

# Literature Review

## Nonconvex Optimization Methods & Applications: Low-Rank Matrix Factorization [1] [2]

### Abstract

While nonconvex optimization methods are not ideal due to local minima and saddle points, various iterative methods work in practice under specific conditions. This is possible given conditions like a benign landscape (as an example) that can provide theoretical convergence guarantees in a provable manner. Moreover, two different methodologies are discussed to solve nonconvex problems under matrix factorization applications.

### 1 Overview

1. Overview & Motivation
2. Canonical Matrix Factorization Problems
3. A Short Discussion on Nonconvexity & Convexity
4. Two-Stage Algorithms
  - Gradient Descent Method
  - Alternating Minimization Method
  - Iterative Hard Thresholding Method
  - Initialization w/ Spectral Method
5. Global Landscape Algorithms
  - Gradient-based Methods
  - Hessian-based Methods
  - Almost-Convex-AGD & Negative-Curvature-Descent Method
6. Conclusion & References
7. Appendices (Various Proofs)

Convex optimization methods, such as gradient descent methods, have been extremely popular to solve optimization problems due to theoretical guarantees of global convergence, fast convergence rate, and ease of use. However, some research problems are inherently nonconvex meaning multiple local minima exist, and thus algorithms may not necessarily converge to the best global solution. Despite this fact, nonconvex algorithms exist that are efficient in a provable manner under certain conditions. This report will include discussion on various nonconvex methods applied to matrix factorization applications. Some examples of matrix factorization problems include problems such as matrix sensing, phase retrieval, matrix completion, blind deconvolution, robust principal component analysis, phase synchronization, and joint alignment. The nonconvex methods are efficient and achieve theoretical convergence guarantees due to the use of statistical models and specific conditions that are met like restricted isometry property. This report will mainly discuss two different approaches to the nonconvex problem: 1) two-part algorithms with a particular initialization step with iterative refinement of the initial estimate; and 2) initialization-free algorithms and optimization landscape analysis. Before doing so, various canonical matrix factorization problems are introduced to further motivate discussion of nonconvex methods.

## 2 Canonical Problems

### 2.1 Matrix Sensing [3]

Assume  $m$  observations are made and denoted as  $y_i \in \mathbb{F}, \forall i \in 1, \dots, m$ . The observation model is described as  $y_i = \langle \mathbf{A}_i, \mathbf{M}_* \rangle = \text{trace}(\mathbf{A}_i \mathbf{M}_*^T)$  where we wish to reconstruct the unprocessed measurement matrix  $\mathbf{M}_*$  (assumed rank  $r$ ) given the set of sensing matrices known beforehand denoted  $\mathbf{A}_i$ . A simple example can be considered where  $\mathbf{M}_*$  is a rank-1 matrix (represented as  $\mathbf{M}_* = \mathbf{X}_* \mathbf{X}_*^T \in \mathbb{R}^{n \times n}$ ). In this case, the minimization problem reduces to the classical least-squares setup

$$\min_{\mathbf{X}} \frac{1}{4m} \sum_i^m (\langle \mathbf{A}_i, \mathbf{X} \mathbf{X}^T \rangle - y_i)^2 \quad (2.1.1)$$

One important observation to make is that the trace operator is unitary invariant and thus multiple solutions exist for the minimization problem since the solution can be multiplied by a unitary matrix. Meaning, this is a nonconvex optimization problem (albeit easily solvable and well understood). Generally, for the rank- $r$  case, we describe  $\mathbf{M}_* = \mathbf{L}_* \mathbf{R}_* = \mathbf{U}_* \boldsymbol{\Sigma}_*^{1/2} \boldsymbol{\Sigma}_*^{1/2} \mathbf{V}_*^T$  where  $\mathbf{U}_*, \mathbf{V}_*, \boldsymbol{\Sigma}_*$  represent the singular value decomposition of  $\mathbf{M}_*$  to solve:

$$\min_{\mathbf{L}, \mathbf{R}} \frac{1}{4m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{L} \mathbf{R}^T \rangle - y_i)^2 \quad (2.1.2)$$

Some examples of practical matrix sensing problems include facial recognition and quantum state tomography. This matrix sensing problem will be used extensively in the coming chapters when discussing theoretical guarantees of different methods. The rest of this chapter is included to show variety of matrix factorization problems but may not necessarily show up in subsequent chapters.

### 2.2 Robust Principle Component Analysis (RPCA) [4]

Classical PCA boils down to taking a singular value decomposition of a matrix in order to create a low-rank approximation of the original matrix. However, under non-ideal conditions like outlier data, the SVD of a matrix will incorporate the outlier data as part of the low-rank construction. To avoid this, we wish to decompose data matrix  $\mathbf{\Gamma}_* = \mathbf{M}_* + \mathbf{S}_*$  where  $\mathbf{M}_*$  is the low-rank matrix to recover but is essentially corrupted by  $\mathbf{S}_*$  the sparse outliers. For the general  $r$ -rank case ( $\mathbf{M}_* = \mathbf{L} \mathbf{R}^T$ ), the minimization problem is the following assuming all entries of  $\mathbf{\Gamma}_*$  are observed (no missing entries):

$$\min_{\mathbf{L}, \mathbf{R}, \mathbf{S}} \frac{1}{4p} \|\mathbf{\Gamma}_* - \mathbf{L} \mathbf{R}^T - \mathbf{S}\|_F^2 \text{ where } \mathbf{S} \in \mathcal{S}_\alpha = \{\mathbf{S} : \|(\mathbf{S}_*)_{i,\cdot}\|_0 \leq \alpha n_2, \|(\mathbf{S}_*)_{\cdot,j}\|_0 \leq \alpha n_1 \forall i, j\}$$

This is in contrast to the convex relaxation version which is written as follows:

$$\min_{\mathbf{M}, \mathbf{S}} \|\mathbf{M}\|_* + \lambda \|\mathbf{S}\|_1 \text{ such that } \mathbf{M} + \mathbf{S} = \mathbf{\Gamma}$$

In this case, sparsity is promoted for the outlier matrix since there is an implicit assumption that outlier data is relatively uncommon and represents a small portion of the overall data matrix. Moreover, the  $\mathbf{L} \mathbf{R}$  term is now replaced by the low-rank matrix  $\mathbf{M}$  that has been relaxed by the nuclear norm constraint. These kinds of RPCA problems appear in computer vision, medical imaging, and surveillance to provide a few examples.

### 2.3 Phase Retrieval [5]

For the rank-1 case, suppose we have  $m$  measurements  $y_i, i \in \{1, \dots, m\}$  and  $M_* = X_* X_*^T$ . The problem is now to reconstruct  $X_*$  given  $y_i$ . These parameters are related via  $y_i = \|\mathbf{a}_i^T X_*\|_2^2 = \mathbf{a}_i M_* \mathbf{a}_i$  where  $\mathbf{a}_i$  is the design vector that is known beforehand. Given this setup, the optimization problem is described as:

$$\min_{\mathbf{X}} \frac{1}{4m} \sum_{i=1}^m (\|\mathbf{a}_i^T \mathbf{X}\|_2^2 - y_i)^2 \quad (2.3.1)$$

An astute observation is the similarity with matrix sensing in that  $\mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^T$  so phase retrieval is a particular setup of the general matrix sensing problem. Phase retrieval shows up in applications such as X-ray crystallography. The main idea to this problem is to figure out phase information given only magnitude information since this is what the transducers are able to detect.

