Efficient Federated Low Rank Matrix Recovery via Alternating GD and Minimization: A Simple Proof

Namrata Vaswani[®], Fellow, IEEE

Abstract—This note provides a significantly simpler and shorter proof of our sample complexity guarantee for solving the low rank column-wise sensing problem using the Alternating Gradient Descent (GD) and Minimization (AltGDmin) algorithm. AltGDmin was developed and analyzed for solving this problem in our recent work. We also provide an improved guarantee.

Index Terms—Low rank column-wise sensing (LRCS), federated learning, multi-task representation learning.

I. INTRODUCTION

E STUDY the low rank column-wise sensing (LRCS) problem which involves recovering a low rank matrix from independent compressive measurements of each of its columns. This problem occurs in dynamic MRI [1] and in multi-task linear representation learning for few-shot learning [2]. The alternating gradient descent (GD) and minimization (AltGDmin) algorithm for solving it in a fast, communication-efficient and private fashion was developed and analyzed in our recent work [3]. This short paper provides a significantly simpler and shorter proof of our sample complexity guarantee for AltGDmin. In fact, it also improves the sample complexity needed by the AltGDmin iterations by a factor of r.

II. PROBLEM STATEMENT, NOTATION, AND ALGORITHM A. Problem Statement and Assumption and Notation

The goal is to recover an $n \times q$ rank-r matrix $\boldsymbol{X}^{\star} = [\boldsymbol{x}_1^{\star}, \boldsymbol{x}_2^{\star}, \dots, \boldsymbol{x}_q^{\star}]$ from m linear projections (sketches) of each of its q columns, i.e. from

$$\boldsymbol{y}_k := \boldsymbol{A}_k \boldsymbol{x}_k^{\star}, \ k \in [q] \tag{1}$$

where each \boldsymbol{y}_k is an m-length vector, $[q] := \{1,2,\ldots,q\}$, and the measurement/sketching matrices \boldsymbol{A}_k are mutually independent and known. The setting of interest is low-rank (LR), $r \ll \min(n,q)$, and undersampled measurements, m < n. Each \boldsymbol{A}_k is assumed to be random-Gaussian: each entry of it is independent and identically distributed (i.i.d.) standard Gaussian. Let $\boldsymbol{X}^\star \stackrel{\mathrm{SVD}}{=} \boldsymbol{U}^\star \boldsymbol{\Sigma}^* \boldsymbol{V}^\star := \boldsymbol{U}^\star \boldsymbol{B}^\star$ denote its

Manuscript received 3 July 2023; revised 26 October 2023; accepted 10 February 2024. Date of publication 1 March 2024; date of current version 18 June 2024. This work was supported in part by NSF under Grant CCF-2115200.

The author is with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA 50011 USA (e-mail: namrata@iastate.edu). Communicated by Y. Chi, Associate Editor for Machine Learning and

Digital Object Identifier 10.1109/TIT.2024.3365795

reduced (rank r) SVD, and $\kappa := \sigma_{\max}^{\star}/\sigma_{\min}^{\star}$ the condition number of Σ^{*} . We let $B^{\star} := \Sigma^{*}V^{\star}$.

Since no measurement y_{ki} is a global function of the entire matrix, X^* , we need the following assumption, borrowed from LR matrix completion literature, to make our problem well-posed (allow for correct interpolation across columns).

Assumption 1 (Incoherence of Right Singular Vectors): Assume that $\|\boldsymbol{b}_k^\star\|^2 \leq \mu^2 r \sigma_{\max}^{\star}^2/q$ for a numerical constant μ . In our discussion of communication complexity and privacy,

we assume a vertically federated setting: different subsets of y_k , A_k are available at different nodes.

1) Notation: We use $\|.\|_F$ to denote the Frobenius norm, $\|.\|$ without a subscript for the (induced) l_2 norm, $^{\top}$ to denote matrix or vector transpose, \boldsymbol{e}_k to denote the k-th canonical basis vector (k-th column of \boldsymbol{I}), and $\boldsymbol{M}^{\dagger} := (\boldsymbol{M}^{\top}\boldsymbol{M})^{-1}\boldsymbol{M}^{\top}$. For two $n \times r$ matrices $\boldsymbol{U}_1, \boldsymbol{U}_2$ that have orthonormal columns, we use

$$\mathrm{SD}_2(\boldsymbol{U}_1, \boldsymbol{U}_2) := \|(\boldsymbol{I} - \boldsymbol{U}_1 \boldsymbol{U}_1^\top) \boldsymbol{U}_2\|$$

as the Subspace Distance (SD) measure. In our previous work [3], we used the Frobenius norm SD,

$$\mathrm{SD}_F(\boldsymbol{U}_1, \boldsymbol{U}_2) := \|(\boldsymbol{I} - \boldsymbol{U}_1 \boldsymbol{U}_1^\top) \boldsymbol{U}_2\|_F.$$

Clearly, $\mathrm{SD}_F(U_1, U_2) \leq \sqrt{r} \mathrm{SD}_2(U_1, U_2)$. We reuse the letters c, C to denote different numerical constants in each use with the convention that c < 1 and $C \geq 1$. We use \sum_k as a shortcut for the summation over k = 1 to q and \sum_{ki} for the summation over i = 1 to m and k = 1 to q. We use whp to refer to "with high probability" and this means that the claim holds with probability (w.p.) at least $1 - n^{-10}$.

B. Review of AltGDmin Algorithm [3]

AltGDmin, summarized in Algorithm 1, imposes the LR constraint by factorizing the unknown matrix \boldsymbol{X} as $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{B}$ with \boldsymbol{U} being an $n \times r$ matrix and \boldsymbol{B} an $r \times q$ matrix. It minimizes $f(\boldsymbol{U}, \boldsymbol{B}) := \sum_{k=1}^q \|\boldsymbol{y}_k - \boldsymbol{U}\boldsymbol{b}_k\|^2$ as follows:

- 1) Truncated spectral initialization: Initialize U (see below).
- 2) At each iteration, update B and U as follows:
 - a) Minimization for B: keeping U fixed, update B by solving $\min_{B} f(U, B)$. Due to the form of the LRCS model, this minimization decouples across columns, making it a cheap least squares problem of recovering q different r length vectors. It is solved as $b_k = (A_k U)^{\dagger} y_k$ for each $k \in [q]$.

0018-9448 © 2024 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

b) Projected-GD for U: keeping B fixed, update U by a GD step, followed by orthonormalizing its columns: $U^+ = QR(U - \eta \nabla_U f(U, B))$. Here QR(.) orthonormalizes the columns of its input.

We initialize U by computing the top r singular vectors of

$$\boldsymbol{X}_0 := \sum_k \boldsymbol{A}_k^\top \boldsymbol{y}_{k,trnc} \boldsymbol{e}_k^\top, \ \boldsymbol{y}_{k,trnc} := \operatorname{trunc}(\boldsymbol{y}_k, \alpha)$$

Here $\alpha := \tilde{C} \sum_k \|\boldsymbol{y}_k\|^2 / mq$ with $\tilde{C} := 9\kappa^2 \mu^2$. The function trunc truncates (zeroes out) all entries of y_k with magnitude greater than $\sqrt{\alpha}$, i.e., for all $j \in [n]$, trunc $(\boldsymbol{y}, \alpha)_j =$ $(y)_j \mathbb{1}_{|y_j| < \sqrt{\alpha}}$, with $\mathbb{1}$ being the indicator function.

Sample-splitting is assumed, i.e., each new update of Uand B uses a new independent set of measurements and measurement matrices, y_k , A_k .

