W_{t,\perp}\|}{\sigma_{\min}(V_X^T U_{t+1})} \\ &\leq \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_t W_t \| \|U_t W_t\| \|U_t W_{t,\perp}\| \|V_{X^\perp}^T V_{U_t W_t}\| \\ &\quad + \mu \frac{\| [(Id - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \| \|U_t W_{t,\perp}\|}{\sigma_{\min}(V_X^T U_{t+1})}. \end{aligned}$$

We note that

$$\begin{aligned} (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_t W_t &= (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_{t+1} W_t \\ &\quad - \mu (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) U_t W_t. \end{aligned}$$

It follows that

$$\begin{aligned} &\| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_t W_t \| \\ &\leq \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_{t+1} W_t \| + \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) U_t W_t \| \\ &\leq \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T U_{t+1} \| + \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) U_t W_t \| \\ &= 1 + \mu \| (V_X^T U_{t+1} U_{t+1}^T V_X)^{-1/2} V_X^T \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) U_t W_t \| \\ &\leq 1 + \mu \frac{\| \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) \| \| U_t W_t \|}{\sigma_{\min}(V_X^T U_{t+1})}. \end{aligned}$$

Hence, we obtain that

$$\begin{aligned} \| W_{t,\perp}^T W_{t+1} \| &\leq \mu \left(1 + \mu \frac{\| \mathcal{A}^* \mathcal{A} (X X^T - U_t U_t^T) \| \| U_t W_t \|}{\sigma_{\min}(V_X^T U_{t+1})} \right) \| U_t W_t \| \| U_t W_{t,\perp} \| \| V_{X^\perp}^T V_{U_t W_t} \| \\ &\quad + \mu \frac{\| (\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T) \| \| U_t W_{t,\perp} \|}{\sigma_{\min}(V_X^T U_{t+1})}. \end{aligned} \tag{95}$$

Next, we are going to show $\sigma_{\min}(V_X^T U_{t+1}) \geq \frac{\sigma_{\min}(U_t W_t)}{2}$. We note that

$$\begin{aligned} \sigma_{\min}(V_X^T U_{t+1}) &\geq \sigma_{\min}(V_X^T U_{t+1} W_t) \\ &= \sigma_{\min}(V_X^T (\text{Id} + \mu [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)]) U_t W_t) \\ &= \sigma_{\min}(V_X^T (\text{Id} + \mu [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)]) V_{U_t W_t} V_{U_t W_t}^T U_t W_t) \\ &\geq \sigma_{\min}(V_X^T (\text{Id} + \mu [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)]) V_{U_t W_t}) \sigma_{\min}(V_{U_t W_t}^T U_t W_t) \\ &= \sigma_{\min}(V_X^T V_{U_t W_t} + \mu V_X^T [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)] V_{U_t W_t}) \sigma_{\min}(U_t W_t) \\ &\geq (\sigma_{\min}(V_X^T V_{U_t W_t}) - \mu \| V_X^T [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)] V_{U_t W_t} \|) \sigma_{\min}(U_t W_t). \end{aligned}$$

We observe that due to our assumption on $\| V_{X^\perp}^T V_{U_t W_t} \|$ we have that

$$\sigma_{\min}(V_X^T V_{U_t W_t}) = \sqrt{1 - \| V_{X^\perp}^T V_{U_t W_t} \|^2} \geq \frac{3}{4}.$$

Next, we note that due to assumptions (43), (45) and $\| U_t \| \leq 3 \| X \|$ we have that

$$\begin{aligned} \mu \| V_X^T [(\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)] V_{U_t W_t} \| &\leq \mu \| (\mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T) \| \\ &\leq 10\mu \| X \|^2 + \mu \| (\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T) \| \leq \frac{1}{4}. \end{aligned}$$

Hence, we have shown that $\sigma_{\min}(V_X^T U_{t+1}) \geq \frac{\sigma_{\min}(U_t W_t)}{2}$. This implies due to inequality (95) that

$$\begin{aligned} \|W_{t,\perp}^T W_{t+1}\| &\leq \mu \left(1 + 2\mu \frac{\|\mathcal{A}^* \mathcal{A}(XX^T - U_t U_t^T)\| \|U_t W_t\|}{\sigma_{\min}(U_t W_t)} \right) \|U_t W_t\| \|U_t W_{t,\perp}\| \|V_{X^\perp}^T V_{U_t W_t}\| \\ &\quad + 2\mu \frac{\|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \|U_t W_{t,\perp}\|}{\sigma_{\min}(U_t W_t)} \\ &\stackrel{(a)}{\leq} \mu \|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_t\| \|U_t W_{t,\perp}\| + 4\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \\ &\quad + 4\mu^2 \|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\|. \end{aligned}$$

In inequality (a) we have used the assumption that $\|U_t W_{t,\perp}\| \leq 2\sigma_{\min}(U_t W_t)$. In order to proceed, we note that

$$\begin{aligned} \|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| &\leq \|XX^T - U_t U_t^T\| + \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \\ &\leq 11\|X\|^2, \end{aligned}$$

where we used the assumption $\|U_t W_t\| \leq 3\|X\|$ and $\|[(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \leq c\sigma_{\min}(X)^2$. Hence, we obtain that

$$\|W_{t,\perp}^T W_{t+1}\| \leq \mu \left(\frac{1}{6400} \sigma_{\min}(X)^2 + \|U_t W_t\| \|U_t W_{t,\perp}\| \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 4\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\|,$$

where we also used $\mu \leq c\kappa^{-2}\|X\|^{-2}$. Hence, we have shown inequality (94).