## 3 A Short Discussion on Nonconvexity & Convexity

Before discussing the nonconvex methods, a short discussion will be provided to set context on the two chapters to follow. Usually, within a matrix factorization problem, it is desirable to solve  $\min_{\mathbf{L}, \mathbf{R}} f(\mathbf{L}, \mathbf{R})$  subject to  $\mathbf{M} = \mathbf{L}\mathbf{R}^T$ . More often than not, the problem is inherently nonconvex and may seem impossible to solve for the global solution. Classically, the approach taken is to use *convex relaxation* to replace  $\mathbf{M} = \mathbf{L}\mathbf{R}^T$  by a nuclear norm restraint as an example and instead solve for  $\mathbf{M}$  instead of  $\mathbf{L}\mathbf{R}$ . These kinds of methods have been incredible successful in practice and have performance guarantees when talking about things like sample complexity and robustness (against noise). However, the computational cost is quite large  $\mathcal{O}(n^3)$  depending on the size of the matrix which may even far exceed the time it takes to read the data. Moreover, extremely large matrices will require large storage allocation both in storage and memory. In this paper, low-rank matrix estimation will be based on nonconvex methods  $(\mathbf{L}, \mathbf{R})$  instead of convex approaches. It may not seem intuitive, but it turns out that the nonconvex approach has low memory/storage requirements, per-iteration cost significantly improves (e.g. no nuclear norm computation), malleability towards parallelization, and scalability to large-scale problems [1]. A fair question one may ask is whether it is even possible to solve a nonconvex problem. On the surface it would appear to be a daunting task. Remarkably, many important nonconvex problems are not difficult to solve. In the case of matrix factorization problems, straightforward first-order methods are guaranteed to succeed with not too many iterations and thus have low computational cost [6]. The main idea is to analyze the estimation problem which could have a benign landscape that would make problematic regions like local minima nonexistent. Specifically, the two chapters to follow will focus on two key ideas:

- **Basin of attraction**

For many problems, one can find a large basin of attraction (e.g.  $L_2$  sphere) centered around the global solution. This set is convex around the minimizer so clearly, if the initialization begins within the basin, then iterative methods like gradient descent are guaranteed to converge and do so quickly. The idea of how to initialize within the basin is discussed in the section about the spectral method.

- **Benign (Harmless) Landscape**

Many problems have a “benign global landscape” when the sample size is large enough. This means that there are no false local solutions (i.e. all local solutions are global solutions) and that the only unwanted critical point situations involve strict saddle points. As will be seen later on, it is surprisingly simple to escape saddle points and continue the descent in many situations.

These two key ideas inspire two different branches of nonconvex methods: 1) two-stage approaches and 2) global landscape analysis methods. The next chapter will focus on two-stage algorithms.

## 4 Two-Stage Algorithms

Two-stage algorithms are inspired by the idea of a cost function having a basin of attraction. Stage one involves initialization which will locate a guess that belongs to the basin. Stage two involves refining the solution iteratively without leaving the basin. Clearly, this is very efficient and essentially reduces to a convex problem (locally). There are three important two-stage algorithms that will be discussed in this paper: gradient descent, alternating minimization, and iterative hard thresholding.

### 4.1 Gradient Descent [7]

Gradient descent in a non-convex setting is expressed similarly to the convex setting. For example, with matrix sensing, using the matrix sensing cost function [EQ 2.1.1](#) and assuming symmetry of apriori sensing matrix ( $\mathbf{A}_i = \mathbf{A}_i^T$ ), the gradient rule update is expressed as the following:

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \alpha_k \underbrace{\frac{1}{m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X}_k \mathbf{X}_k^T \rangle - y_i) \mathbf{A}_i \mathbf{X}_k}_{\nabla f(\mathbf{X}_k)}$$

However, this rule will only converge if the sensing operator  $\mathcal{A}(\cdot)$  satisfies a condition known as RIP condition (will be discussed later on this page). The sensing operator is the “function” that maps the input matrix of interest  $\mathbf{X}$  to the observation domain processed by the sensing matrix. Stated mathematically, for matrix sensing cost function [EQ 2.1.1](#), the operator is defined as:

$$\mathcal{A}(\mathbf{X}) = [m^{-1/2} \langle \mathbf{A}_i, \mathbf{X} \rangle]_{1 \leq i \leq m} \quad (4.1.1)$$

$\mathcal{A}(\mathbf{X})$  must satisfy the RIP condition for convergence. The definition for RIP is written as follows.

**Definition 4.1.1** (Restricted Isometry Property (RIP)). *An operator  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  ([EQ 4.1.1](#)) is said to satisfy the  $r$ -RIP condition with constant  $\delta_r < 1$  if:*

$$(1 - \delta_r) \|\mathbf{X}\|_F^2 \leq \|\mathcal{A}(\mathbf{X})\|_2^2 \leq (1 + \delta_r) \|\mathbf{X}\|_F^2 \quad \forall \mathbf{X} \text{ where } \text{rank}(\mathbf{X}) \leq r \quad (4.1.2)$$

This definition states that the operator preserves the norm of the matrix when dealing with a rank- $r$  matrix. To put it more simply, imagine a triangle on the Cartesian 2D plane. The triangle can move, rotate, and be reflected but the distance between the vertices is preserved. The triangle is not allowed to be scaled or distorted in shape. This is what an isometry describes to give a conceptual example. Surprisingly, many random matrix sensing matrices satisfy this  $r$ -RIP condition if  $\mathbf{A}_i \sim \mathcal{N}(0, 1)$ . This is also usually the case in phase retrieval for  $\mathbf{a}_i$  as well. The RIP condition implies smoothness is well controlled within the basin and local strong convexity is achieved within the basin. Before discussing convergence for matrix sensing, due to the trace operator being unitary invariant, a more suitable error metric is the following:

$$\text{dist}(\mathbf{X}, \mathbf{X}_*) = \min_H \|\mathbf{X} \mathbf{H} - \mathbf{X}_*\|_F \text{ where } \mathbf{H} \text{ is a unitary matrix}$$

Given this metric above, the convergence of gradient descent for the matrix sensing problem is stated below.

