

Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent

Tian Tong

*Department of Electrical and Computer Engineering
Carnegie Mellon University
Pittsburgh, PA 15213, USA*

TTONG1@ANDREW.CMU.EDU

Cong Ma

*Department of Electrical Engineering and Computer Sciences
University of California, Berkeley
Berkeley CA 94720, USA*

CONGM@BERKELEY.EDU

Yuejie Chi

*Department of Electrical and Computer Engineering
Carnegie Mellon University
Pittsburgh, PA 15213, USA*

YUEJIECHI@CMU.EDU

Editor: Qiang Liu

Abstract

Low-rank matrix estimation is a canonical problem that finds numerous applications in signal processing, machine learning and imaging science. A popular approach in practice is to factorize the matrix into two compact low-rank factors, and then optimize these factors directly via simple iterative methods such as gradient descent and alternating minimization. Despite nonconvexity, recent literatures have shown that these simple heuristics in fact achieve linear convergence when initialized properly for a growing number of problems of interest. However, upon closer examination, existing approaches can still be computationally expensive especially for ill-conditioned matrices: the convergence rate of gradient descent depends linearly on the condition number of the low-rank matrix, while the per-iteration cost of alternating minimization is often prohibitive for large matrices.

The goal of this paper is to set forth a competitive algorithmic approach dubbed *Scaled Gradient Descent* (**ScaledGD**) which can be viewed as preconditioned or diagonally-scaled gradient descent, where the preconditioners are adaptive and iteration-varying with a minimal computational overhead. With tailored variants for low-rank matrix sensing, robust principal component analysis and matrix completion, we theoretically show that **ScaledGD** achieves the best of both worlds: it converges linearly at a rate independent of the condition number of the low-rank matrix similar as alternating minimization, while maintaining the low per-iteration cost of gradient descent. Our analysis is also applicable to general loss functions that are restricted strongly convex and smooth over low-rank matrices. To the best of our knowledge, **ScaledGD** is the first algorithm that provably has such properties over a wide range of low-rank matrix estimation tasks. At the core of our analysis is the introduction of a new distance function that takes account of the preconditioners when measuring the distance between the iterates and the ground truth. Finally, numerical examples are provided to demonstrate the effectiveness of **ScaledGD** in accelerating the convergence rate of ill-conditioned low-rank matrix estimation in a wide number of applications.

Keywords: low-rank matrix factorization, scaled gradient descent, ill-conditioned matrix recovery, matrix sensing, robust PCA, matrix completion, general losses

1. Introduction

Low-rank matrix estimation plays a critical role in fields such as machine learning, signal processing, imaging science, and many others. Broadly speaking, one aims to recover a rank- r matrix $\mathbf{X}_\star \in \mathbb{R}^{n_1 \times n_2}$ from a set of observations $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star)$, where the operator $\mathcal{A}(\cdot)$ models the measurement process. It is natural to minimize the least-squares loss function subject to a rank constraint:

$$\underset{\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}}{\text{minimize}} \quad f(\mathbf{X}) := \frac{1}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2 \quad \text{s.t.} \quad \text{rank}(\mathbf{X}) \leq r, \quad (1)$$

which is, however, computationally intractable in general due to the rank constraint. Moreover, as the size of the matrix increases, the costs involved in optimizing over the full matrix space (i.e. $\mathbb{R}^{n_1 \times n_2}$) are prohibitive in terms of both memory and computation. To cope with these challenges, one popular approach is to parametrize $\mathbf{X} = \mathbf{L}\mathbf{R}^\top$ by two low-rank factors $\mathbf{L} \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{R} \in \mathbb{R}^{n_2 \times r}$ that are more memory-efficient, and then to optimize over the factors instead:

$$\underset{\mathbf{L} \in \mathbb{R}^{n_1 \times r}, \mathbf{R} \in \mathbb{R}^{n_2 \times r}}{\text{minimize}} \quad \mathcal{L}(\mathbf{L}, \mathbf{R}) := f(\mathbf{L}\mathbf{R}^\top). \quad (2)$$

Although this leads to a nonconvex optimization problem over the factors, recent breakthroughs have shown that simple algorithms (e.g. gradient descent, alternating minimization), when properly initialized (e.g. via the spectral method), can provably converge to the true low-rank factors under mild statistical assumptions. These benign convergence guarantees hold for a growing number of problems such as low-rank matrix sensing, matrix completion, robust principal component analysis (robust PCA), phase synchronization, and so on.

However, upon closer examination, existing approaches such as gradient descent and alternating minimization are still computationally expensive, especially for ill-conditioned matrices. Take low-rank matrix sensing as an example: although the per-iteration cost is small, the iteration complexity of gradient descent scales linearly with respect to the condition number of the low-rank matrix \mathbf{X}_\star Tu et al. (2016); on the other end, while the iteration complexity of alternating minimization Jain et al. (2013) is independent of the condition number, each iteration requires inverting a linear system whose size is proportional to the dimension of the matrix and thus the per-iteration cost is prohibitive for large-scale problems. These together raise an important open question: *can one design an algorithm with a comparable per-iteration cost as gradient descent, but converges much faster at a rate that is independent of the condition number as alternating minimization in a provable manner for a wide variety of low-rank matrix estimation tasks?*

1.1 Preconditioning helps: scaled gradient descent

In this paper, we answer this question affirmatively by studying the following scaled gradient descent (**ScaledGD**) algorithm to optimize (2). Given an initialization $(\mathbf{L}_0, \mathbf{R}_0)$, **ScaledGD** proceeds as follows

$$\begin{aligned} \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta \nabla_{\mathbf{L}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t) (\mathbf{R}_t^\top \mathbf{R}_t)^{-1}, \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t) (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}, \end{aligned} \quad (3)$$

where $\eta > 0$ is the step size and $\nabla_{\mathbf{L}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t)$ (resp. $\nabla_{\mathbf{R}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t)$) is the gradient of the loss function \mathcal{L} with respect to the factor \mathbf{L}_t (resp. \mathbf{R}_t) at the t -th iteration. Comparing to vanilla gradient descent, the search directions of the low-rank factors $\mathbf{L}_t, \mathbf{R}_t$ in (3) are *scaled* by $(\mathbf{R}_t^\top \mathbf{R}_t)^{-1}$ and $(\mathbf{L}_t^\top \mathbf{L}_t)^{-1}$ respectively. Intuitively, the scaling serves as a preconditioner as in quasi-Newton type algorithms, with the hope of improving the quality of the search direction to allow larger step sizes. Since the computation of the Hessian is extremely expensive, it is necessary to design preconditioners that are both theoretically sound and practically cheap to compute. Such requirements are met by **ScaledGD**, where the preconditioners are computed by inverting two $r \times r$ matrices, whose size is much smaller than the dimension of matrix factors. Therefore, each iteration of **ScaledGD** adds

minimal overhead to the gradient computation and has the order-wise same per-iteration cost as gradient descent. Moreover, the preconditioners are adaptive and iteration-varying. Another key property of `ScaledGD` is that it ensures the iterates are covariant with respect to the parameterization of low-rank factors up to invertible transforms.

While `ScaledGD` and its alternating variants have been proposed in Mishra et al. (2012); Mishra and Sepulchre (2016); Tanner and Wei (2016) for a subset of the problems we studied, none of these prior art provides any theoretical validations to the empirical success. In this work, we confirm *theoretically* that `ScaledGD` achieves linear convergence at a rate *independent of* the condition number of the matrix when initialized properly, e.g. using the standard spectral method, for several canonical problems: low-rank matrix sensing, robust PCA, and matrix completion. Table 1 summarizes the performance guarantees of `ScaledGD` in terms of both statistical and computational complexities with comparisons to prior algorithms using the vanilla gradient method.

- *Low-rank matrix sensing.* As long as the measurement operator satisfies the standard restricted isometry property (RIP) with an RIP constant $\delta_{2r} \lesssim 1/(\sqrt{r}\kappa)$, where κ is the condition number of \mathbf{X}_* , `ScaledGD` reaches ϵ -accuracy in $O(\log(1/\epsilon))$ iterations when initialized by the spectral method. This strictly improves the iteration complexity $O(\kappa \log(1/\epsilon))$ of gradient descent in Tu et al. (2016) under the same sample complexity requirement.
- *Robust PCA.* Under the deterministic corruption model Chandrasekaran et al. (2011), as long as the fraction α of corruptions per row / column satisfies $\alpha \lesssim 1/(\mu r^{3/2} \kappa)$, where μ is the incoherence parameter of \mathbf{X}_* , `ScaledGD` in conjunction with hard thresholding reaches ϵ -accuracy in $O(\log(1/\epsilon))$ iterations when initialized by the spectral method. This strictly improves the iteration complexity of projected gradient descent Yi et al. (2016).
- *Matrix completion.* Under the random Bernoulli observation model, as long as the sample complexity satisfies $n_1 n_2 p \gtrsim (\mu \kappa^2 \vee \log n) \mu n r^2 \kappa^2$ with $n = n_1 \vee n_2$, `ScaledGD` in conjunction with a properly designed projection operator reaches ϵ -accuracy in $O(\log(1/\epsilon))$ iterations when initialized by the spectral method. This improves the iteration complexity of projected gradient descent Zheng and Lafferty (2016) at the expense of requiring a larger sample size.

In addition, `ScaledGD` does not require any explicit regularizations that balance the norms of two low-rank factors as required in Tu et al. (2016); Yi et al. (2016); Zheng and Lafferty (2016), and removed the additional projection that maintains the incoherence properties in robust PCA Yi et al. (2016), thus unveiling the implicit regularization property of `ScaledGD`. To the best of our knowledge, this is the first factored gradient descent algorithm that achieves a fast convergence rate that is independent of the condition number of the low-rank matrix at near-optimal sample complexities without increasing the per-iteration computational cost. Our analysis is also applicable to general loss functions that are restricted strongly convex and smooth over low-rank matrices.

At the core of our analysis, we introduce a new distance metric (i.e. Lyapunov function) that accounts for the preconditioners, and carefully show the contraction of the `ScaledGD` iterates under the new distance metric. We expect that the `ScaledGD` algorithm can accelerate the convergence for other low-rank matrix estimation problems, as well as facilitate the design and analysis of other quasi-Newton first-order algorithms. As a teaser, Figure 1 illustrates the relative error of completing a 1000×1000 incoherent matrix of rank 10 with varying condition numbers from 20% of its entries, using either `ScaledGD` or vanilla GD with spectral initialization. Even for moderately ill-conditioned matrices, the convergence rate of vanilla GD slows down dramatically, while it is evident that `ScaledGD` converges at a rate independent of the condition number and therefore is much more efficient.

Remark 1 (ScaledGD for PSD matrices) *When the low-rank matrix of interest is positive semi-definite (PSD), we factorize the matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$ as $\mathbf{X} = \mathbf{L}\mathbf{L}^\top$, with $\mathbf{L} \in \mathbb{R}^{n \times r}$. The update rule*

	Matrix sensing		Robust PCA		Matrix completion	
Algorithms	sample complexity	iteration complexity	corruption fraction	iteration complexity	sample complexity	iteration complexity
GD	$nr^2\kappa^2$	$\kappa \log \frac{1}{\epsilon}$	$\frac{1}{\mu r^{3/2} \kappa^{3/2} \sqrt{\mu r \kappa^2}}$	$\kappa \log \frac{1}{\epsilon}$	$(\mu \vee \log n) \mu n r^2 \kappa^2$	$\kappa \log \frac{1}{\epsilon}$
ScaledGD (this paper)	$nr^2\kappa^2$	$\log \frac{1}{\epsilon}$	$\frac{1}{\mu r^{3/2} \kappa}$	$\log \frac{1}{\epsilon}$	$(\mu \kappa^2 \vee \log n) \mu n r^2 \kappa^2$	$\log \frac{1}{\epsilon}$

Table 1: Comparisons of ScaledGD with gradient descent (GD) when tailored to various problems (with spectral initialization) Tu et al. (2016); Yi et al. (2016); Zheng and Lafferty (2016), where they have comparable per-iteration costs. Here, we say that the output \mathbf{X} of an algorithm reaches ϵ -accuracy, if it satisfies $\|\mathbf{X} - \mathbf{X}_*\|_F \leq \epsilon \sigma_r(\mathbf{X}_*)$. Here, $n := n_1 \vee n_2 = \max\{n_1, n_2\}$, κ and μ are the condition number and incoherence parameter of \mathbf{X}_* .

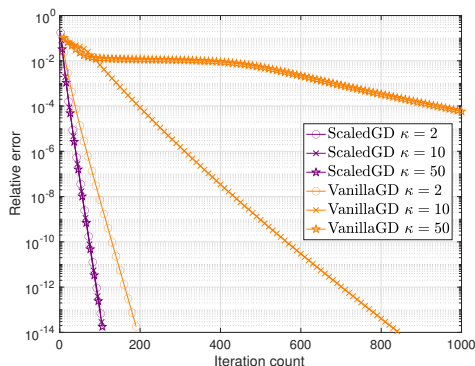


Figure 1: Performance of ScaledGD and vanilla GD for completing a 1000×1000 incoherent matrix of rank 10 with different condition numbers $\kappa = 2, 10, 50$, where each entry is observed independently with probability 0.2. Here, both methods are initialized via the spectral method. It can be seen that ScaledGD converges much faster than vanilla GD even for moderately large condition numbers.

of ScaledGD simplifies to

$$\mathbf{L}_{t+1} = \mathbf{L}_t - \eta \nabla_{\mathbf{L}} \mathcal{L}(\mathbf{L}_t) (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}. \quad (4)$$

We focus on the asymmetric case since the analysis is more involved with two factors. Our theory applies to the PSD case without loss of generality.

1.2 Related work

Our work contributes to the growing literature of design and analysis of provable nonconvex optimization procedures for high-dimensional signal estimation; see e.g. Jain and Kar (2017); Chen and Chi (2018); Chi et al. (2019) for recent overviews. A growing number of problems have been demonstrated to possess benign geometry that is amenable for optimization Mei et al. (2018) either globally or locally under appropriate statistical models. On one end, it is shown that there are no spurious local minima in the optimization landscape of matrix sensing and completion Ge et al. (2016); Bhojanapalli et al. (2016b); Park et al. (2017); Ge et al. (2017), phase retrieval Sun et al. (2018); Davis et al. (2017), dictionary learning Sun et al. (2015), kernel PCA Chen and Li (2019)

and linear neural networks Baldi and Hornik (1989); Kawaguchi (2016). Such landscape analysis facilitates the adoption of generic saddle-point escaping algorithms Nesterov and Polyak (2006); Ge et al. (2015); Jin et al. (2017) to ensure global convergence. However, the resulting iteration complexity is typically high. On the other end, local refinements with carefully-designed initializations often admit fast convergence, for example in phase retrieval Candès et al. (2015); Ma et al. (2019), matrix sensing Jain et al. (2013); Zheng and Lafferty (2015); Wei et al. (2016), matrix completion Sun and Luo (2016); Chen and Wainwright (2015); Ma et al. (2019); Chen et al. (2020a); Zheng and Lafferty (2016); Chen et al. (2020b), blind deconvolution Li et al. (2019); Ma et al. (2019), and robust PCA Netrapalli et al. (2014); Yi et al. (2016); Chen et al. (2020c), to name a few.

Existing approaches for asymmetric low-rank matrix estimation often require additional regularization terms to balance the two factors, either in the form of $\frac{1}{2}\|\mathbf{L}^\top\mathbf{L}-\mathbf{R}^\top\mathbf{R}\|_F^2$ Tu et al. (2016); Park et al. (2017) or $\frac{1}{2}\|\mathbf{L}\|_F^2 + \frac{1}{2}\|\mathbf{R}\|_F^2$ Zhu et al. (2018); Chen et al. (2020b,c), which ease the theoretical analysis but are often unnecessary for the practical success, as long as the initialization is balanced. Some recent work studies the unregularized gradient descent for low-rank matrix factorization and sensing including Charisopoulos et al. (2021); Du et al. (2018); Ma et al. (2021). However, the iteration complexity of all these approaches scales at least linearly with respect to the condition number κ of the low-rank matrix, e.g. $O(\kappa \log(1/\epsilon))$, to reach ϵ -accuracy, therefore they converge slowly when the underlying matrix becomes ill-conditioned. In contrast, **ScaLedGD** enjoys a local convergence rate of $O(\log(1/\epsilon))$, therefore incurring a much smaller computational footprint when κ is large. Last but not least, alternating minimization Jain et al. (2013); Hardt and Wootters (2014) (which alternatively updates \mathbf{L}_t and \mathbf{R}_t) or singular value projection Netrapalli et al. (2014); Jain et al. (2010) (which operates in the matrix space) also converge at the rate $O(\log(1/\epsilon))$, but the per-iteration cost is much higher than **ScaLedGD**. Another notable algorithm is the Riemannian gradient descent algorithm in Wei et al. (2016), which also converges at the rate $O(\log(1/\epsilon))$ under the same sample complexity for low-rank matrix sensing, but requires a higher memory complexity since it operates in the matrix space rather than the factor space.

From an algorithmic perspective, our approach is closely related to the alternating steepest descent (ASD) method in Tanner and Wei (2016) for low-rank matrix completion, which performs the proposed updates (3) for the low-rank factors in an alternating manner. Furthermore, the scaled gradient updates were also introduced in Mishra et al. (2012); Mishra and Sepulchre (2016) for low-rank matrix completion from the perspective of Riemannian optimization. However, none of Tanner and Wei (2016); Mishra et al. (2012); Mishra and Sepulchre (2016) offered any statistical nor computational guarantees for global convergence. Our analysis of **ScaLedGD** can be viewed as providing justifications to these precursors. Moreover, we have systematically extended the framework of **ScaLedGD** to work in a large number of low-rank matrix estimation tasks such as robust PCA.

1.3 Paper organization and notation

The rest of this paper is organized as follows. Section 2 describes the proposed **ScaLedGD** method and details its application to low-rank matrix sensing, robust PCA and matrix completion with theoretical guarantees in terms of both statistical and computational complexities, highlighting the role of a new distance metric. The convergence guarantee of **ScaLedGD** under the general loss function is also presented. In Section 3, we outline the proof for our main results. Section 4 illustrates the excellent empirical performance of **ScaLedGD** in a variety of low-rank matrix estimation problems. Finally, we conclude in Section 5.

Before continuing, we introduce several notation used throughout the paper. First of all, we use boldfaced symbols for vectors and matrices. For a vector \mathbf{v} , we use $\|\mathbf{v}\|_0$ to denote its ℓ_0 counting norm, and $\|\mathbf{v}\|_2$ to denote the ℓ_2 norm. For any matrix \mathbf{A} , we use $\sigma_i(\mathbf{A})$ to denote its i -th largest singular value, and let $\mathbf{A}_{i,\cdot}$ and $\mathbf{A}_{\cdot,j}$ denote its i -th row and j -th column, respectively. In addition, $\|\mathbf{A}\|_{\text{op}}$, $\|\mathbf{A}\|_F$, $\|\mathbf{A}\|_{1,\infty}$, $\|\mathbf{A}\|_{2,\infty}$, and $\|\mathbf{A}\|_\infty$ stand for the spectral norm (i.e. the largest singular value), the Frobenius norm, the $\ell_{1,\infty}$ norm (i.e. the largest ℓ_1 norm of the rows), the $\ell_{2,\infty}$ norm

(i.e. the largest ℓ_2 norm of the rows), and the entrywise ℓ_∞ norm (the largest magnitude of all entries) of a matrix \mathbf{A} . We denote

$$\mathcal{P}_r(\mathbf{A}) = \min_{\tilde{\mathbf{A}}: \text{rank}(\tilde{\mathbf{A}}) \leq r} \|\mathbf{A} - \tilde{\mathbf{A}}\|_{\text{F}}^2 \quad (5)$$

as the rank- r approximation of \mathbf{A} , which is given by the top- r SVD of \mathbf{A} by the Eckart-Young-Mirsky theorem. We also use $\text{vec}(\mathbf{A})$ to denote the vectorization of a matrix \mathbf{A} . For matrices \mathbf{A}, \mathbf{B} of the same size, we use $\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i,j} \mathbf{A}_{i,j} \mathbf{B}_{i,j} = \text{tr}(\mathbf{A}^\top \mathbf{B})$ to denote their inner product. The set of invertible matrices in $\mathbb{R}^{r \times r}$ is denoted by $\text{GL}(r)$. Let $a \vee b = \max\{a, b\}$ and $a \wedge b = \min\{a, b\}$. Throughout, $f(n) \lesssim g(n)$ or $f(n) = O(g(n))$ means $|f(n)|/|g(n)| \leq C$ for some constant $C > 0$ when n is sufficiently large; $f(n) \gtrsim g(n)$ means $|f(n)|/|g(n)| \geq C$ for some constant $C > 0$ when n is sufficiently large. Last but not least, we use the terminology ‘‘with overwhelming probability’’ to denote the event happens with probability at least $1 - c_1 n^{-c_2}$, where $c_1, c_2 > 0$ are some universal constants, whose values may vary from line to line.

2. Scaled Gradient Descent for Low-Rank Matrix Estimation

This section is devoted to introducing **ScaledGD** and establishing its statistical and computational guarantees for various low-rank matrix estimation problems. Before we instantiate tailored versions of **ScaledGD** on concrete low-rank matrix estimation problems, we first pause to provide more insights of the update rule of **ScaledGD**, by connecting it to the quasi-Newton method. Note that the update rule (3) for **ScaledGD** can be equivalently written in a vectorization form as

$$\begin{aligned} \text{vec}(\mathbf{F}_{t+1}) &= \text{vec}(\mathbf{F}_t) - \eta \begin{bmatrix} (\mathbf{R}_t^\top \mathbf{R}_t)^{-1} \otimes \mathbf{I}_{n_1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{L}_t^\top \mathbf{L}_t)^{-1} \otimes \mathbf{I}_{n_2} \end{bmatrix} \text{vec}(\nabla_{\mathbf{F}} \mathcal{L}(\mathbf{F}_t)) \\ &= \text{vec}(\mathbf{F}_t) - \eta \mathbf{H}_t^{-1} \text{vec}(\nabla_{\mathbf{F}} \mathcal{L}(\mathbf{F}_t)), \end{aligned} \quad (6)$$

where we denote $\mathbf{F}_t = [\mathbf{L}_t^\top, \mathbf{R}_t^\top]^\top \in \mathbb{R}^{(n_1+n_2) \times r}$, and by \otimes the Kronecker product. Here, the block diagonal matrix \mathbf{H}_t is set to be

$$\mathbf{H}_t := \begin{bmatrix} (\mathbf{R}_t^\top \mathbf{R}_t) \otimes \mathbf{I}_{n_1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{L}_t^\top \mathbf{L}_t) \otimes \mathbf{I}_{n_2} \end{bmatrix}.$$

The form (6) makes it apparent that **ScaledGD** can be interpreted as a quasi-Newton algorithm, where the inverse of \mathbf{H}_t can be cheaply computed through inverting two rank- r matrices.

2.1 Assumptions and error metric

Denote by $\mathbf{U}_* \mathbf{\Sigma}_* \mathbf{V}_*^\top$ the compact singular value decomposition (SVD) of the rank- r matrix $\mathbf{X}_* \in \mathbb{R}^{n_1 \times n_2}$. Here $\mathbf{U}_* \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{V}_* \in \mathbb{R}^{n_2 \times r}$ are composed of r left and right singular vectors, respectively, and $\mathbf{\Sigma}_* \in \mathbb{R}^{r \times r}$ is a diagonal matrix consisting of r singular values of \mathbf{X}_* organized in a non-increasing order, i.e. $\sigma_1(\mathbf{X}_*) \geq \dots \geq \sigma_r(\mathbf{X}_*) > 0$. Define

$$\kappa := \sigma_1(\mathbf{X}_*) / \sigma_r(\mathbf{X}_*) \quad (7)$$

as the condition number of \mathbf{X}_* . Define the ground truth low-rank factors as

$$\mathbf{L}_* := \mathbf{U}_* \mathbf{\Sigma}_*^{1/2}, \quad \text{and} \quad \mathbf{R}_* := \mathbf{V}_* \mathbf{\Sigma}_*^{1/2}, \quad (8)$$

so that $\mathbf{X}_* = \mathbf{L}_* \mathbf{R}_*^\top$. Correspondingly, denote the stacked factor matrix as

$$\mathbf{F}_* := \begin{bmatrix} \mathbf{L}_* \\ \mathbf{R}_* \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times r}. \quad (9)$$

Next, we are in need of a right metric to measure the performance of the `ScaledGD` iterates $\mathbf{F}_t := [\mathbf{L}_t^\top, \mathbf{R}_t^\top]^\top$. Obviously, the factored representation is not unique in that for any invertible matrix $\mathbf{Q} \in \text{GL}(r)$, one has $\mathbf{L}\mathbf{R}^\top = (\mathbf{L}\mathbf{Q})(\mathbf{R}\mathbf{Q}^{-\top})^\top$. Therefore, the reconstruction error metric needs to take into account this identifiability issue. More importantly, we need a diagonal scaling in the distance error metric to properly account for the effect of preconditioning. To provide intuition, note that the update rule (3) can be viewed as finding the best local quadratic approximation of $\mathcal{L}(\cdot)$ in the following sense:

$$\begin{aligned} (\mathbf{L}_{t+1}, \mathbf{R}_{t+1}) = \underset{\mathbf{L}, \mathbf{R}}{\operatorname{argmin}} & \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t) + \langle \nabla_{\mathbf{L}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t), \mathbf{L} - \mathbf{L}_t \rangle + \langle \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t), \mathbf{R} - \mathbf{R}_t \rangle \\ & + \frac{1}{2\eta} \left(\left\| (\mathbf{L} - \mathbf{L}_t)(\mathbf{R}_t^\top \mathbf{R}_t)^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R} - \mathbf{R}_t)(\mathbf{L}_t^\top \mathbf{L}_t)^{1/2} \right\|_{\mathbb{F}}^2 \right), \end{aligned}$$

where it is different from the common interpretation of gradient descent in the way the quadratic approximation is taken by a scaled norm. When $\mathbf{L}_t \approx \mathbf{L}_\star$ and $\mathbf{R}_t \approx \mathbf{R}_\star$ are approaching the ground truth, the additional scaling factors can be approximated by $\mathbf{L}_t^\top \mathbf{L}_t \approx \Sigma_\star$ and $\mathbf{R}_t^\top \mathbf{R}_t \approx \Sigma_\star$, leading to the following error metric

$$\operatorname{dist}^2(\mathbf{F}, \mathbf{F}_\star) := \inf_{\mathbf{Q} \in \text{GL}(r)} \left\| (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star)\Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star)\Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2. \quad (10)$$

Correspondingly, we define the optimal alignment matrix \mathbf{Q} between \mathbf{F} and \mathbf{F}_\star as

$$\mathbf{Q} := \underset{\mathbf{Q} \in \text{GL}(r)}{\operatorname{argmin}} \left\| (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star)\Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star)\Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2, \quad (11)$$

whenever the minimum is achieved.¹ It turns out that for the `ScaledGD` iterates $\{\mathbf{F}_t\}$, the optimal alignment matrices $\{\mathbf{Q}_t\}$ always exist (at least when properly initialized) and hence are well-defined. The design and analysis of this new distance metric are of crucial importance in obtaining the improved rate of `ScaledGD`; see Appendix A.1 for a collection of its properties. In comparison, the previously studied distance metrics (proposed mainly for GD) either do not include the diagonal scaling Ma et al. (2021); Tu et al. (2016), or only consider the ambiguity class up to orthonormal transforms Tu et al. (2016), which fail to unveil the benefit of `ScaledGD`.

2.2 Matrix sensing

Assume that we have collected a set of linear measurements about a rank- r matrix $\mathbf{X}_\star \in \mathbb{R}^{n_1 \times n_2}$, given as

$$\mathbf{y} = \mathcal{A}(\mathbf{X}_\star) \in \mathbb{R}^m, \quad (12)$$

where $\mathcal{A}(\mathbf{X}) = \{\{\mathbf{A}_k, \mathbf{X}\}\}_{k=1}^m : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}^m$ is the linear map modeling the measurement process. The goal of low-rank matrix sensing is to recover \mathbf{X}_\star from \mathbf{y} , especially when the number of measurements $m \ll n_1 n_2$, by exploiting the low-rank property. This problem has wide applications in medical imaging, signal processing, and data compression Candès and Plan (2011).

Algorithm. Writing $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ into a factored form $\mathbf{X} = \mathbf{L}\mathbf{R}^\top$, we consider the following optimization problem:

$$\underset{\mathbf{F} \in \mathbb{R}^{(n_1+n_2) \times r}}{\operatorname{minimize}} \quad \mathcal{L}(\mathbf{F}) = \frac{1}{2} \left\| \mathcal{A}(\mathbf{L}\mathbf{R}^\top) - \mathbf{y} \right\|_2^2. \quad (13)$$

Here as before, \mathbf{F} denotes the stacked factor matrix $[\mathbf{L}^\top, \mathbf{R}^\top]^\top$. We suggest running `ScaledGD` (3) with the spectral initialization to solve (13), which performs the top- r SVD on $\mathcal{A}^*(\mathbf{y})$, where $\mathcal{A}^*(\cdot)$ is the adjoint operator of $\mathcal{A}(\cdot)$. The full algorithm is stated in Algorithm 1. The low-rank matrix can be estimated as $\mathbf{X}_T = \mathbf{L}_T \mathbf{R}_T^\top$ after running T iterations of `ScaledGD`.

1. If there are multiple minimizers, we can arbitrarily take one to be \mathbf{Q} .

Algorithm 1 ScaledGD for low-rank matrix sensing with spectral initialization

Spectral initialization: Let $\mathbf{U}_0 \mathbf{\Sigma}_0 \mathbf{V}_0^\top$ be the top- r SVD of $\mathcal{A}^*(\mathbf{y})$, and set

$$\mathbf{L}_0 = \mathbf{U}_0 \mathbf{\Sigma}_0^{1/2}, \quad \text{and} \quad \mathbf{R}_0 = \mathbf{V}_0 \mathbf{\Sigma}_0^{1/2}. \quad (14)$$

Scaled gradient updates: for $t = 0, 1, 2, \dots, T - 1$ do

$$\begin{aligned} \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta \mathcal{A}^*(\mathcal{A}(\mathbf{L}_t \mathbf{R}_t^\top) - \mathbf{y}) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1}, \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta \mathcal{A}^*(\mathcal{A}(\mathbf{L}_t \mathbf{R}_t^\top) - \mathbf{y})^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}. \end{aligned} \quad (15)$$

Theoretical guarantees. To understand the performance of ScaledGD for low-rank matrix sensing, we adopt a standard assumption on the sensing operator $\mathcal{A}(\cdot)$, namely the Restricted Isometry Property (RIP).

Definition 2 (RIP Recht et al. (2010)) *The linear map $\mathcal{A}(\cdot)$ is said to obey the rank- r RIP with a constant $\delta_r \in [0, 1)$, if for all matrices $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$ of rank at most r , one has*

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

It is well-known that many measurement ensembles satisfy the RIP property Recht et al. (2010); Candès and Plan (2011). For example, if the entries of \mathbf{A}_i 's are composed of i.i.d. Gaussian entries $\mathcal{N}(0, 1/m)$, then the RIP is satisfied for a constant δ_r as long as m is on the order of $(n_1 + n_2)r/\delta_r^2$. With the RIP condition in place, the following theorem demonstrates that ScaledGD converges linearly — in terms of the new distance metric (cf. (10)) — at a constant rate as long as the sensing operator $\mathcal{A}(\cdot)$ has a sufficiently small RIP constant.

Theorem 3 *Suppose that $\mathcal{A}(\cdot)$ obeys the $2r$ -RIP with $\delta_{2r} \leq 0.02/(\sqrt{r}\kappa)$. If the step size obeys $0 < \eta \leq 2/3$, then for all $t \geq 0$, the iterates of the ScaledGD method in Algorithm 1 satisfy*

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq (1 - 0.6\eta)^t 0.1\sigma_r(\mathbf{X}_\star), \quad \text{and} \quad \|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq (1 - 0.6\eta)^t 0.15\sigma_r(\mathbf{X}_\star).$$

Theorem 3 establishes that the distance $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ contracts linearly at a constant rate, as long as the sample size satisfies $m = O(nr^2\kappa^2)$ with Gaussian random measurements Recht et al. (2010), where we recall that $n = n_1 \vee n_2$. To reach ϵ -accuracy, i.e. $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq \epsilon\sigma_r(\mathbf{X}_\star)$, ScaledGD takes at most $T = O(\log(1/\epsilon))$ iterations, which is *independent* of the condition number κ of \mathbf{X}_\star . In comparison, alternating minimization with spectral initialization (AltMinSense) converges in $O(\log(1/\epsilon))$ iterations as long as $m = O(nr^3\kappa^4)$ Jain et al. (2013), where the per-iteration cost is much higher.² On the other end, gradient descent with spectral initialization in Tu et al. (2016) converges in $O(\kappa \log(1/\epsilon))$ iterations as long as $m = O(nr^2\kappa^2)$. Therefore, ScaledGD converges at a much faster rate than GD at the same sample complexity while requiring a significantly lower per-iteration cost than AltMinSense.

Remark 4 *Tu et al. (2016) suggested that one can employ a more expensive initialization scheme, e.g. performing multiple projected gradient descent steps over the low-rank matrix, to reduce the sample complexity. By seeding ScaledGD with the output of updates of the form $\mathbf{X}_{\tau+1} = \mathcal{P}_r(\mathbf{X}_\tau - \mathcal{A}^*(\mathcal{A}(\mathbf{X}_\tau) - \mathbf{y}))$ after $T_0 \gtrsim \log(\sqrt{r}\kappa)$ iterations, where $\mathcal{P}_r(\cdot)$ is defined in (5), ScaledGD succeeds with the sample size $O(nr)$ which is information theoretically optimal.*

² The exact per-iteration complexity of AltMinSense depends on how the least-squares subproblems are solved with m equations and nr unknowns; see (Luo et al., 2020, Table 1) for detailed comparisons.