In order to finish the proof we note that

$$\begin{aligned} \|W_{t,\perp}^T W_{t+1}\| &\stackrel{(a)}{\leq} \mu \left(\frac{1}{6400} \sigma_{\min}(X)^2 + \|U_t W_t\| \|U_t W_{t,\perp}\| \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 4\mu c \sigma_{\min}(X)^2 \\ &\stackrel{(b)}{\leq} \mu \left(\frac{1}{6400} \sigma_{\min}(X)^2 + 9\|X\|^2 \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 4c\mu \sigma_{\min}(X)^2 \\ &\stackrel{(c)}{\lesssim} c. \end{aligned}$$

In inequality (a) we used the assumption $\|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \leq c\sigma_{\min}(X)^2$. In inequality (b) we used $\|U_t W_{t,\perp}\| \leq 3\|X\|$ and $\|U_t W_t\| \leq 3\|X\|$. To obtain inequality (c) we used the assumption $\mu \leq c\|X\|^{-2}$. By choosing $c > 0$ small enough we obtain that $\|W_{t,\perp}^T W_{t+1}\| \leq 1/2$. Note that this implies that

$$\sigma_{\min}(W_t^T W_{t+1}) = \sqrt{1 - \|W_{t,\perp}^T W_{t+1}\|^2} \geq 1/2,$$

which finishes the proof. \square

Now we have provided all the technical preliminaries to prove Lemma 9.3.

Proof of Lemma 9.3. In order to simplify notation, we define $M := \mathcal{A}^* \mathcal{A}(XX^T - U_t U_t^T)$. Hence, we may write

$$U_{t+1} = (\text{Id} + \mu M) U_t.$$

Now we note that

$$\begin{aligned}
U_{t+1}W_{t+1} &= (\text{Id} + \mu M) U_t W_{t+1} \\
&= (\text{Id} + \mu M) U_t W_t W_t^T W_{t+1} + (\text{Id} + \mu M) U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} \\
&= (\text{Id} + \mu M) V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T W_{t+1} + (\text{Id} + \mu M) U_t W_{t,\perp} W_{t,\perp}^T W_{t+1}.
\end{aligned} \tag{96}$$

Note that $V_{U_t W_t}^T U_t W_t W_t^T W_{t+1}$ is invertible, since $V_{U_t W_t}^T U_t W_t$ is invertible by assumption (44) and $W_t^T W_{t+1}$ is invertible by Lemma B.3. Hence, we see that

$$\begin{aligned}
&(\text{Id} + \mu M) U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} \\
&= (\text{Id} + \mu M) U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} \left(V_{U_t W_t}^T U_t W_t W_t^T W_{t+1} \right)^{-1} V_{U_t W_t}^T U_t W_t W_t^T W_{t+1} \\
&= (\text{Id} + \mu M) \underbrace{U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} \left(V_{U_t W_t}^T U_t W_t W_t^T W_{t+1} \right)^{-1} V_{U_t W_t}^T V_{U_t W_t} U_t W_t W_t^T W_{t+1}}_{=P} \\
&= (\text{Id} + \mu M) P V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T W_{t+1}.
\end{aligned}$$

Hence, by inserting the last equation into equation (96) we obtain that

$$U_{t+1}W_{t+1} = (\text{Id} + \mu M) (\text{Id} + P) V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T W_{t+1}$$

Recall that $V_{U_t W_t}^T U_t W_t W_t^T W_{t+1}$ is an invertible matrix. This implies that the span of the left-singular vectors of

$$Z := (\text{Id} + \mu M) (\text{Id} + P) V_{U_t W_t}$$

is the same as the span of the left-singular vectors of $U_{t+1}W_{t+1}$. Let $V_Z \Sigma_Z W_Z^T$ be the singular value decomposition of Z . From these considerations it follows that

$$\|V_{X^\perp}^T V_{U_{t+1}W_{t+1}}\| = \|V_{X^\perp}^T V_Z\| = \|V_{X^\perp}^T V_Z W_Z^T\|.$$

Next, we note that

$$\begin{aligned}
V_Z W_Z^T &= Z \left(Z^T Z \right)^{-1/2} \\
&= (\text{Id} + \mu M) (\text{Id} + P) V_{U_t W_t} \left(V_{U_t W_t}^T (\text{Id} + P^T) (\text{Id} + \mu M)^2 (\text{Id} + P) V_{U_t W_t} \right)^{-1/2}.
\end{aligned}$$

We note that

$$\begin{aligned}
(\text{Id} + \mu M) (\text{Id} + P) &= \text{Id} + \underbrace{\mu M + P + \mu M P}_{=:B} \\
&= \text{Id} + \underbrace{\mu (X X^T - U_t U_t^T)}_{=:B_1} + \underbrace{\mu (\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)}_{=:B_2} \\
&\quad + \underbrace{U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} \left(V_{U_t W_t}^T U_t W_t W_t^T W_{t+1} \right)^{-1} V_{U_t W_t}^T}_{=:B_3} + \underbrace{\mu M P}_{=:B_4}.
\end{aligned}$$