**Theorem 4.1.2** (Matrix sensing (rank- $r$ ) convergence of GD). *The matrix sensing problem, with cost function [EQ 2.1.1](#), must have the  $\mathcal{A}(\mathbf{X})$  operator ([EQ 4.1.1](#)) satisfy the  $6r$ -RIP condition ([EQ 4.1.2](#)) with constant  $\delta_{6r} \leq \frac{1}{10}$ . Then,  $\exists c_0, c_1 > 0$  such that if  $\text{dist}(\mathbf{X}_0, \mathbf{X}_*) \leq \frac{\sigma_r(\mathbf{M}_*)}{16}$  then with  $\alpha_k = \frac{c_0}{\sigma_1(\mathbf{M}_*)}$  obeys convergence:*

$$\text{dist}(\mathbf{X}_k, \mathbf{X}_*) \leq (1 - \frac{c_1}{\kappa})^k \text{dist}^2(\mathbf{X}_0, \mathbf{X}_*) \quad \forall k = 0, 1, \dots$$

**Proof:** [Appendix A: Matrix Sensing Gradient Method Proof](#)

Thus, the distance between the initial solution  $\mathbf{X}_0$  and the global minima  $\mathbf{X}_*$  becomes increasingly smaller with convergence rate on par with classical gradient descent for convex problems. If the nonconvex problem does not satisfy RIP (e.g. phase retrieval), then the problem becomes challenging since the smoothness condition is not well-controlled and local strong convexity may fail so vanilla GD will get stuck in local minima. To resolve this dilemma, we wish to identify a constrained direction within a sphere that has nice geometry and then apply a regularized GD method to ensure the iterates stay in the constrained region set. This means that the descent direction must be along a specific trajectory to achieve smoothness and strong convexity even though the initializer is inside the local ball. The theorem below discusses desired smoothness conditions for phase retrieval since this usually does not satisfy RIP.

**Lemma 4.1.3** (Restricted smoothness (phase retrieval)). *For  $\mathbf{x} \in \mathbb{R}^n$ ,  $\exists c_0, \dots, c_3 > 0$  such that for arbitrary  $m \geq c_0 n \log(n)$  obeys the following for phase retrieval cost function [EQ 2.3.1](#):*

$$\nabla^2 f(\mathbf{x}) \preceq c_1 \log(n) \mathbf{I}_n$$

*but specifically  $\forall \mathbf{x}$  that meet  $\|\mathbf{x} - \mathbf{x}_*\|_2 \leq c_2$  and  $\max_j |\mathbf{a}_j^T(\mathbf{x} - \mathbf{x}_*)| \leq c_3 \sqrt{\log(n)}$ .*

This means that given an initializer within an  $c_2$ -ball, the hessian  $\nabla^2 f(\mathbf{x})$  will be well-behaved and smooth so long as the iterates are nearly orthogonal ( $\langle \cdot, \cdot \rangle \leq c_3$ ) to the sampling vectors  $\mathbf{a}_j$ . In the literature, this means  $\mathbf{x}$  is “incoherent” with  $\mathbf{a}$  and keeping the iterates nearly orthogonal ( $\leq c_3$ ) to sampling vectors results in convergence. Thus, the iterates must be kept in this set of near-orthogonality. Surprisingly, vanilla GD does this already for the problems mentioned above without any kind of explicit regularization. In other words, vanilla GD will force its iterates to stay incoherent with sampling vectors and matrices. Therefore, the vanilla GD trajectory, with high probability, already belongs to a region that is smooth and achieves strong convexity. However, there does exist explicit regularized GD methods to promote incoherence (via an additional term in the cost function) for various problems to improve computational performance. The basic structure looks like  $\min_{\mathbf{x}} f_{\text{reg}}(\mathbf{x}) = f(\mathbf{x}) + \lambda G(\mathbf{x})$  with  $\lambda$  as the regularization parameter to force incoherence. The function  $G(\mathbf{x})$  is different for each application and there is no “one size fits all” approach. For matrix sensing specifically, the following cost function has been proposed by [8]:

$$G(\mathbf{L}, \mathbf{R}) = G_0(c_1 \|\mathbf{L}\|_F^2) + G_0(c_2 \|\mathbf{R}\|_F^2) + \sum_{i=1}^{\text{rows}} G_0(c_3 \|\mathbf{L}_{i,\cdot}\|_2^2) + \sum_{i=1}^{\text{rows}} G_0(c_4 \|\mathbf{R}_{i,\cdot}\|_2^2)$$

for selected non-zero scalars  $c_1, \dots, c_4$  and  $G_0(x) = \max\{x - 1, 0\}^2$ . By doing this, the incoherence can be encouraged although as mentioned before, it is not required since vanilla GD does a great job already at ensuring incoherence.

## 4.2 Alternating Minimization [9]

Alternating minimization, under a matrix sensing situation, involves solving between the two sub-problems:

$$\mathbf{R}_{k+1} = \min_{\mathbf{R}} f(\mathbf{L}_k, \mathbf{R}) = \|\mathcal{A}(\mathbf{L}_k \mathbf{R}^T - \mathbf{M}_*)\|_F^2 \quad (4.2.1)$$

$$\mathbf{L}_{k+1} = \min_{\mathbf{L}} f(\mathbf{L}, \mathbf{R}_k) = \|\mathcal{A}(\mathbf{L} \mathbf{R}_k^T - \mathbf{M}_*)\|_F^2 \quad (4.2.2)$$

where both matrices are updated sequentially given initialization matrices  $\mathbf{L}_0, \mathbf{R}_0$ . Since minimization occurs with respect to one variable, each substep simplifies to the classical least-squares setup. Each step is then solved using conjugate gradient algorithm. If the reader is not familiar with this method, it is fairly similar to vanilla gradient descent except the current search direction is required to be orthogonal to the last search direction and so forth. Interestingly, if  $\mathbf{L}_0$  is initialized to be based on the singular vectors of  $\mathcal{A}^*(\mathbf{y})$ , then convergence is guaranteed. The following theorem describes this succinctly.

**Theorem 4.2.1** (Matrix Sensing AltMin Convergence). *Under the context of the general rank- $r$  cost function EQ 2.1.2, suppose the sensing operator satisfies 2 $r$ -RIP (EQ 4.1.2) with RIP constant  $\delta_{2r} \leq \frac{1}{100\kappa^2 r}$ . If  $\mathbf{L}_0$  represents the first  $r$  left singular vectors of  $\mathcal{A}^*(\mathbf{y})$ , then AltMin has the following convergence:*

$$\|\mathbf{M}_* - \mathbf{L}_k \mathbf{R}_k^T\|_F \leq \epsilon, \forall k \geq 2 \log\left(\frac{\|\mathbf{M}_*\|_F}{\epsilon}\right)$$

As seen above, AltMin has better iteration complexity compared to vanilla GD. To be clear,  $\kappa$  above is the condition number for matrix  $\mathbf{M}_*$  and  $\mathcal{A}^*(\mathbf{y})$  is the adjoint operator that satisfies  $\langle \mathcal{A}(\mathbf{x}), \mathbf{y} \rangle = \langle \mathbf{x}, \mathcal{A}^*(\mathbf{y}) \rangle$ .