2.3 Robust PCA

Assume that we have observed the data matrix

$$\mathbf{Y} = \mathbf{X}_* + \mathbf{S}_*,$$

which is a superposition of a rank- r matrix \mathbf{X}_* , modeling the clean data, and a sparse matrix \mathbf{S}_* , modeling the corruption or outliers. The goal of robust PCA Candès et al. (2011); Chandrasekaran et al. (2011) is to separate the two matrices \mathbf{X}_* and \mathbf{S}_* from their mixture \mathbf{Y} . This problem finds numerous applications in video surveillance, image processing, and so on.

Following Chandrasekaran et al. (2011); Netrapalli et al. (2014); Yi et al. (2016), we consider a deterministic sparsity model for \mathbf{S}_* , in which \mathbf{S}_* contains at most α -fraction of nonzero entries per row and column for some $\alpha \in [0, 1)$, i.e. $\mathbf{S}_* \in \mathcal{S}_\alpha$, where we denote

$$\mathcal{S}_\alpha := \{\mathbf{S} \in \mathbb{R}^{n_1 \times n_2} : \|\mathbf{S}_{i,\cdot}\|_0 \leq \alpha n_2 \text{ for all } i, \text{ and } \|\mathbf{S}_{\cdot,j}\|_0 \leq \alpha n_1 \text{ for all } j\}. \quad (16)$$

Algorithm. Writing $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ into the factored form $\mathbf{X} = \mathbf{L}\mathbf{R}^\top$, we consider the following optimization problem:

$$\underset{\mathbf{F} \in \mathbb{R}^{(n_1+n_2) \times r}, \mathbf{S} \in \mathcal{S}_\alpha}{\text{minimize}} \quad \mathcal{L}(\mathbf{F}, \mathbf{S}) = \frac{1}{2} \|\mathbf{L}\mathbf{R}^\top + \mathbf{S} - \mathbf{Y}\|_{\mathbf{F}}^2. \quad (17)$$

It is thus natural to alternatively update $\mathbf{F} = [\mathbf{L}^\top, \mathbf{R}^\top]^\top$ and \mathbf{S} , where \mathbf{F} is updated via the proposed ScaledGD algorithm, and \mathbf{S} is updated by hard thresholding, which trims the small entries of the residual matrix $\mathbf{Y} - \mathbf{L}\mathbf{R}^\top$. More specifically, for some truncation level $0 \leq \bar{\alpha} \leq 1$, we define the sparsification operator that only keeps $\bar{\alpha}$ fraction of largest entries in each row and column:

$$(\mathcal{T}_{\bar{\alpha}}[\mathbf{A}])_{i,j} = \begin{cases} \mathbf{A}_{i,j}, & \text{if } |\mathbf{A}|_{i,j} \geq |\mathbf{A}|_{i,(\bar{\alpha}n_2)}, \text{ and } |\mathbf{A}|_{i,j} \geq |\mathbf{A}|_{(\bar{\alpha}n_1),j}, \\ 0, & \text{otherwise} \end{cases}, \quad (18)$$

where $|\mathbf{A}|_{i,(k)}$ (resp. $|\mathbf{A}|_{(k),j}$) denote the k -th largest element in magnitude in the i -th row (resp. j -th column).

The ScaledGD algorithm with the spectral initialization for solving robust PCA is formally stated in Algorithm 2. Note that, comparing with Yi et al. (2016), we do not require a balancing term $\|\mathbf{L}^\top \mathbf{L} - \mathbf{R}^\top \mathbf{R}\|_{\mathbf{F}}^2$ in the loss function (17), nor the projection of the low-rank factors onto the $\ell_{2,\infty}$ ball in each iteration.

Algorithm 2 ScaledGD for robust PCA with spectral initialization

Spectral initialization: Let $\mathbf{U}_0 \mathbf{\Sigma}_0 \mathbf{V}_0^\top$ be the top- r SVD of $\mathbf{Y} - \mathcal{T}_\alpha[\mathbf{Y}]$, and set

$$\mathbf{L}_0 = \mathbf{U}_0 \mathbf{\Sigma}_0^{1/2}, \quad \text{and} \quad \mathbf{R}_0 = \mathbf{V}_0 \mathbf{\Sigma}_0^{1/2}. \quad (19)$$

Scaled gradient updates: for $t = 0, 1, 2, \dots, T-1$ do

$$\begin{aligned} \mathbf{S}_t &= \mathcal{T}_{2\alpha}[\mathbf{Y} - \mathbf{L}_t \mathbf{R}_t^\top], \\ \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta(\mathbf{L}_t \mathbf{R}_t^\top + \mathbf{S}_t - \mathbf{Y}) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1}, \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta(\mathbf{L}_t \mathbf{R}_t^\top + \mathbf{S}_t - \mathbf{Y})^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}. \end{aligned} \quad (20)$$

Theoretical guarantee. Before stating our main result for robust PCA, we introduce the incoherence condition which is known to be crucial for reliable estimation of the low-rank matrix \mathbf{X}_* in robust PCA Chen (2015).

Definition 5 (Incoherence) A rank- r matrix $\mathbf{X}_* \in \mathbb{R}^{n_1 \times n_2}$ with compact SVD as $\mathbf{X}_* = \mathbf{U}_* \mathbf{\Sigma}_* \mathbf{V}_*^\top$ is said to be μ -incoherent if

$$\|\mathbf{U}_*\|_{2,\infty} \leq \sqrt{\frac{\mu}{n_1}} \|\mathbf{U}_*\|_F = \sqrt{\frac{\mu r}{n_1}}, \quad \text{and} \quad \|\mathbf{V}_*\|_{2,\infty} \leq \sqrt{\frac{\mu}{n_2}} \|\mathbf{V}_*\|_F = \sqrt{\frac{\mu r}{n_2}}.$$

The following theorem establishes that **ScaledGD** converges linearly at a constant rate as long as the fraction α of corruptions is sufficiently small.

Theorem 6 Suppose that \mathbf{X}_* is μ -incoherent and that the corruption fraction α obeys $\alpha \leq c/(\mu r^{3/2} \kappa)$ for some sufficiently small constant $c > 0$. If the step size obeys $0.1 \leq \eta \leq 2/3$, then for all $t \geq 0$, the iterates of **ScaledGD** in Algorithm 2 satisfy

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_*) \leq (1 - 0.6\eta)^t 0.02\sigma_r(\mathbf{X}_*), \quad \text{and} \quad \|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*\|_F \leq (1 - 0.6\eta)^t 0.03\sigma_r(\mathbf{X}_*).$$

Theorem 6 establishes that the distance $\text{dist}(\mathbf{F}_t, \mathbf{F}_*)$ contracts linearly at a constant rate, as long as the fraction of corruptions satisfies $\alpha \lesssim 1/(\mu r^{3/2} \kappa)$. To reach ϵ -accuracy, i.e. $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*\|_F \leq \epsilon\sigma_r(\mathbf{X}_*)$, **ScaledGD** takes at most $T = O(\log(1/\epsilon))$ iterations, which is *independent* of κ . In comparison, the **AltProj** algorithm³ with spectral initialization converges in $O(\log(1/\epsilon))$ iterations as long as $\alpha \lesssim 1/(\mu r)$ Netrapalli et al. (2014), where the per-iteration cost is much higher both in terms of computation and memory as it requires the computation of the low-rank SVD of the full matrix. On the other hand, projected gradient descent with spectral initialization in Yi et al. (2016) converges in $O(\kappa \log(1/\epsilon))$ iterations as long as $\alpha \lesssim 1/(\mu r^{3/2} \kappa^{3/2} \vee \mu r \kappa^2)$. Therefore, **ScaledGD** converges at a much faster rate than GD while requesting a significantly lower per-iteration cost than **AltProj**. In addition, our theory suggests that **ScaledGD** maintains the incoherence and balancedness of the low-rank factors without imposing explicit regularizations, which is not captured in previous analysis Yi et al. (2016).

2.4 Matrix completion

Assume that we have observed a subset Ω of entries of \mathbf{X}_* given as $\mathcal{P}_\Omega(\mathbf{X}_*)$, where $\mathcal{P}_\Omega : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}^{n_1 \times n_2}$ is a projection such that

$$(\mathcal{P}_\Omega(\mathbf{X}))_{i,j} = \begin{cases} \mathbf{X}_{i,j}, & \text{if } (i,j) \in \Omega \\ 0, & \text{otherwise} \end{cases}. \quad (21)$$

Here Ω is generated according to the Bernoulli model in the sense that each $(i,j) \in \Omega$ independent with probability p . The goal of matrix completion is to recover the matrix \mathbf{X}_* from its partial observation $\mathcal{P}_\Omega(\mathbf{X}_*)$. This problem has many applications in recommendation systems, signal processing, sensor network localization, and so on Candès and Recht (2009).

Algorithm. Again, writing $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ into the factored form $\mathbf{X} = \mathbf{L} \mathbf{R}^\top$, we consider the following optimization problem:

$$\underset{\mathbf{F} \in \mathbb{R}^{(n_1+n_2) \times r}}{\text{minimize}} \quad \mathcal{L}(\mathbf{F}) := \frac{1}{2p} \|\mathcal{P}_\Omega(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_*)\|_F^2. \quad (22)$$

3. **AltProj** employs a multi-stage strategy to remove the dependence on κ in α , which we do not consider here. The same strategy might also improve the dependence on κ for **ScaledGD**, which we leave for future work.

Similarly to robust PCA, the underlying low-rank matrix \mathbf{X}_\star needs to be incoherent (cf. Definition 5) to avoid ill-posedness. One typical strategy to ensure the incoherence condition is to perform projection after the gradient update, by projecting the iterates to maintain small $\ell_{2,\infty}$ norms of the factor matrices. However, the standard projection operator Chen and Wainwright (2015) is not covariant with respect to invertible transforms, and consequently, needs to be modified when using scaled gradient updates. To that end, we introduce the following new projection operator: for every $\tilde{\mathbf{F}} \in \mathbb{R}^{(n_1+n_2) \times r} = [\tilde{\mathbf{L}}^\top, \tilde{\mathbf{R}}^\top]^\top$,

$$\begin{aligned} \mathcal{P}_B(\tilde{\mathbf{F}}) = & \underset{\mathbf{F} \in \mathbb{R}^{(n_1+n_2) \times r}}{\operatorname{argmin}} \left\| (\mathbf{L} - \tilde{\mathbf{L}})(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2} \right\|_{\mathbf{F}}^2 + \left\| (\mathbf{R} - \tilde{\mathbf{R}})(\tilde{\mathbf{L}}^\top \tilde{\mathbf{L}})^{1/2} \right\|_{\mathbf{F}}^2 \\ \text{s.t. } & \sqrt{n_1} \left\| \mathbf{L}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2} \right\|_{2,\infty} \vee \sqrt{n_2} \left\| \mathbf{R}(\tilde{\mathbf{L}}^\top \tilde{\mathbf{L}})^{1/2} \right\|_{2,\infty} \leq B \end{aligned} \quad (23)$$

which finds a factored matrix that is closest to $\tilde{\mathbf{F}}$ and stays incoherent in a weighted sense. Luckily, the solution to the above scaled projection admits a simple closed-form solution, as stated below.

Proposition 7 *The solution to (23) is given by*

$$\begin{aligned} \mathcal{P}_B(\tilde{\mathbf{F}}) := & \begin{bmatrix} \mathbf{L} \\ \mathbf{R} \end{bmatrix}, \quad \text{where } \mathbf{L}_{i,\cdot} := \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \right) \tilde{\mathbf{L}}_{i,\cdot}, \quad 1 \leq i \leq n_1, \\ & \mathbf{R}_{j,\cdot} := \left(1 \wedge \frac{B}{\sqrt{n_2} \|\tilde{\mathbf{R}}_{j,\cdot} \tilde{\mathbf{L}}^\top\|_2} \right) \tilde{\mathbf{R}}_{j,\cdot}, \quad 1 \leq j \leq n_2. \end{aligned} \quad (24)$$

Proof See Appendix E.1.1. ■

With the new projection operator in place, we propose the scaled projected gradient descent (**ScaledPGD**) method with the spectral initialization for solving matrix completion, formally stated in Algorithm 3.

Algorithm 3 ScaledPGD for matrix completion with spectral initialization

Spectral initialization: Let $\mathbf{U}_0 \mathbf{\Sigma}_0 \mathbf{V}_0^\top$ be the top- r SVD of $\frac{1}{p} \mathcal{P}_\Omega(\mathbf{X}_\star)$, and set

$$\begin{bmatrix} \mathbf{L}_0 \\ \mathbf{R}_0 \end{bmatrix} = \mathcal{P}_B \left(\begin{bmatrix} \mathbf{U}_0 \mathbf{\Sigma}_0^{1/2} \\ \mathbf{V}_0 \mathbf{\Sigma}_0^{1/2} \end{bmatrix} \right). \quad (25)$$

Scaled projected gradient updates: for $t = 0, 1, 2, \dots, T-1$ do

$$\begin{bmatrix} \mathbf{L}_{t+1} \\ \mathbf{R}_{t+1} \end{bmatrix} = \mathcal{P}_B \left(\begin{bmatrix} \mathbf{L}_t - \frac{\eta}{p} \mathcal{P}_\Omega(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1} \\ \mathbf{R}_t - \frac{\eta}{p} \mathcal{P}_\Omega(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star)^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1} \end{bmatrix} \right). \quad (26)$$

Theoretical guarantee. Consider a random observation model, where each index (i, j) belongs to the index set Ω independently with probability $0 < p \leq 1$. The following theorem establishes that ScaledPGD converges linearly at a constant rate as long as the number of observations is sufficiently large.

Theorem 8 *Suppose that \mathbf{X}_\star is μ -incoherent, and that p satisfies $p \geq C(\mu\kappa^2 \vee \log(n_1 \vee n_2))\mu r^2 \kappa^2 / (n_1 \wedge n_2)$ for some sufficiently large constant C . Set the projection radius as $B = C_B \sqrt{\mu r} \sigma_1(\mathbf{X}_\star)$ for some constant $C_B \geq 1.02$. If the step size obeys $0 < \eta \leq 2/3$, then with probability at least $1 - c_1(n_1 \vee n_2)^{-c_2}$, for all $t \geq 0$, the iterates of ScaledPGD in (26) satisfy*

$$\operatorname{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq (1 - 0.6\eta)^t 0.02 \sigma_r(\mathbf{X}_\star), \quad \text{and} \quad \|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_{\mathbf{F}} \leq (1 - 0.6\eta)^t 0.03 \sigma_r(\mathbf{X}_\star).$$

Here $c_1, c_2 > 0$ are two universal constants.

Theorem 8 establishes that the distance $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ contracts linearly at a constant rate, as long as the probability of observation satisfies $p \gtrsim (\mu\kappa^2 \vee \log(n_1 \vee n_2))\mu r^2 \kappa^2 / (n_1 \wedge n_2)$. To reach ϵ -accuracy, i.e. $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq \epsilon \sigma_r(\mathbf{X}_\star)$, ScaledPGD takes at most $T = O(\log(1/\epsilon))$ iterations, which is *independent* of κ . In comparison, projected gradient descent Zheng and Lafferty (2016) with spectral initialization converges in $O(\kappa \log(1/\epsilon))$ iterations as long as $p \gtrsim (\mu \vee \log(n_1 \vee n_2))\mu r^2 \kappa^2 / (n_1 \wedge n_2)$. Therefore, ScaledPGD achieves much faster convergence than its unscaled counterpart, at an expense of higher sample complexity. We believe this higher sample complexity is an artifact of our proof techniques, as numerically we do not observe a degradation in terms of sample complexity.

2.5 Optimizing general loss functions

Last but not least, we generalize our analysis of ScaledGD to minimize a general loss function in the form of (2), where the update rule of ScaledGD is given by

$$\begin{aligned} \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta \nabla f(\mathbf{L}_t \mathbf{R}_t^\top) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1}, \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta \nabla f(\mathbf{L}_t \mathbf{R}_t^\top)^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}. \end{aligned} \quad (27)$$

Two important properties of the loss function $f : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}$ play a key role in the analysis.

Definition 9 (Restricted smoothness) A differentiable function $f : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}$ is said to be rank- r restricted L -smooth for some $L > 0$ if

$$f(\mathbf{X}_2) \leq f(\mathbf{X}_1) + \langle \nabla f(\mathbf{X}_1), \mathbf{X}_2 - \mathbf{X}_1 \rangle + \frac{L}{2} \|\mathbf{X}_2 - \mathbf{X}_1\|_F^2,$$

for any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^{n_1 \times n_2}$ with rank at most r .

Definition 10 (Restricted strong convexity) A differentiable function $f : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}$ is said to be rank- r restricted μ -strongly convex for some $\mu \geq 0$ if

$$f(\mathbf{X}_2) \geq f(\mathbf{X}_1) + \langle \nabla f(\mathbf{X}_1), \mathbf{X}_2 - \mathbf{X}_1 \rangle + \frac{\mu}{2} \|\mathbf{X}_2 - \mathbf{X}_1\|_F^2,$$

for any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^{n_1 \times n_2}$ with rank at most r . When $\mu = 0$, we simply say $f(\cdot)$ is rank- r restricted convex.

Further, when $\mu > 0$, define the condition number of the loss function $f(\cdot)$ over rank- r matrices as

$$\kappa_f := L/\mu. \quad (28)$$

Encouragingly, many problems can be viewed as a special case of optimizing this general loss (27), including but not limited to:

- *low-rank matrix factorization*, where the loss function $f(\mathbf{X}) = \frac{1}{2} \|\mathbf{X} - \mathbf{X}_\star\|_F^2$ in (29) satisfies $\kappa_f = 1$;
- *low-rank matrix sensing*, where the loss function $f(\mathbf{X}) = \frac{1}{2} \|\mathcal{A}(\mathbf{X} - \mathbf{X}_\star)\|_2^2$ in (13) satisfies $\kappa_f \approx 1$ when $\mathcal{A}(\cdot)$ obeys the rank- r RIP with a sufficiently small RIP constant;
- *quadratic sampling*, where the loss function $f(\mathbf{X}) = \frac{1}{2} \sum_{i=1}^m |\langle \mathbf{a}_i \mathbf{a}_i^\top, \mathbf{X} - \mathbf{X}_\star \rangle|^2$ satisfies restricted strong convexity and smoothness when \mathbf{a}_i 's are i.i.d. Gaussian vectors for sufficiently large m Sanghavi et al. (2017); Li et al. (2021);
- *exponential-family PCA*, where the loss function $f(\mathbf{X}) = -\sum_{i,j} \log p(\mathbf{Y}_{i,j} | \mathbf{X}_{i,j})$, where $p(\mathbf{Y}_{i,j} | \mathbf{X}_{i,j})$ is the probability density function of $\mathbf{Y}_{i,j}$ conditional on $\mathbf{X}_{i,j}$, following an exponential-family distribution such as Bernoulli and Poisson distributions. The resulting loss function satisfies restricted strong convexity and smoothness with a condition number $\kappa_f > 1$ depending on the property of the specific distribution Gunasekar et al. (2014); Lafond (2015).

Indeed, the treatment of a general loss function brings the condition number of $f(\cdot)$ under the spotlight, since in our earlier case studies $\kappa_f \approx 1$. Our purpose is thus to understand the interplay of two types of conditioning numbers in the convergence of first-order methods. For simplicity, we assume that $f(\cdot)$ is minimized at the ground truth rank- r matrix \mathbf{X}_* .⁴ The following theorem establishes that as long as properly initialized, then `ScaledGD` converges linearly at a constant rate.

Theorem 11 *Suppose that $f(\cdot)$ is rank- $2r$ restricted L -smooth and μ -strongly convex, of which \mathbf{X}_* is a minimizer, and that the initialization \mathbf{F}_0 satisfies $\text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 0.1\sigma_r(\mathbf{X}_*)/\sqrt{\kappa_f}$. If the step size obeys $0 < \eta \leq 0.4/L$, then for all $t \geq 0$, the iterates of `ScaledGD` in (27) satisfy*

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_*) \leq (1 - 0.7\eta\mu)^t 0.1\sigma_r(\mathbf{X}_*)/\sqrt{\kappa_f}, \quad \text{and} \quad \|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*\|_{\text{F}} \leq (1 - 0.7\eta\mu)^t 0.15\sigma_r(\mathbf{X}_*)/\sqrt{\kappa_f}.$$

Theorem 11 establishes that the distance $\text{dist}(\mathbf{F}_t, \mathbf{F}_*)$ contracts linearly at a constant rate, as long as the initialization \mathbf{F}_0 is sufficiently close to \mathbf{F}_* . To reach ϵ -accuracy, i.e. $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*\|_{\text{F}} \leq \epsilon\sigma_r(\mathbf{X}_*)$, `ScaledGD` takes at most $T = O(\kappa_f \log(1/\epsilon))$ iterations, which depends only on the condition number κ_f of $f(\cdot)$, but is independent of the condition number κ of the matrix \mathbf{X}_* . In contrast, prior theory of vanilla gradient descent Park et al. (2018); Bhojanapalli et al. (2016a) requires $O(\kappa_f \kappa \log(1/\epsilon))$ iterations, which is worse than our rate by a factor of κ .

3. Proof Sketch

In this section, we sketch the proof of the main theorems, highlighting the role of the scaled distance metric (cf. (10)) in these analyses.

3.1 A warm-up analysis: matrix factorization

Let us consider the problem of factorizing a matrix \mathbf{X}_* into two low-rank factors:

$$\underset{\mathbf{F} \in \mathbb{R}^{(n_1+n_2) \times r}}{\text{minimize}} \quad \mathcal{L}(\mathbf{F}) = \frac{1}{2} \|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_*\|_{\text{F}}^2. \quad (29)$$

For this toy problem, the update rule of `ScaledGD` is given as

$$\begin{aligned} \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1}, \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*)^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1}. \end{aligned} \quad (30)$$

To shed light on why `ScaledGD` is robust to ill-conditioning, it is worthwhile to think of `ScaledGD` as a quasi-Newton algorithm: the following proposition (proven in Appendix B.1) reveals that `ScaledGD` is equivalent to approximating the Hessian of the loss function in (29) by only keeping its diagonal blocks.

Proposition 12 *For the matrix factorization problem (29), `ScaledGD` is equivalent to the following update rule*

$$\text{vec}(\mathbf{F}_{t+1}) = \text{vec}(\mathbf{F}_t) - \eta \begin{bmatrix} \nabla_{\mathbf{L}, \mathbf{L}}^2 \mathcal{L}(\mathbf{F}_t) & \mathbf{0} \\ \mathbf{0} & \nabla_{\mathbf{R}, \mathbf{R}}^2 \mathcal{L}(\mathbf{F}_t) \end{bmatrix}^{-1} \text{vec}(\nabla_{\mathbf{F}} \mathcal{L}(\mathbf{F}_t)).$$

Here, $\nabla_{\mathbf{L}, \mathbf{L}}^2 \mathcal{L}(\mathbf{F}_t)$ (resp. $\nabla_{\mathbf{R}, \mathbf{R}}^2 \mathcal{L}(\mathbf{F}_t)$) denotes the second order derivative w.r.t. \mathbf{L} (resp. \mathbf{R}) at \mathbf{F}_t .

The following theorem, whose proof can be found in Appendix B.2, formally establishes that as long as `ScaledGD` is initialized close to the ground truth, $\text{dist}(\mathbf{F}_t, \mathbf{F}_*)$ will contract at a constant linear rate for the matrix factorization problem.

4. In practice, due to the presence of statistical noise, the minimizer of $f(\cdot)$ might be only approximately low-rank, to which our analysis can be extended in a straightforward fashion.

Theorem 13 *Suppose that the initialization \mathbf{F}_0 satisfies $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 0.1\sigma_r(\mathbf{X}_\star)$. If the step size obeys $0 < \eta \leq 2/3$, then for all $t \geq 0$, the iterates of the *ScaledGD* method in (30) satisfy*

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq (1 - 0.7\eta)^t 0.1\sigma_r(\mathbf{X}_\star), \quad \text{and} \quad \|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq (1 - 0.7\eta)^t 0.15\sigma_r(\mathbf{X}_\star).$$

Comparing to the rate of contraction $(1 - 1/\kappa)$ of gradient descent for matrix factorization Ma et al. (2021); Chi et al. (2019), Theorem 13 demonstrates that the preconditioners indeed allow better search directions in the local neighborhood of the ground truth, and hence a faster convergence rate.

3.2 Proof outline for matrix sensing

It can be seen that the update rule (15) of *ScaledGD* in Algorithm 1 closely mimics (30) when $\mathcal{A}(\cdot)$ satisfies the RIP. Therefore, leveraging the RIP of $\mathcal{A}(\cdot)$ and Theorem 13, we can establish the following local convergence guarantee of Algorithm 1, which has a weaker requirement on δ_{2r} than the main theorem (cf. Theorem 3).

Lemma 14 *Suppose that $\mathcal{A}(\cdot)$ obeys the $2r$ -RIP with $\delta_{2r} \leq 0.02$. If the t -th iterate satisfies $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.1\sigma_r(\mathbf{X}_\star)$, then $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$. In addition, if the step size obeys $0 < \eta \leq 2/3$, then the $(t + 1)$ -th iterate \mathbf{F}_{t+1} of the *ScaledGD* method in (15) of Algorithm 1 satisfies*

$$\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star).$$

It then boils to down to finding a good initialization, for which we have the following lemma on the quality of the spectral initialization.

Lemma 15 *Suppose that $\mathcal{A}(\cdot)$ obeys the $2r$ -RIP with a constant δ_{2r} . Then the spectral initialization in (14) for low-rank matrix sensing satisfies*

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 5\delta_{2r}\sqrt{r}\kappa\sigma_r(\mathbf{X}_\star).$$

Therefore, as long as δ_{2r} is small enough, say $\delta_{2r} \leq 0.02/(\sqrt{r}\kappa)$ as specified in Theorem 3, the initial distance satisfies $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 0.1\sigma_r(\mathbf{X}_\star)$, allowing us to invoke Lemma 14 recursively. The proof of Theorem 3 is then complete. The proofs of Lemmas 14-15 can be found in Appendix C.

3.3 Proof outline for robust PCA

As before, we begin with the following local convergence guarantee of Algorithm 2, which has a weaker requirement on α than the main theorem (cf. Theorem 6). The difference with low-rank matrix sensing is that local convergence for robust PCA requires a further incoherence condition on the iterates (cf. (31)), where we recall from (11) that \mathbf{Q}_t is the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star .

Lemma 16 *Suppose that \mathbf{X}_\star is μ -incoherent and $\alpha \leq 10^{-4}/(\mu r)$. If the t -th iterate satisfies $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$ and the incoherence condition*

$$\sqrt{n_1} \left\| (\mathbf{L}_t \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \vee \sqrt{n_2} \left\| (\mathbf{R}_t \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_\star), \quad (31)$$

*then $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$. In addition, if the step size obeys $0.1 \leq \eta \leq 2/3$, then the $(t + 1)$ -th iterate \mathbf{F}_{t+1} of the *ScaledGD* method in (20) of Algorithm 2 satisfies*

$$\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star),$$

and the incoherence condition

$$\sqrt{n_1} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_{t+1} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \vee \sqrt{n_2} \left\| (\mathbf{R}_{t+1} \mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_\star).$$

As long as the initialization is close to the ground truth and satisfies the incoherence condition, Lemma 16 ensures that the iterates of `ScaledGD` remain incoherent and converge linearly. This allows us to remove the unnecessary projection step in Yi et al. (2016), whose main objective is to ensure the incoherence of the iterates.

We are left with checking the initial conditions. The following lemma ensures that the spectral initialization in (19) is close to the ground truth as long as α is sufficiently small.

Lemma 17 *Suppose that \mathbf{X}_\star is μ -incoherent. Then the spectral initialization (19) for robust PCA satisfies*

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 20\alpha\mu r^{3/2}\kappa\sigma_r(\mathbf{X}_\star).$$

As a result, setting $\alpha \leq 10^{-3}/(\mu r^{3/2}\kappa)$, the spectral initialization satisfies $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$. In addition, we need to make sure that the spectral initialization satisfies the incoherence condition, which is provided in the following lemma.

Lemma 18 *Suppose that \mathbf{X}_\star is μ -incoherent and $\alpha \leq 0.1/(\mu r\kappa)$, and that $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$. Then the spectral initialization (19) satisfies the incoherence condition*

$$\sqrt{n_1} \left\| (\mathbf{L}_0 \mathbf{Q}_0 - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \vee \sqrt{n_2} \left\| (\mathbf{R}_0 \mathbf{Q}_0^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_\star).$$

Combining Lemmas 16-18 finishes the proof of Theorem 6. The proofs of the the three supporting lemmas can be found in Section D.

3.4 Proof outline for matrix completion

A key property of the new projection operator. We start with the following lemma that entails a key property of the scaled projection (24), which ensures the scaled projection satisfies both non-expansiveness and incoherence under the scaled metric.

Lemma 19 *Suppose that \mathbf{X}_\star is μ -incoherent, and $\text{dist}(\tilde{\mathbf{F}}, \mathbf{F}_\star) \leq \epsilon\sigma_r(\mathbf{X}_\star)$ for some $\epsilon < 1$. Set $B \geq (1 + \epsilon)\sqrt{\mu r}\sigma_1(\mathbf{X}_\star)$, then $\mathcal{P}_B(\tilde{\mathbf{F}})$ satisfies the non-expansiveness*

$$\text{dist}(\mathcal{P}_B(\tilde{\mathbf{F}}), \mathbf{F}_\star) \leq \text{dist}(\tilde{\mathbf{F}}, \mathbf{F}_\star),$$

and the incoherence condition

$$\sqrt{n_1} \|\mathbf{L}\mathbf{R}^\top\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R}\mathbf{L}^\top\|_{2,\infty} \leq B.$$

It is worth noting that the incoherence condition adopts a slightly different form than that of robust PCA, which is more convenient for matrix completion. The next lemma guarantees the fast local convergence of Algorithm 3 as long as the sample complexity is large enough and the parameter B is set properly.

Lemma 20 *Suppose that \mathbf{X}_\star is μ -incoherent, and $p \geq C(\mu r\kappa^4 \vee \log(n_1 \vee n_2))\mu r/(n_1 \wedge n_2)$ for some sufficiently large constant C . Set the projection radius as $B = C_B\sqrt{\mu r}\sigma_1(\mathbf{X}_\star)$ for some constant $C_B \geq 1.02$. Under an event \mathcal{E} which happens with overwhelming probability (i.e. at least $1 - c_1(n_1 \vee n_2)^{-c_2}$), if the t -th iterate satisfies $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$, and the incoherence condition*

$$\sqrt{n_1} \|\mathbf{L}_t \mathbf{R}_t^\top\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R}_t \mathbf{L}_t^\top\|_{2,\infty} \leq B,$$

then $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_{\text{F}} \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$. In addition, if the step size obeys $0 < \eta \leq 2/3$, then the $(t + 1)$ -th iterate \mathbf{F}_{t+1} of the `ScaledPGD` method in (26) of Algorithm 3 satisfies

$$\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star),$$

and the incoherence condition

$$\sqrt{n_1}\|\mathbf{L}_{t+1}\mathbf{R}_{t+1}^\top\|_{2,\infty} \vee \sqrt{n_2}\|\mathbf{R}_{t+1}\mathbf{L}_{t+1}^\top\|_{2,\infty} \leq B.$$

As long as we can find an initialization that is close to the ground truth and satisfies the incoherence condition, Lemma 20 ensures that the iterates of **ScaledPGD** remain incoherent and converge linearly. The follow lemma ensures that such an initialization can be ensured via the spectral method.

Lemma 21 *Suppose that \mathbf{X}_\star is μ -incoherent, then with overwhelming probability, the spectral initialization before projection $\tilde{\mathbf{F}}_0 := \begin{bmatrix} \mathbf{U}_0 \Sigma_0^{1/2} \\ \mathbf{V}_0 \Sigma_0^{1/2} \end{bmatrix}$ in (25) satisfies*

$$\text{dist}(\tilde{\mathbf{F}}_0, \mathbf{F}_\star) \leq C_0 \left(\frac{\mu r \log(n_1 \vee n_2)}{p\sqrt{n_1 n_2}} + \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \right) 5\sqrt{r}\kappa\sigma_r(\mathbf{X}_\star).$$

Therefore, as long as $p \geq C\mu r^2 \kappa^2 \log(n_1 \vee n_2)/(n_1 \wedge n_2)$ for some sufficiently large constant C , the initial distance satisfies $\text{dist}(\tilde{\mathbf{F}}_0, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$. One can then invoke Lemma 19 to see that $\mathbf{F}_0 = \mathcal{P}_B(\tilde{\mathbf{F}}_0)$ meets the requirements of Lemma 20 due to the non-expansiveness and incoherence properties of the projection operator. The proofs of the the supporting lemmas can be found in Section E.

4. Numerical Experiments

In this section, we provide numerical experiments to corroborate our theoretical findings, with the codes available at

<https://github.com/Titan-Tong/ScaledGD>.

The simulations are performed in Matlab with a 3.6 GHz Intel Xeon Gold 6244 CPU.