Hence, we have that

$$\begin{aligned}
& Z (Z^T Z)^{-1/2} \\
&= (\text{Id} + B) V_{U_t W_t} (V_{U_t W_t}^T (\text{Id} + B + B^T + B^T B) V_{U_t W_t})^{-1/2} \\
&= (\text{Id} + B) V_{U_t W_t} (\text{Id} + V_{U_t W_t}^T B V_{U_t W_t} + V_{U_t W_t}^T B^T V_{U_t W_t} + V_{U_t W_t}^T B^T B V_{U_t W_t})^{-1/2}.
\end{aligned}$$

It follows from Lemma B.2 that

$$\begin{aligned}
& (\text{Id} + V_{U_t W_t}^T B V_{U_t W_t} + V_{U_t W_t}^T B^T V_{U_t W_t} + V_{U_t W_t}^T B^T B V_{U_t W_t})^{-1/2} \\
&= \text{Id} - \frac{1}{2} (V_{U_t W_t}^T B V_{U_t W_t} + V_{U_t W_t}^T B^T V_{U_t W_t} + V_{U_t W_t}^T B^T B V_{U_t W_t}) + C,
\end{aligned}$$

where C is matrix, which satisfies

$$\|C\| \leq 3 \|V_{U_t W_t}^T B V_{U_t W_t} + V_{U_t W_t}^T B^T V_{U_t W_t} + V_{U_t W_t}^T B^T B V_{U_t W_t}\|^2. \quad (97)$$

It follows that

$$\begin{aligned}
& Z (Z^T Z)^{-1/2} \\
&= (\text{Id} + B) V_{U_t W_t} \left(\text{Id} - \frac{1}{2} (V_{U_t W_t}^T B V_{U_t W_t} + V_{U_t W_t}^T B^T V_{U_t W_t} + V_{U_t W_t}^T B^T B V_{U_t W_t}) + C \right) \\
&= V_{U_t W_t} + B V_{U_t W_t} - \frac{1}{2} (\text{Id} + B) V_{U_t W_t} V_{U_t W_t}^T (B + B^T) V_{U_t W_t} - D,
\end{aligned}$$

where we have set

$$D = (\text{Id} + B) V_{U_t W_t} \left(\frac{1}{2} V_{U_t W_t}^T B^T B V_{U_t W_t} - C \right). \quad (98)$$

Hence,

$$\begin{aligned}
& V_{X^\perp}^T Z (Z^T Z)^{-1/2} \\
&= V_{X^\perp}^T \left(\text{Id} + B - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B + B^T) \right) V_{U_t W_t} - \frac{1}{2} V_{X^\perp}^T B V_{U_t W_t} V_{U_t W_t}^T (B + B^T) V_{U_t W_t} - V_{X^\perp}^T D \\
&= \underbrace{V_{X^\perp}^T \left(\text{Id} + B_1 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_1 + B_1^T) \right) V_{U_t W_t}}_{=:(I)} \\
&\quad + \underbrace{V_{X^\perp}^T \left(B_2 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_2 + B_2^T) \right) V_{U_t W_t}}_{=:(II)} \\
&\quad + \underbrace{V_{X^\perp}^T \left(B_3 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_3 + B_3^T) \right) V_{U_t W_t}}_{=:(III)} \\
&\quad + \underbrace{V_{X^\perp}^T \left(B_4 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_4 + B_4^T) \right) V_{U_t W_t}}_{=:(IV)} - \underbrace{\frac{1}{2} V_{X^\perp}^T B V_{U_t W_t} V_{U_t W_t}^T (B + B^T) V_{U_t W_t}}_{=:(V)} - \underbrace{V_{X^\perp}^T D}_{=:(VI)}.
\end{aligned}$$

Estimating (I): We observe that

$$\begin{aligned}
& V_{X^\perp}^T \left(\text{Id} + B_1 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_1 + B_1^T) \right) V_{U_t W_t} \\
&= V_{X^\perp}^T \left(\text{Id} + \mu (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) (X X^T - U_t U_t^T) \right) V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} + \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) (X X^T - U_t U_t^T) V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} + \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) X X^T V_{U_t W_t} - \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) U_t U_t^T V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} - \mu V_{X^\perp}^T V_{U_t W_t} V_{U_t W_t}^T X X^T V_{U_t W_t} - \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) U_t U_t^T V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} - \mu V_{X^\perp}^T V_{U_t W_t} V_{U_t W_t}^T X X^T V_{U_t W_t} - \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} (\text{Id} - \mu V_{U_t W_t}^T X X^T V_{U_t W_t}) - \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{U_t W_t} \\
&= V_{X^\perp}^T V_{U_t W_t} (\text{Id} - \mu V_{U_t W_t}^T X X^T V_{U_t W_t}) - \mu V_{X^\perp}^T (\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{X^\perp} V_{X^\perp}^T V_{U_t W_t}.
\end{aligned}$$

It follows that

$$\begin{aligned}
\|(I)\| &\leq \|V_{X^\perp}^T V_{U_t W_t}\| \left(1 - \mu \sigma_{\min}(V_{U_t W_t}^T X X^T V_{U_t W_t}) \right) + \mu \|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_{t,\perp}\|^2 \\
&\leq \|V_{X^\perp}^T V_{U_t W_t}\| \left(1 - \frac{\mu}{2} \sigma_{\min}(X)^2 \right) + \mu \|V_{X^\perp}^T V_{U_t W_t}\| \|U_t W_{t,\perp}\|^2 \\
&= \|V_{X^\perp}^T V_{U_t W_t}\| \left(1 - \frac{\mu}{2} \sigma_{\min}(X)^2 + \mu \|U_t W_{t,\perp}\|^2 \right) \\
&\leq \|V_{X^\perp}^T V_{U_t W_t}\| \left(1 - \frac{\mu}{3} \sigma_{\min}(X)^2 \right).
\end{aligned}$$

