

---

# Global Optimality of Local Search for Low Rank Matrix Recovery

---

**Srinadh Bhojanapalli**  
srinadh@ttic.edu

**Behnam Neyshabur**  
bneyshabur@ttic.edu

**Nathan Srebro**  
nati@ttic.edu

Toyota Technological Institute at Chicago

## Abstract

We show that there are no spurious local minima in the non-convex factorized parametrization of low-rank matrix recovery from incoherent linear measurements. With noisy measurements we show all local minima are very close to a global optimum. Together with a curvature bound at saddle points, this yields a polynomial time global convergence guarantee for stochastic gradient descent *from random initialization*.

## 1 Introduction

Low rank matrix recovery problem is heavily studied and has numerous applications in collaborative filtering, quantum state tomography, clustering, community detection, metric learning and multi-task learning [21, 12, 9, 27].

We consider the “matrix sensing” problem of recovering a low-rank (or approximately low rank) p.s.d. matrix<sup>1</sup>  $\mathbf{X}^* \in \mathbb{R}^{n \times n}$ , given a linear measurement operator  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$  and noisy measurements  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{w}$ , where  $\mathbf{w}$  is an i.i.d. noise vector. An estimator for  $\mathbf{X}^*$  is given by the rank-constrained, non-convex problem

$$\underset{\mathbf{X} : \text{rank}(\mathbf{X}) \leq r}{\text{minimize}} \quad \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|^2. \quad (1)$$

This matrix sensing problem has received considerable attention recently [30, 29, 26]. This and other rank-constrained problems are common in machine learning and related fields, and have been used for applications discussed above. A typical theoretical approach to low-rank problems, including (1) is to relax the low-rank constraint to a convex constraint, such as the trace-norm of  $\mathbf{X}$ . Indeed, for matrix sensing, Recht et al. [20] showed that if the measurements are noiseless and the measurement operator  $\mathcal{A}$  satisfies a restricted isometry property, then a low-rank  $\mathbf{X}^*$  can be recovered as the unique solution to a convex relaxation of (1). Subsequent work established similar guarantees also for the noisy and approximate case [14, 6].

However, convex relaxations to the rank are not the common approach employed in practice. In this and other low-rank problems, the method of choice is typically unconstrained local optimization (via e.g. gradient descent, SGD or alternating minimization) on the factorized parametrization

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{U}) = \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top) - \mathbf{y}\|^2, \quad (2)$$

---

<sup>1</sup>We study the case where  $\mathbf{X}^*$  is PSD. We believe the techniques developed here can be used to extend results to the general case.

where the rank constraint is enforced by limiting the dimensionality of  $U$ . Problem (2) is a non-convex optimization problem that could have many bad local minima (as we show in Section 5), as well as saddle points. Nevertheless, local optimization seems to work very well in practice. Working on (2) is much cheaper computationally and allows scaling to large-sized problems—the number of optimization variables is only  $O(nr)$  rather than  $O(n^2)$ , and the updates are usually very cheap, especially compared to typical methods for solving the SDP resulting from the convex relaxation. There is therefore a significant disconnect between the theoretically studied and analyzed methods (based on convex relaxations) and the methods actually used in practice.

Recent attempts at bridging this gap showed that, some form of global “initialization”, typically relying on singular value decomposition, yields a solution that is already close enough to  $\mathbf{X}^*$ ; that local optimization from this initializer gets to the global optima (or to a good enough solution). Jain et al. [15], Keshavan [17] proved convergence for alternating minimization algorithm provided the starting point is close to the optimum, while Zheng and Lafferty [30], Zhao et al. [29], Tu et al. [26], Chen and Wainwright [8], Bhojanapalli et al. [2] considered gradient descent methods on the factor space and proved local convergence. But all these studies rely on global initialization followed by local convergence, and do not tackle the question of the existence of spurious local minima or deal with optimization starting from random initialization. There is therefore still a disconnect between this theory and the empirical practice of starting from random initialization and relying *only* on the local search to find the global optimum.

In this paper we show that, under a suitable incoherence condition on the measurement operator  $\mathcal{A}$  (defined in Section 2), with noiseless measurements and with  $\text{rank}(\mathbf{X}^*) \leq r$ , the problem (2) has no spurious local minima (i.e. all local minima are global and satisfy  $\mathbf{X}^* = \mathbf{U}\mathbf{U}^\top$ ). Furthermore, under the same conditions, all saddle points have a direction with significant negative curvature, and so using a recent result of Ge et al. [10] we can establish that stochastic gradient descent from random initialization converges to  $\mathbf{X}^*$  in polynomial number of iterations. We extend the results also to the noisy and approximately-low-rank settings, where we can guarantee that every local minima is close to a global minimum. The incoherence condition we require is weaker than conditions used to establish recovery through local search, and so our results also ensures recovery in polynomial time under milder conditions than what was previously known. In particular, with i.i.d. Gaussian measurements, we ensure no spurious local minima and recovery through local search with the optimal number  $O(nr)$  of measurements.