The use of minimization to update B at each iteration is what helps ensure that we can show exponential error decay with a constant step size. At the same time, due to the columnwise decoupled nature of LRCS, the time complexity for this step is only as much as that of computing one gradient w.r.t. U. Both steps need time¹ of order mqnr. This is only r times more than "linear time" (time needed to read the algorithm inputs, here y_k, A_k 's). To our knowledge, r-times linear-time is the best known time complexity for any algorithm for any LR matrix recovery problem. Moreover, due to the use of the X = UB factorization, AltGDmin is also communicationefficient. Each node needs to only send nr scalars (gradients w.r.t U) at each iteration.

III. NEW GUARANTEE

Let m_0 denote the total number of samples per column needed for initialization and let m_1 denote this number for each GDmin iteration. Then, the total sample complexity per column is $m = m_0 + m_1 T$. Our guarantee given next provides the required minimum value of m.

Theorem 2: Assume that Assumption 1 holds. Set $\eta =$ $0.4/m\sigma_{\rm max}^{\star^2}$ and $T=C\kappa^2\log(1/\epsilon)$. If

$$mq \ge C\kappa^4 \mu^2 (n+q)r(\kappa^4 r + \log(1/\epsilon))$$

and $m \ge C \max(\log n, \log q, r) \log(1/\epsilon)$, then, with probability (w.p.) at least $1 - n^{-10}$,

$$\mathrm{SD}_2(\boldsymbol{U}, \boldsymbol{U}^\star) \leq \epsilon$$
 and $\|\boldsymbol{x}_k - \boldsymbol{x}_k^\star\| \leq \epsilon \|\boldsymbol{x}_k^\star\|$ for all $k \in [q]$.

The time complexity is $mqnr \cdot T = mqnr \cdot \kappa^2 \log(1/\epsilon)$. The communication complexity is $nr \cdot T = nr \cdot \kappa^2 \log(1/\epsilon)$ per node.

Proof: We prove the three results needed for proving this in Sections IV-B, IV-D, and V below. We use these to prove the above result in Appendix A.

Algorithm 1 The AltGD-Min Algorithm

- 1: **Input:** $y_k, A_k, k \in [q]$
- 2: Sample-split: Partition the measurements and measurement matrices into 2T + 2 equal-sized disjoint sets: two sets for initialization and 2T sets for the iterations. Denote these by $\mathbf{y}_{k}^{(\tau)}, \mathbf{A}_{k}^{(\tau)}, \tau = 00, 0, 1, \dots 2T.$

- these by $\boldsymbol{y}_{k}^{\times}$, \boldsymbol{A}_{k} ,

 3: Initialization:

 4: Using $\boldsymbol{y}_{k} \equiv \boldsymbol{y}_{k}^{(00)}$, $\boldsymbol{A}_{k} \equiv \boldsymbol{A}_{k}^{(00)}$,

 5: set $\alpha \leftarrow \tilde{C} \frac{1}{mq} \sum_{ki} |\boldsymbol{y}_{ki}|^{2}$,

 6: Using $\boldsymbol{y}_{k} \equiv \boldsymbol{y}_{k}^{(0)}$, $\boldsymbol{A}_{k} \equiv \boldsymbol{A}_{k}^{(0)}$,

 7: set $\boldsymbol{y}_{k,trunc}(\alpha) \leftarrow \boldsymbol{y}_{k,trnc}$:= trunc $(\boldsymbol{y}_{k},\alpha)$,

 8: set $\boldsymbol{X}_{0} \leftarrow (1/m) \sum_{k \in [q]} \boldsymbol{A}_{k}^{\top} \boldsymbol{y}_{k,trunc}(\alpha) \boldsymbol{e}_{k}^{\top}$
- 10: **GDmin iterations:**
- 11: **for** t = 1 **to** T **do**
- 12:
- 13:
- Let $oldsymbol{U} \leftarrow oldsymbol{U}_{t-1}.$ Using $oldsymbol{y}_k \equiv oldsymbol{y}_k^{(t)}, oldsymbol{A}_k \equiv oldsymbol{A}_k^{(t)},$ set $oldsymbol{b}_k \leftarrow (oldsymbol{A}_k oldsymbol{U})^\dagger oldsymbol{y}_k, \, oldsymbol{x}_k \leftarrow oldsymbol{U} oldsymbol{b}_k$ for all $k \in [q]$ Using $oldsymbol{y}_k \equiv oldsymbol{y}_k^{(T+t)}, oldsymbol{A}_k \equiv oldsymbol{A}_k^{(T+t)},$ compute set $abla_k oldsymbol{V}_k (oldsymbol{U}, oldsymbol{B}) = \sum_k oldsymbol{A}_k^\top (oldsymbol{A}_k oldsymbol{U} oldsymbol{b}_k oldsymbol{y}_k) oldsymbol{b}_k^\top$ 14:
- 15:
- 16:
- set $\hat{\boldsymbol{U}}^+ \leftarrow \boldsymbol{U} (\eta/m) \nabla_{\boldsymbol{U}} f(\boldsymbol{U}, \boldsymbol{B}_t)$. compute $\hat{\boldsymbol{U}}^+ \stackrel{\mathrm{QR}}{=} \boldsymbol{U}^+ \boldsymbol{R}^+$. 17:
- 18:
- Set $U_t \leftarrow U^+$.
- 20: end for

A. Discussion

We use $a \gtrsim b$ to mean that $a \geq C_{\kappa,\mu}b$ where $C_{\kappa,\mu}$ includes terms dependent on κ, μ . Most of our discussion treats κ, μ as numerical constants that do not grow with n, q, r. This assumption is borrowed from the rich past literature on various other LR matrix recovery problems, e.g., see [4] and [5]. Also, below, whp means w.p. at least $1 - n^{-10}$.

In this work (as well as in other works that use matrix factorization to solve LR recovery problems), the goal is to obtain a bound of the form $SD_2(U, U^*) \leq \epsilon$. If a bound on SD_F is needed, one can use $SD_F(U, U^*) \leq \sqrt{r}SD_2(U, U^*) \leq \sqrt{r}\epsilon$. Using Lemma 3 given below and triangle inequality, the bound on SD₂ helps get the bound $\|\boldsymbol{x}_k - \boldsymbol{x}_k^{\star}\| \lesssim \epsilon \|\boldsymbol{x}_k^{\star}\|$.

Our older result from [3] used $\mathrm{SD}_F(\boldsymbol{U},\boldsymbol{U}^\star)$ in its analysis and showed that, to guarantee $SD_F(U, U^*) \leq \epsilon$ whp, we need

$$mq > C\kappa^4 \mu^2 (n+q) r^2 (\kappa^4 + \log(1/\epsilon))$$

Since $\mathrm{SD}_2(\boldsymbol{U},\boldsymbol{U}^\star) \leq \mathrm{SD}_F(\boldsymbol{U},\boldsymbol{U}^\star)$, this implies that we need this same complexity also to guarantee $SD_2(\boldsymbol{U}, \boldsymbol{U}^*) \leq \epsilon$. Our new result needs $mq \geq C\kappa^4\mu^2(n+q)r(\kappa^4r+\log(1/\epsilon))$ to guarantee $\mathrm{SD}_2(\boldsymbol{U},\boldsymbol{U}^\star) \leq \epsilon$. Thus, this new result improves the dependence on r, ϵ from order $r^2 \log(1/\epsilon)$ to order $r \max(r, \log(1/\epsilon))$. This is an improvement over the old result by a factor of $\min(r, \log(1/\epsilon))$. This improvement is obtained because our new guarantee for the GD step (Theorem 6) only needs $mq \gtrsim nr$ at each iteration. On the other hand, the older result needed $mq \gtrsim nr^2$. Both guarantees need $mq \gtrsim nr^2$ for initialization, see Theorem 13.