Bounding (II): We observe that

$$\begin{aligned}
& V_{X^\perp}^T \left(B_2 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_2 + B_2^T) \right) V_{U_t W_t} \\
&= \mu V_{X^\perp}^T ((\text{Id} - V_{U_t W_t} V_{U_t W_t}^T) (\mathcal{A}^* \mathcal{A} - \text{Id}) (X X^T - U_t U_t^T)) V_{U_t W_t}.
\end{aligned}$$

It follows that

$$\|V_{X^\perp}^T \left(B_2 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_2 + B_2^T) \right) V_{U_t W_t}\| \leq \mu \|(\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)\|.$$

Estimating (III): First, we recall that

$$\begin{aligned}
B_3 &= U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} (V_{U_t W_t}^T U_t W_t W_t^T W_{t+1})^{-1} V_{U_t W_t}^T \\
&= U_t W_{t,\perp} W_{t,\perp}^T W_{t+1} (W_t^T W_{t+1})^{-1} (V_{U_t W_t}^T U_t W_t)^{-1} V_{U_t W_t}^T.
\end{aligned}$$

Before we proceed further, we need to understand $W_{t,\perp}^T W_{t+1}$ and $W_t^T W_{t+1}$. By Lemma B.3 it holds that

$$\|W_{t,\perp}^T W_{t+1}\| \leq \mu \left(\frac{1}{800} \sigma_{\min}^2(X) + \|U_t W_t\| \|U_t W_{t,\perp}\| \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 4\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)\|$$

and

$$\sigma_{\min}(W_t^T W_{t+1}) \geq 1/2.$$

It follows that

$$\begin{aligned}
\|B_3\| &\leq \|W_{t,\perp}^T W_{t+1}\| \|U_t W_{t,\perp}\| \| (W_t^T W_{t+1})^{-1} \| \| (V_{U_t W_t}^T U_t W_t)^{-1} \| \\
&= \frac{\|W_{t,\perp}^T W_{t+1}\| \|U_t W_{t,\perp}\|}{\sigma_{\min}(W_t^T W_{t+1}) \sigma_{\min}(V_{U_t W_t}^T U_t W_t)} \\
&= \frac{\|W_{t,\perp}^T W_{t+1}\| \|U_t W_{t,\perp}\|}{\sigma_{\min}(W_t^T W_{t+1}) \sigma_{\min}(U_t W_t)} \\
&\leq 4 \|W_{t,\perp}^T W_{t+1}\|.
\end{aligned}$$

Hence, we can conclude that

$$\begin{aligned}
&\|V_{X^\perp}^T \left(B_3 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_3 + B_3^T) \right) V_{U_t W_t}\| \\
&\leq 2 \|B_3\| \\
&\leq 8 \|W_{t,\perp}^T W_{t+1}\| \\
&\stackrel{(a)}{\leq} \mu \left(\frac{1}{800} \sigma_{\min}^2(X) + 8 \|U_t W_t\| \|U_t W_{t,\perp}\| \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 32\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(X X^T - U_t U_t^T)\| \\
&\stackrel{(b)}{\leq} \frac{1}{400} \mu \cdot \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_t W_t}\| + 32\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(X X^T - U_t U_t^T)\|.
\end{aligned}$$

Inequality (a) follows from Lemma B.3. In (b) we used the assumption $\|U_t W_{t,\perp}\| \leq c\kappa^{-2} \sigma_{\min}(X)$ and $\|U_t W_t\| \leq 3\|X\|$.

Bounding (IV): We start by noticing that

$$\begin{aligned}
\mu \|\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)\| &\leq \mu (\|X X^T - U_t U_t^T\| + \|(\text{Id} - \mathcal{A}^* \mathcal{A})(X X^T - U_t U_t^T)\|) \\
&\leq 11\mu \|X\|^2 \\
&\leq 11c\kappa^{-2},
\end{aligned}$$

where we have used the assumption $\|U\| \leq 3\|X\|$, (43) and (45). Hence, we obtain that

$$\begin{aligned}
&\|V_{X^\perp}^T \left(B_4 - \frac{1}{2} V_{U_t W_t} V_{U_t W_t}^T (B_4 + B_4^T) \right) V_{U_t W_t}\| \\
&\leq 2 \|B_4\| \\
&= 2\mu \|MP\| \\
&\leq 2\mu \|\mathcal{A}^* \mathcal{A}(X X^T - U_t U_t^T)\| \|B_3\| \\
&\leq 22c\kappa^{-2} \|B_3\| \\
&\stackrel{(a)}{\leq} \frac{\mu}{50} \cdot \kappa^{-2} \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_t W_t}\| + 352c\mu \|[(\text{Id} - \mathcal{A}^* \mathcal{A})(X X^T - U_t U_t^T)]\|.
\end{aligned}$$

Inequality (a) follows from similar arguments as when we were bounding (III) and by choosing the constant $c > 0$ small enough.