**Related Work** Our work is heavily inspired by Bandeira et al. [1], who recently showed similar behavior for the problem of community detection—this corresponds to a specific rank-1 problem with a linear objective, ellipsope constraints and a binary solution. Here we take their ideas, extend them and apply them to matrix sensing with general rank- $r$  matrices. In the past several months, similar type of results were also obtained for other non-convex problems (where the source of non-convexity is *not* a rank constraint), specifically complete dictionary learning [24] and phase recovery [25]. A related recent result of a somewhat different nature pertains to rank unconstrained linear optimization on the ellipsope, showing that local minima of the rank-constrained problem approximate well the global optimum of the rank unconstrained convex problem, even though they might *not* be the global minima (in fact, the approximation guarantee for the actual global optimum is better) [18].

Another non-convex low-rank problem long known to not possess spurious local minima is the PCA problem, which can also be phrased as matrix approximation with full observations, namely  $\min_{\text{rank}(\mathbf{X}) \leq r} \|\mathbf{A} - \mathbf{X}\|_F$  (e.g. [23]). Indeed, local search methods such as the power-method are routinely used for this problem. Recently local optimization methods for the PCA problem working more directly on the optimized formulation have also been studied, including SGD [22] and Grassmannian optimization [28]. These results are somewhat orthogonal to ours, as they study a setting in which it is well known there are never any spurious local minima, and the challenge is obtaining satisfying convergence rates.

The seminal work of Burer and Monteiro [3] proposed low-rank factorized optimization for SDPs, and showed that for extremely high rank  $r > \sqrt{m}$  (number of constraints), an Augmented Lagrangian method converges asymptotically to the optimum. It was also shown that (under mild conditions) any rank deficient local minima is a global minima [4, 16], providing a post-hoc verifiable sufficient condition for global optimality. However, this does not establish any a-priori condition, based on problem structure, implying the lack of spurious local minima.

While preparing this manuscript, we also became aware of parallel work [11] studying the same question for the related but different problem of matrix completion. For this problem they obtain a similar guarantee, though with suboptimal dependence on the incoherence parameters and so suboptimal sample complexity, and requiring adding a specific non-standard regularizer to the objective—this is not needed for our matrix sensing results.

We believe our work, together with the parallel work of [11], are the first to establish the lack of spurious local minima and the global convergence of local search from random initialization for a non-trivial rank-constrained problem (beyond PCA with full observations) with rank  $r > 1$ .

**Notation.** For matrices  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}$ , their inner product is  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{trace}(\mathbf{X}^\top \mathbf{Y})$ . We use  $\|\mathbf{X}\|_F$ ,  $\|\mathbf{X}\|_2$  and  $\|\mathbf{X}\|_*$  for the Frobenius, spectral and nuclear norms of a matrix respectively. Given a matrix  $\mathbf{X}$ , we use  $\sigma_i(\mathbf{X})$  to denote singular values of  $\mathbf{X}$  in decreasing order.  $\mathbf{X}_r = \arg \min_{\text{rank}(\mathbf{Y}) \leq r} \|\mathbf{X} - \mathbf{Y}\|_F$  denotes the rank- $r$  approximation of  $\mathbf{X}$ , as obtained via its truncated singular value decomposition. We use plain capitals  $R$  and  $Q$  to denote orthonormal matrices.

## 2 Formulation and Assumptions

We write the linear measurement operator  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$  as  $\mathcal{A}(\mathbf{X})_i = \langle \mathbf{A}_i, \mathbf{X} \rangle$  where  $\mathbf{A}_i \in \mathbb{R}^{n \times n}$ , yielding  $y_i = \langle \mathbf{A}_i, \mathbf{X}^* \rangle + w_i, i = 1, \dots, m$ . We assume  $w_i \sim \mathcal{N}(0, \sigma_w^2)$  is i.i.d Gaussian noise. We are generally interested in the high dimensional regime where the number of measurements  $m$  is usually much smaller than the dimension  $n^2$ .

Even if we know that  $\text{rank}(\mathbf{X}^*) \leq r$ , having many measurements might not be sufficient for recovery if they are not “spread out” enough. E.g., if all measurements only involve the first  $n/2$  rows and columns, we would never have any information on the bottom-right block. A sufficient condition for identifiability of a low-rank  $\mathbf{X}^*$  from linear measurements by Recht et al. [20] is based on restricted isometry property defined below.

**Definition 2.1** (Restricted Isometry Property). *Measurement operator  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$  (with rows  $\mathbf{A}_i, i = 1, \dots, m$ ) satisfies  $(r, \delta_r)$  RIP if for any  $n \times n$  matrix  $\mathbf{X}$  with rank  $\leq r$ ,*

$$(1 - \delta_r) \|\mathbf{X}\|_F^2 \leq \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle^2 \leq (1 + \delta_r) \|\mathbf{X}\|_F^2. \quad (3)$$

In particular,  $\mathbf{X}^*$  of rank  $r$  is identifiable if  $\delta_{2r} < 1$  [see 20, Theorem 3.2]. One situation in which RIP is obtained is for random measurement operators. For example, matrices with i.i.d.  $\mathcal{N}(0, 1)$  entries satisfy  $(r, \delta_r)$ -RIP when  $m = O(\frac{nr}{\delta_r^2})$  [see 6, Theorem 2.3]. This implies identifiability based on i.i.d. Gaussian measurement with  $m = O(nr)$  measurements (coincidentally, the number of degrees of freedom in  $\mathbf{X}^*$ , optimal up to a constant factor).