4.1 Comparison with vanilla GD

To begin, we compare the iteration complexity of **ScaledGD** with vanilla gradient descent (GD). The update rule of vanilla GD for solving (2) is given as

$$\begin{aligned} \mathbf{L}_{t+1} &= \mathbf{L}_t - \eta_{\text{GD}} \nabla_{\mathbf{L}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t), \\ \mathbf{R}_{t+1} &= \mathbf{R}_t - \eta_{\text{GD}} \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{L}_t, \mathbf{R}_t), \end{aligned} \quad (32)$$

where $\eta_{\text{GD}} = \eta/\sigma_1(\mathbf{X}_\star)$ stands for the step size for gradient descent. This choice is often recommended by the theory of vanilla GD Tu et al. (2016); Yi et al. (2016); Ma et al. (2019) and the scaling by $\sigma_1(\mathbf{X}_\star)$ is needed for its convergence. For ease of comparison, we fix $\eta = 0.5$ for both **ScaledGD** and vanilla GD (see Figure 4 for justifications). Both algorithms start from the same spectral initialization. To avoid notational clutter, we work on square *asymmetric* matrices with $n_1 = n_2 = n$. We consider four low-rank matrix estimation tasks:

- *Low-rank matrix sensing.* The problem formulation is detailed in Section 2.2. Here, we collect $m = 5nr$ measurements in the form of $\mathbf{y}_k = \langle \mathbf{A}_k, \mathbf{X}_\star \rangle + \mathbf{w}_k$, in which the measurement matrices \mathbf{A}_k are generated with i.i.d. Gaussian entries with zero mean and variance $1/m$, and $\mathbf{w}_k \sim \mathcal{N}(0, \sigma_w^2)$ are i.i.d. Gaussian noises.
- *Robust PCA.* The problem formulation is stated in Section 2.3. We generate the corruption with a sparse matrix $\mathbf{S}_\star \in \mathcal{S}_\alpha$ with $\alpha = 0.1$. More specifically, we generate a matrix with standard Gaussian entries and pass it through $\mathcal{T}_\alpha[\cdot]$ to obtain \mathbf{S}_\star . The observation is $\mathbf{Y} = \mathbf{X}_\star + \mathbf{S}_\star + \mathbf{W}$, where $\mathbf{W}_{i,j} \sim \mathcal{N}(0, \sigma_w^2)$ are i.i.d. Gaussian noises.

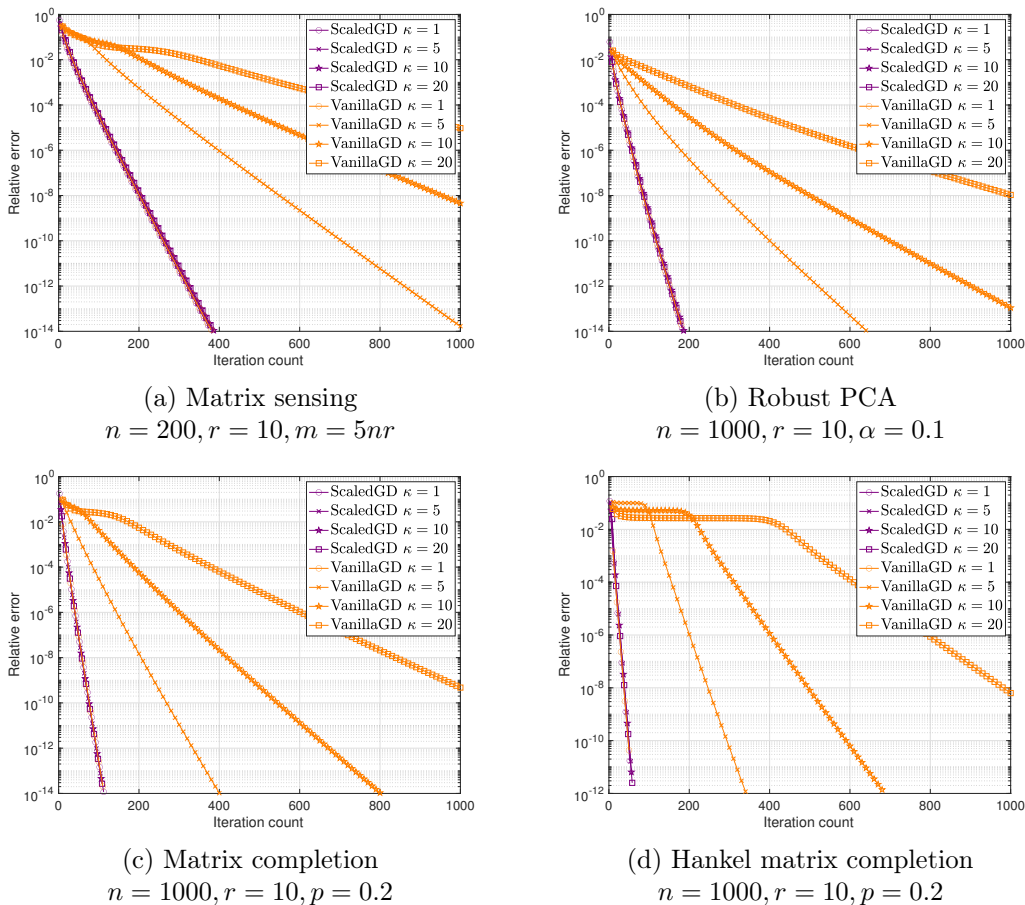


Figure 2: The relative errors of ScaledGD and vanilla GD with respect to the iteration count under different condition numbers $\kappa = 1, 5, 10, 20$ for (a) matrix sensing, (b) robust PCA, (c) matrix completion, and (d) Hankel matrix completion.

- *Matrix completion.* The problem formulation is stated in Section 2.4. We assume random Bernoulli observations, where each entry of \mathbf{X}_* is observed with probability $p = 0.2$ independently. The observation is $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X}_* + \mathbf{W})$, where $\mathbf{W}_{i,j} \sim \mathcal{N}(0, \sigma_w^2)$ are i.i.d. Gaussian noises. Moreover, we perform the scaled gradient updates without projections.
- *Hankel matrix completion.* Briefly speaking, a Hankel matrix shares the same value along each skew-diagonal, and we aim at recovering a low-rank Hankel matrix from observing a few skew-diagonals Chen and Chi (2014); Cai et al. (2018). We assume random Bernoulli observations, where each skew-diagonal of \mathbf{X}_* is observed with probability $p = 0.2$ independently. The loss function is

$$\mathcal{L}(\mathbf{L}, \mathbf{R}) = \frac{1}{2p} \|\mathcal{H}_\Omega(\mathbf{L}\mathbf{R}^\top - \mathbf{Y})\|_F^2 + \frac{1}{2} \|(\mathcal{I} - \mathcal{H})(\mathbf{L}\mathbf{R}^\top)\|_F^2, \quad (33)$$

where $\mathcal{I}(\cdot)$ denotes the identity operator, and the Hankel projection is defined as $\mathcal{H}(\mathbf{X}) := \sum_{k=1}^{2n-1} \langle \mathbf{H}_k, \mathbf{X} \rangle \mathbf{H}_k$, which maps \mathbf{X} to its closest Hankel matrix. Here, the Hankel basis matrix \mathbf{H}_k is the $n \times n$ matrix with the entries in the k -th skew diagonal as $\frac{1}{\sqrt{\omega_k}}$, and all other

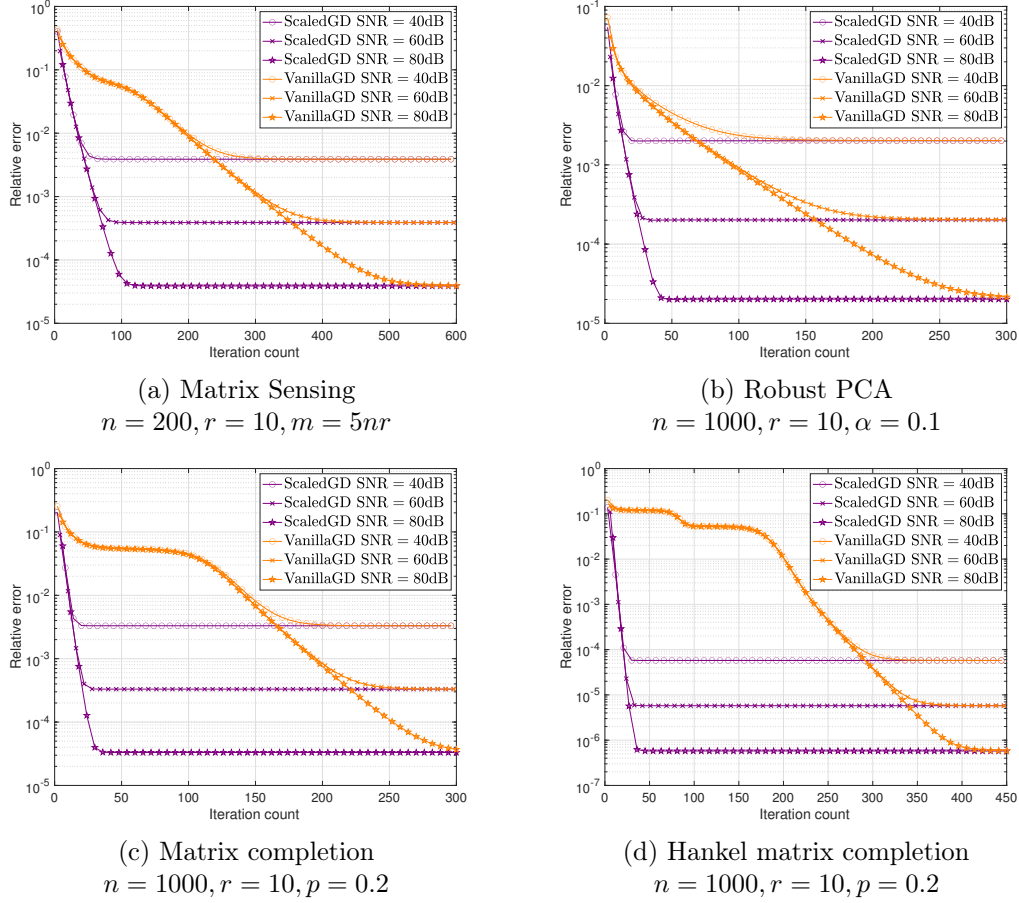


Figure 3: The relative errors of **ScaledGD** and vanilla GD with respect to the iteration count under the condition number $\kappa = 10$ and signal-to-noise ratios SNR = 40, 60, 80dB for (a) matrix sensing, (b) robust PCA, (c) matrix completion, and (d) Hankel matrix completion.

entries as 0, where ω_k is the length of the k -th skew diagonal. Note that \mathbf{X} is a Hankel matrix if and only if $(\mathcal{I} - \mathcal{H})(\mathbf{X}) = \mathbf{0}$. The Hankel projection on the observation index set Ω is defined as $\mathcal{H}_\Omega(\mathbf{X}) := \sum_{k \in \Omega} \langle \mathbf{H}_k, \mathbf{X} \rangle \mathbf{H}_k$. The observation is $\mathbf{Y} = \mathcal{H}_\Omega(\mathbf{X}_* + \mathbf{W})$, where \mathbf{W} is a Hankel matrix whose entries along each skew-diagonal are i.i.d. Gaussian noises $\mathcal{N}(0, \sigma_w^2)$.

For the first three problems, we generate the ground truth matrix $\mathbf{X}_* \in \mathbb{R}^{n \times n}$ in the following way. We first generate an $n \times r$ matrix with i.i.d. random signs, and take its r left singular vectors as \mathbf{U}_* , and similarly for \mathbf{V}_* . The singular values are set to be linearly distributed from 1 to $1/\kappa$. The ground truth is then defined as $\mathbf{X}_* = \mathbf{U}_* \Sigma_* \mathbf{V}_*^\top$ which has the specified condition number κ and rank r . For Hankel matrix completion, we generate \mathbf{X}_* as an $n \times n$ Hankel matrix with entries given as

$$(\mathbf{X}_*)_{i,j} = \sum_{\ell=1}^r \frac{\sigma_\ell}{n} e^{2\pi i(i+j-2)f_\ell}, \quad i, j = 1, \dots, n,$$

where $f_\ell, \ell = 1, \dots, r$ are randomly chosen from $1/n, 2/n, \dots, 1$, and σ_ℓ are linearly distributed from 1 to $1/\kappa$. The Vandermonde decomposition lemma tells that \mathbf{X}_\star has rank r and singular values $\sigma_\ell, \ell = 1, \dots, r$.

We first illustrate the convergence performance under noise-free observations, i.e. $\sigma_w = 0$. We plot the relative reconstruction error $\|\mathbf{X}_t - \mathbf{X}_\star\|_F / \|\mathbf{X}_\star\|_F$ with respect to the iteration count t in Figure 2 for the four problems under different condition numbers $\kappa = 1, 5, 10, 20$. For all these models, we can see that **ScaledGD** has a convergence rate independent of κ , with all curves almost overlay on each other. Under good conditioning $\kappa = 1$, **ScaledGD** converges at the same rate as vanilla GD; under ill conditioning, i.e. when κ is large, **ScaledGD** converges much faster than vanilla GD and leads to significant computational savings.

We next move to demonstrate that **ScaledGD** is robust to small additive noises. Denote the signal-to-noise ratio as $\text{SNR} := 10 \log_{10} \frac{\|\mathbf{X}_\star\|_F^2}{n^2 \sigma_w^2}$ in dB. We plot the reconstruction error $\|\mathbf{X}_t - \mathbf{X}_\star\|_F / \|\mathbf{X}_\star\|_F$ with respect to the iteration count t in Figure 3 under the condition number $\kappa = 10$ and various $\text{SNR} = 40, 60, 80$ dB. We can see that **ScaledGD** and vanilla GD achieve the same statistical error eventually, but **ScaledGD** converges much faster. In addition, the convergence speeds are not influenced by the noise levels.

Careful readers might wonder how sensitivity our comparisons are with respect to the choice of step sizes. To address this, we illustrate the convergence speeds of both **ScaledGD** and vanilla GD under different step sizes η for matrix completion (under the same setting as Figure 2 (c)), where similar plots can be obtained for other problems as well. We run both algorithms for at most 80 iterations, and terminate if the relative error exceeds 10^2 (which happens if the step size is too large and the algorithm diverges). Figure 4 plots the relative error with respect to the step size η for both algorithms, where we can see that **ScaledGD** outperforms vanilla GD over a large range of step sizes, even under optimized values for performance. Hence, our choice of $\eta = 0.5$ in previous experiments renders a typical comparison between **ScaledGD** and vanilla GD.

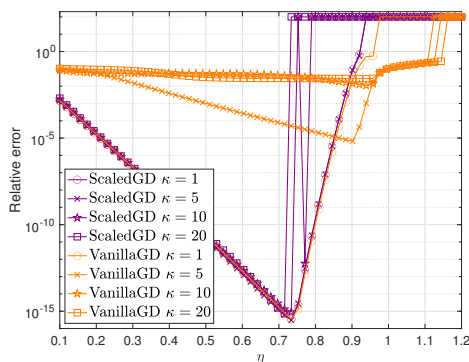


Figure 4: The relative errors of **ScaledGD** and vanilla GD after 80 iterations with respect to different step sizes η from 0.1 to 1.2, for matrix completion with $n = 1000, r = 10, p = 0.2$.

4.2 Run time comparisons

We now compare the run time of **ScaledGD** with vanilla GD and alternating minimization (**AltMin**) Jain et al. (2013). Specifically, for matrix sensing, alternating minimization (**AltMinSense**) updates the factors alternatively as

$$\mathbf{L}_{t+1} = \underset{\mathbf{L}}{\operatorname{argmin}} \|\mathcal{A}(\mathbf{L}\mathbf{R}_t^\top) - \mathbf{y}\|_2^2,$$

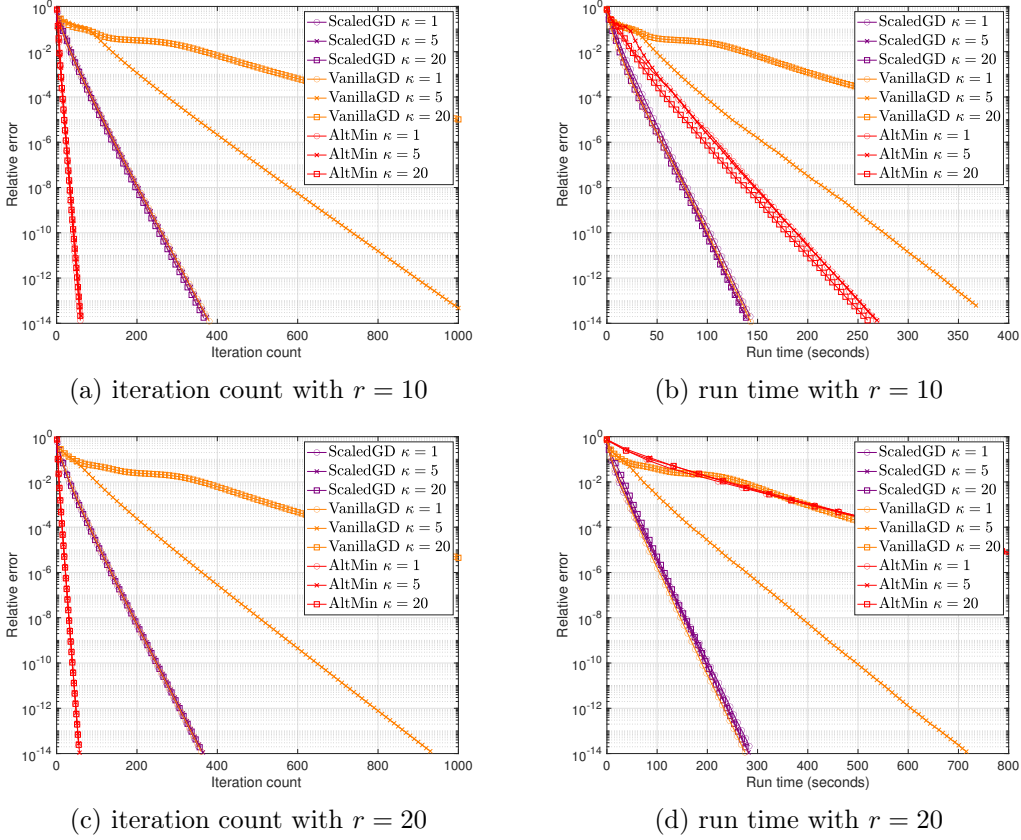


Figure 5: The relative errors of **ScaledGD**, vanilla GD and **AltMin** with respect to the iteration count and run time (in seconds) under different condition numbers $\kappa = 1, 5, 20$ for matrix sensing with $n = 200$, and $m = 5nr$. (a, b): $r = 10$; (c, d): $r = 20$.

$$\mathbf{R}_{t+1} = \underset{\mathbf{R}}{\operatorname{argmin}} \|\mathcal{A}(\mathbf{L}_{t+1}\mathbf{R}^\top) - \mathbf{y}\|_2^2,$$

which corresponds to solving two least-squares problems. For matrix completion, the update rule of alternating minimization proceeds as

$$\begin{aligned} \mathbf{L}_{t+1} &= \underset{\mathbf{L}}{\operatorname{argmin}} \|\mathcal{P}_\Omega(\mathbf{L}\mathbf{R}_t^\top - \mathbf{Y})\|_2^2, \\ \mathbf{R}_{t+1} &= \underset{\mathbf{R}}{\operatorname{argmin}} \|\mathcal{P}_\Omega(\mathbf{L}_{t+1}\mathbf{R}^\top - \mathbf{Y})\|_2^2, \end{aligned}$$

which can be implemented more efficiently since each row of \mathbf{L} (resp. \mathbf{R}) can be updated independently via solving a much smaller least-squares problem due to the decomposable structure of the objective function. It is worth noting that, to the best of our knowledge, this most natural variant of alternating minimization for matrix completion still eludes from a provable performance guarantee, nonetheless, we choose it to compare against due to its popularity and excellent empirical performance.

Figure 5 plots the relative errors of **ScaledGD**, vanilla GD and alternating minimization (**AltMin**) with respect to the iteration count and run time (in seconds) under different condition numbers $\kappa = 1, 5, 20$; and similarly, Figure 6 plots the corresponding results for matrix completion. It

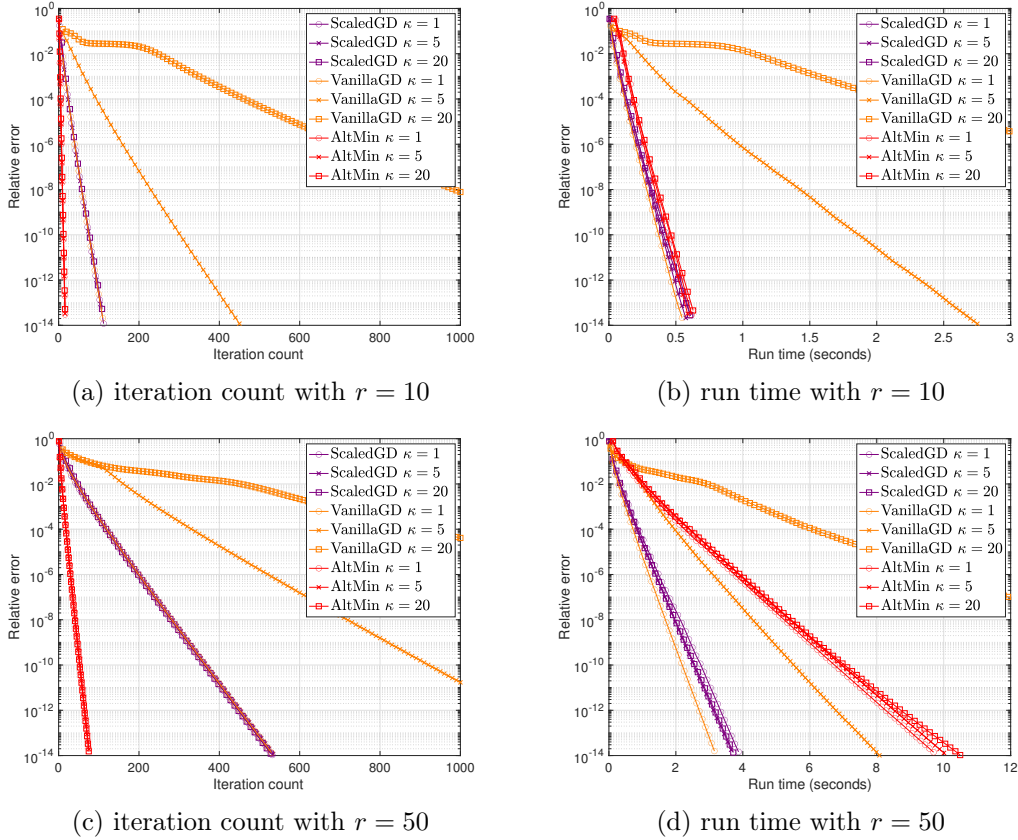


Figure 6: The relative errors of ScaledGD, vanilla GD and AltMin with respect to the iteration count and run time (in seconds) under different condition numbers $\kappa = 1, 5, 20$ for matrix completion with $n = 1000$, and $p = 0.2$. (a, b): $r = 10$; (c, d): $r = 50$.

can be seen that, both ScaledGD and AltMin admit a convergence rate that is independent of the condition number, where the per-iteration complexity of AltMin is much higher than that of ScaledGD. As expected, the run time of ScaledGD only adds a minimal overhead to vanilla GD while being much more robust to ill-conditioning. Noteworthily, AltMin takes much more time and becomes significantly slower than ScaledGD when the rank r is larger. Nonetheless, we emphasize that since the run time is impacted by many factors in terms of problem parameters as well as implementation details, our purpose is to demonstrate the competitive performance of ScaledGD over alternatives, rather than claiming it as the state-of-the-art.

5. Conclusions

This paper proposes scaled gradient descent (ScaledGD) for factored low-rank matrix estimation, which maintains the low per-iteration computational complexity of vanilla gradient descent, but offers significant speed-up in terms of the convergence rate with respect to the condition number κ of the low-rank matrix. In particular, we rigorously establish that for low-rank matrix sensing, robust PCA, and matrix completion, to reach ϵ -accuracy, ScaledGD only takes $O(\log(1/\epsilon))$ iterations without the dependency on the condition number when initialized via the spectral method, under

standard assumptions. The key to our analysis is the introduction of a new distance metric that takes into account the preconditioning and unbalancedness of the low-rank factors, and we have developed new tools to analyze the trajectory of `ScaledGD` under this new metric. This work opens up many venues for future research, as we discuss below.

- *Improved analysis.* In this paper, we have focused on establishing the fast local convergence rate. It is interesting to study if the theory developed herein can be further strengthened in terms of sample complexity and the size of basin of attraction. For matrix completion, it will be interesting to see if a similar guarantee continues to hold in the absence of the projection, which will generalize recent works Ma et al. (2019); Chen et al. (2020a) that successfully removed these projections for vanilla gradient descent.
- *Other low-rank recovery problems.* Besides the problems studied herein, there are many other applications involving the recovery of an ill-conditioned low-rank matrix, such as robust PCA with missing data, quadratic sampling, and so on. It is of interest to establish fast convergence rates of `ScaledGD` that are independent of the condition number for these problems as well. In addition, it is worthwhile to explore if a similar preconditioning trick can be useful to problems beyond low-rank matrix estimation. One recent attempt is to generalize `ScaledGD` for low-rank tensor estimation Tong et al. (2021b).
- *Acceleration schemes?* As it is evident from our analysis of the general loss case, `ScaledGD` may still converge slowly when the loss function is ill-conditioned over low-rank matrices, i.e. κ_f is large. In this case, it might be of interest to combine techniques such as momentum Kyriallidis and Cevher (2012) from the optimization literature to further accelerate the convergence. In our companion paper Tong et al. (2021a), we have extended `ScaledGD` to nonsmooth formulations, which possess better curvatures than their smooth counterparts for certain problems.

Acknowledgments

The work of T. Tong and Y. Chi is supported in part by ONR under the grants N00014-18-1-2142 and N00014-19-1-2404, by ARO under the grant W911NF-18-1-0303, and by NSF under the grants CAREER ECCS-1818571, CCF-1806154 and CCF-1901199.

Appendix A. Technical Lemmas

This section gathers several useful lemmas that will be used in the appendix. Throughout all lemmas, we use \mathbf{X}_\star to denote the ground truth low-rank matrix, with its compact SVD as $\mathbf{X}_\star = \mathbf{U}_\star \mathbf{\Sigma}_\star \mathbf{V}_\star^\top$,

and the stacked factor matrix is defined as $\mathbf{F}_\star = \begin{bmatrix} \mathbf{L}_\star \\ \mathbf{R}_\star \end{bmatrix} = \begin{bmatrix} \mathbf{U}_\star \mathbf{\Sigma}_\star^{1/2} \\ \mathbf{V}_\star \mathbf{\Sigma}_\star^{1/2} \end{bmatrix}$.

A.1 New distance metric

We begin with the investigation of the new distance metric (10), where the matrix \mathbf{Q} that attains the infimum, if exists, is called the optimal alignment matrix between \mathbf{F} and \mathbf{F}_\star ; see (11). Notice that (10) involves a minimization problem over an open set (the set of invertible matrices). Hence the minimizer, i.e. the optimal alignment matrix between \mathbf{F} and \mathbf{F}_\star is not guaranteed to be attained. Fortunately, a simple sufficient condition guarantees the existence of the minimizer; see the lemma below.

Lemma 22 Fix any factor matrix $\mathbf{F} = \begin{bmatrix} \mathbf{L} \\ \mathbf{R} \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times r}$. Suppose that

$$\text{dist}(\mathbf{F}, \mathbf{F}_\star) = \sqrt{\inf_{\mathbf{Q} \in \text{GL}(r)} \left(\left\| (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 \right)} < \sigma_r(\mathbf{X}_\star), \quad (34)$$

then the minimizer of the above minimization problem is attained at some $\mathbf{Q} \in \text{GL}(r)$, i.e. the optimal alignment matrix \mathbf{Q} between \mathbf{F} and \mathbf{F}_\star exists.

Proof In view of the condition (34) and the definition of infimum, one knows that there must exist a matrix $\bar{\mathbf{Q}} \in \text{GL}(r)$ such that

$$\sqrt{\left\| (\mathbf{L}\bar{\mathbf{Q}} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2} \leq \epsilon \sigma_r(\mathbf{X}_\star),$$

for some ϵ obeying $0 < \epsilon < 1$. It further implies that

$$\left\| (\mathbf{L}\bar{\mathbf{Q}} - \mathbf{L}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \vee \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \leq \epsilon.$$

Invoke Weyl's inequality $|\sigma_r(\mathbf{A}) - \sigma_r(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}}$, and use that $\sigma_r(\mathbf{L}_\star \Sigma_\star^{-1/2}) = \sigma_r(\mathbf{U}_\star) = 1$ to obtain

$$\sigma_r(\mathbf{L}\bar{\mathbf{Q}} \Sigma_\star^{-1/2}) \geq \sigma_r(\mathbf{L}_\star \Sigma_\star^{-1/2}) - \left\| (\mathbf{L}\bar{\mathbf{Q}} - \mathbf{L}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \geq 1 - \epsilon. \quad (35)$$

In addition, it is straightforward to verify that

$$\inf_{\mathbf{Q} \in \text{GL}(r)} \left\| (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 \quad (36)$$

$$= \inf_{\mathbf{H} \in \text{GL}(r)} \left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} \mathbf{H}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2. \quad (37)$$

Indeed, if the minimizer of the second optimization problem (cf. (37)) is attained at some \mathbf{H} , then $\bar{\mathbf{Q}}\mathbf{H}$ must be the minimizer of the first problem (36). Therefore, from now on, we focus on proving that the minimizer of the second problem (37) is attained at some \mathbf{H} . In view of (36) and (37), one has

$$\begin{aligned} & \inf_{\mathbf{H} \in \text{GL}(r)} \left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} \mathbf{H}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 \\ & \leq \left\| (\mathbf{L}\bar{\mathbf{Q}} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2, \end{aligned}$$

Clearly, for any $\bar{\mathbf{Q}}\mathbf{H}$ to yield a smaller distance than $\bar{\mathbf{Q}}$, \mathbf{H} must obey

$$\sqrt{\left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} \mathbf{H}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2} \leq \epsilon \sigma_r(\mathbf{X}_\star).$$

It further implies that

$$\left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \vee \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} \mathbf{H}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \leq \epsilon.$$

Invoke Weyl's inequality $|\sigma_1(\mathbf{A}) - \sigma_1(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}}$, and use that $\sigma_1(\mathbf{L}_\star \Sigma_\star^{-1/2}) = \sigma_1(\mathbf{U}_\star) = 1$ to obtain

$$\sigma_1(\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} \Sigma_\star^{-1/2}) \leq \sigma_1(\mathbf{L}_\star \Sigma_\star^{-1/2}) + \left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \leq 1 + \epsilon. \quad (38)$$

Combine (35) and (38), and use the relation $\sigma_r(\mathbf{A})\sigma_1(\mathbf{B}) \leq \sigma_1(\mathbf{AB})$ to obtain

$$\sigma_r(\mathbf{L}\bar{\mathbf{Q}}\Sigma_\star^{-1/2})\sigma_1(\Sigma_\star^{1/2}\mathbf{H}\Sigma_\star^{-1/2}) \leq \sigma_1(\mathbf{L}\bar{\mathbf{Q}}\mathbf{H}\Sigma_\star^{-1/2}) \leq \frac{1+\epsilon}{1-\epsilon}\sigma_r(\mathbf{L}\bar{\mathbf{Q}}\Sigma_\star^{-1/2}).$$

As a result, one has $\sigma_1(\Sigma_\star^{1/2}\mathbf{H}\Sigma_\star^{-1/2}) \leq \frac{1+\epsilon}{1-\epsilon}$.

Similarly, one can show that $\sigma_1(\Sigma_\star^{1/2}\mathbf{H}^{-\top}\Sigma_\star^{-1/2}) \leq \frac{1+\epsilon}{1-\epsilon}$, equivalently, $\sigma_r(\Sigma_\star^{1/2}\mathbf{H}\Sigma_\star^{-1/2}) \geq \frac{1-\epsilon}{1+\epsilon}$. Combining the above two arguments reveals that the minimization problem (37) is equivalent to the constrained problem:

$$\begin{aligned} & \underset{\mathbf{H} \in \text{GL}(r)}{\text{minimize}} && \left\| (\mathbf{L}\bar{\mathbf{Q}}\mathbf{H} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\bar{\mathbf{Q}}^{-\top} \mathbf{H}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 \\ & \text{s.t.} && \frac{1-\epsilon}{1+\epsilon} \leq \sigma_r(\Sigma_\star^{1/2}\mathbf{H}\Sigma_\star^{-1/2}) \leq \sigma_1(\Sigma_\star^{1/2}\mathbf{H}\Sigma_\star^{-1/2}) \leq \frac{1+\epsilon}{1-\epsilon}. \end{aligned}$$

Notice that this is a continuous optimization problem over a compact set. Apply the Weierstrass extreme value theorem to finish the proof. \blacksquare

With the existence of the optimal alignment matrix in place, the following lemma provides the first-order necessary condition for the minimizer.