Bounding (V): First, we want to estimate $\|B\|$. We note that

$$\begin{aligned}
\|B\| &\stackrel{(a)}{\leq} \mu\|M\| + \|P\| + \mu\|MP\| \\
&\stackrel{(b)}{\leq} \mu\|M\| + \|B_3\| + \mu\|M\|\|B_3\| \\
&\stackrel{(c)}{=} \mu\|\mathcal{A}^*\mathcal{A}(XX^T - U_tU_t^T)\| + (1 + \mu\|\mathcal{A}^*\mathcal{A}(XX^T - U_tU_t^T)\|)\|B_3\| \\
&\stackrel{(d)}{\leq} \mu\|\mathcal{A}^*\mathcal{A}(XX^T - U_tU_t^T)\| + 2\|B_3\| \\
&\stackrel{(e)}{\leq} \mu\|\mathcal{A}^*\mathcal{A}(XX^T - U_tU_t^T)\| + \frac{1}{400}\mu \cdot \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\| + 32\mu\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\| \\
&\leq \mu\|XX^T - U_tU_t^T\| + \frac{1}{400}\mu \cdot \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\| + 33\mu\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\|.
\end{aligned}$$

In (a) we used the triangle inequality and in (b) we used $B_3 = P$ and the submultiplicativity of the spectral norm. To obtain equality (c) we inserted the definition of M . For (d) we used that $\mu\|\mathcal{A}^*\mathcal{A}(XX^T - U_tU_t^T)\| \leq 2$, which follows from assumption (43) and (45). For inequality (e) we used our bound for $\|B_3\|$, which we have derived when bounding (III). Hence, we have shown that

$$\|B\| \leq \mu\|XX^T - U_tU_t^T\| + \frac{1}{400}\mu \cdot \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\| + 33\mu\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\|. \quad (99)$$

We obtain that

$$\begin{aligned}
&\frac{1}{2}\|V_{X^\perp}^T B V_{U_tW_t} V_{U_tW_t}^T (B + B^T) V_{U_tW_t}\| \\
&\leq \frac{1}{2}\|B\|\|B + B^T\| \\
&\leq \|B\|^2 \\
&\stackrel{(a)}{\leq} 3\mu^2 \left(\|XX^T - U_tU_t^T\|^2 + \frac{1}{400^2}\sigma_{\min}^4(X) \|V_{X^\perp}^T V_{U_tW_t}\|^2 + 33^2\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\|^2 \right) \\
&\stackrel{(b)}{\leq} 3\mu^2\|XX^T - U_tU_t^T\|^2 + \frac{3}{2}c\frac{\mu\kappa^{-4}}{400^2}\sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\|^2 + \frac{3}{2}\mu^2\sigma_{\min}^2(X) \|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\|^2 \\
&\stackrel{(c)}{\leq} 3\mu^2\|XX^T - U_tU_t^T\|^2 + 3c\frac{\mu\kappa^{-4}}{2 \cdot 400^2}\sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\|^2 + \frac{3}{2}c\mu\kappa^{-4}\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\|,
\end{aligned} \quad (100)$$

where in (a) we used inequality (99) combined with Jensen's inequality. For inequality (b) we used (43) and (45). Inequality (c) follows again from (45).

Bounding (VI): We are first going to show that $\|B\| \leq 1$. Indeed, we have that

$$\begin{aligned}
\|B\| &\stackrel{(a)}{\leq} 10\mu\|X\|^2 + \frac{1}{400}\mu \cdot \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_tW_t}\| + 33\mu\|(\text{Id} - \mathcal{A}^*\mathcal{A})(XX^T - U_tU_t^T)\| \\
&\stackrel{(b)}{\leq} 10\mu\|X\|^2 + \frac{1}{400}\mu \cdot \sigma_{\min}^2(X) + 33c\mu\sigma_{\min}^2(X) \\
&\stackrel{(c)}{\leq} 1,
\end{aligned}$$

where in (a) we used inequality (99) and the assumption $\|U\| \leq 3\|X\|$. For inequality (b) we used (43) and for inequality (c) we used assumption (45). We note that from inequality (98) it follows that

$$\|D\| \leq (1 + \|B\|) \left(\frac{1}{2} \|B\|^2 + \|C\| \right) \leq 2(\|B\|^2 + \|C\|). \quad (101)$$

In the first inequality we used the triangle inequality and in the second inequality we used $\|B\| \leq 1$. Note that from (97) and again $\|B\| \leq 1$ it follows that

$$\|C\| \leq 3(2\|B\| + \|B\|^2)^2 \leq 27\|B\|^2. \quad (102)$$

Hence, we obtain that

$$\begin{aligned} & \|D\| \\ & \stackrel{(a)}{\leq} 56\|B\|^2 \\ & \stackrel{(b)}{\leq} 56 \left(3\mu^2 \|XX^T - U_t U_t^T\|^2 + 3c \frac{\mu\kappa^{-4}}{2 \cdot 400^2} \sigma_{\min}^2(X) \|V_{X^\perp}^T V_{U_t W_t}\|^2 + \frac{3}{2} c \mu \kappa^{-4} \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \right). \end{aligned}$$

Inequality (a) is due to the inequalities (101) and (102) and inequality (b) follows from (100).