## 3 Main Results

We are now ready to present our main result about local minima for the matrix sensing problem (2). We first present the results for noisy sensing of exact low rank matrices, and then generalize the results also to approximately low rank matrices.

Now we will present our result characterizing local minima of  $f(\mathbf{U})$ , for low-rank  $\mathbf{X}^*$ . Recall that measurements are  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{w}$ , where entries of  $\mathbf{w}$  are i.i.d. Gaussian -  $w_i \sim \mathcal{N}(0, \sigma_w^2)$ .

**Theorem 3.1.** *Consider the optimization problem (2) where  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{w}$ ,  $\mathbf{w}$  is i.i.d.  $\mathcal{N}(0, \sigma_w^2)$ ,  $\mathcal{A}$  satisfies  $(4r, \delta_{4r})$ -RIP with  $\delta_{4r} < \frac{1}{10}$ , and  $\text{rank}(\mathbf{X}^*) \leq r$ . Then, with probability  $\geq 1 - \frac{10}{n^2}$  (over the noise), for any local minimum  $\mathbf{U}$  of  $f(\mathbf{U})$ :*

$$\|\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*\|_F \leq 20 \sqrt{\frac{\log(n)}{m}} \sigma_w.$$

In particular, in the noiseless case ( $\sigma_w = 0$ ) we have  $\mathbf{U}\mathbf{U}^\top = \mathbf{X}^*$  and so  $f(\mathbf{U}) = 0$  and every local minima is global. In the noiseless case, we can also relax the RIP requirement to  $\delta_{4r} < 1/5$  (see Theorem 4.1 in Section 4). In the noisy case we cannot expect to ensure we always get to an exact global minima, since the noise might cause tiny fluctuations very close to the global minima possibly

creating multiple very close local minima. But we show that all local minima are indeed very close to some factorization  $\mathbf{U}^*\mathbf{U}^{*\top} = \mathbf{X}^*$  of the true signal, and hence to a global optimum, and this “radius” of local minima decreases as we have more observations.

The proof of the Theorem for the noiseless case is presented in Section 4. The proof for the general setting follows along the same lines and can be found in the Appendix.

So far we have discussed how all local minima are global, or at least very close to a global minimum. Using a recent result by Ge et al. [10] on the convergence of SGD for non-convex functions, we can further obtain a polynomial bound on the number of SGD iterations required to reach the global minima. The main condition that needs to be established in order to ensure this, is that all saddle points of (2) satisfy the “strict saddle point condition”, i.e. have a direction with significant negative curvature:

**Theorem 3.2** (Strict saddle). *Consider the optimization problem (2) in the noiseless case, where  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$ ,  $\mathcal{A}$  satisfies  $(4r, \delta_{4r})$ -RIP with  $\delta_{4r} < \frac{1}{10}$ , and  $\text{rank}(\mathbf{X}^*) \leq r$ . Let  $\mathbf{U}$  be a first order critical point of  $f(\mathbf{U})$  with  $\mathbf{U}\mathbf{U}^\top \neq \mathbf{X}^*$ . Then the smallest eigenvalue of the Hessian satisfies*

$$\lambda_{\min} \left[ \frac{1}{m} \nabla^2(f(\mathbf{U})) \right] \leq \frac{-2}{5} \sigma_r(\mathbf{X}^*).$$

Now consider the stochastic gradient descent updates,

$$\mathbf{U}^+ = \text{Proj}_b \left( \mathbf{U} - \eta \left( \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top \rangle - y_i) \mathbf{A}_i \mathbf{U} + \psi \right) \right), \quad (4)$$

where  $\psi$  is uniformly distributed on the unit sphere and  $\text{Proj}_b$  is a projection onto  $\|\mathbf{U}\|_F \leq b$ . Using Theorem 3.2 and the result of Ge et al. [10] we can establish:

**Theorem 3.3** (Convergence from random initialization). *Consider the optimization problem (2) under the same noiseless conditions as in Theorem 3.2. Using  $b \geq \|\mathbf{U}^*\|_F$ , for some global optimum  $\mathbf{U}^*$  of  $f(\mathbf{U})$ , for any  $\epsilon, c > 0$ , after  $T = \text{poly} \left( \frac{1}{\sigma_r(\mathbf{X}^*)}, \sigma_1(\mathbf{X}^*), b, \frac{1}{\epsilon}, \log(1/c) \right)$  iterations of (4) with an appropriate stepsize  $\eta$ , starting from a random point uniformly distributed on  $\|\mathbf{U}\|_F = b$ , with probability at least  $1 - c$ , we reach an iterate  $\mathbf{U}_T$  satisfying*

$$\|\mathbf{U}_T - \mathbf{U}^*\|_F \leq \epsilon.$$

The above result guarantees convergence of noisy gradient descent to a global optimum. Alternatively, second order methods such as cubic regularization (Nesterov and Polyak [19]) and trust region (Cartis et al. [7]) that have guarantees based on the strict saddle point property can also be used here.