Lemma 23 *For any factor matrix $\mathbf{F} = \begin{bmatrix} \mathbf{L} \\ \mathbf{R} \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times r}$, suppose that the optimal alignment matrix*

$$\mathbf{Q} = \underset{\mathbf{Q} \in \text{GL}(r)}{\text{argmin}} \left\| (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\text{F}}^2$$

between \mathbf{F} and \mathbf{F}_\star exists, then \mathbf{Q} obeys

$$(\mathbf{L}\mathbf{Q})^\top (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star = \Sigma_\star (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star)^\top \mathbf{R}\mathbf{Q}^{-\top}. \quad (39)$$

Proof Expand the squares in the definition of \mathbf{Q} to obtain

$$\mathbf{Q} = \underset{\mathbf{Q} \in \text{GL}(r)}{\text{argmin}} \text{tr}((\mathbf{L}\mathbf{Q} - \mathbf{L}_\star)^\top (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star) + \text{tr}((\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star)^\top (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star) \Sigma_\star).$$

Clearly, the first order necessary condition (i.e. the gradient is zero) yields

$$2\mathbf{L}^\top (\mathbf{L}\mathbf{Q} - \mathbf{L}_\star) \Sigma_\star - 2\mathbf{Q}^{-\top} \Sigma_\star (\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_\star)^\top \mathbf{R}\mathbf{Q}^{-\top} = \mathbf{0},$$

which implies the optimal alignment criterion (39). \blacksquare

Last but not least, we connect the newly proposed distance to the usual Frobenius norm in Lemma 24, the proof of which is a slight modification to (Tu et al., 2016, Lemma 5.4) and (Ge et al., 2017, Lemma 41).

Lemma 24 *For any factor matrix $\mathbf{F} = \begin{bmatrix} \mathbf{L} \\ \mathbf{R} \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times r}$, the distance between \mathbf{F} and \mathbf{F}_\star satisfies*

$$\text{dist}(\mathbf{F}, \mathbf{F}_\star) \leq \left(\sqrt{2} + 1 \right)^{1/2} \|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_{\text{F}}.$$

Proof Suppose that $\mathbf{X} := \mathbf{L}\mathbf{R}^\top$ has compact SVD as $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. Without loss of generality, we can assume that $\mathbf{F} = \begin{bmatrix} \mathbf{U}\mathbf{\Sigma}^{1/2} \\ \mathbf{V}\mathbf{\Sigma}^{1/2} \end{bmatrix}$, since any factorization of $\mathbf{L}\mathbf{R}^\top$ yields the same distance. Introduce two auxiliary matrices $\bar{\mathbf{F}} := \begin{bmatrix} \mathbf{U}\mathbf{\Sigma}^{1/2} \\ -\mathbf{V}\mathbf{\Sigma}^{1/2} \end{bmatrix}$ and $\bar{\mathbf{F}}_\star := \begin{bmatrix} \mathbf{U}_\star\mathbf{\Sigma}_\star^{1/2} \\ -\mathbf{V}_\star\mathbf{\Sigma}_\star^{1/2} \end{bmatrix}$. Apply the dilation trick to obtain

$$2 \begin{bmatrix} \mathbf{0} & \mathbf{X} \\ \mathbf{X}^\top & \mathbf{0} \end{bmatrix} = \mathbf{F}\mathbf{F}^\top - \bar{\mathbf{F}}\bar{\mathbf{F}}^\top, \quad 2 \begin{bmatrix} \mathbf{0} & \mathbf{X}_\star \\ \mathbf{X}_\star^\top & \mathbf{0} \end{bmatrix} = \mathbf{F}_\star\mathbf{F}_\star^\top - \bar{\mathbf{F}}_\star\bar{\mathbf{F}}_\star^\top.$$

As a result, the squared Frobenius norm of $\mathbf{X} - \mathbf{X}_\star$ is given by

$$\begin{aligned} 8\|\mathbf{X} - \mathbf{X}_\star\|_{\mathbb{F}}^2 &= \|\mathbf{F}\mathbf{F}^\top - \bar{\mathbf{F}}\bar{\mathbf{F}}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top + \bar{\mathbf{F}}_\star\bar{\mathbf{F}}_\star^\top\|_{\mathbb{F}}^2 \\ &= \|\mathbf{F}\mathbf{F}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top\|_{\mathbb{F}}^2 + \|\bar{\mathbf{F}}\bar{\mathbf{F}}^\top - \bar{\mathbf{F}}_\star\bar{\mathbf{F}}_\star^\top\|_{\mathbb{F}}^2 - 2\text{tr}((\mathbf{F}\mathbf{F}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top)(\bar{\mathbf{F}}\bar{\mathbf{F}}^\top - \bar{\mathbf{F}}_\star\bar{\mathbf{F}}_\star^\top)) \\ &= 2\|\mathbf{F}\mathbf{F}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top\|_{\mathbb{F}}^2 + 2\|\mathbf{F}^\top\bar{\mathbf{F}}_\star\|_{\mathbb{F}}^2 + 2\|\mathbf{F}_\star^\top\bar{\mathbf{F}}\|_{\mathbb{F}}^2 \\ &\geq 2\|\mathbf{F}\mathbf{F}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top\|_{\mathbb{F}}^2, \end{aligned}$$

where we use the facts that $\|\mathbf{F}\mathbf{F}^\top - \mathbf{F}_\star\mathbf{F}_\star^\top\|_{\mathbb{F}}^2 = \|\bar{\mathbf{F}}\bar{\mathbf{F}}^\top - \bar{\mathbf{F}}_\star\bar{\mathbf{F}}_\star^\top\|_{\mathbb{F}}^2$ and $\mathbf{F}^\top\bar{\mathbf{F}} = \mathbf{F}_\star^\top\bar{\mathbf{F}}_\star = \mathbf{0}$.

Let $\mathbf{O} := \text{sgn}(\mathbf{F}^\top\mathbf{F}_\star)^5$ be the optimal orthonormal alignment matrix between \mathbf{F} and \mathbf{F}_\star . Denote $\mathbf{\Delta} := \mathbf{F}\mathbf{O} - \mathbf{F}_\star$. Follow the same argument as (Tu et al., 2016, Lemma 5.14) and (Ge et al., 2017, Lemma 41) to obtain

$$\begin{aligned} 4\|\mathbf{X} - \mathbf{X}_\star\|_{\mathbb{F}}^2 &\geq \|\mathbf{F}_\star\mathbf{\Delta}^\top + \mathbf{\Delta}\mathbf{F}_\star^\top + \mathbf{\Delta}\mathbf{\Delta}^\top\|_{\mathbb{F}}^2 \\ &= \text{tr}(2\mathbf{F}_\star^\top\mathbf{F}_\star\mathbf{\Delta}^\top\mathbf{\Delta} + (\mathbf{\Delta}^\top\mathbf{\Delta})^2 + 2(\mathbf{F}_\star^\top\mathbf{\Delta})^2 + 4\mathbf{F}_\star^\top\mathbf{\Delta}\mathbf{\Delta}^\top\mathbf{\Delta}) \\ &= \text{tr}\left(2\mathbf{F}_\star^\top\mathbf{F}_\star\mathbf{\Delta}^\top\mathbf{\Delta} + (\mathbf{\Delta}^\top\mathbf{\Delta} + \sqrt{2}\mathbf{F}_\star^\top\mathbf{\Delta})^2 + (4 - 2\sqrt{2})\mathbf{F}_\star^\top\mathbf{\Delta}\mathbf{\Delta}^\top\mathbf{\Delta}\right) \\ &= \text{tr}\left(2(\sqrt{2} - 1)\mathbf{F}_\star^\top\mathbf{F}_\star\mathbf{\Delta}^\top\mathbf{\Delta} + (\mathbf{\Delta}^\top\mathbf{\Delta} + \sqrt{2}\mathbf{F}_\star^\top\mathbf{\Delta})^2 + (4 - 2\sqrt{2})\mathbf{F}_\star^\top\mathbf{F}\mathbf{O}\mathbf{\Delta}^\top\mathbf{\Delta}\right) \\ &\geq \text{tr}\left(4(\sqrt{2} - 1)\mathbf{\Sigma}_\star\mathbf{\Delta}^\top\mathbf{\Delta}\right) = 4(\sqrt{2} - 1)\left\|\mathbf{F}\mathbf{O} - \mathbf{F}_\star\right\|_{\mathbb{F}}^2, \end{aligned}$$

where the last inequality follows from the facts that $\mathbf{F}_\star^\top\mathbf{F}_\star = 2\mathbf{\Sigma}_\star$ and that $\mathbf{F}_\star^\top\mathbf{F}\mathbf{O}$ is positive semi-definite. Therefore we obtain

$$\left\|\mathbf{F}\mathbf{O} - \mathbf{F}_\star\right\|_{\mathbb{F}} \leq (\sqrt{2} + 1)^{1/2} \|\mathbf{X} - \mathbf{X}_\star\|_{\mathbb{F}}.$$

This in conjunction with $\text{dist}(\mathbf{F}, \mathbf{F}_\star) \leq \|\mathbf{F}\mathbf{O} - \mathbf{F}_\star\|_{\mathbb{F}}$ yields the claimed result. \blacksquare

A.2 Matrix perturbation bounds

Lemma 25 For any $\mathbf{L} \in \mathbb{R}^{n_1 \times r}$, $\mathbf{R} \in \mathbb{R}^{n_2 \times r}$, denote $\mathbf{\Delta}_L := \mathbf{L} - \mathbf{L}_\star$ and $\mathbf{\Delta}_R := \mathbf{R} - \mathbf{R}_\star$. Suppose that $\|\mathbf{\Delta}_L\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}} \vee \|\mathbf{\Delta}_R\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}} < 1$, then one has

$$\left\|\mathbf{L}(\mathbf{L}^\top\mathbf{L})^{-1}\mathbf{\Sigma}_\star^{1/2}\right\|_{\text{op}} \leq \frac{1}{1 - \|\mathbf{\Delta}_L\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}}}; \quad (40a)$$

$$\left\|\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\mathbf{\Sigma}_\star^{1/2}\right\|_{\text{op}} \leq \frac{1}{1 - \|\mathbf{\Delta}_R\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}}}; \quad (40b)$$

5. Let $\mathbf{A}\mathbf{S}\mathbf{B}^\top$ be the SVD of $\mathbf{F}^\top\mathbf{F}_\star$, then the matrix sign is $\text{sgn}(\mathbf{F}^\top\mathbf{F}_\star) := \mathbf{A}\mathbf{B}^\top$.

$$\left\| \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Sigma}_*^{1/2} - \mathbf{U}_* \right\|_{\text{op}} \leq \frac{\sqrt{2} \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}}{1 - \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}}; \quad (40c)$$

$$\left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} - \mathbf{V}_* \right\|_{\text{op}} \leq \frac{\sqrt{2} \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}}{1 - \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}}. \quad (40d)$$

Proof We only prove claims (40a) and (40c) on the factor \mathbf{L} , while the claims on the factor \mathbf{R} follow from a similar argument. We start to prove (40a). Notice that

$$\left\| \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}} = \frac{1}{\sigma_r(\mathbf{L} \boldsymbol{\Sigma}_*^{-1/2})}.$$

In addition, invoke Weyl's inequality to obtain

$$\sigma_r(\mathbf{L} \boldsymbol{\Sigma}_*^{-1/2}) \geq \sigma_r(\mathbf{L}_* \boldsymbol{\Sigma}_*^{-1/2}) - \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} = 1 - \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}},$$

where we have used the fact that $\mathbf{U}_* = \mathbf{L}_* \boldsymbol{\Sigma}_*^{-1/2}$ satisfies $\sigma_r(\mathbf{U}_*) = 1$. Combine the preceding two relations to prove (40a).

We proceed to prove (40c). Combine $\mathbf{L}_*^\top \mathbf{U}_* = \boldsymbol{\Sigma}_*^{1/2}$ and $(\mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top) \mathbf{L} = \mathbf{0}$ to obtain the decomposition

$$\mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Sigma}_*^{1/2} - \mathbf{U}_* = -\mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Delta}_L^\top \mathbf{U}_* + (\mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}.$$

The fact that $\mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Delta}_L^\top \mathbf{U}_*$ and $(\mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}$ are orthogonal implies

$$\begin{aligned} \left\| \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Sigma}_*^{1/2} - \mathbf{U}_* \right\|_{\text{op}}^2 &\leq \left\| \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Delta}_L^\top \mathbf{U}_* \right\|_{\text{op}}^2 + \left\| (\mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2} \right\|_{\text{op}}^2 \\ &\leq \left\| \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}}^2 \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2 + \left\| \mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top \right\|_{\text{op}}^2 \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2 \\ &\leq \frac{\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2}{(1 - \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}})^2} + \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2 \\ &\leq \frac{2\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2}{(1 - \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}})^2}, \end{aligned}$$

where we have used (40a) and the fact that $\|\mathbf{I}_{n_1} - \mathbf{L}(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top\|_{\text{op}} \leq 1$ in the third line. \blacksquare

Lemma 26 For any $\mathbf{L} \in \mathbb{R}^{n_1 \times r}$, $\mathbf{R} \in \mathbb{R}^{n_2 \times r}$, denote $\boldsymbol{\Delta}_L := \mathbf{L} - \mathbf{L}_*$ and $\boldsymbol{\Delta}_R := \mathbf{R} - \mathbf{R}_*$, then one has

$$\begin{aligned} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_*\|_{\text{F}} &\leq \|\boldsymbol{\Delta}_L \mathbf{R}_*^\top\|_{\text{F}} + \|\mathbf{L}_* \boldsymbol{\Delta}_R^\top\|_{\text{F}} + \|\boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top\|_{\text{F}} \\ &\leq \left(1 + \frac{1}{2} (\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} \vee \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}) \right) \left(\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} \right). \end{aligned}$$

Proof In light of the decomposition $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_* = \boldsymbol{\Delta}_L \mathbf{R}_*^\top + \mathbf{L}_* \boldsymbol{\Delta}_R^\top + \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top$ and the triangle inequality, one has

$$\begin{aligned} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_*\|_{\text{F}} &\leq \|\boldsymbol{\Delta}_L \mathbf{R}_*^\top\|_{\text{F}} + \|\mathbf{L}_* \boldsymbol{\Delta}_R^\top\|_{\text{F}} + \|\boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top\|_{\text{F}} \\ &= \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} + \|\boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top\|_{\text{F}}, \end{aligned}$$

where we have used the facts that

$$\|\Delta_L \mathbf{R}_*^\top\|_F = \|\Delta_L \Sigma_*^{1/2} \mathbf{V}_*^\top\|_F = \|\Delta_L \Sigma_*^{1/2}\|_F, \quad \text{and} \quad \|\mathbf{L}_* \Delta_R^\top\|_F = \|\mathbf{U}_* \Sigma_*^{1/2} \Delta_R^\top\|_F = \|\Delta_R \Sigma_*^{1/2}\|_F.$$

This together with the simple upper bound

$$\begin{aligned} \|\Delta_L \Delta_R^\top\|_F &= \frac{1}{2} \|\Delta_L \Sigma_*^{1/2} (\Delta_R \Sigma_*^{-1/2})^\top\|_F + \frac{1}{2} \|\Delta_L \Sigma_*^{-1/2} (\Delta_R \Sigma_*^{1/2})^\top\|_F \\ &\leq \frac{1}{2} \|\Delta_L \Sigma_*^{1/2}\|_F \|\Delta_R \Sigma_*^{-1/2}\|_{\text{op}} + \frac{1}{2} \|\Delta_L \Sigma_*^{-1/2}\|_{\text{op}} \|\Delta_R \Sigma_*^{1/2}\|_F \\ &\leq \frac{1}{2} (\|\Delta_L \Sigma_*^{-1/2}\|_{\text{op}} \vee \|\Delta_R \Sigma_*^{-1/2}\|_{\text{op}}) (\|\Delta_L \Sigma_*^{1/2}\|_F + \|\Delta_R \Sigma_*^{1/2}\|_F) \end{aligned}$$

finishes the proof. \blacksquare

Lemma 27 For any $\mathbf{L} \in \mathbb{R}^{n_1 \times r}$, $\mathbf{R} \in \mathbb{R}^{n_2 \times r}$ and any invertible matrices $\mathbf{Q}, \bar{\mathbf{Q}} \in \text{GL}(r)$, suppose that $\|(\mathbf{L}\mathbf{Q} - \mathbf{L}_*) \Sigma_*^{-1/2}\|_{\text{op}} \vee \|(\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_*) \Sigma_*^{-1/2}\|_{\text{op}} < 1$, then one has

$$\begin{aligned} \left\| \Sigma_*^{1/2} \bar{\mathbf{Q}}^{-1} \mathbf{Q} \Sigma_*^{1/2} - \Sigma_* \right\|_{\text{op}} &\leq \frac{\|\mathbf{R}(\bar{\mathbf{Q}}^{-\top} - \mathbf{Q}^{-\top}) \Sigma_*^{1/2}\|_{\text{op}}}{1 - \|(\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_*) \Sigma_*^{-1/2}\|_{\text{op}}}; \\ \left\| \Sigma_*^{1/2} \bar{\mathbf{Q}}^\top \mathbf{Q}^{-\top} \Sigma_*^{1/2} - \Sigma_* \right\|_{\text{op}} &\leq \frac{\|\mathbf{L}(\bar{\mathbf{Q}} - \mathbf{Q}) \Sigma_*^{1/2}\|_{\text{op}}}{1 - \|(\mathbf{L}\mathbf{Q} - \mathbf{L}_*) \Sigma_*^{-1/2}\|_{\text{op}}}. \end{aligned}$$

Proof Insert $\mathbf{R}^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}$, and use the relation $\|\mathbf{A}\mathbf{B}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{op}} \|\mathbf{B}\|_{\text{op}}$ to obtain

$$\begin{aligned} \left\| \Sigma_*^{1/2} \bar{\mathbf{Q}}^{-1} \mathbf{Q} \Sigma_*^{1/2} - \Sigma_* \right\|_{\text{op}} &= \left\| \Sigma_*^{1/2} (\bar{\mathbf{Q}}^{-1} - \mathbf{Q}^{-1}) \mathbf{R}^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{Q} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &\leq \left\| \mathbf{R}(\bar{\mathbf{Q}}^{-\top} - \mathbf{Q}^{-\top}) \Sigma_*^{1/2} \right\|_{\text{op}} \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{Q} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &= \left\| \mathbf{R}(\bar{\mathbf{Q}}^{-\top} - \mathbf{Q}^{-\top}) \Sigma_*^{1/2} \right\|_{\text{op}} \left\| \mathbf{R}\mathbf{Q}^{-\top} ((\mathbf{R}\mathbf{Q}^{-\top})^\top \mathbf{R}\mathbf{Q}^{-\top})^{-1} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &\leq \frac{\|\mathbf{R}(\bar{\mathbf{Q}}^{-\top} - \mathbf{Q}^{-\top}) \Sigma_*^{1/2}\|_{\text{op}}}{1 - \|(\mathbf{R}\mathbf{Q}^{-\top} - \mathbf{R}_*) \Sigma_*^{-1/2}\|_{\text{op}}}, \end{aligned}$$

where the last line uses Lemma 25.

Similarly, insert $\mathbf{L}^\top \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1}$, and use the relation $\|\mathbf{A}\mathbf{B}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{op}} \|\mathbf{B}\|_{\text{op}}$ to obtain

$$\begin{aligned} \left\| \Sigma_*^{1/2} \bar{\mathbf{Q}}^\top \mathbf{Q}^{-\top} \Sigma_*^{1/2} - \Sigma_* \right\|_{\text{op}} &= \left\| \Sigma_*^{1/2} (\bar{\mathbf{Q}}^\top - \mathbf{Q}^\top) \mathbf{L}^\top \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{Q}^{-\top} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &\leq \left\| \mathbf{L}(\bar{\mathbf{Q}} - \mathbf{Q}) \Sigma_*^{1/2} \right\|_{\text{op}} \left\| \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{Q}^{-\top} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &= \left\| \mathbf{L}(\bar{\mathbf{Q}} - \mathbf{Q}) \Sigma_*^{1/2} \right\|_{\text{op}} \left\| \mathbf{L}\mathbf{Q} ((\mathbf{L}\mathbf{Q})^\top \mathbf{L}\mathbf{Q})^{-1} \Sigma_*^{1/2} \right\|_{\text{op}} \\ &\leq \frac{\|\mathbf{L}(\bar{\mathbf{Q}} - \mathbf{Q}) \Sigma_*^{1/2}\|_{\text{op}}}{1 - \|(\mathbf{L}\mathbf{Q} - \mathbf{L}_*) \Sigma_*^{-1/2}\|_{\text{op}}}, \end{aligned}$$

where the last line uses Lemma 25. \blacksquare

A.3 Partial Frobenius norm

We introduce the partial Frobenius norm

$$\|\mathbf{X}\|_{\text{F},r} := \sqrt{\sum_{i=1}^r \sigma_i^2(\mathbf{X})} = \|\mathcal{P}_r(\mathbf{X})\|_{\text{F}} \quad (41)$$

as the ℓ_2 norm of the vector composed of the top- r singular values of the matrix \mathbf{X} , or equivalently as the Frobenius norm of the rank- r approximation $\mathcal{P}_r(\mathbf{X})$ defined in (5). It is straightforward to verify that $\|\cdot\|_{\text{F},r}$ is a norm; see also Mazeika (2016). The following lemma provides several equivalent and useful characterizations of this partial Frobenius norm.

Lemma 28 *For any $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$, one has*

$$\|\mathbf{X}\|_{\text{F},r} = \max_{\tilde{\mathbf{V}} \in \mathbb{R}^{n_2 \times r}: \tilde{\mathbf{V}}^\top \tilde{\mathbf{V}} = \mathbf{I}_r} \|\mathbf{X}\tilde{\mathbf{V}}\|_{\text{F}} \quad (42a)$$

$$= \max_{\tilde{\mathbf{X}} \in \mathbb{R}^{n_1 \times n_2}: \|\tilde{\mathbf{X}}\|_{\text{F}} \leq 1, \text{rank}(\tilde{\mathbf{X}}) \leq r} |\langle \mathbf{X}, \tilde{\mathbf{X}} \rangle| \quad (42b)$$

$$= \max_{\tilde{\mathbf{R}} \in \mathbb{R}^{n_2 \times r}: \|\tilde{\mathbf{R}}\|_{\text{op}} \leq 1} \|\mathbf{X}\tilde{\mathbf{R}}\|_{\text{F}}. \quad (42c)$$

Proof The first representation (42a) follows immediately from the extremal partial trace identity; see (Mazeika, 2016, Proposition 4.4), by noticing the following relation

$$\sum_{i=1}^r \sigma_i^2(\mathbf{X}) = \max_{\mathbb{V} \subseteq \mathbb{R}^{n_2}: \dim(\mathbb{V})=r} \text{tr}(\mathbf{X}^\top \mathbf{X} | \mathbb{V}) = \max_{\tilde{\mathbf{V}} \in \mathbb{R}^{n_2 \times r}: \tilde{\mathbf{V}}^\top \tilde{\mathbf{V}} = \mathbf{I}_r} \|\mathbf{X}\tilde{\mathbf{V}}\|_{\text{F}}^2.$$

Here the partial trace over a vector space \mathbb{V} is defined as

$$\text{tr}(\mathbf{X}^\top \mathbf{X} | \mathbb{V}) := \sum_{i=1}^r \tilde{\mathbf{v}}_i^\top \mathbf{X}^\top \mathbf{X} \tilde{\mathbf{v}}_i,$$

where $\{\tilde{\mathbf{v}}_i\}_{1 \leq i \leq r}$ is any orthonormal basis of \mathbb{V} . The partial trace is invariant to the choice of orthonormal basis and therefore well-defined.

To prove the second representation (42b), for any $\tilde{\mathbf{X}} \in \mathbb{R}^{n_1 \times n_2}$ obeying $\text{rank}(\tilde{\mathbf{X}}) \leq r$ and $\|\tilde{\mathbf{X}}\|_{\text{F}} \leq 1$, denoting $\tilde{\mathbf{X}} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^\top$ as its compact SVD, one has

$$|\langle \mathbf{X}, \tilde{\mathbf{X}} \rangle| = |\langle \mathbf{X}, \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^\top \rangle| = |\langle \mathbf{X}\tilde{\mathbf{V}}, \tilde{\mathbf{U}}\tilde{\Sigma} \rangle| \leq \|\mathbf{X}\tilde{\mathbf{V}}\|_{\text{F}} \|\tilde{\mathbf{U}}\tilde{\Sigma}\|_{\text{F}} \leq \|\mathbf{X}\|_{\text{F},r},$$

where the last inequality follows from (42a). In addition, the maximum in (42b) is attained at $\tilde{\mathbf{X}} = \mathcal{P}_r(\mathbf{X})/\|\mathcal{P}_r(\mathbf{X})\|_{\text{F}}$.

To prove the third representation (42c), for any $\tilde{\mathbf{R}} \in \mathbb{R}^{n_2 \times r}$ obeying $\|\tilde{\mathbf{R}}\|_{\text{op}} \leq 1$, combine the variational representation of the Frobenius norm and (42b) to obtain

$$\begin{aligned} \|\mathbf{X}\tilde{\mathbf{R}}\|_{\text{F}} &= \max_{\tilde{\mathbf{L}} \in \mathbb{R}^{n_1 \times n_2}: \|\tilde{\mathbf{L}}\|_{\text{F}} \leq 1} |\langle \mathbf{X}\tilde{\mathbf{R}}, \tilde{\mathbf{L}} \rangle| \\ &= \max_{\tilde{\mathbf{L}} \in \mathbb{R}^{n_1 \times n_2}: \|\tilde{\mathbf{L}}\|_{\text{F}} \leq 1} |\langle \mathbf{X}, \tilde{\mathbf{L}}\tilde{\mathbf{R}}^\top \rangle| \leq \|\mathbf{X}\|_{\text{F},r}, \end{aligned}$$

where the last inequality follows from (42b). In addition, the maximum in (42c) is attained at $\tilde{\mathbf{R}} = \mathbf{V}$, where \mathbf{V} denotes the top- r right singular vectors of \mathbf{X} . \blacksquare

Remark 29 For self-completeness, we also provide a detailed proof of the first representation (42a). This proof is inductive on r . When $r = 1$, we have

$$\sigma_1(\mathbf{X}) = \|\mathbf{X}\mathbf{v}_1\|_2 = \max_{\tilde{\mathbf{v}} \in \mathbb{R}^{n_2}: \|\tilde{\mathbf{v}}\|_2=1} \|\mathbf{X}\tilde{\mathbf{v}}\|_2,$$

where \mathbf{v}_1 denotes the top right singular vector of \mathbf{X} . Assume that the statement holds for $\|\cdot\|_{\mathbb{F}, r-1}$. Now consider $\|\cdot\|_{\mathbb{F}, r}$. For any $\tilde{\mathbf{V}} \in \mathbb{R}^{n_2 \times r}$ such that $\tilde{\mathbf{V}}^\top \tilde{\mathbf{V}} = \mathbf{I}_r$, we can first pick $\tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r$ as a set of orthonormal vectors in the column space of $\tilde{\mathbf{V}}$ that are orthogonal to \mathbf{v}_1 , and then pick $\tilde{\mathbf{v}}_1$ via the Gram-Schmidt process, so that $\{\tilde{\mathbf{v}}_i\}_{i=1}^r$ provides an orthonormal basis of the column space of $\tilde{\mathbf{V}}$. Further, by the orthogonality of $\tilde{\mathbf{V}}$, there exists an orthonormal matrix \mathbf{O} such that

$$\tilde{\mathbf{V}} = [\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_r] \mathbf{O}.$$

Combining this formula with the induction hypothesis yields

$$\begin{aligned} \|\mathbf{X}\tilde{\mathbf{V}}\|_{\mathbb{F}}^2 &= \|\mathbf{X}[\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_r]\|_{\mathbb{F}}^2 \\ &= \|\mathbf{X}\tilde{\mathbf{v}}_1\|_2^2 + \|\mathbf{X}[\tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r]\|_{\mathbb{F}}^2 \\ &= \|\mathbf{X}\tilde{\mathbf{v}}_1\|_2^2 + \|(\mathbf{X} - \mathcal{P}_1(\mathbf{X}))[\tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r]\|_{\mathbb{F}}^2 \\ &\leq \sigma_1^2(\mathbf{X}) + \|\mathbf{X} - \mathcal{P}_1(\mathbf{X})\|_{\mathbb{F}, r-1}^2 \\ &= \sum_{i=1}^r \sigma_i^2(\mathbf{X}) = \|\mathbf{X}\|_{\mathbb{F}, r}^2, \end{aligned}$$

where the first line holds since \mathbf{O} is orthonormal, the third line holds since $\mathcal{P}_1(\mathbf{X})[\tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r] = \mathbf{0}$, the fourth line follows from the induction hypothesis, and the last line follows from the definition (41). In addition, the maximum in (42a) is attained at $\tilde{\mathbf{V}} = \mathbf{V}$, where \mathbf{V} denotes the top- r right singular vectors of \mathbf{X} . This finishes the proof.

Recall that $\mathcal{P}_r(\mathbf{X})$ denotes the best rank- r approximation of \mathbf{X} under the Frobenius norm. It turns out that $\mathcal{P}_r(\mathbf{X})$ is also the best rank- r approximation of \mathbf{X} under the partial Frobenius norm $\|\cdot\|_{\mathbb{F}, r}$. This claim is formally stated below; see also (Mazeika, 2016, Theorem 4.21).

Lemma 30 Fix any $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ and recall the definition of $\mathcal{P}_r(\mathbf{X})$ in (5). One has

$$\mathcal{P}_r(\mathbf{X}) = \underset{\tilde{\mathbf{X}} \in \mathbb{R}^{n_1 \times n_2}: \text{rank}(\tilde{\mathbf{X}}) \leq r}{\text{argmin}} \|\mathbf{X} - \tilde{\mathbf{X}}\|_{\mathbb{F}, r}.$$

Proof For any $\tilde{\mathbf{X}}$ of rank at most r , invoke Weyl's inequality to obtain $\sigma_{r+i}(\mathbf{X}) \leq \sigma_i(\mathbf{X} - \tilde{\mathbf{X}}) + \sigma_{r+1}(\tilde{\mathbf{X}}) = \sigma_i(\mathbf{X} - \tilde{\mathbf{X}})$, for $i = 1, \dots, r$. Thus one has

$$\|\mathbf{X} - \mathcal{P}_r(\mathbf{X})\|_{\mathbb{F}, r}^2 = \sum_{i=1}^r \sigma_{r+i}^2(\mathbf{X}) \leq \sum_{i=1}^r \sigma_i^2(\mathbf{X} - \tilde{\mathbf{X}}) = \|\mathbf{X} - \tilde{\mathbf{X}}\|_{\mathbb{F}, r}^2.$$

The proof is finished by observing that the rank of $\mathcal{P}_r(\mathbf{X})$ is at most r . ■

Appendix B. Proof for Low-Rank Matrix Factorization

B.1 Proof of Proposition 12

The gradients of $\mathcal{L}(\mathbf{F})$ in (29) with respect to \mathbf{L} and \mathbf{R} are given as

$$\nabla_{\mathbf{L}} \mathcal{L}(\mathbf{F}) = (\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}, \quad \nabla_{\mathbf{R}} \mathcal{L}(\mathbf{F}) = (\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)^\top \mathbf{L},$$

which can be used to compute the Hessian with respect to \mathbf{L} and \mathbf{R} . Writing for the vectorized variables, the Hessians are given as

$$\nabla_{\mathbf{L}, \mathbf{L}}^2 \mathcal{L}(\mathbf{F}) = (\mathbf{R}^\top \mathbf{R}) \otimes \mathbf{I}_{n_1}, \quad \nabla_{\mathbf{R}, \mathbf{R}}^2 \mathcal{L}(\mathbf{F}) = (\mathbf{L}^\top \mathbf{L}) \otimes \mathbf{I}_{n_2}.$$

Viewed in the vectorized form, the ScaledGD update in (3) can be rewritten as

$$\begin{aligned} \text{vec}(\mathbf{L}_{t+1}) &= \text{vec}(\mathbf{L}_t) - \eta((\mathbf{R}_t^\top \mathbf{R}_t)^{-1} \otimes \mathbf{I}_{n_1}) \text{vec}((\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*) \mathbf{R}_t) \\ &= \text{vec}(\mathbf{L}_t) - \eta(\nabla_{\mathbf{L}, \mathbf{L}}^2 \mathcal{L}(\mathbf{F}_t))^{-1} \text{vec}(\nabla_{\mathbf{L}} \mathcal{L}(\mathbf{F}_t)), \\ \text{vec}(\mathbf{R}_{t+1}) &= \text{vec}(\mathbf{R}_t) - \eta((\mathbf{L}_t^\top \mathbf{L}_t)^{-1} \otimes \mathbf{I}_{n_2}) \text{vec}((\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_*)^\top \mathbf{L}_t) \\ &= \text{vec}(\mathbf{R}_t) - \eta(\nabla_{\mathbf{R}, \mathbf{R}}^2 \mathcal{L}(\mathbf{F}_t))^{-1} \text{vec}(\nabla_{\mathbf{R}} \mathcal{L}(\mathbf{F}_t)). \end{aligned}$$

B.2 Proof of Theorem 13

The proof is inductive in nature. More specifically, we intend to show that for all $t \geq 0$,

1. $\text{dist}(\mathbf{F}_t, \mathbf{F}_*) \leq (1 - 0.7\eta)^t \text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 0.1(1 - 0.7\eta)^t \sigma_r(\mathbf{X}_*)$, and
2. the optimal alignment matrix \mathbf{Q}_t between \mathbf{F}_t and \mathbf{F}_* exists.

For the base case, i.e. $t = 0$, the first induction hypothesis trivially holds, while the second also holds true in view of Lemma 22 and the assumption that $\text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 0.1\sigma_r(\mathbf{X}_*)$. We therefore concentrate on the induction step. Suppose that the t -th iterate \mathbf{F}_t obeys the aforementioned induction hypotheses. Our goal is to show that \mathbf{F}_{t+1} continues to satisfy those.