We are able to improve our result because we now use a simpler proof technique that works for the LRCS problem, but

¹The LS step time is $\max(q \cdot mnr, q \cdot mr^2) = mqnr$ (maximum of the time needed for computing $\hat{A_k U}$ for all k, and that for obtaining b_k for all k) while the GD step time is $\max(q \cdot mnr, nr^2) = mqnr$ (maximum of the time needed for computing the gradient w.r.t. U, and time for the QR step).

not for its LR phase retrieval (LRPR) generalization. LRPR involves recovering X^* from $z_k := |y_k|, k \in [q]$. In [3], we were attempting to solve both problems. There are two differences in our new proof compared with that of [3]: (i) we use the 2-norm subspace distance $\mathrm{SD}_2(\boldsymbol{U},\boldsymbol{U}^\star)$ instead of $SD_F(U, U^*)$, and (ii) we do not use the fundamental theorem of calculus [6], [7] for analyzing the GD step, but instead use a much simpler direct approach. If we only use (ii) but not (i), we will still get a simpler proof, but we will not get the sample complexity gain. In [3], we used the Frobenius norm SD because it helped obtain the desired nr^3 guarantee for LRPR. Also, in hindsight, the use of the fundamental theorem of calculus was unnecessary. It has been used in earlier work [7] for analyzing a GD based algorithm for standard PR and for LR matrix completion and that originally motivated us to adapt the same approach for LRCS and LRPR.

Both the current result and our old one have the same dependence on κ . After initialization, the sample complexity grows roughly as κ^4 . A similar dependence on κ also exists in all past sample complexity guarantees for iterative - AlMin or GD – algorithms for LR matrix sensing, LR matrix completion or robust PCA, e.g., see [4] and [5]; and also for our older work on AltMin for a generalization of LRCS, LR phase retrieval [8]. One way to improve this dependence is to develop a stagewise algorithm similar to that introduced in [4] and [9] for LR matrix sensing and completion respectively. Developing this for the current LRCS problem is an open future work question. A second option is to use the convex relaxation approaches which usually depend on lower powers of κ . However, these need a much higher iteration complexity, making them much slower.

IV. NEW PROOF: GDMIN ITERATIONS

A. Definitions and Preliminaries

Let U be the estimate at the t-th iteration. Define

$$egin{aligned} oldsymbol{g}_k &:= oldsymbol{U}^ op oldsymbol{x}_k^\star, k \in [q], \ ext{and} \ oldsymbol{G} := oldsymbol{I} - oldsymbol{U}^ op oldsymbol{U}^{\star op}, \ ext{GradU} &:=
abla_{oldsymbol{U}} f(oldsymbol{U}, oldsymbol{B}) = \sum_k oldsymbol{A}_k^ op (oldsymbol{A}_k oldsymbol{U} oldsymbol{b}_k - oldsymbol{y}_k) oldsymbol{b}_k^ op \ &= \sum_{k:i} (oldsymbol{y}_{ki} - oldsymbol{a}_{ki}^ op oldsymbol{U} oldsymbol{b}_k) oldsymbol{a}_{ki} oldsymbol{b}_k^ op \end{aligned}$$

For an $n_1 \times n_2$ matrix, \boldsymbol{Z} , $\sigma_{\min}(\boldsymbol{Z}) = \sigma_{\min}(\boldsymbol{Z}^{\top}) =$ $\sigma_{\min(n_1,n_2)}(\mathbf{Z})$. Thus, if \mathbf{A} is tall, then $\sigma_{\min}(\mathbf{A}) =$ $\sqrt{\lambda_{\min}(\boldsymbol{A}^{\top}\boldsymbol{A})}$. Using this, it follows that, if $\boldsymbol{A} = \boldsymbol{B}\boldsymbol{C}$ and \boldsymbol{A} and \boldsymbol{B} are tall (or square), then $\sigma_{\min}(\boldsymbol{A}) \geq \sigma_{\min}(\boldsymbol{B})\sigma_{\min}(\boldsymbol{C})$.

B. Minimization Step

Assume $SD_2(\boldsymbol{U}, \boldsymbol{U}^*) \leq \delta_t$ with $\delta_t < 0.02$.

We use the following lemma from [3].

Lemma 3 [3]: Let $g_k := U^{\top} x_k^{\star}$. Then, w.p. at least $1 - \exp(\log q + r - c\epsilon m)$, for all $k \in [q]$,

$$\|oldsymbol{g}_k - oldsymbol{b}_k\| \leq 1.2\epsilon \|\left(oldsymbol{I}_n - oldsymbol{U}oldsymbol{U}^ opoldsymbol{b}_k^\star\|$$

 $^2 \text{Our older guarantee for LRPR [3] needed } mq \gtrsim nr^2 (r + \log(1/\epsilon)).$ If we use the new approach developed here, it will need $mq \gtrsim nr(r^3 + \log(1/\epsilon))$.

Proof: We provide a proof of this lemma in Appendix B to emphasise a slightly more general version of this lemma. In particular, our proof shows that we can replace 0.4 in the previous version of this lemma by any $\epsilon > 0$, and the bound holds w.p. at least $1 - \exp(\log q + r - c\epsilon m)$.

By Lemma 3 with $\epsilon = 0.3$, if $m \gtrsim \max(\log q, \log n, r)$, then, whp, for all $k \in [q]$,

$$\|\boldsymbol{b}_k - \boldsymbol{g}_k\| \le 0.4 \|(\boldsymbol{I} - \boldsymbol{U}\boldsymbol{U}^\top)\boldsymbol{U}^* \boldsymbol{b}_k^*\|$$

Using $SD_2(\boldsymbol{U}, \boldsymbol{U}^*) \leq \delta_t$, this directly implies that

- 1) $\|\boldsymbol{b}_k \boldsymbol{g}_k\| \le 0.4\delta_t \|\boldsymbol{b}_k^{\star}\|$ 2) $\|\boldsymbol{b}_k\| \le \|\boldsymbol{g}_k\| + 0.4 \cdot 0.02 \|\boldsymbol{b}_k^{\star}\| \le 1.1 \|\boldsymbol{b}_k^{\star}\|$
- 3) $\|\boldsymbol{x}_k \boldsymbol{x}_k^{\star}\| \leq 1.4\delta_t \|\boldsymbol{b}_k^{\star}\|$

(to obtain the last bound, we need to add and subtract Ug_k and then use triangle inequality). Using these bounds,

$$\|\boldsymbol{B} - \boldsymbol{G}\|_F \le 0.4\delta_t \sqrt{\sum_k \|\boldsymbol{b}_k^{\star}\|^2} = 0.4\delta_t \sqrt{r} \sigma_{\max}^{\star}$$

and we can get a similar bound on $\|X - X^*\|_F$. Thus,

- 1) $\| \boldsymbol{B} \boldsymbol{G} \|_F \le 0.4 \delta_t \| \boldsymbol{B}^* \|_F \le 0.4 \sqrt{r} \delta_t \sigma_{\text{max}}^*$
- 2) $\|\boldsymbol{X} \boldsymbol{X}^{\star}\|_{F} \leq 1.4\sqrt{r}\delta_{t}\sigma_{\max}^{\star}$

Furthermore,

$$\sigma_{\min}(\boldsymbol{B}) \ge \sigma_{\min}(\boldsymbol{G}) - \|\boldsymbol{B} - \boldsymbol{G}\| \ge \sigma_{\min}(\boldsymbol{G}) - \|\boldsymbol{B} - \boldsymbol{G}\|_F$$

To lower bound $\sigma_{\min}(G)$, observe that

$$\sigma_{\min}(\boldsymbol{G}) = \sigma_{\min}(\boldsymbol{G}^{\top}) \geq \sigma_{\min}^{\star} \sigma_{\min}(\boldsymbol{U}^{\star \top} \boldsymbol{U})$$

and

$$\sigma_{\min}(\boldsymbol{U}^{\star \top} \boldsymbol{U}) = \sqrt{1 - \|\boldsymbol{P}\boldsymbol{U}\|^2} \ge \sqrt{1 - \delta_t^2}.$$