**RIP Requirement:** Our results require  $(4r, 1/10)$ -RIP for the noisy case and  $(4r, 1/5)$ -RIP for the noiseless case. Requiring  $(2r, \delta_{2r})$ -RIP with  $\delta_{2r} < 1$  is sufficient to ensure uniqueness of the global optimum of (1), and thus recovery in the noiseless setting [20], but all known efficient recovery methods require stricter conditions. The best guarantees we are aware of require  $(5r, 1/10)$ -RIP [20] or  $(4r, 0.414)$ -RIP [6] using a convex relaxation. Alternatively,  $(6r, 1/10)$ -RIP is required for global initialization followed by non-convex optimization [26]. In terms of requirements on  $(2r, \delta_{2r})$ -RIP for non-convex methods, the best we are aware of is requiring  $\delta_{2r} < \Omega(1/r)$  [15, 29, 30]—this is a much stronger condition than ours, and it yields a suboptimal required number of spherical Gaussian measurements of  $\Omega(nr^3)$ . So, compared to prior work our requirement is very mild—it ensures efficient recovery, and requires the optimal number of spherical Gaussian measurements (up to a constant factor) of  $O(nr)$ .

**Extension to Approximate Low Rank** We can also obtain similar results that deteriorate gracefully if  $\mathbf{X}^*$  is not exactly low rank, but is close to being low-rank (see proof in the Appendix):

**Theorem 3.4.** *Consider the optimization problem (2) where  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$  and  $\mathcal{A}$  satisfies  $(4r, \delta_{4r})$ -RIP with  $\delta_{4r} < \frac{1}{100}$ . Then, for any local minima  $\mathbf{U}$  of  $f(\mathbf{U})$ :*

$$\|\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*\|_F \leq 4(\|\mathbf{X}^* - \mathbf{X}_r^*\|_F + \delta_{2r} \|\mathbf{X}^* - \mathbf{X}_r^*\|_*),$$

where  $\mathbf{X}_r^*$  is the best rank  $r$  approximation of  $\mathbf{X}^*$ .

This theorem guarantees that any local optimum of  $f(U)$  is close to  $\mathbf{X}^*$  upto an error depending on  $\|\mathbf{X}^* - \mathbf{X}_r^*\|$ . For the low-rank noiseless case we have  $\mathbf{X}^* = \mathbf{X}_r^*$  and the right hand side vanishes. When  $\mathbf{X}^*$  is not exactly low rank, the best recovery error we can hope for is  $\|\mathbf{X}^* - \mathbf{X}_r^*\|_F$ , since  $UU^\top$  is at most rank  $k$ . On the right hand side of Theorem 3.4, we have also a nuclear norm term, which might be higher, but it also gets scaled down by  $\delta_{2r}$ , and so by the number of measurements.

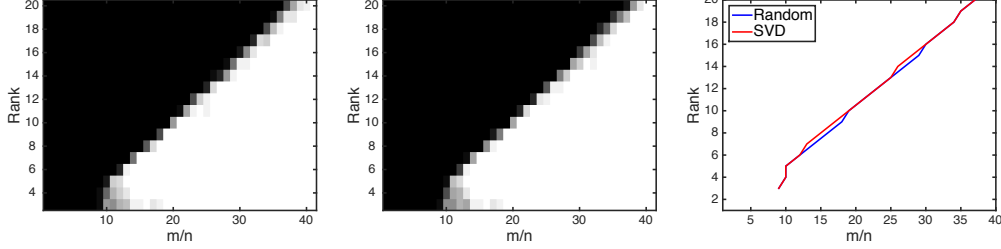


Figure 1: The plots in this figure compare the success probability of gradient descent between (left) random and (center) SVD initialization (suggested in [15]), for problem (2), with increasing number of samples  $m$  and various values of rank  $r$ . Right most plot is the first  $m$  for a given  $r$ , where the probability of success reaches the value 0.5. A run is considered success if  $\|UU^\top - \mathbf{X}^*\|_F / \|\mathbf{X}^*\|_F \leq 1e-2$ . White cells denote success and black cells denote failure of recovery. We set  $n$  to be 100. Measurements  $y_i$  are inner product of entrywise i.i.d Gaussian matrix and a rank- $r$  p.s.d matrix with random subspace. We notice no significant difference between the two initialization methods, suggesting absence of local minima as shown. Both methods have phase transition around  $m = 2 \cdot n \cdot r$ .

#### 4 Proof for the Noiseless Case

In this section we present the proof characterizing the local minima of problem (2). For ease of exposition we first present the results for the noiseless case ( $\mathbf{w} = 0$ ). Proof for the general case can be found in the Appendix.