For notational convenience, denote $\mathbf{L} := \mathbf{L}_t \mathbf{Q}_t$, $\mathbf{R} := \mathbf{R}_t \mathbf{Q}_t^{-\top}$, $\Delta_L := \mathbf{L} - \mathbf{L}_*$, $\Delta_R := \mathbf{R} - \mathbf{R}_*$, and $\epsilon := 0.1$. By the definition of $\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_*)$, one has

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_*) \leq \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \Sigma_*^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_*) \Sigma_*^{1/2} \right\|_F^2, \quad (43)$$

where we recall that \mathbf{Q}_t is the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_* . Utilize the ScaledGD update rule (30) and the decomposition $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_* = \Delta_L \mathbf{R}^\top + \mathbf{L}_* \Delta_R^\top$ to obtain

$$\begin{aligned} (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \Sigma_*^{1/2} &= (\mathbf{L} - \eta(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_*) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_*) \Sigma_*^{1/2} \\ &= (\Delta_L - \eta(\Delta_L \mathbf{R}^\top + \mathbf{L}_* \Delta_R^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}) \Sigma_*^{1/2} \\ &= (1 - \eta) \Delta_L \Sigma_*^{1/2} - \eta \mathbf{L}_* \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_*^{1/2}. \end{aligned}$$

As a result, one can expand the first square in (43) as

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \Sigma_*^{1/2} \right\|_F^2 &= (1 - \eta)^2 \text{tr}(\Delta_L \Sigma_* \Delta_L^\top) - 2\eta(1 - \eta) \underbrace{\text{tr}(\mathbf{L}_* \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_* \Delta_L^\top)}_{\mathfrak{M}_1} \\ &\quad + \eta^2 \underbrace{\left\| \mathbf{L}_* \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_*^{1/2} \right\|_F^2}_{\mathfrak{M}_2}. \end{aligned} \quad (44)$$

The first term $\text{tr}(\Delta_L \Sigma_* \Delta_L^\top)$ is closely related to $\text{dist}(\mathbf{F}_t, \mathbf{F}_*)$, and hence our focus will be on relating \mathfrak{M}_1 and \mathfrak{M}_2 to $\text{dist}(\mathbf{F}_t, \mathbf{F}_*)$. We start with the term \mathfrak{M}_1 . Since \mathbf{L} and \mathbf{R} are aligned with \mathbf{L}_* and \mathbf{R}_* , Lemma 23 tells that $\Sigma_* \Delta_L^\top \mathbf{L} = \mathbf{R}^\top \Delta_R \Sigma_*$. This together with $\mathbf{L}_* = \mathbf{L} - \Delta_L$ allows us to rewrite \mathfrak{M}_1 as

$$\mathfrak{M}_1 = \text{tr}(\mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_* \Delta_L^\top \mathbf{L}_* \Delta_R^\top)$$

$$\begin{aligned}
 &= \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \mathbf{L} \boldsymbol{\Delta}_R^\top \right) - \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right) \\
 &= \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right) - \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right).
 \end{aligned}$$

Moving on to \mathfrak{M}_2 , we can utilize the fact $\mathbf{L}_*^\top \mathbf{L}_* = \boldsymbol{\Sigma}_*$ and the decomposition $\boldsymbol{\Sigma}_* = \mathbf{R}^\top \mathbf{R} - (\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*)$ to obtain

$$\begin{aligned}
 \mathfrak{M}_2 &= \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right) \\
 &= \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right) - \text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} (\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*) (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right).
 \end{aligned}$$

Putting \mathfrak{M}_1 and \mathfrak{M}_2 back to (44) yields

$$\begin{aligned}
 \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 &= (1 - \eta)^2 \text{tr} \left(\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \right) - \eta(2 - 3\eta) \underbrace{\text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right)}_{\mathfrak{F}_1} \\
 &\quad + 2\eta(1 - \eta) \underbrace{\text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right)}_{\mathfrak{F}_2} \\
 &\quad - \eta^2 \underbrace{\text{tr} \left(\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} (\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*) (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top \right)}_{\mathfrak{F}_3}.
 \end{aligned}$$

In what follows, we will control the three terms \mathfrak{F}_1 , \mathfrak{F}_2 and \mathfrak{F}_3 separately.

1. Notice that \mathfrak{F}_1 is the inner product of two positive semi-definite matrices $\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top$ and $\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top$. Consequently we have $\mathfrak{F}_1 \geq 0$.
2. To control \mathfrak{F}_2 , we need certain control on $\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}$ and $\|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}$. The first induction hypothesis

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_*) = \sqrt{\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2} \boldsymbol{\Sigma}_*\|_{\text{F}}^2 + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2} \boldsymbol{\Sigma}_*\|_{\text{F}}^2} \leq \epsilon \sigma_r(\mathbf{X}_*)$$

together with the relation $\|\mathbf{A}\mathbf{B}\|_{\text{F}} \geq \|\mathbf{A}\|_{\text{F}} \sigma_r(\mathbf{B})$ tells that

$$\sqrt{\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{F}}^2 + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{F}}^2} \sigma_r(\mathbf{X}_*) \leq \epsilon \sigma_r(\mathbf{X}_*).$$

In light of the relation $\|\mathbf{A}\|_{\text{op}} \leq \|\mathbf{A}\|_{\text{F}}$, this further implies

$$\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} \vee \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} \leq \epsilon. \tag{45}$$

Invoke Lemma 25 to see

$$\left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}} \leq \frac{1}{1 - \epsilon}.$$

With these consequences, one can bound $|\mathfrak{F}_2|$ by

$$\begin{aligned}
 |\mathfrak{F}_2| &= \left| \text{tr} \left(\boldsymbol{\Sigma}_*^{-1/2} \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2} \right) \right| \\
 &\leq \left\| \boldsymbol{\Sigma}_*^{-1/2} \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}} \text{tr} \left(\boldsymbol{\Sigma}_*^{1/2} \boldsymbol{\Delta}_L^\top \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2} \right) \\
 &\leq \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}} \text{tr} \left(\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \right) \\
 &\leq \frac{\epsilon}{1 - \epsilon} \text{tr} \left(\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top \right).
 \end{aligned}$$

3. Similarly, one can bound $|\mathfrak{F}_3|$ by

$$\begin{aligned} |\mathfrak{F}_3| &\leq \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1}(\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*)(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \right\|_{\text{op}} \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top) \\ &\leq \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{op}}^2 \left\| \boldsymbol{\Sigma}_*^{-1/2}(\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*) \boldsymbol{\Sigma}_*^{-1/2} \right\|_{\text{op}} \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top) \\ &\leq \frac{1}{(1-\epsilon)^2} \left\| \boldsymbol{\Sigma}_*^{-1/2}(\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*) \boldsymbol{\Sigma}_*^{-1/2} \right\|_{\text{op}} \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top). \end{aligned}$$

Further notice that

$$\begin{aligned} \left\| \boldsymbol{\Sigma}_*^{-1/2}(\mathbf{R}^\top \mathbf{R} - \boldsymbol{\Sigma}_*) \boldsymbol{\Sigma}_*^{-1/2} \right\|_{\text{op}} &= \left\| \boldsymbol{\Sigma}_*^{-1/2}(\mathbf{R}_*^\top \boldsymbol{\Delta}_R + \boldsymbol{\Delta}_R^\top \mathbf{R}_* + \boldsymbol{\Delta}_R^\top \boldsymbol{\Delta}_R) \boldsymbol{\Sigma}_*^{-1/2} \right\|_{\text{op}} \\ &\leq 2 \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}} + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{-1/2}\|_{\text{op}}^2 \\ &\leq 2\epsilon + \epsilon^2. \end{aligned}$$

Take the preceding two bounds together to arrive at

$$|\mathfrak{F}_3| \leq \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top).$$

Combining the bounds for $\mathfrak{F}_1, \mathfrak{F}_2, \mathfrak{F}_3$, one has

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 &= \left\| (1-\eta) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2} - \eta \mathbf{L}_* \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 \\ &\leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta(1-\eta) \right) \text{tr}(\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top). \end{aligned} \quad (46)$$

A similarly bound holds for the second square $\|(\mathbf{R}_{t+1} \mathbf{Q}_t - \mathbf{R}_*) \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}}^2$ in (43). Therefore we obtain

$$\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 \leq \rho^2(\eta; \epsilon) \text{dist}^2(\mathbf{F}_t, \mathbf{F}_*),$$

where we identify

$$\text{dist}^2(\mathbf{F}_t, \mathbf{F}_*) = \text{tr}(\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_L^\top) + \text{tr}(\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_* \boldsymbol{\Delta}_R^\top) \quad (47)$$

and the contraction rate $\rho^2(\eta; \epsilon)$ is given by

$$\rho^2(\eta; \epsilon) := (1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta(1-\eta) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2.$$

With $\epsilon = 0.1$ and $0 < \eta \leq 2/3$, one has $\rho(\eta; \epsilon) \leq 1 - 0.7\eta$. Thus we conclude that

$$\begin{aligned} \text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_*) &\leq \sqrt{\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{\text{F}}^2} \\ &\leq (1 - 0.7\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_*) \\ &\leq (1 - 0.7\eta)^{t+1} \text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq (1 - 0.7\eta)^{t+1} 0.1 \sigma_r(\mathbf{X}_*). \end{aligned}$$

This proves the first induction hypothesis. The existence of the optimal alignment matrix \mathbf{Q}_{t+1} between \mathbf{F}_{t+1} and \mathbf{F}_* is assured by Lemma 22, which finishes the proof for the second hypothesis.

So far, we have demonstrated the first conclusion in the theorem. The second conclusion is an easy consequence of Lemma 26 as

$$\begin{aligned} \left\| \mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_* \right\|_{\text{F}} &\leq \left(1 + \frac{\epsilon}{2}\right) \left(\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} + \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{1/2}\|_{\text{F}} \right) \\ &\leq \left(1 + \frac{\epsilon}{2}\right) \sqrt{2} \text{dist}(\mathbf{F}_t, \mathbf{F}_*) \\ &\leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_*). \end{aligned} \quad (48)$$

Here, the second line follows from the elementary inequality $a + b \leq \sqrt{2(a^2 + b^2)}$ and the expression of $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ in (47). The proof is now completed.

Appendix C. Proof for Low-Rank Matrix Sensing

We start by recording a useful lemma.

Lemma 31 (Candès and Plan (2011)) *Suppose that $\mathcal{A}(\cdot)$ obeys the $2r$ -RIP with a constant δ_{2r} . Then for any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^{n_1 \times n_2}$ of rank at most r , one has*

$$|\langle \mathcal{A}(\mathbf{X}_1), \mathcal{A}(\mathbf{X}_2) \rangle - \langle \mathbf{X}_1, \mathbf{X}_2 \rangle| \leq \delta_{2r} \|\mathbf{X}_1\|_F \|\mathbf{X}_2\|_F,$$

which can be stated equivalently as

$$|\text{tr}((\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_1) \mathbf{X}_2^\top)| \leq \delta_{2r} \|\mathbf{X}_1\|_F \|\mathbf{X}_2\|_F. \quad (49)$$

As a simple corollary, one has that for any matrix $\mathbf{R} \in \mathbb{R}^{n_2 \times r}$:

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_1) \mathbf{R}\|_F \leq \delta_{2r} \|\mathbf{X}_1\|_F \|\mathbf{R}\|_{\text{op}}. \quad (50)$$

This is due to the fact that

$$\begin{aligned} \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_1) \mathbf{R}\|_F &= \max_{\tilde{\mathbf{L}}: \|\tilde{\mathbf{L}}\|_F \leq 1} \text{tr}((\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_1) \mathbf{R} \tilde{\mathbf{L}}^\top) \\ &\leq \max_{\tilde{\mathbf{L}}: \|\tilde{\mathbf{L}}\|_F \leq 1} \delta_{2r} \|\mathbf{X}_1\|_F \|\tilde{\mathbf{L}} \mathbf{R}^\top\|_F \\ &\leq \delta_{2r} \|\mathbf{X}_1\|_F \|\mathbf{R}\|_{\text{op}}. \end{aligned}$$

Here, the first line follows from the variational representation of the Frobenius norm, the second line follows from (49), and the last line follows from the relation $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_{\text{op}}$.

C.1 Proof of Lemma 14

The proof mostly mirrors that in Section B.2. First, in view of the condition $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.1\sigma_r(\mathbf{X}_\star)$ and Lemma 22, one knows that \mathbf{Q}_t , the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star exists. Therefore, for notational convenience, denote $\mathbf{L} := \mathbf{L}_t \mathbf{Q}_t$, $\mathbf{R} := \mathbf{R}_t \mathbf{Q}_t^{-\top}$, $\Delta_L := \mathbf{L} - \mathbf{L}_\star$, $\Delta_R := \mathbf{R} - \mathbf{R}_\star$, and $\epsilon := 0.1$. Similar to the derivation in (45), we have

$$\|\Delta_L \Sigma_\star^{-1/2}\|_{\text{op}} \vee \|\Delta_R \Sigma_\star^{-1/2}\|_{\text{op}} \leq \epsilon. \quad (51)$$

The conclusion $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ is a simple consequence of Lemma 26; see (48) for a detailed argument. From now on, we focus on proving the distance contraction.

With these notations in place, we have by the definition of $\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star)$ that

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2. \quad (52)$$

Apply the update rule (15) and the decomposition $\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star = \Delta_L \mathbf{R}^\top + \mathbf{L}_\star \Delta_R^\top$ to obtain

$$\begin{aligned} (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} &= (\mathbf{L} - \eta \mathcal{A}^* \mathcal{A}(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_\star) \Sigma_\star^{1/2} \\ &= (\Delta_L - \eta(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} - \eta(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1}) \Sigma_\star^{1/2} \\ &= (1 - \eta) \Delta_L \Sigma_\star^{1/2} - \eta \mathbf{L}_\star \Delta_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} - \eta(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2}. \end{aligned}$$

This allows us to expand the first square in (52) as

$$\begin{aligned}
 \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2 &= \underbrace{\left\| (1-\eta) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} - \eta \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2}_{\mathfrak{S}_1} \\
 &\quad - 2\eta(1-\eta) \underbrace{\text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right)}_{\mathfrak{S}_2} \\
 &\quad + 2\eta^2 \underbrace{\text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \boldsymbol{\Delta}_R \mathbf{L}_\star^\top \right)}_{\mathfrak{S}_3} \\
 &\quad + \eta^2 \underbrace{\left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2}_{\mathfrak{S}_4}.
 \end{aligned}$$

In what follows, we shall control the four terms separately, of which \mathfrak{S}_1 is the main term, and \mathfrak{S}_2 , \mathfrak{S}_3 and \mathfrak{S}_4 are perturbation terms.

1. Notice that the main term \mathfrak{S}_1 has already been controlled in (46) under the condition (51). It obeys

$$\mathfrak{S}_1 \leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta(1-\eta) \right) \left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \left\| \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2.$$

2. For the second term \mathfrak{S}_2 , decompose $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star = \boldsymbol{\Delta}_L \mathbf{R}_\star^\top + \mathbf{L}_\star \boldsymbol{\Delta}_R^\top + \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top$ and apply the triangle inequality to obtain

$$\begin{aligned}
 |\mathfrak{S}_2| &= \left| \text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\boldsymbol{\Delta}_L \mathbf{R}_\star^\top + \mathbf{L}_\star \boldsymbol{\Delta}_R^\top + \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right) \right| \\
 &\leq \left| \text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\boldsymbol{\Delta}_L \mathbf{R}_\star^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right) \right| \\
 &\quad + \left| \text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\mathbf{L}_\star \boldsymbol{\Delta}_R^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right) \right| \\
 &\quad + \left| \text{tr} \left((\mathcal{A}^* \mathcal{A} - \mathcal{I}) (\boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right) \right|.
 \end{aligned}$$

Invoke Lemma 31 to further obtain

$$\begin{aligned}
 |\mathfrak{S}_2| &\leq \delta_{2r} \left(\left\| \boldsymbol{\Delta}_L \mathbf{R}_\star^\top \right\|_{\mathbb{F}} + \left\| \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} + \left\| \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} \right) \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star \boldsymbol{\Delta}_L^\top \right\|_{\mathbb{F}} \\
 &\leq \delta_{2r} \left(\left\| \boldsymbol{\Delta}_L \mathbf{R}_\star^\top \right\|_{\mathbb{F}} + \left\| \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} + \left\| \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} \right) \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}} \left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}},
 \end{aligned}$$

where the second line follows from the relation $\|\mathbf{A}\mathbf{B}\|_{\mathbb{F}} \leq \|\mathbf{A}\|_{\text{op}} \|\mathbf{B}\|_{\mathbb{F}}$. Take the condition (51) and Lemmas 25 and 26 together to obtain

$$\begin{aligned}
 \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}} &\leq \frac{1}{1-\epsilon}; \\
 \left\| \boldsymbol{\Delta}_L \mathbf{R}_\star^\top \right\|_{\mathbb{F}} + \left\| \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} + \left\| \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} &\leq \left(1 + \frac{\epsilon}{2}\right) \left(\left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} + \left\| \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} \right).
 \end{aligned}$$

These consequences further imply that

$$\begin{aligned}
 |\mathfrak{S}_2| &\leq \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)} \left(\left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} + \left\| \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} \right) \left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} \\
 &= \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)} \left(\left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} \left\| \boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}} \right).
 \end{aligned}$$

For the term $\|\Delta_L \Sigma_\star^{1/2}\|_F \|\Delta_R \Sigma_\star^{1/2}\|_F$, we can apply the elementary inequality $2ab \leq a^2 + b^2$ to see

$$\|\Delta_L \Sigma_\star^{1/2}\|_F \|\Delta_R \Sigma_\star^{1/2}\|_F \leq \frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{1}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2.$$

The preceding two bounds taken collectively yield

$$|\mathfrak{S}_2| \leq \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)} \left(\frac{3}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 \right).$$

3. The third term \mathfrak{S}_3 can be similarly bounded as

$$\begin{aligned} |\mathfrak{S}_3| &\leq \delta_{2r} \left(\|\Delta_L \mathbf{R}_\star^\top\|_F + \|\mathbf{L}_\star \Delta_R^\top\|_F + \|\Delta_L \Delta_R^\top\|_F \right) \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \mathbf{L}_\star^\top \right\|_F \\ &\leq \delta_{2r} \left(\|\Delta_L \mathbf{R}_\star^\top\|_F + \|\mathbf{L}_\star \Delta_R^\top\|_F + \|\Delta_L \Delta_R^\top\|_F \right) \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \|\Delta_R \mathbf{L}_\star^\top\|_F \\ &\leq \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)^2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_R \Sigma_\star^{1/2}\|_F \\ &\leq \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)^2} \left(\frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{3}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right). \end{aligned}$$

4. We are then left with the last term \mathfrak{S}_4 , for which we have

$$\begin{aligned} \sqrt{\mathfrak{S}_4} &= \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F \\ &\leq \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F \\ &\quad + \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{L}_\star \Delta_R^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F \\ &\quad + \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\Delta_L \Delta_R^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F, \end{aligned}$$

where once again we use the decomposition $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star = \Delta_L \mathbf{R}_\star^\top + \mathbf{L}_\star \Delta_R^\top + \Delta_L \Delta_R^\top$. Use (50) to see that

$$\sqrt{\mathfrak{S}_4} \leq \delta_{2r} \left(\|\Delta_L \mathbf{R}_\star^\top\|_F + \|\mathbf{L}_\star \Delta_R^\top\|_F + \|\Delta_L \Delta_R^\top\|_F \right) \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}.$$

Repeating the same argument in bounding \mathfrak{S}_2 yields

$$\sqrt{\mathfrak{S}_4} \leq \frac{\delta_{2r}(2+\epsilon)}{2(1-\epsilon)} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right).$$

We can then take the squares of both sides and use $(a+b)^2 \leq 2a^2 + 2b^2$ to reach

$$\mathfrak{S}_4 \leq \frac{\delta_{2r}^2(2+\epsilon)^2}{2(1-\epsilon)^2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right).$$

Taking the bounds for $\mathfrak{S}_1, \mathfrak{S}_2, \mathfrak{S}_3, \mathfrak{S}_4$ collectively yields

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 &\leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta(1-\eta) \right) \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \\ &\quad + \frac{\delta_{2r}(2+\epsilon)}{1-\epsilon} \eta(1-\eta) \left(\frac{3}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{1}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right) \\ &\quad + \frac{\delta_{2r}^2(2+\epsilon)}{(1-\epsilon)^2} \eta^2 \left(\frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{3}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right) \end{aligned}$$

$$+ \frac{\delta_{2r}^2(2+\epsilon)^2}{2(1-\epsilon)^2} \eta^2 \left(\|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right).$$

Similarly, we can expand the second square in (52) and obtain a similar bound. Combine both to obtain

$$\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2 \leq \rho^2(\eta; \epsilon, \delta_{2r}) \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star),$$

where the contraction rate is given by

$$\rho^2(\eta; \epsilon, \delta_{2r}) := (1-\eta)^2 + \frac{2\epsilon + \delta_{2r}(4+2\epsilon)}{1-\epsilon} \eta(1-\eta) + \frac{2\epsilon + \epsilon^2 + \delta_{2r}(4+2\epsilon) + \delta_{2r}^2(2+\epsilon)^2}{(1-\epsilon)^2} \eta^2.$$

With $\epsilon = 0.1$, $\delta_{2r} \leq 0.02$, and $0 < \eta \leq 2/3$, one has $\rho(\eta; \epsilon, \delta_{2r}) \leq 1 - 0.6\eta$. Thus we conclude that

$$\begin{aligned} \text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) &\leq \sqrt{\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2} \\ &\leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star). \end{aligned}$$

C.2 Proof of Lemma 15

With the knowledge of partial Frobenius norm $\|\cdot\|_{F,r}$, we are ready to establish the claimed result. Invoke Lemma 24 to relate $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star)$ to $\|\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_\star\|_F$, and use that $\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_\star$ has rank at most $2r$ to obtain

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq \sqrt{\sqrt{2}+1} \|\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_\star\|_F \leq \sqrt{2(\sqrt{2}+1)} \|\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_\star\|_{F,r}.$$

Note that $\mathbf{L}_0 \mathbf{R}_0^\top$ is the best rank- r approximation of $\mathcal{A}^* \mathcal{A}(\mathbf{X}_\star)$, and apply the triangle inequality combined with Lemma 30 to obtain

$$\begin{aligned} \|\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_\star\|_{F,r} &\leq \|\mathcal{A}^* \mathcal{A}(\mathbf{X}_\star) - \mathbf{L}_0 \mathbf{R}_0^\top\|_{F,r} + \|\mathcal{A}^* \mathcal{A}(\mathbf{X}_\star) - \mathbf{X}_\star\|_{F,r} \\ &\leq 2 \|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_\star)\|_{F,r} \leq 2\delta_{2r} \|\mathbf{X}_\star\|_F. \end{aligned}$$

Here, the last inequality follows from combining Lemma 28 and (50) as

$$\|(\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_\star)\|_{F,r} = \max_{\tilde{\mathbf{R}} \in \mathbb{R}^{n_2 \times r}: \|\tilde{\mathbf{R}}\|_{\text{op}} \leq 1} \left\| (\mathcal{A}^* \mathcal{A} - \mathcal{I})(\mathbf{X}_\star) \tilde{\mathbf{R}} \right\|_F \leq \delta_{2r} \|\mathbf{X}_\star\|_F.$$

As a result, one has

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 2\sqrt{2(\sqrt{2}+1)} \delta_{2r} \|\mathbf{X}_\star\|_F \leq 5\delta_{2r} \sqrt{r} \kappa \sigma_r(\mathbf{X}_\star).$$

Appendix D. Proof for Robust PCA

We first establish a useful property regarding the truncation operator $\mathcal{T}_{2\alpha}[\cdot]$.

Lemma 32 *Given $\mathbf{S}_\star \in \mathcal{S}_\alpha$ and $\mathbf{S} = \mathcal{T}_{2\alpha}[\mathbf{X}_\star + \mathbf{S}_\star - \mathbf{L}\mathbf{R}^\top]$, one has*

$$\|\mathbf{S} - \mathbf{S}_\star\|_\infty \leq 2\|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_\infty. \quad (53)$$

In addition, for any low-rank matrix $\mathbf{M} = \mathbf{L}_M \mathbf{R}_M^\top \in \mathbb{R}^{n_1 \times n_2}$ with $\mathbf{L}_M \in \mathbb{R}^{n_1 \times r}$, $\mathbf{R}_M \in \mathbb{R}^{n_2 \times r}$, one has

$$\begin{aligned} |\langle \mathbf{S} - \mathbf{S}_\star, \mathbf{M} \rangle| &\leq \sqrt{3\alpha\nu} \left(\|(\mathbf{L} - \mathbf{L}_\star) \Sigma_\star^{1/2}\|_F + \|(\mathbf{R} - \mathbf{R}_\star) \Sigma_\star^{1/2}\|_F \right) \|\mathbf{M}\|_F \\ &\quad + 2\sqrt{\alpha} (\sqrt{n_1} \|\mathbf{L}_M\|_{2,\infty} \|\mathbf{R}_M\|_F \wedge \sqrt{n_2} \|\mathbf{L}_M\|_F \|\mathbf{R}_M\|_{2,\infty}) \|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_F, \end{aligned} \quad (54)$$

where ν obeys

$$\nu \geq \frac{\sqrt{n_1}}{2} \left(\|\mathbf{L}\Sigma_\star^{-1/2}\|_{2,\infty} + \|\mathbf{L}_\star\Sigma_\star^{-1/2}\|_{2,\infty} \right) \vee \frac{\sqrt{n_2}}{2} \left(\|\mathbf{R}\Sigma_\star^{-1/2}\|_{2,\infty} + \|\mathbf{R}_\star\Sigma_\star^{-1/2}\|_{2,\infty} \right).$$

Proof Denote $\Delta_L := \mathbf{L} - \mathbf{L}_\star$, $\Delta_R := \mathbf{R} - \mathbf{R}_\star$, and $\Delta_X := \mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star$. Let Ω, Ω_\star be the support of \mathbf{S} and \mathbf{S}_\star , respectively. As a result, $\mathbf{S} - \mathbf{S}_\star$ is supported on $\Omega \cup \Omega_\star$.

We start with proving the first claim, i.e. (53). For $(i, j) \in \Omega$, by the definition of $\mathcal{T}_{2\alpha}[\cdot]$, we have $(\mathbf{S} - \mathbf{S}_\star)_{i,j} = (-\Delta_X)_{i,j}$. For $(i, j) \in \Omega_\star \setminus \Omega$, one necessarily has $\mathbf{S}_{i,j} = 0$ and therefore $(\mathbf{S} - \mathbf{S}_\star)_{i,j} = (-\mathbf{S}_\star)_{i,j}$. Again by the definition of the operator $\mathcal{T}_{2\alpha}[\cdot]$, we know $|\mathbf{S}_\star - \Delta_X|_{i,j}$ is either smaller than $|\mathbf{S}_\star - \Delta_X|_{i,(2\alpha n_2)}$ or $|\mathbf{S}_\star - \Delta_X|_{(2\alpha n_1),j}$. Furthermore, we know that \mathbf{S}_\star contains at most α -fraction nonzero entries per row and column. Consequently, one has $|\mathbf{S}_\star - \Delta_X|_{i,j} \leq |\Delta_X|_{i,(\alpha n_2)} \vee |\Delta_X|_{(\alpha n_1),j}$. Combining the two cases above, we conclude that

$$|\mathbf{S} - \mathbf{S}_\star|_{i,j} \leq \begin{cases} |\Delta_X|_{i,j}, & (i, j) \in \Omega \\ |\Delta_X|_{i,j} + (|\Delta_X|_{i,(\alpha n_2)} \vee |\Delta_X|_{(\alpha n_1),j}), & (i, j) \in \Omega_\star \setminus \Omega \end{cases}. \quad (55)$$

This immediately implies the ℓ_∞ norm bound (53).

Next, we prove the second claim (54). Recall that $\mathbf{S} - \mathbf{S}_\star$ is supported on $\Omega \cup \Omega_\star$. We then have

$$\begin{aligned} |\langle \mathbf{S} - \mathbf{S}_\star, \mathbf{M} \rangle| &\leq \sum_{(i,j) \in \Omega} |\mathbf{S} - \mathbf{S}_\star|_{i,j} |\mathbf{M}|_{i,j} + \sum_{(i,j) \in \Omega_\star \setminus \Omega} |\mathbf{S} - \mathbf{S}_\star|_{i,j} |\mathbf{M}|_{i,j} \\ &\leq \sum_{(i,j) \in \Omega \cup \Omega_\star} |\Delta_X|_{i,j} |\mathbf{M}|_{i,j} + \sum_{(i,j) \in \Omega_\star \setminus \Omega} (|\Delta_X|_{i,(\alpha n_2)} + |\Delta_X|_{(\alpha n_1),j}) |\mathbf{M}|_{i,j}, \end{aligned}$$

where the second line follows from (55). Let $\beta > 0$ be some positive number, whose value will be determined later. Use $2ab \leq \beta^{-1}a^2 + \beta b^2$ to further obtain

$$|\langle \mathbf{S} - \mathbf{S}_\star, \mathbf{M} \rangle| \leq \underbrace{\sum_{(i,j) \in \Omega \cup \Omega_\star} |\Delta_X|_{i,j} |\mathbf{M}|_{i,j}}_{\mathfrak{A}_1} + \frac{1}{2\beta} \underbrace{\sum_{(i,j) \in \Omega_\star \setminus \Omega} (|\Delta_X|_{i,(\alpha n_2)}^2 + |\Delta_X|_{(\alpha n_1),j}^2)}_{\mathfrak{A}_2} + \beta \underbrace{\sum_{(i,j) \in \Omega_\star \setminus \Omega} |\mathbf{M}|_{i,j}^2}_{\mathfrak{A}_3}.$$

In regard to the three terms $\mathfrak{A}_1, \mathfrak{A}_2$ and \mathfrak{A}_3 , we have the following claims, whose proofs are deferred to the end.