This follows using $\sigma_{\min}^2(\boldsymbol{U}^{\star \top}\boldsymbol{U}) = \lambda_{\min}(\boldsymbol{U}^{\top}\boldsymbol{U}^{\star}\boldsymbol{U}^{\star \top}\boldsymbol{U}) = \lambda_{\min}(\boldsymbol{U}^{\top}(\boldsymbol{I} - \boldsymbol{P})\boldsymbol{U}) = \lambda_{\min}(\boldsymbol{I} - \boldsymbol{U}^{\top}\boldsymbol{P}\boldsymbol{U}) = \lambda_{\min}(\boldsymbol{I} - \boldsymbol{U}^{\top}\boldsymbol{P}^2\boldsymbol{U}) = 1 - \lambda_{\max}(\boldsymbol{U}^{\top}\boldsymbol{P}^2\boldsymbol{U}) = 1 - \|\boldsymbol{P}\boldsymbol{U}\|^2.$

Combining the above three bounds and using the bound on $\|\boldsymbol{B} - \boldsymbol{G}\|_F$, if $\delta_t < 0.02/\sqrt{r}\kappa$, then

$$\sigma_{\min}(\boldsymbol{B}) \ge \sqrt{1 - \delta_t^2} \sigma_{\min}^{\star} - 0.4 \sqrt{r} \delta_t \sigma_{\max}^{\star} \ge 0.9 \sigma_{\min}^{\star}$$

$$\sigma_{\max}(\boldsymbol{B}) \le \|\boldsymbol{G}\| + 0.4\sqrt{r}\delta_t \sigma_{\max}^{\star} \le 1.1\sigma_{\max}^{\star}$$

since $\|G\| \leq \|B^\star\| = \sigma_{\max}^\star$. Thus, we have proved the following claim:

Theorem 4: Assume that $SD_2(\boldsymbol{U}, \boldsymbol{U}^{\star}) \leq \delta_t$. If $\delta_t \leq$ $0.02/\sqrt{r}\kappa$, and if $m \gtrsim \max(\log q, \log n, r)$, then whp,

- 1) $\|\boldsymbol{b}_k \boldsymbol{g}_k\| \le 0.4\delta_t \|\boldsymbol{b}_k^{\star}\|$
- 2) $\|\boldsymbol{b}_{k}\| \le \|\boldsymbol{g}_{k}\| + 0.4 \cdot 0.02 \|\boldsymbol{b}_{k}^{\star}\| \le 1.1 \|\boldsymbol{b}_{k}^{\star}\|$ 3) $\|\boldsymbol{B} \boldsymbol{G}\|_{F} \le 0.4\delta_{t} \|\boldsymbol{B}^{\star}\|_{F} \le 0.4\sqrt{r}\delta_{t}\sigma_{\max}^{\star}$
- 4) $\|\boldsymbol{x}_k \boldsymbol{x}_k^{\star}\| \le 1.4\delta_t \|\boldsymbol{b}_k^{\star}\|$
- 5) $\|\boldsymbol{X} \boldsymbol{X}^{\star}\|_F \le 1.4\sqrt{r}\delta_t \sigma_{\max}^{\star}$
- 6) $\sigma_{\min}(\boldsymbol{B}) \ge 0.9 \sigma_{\min}^{\star}$
- 7) $\sigma_{\max}(\boldsymbol{B}) \leq 1.1 \sigma_{\max}^{\star}$

(only the last two bounds require the upper bound on δ_t).

C. New Bounds on the Expected Gradient and Deviation From it

Using independence of A_k and $\{U, b_k\}$ (due to sample splitting),

$$\mathbb{E}[\operatorname{GradU}] = \sum_k m(\boldsymbol{x}_k - \boldsymbol{x}_k^{\star}) \boldsymbol{b}_k^{\top}$$

Using bounds on ||B|| and $||X^* - X||_F$ from Theorem 4, if $\delta_t < \frac{c}{\sqrt{r}\kappa}$,

$$\|\mathbb{E}[\text{GradU}]\| = \|\sum_{k} m(\boldsymbol{x}_{k} - \boldsymbol{x}_{k}^{\star}) \boldsymbol{b}_{k}^{\top} \| = m \|(\boldsymbol{X} - \boldsymbol{X}^{\star}) \boldsymbol{B}^{\top} \|$$

$$\leq m \|\boldsymbol{X} - \boldsymbol{X}^{\star} \| \cdot \|\boldsymbol{B} \|$$

$$\leq m \|\boldsymbol{X} - \boldsymbol{X}^{\star} \|_{F} \cdot \|\boldsymbol{B} \| \leq 1.1 m \delta_{t} \sqrt{r} \sigma_{\max}^{\star}^{2}.$$

w.p. $1 - \exp(\log q + r - cm)$.

Next, we bound $\|\operatorname{GradU} - \mathbb{E}[\operatorname{GradU}]\|$ $\max_{\|oldsymbol{w}\|=1,\|oldsymbol{z}\|=1} oldsymbol{w}^{ op}(\sum_k \sum_i oldsymbol{a}_{ki} oldsymbol{a}_{ki}^{ op}(oldsymbol{x}_k - oldsymbol{x}_k^{\star}) oldsymbol{b}_k^{ op} - \mathbb{E}[\cdot]) oldsymbol{z}.$ This also uses independence of A_k and $\{U, b_k\}$.

We bound the above for fixed unit norm w, z using subexponential Bernstein inequality (Theorem 2.8.1 of [10]). We extend the bound to all unit norm w, z by using a standard epsilon-net argument. For fixed unit norm w, z, consider

$$\sum_k \sum_i \left((\boldsymbol{w}^\top \boldsymbol{a}_{ki}) (\boldsymbol{b}_k^\top \boldsymbol{z}) \boldsymbol{a}_{ki}^\top (\boldsymbol{x}_k - \boldsymbol{x}_k^\star) - \mathbb{E}[\cdot] \right)$$

Observe that the summands are independent, zero mean, subexponential r.v.s with sub-exponential norm $K_{ki} \leq C \|\boldsymbol{w}\|$. $\|oldsymbol{z}^{ op}oldsymbol{b}_k\|\cdot\|oldsymbol{x}_k-oldsymbol{x}_k^\star\|=C\|oldsymbol{z}^{ op}oldsymbol{b}_k\|\cdot\|oldsymbol{x}_k-oldsymbol{x}_k^\star\|.$ We apply the sub-exponential Bernstein inequality, Theorem 2.8.1 of [10], with $t = \epsilon_1 \delta_t m \sigma_{\min}^*$. We have