Combining the estimates: By combining our results we obtain that for small enough $c > 0$ we have that

$$\begin{aligned} & \|V_{X^\perp}^T V_{U_{t+1} W_{t+1}}\| \\ & \leq \|(I)\| + \|(II)\| + \|(III)\| + \|(IV)\| + \|(V)\| + \|(VI)\| \\ & \leq \left(1 - \frac{\mu}{4} \sigma_{\min}^2(X) \right) \|V_{X^\perp}^T V_{U_t W_t}\| + 100\mu \|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| + 500\mu^2 \|XX^T - U_t U_t^T\|^2. \end{aligned}$$

This finishes the proof. \square

B.4 Proof of Lemma 9.4

Proof of Lemma 9.4. We observe that

$$\begin{aligned} U_{t+1} &= U_t + \mu(XX^T - U_t U_t^T)U_t + \mu[(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)]U_t \\ &= (\text{Id} - \mu U_t U_t^T)U_t + \mu XX^T U_t + \mu[(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)]U_t. \end{aligned}$$

Note that $\|(\text{Id} - \mu U_t U_t^T)U_t\| = (1 - \mu\|U_t\|^2)\|U_t\|$ due to $\mu \leq \frac{1}{27}\|X\|^{-2} \leq \frac{1}{3}\|U_t\|^2$. Hence, by the triangle inequality and submultiplicativity of the spectral norm we obtain that

$$\|U_{t+1}\| \leq (1 - \mu\|U_t\|^2 + \mu\|X\|^2 + \mu\|(\mathcal{A}^* \mathcal{A} - \text{Id})(XX^T - U_t U_t^T)\|)\|U_t\|.$$

Hence, by our assumption on $\|(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\|$ we obtain that

$$\|U_{t+1}\| \leq (1 - \mu\|U_t\|^2 + 2\mu\|X\|^2)\|U_t\|. \quad (103)$$

Now assume that $2\|X\| \leq \|U_t\| \leq 3\|X\|$. Then it follows from the last inequality that $\|U_{t+1}\| \leq \|U_t\|$, which due to the assumption $\|U_t\| \leq 3\|X\|$ implies the claim $\|U_{t+1}\| \leq 3\|X\|$. However, if $\|U_t\| \leq 2\|X\|$ holds, then by combining inequality (103) with the assumption $\mu \leq \frac{\|X\|^{-2}}{27}$ we obtain that $\|U_{t+1}\| \leq 3\|X\|$ as well, which finishes the proof. \square

B.5 Proof of Lemma 9.5

Lemma B.4. *Under the assumptions of Lemma 9.5 it holds that*

$$\|V_{X^\perp}^T U_t U_t^T\| \leq 3\|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|$$

as well as

$$\|X X^T - U_t U_t^T\| \leq 4\|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|.$$

Proof. We notice that by the triangle inequality and submultiplicativity it holds that

$$\begin{aligned} \|V_{X^\perp}^T U_t U_t^T\| &\leq \|V_{X^\perp}^T U_t U_t^T V_X\| + \|V_{X^\perp}^T U_t U_t^T V_{X^\perp}\| \\ &= \|V_{X^\perp}^T (X X^T - U_t U_t^T) V_X\| + \|V_{X^\perp}^T U_t U_t^T V_{X^\perp}\| \\ &= \|V_X^T (X X^T - U_t U_t^T)\| + \|V_{X^\perp}^T U_t U_t^T V_{X^\perp}\|. \end{aligned}$$

In order to bound the second term we compute that

$$\begin{aligned} \|V_{X^\perp}^T U_t U_t^T V_{X^\perp}\| &\leq \|V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp}\| + \|V_{X^\perp}^T U_t W_{t,\perp} W_{t,\perp}^T U_t^T V_{X^\perp}\| \\ &= \|V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp}\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|. \end{aligned}$$

In order to bound the first term we note that

$$\begin{aligned} \|V_{X^\perp}^T U_t W_t W_t^T U_t^T V_{X^\perp}\| &= \|V_{X^\perp}^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T U_t^T V_{X^\perp}\| \\ &= \left\| V_{X^\perp}^T V_{U_t W_t} (V_X^T V_{U_t W_t})^{-1} V_X^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T U_t^T V_{X^\perp} \right\| \\ &\leq \|V_{X^\perp}^T V_{U_t W_t}\| (V_X^T V_{U_t W_t})^{-1} \|V_X^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t W_t^T U_t^T\| \\ &= \frac{\|V_{X^\perp}^T V_{U_t W_t}\|}{\sigma_{\min}(V_X^T V_{U_t W_t})} \|V_X^T U_t U_t^T V_{X^\perp}\| \\ &= \frac{\|V_{X^\perp}^T V_{U_t W_t}\|}{\sigma_{\min}(V_X^T V_{U_t W_t})} \|V_X^T (X X^T - U_t U_t^T) V_{X^\perp}\| \\ &\leq \frac{\|V_{X^\perp}^T V_{U_t W_t}\|}{\sigma_{\min}(V_X^T V_{U_t W_t})} \|V_X^T (X X^T - U_t U_t^T)\| \\ &\leq 2\|V_X^T (X X^T - U_t U_t^T)\|. \end{aligned}$$

Hence we can conclude that

$$\|V_{X^\perp}^T U_t U_t^T\| \leq 3\|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|,$$

which shows the first inequality in the lemma. In order to prove the second inequality, we note that by the triangle inequality and submultiplicativity it holds that

$$\begin{aligned} \|X X^T - U_t U_t^T\| &\leq \|V_X^T (X X^T - U_t U_t^T)\| + \|V_{X^\perp}^T U_t U_t^T\| \\ &\leq 4\|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|, \end{aligned}$$

where in the last line we used the previous inequality. This finishes the proof. \square

After having provided the necessary ingredients, we are in a position to prove Lemma 9.5.