Claim 1 *The first term \mathfrak{A}_1 satisfies*

$$\mathfrak{A}_1 \leq \sqrt{3\alpha\nu} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right) \|\mathbf{M}\|_{\mathbb{F}}.$$

Claim 2 *The second term \mathfrak{A}_2 satisfies*

$$\mathfrak{A}_2 \leq 2\|\Delta_X\|_{\mathbb{F}}^2.$$

Claim 3 *The third term \mathfrak{A}_3 satisfies*

$$\mathfrak{A}_3 \leq \alpha \left(n_1 \|\mathbf{L}_M\|_{2,\infty}^2 \|\mathbf{R}_M\|_{\mathbb{F}}^2 \wedge n_2 \|\mathbf{L}_M\|_{\mathbb{F}}^2 \|\mathbf{R}_M\|_{2,\infty}^2 \right).$$

Combine the pieces to reach

$$\begin{aligned} |\langle \mathbf{S} - \mathbf{S}_\star, \mathbf{M} \rangle| &\leq \sqrt{3\alpha\nu} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right) \|\mathbf{M}\|_{\mathbb{F}} \\ &\quad + \frac{\|\Delta_X\|_{\mathbb{F}}^2}{\beta} + \beta\alpha \left(n_1 \|\mathbf{L}_M\|_{2,\infty}^2 \|\mathbf{R}_M\|_{\mathbb{F}}^2 \wedge n_2 \|\mathbf{L}_M\|_{\mathbb{F}}^2 \|\mathbf{R}_M\|_{2,\infty}^2 \right). \end{aligned}$$

One can then choose β optimally to yield

$$\begin{aligned} |\langle \mathbf{S} - \mathbf{S}_*, \mathbf{M} \rangle| &\leq \sqrt{3\alpha\nu} \left(\|\Delta_L \Sigma_*^{1/2}\|_F + \|\Delta_R \Sigma_*^{1/2}\|_F \right) \|\mathbf{M}\|_F \\ &\quad + 2\sqrt{\alpha} (\sqrt{n_1} \|\mathbf{L}_M\|_{2,\infty} \|\mathbf{R}_M\|_F \wedge \sqrt{n_2} \|\mathbf{L}_M\|_F \|\mathbf{R}_M\|_{2,\infty}) \|\Delta_X\|_F. \end{aligned}$$

This finishes the proof. \blacksquare

Proof [Proof of Claim 1] Use the decomposition $\Delta_X = \Delta_L \mathbf{R}^\top + \mathbf{L}_* \Delta_R^\top = \Delta_L \mathbf{R}_*^\top + \mathbf{L} \Delta_R^\top$ to obtain

$$\begin{aligned} |\Delta_X|_{i,j} &\leq \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2 \|\mathbf{R} \Sigma_*^{-1/2}\|_{2,\infty} + \|\mathbf{L}_* \Sigma_*^{-1/2}\|_{2,\infty} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2, \quad \text{and} \\ |\Delta_X|_{i,j} &\leq \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2 \|\mathbf{R}_* \Sigma_*^{-1/2}\|_{2,\infty} + \|\mathbf{L} \Sigma_*^{-1/2}\|_{2,\infty} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2. \end{aligned}$$

Take the average to yield

$$|\Delta_X|_{i,j} \leq \frac{\nu}{\sqrt{n_2}} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2 + \frac{\nu}{\sqrt{n_1}} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2,$$

where we have used the assumption on ν . With this upper bound on $|\Delta_X|_{i,j}$ in place, we can further control \mathfrak{A}_1 as

$$\begin{aligned} \mathfrak{A}_1 &\leq \sum_{(i,j) \in \Omega \cup \Omega_*} \frac{\nu}{\sqrt{n_2}} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2 |\mathbf{M}|_{i,j} + \sum_{(i,j) \in \Omega \cup \Omega_*} \frac{\nu}{\sqrt{n_1}} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2 |\mathbf{M}|_{i,j} \\ &\leq \left(\sqrt{\sum_{(i,j) \in \Omega \cup \Omega_*} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2^2 / n_2} + \sqrt{\sum_{(i,j) \in \Omega \cup \Omega_*} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2^2 / n_1} \right) \nu \|\mathbf{M}\|_F. \end{aligned}$$

Regarding the first term, one has

$$\begin{aligned} \sum_{(i,j) \in \Omega \cup \Omega_*} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2^2 &= \sum_{i=1}^{n_1} \sum_{j:(i,j) \in \Omega \cup \Omega_*} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2^2 \\ &\leq 3\alpha n_2 \sum_{i=1}^{n_1} \|(\Delta_L \Sigma_*^{1/2})_{i,\cdot}\|_2^2 \\ &= 3\alpha n_2 \|\Delta_L \Sigma_*^{1/2}\|_F^2, \end{aligned}$$

where the second line follows from the fact that $\Omega \cup \Omega_*$ contains at most $3\alpha n_2$ non-zero entries in each row. Similarly, we can show that

$$\sum_{(i,j) \in \Omega \cup \Omega_*} \|(\Delta_R \Sigma_*^{1/2})_{j,\cdot}\|_2^2 \leq 3\alpha n_1 \|\Delta_R \Sigma_*^{1/2}\|_F^2.$$

In all, we arrive at

$$\mathfrak{A}_1 \leq \sqrt{3\alpha\nu} \left(\|\Delta_L \Sigma_*^{1/2}\|_F + \|\Delta_R \Sigma_*^{1/2}\|_F \right) \|\mathbf{M}\|_F,$$

which is the desired claim. \blacksquare

Proof [Proof of Claim 2] Recall that $(\Delta_X)_{i,(\alpha n_2)}$ denotes the (αn_2) -th largest entry in the i -th row of Δ_X . One necessarily has

$$\alpha n_2 |\Delta_X|_{i,(\alpha n_2)}^2 \leq \|(\Delta_X)_{i,\cdot}\|_2^2.$$

As a result, we obtain

$$\begin{aligned}
\sum_{(i,j) \in \Omega_\star \setminus \Omega} |\Delta_X|_{i,(\alpha n_2)}^2 &\leq \sum_{(i,j) \in \Omega_\star} |\Delta_X|_{i,(\alpha n_2)}^2 \\
&\leq \sum_{i=1}^{n_1} \sum_{j:(i,j) \in \Omega_\star} \frac{\|(\Delta_X)_{i,\cdot}\|_2^2}{\alpha n_2} \\
&\leq \sum_{i=1}^{n_1} \|(\Delta_X)_{i,\cdot}\|_2^2 = \|\Delta_X\|_F^2,
\end{aligned}$$

where the last line follows from the fact that Ω_\star contains at most αn_2 nonzero entries in each row. Similarly one can show that

$$\sum_{(i,j) \in \Omega_\star \setminus \Omega} |\Delta_X|_{(\alpha n_1),j}^2 \leq \|\Delta_X\|_F^2.$$

Combining the above two bounds with the definition of \mathfrak{A}_2 completes the proof. \blacksquare

Proof [Proof of Claim 3] By definition, $\mathbf{M} = \mathbf{L}_M \mathbf{R}_M^\top$, and hence one has

$$\mathfrak{A}_3 = \sum_{(i,j) \in \Omega_\star \setminus \Omega} |(\mathbf{L}_M)_{i,\cdot} (\mathbf{R}_M)_{j,\cdot}^\top|^2 \leq \sum_{(i,j) \in \Omega_\star} |(\mathbf{L}_M)_{i,\cdot} (\mathbf{R}_M)_{j,\cdot}^\top|^2.$$

We can further upper bound \mathfrak{A}_3 as

$$\begin{aligned}
\mathfrak{A}_3 &\leq \sum_{(i,j) \in \Omega_\star} \|(\mathbf{L}_M)_{i,\cdot}\|_2^2 \|(\mathbf{R}_M)_{j,\cdot}\|_2^2 \\
&\leq \sum_{i=1}^{n_1} \sum_{j:(i,j) \in \Omega_\star} \|(\mathbf{L}_M)_{i,\cdot}\|_2^2 \|\mathbf{R}_M\|_{2,\infty}^2 \\
&\leq \sum_{i=1}^{n_1} \alpha n_2 \|(\mathbf{L}_M)_{i,\cdot}\|_2^2 \|\mathbf{R}_M\|_{2,\infty}^2 = \alpha n_2 \|\mathbf{L}_M\|_F^2 \|\mathbf{R}_M\|_{2,\infty}^2,
\end{aligned}$$

where the last line follows from the fact that Ω_\star contains at most αn_2 non-zero entries in each row. Similarly, one can obtain

$$\mathfrak{A}_3 \leq \alpha n_1 \|\mathbf{L}_M\|_{2,\infty}^2 \|\mathbf{R}_M\|_F^2,$$

which completes the proof. \blacksquare

D.1 Proof of Lemma 16

We begin with introducing several useful notations and facts. In view of the condition $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$ and Lemma 22, one knows that \mathbf{Q}_t , the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star exists. Therefore, for notational convenience, denote $\mathbf{L} := \mathbf{L}_t \mathbf{Q}_t$, $\mathbf{R} := \mathbf{R}_t \mathbf{Q}_t^{-\top}$, $\Delta_L := \mathbf{L} - \mathbf{L}_\star$, $\Delta_R := \mathbf{R} - \mathbf{R}_\star$, $\mathbf{S} := \mathbf{S}_t = \mathcal{T}_{2\alpha}[\mathbf{X}_\star + \mathbf{S}_\star - \mathbf{L}\mathbf{R}^\top]$, and $\epsilon := 0.02$. Similar to the derivation in (45), we have

$$\|\Delta_L \Sigma_\star^{-1/2}\|_{\text{op}} \vee \|\Delta_R \Sigma_\star^{-1/2}\|_{\text{op}} \leq \epsilon. \quad (56)$$

Moreover, the incoherence condition

$$\sqrt{n_1} \|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\Delta_R \Sigma_\star^{1/2}\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_\star) \quad (57)$$

implies

$$\sqrt{n_1} \|\Delta_L \Sigma_\star^{-1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\Delta_R \Sigma_\star^{-1/2}\|_{2,\infty} \leq \sqrt{\mu r}, \quad (58)$$

which combined with the triangle inequality further implies

$$\sqrt{n_1} \|\mathbf{L} \Sigma_\star^{-1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \leq 2\sqrt{\mu r}. \quad (59)$$

The conclusion $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ is a simple consequence of Lemma 26; see (48) for a detailed argument. In what follows, we shall prove the distance contraction and the incoherence condition separately.

D.1.1 DISTANCE CONTRACTION

By the definition of $\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_\star)$, one has

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2. \quad (60)$$

From now on, we focus on controlling the first square $\|(\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2}\|_F^2$. In view of the update rule (20), one has

$$\begin{aligned} (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} &= (\mathbf{L} - \eta(\mathbf{L} \mathbf{R}^\top + \mathbf{S} - \mathbf{X}_\star - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_\star) \Sigma_\star^{1/2} \\ &= (\Delta_L - \eta(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} - \eta(\mathbf{S} - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}) \Sigma_\star^{1/2} \\ &= (1 - \eta) \Delta_L \Sigma_\star^{1/2} - \eta \mathbf{L}_\star \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} - \eta(\mathbf{S} - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2}. \end{aligned} \quad (61)$$

Here, we use the notation introduced above and the decomposition $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star = \Delta_L \mathbf{R}^\top + \mathbf{L}_\star \Delta_R^\top$. Take the squared Frobenius norm of both sides of (61) to obtain

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 &= \underbrace{\left\| (1 - \eta) \Delta_L \Sigma_\star^{1/2} - \eta \mathbf{L}_\star \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F^2}_{\mathfrak{R}_1} \\ &\quad - 2\eta(1 - \eta) \underbrace{\text{tr} \left((\mathbf{S} - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star \Delta_L^\top \right)}_{\mathfrak{R}_2} \\ &\quad + 2\eta^2 \underbrace{\text{tr} \left((\mathbf{S} - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \mathbf{L}_\star^\top \right)}_{\mathfrak{R}_3} \\ &\quad + \eta^2 \underbrace{\left\| (\mathbf{S} - \mathbf{S}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F^2}_{\mathfrak{R}_4}. \end{aligned}$$

In the sequel, we shall bound the four terms separately, of which \mathfrak{R}_1 is the main term, and $\mathfrak{R}_2, \mathfrak{R}_3$ and \mathfrak{R}_4 are perturbation terms.

1. Notice that the main term \mathfrak{R}_1 has already been controlled in (46) under the condition (56). It obeys

$$\mathfrak{R}_1 \leq \left((1 - \eta)^2 + \frac{2\epsilon}{1 - \epsilon} \eta(1 - \eta) \right) \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{2\epsilon + \epsilon^2}{(1 - \epsilon)^2} \eta^2 \|\Delta_R \Sigma_\star^{1/2}\|_F^2.$$

2. For the second term \mathfrak{R}_2 , set $\mathbf{M} := \Delta_L \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top$ with $\mathbf{L}_M := \Delta_L \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2}$, $\mathbf{R}_M := \mathbf{R} \Sigma_\star^{-1/2}$, and then invoke Lemma 32 with $\nu := 3\sqrt{\mu r}/2$ to see

$$\begin{aligned} |\mathfrak{R}_2| &\leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_L \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top\|_F \\ &\quad + 2\sqrt{\alpha n_2} \left\| \Delta_L \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F \\ &\leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_L \Sigma_\star^{1/2}\|_F \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \\ &\quad + 2\sqrt{\alpha n_2} \|\Delta_L \Sigma_\star^{1/2}\|_F \left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F. \end{aligned}$$

Take the condition (56) and Lemmas 25 and 26 together to obtain

$$\begin{aligned} \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} &\leq \frac{1}{1-\epsilon}; \\ \left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} &= \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \leq \frac{1}{(1-\epsilon)^2}; \\ \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F &\leq \left(1 + \frac{\epsilon}{2}\right) \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right). \end{aligned} \tag{62}$$

These consequences combined with the condition (59) yield

$$\begin{aligned} |\mathfrak{R}_2| &\leq \frac{3\sqrt{3\alpha\mu r}}{2(1-\epsilon)} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_L \Sigma_\star^{1/2}\|_F \\ &\quad + \frac{4\sqrt{\alpha\mu r}}{(1-\epsilon)^2} \|\Delta_L \Sigma_\star^{1/2}\|_F \left(1 + \frac{\epsilon}{2}\right) \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \\ &\leq \sqrt{\alpha\mu r} \frac{3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon}}{2(1-\epsilon)} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_L \Sigma_\star^{1/2}\|_F \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \\ &\leq \sqrt{\alpha\mu r} \frac{3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon}}{2(1-\epsilon)} \left(\frac{3}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{1}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right), \end{aligned}$$

where the last inequality holds since $2ab \leq a^2 + b^2$.

3. The third term \mathfrak{R}_3 can be controlled similarly. Set $\mathbf{M} := \mathbf{L}_\star \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top$ with $\mathbf{L}_M := \mathbf{L}_\star \Sigma_\star^{-1/2}$ and $\mathbf{R}_M := \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \Sigma_\star^{1/2}$, and invoke Lemma 32 with $\nu := 3\sqrt{\mu r}/2$ to arrive at

$$\begin{aligned} |\mathfrak{R}_3| &\leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\mathbf{L}_\star \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top\|_F \\ &\quad + 2\sqrt{\alpha n_1} \|\mathbf{L}_\star \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \Sigma_\star^{1/2} \right\|_F \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F \\ &\leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_R \Sigma_\star^{1/2}\|_F \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \\ &\quad + 2\sqrt{\alpha n_1} \|\mathbf{L}_\star \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \|\Delta_R \Sigma_\star^{1/2}\|_F \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F. \end{aligned}$$

Use the consequences (62) again to obtain

$$|\mathfrak{R}_3| \leq \frac{3\sqrt{3\alpha\mu r}}{2(1-\epsilon)^2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \|\Delta_R \Sigma_\star^{1/2}\|_F$$

$$\begin{aligned}
 & + \frac{2\sqrt{\alpha\mu r}}{(1-\epsilon)^2} \|\Delta_R \Sigma_\star^{1/2}\|_F (1 + \frac{\epsilon}{2}) \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \\
 & \leq \sqrt{\alpha\mu r} \frac{3\sqrt{3} + 2(2+\epsilon)}{2(1-\epsilon)^2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F \|\Delta_R \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right) \\
 & \leq \sqrt{\alpha\mu r} \frac{3\sqrt{3} + 2(2+\epsilon)}{2(1-\epsilon)^2} \left(\frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{3}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right).
 \end{aligned}$$

4. For the last term \mathfrak{R}_4 , utilize the variational representation of the Frobenius norm to see

$$\sqrt{\mathfrak{R}_4} = \text{tr} \left((S - S_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \tilde{\mathbf{L}}^\top \right)$$

for some $\tilde{\mathbf{L}} \in \mathbb{R}^{n_1 \times r}$ obeying $\|\tilde{\mathbf{L}}\|_F = 1$. Setting $\mathbf{M} := \tilde{\mathbf{L}} \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top = \mathbf{L}_M \mathbf{R}_M^\top$ with $\mathbf{L}_M := \tilde{\mathbf{L}} \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2}$ and $\mathbf{R}_M := \mathbf{R} \Sigma_\star^{-1/2}$, we are ready to apply Lemma 32 again with $\nu := 3\sqrt{\mu r}/2$ to see

$$\begin{aligned}
 \sqrt{\mathfrak{R}_4} & \leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \left\| \tilde{\mathbf{L}} \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \right\|_F \\
 & \quad + 2\sqrt{\alpha n_2} \left\| \tilde{\mathbf{L}} \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_F \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F \\
 & \leq \frac{3}{2} \sqrt{3\alpha\mu r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right) \left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \\
 & \quad + 2\sqrt{\alpha n_2} \left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star\|_F.
 \end{aligned}$$

This combined with the consequences (62) and condition (59) yields

$$\sqrt{\mathfrak{R}_4} \leq \sqrt{\alpha\mu r} \frac{3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon}}{2(1-\epsilon)} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F + \|\Delta_R \Sigma_\star^{1/2}\|_F \right).$$

Take the square, and use the elementary inequality $(a+b)^2 \leq 2a^2 + 2b^2$ to reach

$$\mathfrak{R}_4 \leq \alpha\mu r \frac{(3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon})^2}{2(1-\epsilon)^2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right).$$

Taking collectively the bounds for $\mathfrak{R}_1, \mathfrak{R}_2, \mathfrak{R}_3$ and \mathfrak{R}_4 yields the control of $\|(\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2}\|_F^2$ as

$$\begin{aligned}
 \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 & \leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta(1-\eta) \right) \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \\
 & \quad + \sqrt{\alpha\mu r} \frac{3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon}}{1-\epsilon} \eta(1-\eta) \left(\frac{3}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{1}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right) \\
 & \quad + \sqrt{\alpha\mu r} \frac{3\sqrt{3} + 2(2+\epsilon)}{(1-\epsilon)^2} \eta^2 \left(\frac{1}{2} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{3}{2} \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right) \\
 & \quad + \alpha\mu r \frac{(3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon})^2}{2(1-\epsilon)^2} \eta^2 \left(\|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_R \Sigma_\star^{1/2}\|_F^2 \right).
 \end{aligned}$$

Similarly, we can obtain the control of $\|(\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2}\|_F^2$. Combine them together and identify $\text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star) = \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \|\Delta_R \Sigma_\star^{1/2}\|_F^2$ to reach

$$\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2 \leq \rho^2(\eta; \epsilon, \alpha\mu r) \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star),$$

where the contraction rate $\rho^2(\eta; \epsilon, \alpha\mu r)$ is given by

$$\begin{aligned} \rho^2(\eta; \epsilon, \alpha\mu r) &:= (1 - \eta)^2 + \frac{2\epsilon + \sqrt{\alpha\mu r}(6\sqrt{3} + \frac{8(2+\epsilon)}{1-\epsilon})}{1 - \epsilon} \eta(1 - \eta) \\ &\quad + \frac{2\epsilon + \epsilon^2 + \sqrt{\alpha\mu r}(6\sqrt{3} + 4(2 + \epsilon)) + \alpha\mu r(3\sqrt{3} + \frac{4(2+\epsilon)}{1-\epsilon})^2}{(1 - \epsilon)^2} \eta^2. \end{aligned}$$

With $\epsilon = 0.02$, $\alpha\mu r \leq 10^{-4}$, and $0 < \eta \leq 2/3$, one has $\rho(\eta; \epsilon, \alpha\mu r) \leq 1 - 0.6\eta$. Thus we conclude that

$$\begin{aligned} \text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) &\leq \sqrt{\left\| (\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{F}}^2 + \left\| (\mathbf{R}_{t+1}\mathbf{Q}_t^{-\top} - \mathbf{R}_\star)\boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{F}}^2} \\ &\leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star). \end{aligned} \tag{63}$$

D.1.2 INCOHERENCE CONDITION

We start by controlling the term $\|(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}$. We know from (61) that

$$(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\boldsymbol{\Sigma}_\star^{1/2} = (1 - \eta)\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2} - \eta\mathbf{L}_\star\boldsymbol{\Delta}_R^\top\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} - \eta(\mathbf{S} - \mathbf{S}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2}.$$

Apply the triangle inequality to obtain

$$\begin{aligned} \left\| (\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\boldsymbol{\Sigma}_\star^{1/2} \right\|_{2,\infty} &\leq (1 - \eta)\|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} + \eta \underbrace{\left\| \mathbf{L}_\star\boldsymbol{\Delta}_R^\top\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} \right\|_{2,\infty}}_{\mathfrak{T}_1} \\ &\quad + \eta \underbrace{\left\| (\mathbf{S} - \mathbf{S}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} \right\|_{2,\infty}}_{\mathfrak{T}_2}. \end{aligned}$$

The first term $\|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}$ follows from the incoherence condition (57) as

$$\|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n_1}} \sigma_r(\mathbf{X}_\star).$$

In the sequel, we shall bound the terms \mathfrak{T}_1 and \mathfrak{T}_2 .

1. For the term \mathfrak{T}_1 , use the relation $\|\mathbf{A}\mathbf{B}\|_{2,\infty} \leq \|\mathbf{A}\|_{2,\infty}\|\mathbf{B}\|_{\text{op}}$, and combine the condition (56) with the consequences (62) to obtain

$$\begin{aligned} \mathfrak{T}_1 &\leq \|\mathbf{L}_\star\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty} \left\| \boldsymbol{\Sigma}_\star^{1/2}\boldsymbol{\Delta}_R^\top\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}} \\ &\leq \|\mathbf{L}_\star\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{\text{op}} \left\| \mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}} \\ &\leq \frac{\epsilon}{1 - \epsilon} \sqrt{\frac{\mu r}{n_1}} \sigma_r(\mathbf{X}_\star), \end{aligned}$$

2. For the term \mathfrak{T}_2 , use the relation $\|\mathbf{A}\mathbf{B}\|_{2,\infty} \leq \|\mathbf{A}\|_{2,\infty}\|\mathbf{B}\|_{\text{op}}$ to obtain

$$\mathfrak{T}_2 \leq \|\mathbf{S} - \mathbf{S}_\star\|_{2,\infty} \left\| \mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}}.$$

We know from Lemma 32 that $\mathbf{S} - \mathbf{S}_\star$ has at most $3\alpha n_2$ non-zero entries in each row, and $\|\mathbf{S} - \mathbf{S}_\star\|_\infty \leq 2\|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_\infty$. Upper bound the $\ell_{2,\infty}$ norm by the ℓ_∞ norm as

$$\|\mathbf{S} - \mathbf{S}_\star\|_{2,\infty} \leq \sqrt{3\alpha n_2} \|\mathbf{S} - \mathbf{S}_\star\|_\infty \leq 2\sqrt{3\alpha n_2} \|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_\infty.$$

Split $\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star = \mathbf{\Delta}_L\mathbf{R}^\top + \mathbf{L}_\star\mathbf{\Delta}_R^\top$, and take the conditions (57) and (59) to obtain

$$\begin{aligned} \|\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star\|_\infty &\leq \|\mathbf{\Delta}_L\mathbf{R}^\top\|_\infty + \|\mathbf{L}_\star\mathbf{\Delta}_R^\top\|_\infty \\ &\leq \|\mathbf{\Delta}_L\mathbf{\Sigma}_\star^{1/2}\|_{2,\infty}\|\mathbf{R}\mathbf{\Sigma}_\star^{-1/2}\|_{2,\infty} + \|\mathbf{L}_\star\mathbf{\Sigma}_\star^{-1/2}\|_{2,\infty}\|\mathbf{\Delta}_R\mathbf{\Sigma}_\star^{1/2}\|_{2,\infty} \\ &\leq \sqrt{\frac{\mu r}{n_1}}\sigma_r(\mathbf{X}_\star)2\sqrt{\frac{\mu r}{n_2}} + \sqrt{\frac{\mu r}{n_1}}\sqrt{\frac{\mu r}{n_2}}\sigma_r(\mathbf{X}_\star) \\ &= \frac{3\mu r}{\sqrt{n_1 n_2}}\sigma_r(\mathbf{X}_\star). \end{aligned}$$

This combined with the consequences (62) yields

$$\mathfrak{T}_2 \leq \frac{6\sqrt{3\alpha\mu r}}{1-\epsilon}\sqrt{\frac{\mu r}{n_1}}\sigma_r(\mathbf{X}_\star).$$

Taking collectively the bounds for $\mathfrak{T}_1, \mathfrak{T}_2$ yields the control

$$\left\|(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} \leq \left(1 - \eta + \frac{\epsilon + 6\sqrt{3\alpha\mu r}}{1-\epsilon}\eta\right)\sqrt{\frac{\mu r}{n_1}}\sigma_r(\mathbf{X}_\star). \quad (64)$$

The last step is to switch the alignment matrix from \mathbf{Q}_t to \mathbf{Q}_{t+1} . (63) together with Lemma 22 demonstrates the existence of \mathbf{Q}_{t+1} . Apply the triangle inequality to obtain

$$\begin{aligned} \left\|(\mathbf{L}_{t+1}\mathbf{Q}_{t+1} - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} &\leq \left\|(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} + \left\|\mathbf{L}_{t+1}(\mathbf{Q}_{t+1} - \mathbf{Q}_t)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} \\ &\leq \left\|(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} + \|\mathbf{L}_{t+1}\mathbf{Q}_t\mathbf{\Sigma}_\star^{-1/2}\|_{2,\infty}\left\|\mathbf{\Sigma}_\star^{1/2}\mathbf{Q}_t^{-1}\mathbf{Q}_{t+1}\mathbf{\Sigma}_\star^{1/2} - \mathbf{\Sigma}_\star\right\|_{\text{op}}. \end{aligned}$$

We deduct from (64) that

$$\|\mathbf{L}_{t+1}\mathbf{Q}_t\mathbf{\Sigma}_\star^{-1/2}\|_{2,\infty} \leq \|\mathbf{L}_\star\mathbf{\Sigma}_\star^{-1/2}\|_{2,\infty} + \left\|(\mathbf{L}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{-1/2}\right\|_{2,\infty} \leq \left(2 - \eta + \frac{\epsilon + 6\sqrt{3\alpha\mu r}}{1-\epsilon}\eta\right)\sqrt{\frac{\mu r}{n_1}}.$$

Regarding the alignment matrix term, invoke Lemma 27 to obtain

$$\begin{aligned} \left\|\mathbf{\Sigma}_\star^{1/2}\mathbf{Q}_t^{-1}\mathbf{Q}_{t+1}\mathbf{\Sigma}_\star^{1/2} - \mathbf{\Sigma}_\star\right\|_{\text{op}} &\leq \frac{\|(\mathbf{R}_{t+1}(\mathbf{Q}_t^{-\top} - \mathbf{Q}_{t+1}^{-\top})\mathbf{\Sigma}_\star^{1/2})\|_{\text{op}}}{1 - \|(\mathbf{R}_{t+1}\mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}}} \\ &\leq \frac{\|(\mathbf{R}_{t+1}\mathbf{Q}_t^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{1/2}\|_{\text{op}} + \|(\mathbf{R}_{t+1}\mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{1/2}\|_{\text{op}}}{1 - \|(\mathbf{R}_{t+1}\mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}}} \\ &\leq \frac{2\epsilon}{1-\epsilon}\sigma_r(\mathbf{X}_\star), \end{aligned}$$

where we deduct from (63) that the distances using either \mathbf{Q}_t or \mathbf{Q}_{t+1} are bounded by

$$\begin{aligned} \|(\mathbf{R}_{t+1}\mathbf{Q}_t^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{1/2}\|_{\text{op}} &\leq \epsilon\sigma_r(\mathbf{X}_\star); \\ \|(\mathbf{R}_{t+1}\mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{1/2}\|_{\text{op}} &\leq \epsilon\sigma_r(\mathbf{X}_\star); \\ \|(\mathbf{R}_{t+1}\mathbf{Q}_{t+1}^{-\top} - \mathbf{R}_\star)\mathbf{\Sigma}_\star^{-1/2}\|_{\text{op}} &\leq \epsilon. \end{aligned}$$

Combine all pieces to reach

$$\left\|(\mathbf{L}_{t+1}\mathbf{Q}_{t+1} - \mathbf{L}_\star)\mathbf{\Sigma}_\star^{1/2}\right\|_{2,\infty} \leq \left(\frac{1+\epsilon}{1-\epsilon}\left(1 - \eta + \frac{\epsilon + 6\sqrt{3\alpha\mu r}}{1-\epsilon}\eta\right) + \frac{2\epsilon}{1-\epsilon}\right)\sqrt{\frac{\mu r}{n_1}}\sigma_r(\mathbf{X}_\star).$$

With $\epsilon = 0.02$, $\alpha\mu r \leq 10^{-4}$, and $0.1 \leq \eta \leq 2/3$, we get the desired incoherence condition

$$\left\| (\mathbf{L}_{t+1} \mathbf{Q}_{t+1} - \mathbf{L}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n_1}} \sigma_r(\mathbf{X}_*).$$

Similarly, we can prove the other part

$$\left\| (\mathbf{R}_{t+1} \mathbf{Q}_{t+1}^\top - \mathbf{R}_*) \boldsymbol{\Sigma}_*^{1/2} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n_2}} \sigma_r(\mathbf{X}_*).$$

D.2 Proof of Lemma 17

We first record two lemmas from Yi et al. (2016), which are useful for studying the properties of the initialization.

Lemma 33 ((Yi et al., 2016, Section 6.1)) *Given $\mathbf{S}_* \in \mathcal{S}_\alpha$, one has $\|\mathbf{S}_* - \mathcal{T}_\alpha[\mathbf{X}_* + \mathbf{S}_*]\|_\infty \leq 2\|\mathbf{X}_*\|_\infty$.*

Lemma 34 ((Yi et al., 2016, Lemma 1)) *For any matrix $\mathbf{M} \in \mathcal{S}_\alpha$, one has $\|\mathbf{M}\|_{\text{op}} \leq \alpha\sqrt{n_1 n_2} \|\mathbf{M}\|_\infty$.*

With these two lemmas in place, we are ready to establish the claimed result. Invoke Lemma 24 to obtain

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq \sqrt{\sqrt{2} + 1} \left\| \mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_* \right\|_{\text{F}} \leq \sqrt{(\sqrt{2} + 1)2r} \left\| \mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_* \right\|_{\text{op}},$$

where the last relation uses the fact that $\mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_*$ has rank at most $2r$. We can further apply the triangle inequality to see

$$\begin{aligned} \left\| \mathbf{L}_0 \mathbf{R}_0^\top - \mathbf{X}_* \right\|_{\text{op}} &\leq \left\| \mathbf{Y} - \mathcal{T}_\alpha[\mathbf{Y}] - \mathbf{L}_0 \mathbf{R}_0^\top \right\|_{\text{op}} + \left\| \mathbf{Y} - \mathcal{T}_\alpha[\mathbf{Y}] - \mathbf{X}_* \right\|_{\text{op}} \\ &\leq 2 \left\| \mathbf{Y} - \mathcal{T}_\alpha[\mathbf{Y}] - \mathbf{X}_* \right\|_{\text{op}} = 2 \left\| \mathbf{S}_* - \mathcal{T}_\alpha[\mathbf{X}_* + \mathbf{S}_*] \right\|_{\text{op}}. \end{aligned}$$

Here the second inequality hinges on the fact that $\mathbf{L}_0 \mathbf{R}_0^\top$ is the best rank- r approximation of $\mathbf{Y} - \mathcal{T}_\alpha[\mathbf{Y}]$, and the last identity arises from $\mathbf{Y} = \mathbf{X}_* + \mathbf{S}_*$. Follow the same argument as (Yi et al., 2016, Section 6.1), combining Lemmas 33 and 34 to reach

$$\begin{aligned} \left\| \mathbf{S}_* - \mathcal{T}_\alpha[\mathbf{X}_* + \mathbf{S}_*] \right\|_{\text{op}} &\leq 2\alpha\sqrt{n_1 n_2} \left\| \mathbf{S}_* - \mathcal{T}_\alpha[\mathbf{X}_* + \mathbf{S}_*] \right\|_\infty \\ &\leq 4\alpha\sqrt{n_1 n_2} \|\mathbf{X}_*\|_\infty \leq 4\alpha\mu r \kappa \sigma_r(\mathbf{X}_*), \end{aligned}$$

where the last inequality follows from the incoherence assumption

$$\|\mathbf{X}_*\|_\infty \leq \|\mathbf{U}_*\|_{2,\infty} \|\boldsymbol{\Sigma}_*\|_{\text{op}} \|\mathbf{V}_*\|_{2,\infty} \leq \frac{\mu r}{\sqrt{n_1 n_2}} \kappa \sigma_r(\mathbf{X}_*). \quad (65)$$

Take the above inequalities together to arrive at

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 8\sqrt{2(\sqrt{2} + 1)} \alpha\mu r^{3/2} \kappa \sigma_r(\mathbf{X}_*) \leq 20\alpha\mu r^{3/2} \kappa \sigma_r(\mathbf{X}_*).$$

D.3 Proof of Lemma 18

In view of the condition $\text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 0.02\sigma_r(\mathbf{X}_*)$ and Lemma 22, one knows that \mathbf{Q}_0 , the optimal alignment matrix between \mathbf{F}_0 and \mathbf{F}_* exists. Therefore, for notational convenience, denote $\mathbf{L} := \mathbf{L}_0 \mathbf{Q}_0$, $\mathbf{R} := \mathbf{R}_0 \mathbf{Q}_0^\top$, $\boldsymbol{\Delta}_L := \mathbf{L} - \mathbf{L}_*$, $\boldsymbol{\Delta}_R := \mathbf{R} - \mathbf{R}_*$, and $\epsilon := 0.02$. Our objective is then translated to demonstrate

$$\sqrt{n_1} \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_*^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_*^{1/2}\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_*).$$

From now on, we focus on bounding $\|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty}$. Since $U_0 \Sigma_0 V_0^\top$ is the top- r SVD of $Y - \mathcal{T}_\alpha[Y]$, and recall that $Y = X_\star + S_\star$, we have the relation

$$(X_\star + S_\star - \mathcal{T}_\alpha[X_\star + S_\star])V_0 = U_0 \Sigma_0,$$

which further implies the following decomposition of $\Delta_L \Sigma_\star^{1/2}$.

Claim 4 *One has*

$$\Delta_L \Sigma_\star^{1/2} = (S_\star - \mathcal{T}_\alpha[X_\star + S_\star])R(R^\top R)^{-1}\Sigma_\star^{1/2} - L_\star \Delta_R^\top R(R^\top R)^{-1}\Sigma_\star^{1/2}.$$

Combining Claim 4 with the triangle inequality yields

$$\|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty} \leq \underbrace{\|L_\star \Delta_R^\top R(R^\top R)^{-1}\Sigma_\star^{1/2}\|_{2,\infty}}_{\mathfrak{J}_1} + \underbrace{\|(S_\star - \mathcal{T}_\alpha[X_\star + S_\star])R(R^\top R)^{-1}\Sigma_\star^{1/2}\|_{2,\infty}}_{\mathfrak{J}_2}.$$

In what follows, we shall control \mathfrak{J}_1 and \mathfrak{J}_2 in turn.