### 4.3 Iterative Hard Thresholding [10]

Iterative Hard Thresholding (IHT), while similar to gradient descent, iteratively updates in the full matrix space  $\mathbf{M}$  instead of  $\mathbf{L}, \mathbf{R}^T$  (where  $\mathbf{M} = \mathbf{L}\mathbf{R}^T$ ) and then performs hard threshold on the singular values to retain low-rank conditions. Notice that this problem remains nonconvex since  $\|\sigma(\mathbf{M})\|_0 = r$ . Explicitly, we have that:

$$\mathbf{M}_{k+1} = \mathcal{P}_r(\mathbf{M}_k - \alpha_k \nabla f(\mathbf{M}_k)) \quad \forall k = 0, 1, \dots$$

and for matrix sensing specifically we have that  $f(\mathbf{M}) = \frac{1}{2} \|\mathcal{A}(\mathbf{M}) - \mathcal{A}(\mathbf{M}_*)\|_F^2$ . The  $\mathcal{P}_r(\cdot)$  operator will return the best  $r$ -rank approximation by setting  $\sigma_{r+1}, \dots, \sigma_n = 0$ . One advantage is that the iterates are low rank so  $\mathbf{M}_k$  can be stored efficiently (memory size) through the use of a compact SVD. One important note to make is that a full SVD is not performed here when thresholding the singular values. Instead, in practice, a partial SVD is computed via Krylov-based subspace methods like Lanczos algorithm or other linear algebra algorithms. This is done for efficiency as computing a full SVD can be computationally expensive. Moreover, for matrix sensing, there are convergence guarantees for this IHT method.

**Theorem 4.3.1** (Matrix sensing convergence for IHT). *For the matrix sensing problem EQ 2.1.2, assume the sensing operator  $\mathcal{A}$  (EQ 4.1.1) satisfies 2 $r$ -RIP (EQ 4.1.2) with constant  $\delta_{2r} \leq \frac{1}{3}$ . For initial iterate  $\mathbf{M}_0 = \mathbf{0}$  and step size  $\alpha = \frac{1}{(1+\delta_{2r})}$ , then IHT achieves the following:*

$$\|\mathbf{M}_k - \mathbf{M}_*\|_F \leq \epsilon \text{ but only when } k \geq c_1 \log(\|\mathbf{M}_*\|_F/\epsilon) \text{ for some } c_1$$

Thus, it is possible to work with the full matrix space instead of each  $\mathbf{L}, \mathbf{R}^T$  matrix individually in the nonconvex setting. Now that these two two-stage algorithms have been discussed, the question remains on how to initialize such that the first guess is inside the basin of attraction. The next section will discuss this important question.

### 4.4 Spectral Method (Initialization Technique) [11]

Given the methods in the previous section, initialization plays a key role for basin of attraction methods. One well known approach involves the spectral method. Given a data sample matrix  $\mathbf{Y}$ , this is decomposed to the following form  $\mathbf{Y} = \mathbf{Y}_* + \Delta$  where  $\mathbf{Y}_*$  represents data whose eigenspace represents the truth and  $\Delta$  contains fluctuations due to finite-sample effect since the matrix is of finite size. The goal is to estimate the truth by analyzing the subspace information of  $\mathbf{Y}$  assuming the fluctuations are “well controlled”. Assume  $\mathbf{Y}$  is the data matrix with samples and  $\mathbf{Y}_*$  is a symmetric matrix with real eigenvalues such that  $\text{SVD}(\mathbf{Y}_*) \rightarrow \mathbf{U}_*$  the true subspace of the global minimizer. Say  $\mathbf{Y}_*$  is rank- $r$ , then it is assumed that the spectral-gap  $\lambda_r(\mathbf{Y}_*) - \lambda_{r+1}(\mathbf{Y}_*)$  is nonzero. As an example, if  $\mathbf{Y}_*$  is positive semi definite, then the spectral gap will clearly be nonzero for some  $r$ . It turns out that the subspace  $\mathbf{U}_Y$  is close to subspace  $\mathbf{U}_{Y_*}$  by the following theorem.

**Theorem 4.4.1** (Davis-Kahan sin  $\Theta$  Theorem). *Given the positive spectral gap, the error between subspaces for the true matrix and perturbed true matrix are expressed as follows:*

$$\text{dist}(U, U_*) = \|UU^T - U_*U_*^T\|_F \leq \frac{\|\Delta\|_2}{\lambda_r(\mathbf{Y}_*) - \lambda_{r+1}(\mathbf{Y}_*)}$$

Therefore, the perturbed subspace  $U$  of  $\mathbf{Y}$  is bounded by the spectral norm of the fluctuation matrix  $\Delta$  and the spectral gap. This is great news because it provides sound theory as to why the subspace of the data matrix would work as an initialization for two-stage algorithms since it is close to the true subspace  $U_*$  of the global minimizer. To provide an example, the spectral theory can be applied to matrix sensing. Let the surrogate matrix be expressed as  $\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{A}_i$  where the  $\mathbf{A}_i$  are the sensing matrices known beforehand and  $y_i$  is the observation data. Recall that we wish to recover the set of measurements  $\mathbf{M}_*$  but this data is “hidden” by the sensing matrices and only  $y_i$  is collected. This surrogate matrix  $\mathbf{Y}$  is decomposed via SVD such that  $\mathbf{L}_0 = \mathbf{U}\Sigma^{1/2}$  and  $\mathbf{R}_0 = \mathbf{V}\Sigma^{1/2}$ . Then, the end goal is to find  $\mathbf{L}_*\mathbf{R}_* = \mathbf{M}_*$ . If the sensing operator satisfies 2r-RIP then it is possible to bound the initializer with the global solution.

**Lemma 4.4.2** (Matrix Sensing Spectral Initialization). *Given that  $\mathbf{M}_*$  is rank- $r$  matrix and  $\mathcal{A}$  (EQ 4.1.1) satisfies 2r-RIP condition (EQ 4.1.2) with constant  $\delta_{2r} < 1$ , then the following holds:*

$$\|\mathbf{Y} - \mathbf{M}_*\| \leq \delta_{2r} \|\mathbf{M}_*\|_F \leq \delta_{2r} \sqrt{r} \|\mathbf{M}_*\|$$

**Proof:** [Appendix B: Matrix Sensing Spectral Method Proof](#)

Thus, the surrogate matrix  $\mathbf{Y}$  is bounded in distance to  $\mathbf{M}_*$ . Moreover, many surrogate matrices have been designed for various matrix factorization problems where the sensing operator satisfies RIP condition (EQ 4.1.2). For the case where RIP condition is not satisfied, it turns out that the subspace of a surrogate matrix can still be useful for initialization. Recall that phase retrieval (a specific matrix sensing problem) does not satisfy RIP but it is possible to initialize regardless so this will be used as an example as to why the method can still work. For phase retrieval, the sensing matrix is  $\mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^T$  so let the surrogate matrix be expressed similarly in that  $\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m (y_i \mathbf{a}_i \mathbf{a}_i^T)$ . For simplicity, assume the measurement matrix is rank-1  $\mathbf{M}_* = \mathbf{x}_* \mathbf{x}_*^T$  such that our initializer is a vector (specifically the first eigenvector of  $\mathbf{Y}$ ). The goal now becomes to figure out whether this surrogate is bounded with the final solution  $\mathbf{x}_*$ . The surrogate matrix, under a different interpretation, is really an average of  $m$  i.i.d. random rank-1 matrices. When the number of observations is large, using probability theory, the following is achieved:

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[y_i \mathbf{a}_i \mathbf{a}_i^T] = 2\mathbf{x}_* \mathbf{x}_*^T + \|\mathbf{x}_*\|_2^2 \mathbf{I}_n$$

This is fantastic since this tells us that the surrogate matrix  $\mathbf{Y}$  contains information about  $\mathbf{x}_*$  on average even though  $\mathbf{Y}$  does not contain any information about  $\mathbf{x}$  in the surrogate expression written above. Therefore,  $\mathbf{x}_0 = \sqrt{\lambda_1/3} \mathbf{u}_1$  is a valid initializer since first eigenvector  $\mathbf{u}_1$  of  $\mathbf{Y}$  contains information about  $\mathbf{x}_*$ .

## 5 Global Landscape Algorithms

Previously, the goal was to find a basin of attraction and initialize within it. In the global landscape setting, the entire parameter space is considered and the curvature properties of the cost function are considered. This is done by avoiding two kinds of situations: 1) local minima that are not the same value as the global minima and 2) saddle points that confuse convex algorithms. In order to discuss nonconvex methods, the notion of a benign landscape is introduced to differentiate from badly behaved landscapes.

**Definition 5.0.1** (Benign Landscape). *A function  $f(\cdot)$  is said to have a benign (meaning harmless) landscape iff all the critical points of  $f(\cdot)$  are either global minima or strict saddle points.*

The reason why only strict saddle points are considered (instead of degenerate saddle points) is because they are relatively easy to escape during the descent process with only second order derivative information whereas degenerate saddle points require higher order information. Besides matrix factorization, there are many popular problems where benign landscapes appear. For example, dictionary learning and statistical estimation have a nonconvex but benign landscape. An example of a nonconvex problem that does **not** have a benign landscape would be neural networks with nonlinear activation functions. In the matrix factorization context, the nonconvexity usually comes from symmetry of some kind. For example, matrix sensing operator involves the Frobenious norm (trace operation) which is unitary invariant so we have a “ring” of global minima versus one single minima. To take this further, matrix sensing convergence is guaranteed in a global landscape setting by the following theorem:

**Theorem 5.0.2** (Matrix Sensing Global Landscape Theorem). *Consider the matrix sensing operator  $A$  from EQ 4.1.1. If  $A$  satisfies  $2r$ -RIP condition (EQ 4.1.2) with constant  $\delta_{2r} < 1/10$ , then the global landscape meets the following:*

- All local minima are global minima
- All critical points that are not minima are strict saddle points, specifically for critical point  $\mathbf{X}$ ,  $\lambda_{\min}(\nabla^2 f(\mathbf{X})) \leq -4\sigma_r(\mathbf{M}_*)/5$

**Proof:** [Appendix C: Matrix Sensing Benign Landscape Proof](#)

In other words, the matrix sensing problem that satisfies RIP condition (EQ 4.1.2) will have a benign landscape (DF 5.0.1). If matrix sensing has a benign landscape, then the main goal becomes how to escape saddle points efficiently. Let us consider the famous gradient descent method.

## 5.1 Gradient-based Methods [12]

Before discussing gradient theory, it is important to consider what is known as the strict saddle property. Putting it simply, it states that every point in some function  $f(\cdot)$  either has a large gradient, a negative curvature, or a local minimum when there is a benign landscape. Mathematically, the strict saddle property is defined as follows:

**Definition 5.1.1** (Strict Saddle Property). *The function  $f(\cdot)$  satisfies  $(\epsilon, \gamma, \zeta)$ -strict saddle property for said constants if one of the following is met for each  $\mathbf{x}$ :*

- (strong gradient magnitude)  $\|\nabla f(\mathbf{x})\|_2 \geq \epsilon$
- (negative curvature)  $\lambda_{\min}(\nabla^2 f(\mathbf{x})) \leq -\gamma$
- (local minimum) For local minimizer  $\mathbf{x}_*$ ,  $\|\mathbf{x} - \mathbf{x}_*\|_2 \leq \zeta$

It turns out that gradient descent can converge for functions that meet this condition but only if initialized randomly. If the initialization occurs at a location where  $\nabla f(\mathbf{x}) = \mathbf{0}$  then gradient descent will get stuck and unable to escape. However, if the initialization is random, then it can avoid saddle points during the descent since the probability of this event is zero. This powerful idea is expressed in the following theorem.

**Theorem 5.1.2** (GD Convergence - Random Initialization Case). *Assume  $f(\cdot)$  (EQ 2.1.1) satisfies the strict saddle property (DF 5.1.1). This implies a twice continuously differentiable function meaning it is  $\beta$ -smooth for some  $\beta$ . If the step size is such that  $\alpha_k < 1/\beta$  then GD under random initialization converges (in the almost surely sense) to a local minimizer or  $-\infty$ .*



Conceptually, GD iterations will not even reach the saddle point (under random initialization) since there is always negative curvature around that area that forces the GD trajectory away from any saddle point. Since phase retrieval does not satisfy RIP condition (EQ 4.1.2) but satisfies the strict saddle property (DF 5.1.1), then convergence is guaranteed under the random initialization setting.

**Theorem 5.1.3** (Phase retrieval convergence for randomized GD method). *For a sufficiently large iteration number  $k$ , under a random initialization  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \frac{\|\mathbf{x}_*\|_2^2}{n} \mathbf{I}_n)$ , and step size  $\alpha = \frac{1}{c_1 \|\mathbf{x}_*\|_2^2}$ , the following bound is achieved (for cost function EQ 2.3.1):*

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq (1 - c_2)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \text{ with a very high probability for some constants } c_1, c_2 > 0$$

This particular GD method is useful for large-scale problems due to fast convergence. Designing other minimization algorithms besides gradient descent is a topic for the next section.