Proof of Lemma 9.5. Recall that

$$U_{t+1} = U_t + \mu [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t.$$

Next, we compute that

$$\begin{aligned} XX^T - U_{t+1} U_{t+1}^T &= XX^T - U_t U_t^T - \mu [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T - \mu U_t U_t^T [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \\ &\quad - \mu^2 [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \\ &= XX^T - U_t U_t^T - \mu (XX^T - U_t U_t^T) U_t U_t^T - \mu U_t U_t^T (XX^T - U_t U_t^T) \\ &\quad + \mu [(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T + \mu U_t U_t^T [(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \\ &\quad - \mu^2 [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] \\ &= \underbrace{(\text{Id} - \mu U_t U_t^T)(XX^T - U_t U_t^T)(\text{Id} - \mu U_t U_t^T)}_{=(I)} + \underbrace{\mu [(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T}_{=(II)} \\ &\quad + \underbrace{\mu U_t U_t^T [(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]}_{=(III)} - \underbrace{\mu^2 U_t U_t^T (XX^T - U_t U_t^T) U_t U_t^T}_{=(IV)} \\ &\quad - \underbrace{\mu^2 [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)] U_t U_t^T [(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]}_{=(V)}. \end{aligned}$$

We are going to deal with each summand individually.

Estimation of (I): We note that

$$\begin{aligned} &V_X^T (\text{Id} - \mu U_t U_t^T) (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T) \\ &= V_X^T (\text{Id} - \mu U_t U_t^T) V_X V_X^T (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T) \\ &\quad + V_X^T (\text{Id} - \mu U_t U_t^T) V_{X^\perp} V_{X^\perp}^T (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T) \\ &= V_X^T (\text{Id} - \mu U_t U_t^T) V_X V_X^T (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T) + \mu V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T (\text{Id} - \mu U_t U_t^T) \\ &= (\text{Id} - \mu V_X^T U_t U_t^T V_X) V_X^T (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T) + \mu V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T (\text{Id} - \mu U_t U_t^T). \end{aligned}$$

Hence, we obtain that

$$\begin{aligned} &\|(\text{Id} - \mu V_X^T U_t U_t^T V_X) V_X^T (XX^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T)\| \\ &\leq \|(\text{Id} - \mu V_X^T U_t U_t^T V_X)\| \|V_X^T (XX^T - U_t U_t^T)\| \|(\text{Id} - \mu U_t U_t^T)\| \\ &\leq \|(\text{Id} - \mu V_X^T U_t U_t^T V_X)\| \|V_X^T (XX^T - U_t U_t^T)\| \\ &= (1 - \mu \sigma_{\min}(V_X^T U_t U_t^T V_X)) \|V_X^T (XX^T - U_t U_t^T)\| \\ &\leq (1 - \mu \sigma_{\min}^2(V_X^T U_t U_t^T)) \|V_X^T (XX^T - U_t U_t^T)\|. \end{aligned}$$

Next, we note that

$$\begin{aligned}
\sigma_{\min}^2(V_X^T U_t W_t) &= \sigma_{\min}^2(V_X^T V_{U_t W_t} V_{U_t W_t}^T U_t W_t) \\
&\geq \sigma_{\min}^2(V_X^T V_{U_t W_t}) \sigma_{\min}^2(U_t W_t) \\
&\geq \frac{1}{2} \sigma_{\min}^2(U_t W_t) \\
&\geq \frac{1}{20} \sigma_{\min}^2(X),
\end{aligned}$$

where in the last line we used the assumption $\sigma_{\min}^2(U_t W_t) \geq \frac{1}{10} \sigma_{\min}^2(X)$. Hence, we have shown that

$$\|(\text{Id} - \mu V_X^T U_t U_t^T V_X) V_X^T (X X^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T)\| \leq \left(1 - \frac{\mu}{20} \sigma_{\min}^2(X)\right) \|V_X^T (X X^T - U_t U_t^T)\|.$$

Next, we note that

$$\begin{aligned}
&\|V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T (\text{Id} - \mu U_t U_t^T)\| \\
&\leq \|V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T\| \|\text{Id} - \mu U_t U_t^T\| \\
&\leq \|V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T\| \\
&\leq \|V_X^T U_t W_t W_t^T U_t V_{X^\perp} V_{X^\perp}^T U_t U_t^T\| \\
&\leq \|V_X^T U_t W_t\| \|V_{X^\perp}^T U_t W_t\| \|V_{X^\perp}^T U_t U_t^T\| \\
&\leq \|U_t W_t\|^2 \|V_{X^\perp}^T V_{U_t W_t}\| \|V_{X^\perp}^T U_t U_t^T\| \\
&\leq 9 \|X\|^2 \|V_{X^\perp}^T V_{U_t W_t}\| \|V_{X^\perp}^T U_t U_t^T\| \\
&\leq 9 \|X\|^2 \|V_{X^\perp}^T V_{U_t W_t}\| (3 \|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t\|),
\end{aligned}$$

where in the last line we used Lemma B.4. Then, using the assumption $\|V_{X^\perp}^T V_{U_t W_t}\| \leq c\kappa^{-2}$ it follows that

$$\begin{aligned}
&\|V_X^T U_t U_t^T V_{X^\perp} V_{X^\perp}^T U_t U_t^T (\text{Id} - U_t U_t^T)\| \\
&\leq \frac{1}{100} \sigma_{\min}^2(X) \|V_X^T (X X^T - U_t U_t^T)\| + \frac{\sigma_{\min}^2(X)}{400} \|U_t W_{t,\perp} W_{t,\perp}^T U_t\|.
\end{aligned}$$