**Theorem 4.1.** *Consider the optimization problem (2) where  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$ ,  $\mathcal{A}$  satisfies  $(4r, \delta_{4r})$ -RIP with  $\delta_{4r} < \frac{1}{5}$ , and  $\text{rank}(\mathbf{X}^*) \leq r$ . Then, for any local minimum  $\mathbf{U}$  of  $f(\mathbf{U})$ :*

$$UU^\top = \mathbf{X}^*.$$

For the proof of this theorem we first discuss the implications of the first and second order optimality conditions and then show how to combine them to yield the result.

Invariance of  $f(\mathbf{U})$  over  $r \times r$  orthonormal matrices introduces additional challenges in comparing a given stationary point to a global optimum. We have to find the best orthonormal matrix  $R$  to align a given stationary point  $\mathbf{U}$  to a global optimum  $\mathbf{U}^*$ , where  $\mathbf{U}^* \mathbf{U}^{*\top} = \mathbf{X}^*$ , to combine results from the first and second order conditions, without degrading the isometry constants.

Consider a local optimum  $\mathbf{U}$  that satisfies first and second order optimality conditions of problem (2). In particular  $\mathbf{U}$  satisfies  $\nabla f(\mathbf{U}) = 0$  and  $z^\top \nabla^2 f(\mathbf{U}) z \geq 0$  for any  $z \in \mathbb{R}^{n \cdot r}$ . Now we will see how these two conditions constrain the error  $UU^\top - \mathbf{U}^* \mathbf{U}^{*\top}$ .

First we present the following consequence of the RIP assumption [see 5, Lemma 2.1].

**Lemma 4.1.** *Given two  $n \times n$  rank- $r$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ , and a  $(4r, \delta)$ -RIP measurement operator  $\mathcal{A}$ , the following holds:*

$$\left| \frac{1}{m} \sum_{i=1}^m \langle \mathcal{A}_i, \mathbf{X} \rangle \langle \mathcal{A}_i, \mathbf{Y} \rangle - \langle \mathbf{X}, \mathbf{Y} \rangle \right| \leq \delta \|\mathbf{X}\|_F \|\mathbf{Y}\|_F. \quad (5)$$

##### 4.1 First order optimality

First we will consider the first order condition,  $\nabla f(\mathbf{U}) = 0$ . For any stationary point  $\mathbf{U}$  this implies

$$\sum_i \langle \mathcal{A}_i, UU^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \mathcal{A}_i \mathbf{U} = 0. \quad (6)$$

Now using the isometry property of  $\mathbf{A}_i$  gives us the following result.

**Lemma 4.2.** [First order condition] For any first order stationary point  $\mathbf{U}$  of  $f(\mathbf{U})$ , and  $\mathcal{A}$  satisfying the  $(4r, \delta)$ -RIP (3), the following holds:

$$\|(\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top})\mathbf{Q}\mathbf{Q}^\top\|_F \leq \delta \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F,$$

where  $\mathbf{Q}$  is an orthonormal matrix that spans the column space of  $\mathbf{U}$ .

This lemma states that any stationary point of  $f(\mathbf{U})$  is close to a global optimum  $\mathbf{U}^*$  in the subspace spanned by columns of  $\mathbf{U}$ . Notice that the error along the orthogonal subspace  $\mathbf{Q}_\perp$ ,  $\|\mathbf{X}^*\mathbf{Q}_\perp\mathbf{Q}_\perp^\top\|_F$  can still be large making the distance between  $\mathbf{X}$  and  $\mathbf{X}^*$  arbitrarily far.

*Proof of Lemma 4.2.* Let  $\mathbf{U} = \mathbf{Q}\mathbf{R}$ , for some orthonormal  $\mathbf{Q}$ . Consider any matrix of the form  $\mathbf{Z}\mathbf{Q}\mathbf{R}^\dagger{}^2$ . The first order optimality condition then implies,

$$\sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{U}\mathbf{R}^\dagger\mathbf{Q}^\top\mathbf{Z}^\top \rangle = 0$$

The above equation together with Restricted Isometry Property (equation (5)) gives us the following inequality:

$$\left| \langle \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}, \mathbf{Q}\mathbf{Q}^\top\mathbf{Z}^\top \rangle \right| \leq \delta \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F \|\mathbf{Q}\mathbf{Q}^\top\mathbf{Z}^\top\|_F.$$

Note that for any matrix  $\mathbf{A}$ ,  $\langle \mathbf{A}, \mathbf{Q}\mathbf{Q}^\top\mathbf{Z} \rangle = \langle \mathbf{Q}\mathbf{Q}^\top\mathbf{A}, \mathbf{Z} \rangle$ . Furthermore, for any matrix  $\mathbf{A}$ ,  $\sup_{\{\mathbf{Z}: \|\mathbf{Z}\|_F \leq 1\}} \langle \mathbf{A}, \mathbf{Z} \rangle = \|\mathbf{A}\|_F$ . Hence the above inequality implies the lemma statement.  $\square$

## 4.2 Second order optimality

We now consider the second order condition to show that the error along  $\mathbf{Q}_\perp\mathbf{Q}_\perp^\top$  is indeed bounded well. Let  $\nabla^2 f(\mathbf{U})$  be the hessian of the objective function. Note that this is an  $n \cdot r \times n \cdot r$  matrix. Fortunately for our result we need to only evaluate the Hessian along  $\text{vec}(\mathbf{U} - \mathbf{U}^*\mathbf{R})$  for some orthonormal matrix  $\mathbf{R}$ . Here  $\text{vec}(\cdot)$  denotes writing a matrix in vector form.