1. For the term \mathfrak{J}_1 , use the relation $\|AB\|_{2,\infty} \leq \|A\|_{2,\infty}\|B\|_{\text{op}}$ to obtain

$$\mathfrak{J}_1 \leq \|L_\star \Sigma_\star^{-1/2}\|_{2,\infty} \|\Delta_R \Sigma_\star^{1/2}\|_{\text{op}} \left\| R(R^\top R)^{-1}\Sigma_\star^{1/2} \right\|_{\text{op}}.$$

The incoherence assumption tells $\|L_\star \Sigma_\star^{-1/2}\|_{2,\infty} = \|U_\star\|_{2,\infty} \leq \sqrt{\mu r/n_1}$. In addition, the assumption $\text{dist}(F_0, F_\star) \leq \epsilon \sigma_r(X_\star)$ entails the bound $\|\Delta_R \Sigma_\star^{1/2}\|_{\text{op}} \leq \epsilon \sigma_r(X_\star)$. Finally, repeating the argument for obtaining (56) yields $\|\Delta_R \Sigma_\star^{-1/2}\|_{\text{op}} \leq \epsilon$, which together with Lemma 25 reveals

$$\left\| R(R^\top R)^{-1}\Sigma_\star^{1/2} \right\|_{\text{op}} \leq \frac{1}{1-\epsilon}.$$

In all, we arrive at

$$\mathfrak{J}_1 \leq \frac{\epsilon}{1-\epsilon} \sqrt{\frac{\mu r}{n_1}} \sigma_r(X_\star).$$

2. Proceeding to the term \mathfrak{J}_2 , use the relations $\|AB\|_{2,\infty} \leq \|A\|_{1,\infty}\|B\|_{2,\infty}$ and $\|AB\|_{2,\infty} \leq \|A\|_{2,\infty}\|B\|_{\text{op}}$ to obtain

$$\begin{aligned} \mathfrak{J}_2 &\leq \|S_\star - \mathcal{T}_\alpha[X_\star + S_\star]\|_{1,\infty} \left\| R(R^\top R)^{-1}\Sigma_\star^{1/2} \right\|_{2,\infty} \\ &\leq \|S_\star - \mathcal{T}_\alpha[X_\star + S_\star]\|_{1,\infty} \|R \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \Sigma_\star^{1/2} (R^\top R)^{-1}\Sigma_\star^{1/2} \right\|_{\text{op}}. \end{aligned}$$

Regarding $S_\star - \mathcal{T}_\alpha[X_\star + S_\star]$, Lemma 33 tells that $S_\star - \mathcal{T}_\alpha[X_\star + S_\star]$ has at most $2\alpha n_2$ non-zero entries in each row, and $\|S_\star - \mathcal{T}_\alpha[X_\star + S_\star]\|_\infty \leq 2\|X_\star\|_\infty$. Consequently, we can upper bound the $\ell_{1,\infty}$ norm by the ℓ_∞ norm as

$$\begin{aligned} \|S_\star - \mathcal{T}_\alpha[X_\star + S_\star]\|_{1,\infty} &\leq 2\alpha n_2 \|S_\star - \mathcal{T}_\alpha[X_\star + S_\star]\|_\infty \\ &\leq 4\alpha n_2 \|X_\star\|_\infty \\ &\leq 4\alpha n_2 \frac{\mu r}{\sqrt{n_1 n_2}} \kappa \sigma_r(X_\star). \end{aligned}$$

Here the last inequality follows from the incoherence assumption (65). For the term $\|\mathbf{R}\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty}$, one can apply the triangle inequality to see

$$\|\mathbf{R}\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty} \leq \|\mathbf{R}_\star\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty} + \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{-1/2}\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n_2}} + \frac{\|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}}{\sigma_r(\mathbf{X}_\star)}.$$

Last but not least, repeat the argument for (62) to obtain

$$\left\| \boldsymbol{\Sigma}_\star^{1/2}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}} = \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\text{op}}^2 \leq \frac{1}{(1-\epsilon)^2}.$$

Taking together the above bounds yields

$$\mathfrak{J}_2 \leq \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \sqrt{\frac{\mu r}{n_1}} \sigma_r(\mathbf{X}_\star) + \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \sqrt{\frac{n_2}{n_1}} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}.$$

Combine the bounds on \mathfrak{J}_1 and \mathfrak{J}_2 to reach

$$\sqrt{n_1} \|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \leq \left(\frac{\epsilon}{1-\epsilon} + \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \right) \sqrt{\mu r} \sigma_r(\mathbf{X}_\star) + \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \sqrt{n_2} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}.$$

Similarly, we have

$$\sqrt{n_2} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \leq \left(\frac{\epsilon}{1-\epsilon} + \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \right) \sqrt{\mu r} \sigma_r(\mathbf{X}_\star) + \frac{4\alpha\mu r\kappa}{(1-\epsilon)^2} \sqrt{n_1} \|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}.$$

Taking the maximum and solving for $\sqrt{n_1} \|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty}$ yield the relation

$$\sqrt{n_1} \|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \leq \frac{\epsilon(1-\epsilon) + 4\alpha\mu r\kappa}{(1-\epsilon)^2 - 4\alpha\mu r\kappa} \sqrt{\mu r} \sigma_r(\mathbf{X}_\star).$$

With $\epsilon = 0.02$ and $\alpha\mu r\kappa \leq 0.1$, we get the desired conclusion

$$\sqrt{n_1} \|\boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\boldsymbol{\Delta}_R\boldsymbol{\Sigma}_\star^{1/2}\|_{2,\infty} \leq \sqrt{\mu r} \sigma_r(\mathbf{X}_\star).$$

Proof [Proof of Claim 4] Identify \mathbf{U}_0 (resp. \mathbf{V}_0) with $\mathbf{L}_0\boldsymbol{\Sigma}_0^{-1/2}$ (resp. $\mathbf{R}_0\boldsymbol{\Sigma}_0^{-1/2}$) to yield

$$(\mathbf{X}_\star + \mathbf{S}_\star - \mathcal{T}_\alpha[\mathbf{X}_\star + \mathbf{S}_\star])\mathbf{R}_0\boldsymbol{\Sigma}_0^{-1} = \mathbf{L}_0,$$

which is equivalent to $(\mathbf{X}_\star + \mathbf{S}_\star - \mathcal{T}_\alpha[\mathbf{X}_\star + \mathbf{S}_\star])\mathbf{R}_0(\mathbf{R}_0^\top \mathbf{R}_0)^{-1} = \mathbf{L}_0$ since $\boldsymbol{\Sigma}_0 = \mathbf{R}_0^\top \mathbf{R}_0$. Multiply both sides by $\mathbf{Q}_0\boldsymbol{\Sigma}_\star^{1/2}$ to obtain

$$(\mathbf{X}_\star + \mathbf{S}_\star - \mathcal{T}_\alpha[\mathbf{X}_\star + \mathbf{S}_\star])\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} = \mathbf{L}\boldsymbol{\Sigma}_\star^{1/2},$$

where we recall that $\mathbf{L} = \mathbf{L}_0\mathbf{Q}_0$ and $\mathbf{R} = \mathbf{R}_0\mathbf{Q}_0^{-\top}$. In the end, subtract $\mathbf{X}_\star\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2}$ from both sides to reach

$$\begin{aligned} (\mathbf{S}_\star - \mathcal{T}_\alpha[\mathbf{X}_\star + \mathbf{S}_\star])\mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} &= \mathbf{L}\boldsymbol{\Sigma}_\star^{1/2} - \mathbf{L}_\star\mathbf{R}_\star^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \\ &= (\mathbf{L} - \mathbf{L}_\star)\boldsymbol{\Sigma}_\star^{1/2} + \mathbf{L}_\star(\mathbf{R} - \mathbf{R}_\star)^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \\ &= \boldsymbol{\Delta}_L\boldsymbol{\Sigma}_\star^{1/2} + \mathbf{L}_\star\boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2}. \end{aligned}$$

This finishes the proof. ■

Appendix E. Proof for Matrix Completion

E.1 New projection operator

E.1.1 PROOF OF PROPOSITION 7

First, notice that the optimization of \mathbf{L} and \mathbf{R} in (23) can be decomposed and done in parallel, hence we focus on the optimization of \mathbf{L} below:

$$\mathbf{L} = \operatorname{argmin}_{\mathbf{L} \in \mathbb{R}^{n_1 \times r}} \left\| (\mathbf{L} - \tilde{\mathbf{L}})(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2} \right\|_{\mathbb{F}}^2 \quad \text{s.t.} \quad \sqrt{n_1} \left\| \mathbf{L}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2} \right\|_{2, \infty} \leq B.$$

By a change of variables as $\mathbf{G} := \mathbf{L}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2}$ and $\tilde{\mathbf{G}} := \tilde{\mathbf{L}}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{1/2}$, we rewrite the above problem equivalently as

$$\mathbf{G} = \operatorname{argmin}_{\mathbf{G} \in \mathbb{R}^{n_1 \times r}} \left\| \mathbf{G} - \tilde{\mathbf{G}} \right\|_{\mathbb{F}}^2 \quad \text{s.t.} \quad \sqrt{n_1} \left\| \mathbf{G} \right\|_{2, \infty} \leq B,$$

whose solution is given as Chen and Wainwright (2015)

$$\mathbf{G}_{i,\cdot} = \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{G}}_{i,\cdot}\|_2} \right) \tilde{\mathbf{G}}_{i,\cdot}, \quad 1 \leq i \leq n_1.$$

By applying again the change of variable $\mathbf{L} = \mathbf{G}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1/2}$ and $\tilde{\mathbf{L}} = \tilde{\mathbf{G}}(\tilde{\mathbf{R}}^\top \tilde{\mathbf{R}})^{-1/2}$, we obtain the claimed solution.

E.1.2 PROOF OF LEMMA 19

We begin with proving the non-expansiveness property. Denote the optimal alignment matrix between $\tilde{\mathbf{F}}$ and \mathbf{F}_\star as $\tilde{\mathbf{Q}}$, whose existence is guaranteed by Lemma 22. Denoting $\mathcal{P}_B(\tilde{\mathbf{F}}) = [\mathbf{L}^\top, \mathbf{R}^\top]^\top$, by the definition of $\operatorname{dist}(\mathcal{P}_B(\tilde{\mathbf{F}}), \mathbf{F}_\star)$, we know that

$$\operatorname{dist}^2(\mathcal{P}_B(\tilde{\mathbf{F}}), \mathbf{F}_\star) \leq \sum_{i=1}^{n_1} \left\| \mathbf{L}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_\star^{1/2} - (\mathbf{L}_\star \Sigma_\star^{1/2})_{i,\cdot} \right\|_2^2 + \sum_{j=1}^{n_2} \left\| \mathbf{R}_{j,\cdot} \tilde{\mathbf{Q}}^{-\top} \Sigma_\star^{1/2} - (\mathbf{R}_\star \Sigma_\star^{1/2})_{j,\cdot} \right\|_2^2. \quad (66)$$

Recall that the condition $\operatorname{dist}(\tilde{\mathbf{F}}, \mathbf{F}_\star) \leq \epsilon \sigma_r(\mathbf{X}_\star)$ implies

$$\left\| (\tilde{\mathbf{L}} \tilde{\mathbf{Q}} - \mathbf{L}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \vee \left\| (\tilde{\mathbf{R}} \tilde{\mathbf{Q}}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \leq \epsilon,$$

which, together with $\mathbf{R}_\star \Sigma_\star^{-1/2} = \mathbf{V}_\star$, further implies that

$$\begin{aligned} \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top \right\|_2 &\leq \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_\star^{1/2} \right\|_2 \left\| \tilde{\mathbf{R}} \tilde{\mathbf{Q}}^{-\top} \Sigma_\star^{-1/2} \right\|_{\text{op}} \\ &\leq \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_\star^{1/2} \right\|_2 \left(\|\mathbf{V}_\star\|_{\text{op}} + \left\| (\tilde{\mathbf{R}} \tilde{\mathbf{Q}}^{-\top} - \mathbf{R}_\star) \Sigma_\star^{-1/2} \right\|_{\text{op}} \right) \leq (1 + \epsilon) \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_\star^{1/2} \right\|_2. \end{aligned}$$

In addition, the μ -incoherence of \mathbf{X}_\star yields

$$\sqrt{n_1} \left\| (\mathbf{L}_\star \Sigma_\star^{1/2})_{i,\cdot} \right\|_2 \leq \sqrt{n_1} \|\mathbf{U}_\star\|_{2, \infty} \|\Sigma_\star\|_{\text{op}} \leq \sqrt{\mu r} \sigma_1(\mathbf{X}_\star) \leq \frac{B}{1 + \epsilon},$$

where the last inequality follows from the choice of B . Take the above two relations collectively to reach

$$\frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \geq \frac{\left\| (\mathbf{L}_\star \Sigma_\star^{1/2})_{i,\cdot} \right\|_2}{\left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_\star^{1/2} \right\|_2}.$$

We claim that performing the following projection yields a contraction on each row; see also (Zheng and Lafferty, 2016, Lemma 11).

Claim 5 For vectors $\mathbf{u}, \mathbf{u}_* \in \mathbb{R}^n$ and $\lambda \geq \|\mathbf{u}_*\|_2/\|\mathbf{u}\|_2$, it holds that

$$\|(1 \wedge \lambda)\mathbf{u} - \mathbf{u}_*\|_2 \leq \|\mathbf{u} - \mathbf{u}_*\|_2.$$

Apply Claim 5 with $\mathbf{u} := \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_*^{1/2}$, $\mathbf{u}_* := (\mathbf{L}_* \Sigma_*^{1/2})_{i,\cdot}$, and $\lambda := B/(\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2)$ to obtain

$$\begin{aligned} \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_*^{1/2} - (\mathbf{L}_* \Sigma_*^{1/2})_{i,\cdot} \right\|_2^2 &= \left\| \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \right) \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_*^{1/2} - (\mathbf{L}_* \Sigma_*^{1/2})_{i,\cdot} \right\|_2^2 \\ &\leq \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_*^{1/2} - (\mathbf{L}_* \Sigma_*^{1/2})_{i,\cdot} \right\|_2^2. \end{aligned}$$

Following a similar argument for \mathbf{R} , and plugging them back to (66), we conclude that

$$\text{dist}^2(\mathcal{P}_B(\tilde{\mathbf{F}}), \mathbf{F}_*) \leq \sum_{i=1}^{n_1} \left\| \tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{Q}} \Sigma_*^{1/2} - (\mathbf{L}_* \Sigma_*^{1/2})_{i,\cdot} \right\|_2^2 + \sum_{j=1}^{n_2} \left\| \tilde{\mathbf{R}}_{j,\cdot} \tilde{\mathbf{Q}}^{-\top} \Sigma_*^{1/2} - (\mathbf{R}_* \Sigma_*^{1/2})_{j,\cdot} \right\|_2^2 = \text{dist}^2(\tilde{\mathbf{F}}, \mathbf{F}_*).$$

We move on to the incoherence condition. For any $1 \leq i \leq n_1$, one has

$$\begin{aligned} \|\mathbf{L}_{i,\cdot} \mathbf{R}^\top\|_2^2 &= \sum_{j=1}^{n_2} \langle \mathbf{L}_{i,\cdot}, \mathbf{R}_{j,\cdot} \rangle^2 = \sum_{j=1}^{n_2} \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \right)^2 \langle \tilde{\mathbf{L}}_{i,\cdot}, \tilde{\mathbf{R}}_{j,\cdot} \rangle^2 \left(1 \wedge \frac{B}{\sqrt{n_2} \|\tilde{\mathbf{R}}_{j,\cdot} \tilde{\mathbf{L}}^\top\|_2} \right)^2 \\ &\stackrel{(i)}{\leq} \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \right)^2 \sum_{j=1}^{n_2} \langle \tilde{\mathbf{L}}_{i,\cdot}, \tilde{\mathbf{R}}_{j,\cdot} \rangle^2 = \left(1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \right)^2 \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2^2 \\ &\stackrel{(ii)}{\leq} \frac{B^2}{n_1}. \end{aligned}$$

where (i) follows from $1 \wedge \frac{B}{\sqrt{n_2} \|\tilde{\mathbf{R}}_{j,\cdot} \tilde{\mathbf{L}}^\top\|_2} \leq 1$, and (ii) follows from $1 \wedge \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2} \leq \frac{B}{\sqrt{n_1} \|\tilde{\mathbf{L}}_{i,\cdot} \tilde{\mathbf{R}}^\top\|_2}$. Similarly, one has $\|\mathbf{R}_{j,\cdot} \mathbf{L}^\top\|_2^2 \leq B^2/n_2$. Combining these two bounds completes the proof.

Proof [Proof of Claim 5] When $\lambda > 1$, the claim holds as an identity. Otherwise $\lambda \leq 1$. Denote $h(\bar{\lambda}) := \|\bar{\lambda}\mathbf{u} - \mathbf{u}_*\|_2^2$. Calculate its derivative to conclude that $h(\bar{\lambda})$ is monotonically increasing when $\bar{\lambda} \geq \lambda_* := \langle \mathbf{u}, \mathbf{u}_* \rangle / \|\mathbf{u}\|_2^2$. Note that $\lambda \geq \|\mathbf{u}_*\|_2/\|\mathbf{u}\|_2 \geq \lambda_*$, thus $h(\lambda) \leq h(1)$, i.e. the claim holds. ■

E.2 Proof of Lemma 20

We first record two useful lemmas regarding the projector $\mathcal{P}_\Omega(\cdot)$.

Lemma 35 ((Zheng and Lafferty, 2016, Lemma 10)) Suppose that \mathbf{X}_* is μ -incoherent, and $p \gtrsim \mu r \log(n_1 \vee n_2)/(n_1 \wedge n_2)$. With overwhelming probability, one has

$$\begin{aligned} &\left| \langle (p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}_* \mathbf{R}_A^\top + \mathbf{L}_A \mathbf{R}_*^\top), \mathbf{L}_* \mathbf{R}_B^\top + \mathbf{L}_B \mathbf{R}_*^\top \rangle \right| \\ &\leq C_1 \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\mathbf{L}_* \mathbf{R}_A^\top + \mathbf{L}_A \mathbf{R}_*^\top\|_F \|\mathbf{L}_* \mathbf{R}_B^\top + \mathbf{L}_B \mathbf{R}_*^\top\|_F, \end{aligned}$$

simultaneously for all $\mathbf{L}_A, \mathbf{L}_B \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{R}_A, \mathbf{R}_B \in \mathbb{R}^{n_2 \times r}$, where $C_1 > 0$ is some universal constant.

Lemma 36 ((Chen and Li, 2019, Lemma 8), (Chen et al., 2020a, Lemma 12)) *Suppose that $p \gtrsim \log(n_1 \vee n_2)/(n_1 \wedge n_2)$. With overwhelming probability, one has*

$$\begin{aligned} & \left| \langle (p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}_A \mathbf{R}_A^\top, \mathbf{L}_B \mathbf{R}_B^\top) \rangle \right| \\ & \leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} (\|\mathbf{L}_A\|_F \|\mathbf{L}_B\|_{2,\infty} \wedge \|\mathbf{L}_A\|_{2,\infty} \|\mathbf{L}_B\|_F) (\|\mathbf{R}_A\|_F \|\mathbf{R}_B\|_{2,\infty} \wedge \|\mathbf{R}_A\|_{2,\infty} \|\mathbf{R}_B\|_F), \end{aligned}$$

simultaneously for all $\mathbf{L}_A, \mathbf{L}_B \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{R}_A, \mathbf{R}_B \in \mathbb{R}^{n_2 \times r}$, where $C_2 > 0$ is some universal constant.

In view of the above two lemmas, define the event \mathcal{E} as the intersection of the events that the bounds in Lemmas 35 and 36 hold, which happens with overwhelming probability. The rest of the proof is then performed under the event that \mathcal{E} holds.

By the condition $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.02\sigma_r(\mathbf{X}_\star)$ and Lemma 22, one knows that \mathbf{Q}_t , the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star exists. Therefore, for notational convenience, we denote $\mathbf{L} := \mathbf{L}_t \mathbf{Q}_t$, $\mathbf{R} := \mathbf{R}_t \mathbf{Q}_t^{-\top}$, $\Delta_L := \mathbf{L} - \mathbf{L}_\star$, $\Delta_R := \mathbf{R} - \mathbf{R}_\star$, and $\epsilon := 0.02$. In addition, denote $\tilde{\mathbf{F}}_{t+1}$ as the update before projection as

$$\tilde{\mathbf{F}}_{t+1} := \begin{bmatrix} \tilde{\mathbf{L}}_{t+1} \\ \tilde{\mathbf{R}}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_t - \eta p^{-1} \mathcal{P}_\Omega(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star) \mathbf{R}_t (\mathbf{R}_t^\top \mathbf{R}_t)^{-1} \\ \mathbf{R}_t - \eta p^{-1} \mathcal{P}_\Omega(\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star)^\top \mathbf{L}_t (\mathbf{L}_t^\top \mathbf{L}_t)^{-1} \end{bmatrix},$$

and therefore $\mathbf{F}_{t+1} = \mathcal{P}_B(\tilde{\mathbf{F}}_{t+1})$. Note that in view of Lemma 19, it suffices to prove the following relation

$$\text{dist}(\tilde{\mathbf{F}}_{t+1}, \mathbf{F}_\star) \leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star). \quad (67)$$

The conclusion $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_F \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ is a simple consequence of Lemma 26; see (48) for a detailed argument. In what follows, we concentrate on proving (67).

To begin with, we list a few easy consequences under the assumed conditions.

Claim 6 *Under conditions $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq \epsilon \sigma_r(\mathbf{X}_\star)$ and $\sqrt{n_1} \|\mathbf{L} \mathbf{R}^\top\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R} \mathbf{L}^\top\|_{2,\infty} \leq C_B \sqrt{\mu r} \sigma_1(\mathbf{X}_\star)$, one has*

$$\|\Delta_L \Sigma_\star^{-1/2}\|_{\text{op}} \vee \|\Delta_R \Sigma_\star^{-1/2}\|_{\text{op}} \leq \epsilon; \quad (68a)$$

$$\left\| \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \leq \frac{1}{1 - \epsilon}; \quad (68b)$$

$$\left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} \leq \frac{1}{(1 - \epsilon)^2}; \quad (68c)$$

$$\sqrt{n_1} \|\mathbf{L} \Sigma_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R} \Sigma_\star^{1/2}\|_{2,\infty} \leq \frac{C_B}{1 - \epsilon} \sqrt{\mu r} \sigma_1(\mathbf{X}_\star); \quad (68d)$$

$$\sqrt{n_1} \|\mathbf{L} \Sigma_\star^{-1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \leq \frac{C_B \kappa}{1 - \epsilon} \sqrt{\mu r}; \quad (68e)$$

$$\sqrt{n_1} \|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty} \vee \sqrt{n_2} \|\Delta_R \Sigma_\star^{1/2}\|_{2,\infty} \leq \left(1 + \frac{C_B}{1 - \epsilon}\right) \sqrt{\mu r} \sigma_1(\mathbf{X}_\star). \quad (68f)$$

Now we are ready to embark on the proof of (67). By the definition of $\text{dist}(\tilde{\mathbf{F}}_{t+1}, \mathbf{F}_\star)$, one has

$$\text{dist}^2(\tilde{\mathbf{F}}_{t+1}, \mathbf{F}_\star) \leq \left\| (\tilde{\mathbf{L}}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_F^2 + \left\| (\tilde{\mathbf{R}}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_F^2, \quad (69)$$

where we recall that \mathbf{Q}_t is the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star . Plug in the update rule (26) and the decomposition $\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star = \Delta_L \mathbf{R}^\top + \mathbf{L}_\star \Delta_R^\top$ to obtain

$$(\tilde{\mathbf{L}}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} = (\mathbf{L} - \eta p^{-1} \mathcal{P}_\Omega(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_\star) \Sigma_\star^{1/2}$$

$$\begin{aligned}
&= \Delta_L \Sigma_\star^{1/2} - \eta(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2} - \eta(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2} \\
&= (1-\eta)\Delta_L \Sigma_\star^{1/2} - \eta\mathbf{L}_\star\Delta_R^\top\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2} - \eta(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2}.
\end{aligned}$$

This allows us to expand the first square in (69) as

$$\begin{aligned}
\left\|(\tilde{\mathbf{L}}_{t+1}\mathbf{Q}_t - \mathbf{L}_\star)\Sigma_\star^{1/2}\right\|_F^2 &= \underbrace{\left\|(1-\eta)\Delta_L \Sigma_\star^{1/2} - \eta\mathbf{L}_\star\Delta_R^\top\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2}\right\|_F^2}_{\mathfrak{P}_1} \\
&\quad - 2\eta(1-\eta)\underbrace{\text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star\Delta_L^\top\right)}_{\mathfrak{P}_2} \\
&\quad + 2\eta^2\underbrace{\text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star(\mathbf{R}^\top\mathbf{R})^{-1}\mathbf{R}^\top\Delta_R\mathbf{L}_\star^\top\right)}_{\mathfrak{P}_3} \\
&\quad + \eta^2\underbrace{\left\|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star^{1/2}\right\|_F^2}_{\mathfrak{P}_4}.
\end{aligned}$$

In the sequel, we shall control the four terms separately, of which \mathfrak{P}_1 is the main term, and $\mathfrak{P}_2, \mathfrak{P}_3$ and \mathfrak{P}_4 are perturbation terms.

1. Notice that the main term \mathfrak{P}_1 has already been controlled in (46) under the condition (68a). It obeys

$$\mathfrak{P}_1 \leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon}\eta(1-\eta) \right) \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2}\eta^2 \|\Delta_R \Sigma_\star^{1/2}\|_F^2.$$

2. For the second term \mathfrak{P}_2 , decompose $\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star = \Delta_L \mathbf{R}_\star^\top + \mathbf{L}\Delta_R^\top$ and apply the triangle inequality to obtain

$$\begin{aligned}
|\mathfrak{P}_2| &= \left| \text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top + \mathbf{L}\Delta_R^\top)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star\Delta_L^\top\right) \right| \\
&\leq \underbrace{\left| \text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top)\mathbf{R}_\star(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star\Delta_L^\top\right) \right|}_{\mathfrak{P}_{2,1}} \\
&\quad + \underbrace{\left| \text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top)\Delta_R(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star\Delta_L^\top\right) \right|}_{\mathfrak{P}_{2,2}} \\
&\quad + \underbrace{\left| \text{tr}\left((p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}\Delta_R^\top)\mathbf{R}(\mathbf{R}^\top\mathbf{R})^{-1}\Sigma_\star\Delta_L^\top\right) \right|}_{\mathfrak{P}_{2,3}}.
\end{aligned}$$

For the first term $\mathfrak{P}_{2,1}$, under the event \mathcal{E} , we can invoke Lemma 35 to obtain

$$\begin{aligned}
\mathfrak{P}_{2,1} &\leq C_1 \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_L \mathbf{R}_\star^\top\|_F \|\Delta_L \Sigma_\star (\mathbf{R}^\top\mathbf{R})^{-1} \mathbf{R}_\star^\top\|_F \\
&\leq C_1 \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_F^2 \|\Sigma_\star^{1/2} (\mathbf{R}^\top\mathbf{R})^{-1} \Sigma_\star^{1/2}\|_{\text{op}},
\end{aligned}$$

where the second line follows from the relation $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_{\text{op}}\|\mathbf{B}\|_F$. Use the condition (68c) to obtain

$$\mathfrak{P}_{2,1} \leq \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_F^2.$$

Regarding the remaining terms $\mathfrak{P}_{2,2}$ and $\mathfrak{P}_{2,3}$, our main hammer is Lemma 36. Invoking Lemma 36 under the event \mathcal{E} with $L_A := \Delta_L \Sigma_\star^{1/2}$, $R_A := R_\star \Sigma_\star^{-1/2}$, $L_B := \Delta_L \Sigma_\star^{1/2}$, and $R_B := \Delta_R (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2}$, we arrive at

$$\begin{aligned} \mathfrak{P}_{2,2} &\leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|R_\star \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \Delta_R (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\mathbb{F}} \\ &\leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\Delta_L \Sigma_\star^{1/2}\|_{2,\infty} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|R_\star \Sigma_\star^{-1/2}\|_{2,\infty} \|\Delta_R \Sigma_\star^{-1/2}\|_{\mathbb{F}} \left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}. \end{aligned}$$

Similarly, with the help of Lemma 36, one has

$$\mathfrak{P}_{2,3} \leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\mathbf{L} \Sigma_\star^{-1/2}\|_{2,\infty} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \|R \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \Sigma_\star^{1/2} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}.$$

Utilizing the consequences in Claim 6, we arrive at

$$\begin{aligned} \mathfrak{P}_{2,2} &\leq \frac{C_2 \kappa}{(1-\epsilon)^2} \left(1 + \frac{C_B}{1-\epsilon}\right) \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}; \\ \mathfrak{P}_{2,3} &\leq \frac{C_2 C_B^2 \kappa^2}{(1-\epsilon)^4} \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}. \end{aligned}$$

We then combine the bounds for $\mathfrak{P}_{2,1}$, $\mathfrak{P}_{2,2}$ and $\mathfrak{P}_{2,3}$ to see

$$\begin{aligned} \mathfrak{P}_2 &\leq \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \\ &\quad + \frac{C_2 \kappa}{(1-\epsilon)^2} \left(1 + \frac{C_B}{1-\epsilon} + \frac{C_B^2 \kappa}{(1-\epsilon)^2}\right) \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &= \delta_1 \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \delta_2 \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &\leq (\delta_1 + \frac{\delta_2}{2}) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \frac{\delta_2}{2} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2, \end{aligned}$$

where we denote

$$\delta_1 := \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}}, \quad \text{and} \quad \delta_2 := \frac{C_2 \kappa}{(1-\epsilon)^2} \left(1 + \frac{C_B}{1-\epsilon} + \frac{C_B^2 \kappa}{(1-\epsilon)^2}\right) \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}}.$$

3. Following a similar argument for controlling \mathfrak{P}_2 (i.e. repeatedly using Lemmas 35 and 36), we can obtain the following bounds for \mathfrak{P}_3 and \mathfrak{P}_4 , whose proof are deferred to the end of this section.

Claim 7 *Under the event \mathcal{E} , one has*

$$\begin{aligned} \mathfrak{P}_3 &\leq \frac{\delta_2}{2} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + (\delta_1 + \frac{\delta_2}{2}) \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2; \\ \mathfrak{P}_4 &\leq \delta_1 (\delta_1 + \delta_2) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \delta_2 (\delta_1 + \delta_2) \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2. \end{aligned}$$

Taking the bounds for \mathfrak{P}_1 , \mathfrak{P}_2 , \mathfrak{P}_3 and \mathfrak{P}_4 collectively yields

$$\begin{aligned} \left\| (\tilde{\mathbf{L}}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 &\leq \left((1-\eta)^2 + \frac{2\epsilon}{1-\epsilon} \eta (1-\eta) \right) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \\ &\quad + \eta (1-\eta) \left((2\delta_1 + \delta_2) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \delta_2 \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right) \end{aligned}$$

$$\begin{aligned}
 & + \eta^2 \left(\delta_2 \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + (2\delta_1 + \delta_2) \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right) \\
 & + \eta^2 \left(\delta_1 (\delta_1 + \delta_2) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \delta_2 (\delta_1 + \delta_2) \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right).
 \end{aligned}$$

A similar upper bound holds for the second square in (69). As a result, we reach the conclusion that

$$\left\| (\tilde{\mathbf{L}}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\tilde{\mathbf{R}}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 \leq \rho^2(\eta; \epsilon, \delta_1, \delta_2) \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star),$$

where the contraction rate $\rho^2(\eta; \epsilon, \delta_1, \delta_2)$ is given by

$$\rho^2(\eta; \epsilon, \delta_1, \delta_2) := (1 - \eta)^2 + \left(\frac{2\epsilon}{1 - \epsilon} + 2(\delta_1 + \delta_2) \right) \eta(1 - \eta) + \left(\frac{2\epsilon + \epsilon^2}{(1 - \epsilon)^2} + 2(\delta_1 + \delta_2) + (\delta_1 + \delta_2)^2 \right) \eta^2.$$

As long as $p \geq C(\mu r \kappa^4 \vee \log(n_1 \vee n_2)) \mu r / (n_1 \wedge n_2)$ for some sufficiently large constant C , one has $\delta_1 + \delta_2 \leq 0.1$ under the setting $\epsilon = 0.02$. When $0 < \eta \leq 2/3$, one further has $\rho(\eta; \epsilon, \delta_1, \delta_2) \leq 1 - 0.6\eta$. Thus we conclude that

$$\begin{aligned}
 \text{dist}(\tilde{\mathbf{F}}_{t+1}, \mathbf{F}_\star) & \leq \sqrt{\left\| (\tilde{\mathbf{L}}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\tilde{\mathbf{R}}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2} \\
 & \leq (1 - 0.6\eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star),
 \end{aligned}$$

which is exactly the upper bound we are after; see (67). This finishes the proof.