## 5.2 Hessian-based Methods [13] [14]

Before talking about trust-region methods, the reader should know that both trust-region methods and cubic regularization fall under a broad umbrella of algorithms known as generic saddle-escaping algorithms. To explain what these methods do, it is important to consider the Taylor expansion of any saddle point. Given a saddle point  $\mathbf{x}$ , the following Taylor expansion

$$f(\mathbf{x} + \boldsymbol{\epsilon}) \approx f(\mathbf{x}) + \frac{1}{2} \boldsymbol{\epsilon}^T \nabla^2 f(\mathbf{x}) \boldsymbol{\epsilon}$$

shows that  $f(\mathbf{x} + \boldsymbol{\epsilon}) < f(\mathbf{x})$  since the  $\nabla^2 f(\mathbf{x})$  term will be negative. This is remarkable since it is possible to identify a direction where the curvature is negative and continue to decrease the objective function (therefore escape saddle points). One way to use this property is to use an algorithm that involves the Hessian  $\nabla^2 f(\mathbf{x})$  which is the basis for trust-region methods. The trust-region method is written as follows for some  $\mathbf{x}$ :

$$\mathbf{x}_{k+1} = \min_{\mathbf{z}: \|\mathbf{z} - \mathbf{x}_k\| \leq \Delta} \{ \langle \nabla f(\mathbf{x}_k), \mathbf{z} - \mathbf{x}_k \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_k)(\mathbf{z} - \mathbf{x}_k), \mathbf{z} - \mathbf{x}_k \rangle \}$$

where  $\Delta$  is a user chosen radius known as the “trust region”. From the Taylor expansion written above, we can expect a reasonable approximation only if the trust region is small. If the current iterate happens to be the saddle point then the method above will move in the direction of the negative curvature since  $\nabla^2 f(\mathbf{x})(\mathbf{z} - \mathbf{x})$  is negative. Thus, the method will escape a saddle point more efficiently than GD would if the current iterate happens to meet the condition. On the other hand, cubic regularization explicitly adds a cubic term to the cost function so that the minimization problem becomes unconstrained with respect to  $\mathbf{z}$ . This cubic regularized minimization expression can be written as follows:

$$\mathbf{x}_{k+1} = \min_{\mathbf{z}} \{ \langle \nabla f(\mathbf{x}_k), \mathbf{z} - \mathbf{x}_k \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_k)(\mathbf{z} - \mathbf{x}_k), \mathbf{z} - \mathbf{x}_k \rangle + \frac{L_2}{6} \|\mathbf{z} - \mathbf{x}_k\|_2^3 \}$$

where  $L_2$  corresponds to the Lipschitz constant of the Hessian of the cost function. In this case, the Lipschitz constant is selected since it will allow the new cubic-regularized function to *majorize* the objective function  $f(\mathbf{x})$ . The majorization definition is written below.

**Definition 5.2.1** (Majorization). *For vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ ,  $\mathbf{a}$  is said to majorize  $\mathbf{b}$  iff  $\sum_{i=1}^k \mathbf{a}_i \geq \sum_{i=1}^k \mathbf{b}_i$  for  $k = 1, \dots, d$ .*

Without diving into majorization theory, functions that preserve the ordering of majorization (e.g. this cubic-regularized function) are Schur-convex functions. While a Schur-convex function does not imply a convex function, gradient descent can be used on Schur-convex functions to efficiently solve the minimization problem.

### 5.3 Almost-Convex-AGD & Negative-Curvature-Descent [15]

Both Almost-Convex-AGD and Negative-Curvature-Descent make up one complete algorithm that consists of the two subroutines. Almost-Convex-AGD uses Nesterov’s accelerated gradient method to minimize an almost convex function (hence the name) that combines the objective function with a distance term. Negative-Curvature-Descent will find a direction that decreases the function value by calculating the smallest eigenvector of the Hessian and thus move in the negative curvature direction and escape saddle points. The overall scheme of the algorithm is that for each iterate, it will first calculate the smallest eigenvalue for the Hessian and decide whether to move in the direction computed by Negative-Curvature-Descent (to avoid saddle points) or either to apply Almost-Convex-AGD to optimize an almost convex function. Almost-Convex-AGD is exactly the same as Nesterov’s method except the cost function becomes  $f(\mathbf{z}) + \gamma\|\mathbf{z} - \mathbf{z}_i\|^2$  where  $\gamma < L_1$  (Lipschitz constant of the gradient) and accelerated GD is performed on the modified (“convex”) cost function. This is done by coming up with a momentum term  $\mathbf{p}_{k+1} = \mathbf{z}_k + \beta_k(\mathbf{z} - \mathbf{z}_{k-1})$  and using that in vanilla GD such that  $\mathbf{z}_{k+1} = \mathbf{p}_{k+1} - \tau\nabla f(\mathbf{p}_{k+1})$  where  $\beta$  is the momentum weight and  $\tau$  is the step size. For the other subroutine, the goal of Negative-Curvature-Descent is to find a vector  $\mathbf{v}_j$  where  $\|\mathbf{v}_j\| \leq 1$  and meets the following condition for current iterate  $\mathbf{x}_k$ :

$$\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) \geq \mathbf{v}_j^T \nabla^2 f(\mathbf{x}_k) \mathbf{v}_j$$

Then if  $\mathbf{v}_j^T \nabla^2 f(\mathbf{x}_k) \mathbf{v}_j < 0$  we will either take a step in the negative direction computed (i.e.  $\mathbf{v}_j$ ) of size proportional to  $\mathbf{v}_j^T \nabla^2 f(\mathbf{x}_k) \mathbf{v}_j$  or otherwise retain current iterate  $\mathbf{x}_k$ . In other words, it will stay at the current iterate if a negative descent direction cannot be calculated (since it is away from saddle points) or descend based on negative curvature if near a saddle point. Specifically, the new iterate is written as follows:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{2|\mathbf{v}_j^T \nabla^2 f(\mathbf{x}_k) \mathbf{v}_j|}{L_2} \text{sign}(\mathbf{v}_j^T \nabla^2 f(\mathbf{x}_k) \mathbf{v}_j) \mathbf{v}_j$$

where  $L_2$  represents the Lipschitz constant of the Hessian. It is important to note that, unlike trust region method and cubic regularization, this method does not use the Hessian directly. Instead, it approximates it by calculating the following:

$$\nabla^2 f(\mathbf{x}) \mathbf{v} = \lim_{h \rightarrow 0} \frac{\nabla f(\mathbf{x} + h\mathbf{v}) - \nabla f(\mathbf{x})}{h}$$

## 6 Conclusion

In this paper, nonconvex optimization algorithms are discussed within the context of matrix factorization problems. Canonical problems like matrix sensing, phase retrieval, and robust PCA were introduced to show a representative sample of the world of matrix factorization. For the sake of simplicity and consistency, special attention was placed to matrix sensing for each topic. Thereafter, two-stage algorithms and global landscape algorithms were discussed that take advantage of two key ideas: basin of attraction and benign landscape. Specifically, for two-stage algorithms (basin of attraction), gradient descent method and alternating minimization methods are useful in practice so long as the initialization process is carefully designed via the spectral method. Moreover, for global landscape algorithms (benign landscape), gradient methods were discussed for strict saddle functions and generic saddle-escaping algorithms such as trust-region methods and AGD-NCD. Throughout these sections, theorems are also included for theoretical guarantees under a matrix factorization setting. Nonconvex methods are extremely important for a wide range of problems as shown in this paper. However, some attention should be placed on whether it is optimal for the specific problem. It may be that a convex relaxation approach would be more appropriate instead of a nonconvex formulation depending on the situation. Each specific problem has different conditions that determine whether or not a nonconvex approach would be ideal so individual analysis is required before applying any such methods discussed.