**Lemma 4.3.** [Hessian computation] Let  $\mathbf{U}$  be a first order critical point of  $f(\mathbf{U})$ . Then for any  $r \times r$  orthonormal matrix  $\mathbf{R}$  and  $\Delta_j = \Delta e_j e_j^\top$  ( $\Delta = \mathbf{U} - \mathbf{U}^*\mathbf{R}$ ),

$$\sum_{j=1}^r \text{vec}(\Delta_j)^\top [\nabla^2 f(\mathbf{U})] \text{vec}(\Delta_j) = \sum_{i=1}^m \left( \sum_{j=1}^r 4 \langle \mathbf{A}_i, \mathbf{U}\Delta_j^\top \rangle^2 - 2 \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle^2 \right),$$

Hence from second order optimality of  $\mathbf{U}$  we get,

**Corollary 4.1.** [Second order optimality] Let  $\mathbf{U}$  be a local minimum of  $f(\mathbf{U})$ . For any  $r \times r$  orthonormal matrix  $\mathbf{R}$ ,

$$\sum_{j=1}^r \sum_{i=1}^m 4 \langle \mathbf{A}_i, \mathbf{U}\Delta_j^\top \rangle^2 \geq \frac{1}{2} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle^2, \quad (7)$$

Further for  $\mathcal{A}$  satisfying  $(2r, \delta)$ -RIP (equation (3)) we have,

$$\sum_{j=1}^r \|\mathbf{U}e_j e_j^\top (\mathbf{U} - \mathbf{U}^*\mathbf{R})^\top\|_F^2 \geq \frac{1 - \delta}{2(1 + \delta)} \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F^2. \quad (8)$$

The proof of this result follows simply by applying Lemma 4.3. The above Lemma gives a bound on the distance in the factor ( $\mathbf{U}$ ) space  $\|\mathbf{U}(\mathbf{U} - \mathbf{U}^*\mathbf{R})^\top\|_F^2$ . To be able to compare the second order condition to the first order condition we need a relation between  $\|\mathbf{U}(\mathbf{U} - \mathbf{U}^*\mathbf{R})^\top\|_F^2$  and  $\|\mathbf{X} - \mathbf{X}^*\|_F^2$ . Towards this we show the following result.

---

<sup>2</sup> $\mathbf{R}^\dagger$  is the pseudo inverse of  $\mathbf{R}$

**Lemma 4.4.** Let  $\mathbf{U}$  and  $\mathbf{U}^*$  be two  $n \times r$  matrices, and  $\mathbf{Q}$  is an orthonormal matrix that spans the column space of  $\mathbf{U}$ . Then there exists an  $r \times r$  orthonormal matrix  $\mathbf{R}$  such that for any first order stationary point  $\mathbf{U}$  of  $f(\mathbf{U})$ , the following holds:

$$\sum_{j=1}^r \|\mathbf{U} e_j e_j^\top (\mathbf{U} - \mathbf{U}^* \mathbf{R})^\top\|_F^2 \leq \frac{1}{8} \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 + \frac{34}{8} \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{Q} \mathbf{Q}^\top\|_F^2.$$

This Lemma bounds the distance in the factor space ( $\|(\mathbf{U} - \mathbf{U}^* \mathbf{R}) \mathbf{U}^\top\|_F^2$ ) with  $\|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2$  and  $\|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{Q} \mathbf{Q}^\top\|_F^2$ . Combining this with the result from second order optimality (Corollary 4.1) shows  $\|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2$  is bounded by a constant factor of  $\|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{Q} \mathbf{Q}^\top\|_F^2$ . This implies  $\|\mathbf{X}^* \mathbf{Q}_\perp \mathbf{Q}_\perp^\top\|_F$  is bounded, opposite to what the first order condition implied (Lemma 4.2). The proof of the above lemma is in Section B. Hence from the above optimality conditions we get the proof of Theorem 4.1.

*Proof of Theorem 4.1.* Assuming  $\mathbf{U} \mathbf{U}^\top \neq \mathbf{U}^* \mathbf{U}^{*\top}$ , from Lemmas 4.2, 4.4 and Corollary 4.1 we get,

$$\left( \frac{1 - \delta}{2(1 + \delta)} - \frac{1}{8} \right) \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 \leq \frac{34}{8} \delta^2 \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top})\|_F^2.$$

If  $\delta \leq \frac{1}{5}$  the above inequality holds only if  $\mathbf{U} \mathbf{U}^\top = \mathbf{U}^* \mathbf{U}^{*\top}$ .  $\square$

## 5 Necessity of RIP

We showed that there are no spurious local minima only under a restricted isometry assumption. A natural question is whether this is necessary, or whether perhaps the problem (2) never has any spurious local minima, perhaps similarly to the non-convex PCA problem  $\min_{\mathbf{U}} \|\mathbf{A} - \mathbf{U} \mathbf{U}^\top\|$ .