$$\begin{split} \frac{t^2}{\sum_{ki} K_{ki}^2} &\geq \frac{\epsilon_1^2 \delta_t^2 m^2 {\sigma_{\min}^{\star}}^4}{m \sum_k \| \boldsymbol{x}_k - \boldsymbol{x}_k^{\star} \|^2 (\boldsymbol{z}^{\top} \boldsymbol{b}_k)^2} \\ &\geq \frac{\epsilon_1^2 \delta_t^2 m^2 {\sigma_{\min}^{\star}}^4}{m \max_k \| \boldsymbol{x}_k - \boldsymbol{x}_k^{\star} \|^2 \sum_k (\boldsymbol{z}^{\top} \boldsymbol{b}_k)^2} \\ &= \frac{\epsilon_1^2 \delta_t^2 m^2 {\sigma_{\min}^{\star}}^4}{m \max_k \| \boldsymbol{x}_k - \boldsymbol{x}_k^{\star} \|^2 \| \boldsymbol{z}^{\top} \boldsymbol{B} \|^2} \\ &\geq \frac{\epsilon_1^2 \delta_t^2 m^2 {\sigma_{\min}^{\star}}^4 q}{1.4^2 m \mu^2 \delta_t^2 r {\sigma_{\max}^{\star}}^2 \| \boldsymbol{B} \|^2} \\ &\geq \frac{\epsilon_1^2 \delta_t^2 m^2 {\sigma_{\min}^{\star}}^4 q}{1.4^2 m \mu^2 \delta_t^2 r {\sigma_{\max}^{\star}}^2 1.1 {\sigma_{\max}^{\star}}^2} = c \frac{\epsilon_1^2 m q}{\kappa^4 \mu^2 r} \\ \frac{t}{\max_k K_{ki}} &\geq \frac{\epsilon_1 \delta_t m {\sigma_{\min}^{\star}}^2}{\max_k \| \boldsymbol{b}_k \| \max_k \| \boldsymbol{x}_k - \boldsymbol{x}_k^{\star} \|} \geq c \frac{\epsilon_1 m q}{\kappa^2 \mu^2 r} \end{split}$$

In the above, we used (i) $\sum_k (oldsymbol{z}^ op oldsymbol{b}_k)^2 = \|oldsymbol{z}^ op oldsymbol{B}\|^2 \leq \|oldsymbol{B}\|^2$ since z is unit norm, (ii) Theorem 4 to bound $\|B\| \le$ $1.1\sigma_{\mathrm{max}}^{\star}$, and (iii) Theorem 4 followed by Assumption 1 (right incoherence) to bound $\|m{x}_k - m{x}_k^\star\| \leq \delta_t \cdot \mu \sigma_{\max}^\star \sqrt{r/q}$ and $|\boldsymbol{z}^{\top}\boldsymbol{b}_{k}| \leq \|\boldsymbol{b}_{k}\| \leq 1.1\|\boldsymbol{b}_{k}^{\star}\| \leq 1.1\mu\sigma_{\max}^{\star}\sqrt{r/q}.$

For $\epsilon_1 < 1$, the first term above is smaller (since $1/\kappa^4 \le 1/\kappa^2$), i.e., $\min(\frac{t^2}{\sum_{ki} K_{ki}^2}, \frac{t}{\max_{ki} K_{ki}}) = c\frac{\epsilon_1^2 mq}{\kappa^4 \mu^2 r}$. Thus, by sub-exponential Bernstein, w.p. at least $1-\exp(-c\frac{\epsilon_1^2 mq}{\kappa^4 u^2 r})$ — $\exp(\log q + r - cm)$, for a given w, z,

$$\boldsymbol{w}^{\top}(\operatorname{GradU} - \mathbb{E}[\operatorname{GradU}])\boldsymbol{z} \leq \epsilon_1 \delta_t m \sigma_{\min}^{\star}^{2}$$

Using a standard epsilon-net argument to bound the maximum of the above over all unit norm w, z, e.g., using [3, Proposition 4.7], we can conclude that

$$\|\operatorname{GradU} - \mathbb{E}[\operatorname{GradU}]\| \le 1.1\epsilon_1 \delta_t m \sigma_{\min}^{\star^2}$$

w.p. at least $1-\exp(C(n+r)-c\frac{\epsilon_1^2mq}{\kappa^4\mu^2r})-\exp(\log q+r-cm)$. The factor of $\exp(C(n+r))$ is due to the epsilon-net over ${\pmb w}$ and that over z: w is an n-length unit norm vector while z is an r-length unit norm vector. The smallest epsilon net covering the hyper-sphere of all ws is of size $(1+2/\epsilon_{net})^n=C^n$ with $\epsilon_{net} = c$ while that for z is of size C^r . Union bounding over both thus gives a factor of C^{n+r} . By replacing ϵ_1 by $\epsilon_1/1.1$, our bound becomes simpler (and $1/1.1^2$ gets incorporated into the factor c). We have thus proved the following.

Lemma 5: Assume that $SD_2(U, U^*) \leq \delta_t$. The following

- 1) $\mathbb{E}[\text{GradU}] = m(\boldsymbol{X} \boldsymbol{X}^*)\boldsymbol{B}^{\top} = m(\boldsymbol{U}\boldsymbol{B}\boldsymbol{B}^{\top} \boldsymbol{X}^*\boldsymbol{B}^{\top})$ 2) $\|\mathbb{E}[\text{GradU}]\| \le 1.1m\delta_t\sqrt{r}\sigma_{\max}^*$
- 3) If $\delta_t < \frac{c}{\sqrt{r}\kappa}$, then, w.p. at least $1 - \exp(C(n+r) - c\frac{\epsilon_1^2 mq}{\kappa^4 u^2 r}) - \exp(\log q + r - cm),$

$$\|\operatorname{GradU} - \mathbb{E}[\operatorname{GradU}]\| \le \epsilon_1 \delta_t m \sigma_{\min}^{\star}^2$$

The above lemma is an improvement over the bounds given in [3] because δ_t is now the bound on the 2-norm SD, and still it only needs $mq \gtrsim nr/\epsilon_1^2$.

D. GD Step

Assume that $SD_2(\boldsymbol{U}, \boldsymbol{U}^*) \leq \delta_t$ with $\delta_t < 0.02$. Recall the Projected GD step for U:

$$\widetilde{m{U}}^+ = m{U} - \eta ext{GradU}$$
 and $\widetilde{m{U}}^+ \stackrel{ ext{QR}}{=} m{U}^+ m{R}^+$

Since $U^+ = \widetilde{U}^+(R^+)^{-1}$ and since $\|(R^+)^{-1}\|$ $1/\sigma_{\min}(\boldsymbol{R}^+) = 1/\sigma_{\min}(\widetilde{\boldsymbol{U}}^+)$, thus, $\mathrm{SD}_2(\boldsymbol{U}^+, \boldsymbol{U}^\star)$ $\| \boldsymbol{P} \boldsymbol{U}^{+} \|$ can be bounded as

$$\mathrm{SD}_2(\boldsymbol{U}^+, \boldsymbol{U}^\star) \leq \frac{\|\boldsymbol{P}\widetilde{\boldsymbol{U}}^+\|}{\sigma_{\min}(\widetilde{\boldsymbol{U}}^+)} \leq \frac{\|\boldsymbol{P}\widetilde{\boldsymbol{U}}^+\|}{\sigma_{\min}(\boldsymbol{U}) - \eta\|\mathrm{GradU}\|}$$
 (2)

Consider the numerator. Adding/subtracting $\mathbb{E}[GradU]$, left multiplying both sides by P, and using Lemma 5 (first part),

$$\tilde{\boldsymbol{U}}^+ = \boldsymbol{U} - \eta \mathbb{E}[\operatorname{GradU}] + \eta(\mathbb{E}[\operatorname{GradU}] - \operatorname{GradU}), \text{ thus,}$$

$$\boldsymbol{P}\tilde{\boldsymbol{U}}^+ = \boldsymbol{P}\boldsymbol{U} - \eta m \boldsymbol{P}\boldsymbol{U} \boldsymbol{B} \boldsymbol{B}^\top + \eta \boldsymbol{P}(\mathbb{E}[\operatorname{GradU}] - \operatorname{GradU})$$

The last row used $PX^* = 0$. Thus,

$$\|\boldsymbol{P}\widetilde{\boldsymbol{U}}^{\dagger}\| \le \|\boldsymbol{P}\boldsymbol{U}\|\|\boldsymbol{I} - \eta m\boldsymbol{B}\boldsymbol{B}^{\top}\| + \eta\|\mathbb{E}[\operatorname{GradU}] - \operatorname{GradU}\|$$
(3)