## References

- [1] Y. Chi, Y. M. Lu, and Y. Chen, “Nonconvex optimization meets low-rank matrix factorization: An overview,” *CoRR*, vol. abs/1809.09573, 2018.
- [2] J. Wright and Y. Ma, *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. Cambridge University Press, 2022.
- [3] E. J. Candès and Y. Plan, “Tight oracle bounds for low-rank matrix recovery from a minimal number of random measurements,” *CoRR*, vol. abs/1001.0339, 2010.
- [4] N. Vaswani, T. Bouwmans, S. Javed, and P. Narayanamurthy, “Robust PCA and robust subspace tracking,” *CoRR*, vol. abs/1711.09492, 2017.
- [5] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev, “Phase retrieval with application to optical imaging,” *CoRR*, vol. abs/1402.7350, 2014.
- [6] Y. Chen and E. J. Candès, “Solving random quadratic systems of equations is nearly as easy as solving linear systems,” *CoRR*, vol. abs/1505.05114, 2015.
- [7] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht, “Low-rank solutions of linear matrix equations via procrustes flow,” in *International Conference on Machine Learning*, pp. 964–973, PMLR, 2016.
- [8] R. Sun and Z.-Q. Luo, “Guaranteed matrix completion via non-convex factorization,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.
- [9] P. Jain, P. Netrapalli, and S. Sanghavi, “Low-rank matrix completion using alternating minimization,” in *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pp. 665–674, 2013.
- [10] P. Jain, R. Meka, and I. Dhillon, “Guaranteed rank minimization via singular value projection,” *Advances in Neural Information Processing Systems*, vol. 23, 2010.
- [11] C. Davis and W. M. Kahan, “The rotation of eigenvectors by a perturbation. iii,” *SIAM Journal on Numerical Analysis*, vol. 7, no. 1, pp. 1–46, 1970.
- [12] J. Sun, *When are nonconvex optimization problems not scary?* Columbia University, 2016.
- [13] A. R. Conn, N. I. Gould, and P. L. Toint, *Trust region methods*. SIAM, 2000.
- [14] Y. Nesterov and B. T. Polyak, “Cubic regularization of newton method and its global performance,” *Mathematical Programming*, vol. 108, no. 1, pp. 177–205, 2006.
- [15] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford, “Accelerated methods for nonconvex optimization,” *SIAM Journal on Optimization*, vol. 28, no. 2, pp. 1751–1772, 2018.
- [16] S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Global optimality of local search for low rank matrix recovery,” *Advances in Neural Information Processing Systems*, vol. 29, 2016.

## 7 Appendix A: Matrix Sensing Gradient Method Proof [1]

**Theorem 7.0.1** (Matrix sensing (rank-r) convergence of GD). *The matrix sensing problem (EQ 2.1.2) must have the  $\mathcal{A}(\mathbf{X})$  operator (EQ 4.1.1) satisfy the 6r-RIP condition (EQ 4.1.2) with constant  $\delta_{6r} \leq \frac{1}{10}$ . Then,  $\exists c_0, c_1 > 0$  such that if  $\text{dist}(\mathbf{X}_0, \mathbf{X}_*) \leq \frac{\sigma_r(M_*)}{16}$  then with  $\alpha_k = \frac{c_0}{\sigma_1(M_*)}$  obeys convergence:*

$$\text{dist}(\mathbf{X}_k, \mathbf{X}_*) \leq \left(1 - \frac{c_1}{\kappa}\right)^k \text{dist}^2(\mathbf{X}_0, \mathbf{X}_*) \quad \forall k = 0, 1, \dots$$

**Proof (rank-1 case):**

For arbitrary vectors  $\mathbf{x}$  and  $\Delta \in \mathbb{R}^n$  and  $\mathbf{D} = \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T$  we have that the quadratic form of matrix sensing cost function stated as:

$$\Delta^T \nabla^2 f(\mathbf{x}) \Delta = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{D} \rangle (\Delta^T \mathbf{A}_i \Delta) + 2(\Delta^T \mathbf{A}_i \mathbf{x})^2$$

This is a lemma from [16]. Introducing the sensing operator  $\mathcal{A}(\cdot)$ , the expression becomes the following:

$$\Delta^T \nabla^2 f(\mathbf{x}) \Delta = \langle \mathcal{A}(\mathbf{D}), \mathcal{A}(\Delta \Delta^T) \rangle + \frac{1}{2} \langle \mathcal{A}(\Delta \mathbf{x}^T + \mathbf{x} \Delta^T), \mathcal{A}(\Delta \mathbf{x}^T + \mathbf{x} \Delta^T) \rangle$$

The goal is to show that

$$\Delta^T \nabla^2 f(\mathbf{x}) \Delta \approx \langle \mathbf{D}, \Delta \Delta^T \rangle + \frac{1}{2} \|\Delta \mathbf{x}^T + \mathbf{x} \Delta^T\|_F^2 = g(\mathbf{x}, \Delta)$$

The reason is that  $g(\cdot)$  has local strong convexity and is smooth. This means that convergence is guaranteed. This is expressed by the following lemma.

**Lemma 7.0.2** (GD Convergence (Strong Convexity & Smoothness)). *Suppose that function  $f(\cdot)$  is  $\alpha$ -strongly convex and  $\beta$ -smooth within an L2 sphere  $\mathcal{B}[\mathbf{x}_*, \zeta]$ . If  $\mathbf{x}_0 \in \mathcal{B}[\cdot]$  and  $\tau = \frac{1}{\beta}$ , then the following is achieved:*

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \quad \forall k = 0, 1, \dots$$

Going back to the problem at hand, when  $\mathcal{A}(\cdot)$  satisfies 4-RIP (EQ 4.1.2), then by definition, we have that

$$|\langle \mathcal{A}(\mathbf{D}), \mathcal{A}(\Delta \Delta^T) \rangle - \langle \mathbf{D}, \Delta \Delta^T \rangle| \leq \delta \|\mathbf{D}\|_F \|\Delta \Delta^T\|_F \leq 3\delta \|\mathbf{x}_*\|_2^2 \|\Delta\|_2^2 \text{ assuming } \|\mathbf{x} - \mathbf{x}_*\|_2 \leq \|\mathbf{x}_*\|_2$$

Likewise, applying the same logic, we can get another bound:

$$|\langle \mathcal{A}(\Delta \mathbf{x}^T + \mathbf{x} \Delta^T), \mathcal{A}(\Delta \mathbf{x}^T + \mathbf{x} \Delta^T) \rangle - \|\Delta \mathbf{x}^T + \mathbf{x} \Delta^T\|_F^2| \leq 4\delta \|\mathbf{x}\|_2^2 \|\Delta\|_2^2 \leq 16\delta \|\mathbf{x}_*\|_2^2 \|\Delta\|_2^2$$

So long as  $\delta \leq \frac{1}{44}$ , then with some norm inequality properties, it is possible to get the result:

$$\frac{1}{4} \|\Delta\|_2^2 \|\mathbf{x}_*\|_2^2 \leq \Delta^T \nabla^2 f(\mathbf{x}) \Delta \leq 3 \|\Delta\|_2^2 \|\mathbf{x}_*\|_2^2 \quad \forall \Delta$$

Thus, we have bounds on  $(\alpha, \beta)$  the smoothness and convexity parameters. Because of this, through the lemma stated above, matrix sensing is guaranteed to converge using gradient descent method so long as the sensing operator satisfies the RIP condition (EQ 4.1.2).