A good indication that this is not the case is that (2) is NP-hard, even in the noiseless case when  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$  for  $\text{rank}(\mathbf{X}^*) \leq k$  [20] (if we don't require RIP, we can have each  $\mathbf{A}_i$  be non-zero on a single entry in which case (2) becomes a matrix completion problem, for which hardness has been shown even under fairly favorable conditions [13])<sup>3</sup>. That is, we are unlikely to have a poly-time algorithm that succeeds for any linear measurement operator. Although this doesn't formally preclude the possibility that there are no spurious local minima, but it just takes a very long time to find a local minima, this scenario seems somewhat unlikely.

To resolve the question, we present an explicit example of a measurement operator  $\mathcal{A}$  and  $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$  (i.e.  $f(\mathbf{X}^*) = 0$ ), with  $\text{rank}(\mathbf{X}^*) = r$ , for which (1), and so also (2), have a non-global local minima.

*Example 1:* Let  $f(\mathbf{X}) = (X_{11} + X_{22} - 1)^2 + (X_{11} - 1)^2 + X_{12}^2 + X_{21}^2$  and consider (1) with  $r = 1$  (i.e. a rank-1 constraint). For  $\mathbf{X}^* = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$  we have  $f(\mathbf{X}^*) = 0$  and  $\text{rank}(\mathbf{X}^*) = 1$ . But  $\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$  is a rank 1 local minimum with  $f(\mathbf{X}) = 1$ .

We can be extended the construction to any rank  $r$  by simply adding  $\sum_{i=3}^{r+2} (X_{ii} - 1)^2$  to the objective, and padding both the global and local minimum with a diagonal beneath the leading  $2 \times 2$  block.

In Example 1, we had a rank- $r$  problem, with a rank- $r$  exact solution, and a rank- $r$  local minima. Another question we can ask is what happens if we allow a larger rank than the rank of the optimal solution. That is, if we have  $f(\mathbf{X}^*) = 0$  with low  $\text{rank}(\mathbf{X}^*)$ , even  $\text{rank}(\mathbf{X}^*) = 1$ , but consider (1) or (2) with a high  $r$ . Could we still have non-global local minima? The answer is yes...

*Example 2:* Let  $f(\mathbf{X}) = (X_{11} + X_{22} + X_{33} - 1)^2 + (X_{11} - 1)^2 + (X_{22} - X_{33})^2 + \sum_{i,j:i \neq j} X_{ij}^2$  and consider the problem (1) with a rank  $r = 2$  constraint. We can verify that  $\mathbf{X}^* = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

is a rank=1 global minimum with  $f(\mathbf{X}^*) = 0$ , but  $\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1/2 & 0 \\ 0 & 0 & 1/2 \end{bmatrix}$  is a local minimum with

<sup>3</sup>Note that matrix completion is tractable under incoherence assumptions, similar to RIP

$f(\mathbf{X}) = 1$ . Also for an arbitrary large rank constraint  $r > 1$  (taking  $r$  to be odd for simplicity), extend the objective to  $f(\mathbf{X}) = (X_{11} - 1)^2 + \sum_{i=1}^{(r-1)/2} [(X_{11} + X_{2i,2i} + X_{(2i+1),(2i+1)} - 1)^2 + (X_{2i,2i} - X_{(2i+1),(2i+1)})^2]$ . We still have a rank-1 global minimum  $\mathbf{X}^*$  with a single non-zero entry  $\mathbf{X}_{11}^* = 1$ , while  $\mathbf{X} = (I - \mathbf{X}^*)/2$  is a local minimum with  $f(\mathbf{X}) = 1$ .

## 6 Conclusion

We established that under conditions similar to those required for convex relaxation recovery guarantees, the non-convex formulation of matrix sensing (2) does not exhibit any spurious local minima (or, in the noisy and approximate settings, at least not outside some small radius around a global minima), and we can obtain theoretical guarantees on the success of optimizing it using SGD from *random initialization*. This matches the methods frequently used in practice, and can explain their success. This guarantee is very different in nature from other recent work on non-convex optimization for low-rank problems, which relied heavily on initialization to get close to the global optimum, and on local search just for the final local convergence to the global optimum. We believe this is the first result, together with the parallel work of Ge et al. [11], on the global convergence of local search for common rank-constrained problems that are worst-case hard.

Our result suggests that SVD initialization is not necessary for global convergence, and random initialization would succeed under similar conditions (in fact, our conditions are even weaker than in previous work that used SVD initialization). To investigate empirically whether SVD initialization is indeed helpful for ensuring global convergence, in Figure 1 we compare recovery probability of random rank- $k$  matrices for random and SVD initialization—there is no significant difference between the two.