Hence, we have shown that

$$\begin{aligned}
&\|V_X^T (\text{Id} - \mu U_t U_t^T) (X X^T - U_t U_t^T) (\text{Id} - \mu U_t U_t^T)\| \\
&\leq \left(1 - \frac{\mu}{40} \sigma_{\min}^2(X)\right) \|V_X^T (X X^T - U_t U_t^T)\| + \mu \frac{\sigma_{\min}^2(X)}{400} \|U_t W_{t,\perp} W_{t,\perp}^T U_t\|.
\end{aligned}$$

Estimation of (II): We note that

$$\begin{aligned}
\|V_X^T [(\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)] U_t U_t^T\| &\leq \|(\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)\| \|U_t\|^2 \\
&\lesssim \|(\text{Id} - \mathcal{A}^* \mathcal{A}) (X X^T - U_t U_t^T)\| \|X\|^2 \\
&\lesssim c \sigma_{\min}^2(X) \|X X^T - U_t U_t^T\| \\
&\lesssim c \sigma_{\min}^2(X) (\|V_X^T (X X^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t\|).
\end{aligned}$$

In the second inequality we used the assumption $\|U\| \leq 3\|X\|$ and in the third inequality we used assumption (47). In the fourth inequality we applied Lemma B.4. Hence, by choosing the constant $c > 0$ small enough, we obtain that

$$\|V_X^T[(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - UU_t^T)]U_t U_t^T\| \leq \frac{1}{1000} \sigma_{\min}^2(X) (\|V_X^T(XX^T - UU_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|).$$

Estimation of (III): In the analogous way as in the estimation of (II) we derive that

$$\|V_X^T U_t U_t^T [(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \leq \frac{1}{1000} \sigma_{\min}^2(X) (\|V_X^T(XX^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|).$$

Estimation of (IV): We note that it follows from submultiplicativity of the spectral norm that

$$\begin{aligned} \|V_X^T U_t U_t^T (XX^T - U_t U_t^T) U_t U_t^T\| &\leq \|U\|^4 \|XX^T - U_t U_t^T\| \\ &\lesssim \|X\|^4 \|XX^T - U_t U_t^T\| \\ &\lesssim \|X\|^4 \|V_X^T(XX^T - U_t U_t^T)\| + \|X\|^4 \|U_t W_{t,\perp}\|^2, \end{aligned}$$

where in the second line we used the assumption $\|U_t\| \leq 3\|X\|$. In the third line we used Lemma B.4. Then using the assumption $\mu \leq c\kappa^{-2}\|X\|^{-2}$ it follows that

$$\mu^2 \|V_X^T U_t U_t^T (XX^T - U_t U_t^T) U_t U_t^T\| \leq \frac{\mu}{200} \sigma_{\min}^2(X) \|V_X^T(XX^T - U_t U_t^T)\| + \mu \frac{\sigma_{\min}^2(X)}{1000} \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|.$$

Estimation of (V): We first note that

$$\begin{aligned} \|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| &\leq \|XX^T - U_t U_t^T\| + \|[(\text{Id} - \mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \\ &\leq (1 + c\kappa^{-2}) \|XX^T - U_t U_t^T\| \\ &\leq 2\|XX^T - U_t U_t^T\|, \end{aligned}$$

where we have used Assumption (47). In a similar manner, again using Assumption (47), we can show that

$$\|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \leq 2\|XX^T - U_t U_t^T\| \leq 2(\|X\|^2 + \|U_t\|^2).$$

Hence, it follows that

$$\begin{aligned} &\|V_X^T[(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]U_t U_t^T[(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \\ &\leq \|[(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \|U_t\|^2 \|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \\ &\lesssim \|XX^T - U_t U_t^T\| \|U_t\|^2 \|(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)\| \\ &\lesssim \|XX^T - U_t U_t^T\| \|U_t\|^2 (\|X\|^2 + \|U_t\|^2) \\ &\lesssim \|XX^T - U_t U_t^T\| \|X\|^4 \\ &\lesssim (\|V_X^T(XX^T - U_t U_t^T)\| + \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|) \|X\|^4, \end{aligned}$$

where in the third and fourth line we used the estimates from above. The fifth line is due to assumption $\|U_t\| \leq 3\|X\|$. In the last line we used Lemma B.4. Hence, it follows that

$$\begin{aligned} &\mu^2 \|V_X^T[(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]U_t U_t^T[(\mathcal{A}^* \mathcal{A})(XX^T - U_t U_t^T)]\| \\ &\leq \frac{\mu}{1000} \sigma_{\min}^2(X) \|V_X^T(XX^T - U_t U_t^T)\| + \frac{\mu}{400} \sigma_{\min}^2(X) \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|, \end{aligned}$$

where the last inequality is due to the assumption $\mu \leq c\kappa^{-2}\|X\|^{-2}$ for a sufficiently small constant $c > 0$.

Combining the estimates: By combining the estimates, it follows that

$$\|V_X^T (XX^T - U_{t+1}U_{t+1}^T)\| \leq \left(1 - \frac{\mu}{200}\sigma_{\min}^2(X)\right) \|V_X^T (XX^T - U_tU_t^T)\| + \mu \frac{\sigma_{\min}^2(X)}{100} \|U_t W_{t,\perp} W_{t,\perp}^T U_t^T\|,$$

which finishes the proof. \square