Using Theorem 4, we get

$$\lambda_{\min}(\boldsymbol{I} - \eta m \boldsymbol{B} \boldsymbol{B}^{\top}) = 1 - \eta m \|\boldsymbol{B}\|^2 \ge 1 - 1.2 \eta m \sigma_{\max}^{\star 2}$$

Thus, if $\eta < 0.5/m\sigma_{\rm max}^{\star~2}$, then the above matrix is p.s.d. This along with Theorem 4 then implies that

$$\|\boldsymbol{I} - \eta m \boldsymbol{B} \boldsymbol{B}^{\top}\| = \lambda_{\max}(\boldsymbol{I} - \eta m \boldsymbol{B} \boldsymbol{B}^{\top}) \le 1 - 0.9 \eta m \sigma_{\min}^{\star 2}$$

Using the above, (3), and the bound on $||\mathbb{E}[\text{GradU}]|$ – GradU|| from Lemma 5, we conclude the following: If $\eta \leq$ $0.5/m\sigma_{\max}^{\star^2}$, and $\delta_t \leq c/\sqrt{r}\kappa$, then

$$\|\boldsymbol{P}\tilde{\boldsymbol{U}}^{+}\| \leq \|\boldsymbol{P}\boldsymbol{U}\| \|\boldsymbol{I} - \eta m\boldsymbol{B}\boldsymbol{B}^{\top}\| + !\eta \|\mathbb{E}[\operatorname{GradU}] - \operatorname{GradU}\|$$

$$\leq \delta_{t}(1 - 0.9 \ \eta m\sigma_{\min}^{\star}^{2}) + \eta m\epsilon_{1}\delta_{t}\sigma_{\min}^{\star}^{2}$$
(4)

w.p. at least $1-\exp(C(n+r)-c\frac{\epsilon_1^2mq}{\kappa^4\mu^2r})-\exp(\log q+r-cm)$. This probability is at least $1-n^{-10}$ if $mq\gtrsim \kappa^4\mu^2nr/\epsilon_1^2$ and $m \gtrsim \max(\log n, \log q, r).$

Next we use (4) with $\epsilon_1 = 0.1$ and Lemma 5 to bound the right hand side of (2). Set $\eta = c_{\eta}/m\sigma_{\max}^{\star}^{2}$. If $c_{\eta} \leq 0.5$, if $\delta_t \leq c/\sqrt{r}\kappa^2$, and lower bounds on m from above hold, (2) implies that, whp,

$$\begin{split} &\mathrm{SD}_{2}(\boldsymbol{U}^{+}, \boldsymbol{U}^{\star}) \\ &\leq \frac{\|\boldsymbol{P}\widetilde{\boldsymbol{U}}^{+}\|}{\sigma_{\min}(\boldsymbol{U}) - \eta\|\mathrm{GradU}\|} \\ &\leq \frac{\|\boldsymbol{P}\boldsymbol{U}^{+}\|}{\sigma_{\min}(\boldsymbol{U}) - \eta\|\mathrm{GradU}\|} \\ &\leq \frac{\|\boldsymbol{P}\boldsymbol{U}\|\|\boldsymbol{I} - \eta m\boldsymbol{B}\boldsymbol{B}^{\top}\| + \eta\|\mathbb{E}[\mathrm{GradU}] - \mathrm{GradU}\|}{1 - \eta\|\mathbb{E}[\mathrm{GradU}]\| - \eta\|\mathrm{GradU} - \mathbb{E}[\mathrm{GradU}]} \\ &\leq \frac{\delta_{t}(1 - \eta m\sigma_{\min}^{\star}{}^{2}(0.9 - 0.1))}{1 - \eta\|\mathbb{E}[\mathrm{GradU}]\| - \eta\|\mathrm{GradU} - \mathbb{E}[\mathrm{GradU}]\|} \\ &\leq \frac{\delta_{t}(1 - 0.8\eta m\sigma_{\min}^{\star}{}^{2})}{1 - \eta m\delta_{t}\sqrt{r}\sigma_{\max}^{\star}{}^{2}(1.4 + \frac{0.1}{\kappa^{2}\sqrt{r}})} \\ &\leq \frac{\delta_{t}(1 - 0.8\eta m\sigma_{\min}^{\star}{}^{2})}{1 - 1.5\eta m\delta_{t}\sqrt{r}\sigma_{\max}^{\star}{}^{2}} \\ &\leq \delta_{t}(1 - 0.8\eta m\sigma_{\min}^{\star}{}^{2})(1 + 1.5\eta m\delta_{t}\sqrt{r}\sigma_{\max}^{\star}{}^{2}) \\ &\leq \delta_{t}(1 - \eta m\sigma_{\min}^{\star}{}^{2}(0.8 - 1.5\delta_{t}\sqrt{r}\kappa^{2})) \\ &\leq \delta_{t}(1 - \eta m\sigma_{\min}^{\star}{}^{2}(0.8 - 0.15)) \\ &\leq \delta_{t}(1 - 0.6\eta m\sigma_{\max}^{\star}{}^{2}/\kappa^{2}) = \delta_{t}(1 - 0.6c_{\eta}/\kappa^{2}) \end{split}$$

In the above we used $\kappa^2 \sqrt{r} > 1$, $(1-x)^{-1} < (1+x)$ if |x| < 1, and $\delta_t < 0.1/\sqrt{r}\kappa^2$ (used in second-last inequality). Thus, we have proved the following result.

Theorem 6: Assume that Assumption 1 holds and $\mathrm{SD}_2(\boldsymbol{U},\boldsymbol{U}^\star) \leq \delta_t$. If $\delta_t \leq 0.02/\sqrt{r}\kappa^2$, if $\eta = c_\eta/m\sigma_{\mathrm{max}}^\star$ with $c_\eta \leq 0.5$, and if $mq \gtrsim \kappa^4\mu^2 nr$ and $m \gtrsim \max(\log n, \log q, r)$, then, whp

$$SD_2(U^+, U^*) \le \delta_{t+1} := \delta_t (1 - 0.6c_n/\kappa^2)$$

V. NEW PROOF: INITIALIZATION

We need a new result for the initialization step because we need a tight bound on $SD_2(U_0, U^*)$.

A. Results Taken From [3]

Recall from Algorithm 1 that α uses a different set of measurements that is independent of those used for X_0 . We use the following four results from [3].

Lemma 7 [3]: Conditioned on α , we have the following conclusions. Let ζ be a scalar standard Gaussian r.v.. Define

$$\beta_k(\alpha) := \mathbb{E}[\zeta^2 \mathbb{1}_{\{\|\boldsymbol{x}_k^{\star}\|^2 \zeta^2 < \alpha\}}].$$

Then,

$$\mathbb{E}[\boldsymbol{X}_0|\alpha] = \boldsymbol{X}^*\boldsymbol{D}(\alpha),$$
where $\boldsymbol{D}(\alpha) := diagonal(\beta_k(\alpha), k \in [q])$

i.e. $D(\alpha)$ is a diagonal matrix of size $q \times q$ with diagonal entries β_k defined above.