Beyond the implications for matrix sensing, we are hoping these type of results could be a first step and serve as a model for understanding local search in deep networks. Matrix factorization, such as in (2), is a depth-two neural network with linear transfer—an extremely simple network, but already non-convex and arguably the most complicated network we have a good theoretical understanding of. Deep networks are also hard to optimize in the worst case, but local search seems to do very well in practice. Our ultimate goal is to use the study of matrix recovery as a guide in understating the conditions that enable efficient training of deep networks.

## Acknowledgements

Authors would like to thank Afonso Bandeira for discussions, Jason Lee and Tengyu Ma for sharing and discussing their work. This research was supported in part by an NSF RI/AF grant 1302662.

## References

- [1] A. S. Bandeira, N. Boumal, and V. Voroninski. On the low-rank approach for semidefinite programs arising in synchronization and community detection. *arXiv preprint arXiv:1602.04426*, 2016.
- [2] S. Bhojanapalli, A. Kyriillidis, and S. Sanghavi. Dropping convexity for faster semi-definite optimization. *arXiv preprint arXiv:1509.03917*, 2015.
- [3] S. Burer and R. D. Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, 2003.
- [4] S. Burer and R. D. Monteiro. Local minima and convergence in low-rank semidefinite programming. *Mathematical Programming*, 103(3):427–444, 2005.
- [5] E. J. Candès. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9):589–592, 2008.
- [6] E. J. Candès and Y. Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *Information Theory, IEEE Transactions on*, 57(4):2342–2359, 2011.
- [7] C. Cartis, N. I. Gould, and P. L. Toint. Complexity bounds for second-order optimality in unconstrained optimization. *Journal of Complexity*, 28(1):93–108, 2012.
- [8] Y. Chen and M. J. Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.



- [9] S. Flammia, D. Gross, Y.-K. Liu, and J. Eisert. Quantum tomography via compressed sensing: Error bounds, sample complexity and efficient estimators. *New Journal of Physics*, 14(9):095022, 2012.
- [10] R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Proceedings of The 28th Conference on Learning Theory*, pages 797–842, 2015.
- [11] R. Ge, J. Lee, and T. Ma. Matrix completion has no spurious local minimum. *arXiv preprint arXiv:1605.07272*, 2016.
- [12] D. Gross, Y.-K. Liu, S. T. Flammia, S. Becker, and J. Eisert. Quantum state tomography via compressed sensing. *Physical review letters*, 105(15):150401, 2010.
- [13] M. Hardt, R. Meka, P. Raghavendra, and B. Weitz. Computational limits for matrix completion. In *Proceedings of The 27th Conference on Learning Theory*, pages 703–725, 2014.
- [14] P. Jain, R. Meka, and I. S. Dhillon. Guaranteed rank minimization via singular value projection. In *Advances in Neural Information Processing Systems*, pages 937–945, 2010.
- [15] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the 45th annual ACM Symposium on theory of computing*, pages 665–674. ACM, 2013.
- [16] M. Journée, F. Bach, P.-A. Absil, and R. Sepulchre. Low-rank optimization on the cone of positive semidefinite matrices. *SIAM Journal on Optimization*, 20(5):2327–2351, 2010.
- [17] R. H. Keshavan. *Efficient algorithms for collaborative filtering*. PhD thesis, STANFORD, 2012.
- [18] A. Montanari. A Grothendieck-type inequality for local maxima. *arXiv preprint arXiv:1603.04064*, 2016.
- [19] Y. Nesterov and B. T. Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [20] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010.
- [21] J. D. Rennie and N. Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *Proceedings of the 22nd international conference on Machine learning*, pages 713–719. ACM, 2005.
- [22] C. D. Sa, C. Re, and K. Olukotun. Global convergence of stochastic gradient descent for some non-convex matrix problems. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2332–2341, 2015.
- [23] N. Srebro and T. Jaakkola. Weighted low-rank approximations. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 720–727, 2003.
- [24] J. Sun, Q. Qu, and J. Wright. Complete dictionary recovery using nonconvex optimization. In *Proceedings of The 32nd International Conference on Machine Learning*, pages 2351–2360, 2015.
- [25] J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *preprint arXiv:1602.06664*, 2016.
- [26] S. Tu, R. Boczar, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via Procrustes flow. *arXiv preprint arXiv:1507.03566*, 2015.
- [27] H.-F. Yu, P. Jain, P. Kar, and I. Dhillon. Large-scale multi-label learning with missing labels. In *Proceedings of The 31st International Conference on Machine Learning*, pages 593–601, 2014.
- [28] D. Zhang and L. Balzano. Global convergence of a grassmannian gradient descent algorithm for subspace estimation. *arXiv preprint arXiv:1506.07405*, 2015.
- [29] T. Zhao, Z. Wang, and H. Liu. A nonconvex optimization framework for low rank matrix estimation. In *Advances in Neural Information Processing Systems*, pages 559–567, 2015.
- [30] Q. Zheng and J. 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.