Proof [Proof of Claim 6] First, repeating the derivation for (45) obtains (68a). Second, take the condition (68a) and Lemma 25 together to obtain (68b) and (68c). Third, take the incoherence condition $\sqrt{n_1} \|\mathbf{L} \mathbf{R}^\top\|_{2,\infty} \vee \sqrt{n_2} \|\mathbf{R} \mathbf{L}^\top\|_{2,\infty} \leq C_B \sqrt{\mu r} \sigma_1(\mathbf{X}_\star)$ together with the relations

$$\begin{aligned}
 \|\mathbf{L} \mathbf{R}^\top\|_{2,\infty} & \geq \sigma_r(\mathbf{R} \Sigma_\star^{-1/2}) \|\mathbf{L} \Sigma_\star^{1/2}\|_{2,\infty} \\
 & \geq \left(\sigma_r(\mathbf{R}_\star \Sigma_\star^{-1/2}) - \|\Delta_R \Sigma_\star^{-1/2}\|_{\text{op}} \right) \|\mathbf{L} \Sigma_\star^{1/2}\|_{2,\infty} \\
 & \geq (1 - \epsilon) \|\mathbf{L} \Sigma_\star^{1/2}\|_{2,\infty}; \\
 \|\mathbf{R} \mathbf{L}^\top\|_{2,\infty} & \geq \sigma_r(\mathbf{L} \Sigma_\star^{-1/2}) \|\mathbf{R} \Sigma_\star^{1/2}\|_{2,\infty} \\
 & \geq \left(\sigma_r(\mathbf{L}_\star \Sigma_\star^{-1/2}) - \|\Delta_L \Sigma_\star^{-1/2}\|_{\text{op}} \right) \|\mathbf{R} \Sigma_\star^{1/2}\|_{2,\infty} \\
 & \geq (1 - \epsilon) \|\mathbf{R} \Sigma_\star^{1/2}\|_{2,\infty}
 \end{aligned}$$

to obtain (68d) and (68e). Finally, apply the triangle inequality together with incoherence assumption to obtain (68f). \blacksquare

Proof [Proof of Claim 7] We start with the term \mathfrak{P}_3 , for which we have

$$\begin{aligned}
 |\mathfrak{P}_3| & \leq \left| \underbrace{\text{tr} \left((p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\mathbf{L}_\star \Delta_R^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \mathbf{L}_\star^\top \right)}_{\mathfrak{P}_{3,1}} \right| \\
 & \quad + \left| \underbrace{\text{tr} \left((p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}^\top) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \mathbf{L}_\star^\top \right)}_{\mathfrak{P}_{3,2}} \right|.
 \end{aligned}$$

Invoke Lemma 35 to bound $\mathfrak{P}_{3,1}$ as

$$\mathfrak{P}_{3,1} \leq C_1 \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\mathbf{L}_\star \Delta_R^\top\|_{\mathbb{F}} \|\mathbf{L}_\star \Delta_R^\top \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top\|_{\mathbb{F}}$$

$$\leq C_1 \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2.$$

The condition (68b) allows us to obtain a simplified bound

$$\mathfrak{P}_{3,1} \leq \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2.$$

In regard to $\mathfrak{P}_{3,2}$, we apply Lemma 36 with $\mathbf{L}_A := \Delta_L \Sigma_\star^{1/2}$, $\mathbf{R}_A := \mathbf{R} \Sigma_\star^{-1/2}$, $\mathbf{L}_B := \mathbf{L}_\star \Sigma_\star^{-1/2}$, and $\mathbf{R}_B := \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \Sigma_\star^{1/2}$ to see

$$\begin{aligned} \mathfrak{P}_{3,2} &\leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\mathbf{L}_\star \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \Delta_R \Sigma_\star^{1/2} \right\|_{\mathbb{F}} \\ &\leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\mathbf{L}_\star \Sigma_\star^{-1/2}\|_{2,\infty} \|\mathbf{R} \Sigma_\star^{-1/2}\|_{2,\infty} \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}. \end{aligned}$$

Again, use the consequences in Claim 6 to reach

$$\begin{aligned} \mathfrak{P}_{3,2} &\leq C_2 \sqrt{\frac{n_1 \vee n_2}{p}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \sqrt{\frac{\mu r}{n_1}} \frac{C_B \kappa}{1-\epsilon} \sqrt{\frac{\mu r}{n_2}} \frac{1}{(1-\epsilon)^2} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &= \frac{C_2 C_B \kappa}{(1-\epsilon)^3} \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}. \end{aligned}$$

Combine the bounds of $\mathfrak{P}_{3,1}$ and $\mathfrak{P}_{3,2}$ to reach

$$\begin{aligned} \mathfrak{P}_3 &\leq \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \\ &\quad + \frac{C_2 C_B \kappa}{(1-\epsilon)^3} \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &\leq \delta_1 \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \delta_2 \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &\leq \frac{\delta_2}{2} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + (\delta_1 + \frac{\delta_2}{2}) \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2. \end{aligned}$$

Moving on to the term \mathfrak{P}_4 , we have

$$\begin{aligned} \sqrt{\mathfrak{P}_4} &= \left\| (p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\mathbf{L} \mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\mathbb{F}} \\ &\leq \underbrace{\left| \text{tr} \left((p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top) \mathbf{R}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \tilde{\mathbf{L}}^\top \right) \right|}_{\mathfrak{P}_{4,1}} \\ &\quad + \underbrace{\left| \text{tr} \left((p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\Delta_L \mathbf{R}_\star^\top) \Delta_R (\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \tilde{\mathbf{L}}^\top \right) \right|}_{\mathfrak{P}_{4,2}} \\ &\quad + \underbrace{\left| \text{tr} \left((p^{-1} \mathcal{P}_\Omega - \mathcal{I})(\mathbf{L} \Delta_R^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \tilde{\mathbf{L}}^\top \right) \right|}_{\mathfrak{P}_{4,3}}, \end{aligned}$$

where we have used the variational representation of the Frobenius norm for some $\tilde{\mathbf{L}} \in \mathbb{R}^{n_1 \times r}$ obeying $\|\tilde{\mathbf{L}}\|_{\mathbb{F}} = 1$. Note that the decomposition of $\sqrt{\mathfrak{P}_4}$ is extremely similar to that of \mathfrak{P}_2 . Therefore we can follow a similar argument (i.e. applying Lemmas 35 and 36) to control these terms as

$$\mathfrak{P}_{4,1} \leq \frac{C_1}{(1-\epsilon)^2} \sqrt{\frac{\mu r \log(n_1 \vee n_2)}{p(n_1 \wedge n_2)}} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}};$$

$$\begin{aligned}\mathfrak{P}_{4,2} &\leq \frac{C_2\kappa}{(1-\epsilon)^2} \left(1 + \frac{C_B}{1-\epsilon}\right) \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_R \Sigma_\star^{1/2}\|_F; \\ \mathfrak{P}_{4,3} &\leq \frac{C_2 C_B^2 \kappa^2}{(1-\epsilon)^4} \frac{\mu r}{\sqrt{p(n_1 \wedge n_2)}} \|\Delta_R \Sigma_\star^{1/2}\|_F.\end{aligned}$$

For conciseness, we omit the details for bounding each term. Combine them to reach

$$\sqrt{\mathfrak{P}_4} \leq \delta_1 \|\Delta_L \Sigma_\star^{1/2}\|_F + \delta_2 \|\Delta_R \Sigma_\star^{1/2}\|_F.$$

Finally take the square on both sides and use $2ab \leq a^2 + b^2$ to obtain the upper bound

$$\mathfrak{P}_4 \leq \delta_1(\delta_1 + \delta_2) \|\Delta_L \Sigma_\star^{1/2}\|_F^2 + \delta_2(\delta_1 + \delta_2) \|\Delta_R \Sigma_\star^{1/2}\|_F^2. \quad \blacksquare$$

E.3 Proof of Lemma 21

We start by recording a useful lemma below.

Lemma 37 ((Chen, 2015, Lemma 2), (Chen et al., 2020a, Lemma 4)) *For any fixed $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$, with overwhelming probability, one has*

$$\|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{X})\|_{\text{op}} \leq C_0 \frac{\log(n_1 \vee n_2)}{p} \|\mathbf{X}\|_\infty + C_0 \sqrt{\frac{\log(n_1 \vee n_2)}{p}} (\|\mathbf{X}\|_{2,\infty} \vee \|\mathbf{X}^\top\|_{2,\infty}),$$

where $C_0 > 0$ is some universal constant that does not depend on \mathbf{X} .

In view of Lemma 24, one has

$$\text{dist}(\tilde{\mathbf{F}}_0, \mathbf{F}_\star) \leq \sqrt{\sqrt{2} + 1} \|\mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top - \mathbf{X}_\star\|_F \leq \sqrt{(\sqrt{2} + 1)2r} \|\mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top - \mathbf{X}_\star\|_{\text{op}}, \quad (70)$$

where the last relation uses the fact that $\mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top - \mathbf{X}_\star$ has rank at most $2r$. Applying the triangle inequality, we obtain

$$\begin{aligned}\|\mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top - \mathbf{X}_\star\|_{\text{op}} &\leq \|p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star) - \mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top\|_{\text{op}} + \|p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star) - \mathbf{X}_\star\|_{\text{op}} \\ &\leq 2 \|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{X}_\star)\|_{\text{op}}.\end{aligned} \quad (71)$$

Here the second inequality hinges on the fact that $\mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top$ is the best rank- r approximation to $p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star)$, i.e.

$$\|p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star) - \mathbf{U}_0 \Sigma_0 \mathbf{V}_0^\top\|_{\text{op}} \leq \|p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star) - \mathbf{X}_\star\|_{\text{op}}.$$

Combining (70) and (71) yields

$$\text{dist}(\tilde{\mathbf{F}}_0, \mathbf{F}_\star) \leq 2\sqrt{(\sqrt{2} + 1)2r} \|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{X}_\star)\|_{\text{op}} \leq 5\sqrt{r} \|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{X}_\star)\|_{\text{op}}.$$

It then boils down to controlling $\|p^{-1}\mathcal{P}_\Omega(\mathbf{X}_\star) - \mathbf{X}_\star\|_{\text{op}}$, which is readily supplied by Lemma 37 as

$$\|(p^{-1}\mathcal{P}_\Omega - \mathcal{I})(\mathbf{X}_\star)\|_{\text{op}} \leq C_0 \frac{\log(n_1 \vee n_2)}{p} \|\mathbf{X}_\star\|_\infty + C_0 \sqrt{\frac{\log(n_1 \vee n_2)}{p}} (\|\mathbf{X}_\star\|_{2,\infty} \vee \|\mathbf{X}_\star^\top\|_{2,\infty}),$$

which holds with overwhelming probability. The proof is finished by plugging the following bounds from incoherence assumption of \mathbf{X}_\star :

$$\begin{aligned}\|\mathbf{X}_\star\|_\infty &\leq \|\mathbf{U}_\star\|_{2,\infty}\|\boldsymbol{\Sigma}_\star\|_{\text{op}}\|\mathbf{V}_\star\|_{2,\infty} \leq \frac{\mu r}{\sqrt{n_1 n_2}}\kappa\sigma_r(\mathbf{X}_\star); \\ \|\mathbf{X}_\star\|_{2,\infty} &\leq \|\mathbf{U}_\star\|_{2,\infty}\|\boldsymbol{\Sigma}_\star\|_{\text{op}}\|\mathbf{V}_\star\|_{\text{op}} \leq \sqrt{\frac{\mu r}{n_1}}\kappa\sigma_r(\mathbf{X}_\star); \\ \|\mathbf{X}_\star^\top\|_{2,\infty} &\leq \|\mathbf{U}_\star\|_{\text{op}}\|\boldsymbol{\Sigma}_\star\|_{\text{op}}\|\mathbf{V}_\star\|_{2,\infty} \leq \sqrt{\frac{\mu r}{n_2}}\kappa\sigma_r(\mathbf{X}_\star).\end{aligned}$$

Appendix F. Proof for General Loss Functions

We first present a useful property of restricted smooth and convex functions.

Lemma 38 *Suppose that $f : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}$ is rank- $2r$ restricted L -smooth and rank- $2r$ restricted convex. Then for any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^{n_1 \times n_2}$ of rank at most r , one has*

$$\langle \nabla f(\mathbf{X}_1) - \nabla f(\mathbf{X}_2), \mathbf{X}_1 - \mathbf{X}_2 \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{X}_1) - \nabla f(\mathbf{X}_2)\|_{\mathbb{F},r}^2.$$

Proof Since $f(\cdot)$ is rank- $2r$ restricted L -smooth and convex, it holds for any $\bar{\mathbf{X}} \in \mathbb{R}^{n_1 \times n_2}$ with rank at most $2r$ that

$$f(\mathbf{X}_1) + \langle \nabla f(\mathbf{X}_1), \bar{\mathbf{X}} - \mathbf{X}_1 \rangle \leq f(\bar{\mathbf{X}}) \leq f(\mathbf{X}_2) + \langle \nabla f(\mathbf{X}_2), \bar{\mathbf{X}} - \mathbf{X}_2 \rangle + \frac{L}{2} \|\bar{\mathbf{X}} - \mathbf{X}_2\|_{\mathbb{F}}^2.$$

Reorganize the terms to yield

$$f(\mathbf{X}_1) + \langle \nabla f(\mathbf{X}_1), \mathbf{X}_2 - \mathbf{X}_1 \rangle \leq f(\mathbf{X}_2) + \langle \nabla f(\mathbf{X}_2) - \nabla f(\mathbf{X}_1), \bar{\mathbf{X}} - \mathbf{X}_2 \rangle + \frac{L}{2} \|\bar{\mathbf{X}} - \mathbf{X}_2\|_{\mathbb{F}}^2.$$

Take $\bar{\mathbf{X}} = \mathbf{X}_2 - \frac{1}{L} \mathcal{P}_r(\nabla f(\mathbf{X}_2) - \nabla f(\mathbf{X}_1))$, whose rank is at most $2r$, to see

$$f(\mathbf{X}_1) + \langle \nabla f(\mathbf{X}_1), \mathbf{X}_2 - \mathbf{X}_1 \rangle + \frac{1}{2L} \|\nabla f(\mathbf{X}_2) - \nabla f(\mathbf{X}_1)\|_{\mathbb{F},r}^2 \leq f(\mathbf{X}_2).$$

We can further switch the roles of \mathbf{X}_1 and \mathbf{X}_2 to obtain

$$f(\mathbf{X}_2) + \langle \nabla f(\mathbf{X}_2), \mathbf{X}_1 - \mathbf{X}_2 \rangle + \frac{1}{2L} \|\nabla f(\mathbf{X}_2) - \nabla f(\mathbf{X}_1)\|_{\mathbb{F},r}^2 \leq f(\mathbf{X}_1).$$

Adding the above two inequalities yields the desired bound. \blacksquare

F.1 Proof of Theorem 11

Suppose that the t -th iterate \mathbf{F}_t obeys the condition $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq 0.1\sigma_r(\mathbf{X}_\star)/\sqrt{\kappa\bar{f}}$. In view of Lemma 22, one knows that \mathbf{Q}_t , the optimal alignment matrix between \mathbf{F}_t and \mathbf{F}_\star exists. Therefore, for notational convenience, denote $\mathbf{L} := \mathbf{L}_t \mathbf{Q}_t$, $\mathbf{R} := \mathbf{R}_t \mathbf{Q}_t^{-\top}$, $\boldsymbol{\Delta}_L := \mathbf{L} - \mathbf{L}_\star$, $\boldsymbol{\Delta}_R := \mathbf{R} - \mathbf{R}_\star$, and $\epsilon := 0.1/\sqrt{\kappa\bar{f}}$. Similar to the derivation in (45), we have

$$\|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{-1/2}\|_{\text{op}} \vee \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{-1/2}\|_{\text{op}} \leq \epsilon. \quad (72)$$

The conclusion $\|\mathbf{L}_t \mathbf{R}_t^\top - \mathbf{X}_\star\|_{\mathbb{F}} \leq 1.5 \text{dist}(\mathbf{F}_t, \mathbf{F}_\star)$ is a simple consequence of Lemma 26; see (48) for a detailed argument. From now on, we focus on proving the distance contraction.

By the definition of $\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star)$, one has

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2. \quad (73)$$

Introduce an auxiliary function

$$f_\mu(\mathbf{X}) = f(\mathbf{X}) - \frac{\mu}{2} \|\mathbf{X} - \mathbf{X}_\star\|_{\mathbb{F}}^2,$$

which is rank- $2r$ restricted $(L - \mu)$ -smooth and rank- $2r$ restricted convex. Using the ScaledGD update rule (27) and the decomposition $\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star = \boldsymbol{\Delta}_L \mathbf{R}^\top + \mathbf{L}_\star \boldsymbol{\Delta}_R^\top$, we obtain

$$\begin{aligned} (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} &= (\mathbf{L} - \eta \nabla f(\mathbf{L}\mathbf{R}^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} \\ &= (\mathbf{L} - \eta \mu (\mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} - \eta \nabla f_\mu(\mathbf{L}\mathbf{R}^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} \\ &= (1 - \eta \mu) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} - \eta \mu \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} - \eta \nabla f_\mu(\mathbf{L}\mathbf{R}^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2}. \end{aligned}$$

As a result, one can expand the first square in (73) as

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2 &= \underbrace{\left\| (1 - \eta \mu) \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2} - \eta \mu \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2}_{\mathfrak{G}_1} \\ &\quad - 2\eta(1 - \eta \mu) \underbrace{\left\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top - \boldsymbol{\Delta}_L \mathbf{R}_\star^\top - \frac{1}{2} \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\rangle}_{\mathfrak{G}_2} \\ &\quad - 2\eta(1 - \eta \mu) \left\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \boldsymbol{\Delta}_L \mathbf{R}_\star^\top + \frac{1}{2} \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\rangle \\ &\quad + 2\eta^2 \mu \underbrace{\left\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \mathbf{L}_\star \boldsymbol{\Delta}_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \right\rangle}_{\mathfrak{G}_3} \\ &\quad + \eta^2 \underbrace{\left\| \nabla f_\mu(\mathbf{L}\mathbf{R}^\top) \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} \right\|_{\mathbb{F}}^2}_{\mathfrak{G}_4}. \end{aligned}$$

In the sequel, we shall bound the four terms separately.

1. Notice that the main term \mathfrak{G}_1 has already been controlled in (46) under the condition (72). It obeys

$$\mathfrak{G}_1 \leq \left((1 - \eta \mu)^2 + \frac{2\epsilon}{1 - \epsilon} \eta \mu (1 - \eta \mu) \right) \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2}\|_{\mathbb{F}}^2 + \frac{2\epsilon + \epsilon^2}{(1 - \epsilon)^2} \eta^2 \mu^2 \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{1/2}\|_{\mathbb{F}}^2,$$

as long as $\eta \mu \leq 2/3$.

2. For the second term \mathfrak{G}_2 , note that $\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top - \boldsymbol{\Delta}_L \mathbf{R}_\star^\top - \frac{1}{2} \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top$ has rank at most r . Hence we can invoke Lemma 28 to obtain

$$\begin{aligned} |\mathfrak{G}_2| &\leq \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F}, r} \left\| \boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top - \boldsymbol{\Delta}_L \mathbf{R}_\star^\top - \frac{1}{2} \boldsymbol{\Delta}_L \boldsymbol{\Delta}_R^\top \right\|_{\mathbb{F}} \\ &\leq \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F}, r} \|\boldsymbol{\Delta}_L \boldsymbol{\Sigma}_\star^{1/2}\|_{\mathbb{F}} \left(\left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \boldsymbol{\Sigma}_\star^{1/2} - \mathbf{V}_\star \right\|_{\text{op}} + \frac{1}{2} \|\boldsymbol{\Delta}_R \boldsymbol{\Sigma}_\star^{-1/2}\|_{\text{op}} \right), \end{aligned}$$

where the second line uses $\mathbf{R}_\star = \mathbf{V}_\star \Sigma_\star^{1/2}$. Take the condition (72) and Lemma 25 together to obtain

$$\begin{aligned} \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}} &\leq \frac{1}{1-\epsilon}; \\ \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} - \mathbf{V}_\star \right\|_{\text{op}} &\leq \frac{\sqrt{2}\epsilon}{1-\epsilon}. \end{aligned}$$

These consequences further imply that

$$|\mathfrak{G}_2| \leq \left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}.$$

3. As above, the third term \mathfrak{G}_3 can be similarly bounded as

$$\begin{aligned} |\mathfrak{G}_3| &\leq \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \left\| \mathbf{L}_\star \Delta_R^\top \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star (\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \right\|_{\mathbb{F}} \\ &\leq \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \\ &\leq \frac{1}{(1-\epsilon)^2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}. \end{aligned}$$

4. For the last term \mathfrak{G}_4 , invoke Lemma 28 to obtain

$$\mathfrak{G}_4 \leq \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2 \left\| \mathbf{R}(\mathbf{R}^\top \mathbf{R})^{-1} \Sigma_\star^{1/2} \right\|_{\text{op}}^2 \leq \frac{1}{(1-\epsilon)^2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2.$$

Taking collectively the bounds for $\mathfrak{G}_1, \mathfrak{G}_2, \mathfrak{G}_3$ and \mathfrak{G}_4 yields

$$\begin{aligned} \left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 &\leq \left((1-\eta\mu)^2 + \frac{2\epsilon}{1-\epsilon} \eta\mu(1-\eta\mu) \right) \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \mu^2 \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \\ &\quad + 2\eta \left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) (1-\eta\mu) \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} \\ &\quad - 2\eta(1-\eta\mu) \left\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \Delta_L \mathbf{R}_\star^\top + \frac{1}{2} \Delta_L \Delta_R^\top \right\rangle \\ &\quad + \frac{2\eta^2 \mu}{(1-\epsilon)^2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} + \frac{\eta^2}{(1-\epsilon)^2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2. \end{aligned}$$

Similarly, we can obtain the control of $\|(\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2}\|_{\mathbb{F}}^2$. Combine them together to reach

$$\begin{aligned} &\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 \\ &\leq \left((1-\eta\mu)^2 + \frac{2\epsilon}{1-\epsilon} \eta\mu(1-\eta\mu) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \mu^2 \right) \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right) \\ &\quad + 2\eta \left(\left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) (1-\eta\mu) + \frac{\eta\mu}{(1-\epsilon)^2} \right) \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right) \\ &\quad - 2\eta(1-\eta\mu) \left\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \Delta_L \mathbf{R}_\star^\top + \mathbf{L}_\star \Delta_R^\top + \Delta_L \Delta_R^\top \right\rangle + \frac{2\eta^2}{(1-\epsilon)^2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2 \\ &\leq \left((1-\eta\mu)^2 + \frac{2\epsilon}{1-\epsilon} \eta\mu(1-\eta\mu) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \mu^2 \right) \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right) \end{aligned}$$

$$\begin{aligned}
 & + 2\eta \underbrace{\left(\left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) (1-\eta\mu) + \frac{\eta\mu}{(1-\epsilon)^2} \right)}_{\mathfrak{C}_1} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right) \\
 & - 2\eta \underbrace{\left(\frac{1-\eta\mu}{L-\mu} - \frac{\eta}{(1-\epsilon)^2} \right)}_{\mathfrak{C}_2} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2,
 \end{aligned}$$

where the last line follows from Lemma 38 (notice that $\nabla f_\mu(\mathbf{X}_\star) = \mathbf{0}$) as

$$\langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \Delta_L \mathbf{R}_\star^\top + \mathbf{L}_\star \Delta_R^\top + \Delta_L \Delta_R^\top \rangle = \langle \nabla f_\mu(\mathbf{L}\mathbf{R}^\top), \mathbf{L}\mathbf{R}^\top - \mathbf{X}_\star \rangle \geq \frac{1}{L-\mu} \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2.$$

Notice that $\mathfrak{C}_2 > 0$ as long as $\eta \leq (1-\epsilon)^2/L$. Maximizing the quadratic function of $\|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}$ yields

$$\begin{aligned}
 \mathfrak{C}_1 \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right) - \mathfrak{C}_2 \|\nabla f_\mu(\mathbf{L}\mathbf{R}^\top)\|_{\mathbb{F},r}^2 & \leq \frac{\mathfrak{C}_1^2}{4\mathfrak{C}_2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}} + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}} \right)^2 \\
 & \leq \frac{\mathfrak{C}_1^2}{2\mathfrak{C}_2} \left(\|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 \right),
 \end{aligned}$$

where the last inequality holds since $(a+b)^2 \leq 2(a^2+b^2)$. Identify $\text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star) = \|\Delta_L \Sigma_\star^{1/2}\|_{\mathbb{F}}^2 + \|\Delta_R \Sigma_\star^{1/2}\|_{\mathbb{F}}^2$ to obtain

$$\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 \leq \rho^2(\eta; \epsilon, \mu, L) \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star),$$

where the contraction rate is given by

$$\rho^2(\eta; \epsilon, \mu, L) := (1-\eta\mu)^2 + \frac{2\epsilon}{1-\epsilon} \eta\mu(1-\eta\mu) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \mu^2 + \frac{\left(\left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) (1-\eta\mu) + \frac{\eta\mu}{(1-\epsilon)^2} \right)^2}{1-\eta\mu - \frac{\eta(L-\mu)}{(1-\epsilon)^2}} \eta(L-\mu).$$

With $\epsilon = 0.1/\sqrt{\kappa_f}$ and $0 < \eta \leq 0.4/L$, one has $\rho(\eta; \epsilon, \mu, L) \leq 1 - 0.7\eta\mu$. Thus we conclude that

$$\begin{aligned}
 \text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) & \leq \sqrt{\left\| (\mathbf{L}_{t+1} \mathbf{Q}_t - \mathbf{L}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2 + \left\| (\mathbf{R}_{t+1} \mathbf{Q}_t^{-\top} - \mathbf{R}_\star) \Sigma_\star^{1/2} \right\|_{\mathbb{F}}^2} \\
 & \leq (1 - 0.7\eta\mu) \text{dist}(\mathbf{F}_t, \mathbf{F}_\star),
 \end{aligned}$$

which is the desired claim.

Remark 39 We provide numerical details for the contraction rate. For simplicity, we shall prove $\rho(\eta; \epsilon, \mu, L) \leq 1 - 0.7\eta\mu$ under a stricter condition $\epsilon = 0.02/\sqrt{\kappa_f}$. The stronger result under the condition $\epsilon = 0.1/\sqrt{\kappa_f}$ can be verified through a subtler analysis.

With $\epsilon = 0.02/\sqrt{\kappa_f}$ and $0 < \eta \leq 0.4/L$, one can bound the terms in $\rho^2(\eta; \epsilon, \mu, L)$ as

$$\begin{aligned}
 (1-\eta\mu)^2 + \frac{2\epsilon}{1-\epsilon} \eta\mu(1-\eta\mu) + \frac{2\epsilon + \epsilon^2}{(1-\epsilon)^2} \eta^2 \mu^2 & \leq 1 - 1.959\eta\mu + 1.002\eta^2 \mu^2; \tag{74} \\
 \frac{\left(\left(\frac{\sqrt{2}\epsilon}{1-\epsilon} + \frac{\epsilon}{2} \right) (1-\eta\mu) + \frac{\eta\mu}{(1-\epsilon)^2} \right)^2}{1-\eta\mu - \frac{\eta(L-\mu)}{(1-\epsilon)^2}} \eta(L-\mu) & \leq \frac{\frac{0.0016}{\kappa_f} + 0.078\eta\mu + 1.005\eta^2 \mu^2}{1 - 1.042\eta L} \eta L \\
 & \leq \frac{0.0016\eta \frac{L}{\kappa_f} + 0.4 \times (0.078\eta\mu + 1.005\eta^2 \mu^2)}{1 - 0.4 \times 1.042}
 \end{aligned}$$

$$\leq 0.057\eta\mu + 0.69\eta^2\mu^2, \quad (75)$$

where the last line uses the definition (28) of κ_f . Putting (74) and (75) together further implies

$$\rho^2(\eta; \epsilon, \mu, L) \leq 1 - 1.9\eta\mu + 1.7\eta^2\mu^2 \leq (1 - 0.7\eta\mu)^2,$$

as long as $0 < \eta\mu \leq 0.4$.

References

- Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Srinadh Bhojanapalli, Anastasios Kyrillidis, and Sujay Sanghavi. Dropping convexity for faster semi-definite optimization. In *Conference on Learning Theory*, pages 530–582. PMLR, 2016a.
- Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016b.
- Jian-Feng Cai, Tianming Wang, and Ke Wei. Spectral compressed sensing via projected gradient descent. *SIAM Journal on Optimization*, 28(3):2625–2653, 2018.
- Emmanuel Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow: Theory and algorithms. *Information Theory, IEEE Transactions on*, 61(4):1985–2007, 2015.
- Emmanuel J Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359, 2011.
- Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- Emmanuel J. Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM*, 58(3):11:1–11:37, 2011.
- Venkat Chandrasekaran, Sujay Sanghavi, Pablo Parrilo, and Alan Willsky. Rank-sparsity incoherence for matrix decomposition. *SIAM Journal on Optimization*, 21(2):572–596, 2011.
- Vasileios Charisopoulos, Yudong Chen, Damek Davis, Mateo Díaz, Lijun Ding, and Dmitriy Drusvyatskiy. Low-rank matrix recovery with composite optimization: good conditioning and rapid convergence. *Foundations of Computational Mathematics*, pages 1–89, 2021.
- Ji Chen and Xiaodong Li. Model-free nonconvex matrix completion: Local minima analysis and applications in memory-efficient kernel PCA. *Journal of Machine Learning Research*, 20(142):1–39, 2019.
- Ji Chen, Dekai Liu, and Xiaodong Li. Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization. *IEEE Transactions on Information Theory*, 66(9):5806–5841, 2020a.
- Yudong Chen. Incoherence-optimal matrix completion. *IEEE Transactions on Information Theory*, 61(5):2909–2923, 2015.
- Yudong Chen and Yuejie Chi. Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine*, 35(4):14 – 31, 2018.

- Yudong Chen and Martin J Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- Yuxin Chen and Yuejie Chi. Robust spectral compressed sensing via structured matrix completion. *IEEE Transactions on Information Theory*, 60(10):6576–6601, 2014.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, and Yuling Yan. Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM Journal on Optimization*, 30(4):3098–3121, 2020b.
- Yuxin Chen, Jianqing Fan, Cong Ma, and Yuling Yan. Bridging convex and nonconvex optimization in robust PCA: Noise, outliers, and missing data. *arXiv preprint arXiv:2001.05484*, 2020c.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Damek Davis, Dmitriy Drusvyatskiy, and Courtney Paquette. The nonsmooth landscape of phase retrieval. *arXiv preprint arXiv:1711.03247*, 2017.
- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems*, pages 384–395, 2018.
- Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points-online stochastic gradient for tensor decomposition. In *Conference on Learning Theory (COLT)*, pages 797–842, 2015.
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. In *Advances in Neural Information Processing Systems*, pages 2973–2981, 2016.
- Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, pages 1233–1242, 2017.
- Suriya Gunasekar, Pradeep Ravikumar, and Joydeep Ghosh. Exponential family matrix completion under structural constraints. In *International Conference on Machine Learning*, pages 1917–1925, 2014.
- Moritz Hardt and Mary Wootters. Fast matrix completion without the condition number. In *Proceedings of The 27th Conference on Learning Theory*, pages 638–678, 2014.
- Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Foundations and Trends® in Machine Learning*, 10(3-4):142–336, 2017.
- Prateek Jain, Raghu Meka, and Inderjit S Dhillon. Guaranteed rank minimization via singular value projection. In *Advances in Neural Information Processing Systems*, pages 937–945, 2010.
- Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 665–674. ACM, 2013.
- Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. In *International Conference on Machine Learning*, pages 1724–1732, 2017.
- Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in neural information processing systems*, pages 586–594, 2016.

- Anastasios Kyrillidis and Volkan Cevher. Matrix ALPS: Accelerated low rank and sparse matrix reconstruction. In *2012 IEEE Statistical Signal Processing Workshop (SSP)*, pages 185–188. IEEE, 2012.
- Jean Lafond. Low rank matrix completion with exponential family noise. In *Conference on Learning Theory*, pages 1224–1243, 2015.
- Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Applied and computational harmonic analysis*, 47(3): 893–934, 2019.
- Y. Li, C. Ma, Y. Chen, and Y. Chi. Nonconvex matrix factorization from rank-one measurements. *IEEE Transactions on Information Theory*, 67(3):1928–1950, 2021.
- Yuetian Luo, Wen Huang, Xudong Li, and Anru R Zhang. Recursive importance sketching for rank constrained least squares: Algorithms and high-order convergence. *arXiv preprint arXiv:2011.08360*, 2020.
- 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. *Foundations of Computational Mathematics*, pages 1–182, 2019.
- Cong Ma, Yuanxin Li, and Yuejie Chi. Beyond Procrustes: Balancing-free gradient descent for asymmetric low-rank matrix sensing. *IEEE Transactions on Signal Processing*, 69:867–877, 2021.
- Mantas Mazeika. The singular value decomposition and low rank approximation. Technical report, University of Chicago, 2016.
- Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774, 2018.
- Bamdev Mishra and Rodolphe Sepulchre. Riemannian preconditioning. *SIAM Journal on Optimization*, 26(1):635–660, 2016.
- Bamdev Mishra, K Adithya Apuroop, and Rodolphe Sepulchre. A Riemannian geometry for low-rank matrix completion. *arXiv preprint arXiv:1211.1550*, 2012.
- Yurii Nesterov and Boris T Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- Praneeth Netrapalli, UN Niranjan, Sujay Sanghavi, Animashree Anandkumar, and Prateek Jain. Non-convex robust PCA. In *Advances in Neural Information Processing Systems*, pages 1107–1115, 2014.
- Dohyung Park, Anastasios Kyrillidis, Constantine Carmanis, and Sujay Sanghavi. Non-square matrix sensing without spurious local minima via the Burer-Monteiro approach. In *Artificial Intelligence and Statistics*, pages 65–74, 2017.
- Dohyung Park, Anastasios Kyrillidis, Constantine Caramanis, and Sujay Sanghavi. Finding low-rank solutions via nonconvex matrix factorization, efficiently and provably. *SIAM Journal on Imaging Sciences*, 11(4):2165–2204, 2018.
- Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010.
- Sujay Sanghavi, Rachel Ward, and Chris D White. The local convexity of solving systems of quadratic equations. *Results in Mathematics*, 71(3-4):569–608, 2017.

- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery using nonconvex optimization. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 2351–2360, 2015.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18(5):1131–1198, 2018.
- Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- Jared Tanner and Ke Wei. Low rank matrix completion by alternating steepest descent methods. *Applied and Computational Harmonic Analysis*, 40(2):417–429, 2016.
- Tian Tong, Cong Ma, and Yuejie Chi. Low-rank matrix recovery with scaled subgradient methods: Fast and robust convergence without the condition number. *IEEE Transactions on Signal Processing*, 2021a.
- Tian Tong, Cong Ma, Ashley Prater-Bennette, Erin Tripp, and Yuejie Chi. Scaling and scalability: Provable nonconvex low-rank tensor estimation from incomplete measurements. *arXiv preprint arXiv:2104.14526*, 2021b.
- Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Benjamin Recht. Low-rank solutions of linear matrix equations via Procrustes flow. In *International Conference Machine Learning*, pages 964–973, 2016.
- Ke Wei, Jian-Feng Cai, Tony F Chan, and Shingyu Leung. Guarantees of Riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- Xinyang Yi, Dohyung Park, Yudong Chen, and Constantine Caramanis. Fast algorithms for robust PCA via gradient descent. In *Advances in neural information processing systems*, pages 4152–4160, 2016.
- Qinqing Zheng and John Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In *Advances in Neural Information Processing Systems*, pages 109–117, 2015.
- Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using Burer-Monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.
- Zhihui Zhu, Qiuwei Li, Gongguo Tang, and Michael B Wakin. Global optimality in low-rank matrix optimization. *IEEE Transactions on Signal Processing*, 66(13):3614–3628, 2018.