Fact 8 [3]: Let

$$\mathcal{E} := \left\{ \tilde{C}(1 - \epsilon_1) \frac{\|\boldsymbol{X}^{\star}\|_F^2}{q} \leq \alpha \leq \tilde{C}(1 + \epsilon_1) \frac{\|\boldsymbol{X}^{\star}\|_F^2}{q} \right\}.$$

$$\begin{split} & \Pr(\alpha \in \mathcal{E}) \geq 1 - \exp(-\tilde{c}mq\epsilon_1^2). \text{ Here } \tilde{c} = c/\tilde{C} = c/\kappa^2\mu^2. \\ & \textit{Fact 9} \quad \textit{[3]:} \quad \text{For any} \quad \epsilon_1 \leq 0.1, \\ & \min_k \mathbb{E} \left[\zeta^2 \mathbb{I}_{\left\{ |\zeta| \leq \tilde{C} \frac{\sqrt{1-\epsilon_1} \|\mathbf{X}^\star\|_F}{\sqrt{q} \|\boldsymbol{x}_k^\star\|_F} \right\}} \right] \geq 0.92. \\ & \textit{Lemma 10} \quad \textit{[3]:} \quad \text{Fix } 0 < \epsilon_1 < 1. \text{ Then, w.p. at least } 1 - \epsilon_1 < \epsilon_2 < \epsilon_3 < \epsilon_4 < 1. \end{split}$$

 $\exp\left[(n+q)-c\epsilon_1^2mq/\mu^2\kappa^2\right]$, conditioned on α , for an $\alpha\in\mathcal{E}$,

$$\|\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha]\| \le 1.1\epsilon_1 \|\boldsymbol{X}^{\star}\|_F$$

Facts 8 and 9 together imply that, w.p. at least 1 - $\exp(-\tilde{c}mq\epsilon_1^2),$

$$\min_{k} \beta_{k}(\alpha) \ge \min_{k} \mathbb{E} \left[\zeta^{2} \mathbb{1}_{\left\{ |\zeta| \le \tilde{C} \frac{\sqrt{1-\epsilon_{1}} \|\mathbf{X}^{\star}\|_{F}}{\sqrt{q} \|\mathbf{x}^{\star}_{k}\|_{F}} \right\}} \right] \ge 0.92. \quad (5)$$

The first inequality is an immediate consequence of Fact 8 $(\alpha \in \mathcal{E})$ and the second follows by Fact 9.

By setting $\epsilon_1 = \epsilon_0/1.1\sqrt{r}\kappa$ in Lemma 10, and using Fact 8, we get the following corollary.

Corollary 11: Fix $0 < \epsilon_1 < 1$. Then, w.p. at least $1 - \epsilon_1 < \epsilon_2 < \epsilon_2 < \epsilon_3 < \epsilon_4$ $\exp\left[(n+q) - c\frac{\epsilon_0^2 mq}{\mu^2 \kappa^4 r}\right] - \exp(-cmq\epsilon_1^2/\kappa^2 \mu^2),$

$$\|\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha]\| \le \epsilon_0 \sigma_{\min}^{\star}$$

with $\mathbb{E}[X_0|\alpha]$ being as given on Lemma 7.

B. Obtaining the SD Error Bound for Initialization

By Lemma 7, $\mathbb{E}[X_0|\alpha] = X^*D(\alpha)$ with $D(\alpha)$ as defined there. Clearly, its rank is r or less. To obtain the bound, we apply Wedin's $\sin \theta$ theorem [11] for SD₂ with $M = X_0$, $M^* = \mathbb{E}[X_0|\alpha] = X^*D$. Recall that U_0 is the matrix of top r singular vectors of X_0 . Also, the span of top r singular vectors of $\mathbb{E}[X_0|\alpha] = X^*D$ equals the column span of U^* . Thus applying Wedin will help us bound $SD_2(U_0, U^*)$. To do this, we need to define the SVD of $\mathbb{E}[X_0|\alpha]$. Let

$$\mathbb{E}[\boldsymbol{X}_0|\alpha] = \boldsymbol{X}^{\star}\boldsymbol{D}(\alpha) \overset{\text{SVD}}{=} (\boldsymbol{U}^{\star}\boldsymbol{Q}) \check{\boldsymbol{\Sigma}^{\star}} \check{\boldsymbol{V}}$$

where Q is a $r \times r$ unitary matrix, $\check{\Sigma}^*$ is an $r \times r$ diagonal matrix with non-negative entries (singular values) and $\check{m{V}}$ is an $r \times q$ matrix with orthonormal rows, i.e. $\check{\boldsymbol{V}}\check{\boldsymbol{V}}^{\top} = \boldsymbol{I}$. Observe also that $\sigma_{r+1}(\mathbb{E}[\boldsymbol{X}_0|\alpha]) = 0$ since it is a rank r matrix. Also, from above,

$$\sigma_r(\mathbb{E}[\boldsymbol{X}_0|\alpha]) \!=\! \sigma_{\min}(\check{\boldsymbol{\Sigma}}^*) \geq \sigma_{\min}^{\star} \sigma_{\min}(\boldsymbol{D}) = \sigma_{\min}^{\star} \min_{k} \beta_k(\alpha)$$

This follows since $\check{\mathbf{\Sigma}}^* = \mathbf{Q}^{\top} \mathbf{B}^* \mathbf{D} \check{\mathbf{V}}^{\top}$ and so $\sigma_{\min}(\check{\mathbf{\Sigma}}^*) \geq \sigma_{\min}(\mathbf{B}^* \mathbf{D} \check{\mathbf{V}}^{\top}) \geq \sigma_{\min}(\mathbf{B}^* \mathbf{D}) \cdot 1 \geq \sigma_{\min}(\mathbf{D} \mathbf{B}^{*\top}) \geq \sigma_{\min}(\mathbf{D}) \cdot \sigma_{\min}^*$ and $\sigma_{\min}(\mathbf{D}) = \min_k \beta_k$. The last equality follows since \mathbf{D} is diagonal with entries β_k . Thus,

Fact 12: $\sigma_r(\mathbb{E}[\boldsymbol{X}_0|\alpha]) \geq \sigma_{\min}^* \min_k \beta_k(\alpha), \quad \sigma_{r+1}(\mathbb{E}[\boldsymbol{X}_0|\alpha]) = 0.$

Applying Wedin's $\sin \theta$ theorem, and then using Corollary 11, Fact 12 and (5), we get

 $SD_2(\boldsymbol{U}_0, \boldsymbol{U}^{\star})$

$$\leq \sqrt{2} \frac{\max(\|(\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha])^{\top} \boldsymbol{U}^{\star}\|, \|(\boldsymbol{X}_0 - \mathbb{E}[X_0|\alpha]) \boldsymbol{\check{V}}\|)}{\sigma_{\min}^{\star} \min_{k} \beta_k(\alpha) - 0 - \|\boldsymbol{X}_0 - \mathbb{E}[X_0|\alpha]\|} \\ \lesssim \sqrt{2} \frac{\epsilon_0 \sigma_{\min}^{\star}}{0.92 \sigma_{\min}^{\star} - \epsilon_0 \sigma_{\min}^{\star}} \lesssim \sqrt{2} \frac{\epsilon_0}{0.9} < 1.6 \epsilon_0$$

if $\epsilon_0 < 0.02$ and $mq \gtrsim (n+q)r/\epsilon_0^2$. In the above we used $\|(\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha])^\top \boldsymbol{U}^\star\| \le \|\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha]\| \cdot 1$ and $\|(\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha])\dot{\boldsymbol{V}}\| \le \|\boldsymbol{X}_0 - \mathbb{E}[\boldsymbol{X}_0|\alpha]\| \cdot 1$. Setting $\epsilon_0 = 0.5\delta_0$ with $\delta_0 < 0.02$, we obtain our desired result.