## 8 Appendix B: Matrix Sensing Spectral Method Proof [1]

**Lemma 8.0.1** (Matrix Sensing Spectral Initialization). *Given that  $M_*$  is rank- $r$  matrix and  $\mathcal{A}$  (EQ 4.1.1) satisfies 2 $r$ -RIP condition (EQ 4.1.2) with constant  $\delta_{2r} < 1$ , then the following holds for surrogate matrix  $\mathbf{Y}$ :*

$$\|\mathbf{Y} - M_*\| \leq \delta_{2r} \|M_*\|_F \leq \delta_{2r} \sqrt{r} \|M_*\|_F$$

**Proof:**

For arbitrary vectors  $\mathbf{x}$  and  $\mathbf{y}$  we have that,

$$\begin{aligned} & \mathbf{x}^T (M_* - \mathbf{Y}) \mathbf{y} \\ &= \mathbf{x}^T (M_* - \mathcal{A}^* \mathcal{A}(M_*)) \mathbf{y} && \text{(per definition of the surrogate matrix)} \\ &= \langle \mathbf{x} \mathbf{y}^T, M_* - \mathcal{A}^* \mathcal{A}(M_*) \rangle && \text{(per inner product definition)} \\ &= \langle \mathbf{x} \mathbf{y}^T, M_* \rangle - \langle \mathcal{A}(\mathbf{x} \mathbf{y}^T), \mathcal{A}(M_*) \rangle && \text{(per inner product properties)} \\ &\leq \delta_{2r} (\|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2) \|M_*\|_F && \text{(per 2r-RIP property definition)} \\ &\implies \|M_* - \mathbf{Y}\| = \max_{\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1} \mathbf{x}^T (M_* - \mathbf{Y}) \mathbf{y} && \text{(spectral norm definition)} \\ &\leq \delta_{2r} \|M_*\|_F \leq \delta_{2r} \sqrt{r} \|M_*\| && \text{(using above bound)} \\ &\implies \|\mathbf{Y} - M_*\| \leq \delta_{2r} \|M_*\|_F \leq \delta_{2r} \sqrt{r} \|M_*\| \quad \square \end{aligned}$$

Thus, the surrogate matrix  $\mathbf{Y}$  is bounded to the minimizer  $M_*$  by some distance proportional to  $\|M_*\|_F$ . Because of this, the surrogate matrix is a valid initialization for any two-stage algorithm.

## 9 Appendix C: Matrix Sensing Benign Landscape Proof [1]

**Theorem 9.0.1** (Matrix Sensing Global Landscape Theorem). *Consider the matrix sensing operator  $\mathcal{A}$  (EQ 4.1.1). If  $\mathcal{A}$  satisfies 2r-RIP condition (EQ 4.1.2) with constant  $\delta_{2r} < 1/10$ , then the global landscape is benign (DF 5.0.1) and meets the following conditions for cost function EQ 2.1.1:*

- All local minima are global minima
- All critical points that are not minima are strict saddle points, specifically for critical point  $\mathbf{X}$ ,  $\lambda_{\min}(\nabla^2 f(\mathbf{X})) \leq -4\sigma_r(\mathbf{M}_*)/5$

**Proof:**

For rank-1 case with arbitrary vector  $\mathbf{x}$  that is a critical point for the matrix sensing problem we have that,

$$\nabla f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T \rangle \mathbf{A}_i \mathbf{x} = \mathbf{0} \quad (\text{first order optimality condition})$$

$$\implies \|(\mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T)\mathbf{x}\mathbf{x}^T\|_F \leq \delta \|\mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T\|_F \cdot \|\mathbf{x}\|_2^2 \quad (\text{direct consequence from above})$$

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m 2\mathbf{A}_i \mathbf{x}\mathbf{x}^T \mathbf{A}_i + \langle \mathbf{A}_i, \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T \rangle \mathbf{A}_i \succeq \mathbf{0} \quad (\text{second order optimality condition})$$

Let  $\Delta = \mathbf{x} - \mathbf{x}_*$  where  $\mathbf{x}_*$  is the minimizer. Then,

$$\Delta^T \nabla^2 f(\mathbf{x}) \Delta \quad (\text{quadratic form of function})$$

$$= \frac{1}{m} \sum_{i=1}^m 2\langle \mathbf{A}_i, \mathbf{x}\Delta^T \rangle^2 + \langle \mathbf{A}_i, \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T \rangle \langle \mathbf{A}_i, \Delta\Delta^T \rangle$$

$$= \frac{1}{m} \sum_{i=1}^m 2\langle \mathbf{A}_i, \mathbf{x}\Delta^T \rangle^2 - \langle \mathbf{A}_i, \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T \rangle^2$$

$$\implies \Delta^T \nabla^2 f(\mathbf{x}) \Delta \leq 2(1 + \delta) \|\mathbf{x}\Delta^T\|_F^2 - (1 - \delta) \|\mathbf{D}\|_F^2 \quad (\text{using 2r-RIP condition})$$

Where  $\mathbf{D} = \mathbf{x}\mathbf{x}^T - \mathbf{x}_*\mathbf{x}_*^T$ . Without proof, the bound is stated:

$$\|\mathbf{x}\Delta^T\|_F^2 \leq \frac{1}{8} \|\mathbf{D}\|_F^2 + \frac{34}{8\|\mathbf{x}\|_2^2} \|\mathbf{D}\mathbf{x}\mathbf{x}^T\|_F^2 \quad (\text{Lemma from [16]})$$

$$\implies \Delta^T \nabla^2 f(\mathbf{x}) \Delta$$

$$\leq \frac{5\delta - 3}{4} \|\mathbf{D}\|_F^2 + \frac{34 + 34\delta}{4\|\mathbf{x}\|_2^2} \|\mathbf{D}\mathbf{x}\mathbf{x}^T\|_F^2 \quad (\text{using bound above on quadratic form})$$

$$\leq -\gamma \|\mathbf{D}\|_F^2 \text{ where } 0 < \gamma < 1 \quad (\text{using first order optimality condition})$$

$$\implies \Delta^T \nabla^2 f(\mathbf{x}) \Delta < 0 \quad \forall \Delta \text{ unless } \Delta = \mathbf{0}$$

$$\Delta = \mathbf{0} \implies \mathbf{x} = \mathbf{x}_* \quad \square$$

Therefore, either  $\mathbf{x}$  (a critical point) is the minimizer  $\mathbf{x}_*$  or it is a strict saddle point since the quadratic form is negative as shown above.