Theorem 13: Assume that Assumption 1 holds. Pick a $\delta_0 \le 0.02$. If $mq \gtrsim \kappa^4 \mu^2 (n+q) r/\delta_0^2$ then whp

$$SD_2(\boldsymbol{U}_0, \boldsymbol{U}^{\star}) \leq \delta_0$$

APPENDIX A PROOF OF THEOREM 2

Theorem 13 tells us that $\mathrm{SD}_2(\boldsymbol{U}_0, \boldsymbol{U}^\star) \leq \delta_0$ whp if $m_0 q \gtrsim (n+q)r/\delta_0^2$. Theorem 6 tells us that if $\delta_t \leq 0.02/\sqrt{r}\kappa^2$, and if $m_1 q \gtrsim nr$, then, whp, δ_t reduces by a factor of $(1-0.6c_\eta/\kappa^2)$ at each iteration. In particular, this implies that $\delta_t \leq \delta_0$. Thus, in order to apply Theorem 6, it suffices to require $\delta_0 = 0.02/\sqrt{r}\kappa^2$.

Combining both results, we have shown that if $m_0 q \ge C\kappa^8\mu^2(n+q)r^2$ and if $m_1 q \ge C\kappa^4mu^2nr$ and $m_1 \ge C\max(r,\log n,\log r)$, and if $\eta = c_\eta/(m\sigma_{\max}^*{}^2)$ with $c_\eta < 0.5$, then, whp, at each $t \ge 0$,

$$SD_2(\boldsymbol{U}_t, \boldsymbol{U}^*) \le \delta_t := (1 - 0.6 \frac{c_{\eta}}{\kappa^2})^t \delta_0 = (1 - 0.6 \frac{c_{\eta}}{\kappa^2})^t \frac{0.02}{\sqrt{r}\kappa^2}$$

and all bounds of Theorem 4 hold with δ_t as above.

Thus, to guarantee $\mathrm{SD}_2(\boldsymbol{U}_T, \boldsymbol{U}^\star) \leq \epsilon$, we need

$$T \ge C \frac{\kappa^2}{c_n} \log(1/\epsilon)$$

This follows by using $\log(1-x) < -x$ for |x| < 1 and using $\kappa^2 \sqrt{r} \ge 1$. Thus, setting $c_{\eta} = 0.4$, our sample complexity $m = m_0 + m_1 T$ becomes $mq \ge C \kappa^8 \mu^2 (n+q) r (1+\log(1/\epsilon))$, and $m \ge C \max(r, \log q, \log n) \log(1/\epsilon)$.

APPENDIX B PROOF OF LEMMA 3

Using the expression for b_k and simplifying it,

$$oldsymbol{b}_k - oldsymbol{g}_k = oldsymbol{M}^{-1} oldsymbol{U}^ op oldsymbol{A}_k^ op oldsymbol{A}_k (oldsymbol{I} - oldsymbol{U} oldsymbol{U}^ op) oldsymbol{U}^\star oldsymbol{b}_k^\star$$

with $\boldsymbol{M} := \boldsymbol{U}^{\top} \boldsymbol{A}_k^{\top} \boldsymbol{A}_k \boldsymbol{U}$. Clearly,

$$\mathbb{E}[\boldsymbol{M}] = m\boldsymbol{U}^{\top}\boldsymbol{U} = m\boldsymbol{I}_r, \ \mathbb{E}[\boldsymbol{U}^{\top}\boldsymbol{A}_k^{\top}\boldsymbol{A}_k(\boldsymbol{I} - \boldsymbol{U}\boldsymbol{U}^{\top})\boldsymbol{U}^{\star}\boldsymbol{b}_k^{\star}] = 0$$

Thus, using the standard technique (sub-expo Bern ineq followed by an epsilon-net argument), one can show that

1) w.p. at least $1 - \exp(r - \epsilon_2 m)$,

$$\|\boldsymbol{M} - \boldsymbol{I}_r\| \le 1.1\epsilon_2 m$$

This implies of course that $\sigma_{\min}^{\star}(\boldsymbol{M}) \geq (1 - 1.1\epsilon_2)m$ and thus $\|\boldsymbol{M}^{-1}\| \leq 1/((1 - 1.1\epsilon_2)m)$.

2) w.p. at least $1 - \exp(r - \epsilon_3 m)$,

$$\|\boldsymbol{U}^{\top}\boldsymbol{A}_{k}^{\top}\boldsymbol{A}_{k}(\boldsymbol{I}-\boldsymbol{U}\boldsymbol{U}^{\top})\boldsymbol{U}^{\star}\boldsymbol{b}_{k}^{\star}\|$$

$$\leq 1.1\epsilon_{3}m\|(\boldsymbol{I}-\boldsymbol{U}\boldsymbol{U}^{\top})\boldsymbol{U}^{\star}\boldsymbol{b}_{k}^{\star}\|$$

Thus, setting $\epsilon_2 = 0.1$, and taking union bound over all q vectors, w.p. at least $1 - q \exp(r - \epsilon_3 m)$,

$$\|\boldsymbol{b}_k - \boldsymbol{g}_k\| \le 1.2\epsilon_3 \|(\boldsymbol{I} - \boldsymbol{U}\boldsymbol{U}^{\top})\boldsymbol{U}^{\star}\boldsymbol{b}_k^{\star}\|$$
 for all $k \in [q]$

REFERENCES

- [1] S. Babu, S. G. Lingala, and N. Vaswani, "Fast low rank compressive sensing for accelerated dynamic MRI," *IEEE Trans. Comput. Imag.*, vol. 9, pp. 409–424, 2023, doi: 10.1109/TCI.2023.3263810.
- [2] S. S. Du, W. Hu, S. M. Kakade, J. D. Lee, and Q. Lei, "Few-shot learning via learning the representation, provably," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2021.
- [3] S. Nayer and N. Vaswani, "Fast and sample-efficient federated low rank matrix recovery from column-wise linear and quadratic projections," *IEEE Trans. Inf. Theory*, vol. 69, no. 2, pp. 1177–1202, Feb. 2023.
- [4] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Proc. 45th Annu. ACM Symp. Theory Comput.*, Jun. 2013.
- [5] X. Yi, D. Park, Y. Chen, and C. Caramanis, "Fast algorithms for robust PCA via gradient descent," in *Proc. NeurIPS*, 2016.
- [6] S. Lang, Real and Functional Analysis, vol. 10. New York, NY, USA: Springer, 1993, pp. 11–13.
- [7] C. Ma, K. Wang, Y. Chi, and Y. Chen, "Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018.
- [8] S. Nayer and N. Vaswani, "Sample-efficient low rank phase retrieval," in Proc. IEEE Int. Symp. Inf. Theory (ISIT), Jul. 2021, pp. 2244–2249.
- [9] Y. Cherapanamjeri, K. Gupta, and P. Jain, "Nearly-optimal robust matrix completion," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2016.
- [10] R. Vershynin, High-Dimensional Probability: An Introduction With Applications in Data Science, vol. 47. Cambridge, U.K.: Cambridge Univ. Press, 2018.
- [11] P. Wedin, "Perturbation bounds in connection with singular value decomposition," *BIT Numer. Math.*, vol. 12, no. 1, pp. 99–111, 1972.

Namrata Vaswani (Fellow, IEEE) received the B.Tech. degree from IIT Delhi, India, in 1999, and the Ph.D. degree from the University of Maryland, College Park, in 2004.

She is currently the Anderlik Professor in electrical and computer engineering with Iowa State University, where she is also the Director of the CyMath Program, in which graduate students provided school-year-long math tutoring for under-served grade school students. Her research interests are in data science, with a particular focuses on statistical machine learning for signal processing and computational imaging. She was a recipient of the 2014 IEEE Signal Processing Society Best Paper Award, the Iowa State Early Career Engineering Faculty Research Award in 2014, and the Iowa State Mid-Career Achievement in Research Award in 2019. She has served as an Associate Editor or an Area Editor for IEEE TRANSACTIONS ON INFORMATION THEORY, IEEE TRANSACTIONS ON SIGNAL PROCESSING, and the Signal Processing Magazine. She has also guest-edited special issues for PROCEEDINGS OF THE IEEE and the IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING.