
Rank-1 Matrix Completion with Gradient Descent and Small Random Initialization

Daesung Kim
Samsung Electronics
dskim95phd@gmail.com

Hye Won Chung
School of Electrical Engineering
KAIST
hwchung@kaist.ac.kr

Abstract

The nonconvex formulation of the matrix completion problem has received significant attention in recent years due to its affordable complexity compared to the convex formulation. Gradient Descent (GD) is a simple yet efficient baseline algorithm for solving nonconvex optimization problems. The success of GD has been witnessed in many different problems in both theory and practice when it is combined with random initialization. However, previous works on matrix completion require either careful initialization or regularizers to prove the convergence of GD. In this paper, we study the rank-1 symmetric matrix completion and prove that GD converges to the ground truth when small random initialization is used. We show that in a logarithmic number of iterations, the trajectory enters the region where local convergence occurs. We provide an upper bound on the initialization size that is sufficient to guarantee the convergence, and show that a larger initialization can be used as more samples are available. We observe that the implicit regularization effect of GD plays a critical role in the analysis, and for the entire trajectory, it prevents each entry from becoming much larger than the others.

1 Introduction

Recovering a low-rank matrix from a set of linear measurements is at the heart of many statistical learning problems. Depending on the structure of the matrix and the linear measurements, it reduces to various problems such as phase retrieval [1], blind deconvolution [2], and matrix sensing [3]. Matrix completion [4] is also one such type of problem where each measurement provides an entry of the matrix, and the goal is to recover the low-rank matrix from a partial, usually very sparse, observation of the entries. One of the most notable applications of matrix completion is collaborative filtering [5], which aims to predict user preferences for items based on a highly incomplete observation of user-item ratings. There are also a number of different applications, such as principal component analysis [6] and image reconstruction [7], just to name a few.

Extensive amount of work has been dedicated to provide an efficient recovery algorithm for matrix completion with theoretical guarantees [8]. The convex relaxation based nuclear norm minimization [4, 9] was the first algorithm proven to recover the matrix with near optimal sample complexity. Despite its theoretical success, the convex algorithm was found hard to be used in practical scenarios due to its unaffordable computational complexity and memory size. Therefore, the nonconvex formulation of matrix completion with quadratic loss has received significant attention in recent years. Many different algorithms have been proposed for the nonconvex problem, and their convergence toward the ground truth has been analyzed. Examples include optimization on Grassmann manifolds [10], alternating minimization [11], projected gradient descent [12], gradient descent with regularizer [13], and (vanilla) gradient descent [14, 15].

Gradient descent (GD) has served as a baseline algorithm for solving nonconvex optimization problems. However, the convergence of GD to global minimizers is not guaranteed, and it can take exponential time to escape saddle points [16]. Nevertheless, GD with random initialization has been shown to successfully recover the global minimum in many different problems such as phase retrieval [1], matrix sensing [17], matrix factorization [18], and neural network training [19]. Previous work on matrix completion [14, 15] proved the convergence of GD under the spectral initialization, which locates the initial point in the local region of the minima. However, the role of random initialization in solving matrix completion with GD is not fully understood yet, although its success is observed in practice. Therefore, we aim to answer the following question:

Can GD with random initialization solve the nonconvex matrix completion problem?

We answer this question affirmatively and show that GD with small random initialization successfully converges to the ground truth for rank-1 symmetric matrix completion. In the analysis, we use vanilla GD, which does not incorporate any modifications, such as regularization or truncation, into the GD algorithm. We also characterize the entire trajectory that GD follows by showing that the trajectory is well approximated by the fully observed case. The small initialization plays a critical role in analyzing the trajectory of the early stages, where the randomly initialized vector is nearly orthogonal to the first eigenvector of the ground truth matrix. We provide a bound on the required initialization size for the algorithm to converge, and our bound suggests that one can use a larger initialization to improve the convergence speed as more samples are provided. However, in any case, GD with a small random initialization takes only logarithmic amount of time (with respect to the matrix dimension) to reach the point where local convergence can begin. To the best of our knowledge, this is the first result on matrix completion that proves the convergence of vanilla GD without a carefully designed initialization.

Although our result is restricted to the rank-1 case, we believe that this work provides an important evidence for understanding the more general rank- r case. At the end of this paper, we will discuss some technical difficulties that the rank- r case naturally has, and provide some empirical results related to them. However, studying the rank-1 matrix completion problem is not only motivated by theoretical interest, but the problem itself also appears in some practical problems such as crowdsourcing [20, 21].

Related Works This work is motivated by the recent success of small initialization in matrix factorization and matrix sensing. It was first conjectured in [22] that sufficiently small step sizes and initialization lead GD to converge to the minimum nuclear norm solution of a full-dimensional matrix sensing problem. The conjecture was proved in [17] for the fully overparameterized matrix sensing under the standard restricted isometry property (RIP). A recent study by [23] provided more general results by showing that the early iterations of GD with small initialization have spectral bias. Many other works such as [24, 25, 26] have also studied how GD or gradient flow with small initialization implicitly forces the recovered matrix to be low-rank. However, the recovery guarantee for matrix completion has not been provided by any work.

For the matrix sensing where RIP holds, the loss function has global benign geometry in that it does not contain any spurious local minima or non-strict saddle points [27]. In the case of matrix completion, a similar result was obtained but with a regularizer that penalizes the matrices with large rows [28]. Controlling the norm of each row (absolute value of each entry in the case of rank-1) is the biggest hurdle in the analysis of matrix completion. In the local convergence analysis of [14], it was proved that GD implicitly regularizes the largest ℓ_2 -norm of the rows of error matrices, showing that explicit regularization is unnecessary. In this paper, we also prove that such an implicit regularization is induced by GD when it starts from a point of small size. We show that the trajectory is close to the fully observed case in both ℓ_2 and ℓ_∞ norms. Thus, the trajectory is confined to the region where it has benign geometry, and GD can converge without an explicit regularizer.

Notations We denote vectors with lowercase bold letters and matrices with uppercase bold letters. The components or entries of them are written without bold. We use $\|\cdot\|_2$ and $\|\cdot\|_\infty$ to denote ℓ_2 and ℓ_∞ -norm of vectors, respectively, and $\|\cdot\|_F$ is used for Frobenius norm of matrices. For any norm $\|\cdot\|$ and two vectors \mathbf{x}, \mathbf{y} , we let $\|\mathbf{x} \pm \mathbf{y}\| = \min\{\|\mathbf{x} + \mathbf{y}\|, \|\mathbf{x} - \mathbf{y}\|\}$. Asymptotic dependencies with respect to the matrix dimension are denoted with the standard big O notations, or with the symbols, \lesssim , \asymp and \gtrsim .

2 Problem Formulation

The matrix completion problem aims to reconstruct a low-rank matrix from partially observed entries. In this paper, we focus on the case where the ground truth matrix, denoted by $\mathbf{M}^* \in \mathbb{R}^{n \times n}$, is a rank-1 positive semidefinite matrix. Thus, the ground truth matrix is decomposed as $\mathbf{M}^* = \lambda^* \mathbf{u}^* \mathbf{u}^{*\top}$ with $\lambda^* > 0$ and a unit vector \mathbf{u}^* . We define $\mathbf{x}^* = \sqrt{\lambda^*} \mathbf{u}^*$ so that $\mathbf{M}^* = \mathbf{x}^* \mathbf{x}^{*\top}$. To follow the standard incoherence assumption, we let $\|\mathbf{u}^*\|_\infty = \sqrt{\frac{\mu}{n}}$ and allow μ to be as large as $\text{poly}(\log n)$. We consider a random sampling model that is also symmetric as \mathbf{M}^* . Each entry in the diagonal and the upper (or lower) triangular part of \mathbf{M}^* is independently revealed with probability $0 < p \leq 1$. We consider the noisy case where Gaussian noise is added to each observation. Formally, we get as an observation the matrix \mathbf{M}° whose (i, j) th entry is $\frac{1}{p} \delta_{ij} (M_{ij}^* + E_{ij})$, where $[\delta_{ij}]_{1 \leq i \leq j \leq n}$ are independent Bernoulli random variables with expectation p and $[E_{ij}]_{1 \leq i \leq j \leq n}$ are independent Gaussian random variables with the distribution $\mathcal{N}(0, \sigma^2)$. They are both symmetric in the sense that $\delta_{ij} = \delta_{ji}$ and $E_{ij} = E_{ji}$ for all $1 \leq i \leq j \leq n$. We use \mathbf{E} to denote the symmetric matrix whose entries are E_{ij} . We denote the set of observed entries as $\Omega := \{(i, j) \mid \delta_{ij} = 1\}$, and define an operator \mathcal{P}_Ω on matrices that sets the entries not contained in Ω to zero. (e.g. $\mathbf{M}^\circ = \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^* + \mathbf{E})$)

To recover the matrix \mathbf{M}^* , we find $\mathbf{x} \in \mathbb{R}^n$ that minimizes the nonconvex loss function $f(\mathbf{x})$, which is the sum of the squared differences on the observed entries. It is explicitly written as $f(\mathbf{x}) := \frac{1}{4p} \sum_{(i,j) \in \Omega} (x_i x_j - x_i^* x_j^* - E_{ij})^2$. We apply vanilla GD to solve the optimization problem starting from a small randomly initialized vector $\mathbf{x}^{(0)}$. Each entry of $\mathbf{x}^{(0)}$ is sampled independently from the Gaussian distribution $\mathcal{N}(0, \frac{1}{n} \beta_0^2)$, so that the squared norm of $\mathbf{x}^{(0)}$ is expected to be β_0^2 . The update rule of GD is written as

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla f(\mathbf{x}^{(t)}) = \mathbf{x}^{(t)} - \frac{\eta}{p} \mathcal{P}_\Omega(\mathbf{x}^{(t)} \mathbf{x}^{(t)\top}) \mathbf{x}^{(t)} + \eta \mathbf{M}^\circ \mathbf{x}^{(t)}, \quad (1)$$

where $\eta > 0$ is the step size.

We define F as the loss function f when all entries of \mathbf{M}^* are observed without noise, i.e., $F(\mathbf{x}) := \frac{1}{4} \|\mathbf{x} \mathbf{x}^\top - \mathbf{M}^*\|_F^2$. We also define $\tilde{\mathbf{x}}^{(t)}$ as the trajectory of GD when it is applied to F with the same initial point $\mathbf{x}^{(0)}$, i.e., $\tilde{\mathbf{x}}^{(t)}$ is the trajectory of the fully observed case. Specifically, it evolves with

$$\tilde{\mathbf{x}}^{(t+1)} = \tilde{\mathbf{x}}^{(t)} - \eta \nabla F(\tilde{\mathbf{x}}^{(t)}) = \tilde{\mathbf{x}}^{(t)} - \eta \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2^2 \tilde{\mathbf{x}}^{(t)} + \eta \mathbf{M}^* \tilde{\mathbf{x}}^{(t)} \quad (2)$$

from the same starting point $\tilde{\mathbf{x}}^{(0)} = \mathbf{x}^{(0)}$.

Lastly, we introduce the so-called *leave-one-out* sequences. These were the main ingredient in controlling the ℓ_∞ -norm of trajectory in [14]. We use them for a similar purpose. For each $l \in [n]$, we define an operator $\mathcal{P}_\Omega^{(l)}$ such that $\mathcal{P}_\Omega^{(l)}(\mathbf{X})$ is equal to \mathbf{X} on the l th row and column, and equal to $\frac{1}{p} \mathcal{P}_\Omega(\mathbf{X})$ otherwise. The l th leave-one-out sequence, $\mathbf{x}^{(t,l)}$, evolves with

$$\mathbf{x}^{(t+1,l)} = \mathbf{x}^{(t,l)} - \eta \mathcal{P}_\Omega^{(l)}(\mathbf{x}^{(t,l)} \mathbf{x}^{(t,l)\top}) \mathbf{x}^{(t,l)} + \eta \mathbf{M}^{(l)} \mathbf{x}^{(t,l)}, \quad (3)$$

for $\mathbf{x}^{(0,l)} = \mathbf{x}^{(0)}$, where $\mathbf{M}^{(l)} = \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) + \mathbf{E}^{(l)}$, and $\mathbf{E}^{(l)}$ is obtained by zeroing out the l th row and column of $\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E})$.

3 Main Results

In this section, we present our main results. The first main result concerns the global convergence of GD with small random initialization.

Theorem 3.1. *Let us consider a rank-1 matrix completion problem that recovers the matrix $\mathbf{M}^* = \mathbf{x}^* \mathbf{x}^{*\top} \in \mathbb{R}^{n \times n}$ such that $\|\mathbf{x}^*\|_2 = \sqrt{\lambda^*}$ and $\|\mathbf{x}^*\|_\infty = \sqrt{\frac{\mu}{n}} \|\mathbf{x}^*\|_2$, where $\mu = O(\text{poly}(\log n))$. Let the initial point $\mathbf{x}^{(0)} \in \mathbb{R}^n$ be sampled from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \frac{1}{n} \beta_0^2 \mathbf{I})$ and $\mathbf{x}^{(t)}$ be updated with (1). Suppose that a small step size with $\eta \lambda^* < 0.1$ is used and the sample complexity satisfies $n^2 p \gtrsim \mu^5 n \log^{22} n$. Then, there exists $T^* = (1 + o(1)) \frac{1}{\eta \lambda^*} \log \frac{\sqrt{\lambda^* n}}{\beta_0}$ such that*

$$\left\| \mathbf{x}^{(t)} \pm \mathbf{x}^* \right\|_2 \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_2, \quad (4) \quad \max_{1 \leq l \leq n} \left\| \mathbf{x}^{(t)} - \mathbf{x}^{(t,l)} \right\|_2 \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_\infty, \quad (6)$$

$$\left\| \mathbf{x}^{(t)} \pm \mathbf{x}^* \right\|_\infty \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_\infty, \quad (5) \quad \max_{1 \leq l \leq n} \left| (\mathbf{x}^{(t,l)} - \mathbf{x}^*)_l \right| \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_\infty \quad (7)$$

hold at $t = T^*$ with probability at least $1 - o(1/\sqrt{\log n})$, if a sufficiently small initialization with

$$\sqrt{\lambda^* n^{-10}} \lesssim \beta_0 \lesssim \sqrt{\lambda^*} \sqrt[4]{\frac{np}{\mu^5 \log^{26} n}} \frac{1}{\sqrt[4]{n}} \quad (8)$$

is used and the noise satisfies $\sigma \lesssim \frac{\lambda^* \mu}{n} \sqrt{\log n}$.

Theorem 3.1 proves that, starting from a small random initialization, the trajectory of GD eventually enters the local region of the global minimizers $\pm \mathbf{x}^*$ in terms of both ℓ_2 and ℓ_∞ norms. Combined with the result of [14], GD starts to converge linearly to either \mathbf{x}^* or $-\mathbf{x}^*$ after $t = T^*$, as stated in the corollary below.

Corollary 3.2. *Suppose that the conditions in Theorem 3.1 are satisfied, and let ρ be a constant such that $1 - \frac{\eta}{10} \leq \rho < 1$. Then, with probability at least $1 - o(1/\sqrt{\log n})$, we have*

$$\left\| \mathbf{x}^{(t)} \pm \mathbf{x}^* \right\|_2 \lesssim \left(\frac{1}{\sqrt{\log n}} \rho^{t-T^*} + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}} \right) \|\mathbf{x}^*\|_2, \quad (9)$$

$$\left\| \mathbf{x}^{(t)} \pm \mathbf{x}^* \right\|_\infty \lesssim \left(\frac{1}{\sqrt{\log n}} \rho^{t-T^*} + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}} \right) \|\mathbf{x}^*\|_\infty, \quad (10)$$

for all $T^* \leq t \leq T = O(n^5)$.

The desired global convergence result is provided by Corollary 3.2. Several remarks about Theorem 3.1 and Corollary 3.2 are in order.

Matrix Recovery Suppose $\mathbf{x}^{(t)}$ converges to a global minimum \mathbf{y}^* of the function f , which is different from $\pm \mathbf{x}^*$. In such a case, despite achieving global convergence, the reconstructed matrix $\mathbf{y}^* \mathbf{y}^{*\top}$ deviates from the ground truth matrix \mathbf{M}^* . However, Theorem 3.1 establishes that $\mathbf{x}^{(t)}$ converges exclusively to the correct global minima $\pm \mathbf{x}^*$, so that the matrix \mathbf{M}^* is recovered with high probability.

Leave-one-out Sequence To apply the local convergence result of [14], in addition to (4) and (5), the existence of leave-one-out sequences $\{\mathbf{x}^{(t,l)}\}_{l \in [n]}$ satisfying (6) and (7) is required. Leave-one-out sequences also play a critical role and appear naturally in the proof of Theorem 3.1.

Sample Complexity The required sample complexity for Theorem 3.1 to hold is optimal up to a logarithmic factor compared to the statistical lower bound of $\Omega(n \log n)$. We have not done our best to optimize the log factors, and about half of them can be reduced with more delicate analysis. We will discuss this briefly in Section 6.

Convergence Time Considering that β_0^{-1} is at most polynomial in n (due to the lower bound of (8)), only $O(\log n)$ iterations are required for GD to enter the local region. It takes $O(\log(\frac{1}{\epsilon}))$ more iterations to achieve ϵ -accuracy in the local region, so the total iteration complexity is given by $O(\log n) + O(\log(\frac{1}{\epsilon}))$.

Initialization Size Although small initialization provides a good geometry to GD, a larger initialization is preferred because the convergence time, T^* , is inversely proportional to β_0 . When the sample complexity is optimal, i.e., $n^2 p \asymp n \text{ poly}(\log n)$, an upper bound on the initialization size given by Theorem 3.1 is $n^{-\frac{1}{4}}$, ignoring the log factors. However, as more samples are provided, we are allowed to use a larger initialization to reduce the convergence time. When the sample complexity satisfies $n^2 p \asymp n^{1+a}$, the bound is $n^{-\frac{1}{4}(1-a)}$ ignoring the log factors. The bound becomes nearly constant as a approaches 1, namely the fully observed case, and this is consistent with the previous result that small initialization is unnecessary for the fully observed case [18]. We also note that the lower bound of (8) is necessary in the proof of Theorem 3.1, since we derive probabilistic bounds for all iterations, and the lower bound limits the maximum number of iterations. However, we can further reduce the lower bound n^{-10} to n^{-c} for any constant $c > 10$ by tuning some constant factors during the proof.

Noise Size From the incoherence assumption, the maximum absolute value of entries of \mathbf{M}^* is bounded by $\frac{\lambda^* \mu}{n}$. The condition $\sigma \lesssim \frac{\lambda^* \mu}{n} \sqrt{\log n}$ in Theorem 3.1 allows the standard deviation of the Gaussian noise to be much larger than the maximum entry. It also implies $\frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}} \lesssim \mu \sqrt{\frac{\log n}{np}}$, so that the upper bounds in Corollary 3.2 are dominated by the first terms at $t = T^*$ and they eventually converge to the second terms as t increases.

Estimation Error The current estimation bounds (4) to (7) are all proportional to $\frac{1}{\sqrt{\log n}}$ times the norms of \mathbf{x}^* . However, if we do not allow the initialization size to grow with the sample complexity, we are able to obtain tighter bounds; if we use the fixed initialization size $n^{-\frac{1}{4}}$ regardless of the sample complexity, in Theorem 3.1, the factor $\frac{1}{\sqrt{\log n}}$ is improved to $\frac{1}{\sqrt{np}} + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}}$, and the upper bound on noise size is also improved to $\sigma \lesssim \frac{\lambda^* \mu}{n} \sqrt{np}$ (not being precise on the factors of μ and $\log n$ here). Then, the estimation error in Corollary 3.2 is improved to $\frac{1}{\sqrt{np}} \rho^t + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}}$ to match the result of [14] which uses spectral initialization. Thus, we have a tradeoff between estimation error and initialization size.

The next main result concerns the trajectory of GD before it enters the local region. The theorem states that for all $t \leq T^*$, $\mathbf{x}^{(t)}$ stays close to the fully observed case $\tilde{\mathbf{x}}^{(t)}$ in both ℓ_2 and ℓ_∞ -norm.

Theorem 3.3. *Suppose that the conditions of Theorem 3.1 hold, and T^* is defined as in Theorem 3.1. Then, for all $t \leq T^*$, we have*

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2, \quad (11) \quad \left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty \quad (12)$$

with probability at least $1 - o(1/\sqrt{\log n})$.

Trajectory of GD The sequence $\tilde{\mathbf{x}}^{(t)}$ is a linear combination of $\mathbf{x}^{(0)}$ and \mathbf{u}^* (see (C.1) in the appendix), and it is easy to analyze how $\tilde{\mathbf{x}}^{(t)}$ evolves. By showing that $\mathbf{x}^{(t)}$ stays close to $\tilde{\mathbf{x}}^{(t)}$ for all iterations, we not only show the convergence of GD with small initialization as in Theorem 3.1, but also characterize the exact trajectory that GD follows by Theorem 3.3.

Implicit Regularization One can prove that $\tilde{\mathbf{x}}^{(t)}$ is incoherent up to some log factors over all iterations, and from (11) and (12), the incoherence of $\mathbf{x}^{(t)}$ is bounded by that of $\tilde{\mathbf{x}}^{(t)}$. Thus, Theorem 3.3 shows that the incoherence of $\mathbf{x}^{(t)}$ is *implicitly* controlled by GD without any regularizer. This is an improvement over the previous result on the global convergence of GD for matrix completion [28], where an explicit regularizer was used to control the ℓ_∞ -norm of $\mathbf{x}^{(t)}$, although no small initialization was used in that work.

4 Fully Observed Case and Proof Sketch

Before we explain the proof of Theorems 3.1 and 3.3, we describe the trajectory of the fully observed case. We characterize $\tilde{\mathbf{x}}^{(t)}$ with three variables: $\tilde{\alpha}_t = |\mathbf{u}^{*\top} \tilde{\mathbf{x}}^{(t)}|$, $\tilde{\beta}_t = \|\tilde{\mathbf{x}}^{(t)}\|_2$, and $\tilde{\gamma}_t = \|\tilde{\mathbf{x}}_\perp^{(t)}\|_2$, where $\tilde{\mathbf{x}}_\perp^{(t)} = \tilde{\mathbf{x}}^{(t)} - \mathbf{u}^* \mathbf{u}^{*\top} \tilde{\mathbf{x}}^{(t)}$. According to (2), the three variables are updated with

$$\begin{aligned} \tilde{\alpha}_{t+1} &= (1 - \eta \tilde{\beta}_t^2 + \eta \lambda^*) \tilde{\alpha}_t; & \tilde{\gamma}_{t+1} &= (1 - \eta \tilde{\beta}_t^2) \tilde{\gamma}_t; \\ \tilde{\beta}_t^2 &= \tilde{\alpha}_t^2 + \tilde{\gamma}_t^2. \end{aligned}$$

At $t = 0$, due to random initialization, the initial vector is nearly orthogonal to \mathbf{u}^* , and we have $\tilde{\alpha}_0 \approx \frac{1}{\sqrt{n}} \beta_0$ and $\tilde{\gamma}_0 \approx \tilde{\beta}_0 = \beta_0$. Also, due to the small initialization, the term $\eta \tilde{\beta}_t^2$ is ignorable until $\tilde{\beta}_t$ becomes sufficiently large, so $\tilde{\alpha}_t$ grows exponentially at the rate of $1 + \eta \lambda^*$, while $\tilde{\gamma}_t$ remains still. Thus, in the early iterations where $(1 + \eta \lambda^*)^t$ is still much less than \sqrt{n} , $\tilde{\beta}_t$ is kept close to its initial

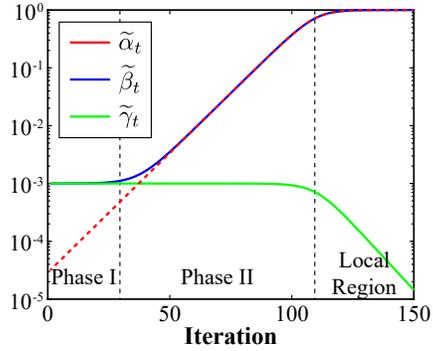


Figure 1: Evolution of the quantities $\tilde{\alpha}_t$, $\tilde{\beta}_t$, and $\tilde{\gamma}_t$ simulated with $\tilde{\alpha}_0 = \frac{1}{\sqrt{n}} \beta_0$, $\beta_0 = \frac{1}{n}$, $\lambda^* = 1$, and $n = 1000$.

value β_0 while the trajectory becomes more parallel to \mathbf{u}^* as $\tilde{\alpha}_t$ increases. When $(1 + \eta\lambda^*)^t$ becomes much larger than \sqrt{n} , the trajectory becomes almost parallel to \mathbf{u}^* in that $\tilde{\beta}_t \approx \tilde{\alpha}_t \gg \tilde{\gamma}_t$. Until $\tilde{\beta}_t$ (asymptotically) reaches $\frac{\sqrt{\lambda^*}}{\sqrt{\log n}}$, we can consider $\tilde{\alpha}_t$ as increasing at a rate of $(1 + \eta\lambda^*)$, and it takes about $\frac{1}{\log(1+\eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0}$ steps to reach this point. After that, we can no longer ignore the term $\eta\tilde{\beta}_t^2$, and $\tilde{\alpha}_t$ increases at a slower rate as $\tilde{\beta}_t$ increases. We can show that $\tilde{\beta}_t^2$ becomes sufficiently close to λ^* within $O(\log \log n)$ additional iterations, as stated in the following lemma.

Lemma 4.1. *Let T'_2 be the largest t such that $\tilde{\beta}_t^2 \leq \frac{\lambda^*}{64 \log n}$. At $t = T'_2 + \frac{6 \log \log n}{\log(1+\eta\lambda^*)}$, we have $\tilde{\beta}_t^2 \geq \lambda^* \left(1 - \frac{1}{\log n}\right)$.*

Finally, local convergence to \mathbf{u}^* occurs in that $\tilde{\alpha}_t$ approaches λ^* and $\tilde{\gamma}_t$ decreases exponentially with the rate $(1 - \eta\lambda^*)$. The actual behavior of quantities $\tilde{\alpha}_t, \tilde{\beta}_t, \tilde{\gamma}_t$ are plotted in Figure 1.

We define the iterates before $(1 + \eta\lambda^*)^t$ reaches $\frac{1}{\sqrt{np}}\sqrt{n}$, within some logarithmic factors, as Phase I, and the next iterates before $\tilde{\beta}_t^2$ reaches $\lambda^* \left(1 - \frac{1}{\log n}\right)$ as Phase II. Different techniques are used for each phase to prove that $\mathbf{x}^{(t)}$ stays close to $\tilde{\mathbf{x}}^{(t)}$. At the end of Phase I, $\tilde{\alpha}_t$ is increased to $\frac{1}{\sqrt{np}}\beta_0$ from its initial scale $\frac{1}{\sqrt{n}}\beta_0$, but it is still not dominant over β_0 . Therefore, the magnitudes of both $\mathbf{x}^{(t)}$ and $\tilde{\mathbf{x}}^{(t)}$ are kept close to β_0 throughout Phase I, and we take advantage of the small random initialization to show that the deviation of $\mathbf{x}^{(t)}$ from $\tilde{\mathbf{x}}^{(t)}$ does not increase much, and is kept at $\sqrt{\frac{1}{np}}$ times the norms of $\mathbf{x}^{(t)}$. In Phase II, we show that $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ expands at a rate of at most $(1 + \eta\lambda^*)$. Since the norms of $\mathbf{x}^{(t)}$ also grows at a rate of $(1 + \eta\lambda^*)$ during most of Phase II, the norms of $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ remain negligible compared to those of $\mathbf{x}^{(t)}$. The next two sections give the main lemmas of Phase I and II, respectively, which are used to prove Theorems 3.1 and 3.3. For a visual representation of the results in the following two sections, please refer to Figure 2.

5 Phase I: Finding Direction

We provide detailed results and proof ideas for Phase I. Our main goal is to analyze the deviation of $\mathbf{x}^{(t)}$ from $\tilde{\mathbf{x}}^{(t)}$. First, if we look at the update equations (1) and (2), the second term is proportional to the third power of $\|\mathbf{x}^{(t)}\|_2$, while the other terms depend linearly on $\|\mathbf{x}^{(t)}\|_2$. Thus, the second term is almost negligible due to the small initialization. Without the second terms, the difference between $\mathbf{x}^{(t)}$ and $\tilde{\mathbf{x}}^{(t)}$ at $t = 1$ is $\eta(\mathbf{M}^\circ - \mathbf{M}^*)\mathbf{x}^{(0)}$. From concentration inequalities, one can see that the ℓ_2 and ℓ_∞ norms of $\eta(\mathbf{M}^\circ - \mathbf{M}^*)\mathbf{x}^{(0)}$ are about $\frac{1}{\sqrt{np}}$ times smaller than those of $\tilde{\mathbf{x}}^{(1)}$.

Due to the third terms of (1) and (2), the norms of $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ can grow exponentially at a rate of $(1 + \eta\lambda^*)$ in the worst case where $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ is parallel to \mathbf{u}^* . In such a case, the norms of $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ would be larger than those of $\tilde{\mathbf{x}}^{(t)}$ at the end of Phase I, since those of $\tilde{\mathbf{x}}^{(t)}$ remain still in Phase I. However, we overcome this problem by proving that the bounds grow at most *polynomially* with respect to t , and since t is at most $O(\log n)$, the bounds remain $\frac{1}{\sqrt{np}}$ times smaller than the norms of $\tilde{\mathbf{x}}^{(t)}$ up to logarithmic factors throughout Phase I.

Lemma 5.1. *Let T_1 be the largest t such that $(1 + \eta\lambda^*)^t \leq \sqrt{\frac{\mu^4 \log^{21} n}{np}} \sqrt{n}$. Under the conditions of Theorem 3.1, with probability at least $1 - o(1/\sqrt{\log n})$, for all $t \leq T_1$, we have*

$$\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2 \lesssim \mu \sqrt{\frac{\log n}{np}} \beta_0 t, \quad (13) \quad \|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_\infty \lesssim \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t^2. \quad (14)$$

T_1 is defined to be the end of Phase I. Lemma 5.1 proves Theorem 3.3 for Phase I.

Proof of (13) We will first demonstrate how to obtain the ℓ_2 -norm bound of Lemma 5.1. Let us define a sequence $\hat{\mathbf{x}}^{(t)}$ that is updated as

$$\hat{\mathbf{x}}^{(t+1)} = \hat{\mathbf{x}}^{(t)} - \eta \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2^2 \hat{\mathbf{x}}^{(t)} + \eta \mathbf{M}^\circ \hat{\mathbf{x}}^{(t)}; \quad \hat{\mathbf{x}}^{(0)} = \mathbf{x}^{(0)}. \quad (15)$$

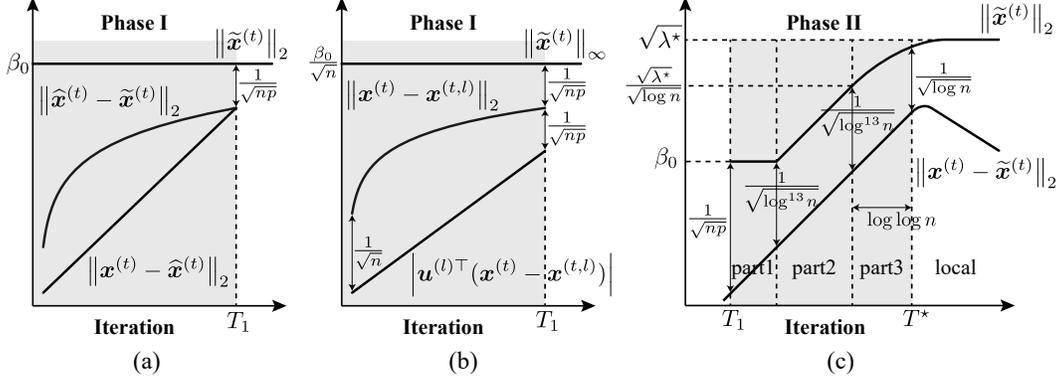


Figure 2: An illustrative description of trajectory of various quantities compared to the norms of $\mathbf{x}^{(t)}$ on logarithmic scales. Arrows between lines represent the ratio between them. The quantities depicted are not precise, and only the key factors are shown for simplicity. (a) In Phase I, $\|\widehat{\mathbf{x}}^{(t)} - \widetilde{\mathbf{x}}^{(t)}\|_2$ increases linearly and $\|\mathbf{x}^{(t)} - \widehat{\mathbf{x}}^{(t)}\|_2$ increases exponentially with the rate of $(1 + \eta\lambda^*)$. They have the same scale at the end of Phase I. (b) In Phase I, even if the \mathbf{u}^* component of $\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}$ grows exponentially, it remains almost orthogonal to \mathbf{u}^* throughout the phase. (c) Phase II is divided into three parts according to the growth speed of $\mathbf{x}^{(t)}$ and $\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)}$. The ratio between them at the start and the end of each part is described, and it is at most $\frac{1}{\sqrt{\log n}}$ in Phase II.

Note that the norm of $\widetilde{\mathbf{x}}^{(t)}$ is used in the second term of (15). The update equation of $\widehat{\mathbf{x}}^{(t)}$ differs from $\widetilde{\mathbf{x}}^{(t)}$ in the third term and from $\mathbf{x}^{(t)}$ in the second term. We use $\widehat{\mathbf{x}}^{(t)}$ as a proxy for bounding $\|\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)}\|_2$. We first show that $\|\widehat{\mathbf{x}}^{(t)} - \widetilde{\mathbf{x}}^{(t)}\|_2$ grows at most linearly with respect to t .

Lemma 5.2. *With probability at least $1 - o(1/\sqrt{\log n})$, for all $t \leq T_1$, we have*

$$\|\widehat{\mathbf{x}}^{(t)} - \widetilde{\mathbf{x}}^{(t)}\|_2 \lesssim \mu \sqrt{\frac{\log n}{np}} \beta_0 t. \quad (16)$$

The proof of this lemma is based on the fact that $\widehat{\mathbf{x}}^{(t)}$ is a product of $\mathbf{x}^{(0)}$ and a matrix polynomial of \mathbf{I} and \mathbf{M}° , while $\widetilde{\mathbf{x}}^{(t)}$ is a product between $\mathbf{x}^{(0)}$ and a matrix polynomial of \mathbf{I} and \mathbf{M}^* . We prove the lemma by comparing the two matrix polynomials. We remark that Lemma 5.2 holds regardless of the small initialization, but it relies on the randomness of $\mathbf{x}^{(0)}$.

Since $\mathbf{x}^{(t)}$ and $\widehat{\mathbf{x}}^{(t)}$ differ only in the second term, their initial difference is proportional to β_0^3 . More precisely, it is $\frac{1}{\sqrt{np}}\beta_0^3$. We show that the difference grows exponentially at a rate of $(1 + \eta\lambda^*)$.

Lemma 5.3. *If (14) holds for all $t \leq T_1$, we have*

$$\|\mathbf{x}^{(t)} - \widehat{\mathbf{x}}^{(t)}\|_2 \lesssim \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta\lambda^*)^t \beta_0^3 \quad (17)$$

for all $t \leq T_1$ with probability at least $1 - o(1/\sqrt{\log n})$.

The upper bound in (17) becomes smaller than that of (16) if $(1 + \eta\lambda^*)^t \beta_0^2 \leq \lambda^* \sqrt{\frac{1}{\mu \log^2 n}}$. One can check that this condition is satisfied from the definition of T_1 given in Lemma 5.1 and the bound on the initialization size (8). Thus, (13) is proved by (16) and (17).

Proof of (14) We control the l th component of $\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)}$ using the l th leave-one-out sequence. Leave-one-out sequences have two important properties. First, because they are defined without only one row/column, they are extremely close to $\mathbf{x}^{(t)}$, and at $t = 1$, $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$ is about $\frac{1}{\sqrt{np}} \frac{\beta_0}{\sqrt{n}}$. Second, the l th component of the l th leave-one-out sequence evolves similarly to that of $\widetilde{\mathbf{x}}^{(t)}$ and is easy to analyze. With these two properties, we bound the l th component of $\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)}$ as

$$\left| \left(\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)} \right)_l \right| \leq \left\| \mathbf{x}^{(t)} - \mathbf{x}^{(t,l)} \right\|_2 + \left| \left(\mathbf{x}^{(t,l)} - \widetilde{\mathbf{x}}^{(t)} \right)_l \right|. \quad (18)$$

We claim that both $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$ and $|(\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l|$ increase at most polynomially with respect to t from the initial scale $\frac{1}{\sqrt{np}} \frac{\beta_0}{\sqrt{n}}$.

Lemma 5.4. *With probability at least $1 - o(1/\sqrt{\log n})$, for all $t \leq T_1$, we have*

$$\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2 \lesssim \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t, \quad (19) \quad |(\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l| \lesssim \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t^2. \quad (20)$$

As explained for $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$, due to the third terms of (1) and (3), $\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}$ can also grow exponentially at the rate of $(1 + \eta\lambda^*)$ in the worst case where $\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}$ is parallel to \mathbf{u}^* . This contradicts our result (19) that $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$ grows only linearly. We show that $\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}$ remains nearly orthogonal to \mathbf{u}^* in Phase I, and thus the worst case does not occur.

Lemma 5.5. *For all $l \in [n]$ and $t \leq T_1$, we have*

$$|\mathbf{u}^{(l)\top}(\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)})| \lesssim \sqrt{\frac{\mu^3 \log^2 n}{np}} (1 + \eta\lambda^*)^t \frac{\beta_0}{n}$$

with probability at least $1 - o(1/\sqrt{\log n})$, where $\mathbf{u}^{(l)}$ is the first eigenvector of $\mathbf{M}^{(l)}$.

Note that $\mathbf{u}^{(l)}$ is almost parallel to \mathbf{u}^* (see Lemma A.5 in the appendix). The $\mathbf{u}^{(l)}$ component of $\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}$ is initialized to the order of $\frac{1}{\sqrt{np}} \frac{\beta_0}{n}$, which is $\frac{1}{\sqrt{n}}$ times smaller than $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$.

Although it is increased exponentially, from the definition of T_1 , the $\mathbf{u}^{(l)}$ component remains much smaller than $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$ in Phase I.

One can see that $|(\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l|$ increases by $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2 \|\mathbf{u}^*\|_\infty$ at each step, and summing the bound (13) up to t gives (20). Finally, (14) is obtained by putting (19) and (20) into (18).

6 Phase II: Expansion

In the next phase, we show that the bounds obtained in Phase I are increased at a rate of $(1 + \eta\lambda^*)$.

Lemma 6.1. *Let T_2 be the largest t such that $\tilde{\beta}_t^2 \leq \lambda^* \left(1 - \frac{1}{\log n}\right)$. Then, for all $T_1 < t \leq T_2$, we have*

$$\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2 \lesssim \mu \sqrt{\frac{\log^3 n}{np}} \beta_0 (1 + \eta\lambda^*)^{t-T_1}, \quad (21)$$

$$\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2 \lesssim \mu \sqrt{\frac{\log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}, \quad (22)$$

$$\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_\infty \lesssim \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}, \quad (23)$$

$$|(\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l| \lesssim \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}, \quad (24)$$

with probability at least $1 - o(1/\sqrt{\log n})$.

T_2 is defined as the end of Phase II. We will explain how Lemma 6.1 leads to Theorem 3.3 in Phase II. Let us first focus on (21) and (11). We can divide Phase II into three parts according to the behavior of $\|\tilde{\mathbf{x}}^{(t)}\|_2$. First, $\|\tilde{\mathbf{x}}^{(t)}\|_2$ is kept close to β_0 until $(1 + \eta\lambda^*)^t$ becomes \sqrt{n} , or $(1 + \eta\lambda^*)^{t-T_1}$ becomes \sqrt{np} . In this part, although the bounds increase exponentially with the rate of $(1 + \eta\lambda^*)$, the factor $\frac{1}{\sqrt{np}}$, which was already present in (13) of Phase I, compensates for this increase. At the end of the first part, $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2$ is smaller than $\|\tilde{\mathbf{x}}^{(t)}\|_2$ by some log factors. Next, $\|\tilde{\mathbf{x}}^{(t)}\|_2$ grows at the rate of $(1 + \eta\lambda^*)$ until it reaches $\frac{\sqrt{\lambda^*}}{8\sqrt{\log n}}$. Since both $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2$ and $\|\tilde{\mathbf{x}}^{(t)}\|_2$ increase with $(1 + \eta\lambda^*)$,

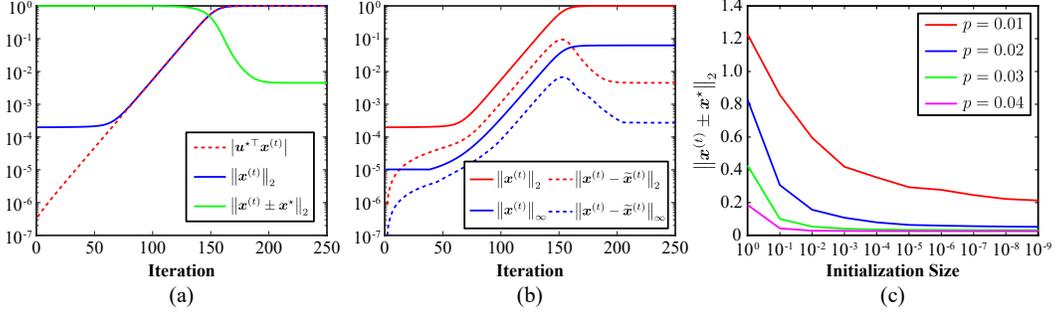


Figure 3: (a) Evolution of the quantities $|\mathbf{u}^{\star\top} \mathbf{x}^{(t)}|$ and $\|\mathbf{x}^{(t)}\|_2$, which behave similarly to $\tilde{\alpha}_t$ and $\tilde{\beta}_t$, respectively, and $\|\mathbf{x}^{(t)} \pm \mathbf{x}^*\|_2$, which shows local convergence. (b) Comparison between the norms of $\mathbf{x}^{(t)}$ and $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$. (c) Convergence of GD with respect to the initialization size and sampling probability. $\|\mathbf{x}^{(t)} \pm \mathbf{x}^*\|_2$ was measured at $t = \frac{1}{\log(1+\eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0} + 100$ and averaged over 1000 trials.

the ratio between them is maintained in the second part. Finally, in the remaining iterations, $\|\tilde{\mathbf{x}}^{(t)}\|_2$ increases with $(1 - \eta\tilde{\beta}_t^2 + \eta\lambda^*)$ at each step, and the increment becomes smaller as it converges to $\sqrt{\lambda^*}$. Thus, as in the first part, $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2$ increases faster than $\|\tilde{\mathbf{x}}^{(t)}\|_2$. However, from Lemma 4.1, the length of this part is $O(\log \log n)$, and the ratio between $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2$ and $\|\tilde{\mathbf{x}}^{(t)}\|_2$ increases only by $\log^6 n$. We prove that the log factors already present at the end of the second part compensate this, and finally (11) holds for all t in Phase II. A more delicate analysis may prove that $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2$ grows at the same rate as $\|\tilde{\mathbf{x}}^{(t)}\|_2$ in the third part, and this will reduce the required sample complexity by at most $\log^{12} n$. A similar argument can be used to prove that the bounds for $\|\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_\infty$, $\|\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}\|_2$, and $|\langle \mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)}, \mathbf{u} \rangle|$ are smaller than $\|\tilde{\mathbf{x}}^{(t)}\|_\infty$ by some log factors throughout Phase II.

At the end of Phase II, $\tilde{\mathbf{x}}^{(t)}$ is very close to $\pm \mathbf{x}^*$ in both ℓ_2 and ℓ_∞ norms (see Corollary C.3 in the appendix), so one can replace $\tilde{\mathbf{x}}^{(t)}$ of Lemma 6.1 with $\pm \mathbf{x}^*$ to prove (4) to (7) of Theorem 3.1. Hence, we can let $T^* = T_2$, and as explained in Section 4, T_2 is approximately given by $\frac{1}{\log(1+\eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0} + O(\log \log n)$.

7 Simulation

In this section, we present some simulation results that support our theoretical findings.

Trajectory of GD With the dimension $n = 5000$, we constructed the ground truth vector \mathbf{u}^* by sampling it from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \frac{1}{n} \mathbf{I})$ and normalizing it to have unit norm. We let $\lambda^* = 1$ so that the matrix \mathbf{M}^* is given by $\mathbf{u}^* \mathbf{u}^{\star\top}$, and we randomly sampled the matrix symmetrically with a sampling rate of $p = 0.1$ and Gaussian noise of $\sigma = \frac{0.1}{n}$. The initialization size was set to $\beta_0 = \frac{1}{n}$ and a step size of 0.1 was used for GD. Figure 3 (a) and (b) represent one trial of the experiment, but similar graphs were obtained in each repetition of the experiment. The evolution of some important quantities such as $\|\mathbf{x}^{(t)}\|_2$ and $|\mathbf{u}^{\star\top} \mathbf{x}^{(t)}|$ is shown in Figure 3(a). As in the fully observed case, the signal component $|\mathbf{u}^{\star\top} \mathbf{x}^{(t)}|$ increases at the rate of $(1 + \eta\lambda^*)$ until it approaches $\sqrt{\lambda^*}$, and a local convergence to \mathbf{x}^* occurs, where $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2$ decreases exponentially and saturates at the level determined by the noise size σ . In Figure 3(b), we describe the deviation of $\mathbf{x}^{(t)}$ from $\tilde{\mathbf{x}}^{(t)}$ in both ℓ_2 and ℓ_∞ norms. The solid lines represent the norms of $\mathbf{x}^{(t)}$ and the dotted lines represent those of $\mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)}$. We can see that there is a gap between the solid and the dotted lines during the whole iterations. Thus, $\mathbf{x}^{(t)}$ stays close to the trajectory of the fully observed case, as we proved in Theorem 3.3.

Small Initialization In the next experiment, we investigated the importance of a small initialization for the convergence of GD. We used the same conditions as in the previous experiment except

$n = 500$. We measured $\|\mathbf{x}^{(t)} \pm \mathbf{x}^*\|_2$ at $t = \frac{1}{\log(1+\eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0} + 100$ and averaged it over 1000 trials. We repeated the experiment while changing the initialization size from 10^0 to 10^{-9} and the sampling probability from 0.01 to 0.04. The result is summarized in Figure 3(c). For all sampling probabilities, the small initialization improves the convergence of GD. Also, the performance starts to saturate at much larger initialization sizes as the sampling probability increases, and this is consistent with our finding (8) that a larger initialization is possible as more samples are available.

8 Discussion

In this paper, we showed that for rank-1 symmetric matrix completion with ℓ_2 loss, GD can converge to the ground truth starting from a small random initialization. Ignoring log factors, the bound on the initialization size is $n^{-\frac{1}{4}}$ when the optimal $n \text{ poly}(\log n)$ samples are provided, and the bound becomes larger as more samples are provided. The result is interesting because the loss function does not have global benign geometry if no regularizer is applied. Our result does not use any explicit regularizer and relies only on the implicit regularizing effect of GD.

The most important future work is an extension to the rank- r case. Suppose that M^* is a rank- r matrix and its eigendecomposition is given by $U^* \Sigma^* U^{*\top} = X^* X^{*\top}$, where $\Sigma^* = \text{diag}(\lambda_1^*, \dots, \lambda_r^*)$ and $X^* = U^* \Sigma^{*\frac{1}{2}}$. Then, the trajectory of GD becomes an $n \times r$ matrix $X^{(t)}$, which is updated as

$$X^{(t+1)} = X^{(t)} - \frac{\eta}{p} \mathcal{P}_\Omega \left(X^{(t)} X^{(t)\top} \right) X^{(t)} + \eta M^\circ X^{(t)}.$$

Each entry of $X^{(0)}$ is sampled independently from the Gaussian distribution $\mathcal{N}\left(0, \frac{1}{n} \beta_0^2\right)$ as in the rank-1 case.

An instance of $X^{(t)}$ is shown in Figure 4. The same conditions as in Figure 3 are used, except that the ground truth matrix is a rank-3 matrix with non-zero eigenvalues 1, 0.75, 0.5. The singular values of $X^{(t)}$ behave similarly to $\|\mathbf{x}^{(t)}\|_2$ in the rank-1 case. In the early iterations, where the orthogonal components dominate, the singular values stay close to their initial scale β_0 . After that, each singular value $\sigma_i(X^{(t)})$ increases at a rate of $(1 + \eta\lambda_i^*)$ and saturates at $\sqrt{\lambda_i^*}$. We use $\|X - Y\|_R$ to denote the Frobenius norm between X and Y under best rotational alignment. $\|X^{(t)} - X^*\|_R$ decreases exponentially and saturates at the level determined by the noise size σ , after all singular values have saturated, as local convergence begins.

To extend the results of the rank-1 case, we need to show that $\|X^{(t)} - \widetilde{X}^{(t)}\|_R$ remains much smaller than $\|X^{(t)}\|_F$ throughout the iterations, where $\widetilde{X}^{(t)}$ is the trajectory of the fully observed case. Before $\sigma_1(X^{(t)})$ saturates around $\sqrt{\lambda_1^*}$, it behaves similarly to $\|\mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)}\|_2$ of the rank-1 case, i.e., it expands at a rate of $(1 + \eta\lambda_1^*)$ along with $\|X^{(t)}\|_F$ after the early iterations. However, because each singular value grows at a different rate, a different phenomenon is observed for the rank- r case. During the iterations before $\sigma_{i+1}(X^{(t)})$ saturates after $\sigma_i(X^{(t)})$ does, both $\|X^{(t)} - \widetilde{X}^{(t)}\|_R$ and $\|X^{(t)}\|_F$ do not increase much. Our current theory can only show that $\|X^{(t)} - \widetilde{X}^{(t)}\|_R$ increases at a rate less than $(1 + \eta\lambda_{i+1}^*)$, and in order for $\|X^{(t)} - \widetilde{X}^{(t)}\|_R$ to remain much smaller than $\|X^{(t)}\|_F$, additional sample complexity is required to compensate for the exponential increases. Therefore, we expect that the convergence of GD for the case of rank- r can be proved with the techniques developed in this paper if $n^{1+\Theta(\kappa-1)} \text{poly}(\kappa, r, \log n)$ samples are provided, where $\kappa = \frac{\lambda_1^*}{\lambda_r^*}$ is the condition number. Nevertheless, whether GD can converge with the optimal $n \text{ poly}(\kappa, r, \log n)$ samples for the rank- r matrix completion problem remains an open problem.

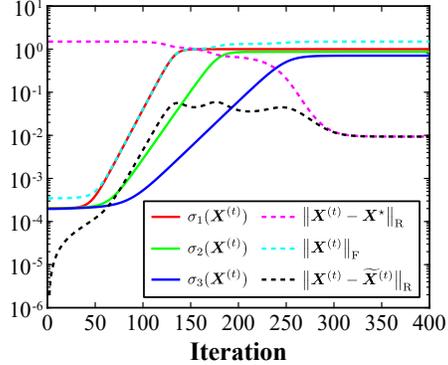


Figure 4: Trajectory of GD obtained for a rank-3 matrix with non-zero eigenvalues 1, 0.75, 0.5. The same parameters were used as in Figure 3.

Acknowledgments and Disclosure of Funding

This research was supported by the National Research Foundation of Korea under grant 2021R1C1C11008539.

References

- [1] Y. Chen, Y. Chi, J. Fan, and C. Ma, “Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval,” *Mathematical Programming*, vol. 176, no. 1-2, pp. 5–37, 2019.
- [2] A. Ahmed, B. Recht, and J. Romberg, “Blind deconvolution using convex programming,” *IEEE Transactions on Information Theory*, vol. 60, no. 3, pp. 1711–1732, 2014.
- [3] 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, 2016.
- [4] E. J. Candès and B. Recht, “Exact matrix completion via convex optimization,” *Foundations of Computational Mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [5] D. F. Gleich and L.-h. Lim, “Rank aggregation via nuclear norm minimization,” in *Proceedings of the 17th ACM international conference on Knowledge discovery and data mining*, pp. 60–68, ACM, 2011.
- [6] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?,” *Journal of the ACM*, vol. 58, no. 3, pp. 11:1–11:37, 2011.
- [7] Y. Hu, X. Liu, and M. Jacob, “A generalized structured low-rank matrix completion algorithm for mr image recovery,” *IEEE Transactions on Medical Imaging*, vol. 38, no. 8, pp. 1841–1851, 2019.
- [8] Y. Chi, Y. M. Lu, and Y. Chen, “Nonconvex optimization meets low-rank matrix factorization: An overview,” *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5239–5269, 2019.
- [9] E. J. Candès and T. Tao, “The power of convex relaxation: Near-optimal matrix completion,” *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [10] R. H. Keshavan, A. Montanari, and S. Oh, “Matrix completion from a few entries,” *IEEE Transactions on Information Theory*, vol. 56, no. 6, pp. 2980–2998, 2010.
- [11] 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, ACM, 2013.
- [12] 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.
- [13] 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.
- [14] 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,” *Foundations of Computational Mathematics*, vol. 20, no. 3, pp. 451–632, 2020.
- [15] J. Chen, D. Liu, and X. Li, “Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization,” *IEEE Transactions on Information Theory*, vol. 66, no. 9, pp. 5806–5841, 2020.
- [16] S. S. Du, C. Jin, J. D. Lee, M. I. Jordan, A. Singh, and B. Póczos, “Gradient descent can take exponential time to escape saddle points,” in *Advances in Neural Information Processing Systems*, 2017.
- [17] Y. Li, T. Ma, and H. Zhang, “Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations,” in *Conference on Learning Theory*, pp. 2–47, 2018.
- [18] T. Ye and S. S. Du, “Global convergence of gradient descent for asymmetric low-rank matrix factorization,” in *Advances in Neural Information Processing Systems*, pp. 1429–1439, 2021.

- [19] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai, “Gradient descent finds global minima of deep neural networks,” in *International Conference on Machine Learning*, pp. 1675–1685, 2019.
- [20] Y. Ma, A. Olshevsky, C. Szepesvari, and V. Saligrama, “Gradient descent for sparse rank-one matrix completion for crowd-sourced aggregation of sparsely interacting workers,” in *International Conference on Machine Learning*, pp. 3335–3344, 2018.
- [21] Q. Ma and A. Olshevsky, “Adversarial crowdsourcing through robust rank-one matrix completion,” in *Advances in Neural Information Processing Systems*, pp. 21841–21852, 2020.
- [22] S. Gunasekar, B. E. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Implicit regularization in matrix factorization,” in *Advances in Neural Information Processing Systems*, 2017.
- [23] D. Stöger and M. Soltanolkotabi, “Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction,” in *Advances in Neural Information Processing Systems*, pp. 23831–23843, 2021.
- [24] S. Arora, N. Cohen, W. Hu, and Y. Luo, “Implicit regularization in deep matrix factorization,” in *Advances in Neural Information Processing Systems*, pp. 7413–7424, 2019.
- [25] N. Razin and N. Cohen, “Implicit regularization in deep learning may not be explainable by norms,” in *Advances in Neural Information Processing Systems*, pp. 21174–21187, 2020.
- [26] Z. Li, Y. Luo, and K. L. Lyu, “Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning,” in *In International Conference on Learning Representations*, 2021.
- [27] S. Bhojanapalli, B. Neyshabur, and N. Srebro, “Global optimality of local search for low rank matrix recovery,” in *Advances in Neural Information Processing Systems*, 2016.
- [28] R. Ge, J. D. Lee, and T. Ma, “Matrix completion has no spurious local minimum,” in *Advances in Neural Information Processing Systems*, 2016.

Detailed proofs for the results explained in the main text are provided in this appendix. We say that an event happens *with high probability* if it happens with probability at least $1 - \frac{1}{n^C}$ for a constant $C > 0$ and C can be made arbitrary large by controlling constant factors. A union of $\text{poly}(n)$ number of events that happens with high probability still happens with high probability. For a matrix \mathbf{A} , we denote the spectral norm by $\|\mathbf{A}\|$ and the maximum absolute value of entries by $\|\mathbf{A}\|_\infty$. Also, the largest ℓ_2 -norm of rows of \mathbf{A} is denoted as $\|\mathbf{A}\|_{2,\infty}$.

A Spectral Analysis

We introduce some spectral bounds related to random sampling and Gaussian noise.

Lemma A.1. *If $n^2 p \gtrsim n \log n$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathbf{M}^* \right\| \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}}$$

with high probability.

Lemma A.2. *If $n^2 p \gtrsim \mu n \log n$, for all $l \in [n]$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right\| \lesssim \lambda^* \sqrt{\frac{\mu}{np}}$$

with high probability.

Lemma A.3. *If $n^2 p \gtrsim n \log^2 n$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \right\| \lesssim \sigma \sqrt{\frac{n}{p}}$$

with high probability.

Note that Lemma A.3 also implies that $\|\mathbf{E}^{(l)}\| \lesssim \sigma \sqrt{\frac{n}{p}}$ for all $l \in [n]$ with high probability.

Combined with the condition $\sigma \lesssim \frac{\lambda^* \mu}{n} \sqrt{\log n}$, we have

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \right\| \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}}, \quad \|\mathbf{E}^{(l)}\| \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}}$$

for all $l \in [n]$. Proofs for Lemmas A.1 and A.2 are provided in Appendix G. Check Lemma 11 of [12] for the proof of Lemma A.3.

Next, we state bounds on the eigenvalues of \mathbf{M}° and $\mathbf{M}^{(l)}$. The first eigenvalues of \mathbf{M}° and $\mathbf{M}^{(l)}$ are denoted as λ° and $\lambda^{(l)}$, respectively. The following lemma is derived from Lemmas A.1 and A.2 with Weyl's Theorem.

Lemma A.4. *If $n^2 p \gtrsim \mu n \log n$, we have*

$$|\lambda^\circ - \lambda^*| \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}}, \tag{A.1}$$

$$|\lambda^{(l)} - \lambda^*| \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \tag{A.2}$$

for all $l \in [n]$ with high probability.

Lastly, Lemmas A.1 and A.2 with Davis-Kahan Theorem give the following lemma.

Lemma A.5. *If $n^2 p \gtrsim \mu n \log n$, we have*

$$\|\mathbf{u}^{(l)} - \mathbf{u}^*\|_2 \lesssim \mu \sqrt{\frac{\log n}{np}}$$

for all $l \in [n]$ with high probability.

B Initialization

In this section, we introduce some properties that the initialization vector $\mathbf{x}^{(0)}$ satisfies. Recall that each entry of $\mathbf{x}^{(0)}$ is sampled from $\mathcal{N}(0, \frac{1}{n}\beta_0^2)$ independently. We use \mathbf{H} to denote the perturbation $\mathbf{M}^\circ - \mathbf{M}^*$.

Lemma B.1. *The initialization vector $\mathbf{x}^{(0)}$ satisfies*

$$\frac{1}{2}\beta_0 \leq \|\mathbf{x}^{(0)}\|_2 \leq \frac{3}{2}\beta_0 \quad (\text{B.1})$$

with probability at least $1 - e^{-n/32}$, and

$$\|\mathbf{x}^{(0)}\|_\infty \leq 2\sqrt{\log n} \frac{\beta_0}{\sqrt{n}}, \quad (\text{B.2})$$

$$\left| \mathbf{u}^{\star\top} \mathbf{H}^s \mathbf{x}^{(0)} \right| \leq 2\sqrt{\log n} \frac{\beta_0}{\sqrt{n}} \|\mathbf{H}\|^s, \quad \forall s \leq 30 \log n, \quad (\text{B.3})$$

with probability at least $1 - \frac{1}{n} - \frac{30 \log n}{n^2}$. It also satisfies

$$\frac{1}{\sqrt{\log n}} \frac{\beta_0}{\sqrt{n}} \leq \left| \mathbf{u}^{\star\top} \mathbf{x}^{(0)} \right| \quad (\text{B.4})$$

with probability at least $1 - \frac{1}{2\sqrt{\log n}}$.

Proof. To bound $\|\mathbf{x}^{(0)}\|_2$, we use the following basic concentration inequality that holds for i.i.d. standard normal variables $\{X_i\}_{i \in [n]}$.

$$\mathbb{P} \left[\left| \frac{1}{n} \sum_{i=1}^n X_i^2 - 1 \right| \geq t \right] \leq 2e^{-nt^2/8}$$

If we put $t = \frac{1}{2}$, with probability at least $1 - e^{-n/32}$, we have

$$\frac{1}{2}\beta_0^2 \leq \|\mathbf{x}^{(0)}\|_2^2 \leq \frac{3}{2}\beta_0^2,$$

and this implies (B.1).

For a centered Gaussian random variable with standard deviation σ , we have

$$\mathbb{P}[|X| \geq t] \leq e^{-\frac{t^2}{2\sigma^2}}.$$

Hence, an entry of $\mathbf{x}^{(0)}$ is less than $2\sqrt{\log n} \frac{\beta_0}{\sqrt{n}}$ with probability at least $1 - \frac{1}{n^2}$, and all entries of $\mathbf{x}^{(0)}$ are less than $2\sqrt{\log n} \frac{\beta_0}{\sqrt{n}}$ with probability at least $1 - \frac{1}{n}$. $\mathbf{u}^{\star\top} \mathbf{H}^s \mathbf{x}^{(0)}$ follows a centered Gaussian distribution with standard deviation $\|\mathbf{H}^s \mathbf{u}^*\|_2 \leq \|\mathbf{H}\|^s \|\mathbf{u}^*\|_2$ for all s , and (B.3) holds with probability at least $1 - \frac{10 \log n}{n^2}$.

For a random variable X that is sampled from $\mathcal{N}(0, \sigma^2)$, we have

$$\mathbb{P}[|X| \leq t] \leq \frac{t}{\sqrt{2\pi}\sigma^2}.$$

Hence, we have

$$\frac{1}{\sqrt{\log n}} \frac{\beta_0}{\sqrt{n}} \leq \left| \mathbf{u}^{\star\top} \mathbf{x}^{(0)} \right|$$

with probability at least $1 - \frac{1}{2\sqrt{\log n}}$. □

Lemma B.1 implies that the \mathbf{u}^* component of $\mathbf{x}^{(0)}$ is in the range

$$\frac{1}{\sqrt{\log n}} \frac{\beta_0}{\sqrt{n}} \leq \left| \mathbf{u}^{\star\top} \mathbf{x}^{(0)} \right| \leq 2\sqrt{\log n} \frac{\beta_0}{\sqrt{n}}. \quad (\text{B.5})$$

Lemma B.2. *We have*

$$\left\| a\mathbf{x}^{(0)} + b\mathbf{u}^* \right\|_\infty \geq (|a|\beta_0 + |b|) \frac{1}{\sqrt{n}} \quad (\text{B.6})$$

for all a, b , with probability at least $1 - \exp\left(-\frac{n}{2\mu}\right)$.

Proof. The probability that an entry of $\mathbf{x}^{(0)}$ is less than $\frac{\beta_0}{\sqrt{n}}$ is bounded by $\frac{1}{\sqrt{2\pi}}$. Without loss of generality, let us assume that all entries of \mathbf{u}^* are not negative and $a, b \geq 0$. There are at least $\frac{n}{\mu}$ entries of \mathbf{u}^* that are larger than $\frac{1}{\sqrt{n}}$. For such entries, the probability that all entries of $\mathbf{x}^{(0)}$ is less than $\frac{\beta_0}{\sqrt{n}}$ is bounded by $\left(\frac{1}{\sqrt{2\pi}}\right)^{\frac{n}{\mu}} \leq \exp\left(-\frac{n}{2\mu}\right)$. Hence, for at least one position, both entries of $\mathbf{x}^{(0)}$ and \mathbf{u}^* are larger than $\frac{\beta_0}{\sqrt{n}}$ and $\frac{1}{\sqrt{n}}$, respectively, with probability at least $1 - \exp\left(-\frac{n}{2\mu}\right)$. \square

In the following sections, we assume that we are given an initialization vector $\mathbf{x}^{(0)}$ that satisfies (B.1) to (B.6).

C Fully Observed Case

We provide some lemmas related to $\tilde{\mathbf{x}}^{(t)}$ in this section. We first note that $\tilde{\mathbf{x}}^{(t)}$ is explicitly written as

$$\tilde{\mathbf{x}}^{(t)} = \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2) \mathbf{x}^{(0)} + \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2 + \eta\lambda^*) (\mathbf{u}^{*\top} \mathbf{x}^{(0)}) \mathbf{u}^* := A^{(t)} \mathbf{x}^{(0)} + B^{(t)} \mathbf{u}^*. \quad (\text{C.1})$$

Let us define T'_2 as the last t such that $\tilde{\beta}_t^2 \leq \frac{\lambda^*}{64 \log n}$. We claim that $T'_2 \leq \frac{64 \log n}{\eta\lambda^*}$ and prove this later. Then, for all $t \leq T'_2$, we have

$$\begin{aligned} \frac{1}{4} (1 + \eta\lambda^*)^t &\leq \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2 + \eta\lambda^*) \leq (1 + \eta\lambda^*)^t \\ \frac{1}{4} &\leq \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2) \leq 1 \end{aligned} \quad (\text{C.2})$$

because

$$\prod_{s=1}^{T'_2} \left(\frac{1 + \eta\lambda^* - \eta\tilde{\beta}_s^2}{1 + \eta\lambda^*} \right) \geq \prod_{s=1}^{T'_2} (1 - \eta\tilde{\beta}_s^2) \geq \left(1 - \frac{\eta\lambda^*}{64 \log n} \right)^{\frac{64 \log n}{\eta\lambda^*}} \geq \frac{1}{4}$$

if $\frac{\eta\lambda^*}{64 \log n} \leq \frac{1}{2}$. Note that the upper bounds in (C.2) hold even if $t > T_2$.

From (C.2), we have the approximation $\tilde{\mathbf{x}}^{(t)} \approx \mathbf{x}^{(0)} + (1 + \eta\lambda^*)^t (\mathbf{u}^{*\top} \mathbf{x}^{(0)}) \mathbf{u}^*$ for all $t \leq T'_2$ and the ℓ_2 -norm of $\tilde{\mathbf{x}}^{(t)}$ is also approximately given by $\left(1 + \frac{(1 + \eta\lambda^*)^t}{\sqrt{n}}\right) \beta_0$. The ℓ_∞ -norm is about $\frac{1}{\sqrt{n}}$ times smaller than the ℓ_2 -norm. We make this observation rigorous with the following lemma.

Lemma C.1. *For all $t \leq T'_2$, we have*

$$\begin{aligned} \frac{1}{8} \frac{1}{\sqrt{\log n}} \left(1 + \frac{(1 + \eta\lambda^*)^t}{\sqrt{n}} \right) \beta_0 &\leq \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2 \leq 2\sqrt{\log n} \left(1 + \frac{(1 + \eta\lambda^*)^t}{\sqrt{n}} \right) \beta_0, \\ \frac{1}{4} \frac{1}{\sqrt{\log n}} \left(1 + \frac{(1 + \eta\lambda^*)^t}{\sqrt{n}} \right) \frac{\beta_0}{\sqrt{n}} &\leq \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty \leq 2\sqrt{\log n} \left(1 + (1 + \eta\lambda^*)^t \sqrt{\frac{\mu}{n}} \right) \frac{\beta_0}{\sqrt{n}}. \end{aligned}$$

Proof. For brevity, let us drop the superscript (t) and write $\tilde{\mathbf{x}} = A\mathbf{x}^{(0)} + B\mathbf{u}^*$. For the upper bounds, we may use the triangle inequality

$$\left\| A\mathbf{x}^{(0)} + B\mathbf{u}^* \right\| \leq A \left\| \mathbf{x}^{(0)} \right\| + |B| \left\| \mathbf{u}^* \right\|.$$

If we use (B.5) and (C.2), we get the upper bounds for A and B . We have $\|\mathbf{x}^{(0)}\|_2 \leq 2\beta_0$ by (B.1), and the ℓ_∞ -norm of $\mathbf{x}^{(0)}$ is controlled through (B.2). These finish the proof for the upper bounds.

From the definition of B , we have $B(\mathbf{u}^{\star\top} \mathbf{x}^{(0)}) \geq 0$, and

$$\begin{aligned} \|A\mathbf{x}^{(0)} + B\mathbf{u}^*\|_2^2 &= A^2 \|\mathbf{x}^{(0)}\|_2^2 + B^2 + 2AB(\mathbf{u}^{\star\top} \mathbf{x}^{(0)}) \\ &\geq A^2 \|\mathbf{x}^{(0)}\|_2^2 + B^2 \\ &\geq \frac{1}{4} \left(A \|\mathbf{x}^{(0)}\|_2 + |B| \right)^2. \end{aligned}$$

(B.1) and (B.4) together with the lower bound in (C.2) give the desired lower bound for $\|\tilde{\mathbf{x}}\|_2$. The lower bound for ℓ_∞ -norm is directly implied from Lemma B.2 together with (B.4) and (C.2). \square

With Lemma C.1, we have

$$\frac{1}{8} \frac{1}{\sqrt{\log n}} \frac{(1 + \eta\lambda^*)^{T'_2}}{\sqrt{n}} \beta_0 \leq \frac{1}{8} \frac{1}{\sqrt{\log n}} \left(1 + \frac{(1 + \eta\lambda^*)^{T'_2}}{\sqrt{n}} \right) \beta_0 \leq \|\tilde{\mathbf{x}}^{(T'_2)}\|_2 \leq \frac{\sqrt{\lambda^*}}{8\sqrt{\log n}},$$

and thus

$$T'_2 \leq \frac{1}{\log(1 + \eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0} \leq \frac{11 \log n}{\log(1 + \eta\lambda^*)} \leq \frac{64 \log n}{\eta\lambda^*}.$$

For $t \leq T_1$ where $(1 + \eta\lambda^*)^t$ is not big, the bounds in Lemma C.1 are simplified to

$$\frac{1}{\sqrt{\log n}} \beta_0 \lesssim \|\tilde{\mathbf{x}}^{(t)}\|_2 \lesssim \sqrt{\log n} \beta_0, \quad (\text{C.3})$$

$$\frac{1}{\sqrt{\log n}} \frac{\beta_0}{\sqrt{n}} \lesssim \|\tilde{\mathbf{x}}^{(t)}\|_\infty \lesssim \sqrt{\log n} \frac{\beta_0}{\sqrt{n}}. \quad (\text{C.4})$$

After $\tilde{\mathbf{x}}^{(t)}$ becomes almost parallel to \mathbf{u}^* and before T'_2 , we could approximate $\tilde{\beta}_t$ as increasing with the rate $(1 + \eta\lambda^*)$. However, after T'_2 , this approximation is invalid, and $\tilde{\beta}_t$ grows at a slower rate as it increases and it eventually converges to $\sqrt{\lambda^*}$. How much iterations will be required for it to reach $\sqrt{\lambda^*} \sqrt{1 - \frac{1}{\log n}}$ after T'_2 ? With Lemma C.2, we will prove that $O(\log \log n)$ iterations are required after T'_2 .

Lemma C.2. *At $t = T'_2 + \frac{6 \log \log n}{\log(1 + \eta\lambda^*)}$, we have $\tilde{\beta}_t^2 \geq \lambda^* \left(1 - \frac{1}{\log n}\right)$.*

Proof. From the decomposition (C.1), we have

$$\begin{aligned} \left| \|\tilde{\mathbf{x}}^{(t)}\|_2 - |B^{(t)}| \right| &\leq \|\mathbf{x}^{(0)}\|_2 \\ \left| \mathbf{u}^{\star\top} \mathbf{x}^{(0)} - |B^{(t)}| \right| &\leq |A^{(t)}| \|\mathbf{u}^{\star\top} \mathbf{x}^{(0)}\| \leq \|\mathbf{x}^{(0)}\|_2, \end{aligned}$$

and thus

$$\left| \tilde{\alpha}_t - \tilde{\beta}_t \right| \leq 2 \|\mathbf{x}^{(0)}\|_2 \leq \frac{\sqrt{\lambda^*}}{3 \log^2 n}. \quad (\text{C.5})$$

holds for all t . Because $\tilde{\beta}_t \gtrsim \frac{1}{\log n}$ for all $t \geq T'_2$, (C.5) implies that $\tilde{\beta}_t$ is well approximated by $\tilde{\alpha}_t$. Hence, we will focus on $\tilde{\alpha}_t$, which is an increasing sequence that evolves with

$$\tilde{\alpha}_{t+1} = (1 - \eta\tilde{\beta}_t^2 + \eta\lambda^*) \tilde{\alpha}_t.$$

For all $i \geq 1$, let N_i be the last t such that $\lambda^* - \tilde{\alpha}_t^2 \geq \frac{\lambda^*}{e^i}$. Then, we have

$$\lambda^* - \tilde{\alpha}_{N_i}^2 \geq \frac{\lambda^*}{e^i} > \lambda^* - \tilde{\alpha}_{N_i+1}^2. \quad (\text{C.6})$$

Let $i \geq 2$. For all $N_{i-1} < t \leq N_i$,

$$\frac{\tilde{\alpha}_{t+1}}{\tilde{\alpha}_t} = 1 - \eta \tilde{\beta}_t^2 + \eta \lambda^* = 1 + \eta(\lambda^* - \tilde{\alpha}_t^2) + \eta(\tilde{\alpha}_t^2 - \tilde{\beta}_t^2) \geq 1 + \frac{\eta \lambda^*}{e^i} - \frac{\eta \lambda^*}{\log^2 n} \geq 1 + 0.99 \frac{\eta \lambda^*}{e^i}.$$

We used (C.5), (C.6), and the fact that $\tilde{\alpha}_t, \tilde{\beta}_t \leq \sqrt{\lambda^*}$ for all t . This implies

$$\left(1 + 0.99 \frac{\eta \lambda^*}{e^i}\right)^{N_i - N_{i-1} - 1} x_{N_{i-1}+1} \leq x_{N_i}.$$

From the lower and upper bounds provided by (C.6), we have

$$\begin{aligned} \sqrt{\lambda^*} \sqrt{1 - \frac{1}{e^{i-1}}} \left(1 + 0.99 \frac{\eta \lambda^*}{e^i}\right)^{N_i - N_{i-1} - 1} &\leq \sqrt{\lambda^*} \sqrt{1 - \frac{1}{e^i}}, \\ \left(1 + 0.99 \frac{\eta \lambda^*}{e^i}\right)^{2(N_i - N_{i-1} - 1)} &\leq \frac{e^i - 1}{e^i} \frac{e^{i-1}}{e^{i-1} - 1} = \frac{e^{i-1} - \frac{1}{e}}{e^{i-1} - 1} \leq 1 + \frac{1}{e^{i-1}}. \end{aligned}$$

Taking log on both sides and using the inequality $\frac{1}{2}x < \log(1+x) < x$ that holds for $0 < x < 1$, we get

$$N_i - N_{i-1} \leq 1 + \frac{1}{2} \frac{\log\left(1 + \frac{1}{e^{i-1}}\right)}{\log\left(1 + 0.99 \frac{\eta \lambda^*}{e^i}\right)} \leq 1 + \frac{e}{0.99 \eta \lambda^*}.$$

For $t \leq N_1$, we have

$$\frac{\tilde{\alpha}_{t+1}}{\tilde{\alpha}_t} \geq 1 + 0.99 \frac{\eta \lambda^*}{e},$$

and thus

$$\sqrt{\lambda^*} \sqrt{1 - \frac{1}{e}} \geq \alpha_{N_1} \geq \left(1 + 0.99 \frac{\eta \lambda^*}{e}\right)^{N_1} \tilde{\alpha}_{T_2} = \left(1 + 0.99 \frac{\eta \lambda^*}{e}\right)^{N_1} \sqrt{\frac{\lambda^*}{21 \log n}}.$$

Taking log on both sides we get

$$N_1 \leq 3 \frac{\log \log n}{\eta \lambda^*}.$$

Hence, we have

$$\begin{aligned} N_{\log \log n+1} + 1 &\leq \left(1 + \frac{e}{0.99 \eta \lambda^*}\right) \log \log n + N_1 + 1 \\ &\leq \left(1 + \frac{e}{0.99 \eta \lambda^*}\right) \log \log n + 3 \frac{\log \log n}{\eta \lambda^*} + 1 \\ &\leq \frac{6 \log \log n}{\log(1 + \eta \lambda^*)}, \end{aligned}$$

but at $t = N_{\log \log n+1} + 1$, it holds that

$$\tilde{\alpha}_t^2 > \lambda^* \left(1 - \frac{1}{e \log n}\right),$$

and we have

$$\tilde{\beta}_t^2 > \lambda^* \left(1 - \frac{1}{\log n}\right)$$

as desired. Note that $\tilde{\beta}_t$ is also an increasing sequence as $\tilde{\alpha}_t$. \square

It is implied from Lemma C.2 that $T_2 \leq \frac{1}{\log(1+\eta\lambda^*)} \log \frac{\sqrt{\lambda^* n}}{\beta_0} + \frac{6 \log \log n}{\log(1+\eta\lambda^*)} = (1+o(1)) \frac{1}{\eta \lambda^*} \log \frac{\sqrt{\lambda^* n}}{\beta_0}$. The following corollary shows that $\tilde{\mathbf{x}}^{(t)}$ is sufficiently close to \mathbf{x}^* at $t = T_2$.

Corollary C.3. *At $t = T_2$, we have*

$$\min \left\{ \left\| \tilde{\mathbf{x}}^{(t)} - \mathbf{x}^* \right\|_2, \left\| \tilde{\mathbf{x}}^{(t)} + \mathbf{x}^* \right\|_2 \right\} \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_2, \quad (\text{C.7})$$

$$\min \left\{ \left\| \tilde{\mathbf{x}}^{(t)} - \mathbf{x}^* \right\|_\infty, \left\| \tilde{\mathbf{x}}^{(t)} + \mathbf{x}^* \right\|_\infty \right\} \lesssim \frac{1}{\sqrt{\log n}} \|\mathbf{x}^*\|_\infty. \quad (\text{C.8})$$

Proof. When $B^{(t)} > 0$, from the decomposition

$$\mathbf{x}^{(t)} - \mathbf{x}^* = A^{(t)}\mathbf{x}^{(0)} + (B^{(t)} - \tilde{\beta}_t)\mathbf{u}^* + (\tilde{\beta}_t - \sqrt{\lambda^*})\mathbf{u}^*,$$

we have

$$\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq 2\|\mathbf{x}^{(0)}\|_2 + \frac{\sqrt{\lambda^*}}{\sqrt{\log n}} \leq \frac{2\sqrt{\lambda^*}}{\sqrt{\log n}} = \frac{2}{\sqrt{\log n}}\|\mathbf{x}^*\|_2.$$

For the cases $B^{(t)} < 0$ and ℓ_∞ -norm, we may use similar technique. \square

D Phase I

D.1 Proof of Lemma 5.2

In this subsection, we provide a proof to the following lemma, which is a formal statement of Lemma 5.2.

Lemma D.1. *With high probability, there exists a universal constant $c_0 > 0$ such that*

$$\|\hat{\mathbf{x}}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2 \leq c_0\mu\sqrt{\frac{\log n}{np}}\beta_0t \quad (\text{D.1})$$

for all $t \leq T_1$, if $n^2p \gtrsim \mu^4n \log^{21} n$ and the initialization point $\mathbf{x}^{(0)}$ satisfies (B.1) to (B.6).

Proof. Let us rewrite the update equations (2) and (15) as

$$\begin{aligned} \tilde{\mathbf{x}}^{(t+1)} &= \left(\mathbf{I} - \eta\tilde{\beta}_t^2 + \eta\mathbf{M}^*\right)\tilde{\mathbf{x}}^{(t)}, \\ \hat{\mathbf{x}}^{(t+1)} &= \left(\mathbf{I} - \eta\tilde{\beta}_t^2 + \eta\mathbf{M}^\circ\right)\hat{\mathbf{x}}^{(t)}, \end{aligned}$$

where $\tilde{\beta}_t = \|\tilde{\mathbf{x}}^{(t)}\|_2^2$. Then, $\hat{\mathbf{x}}^{(t)} - \tilde{\mathbf{x}}^{(t)}$ is a product between $\mathbf{x}^{(0)}$ and $P^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H})$, which is a matrix polynomial of $\mathbf{I}, \mathbf{M}^*, \mathbf{H}$, where $\mathbf{H} = \mathbf{M}^\circ - \mathbf{M}^*$.

$$P^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H}) := \left(\prod_{s=1}^t \left((1 - \eta\tilde{\beta}_s^2)\mathbf{I} + \eta\mathbf{M}^* + \eta\mathbf{H} \right) - \prod_{s=1}^t \left((1 - \eta\tilde{\beta}_s^2)\mathbf{I} + \eta\mathbf{M}^* \right) \right) \quad (\text{D.2})$$

$$\hat{\mathbf{x}}^{(t)} - \tilde{\mathbf{x}}^{(t)} = P^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H})\mathbf{x}^{(0)} \quad (\text{D.3})$$

We classify the terms that appear after expanding the matrix polynomial $P^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H})$ into two types; 1) the terms that contain \mathbf{H} but not \mathbf{M}^* , 2) the terms that contain both \mathbf{H} and \mathbf{M}^* . We define $P_1^{(t)}(\mathbf{I}, \mathbf{H})$ to be a matrix polynomial of \mathbf{I} and \mathbf{H} , which is equal to summation of the first type, and it is explicitly written as

$$P_1^{(t)}(\mathbf{I}, \mathbf{H}) = \prod_{s=1}^t \left((1 - \eta\tilde{\beta}_s^2)\mathbf{I} + \eta\mathbf{H} \right) - \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2)\mathbf{I}.$$

We correspondingly define $P_2^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H})$ to be summation of the second type, and it is equal to

$$P_2^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H}) = P^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H}) - P_1^{(t)}(\mathbf{I}, \mathbf{H}).$$

For $x, y \in \mathbb{R}$, we define $P_1^{(t)}(x, y)$ as the value that is obtained by substituting x, y instead of \mathbf{I}, \mathbf{H} , respectively. For example, $P_1^{(t)}(1, 2) = \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2 + 2\eta) - \prod_{s=1}^t (1 - \eta\tilde{\beta}_s^2)$. For $x, y, z \in \mathbb{R}$, $P_2^{(t)}(x, y, z)$ is defined in a similar manner.

We bound the contribution of each type separately because the triangle inequality gives

$$\|\hat{\mathbf{x}}^{(t)} - \tilde{\mathbf{x}}^{(t)}\|_2 \leq \|P_1^{(t)}(\mathbf{I}, \mathbf{H})\mathbf{x}^{(0)}\|_2 + \|P_2^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H})\mathbf{x}^{(0)}\|_2.$$

Every term in $P_1^{(t)}(\mathbf{I}, \mathbf{H})$ is \mathbf{H}^s times a constant. We have $\|\mathbf{H}^s\mathbf{x}^{(0)}\|_2 \leq \|\mathbf{H}\|^s\|\mathbf{x}^{(0)}\|_2$, and hence with triangle inequality

$$\|P_1^{(t)}(\mathbf{I}, \mathbf{H})\mathbf{x}^{(0)}\|_2 \leq P_1^{(t)}(1, \|\mathbf{H}\|)\beta_0.$$

If $n^2 p \gtrsim \mu^2 n \log^3 n$, we can further bound $P_1^{(t)}(1, \|\mathbf{H}\|)$ as

$$\begin{aligned}
P_1^{(t)}(1, \|\mathbf{H}\|) &= \prod_{s=1}^t (1 - \eta \tilde{\beta}_s^2 + \eta \|\mathbf{H}\|) - \prod_{s=1}^t (1 - \eta \tilde{\beta}_s^2) \\
&= \prod_{s=1}^t (1 - \eta \tilde{\beta}_s^2) \left(\prod_{s=1}^t \left(1 + \frac{\eta \|\mathbf{H}\|}{1 - \eta \tilde{\beta}_s^2} \right) - 1 \right) \\
&\leq \left(1 + \frac{\eta}{1 - \eta \lambda^*} \|\mathbf{H}\| \right)^t - 1 \\
&\leq \left(\exp \left(\frac{\eta}{1 - \eta \lambda^*} \|\mathbf{H}\| t \right) - 1 \right) \\
&\leq \frac{2\eta}{1 - \eta \lambda^*} \|\mathbf{H}\| t
\end{aligned}$$

The third line uses the fact that $\tilde{\beta}_t^2 \leq \lambda^*$ for all $t \leq T_1$. The fourth and fifth lines are derived from an elementary inequality $1 + x \leq e^x \leq 1 + 2x$, which holds for small $x > 0$. Note that $\eta \|\mathbf{H}\| t \lesssim \eta \lambda^* \mu \sqrt{\frac{\log n}{np}} \ll 1$ from Lemmas A.1 and A.3, and the fact that $T_1 \lesssim \log n$.

We can decompose every term of second type as a product of η , λ^* , \mathbf{u}^* , $\mathbf{H}^s \mathbf{u}^*$, $\mathbf{u}^{*\top} \mathbf{H}^s \mathbf{x}^{(0)}$, $\mathbf{u}^{*\top} \mathbf{H}^s \mathbf{u}^*$, and $\mathbf{u}^{*\top} \mathbf{x}^{(0)}$. We describe this with some examples.

$$\begin{aligned}
(\eta \mathbf{H})^{s_1} (\eta \mathbf{M}^*) (\eta \mathbf{H})^{s_2} \mathbf{x}^{(0)} &= \eta^{s_1 + s_2 + 1} \lambda^* (\mathbf{H}^{s_1} \mathbf{u}^*) (\mathbf{u}^{*\top} \mathbf{H}^{s_2} \mathbf{x}^{(0)}) \\
(\eta \mathbf{H})^s (\eta \mathbf{M}^*) \mathbf{x}^{(0)} &= \eta^{s+1} (\mathbf{H}^s \mathbf{u}^*) (\mathbf{u}^{*\top} \mathbf{x}^{(0)}) \\
(\eta \mathbf{M}^*) (\eta \mathbf{H})^s (\eta \mathbf{M}^*) \mathbf{x}^{(0)} &= \eta^{s+2} \lambda^{*2} \mathbf{u}^* (\mathbf{u}^{*\top} \mathbf{H} \mathbf{u}^*) (\mathbf{u}^{*\top} \mathbf{x}^{(0)}) \\
(\eta \mathbf{M}^*) (\eta \mathbf{H})^s \mathbf{x}^{(0)} &= \eta^{s+1} \lambda^* \mathbf{u}^* (\mathbf{u}^{*\top} \mathbf{H}^s \mathbf{x}^{(0)})
\end{aligned}$$

The terms $\mathbf{H}^s \mathbf{u}^*$ and $\mathbf{u}^{*\top} \mathbf{H}^s \mathbf{u}^*$ are bounded with

$$\|\mathbf{H}^s \mathbf{u}^*\|_2 \leq \|\mathbf{H}\|^s, \quad |\mathbf{u}^{*\top} \mathbf{H}^s \mathbf{u}^*| \leq \|\mathbf{H}\|^s, \quad (\text{D.4})$$

and the terms that contain $\mathbf{x}^{(0)}$ are bounded with (B.3). For every term of second type that includes s_1 times of $\eta \mathbf{M}^*$ and s_2 times of $\eta \mathbf{H}$, the bounds (D.4) and (B.3) imply that ℓ_2 -norm of the term multiplied by $\mathbf{x}^{(0)}$ is at most

$$(\eta \lambda^*)^{s_1} (\eta \|\mathbf{H}\|)^{s_2} 2 \sqrt{\frac{\log n}{n}} \beta_0.$$

Hence, similar to the first type, we have

$$\left\| P_2^{(t)}(\mathbf{I}, \mathbf{M}^*, \mathbf{H}) \mathbf{x}^{(0)} \right\|_2 \leq P_2^{(t)}(1, \lambda^*, \|\mathbf{H}\|) 2 \sqrt{\frac{\log n}{n}} \beta_0.$$

If $n^2 p \gtrsim \mu^2 n \log^3 n$, we can further bound $P_2^{(t)}(1, \lambda^*, \|\mathbf{H}\|)$ as

$$\begin{aligned}
P_2^{(t)}(1, \lambda^*, \|\mathbf{H}\|) &= \prod_{s=1}^t (1 - \eta \beta_s^2 + \eta \lambda^* + \eta \|\mathbf{H}\|) - \prod_{s=1}^t (1 - \eta \beta_s^2 + \eta \lambda^*) - P_1^{(t)}(1, \|\mathbf{H}\|) \\
&\leq \prod_{s=1}^t (1 - \eta \beta_s^2 + \eta \lambda^* + \eta \|\mathbf{H}\|) - \prod_{s=1}^t (1 - \eta \beta_s^2 + \eta \lambda^*) \\
&\leq \left(\prod_{s=1}^t (1 - \eta \beta_s^2 + \eta \lambda^*) \right) \left(\prod_{s=1}^t \left(1 + \frac{\eta}{1 - \eta \beta_s^2 + \eta \lambda^*} \|\mathbf{H}\| \right) - 1 \right) \\
&\leq (1 + \eta \lambda^*)^t \left((1 + \eta \|\mathbf{H}\|)^t - 1 \right) \\
&\leq 2\eta \|\mathbf{H}\| t (1 + \eta \lambda^*)^t.
\end{aligned}$$

Combining all, we have

$$\begin{aligned}
\left\| \widehat{\mathbf{x}}^{(t)} - \widetilde{\mathbf{x}}^{(t)} \right\|_2 &\leq 4\eta \|\mathbf{H}\| t \left(1 + \sqrt{\frac{\log n}{n}} (1 + \eta\lambda^*)^t \right) \beta_0 \\
&\leq 4\eta \|\mathbf{H}\| t \left(1 + \sqrt{\frac{\log n}{n}} (1 + \eta\lambda^*)^{T_1} \right) \beta_0 \\
&\leq 4\eta \|\mathbf{H}\| t \left(1 + \sqrt{\frac{\mu^4 \log^{22} n}{np}} \right) \beta_0 \\
&\leq c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 t
\end{aligned}$$

for all $t \leq T_1$ for some constant $c_0 > 0$ if $n^2 p \gtrsim \mu^4 n \log^{22} n$. \square

D.2 Proof of Lemmas 5.3 to 5.5

We prove Lemmas 5.3 to 5.5 all together in an inductive manner.

Lemma D.2. *Suppose that the initialization point $\mathbf{x}^{(0)}$ satisfies (B.1) to (B.6). If $n^2 p \gtrsim \mu^5 n \log^{22} n$, for all $t \leq T_1$, we have*

$$\left\| \mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)} \right\|_2 \leq 2c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 t, \quad (\text{D.5})$$

$$\left\| \mathbf{x}^{(t)} - \widetilde{\mathbf{x}}^{(t)} \right\|_\infty \leq (3c_0 + c_5) \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t^2, \quad (\text{D.6})$$

$$\left\| \mathbf{x}^{(t)} - \widehat{\mathbf{x}}^{(t)} \right\|_2 \leq c_2 (3c_0 + c_5 + 1) \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta\lambda^*)^t \beta_0^3, \quad (\text{D.7})$$

$$\left\| \mathbf{x}^{(t)} - \mathbf{x}^{(t,l)} \right\|_2 \leq c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t, \quad (\text{D.8})$$

$$\left| \mathbf{u}^{(l)\top} (\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}) \right| \leq c_6 \sqrt{\frac{\mu^3 \log^2 n}{np}} (1 + \eta\lambda^*)^t \frac{\beta_0}{n}, \quad (\text{D.9})$$

$$\left| \left(\mathbf{x}^{(t,l)} - \widetilde{\mathbf{x}}^{(t)} \right)_l \right| \leq 3c_0 \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} t^2, \quad (\text{D.10})$$

with high probability, where c_2, c_5, c_6 are positive constants.

Before we start the proof, we introduce some notations. For $\mathbf{x} \in \mathbb{R}^n$, let us define

$$\|\mathbf{x}\|_{2,i} = \sqrt{\frac{1}{p} \sum_{j=1}^n \delta_{ij} x_j^2}, \quad \mathbf{I}_{\mathbf{x}} = \frac{1}{\|\mathbf{x}\|_2^2} \text{diag}(\|\mathbf{x}\|_{2,1}^2, \dots, \|\mathbf{x}\|_{2,n}^2).$$

$\|\mathbf{x}\|_{2,i}$ is the ℓ_2 -norm of \mathbf{x} estimated with sampling of the i th row. With this notation, we can write the gradient of f as

$$\nabla f(\mathbf{x}) = \|\mathbf{x}\|_2^2 \mathbf{I}_{\mathbf{x}} \mathbf{x} - \mathbf{M}^\circ \mathbf{x}.$$

The function g is defined as

$$g(\mathbf{x}) = \frac{1}{4p} \left\| \mathcal{P}_\Omega(\mathbf{x} \mathbf{x}^\top) \right\|_{\text{F}}^2,$$

and its gradient satisfies

$$\nabla f(\mathbf{x}) = \nabla g(\mathbf{x}) - \mathbf{M}^\circ \mathbf{x}.$$

The Hessian of $g(\mathbf{x})$ is equal to

$$\nabla^2 g(\mathbf{x}) = \|\mathbf{x}\|_2^2 \mathbf{I}_x + \frac{2}{p} \mathcal{P}_\Omega(\mathbf{x} \mathbf{x}^\top).$$

The base case $t = 0$ for induction hypotheses (D.5) to (D.10) trivially hold because all three sequences $\mathbf{x}^{(t)}$, $\widehat{\mathbf{x}}^{(t)}$, $\widetilde{\mathbf{x}}^{(t)}$ start from the same point. Now, we assume that the hypotheses hold up to the t th iteration and show that they hold at the $(t+1)$ st iteration. For brevity, we drop the superscript (t) from $\mathbf{x}^{(t)}$, $\mathbf{x}^{(t,l)}$, $\widehat{\mathbf{x}}^{(t)}$, $\widetilde{\mathbf{x}}^{(t)}$ and denote them as \mathbf{x} , $\mathbf{x}^{(l)}$, $\widehat{\mathbf{x}}$, $\widetilde{\mathbf{x}}$, respectively. Also, recall that T_1 is defined to be the last t such that $(1 + \eta\lambda^*)^t \leq \sqrt{\frac{\mu^4 \log^{21} n}{np}} \sqrt{n}$, and the magnitude of initialization satisfies $\beta_0^2 \lesssim \lambda^* \sqrt{\frac{np}{\mu^5 \log^{26} n}} \frac{1}{\sqrt{n}}$ so that there exists a constant $c_1 > 0$ such that $(1 + \eta\lambda^*)^t \beta_0^2 \leq c_1 \frac{\lambda^*}{\sqrt{\mu \log^5 n}}$.

(D.7) at $(t+1)$ We first decompose $\mathbf{x}^{(t+1)} - \widehat{\mathbf{x}}^{(t+1)}$ as

$$\begin{aligned} \mathbf{x}^{(t+1)} - \widehat{\mathbf{x}}^{(t+1)} &= (\mathbf{I} + \eta \mathbf{M}^\circ)(\mathbf{x} - \widehat{\mathbf{x}}) - \eta \|\mathbf{x}\|_2^2 \mathbf{I}_x \mathbf{x} + \eta \|\widetilde{\mathbf{x}}\|_2^2 \widehat{\mathbf{x}} \\ &= (\mathbf{I} - \eta \|\widetilde{\mathbf{x}}\|_2^2 \mathbf{I} + \eta \mathbf{M}^\circ)(\mathbf{x} - \widehat{\mathbf{x}}) - \eta (\|\mathbf{x}\|_2^2 \mathbf{I}_x - \|\widetilde{\mathbf{x}}\|_2^2 \mathbf{I}) \mathbf{x}. \end{aligned} \quad (\text{D.11})$$

With the help of Lemma G.8, we bound the maximum entry of a diagonal matrix $\|\mathbf{x}\|_2^2 \mathbf{I}_x - \|\widetilde{\mathbf{x}}\|_2^2 \mathbf{I}$. We have

$$\begin{aligned} \max_{i \in [n]} \left| \|\mathbf{x}\|_{2,i}^2 - \|\widetilde{\mathbf{x}}\|_2^2 \right| &\lesssim n \|\mathbf{x} - \widetilde{\mathbf{x}}\|_\infty \|\mathbf{x} + \widetilde{\mathbf{x}}\|_\infty + \sqrt{\frac{\log n}{p}} \|\widetilde{\mathbf{x}}\|_2 \|\widetilde{\mathbf{x}}\|_\infty + \frac{\log n}{p} \|\widetilde{\mathbf{x}}\|_\infty^2 \\ &\lesssim (3c_0 + c_5) \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^2 t^2 + \sqrt{\frac{\log^2 n}{np}} \beta_0^2 + \frac{\log^2 n}{np} \beta_0^2 \\ &\lesssim \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^2 ((3c_0 + c_5)t^2 + 1) \end{aligned}$$

if $n^2 p \gtrsim n \log^2 n$. Hence, there exists a universal constant $c_2 > 0$ that is independent of t such that

$$\left(\max_{i \in [n]} \left| \|\mathbf{x}^{(t)}\|_{2,i}^2 - \|\widetilde{\mathbf{x}}^{(t)}\|_2^2 \right| \right) \|\mathbf{x}^{(t)}\|_2 \leq \frac{1}{2} c_2 \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^3 ((3c_0 + c_5)t^2 + 1)$$

for all $t \leq T_1$. With the decomposition (D.11), we have

$$\begin{aligned} \left\| \mathbf{x}^{(t+1)} - \widehat{\mathbf{x}}^{(t+1)} \right\|_2 &\leq (1 - \eta \|\widetilde{\mathbf{x}}\|_2^2 + \eta \lambda^\circ) \|\mathbf{x} - \widehat{\mathbf{x}}\|_2 + \eta \left(\max_{i \in [n]} \left| \|\mathbf{x}\|_{2,i}^2 - \|\widetilde{\mathbf{x}}\|_2^2 \right| \right) \|\mathbf{x}\|_2 \\ &\leq (1 + \eta \lambda^\circ) \|\mathbf{x} - \widehat{\mathbf{x}}\|_2 + \frac{1}{2} c_2 (\eta \lambda^*)^3 \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^3 ((3c_0 + c_5)t^2 + 1). \end{aligned}$$

From (A.1), there exists a universal constant $c_3 > 0$ such that $\eta \lambda^\circ \leq \eta \lambda^* + \frac{c_3}{\log^{\frac{3}{2}} n}$ if $n^2 p \gtrsim \mu^2 n \log^5 n$. Combining all, for all $s \leq t$, we have

$$\begin{aligned} \left\| \mathbf{x}^{(s+1)} - \widehat{\mathbf{x}}^{(s+1)} \right\|_2 &\leq \left(1 + \eta \lambda^* + \frac{c_3}{\log^2 n} \right) \left\| \mathbf{x}^{(s)} - \widehat{\mathbf{x}}^{(s)} \right\|_2 \\ &\quad + \frac{1}{2} c_2 (\eta \lambda^*)^3 \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^3 ((3c_0 + c_5)s^2 + 1). \end{aligned}$$

An analysis on the recursive equation

$$x_{s+1} = \left(1 + \eta \lambda^* + \frac{c_3}{\log^2 n} \right) x_s + \frac{1}{2} c_2 (\eta \lambda^*)^3 \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} \beta_0^3 ((3c_0 + c_5)s^2 + 1), \quad x_0 = 0,$$

proves that

$$\left\| \mathbf{x}^{(t+1)} - \widehat{\mathbf{x}}^{(t+1)} \right\|_2 \leq c_2 (3c_0 + c_5 + 1) \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta \lambda^*)^{t+1} \beta_0^3.$$

(D.8) at $(t+1)$ We decompose $\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)}$ as

$$\begin{aligned}
\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} &= (1 - \eta \|\tilde{\mathbf{x}}\|_2^2) (\mathbf{x} - \mathbf{x}^{(l)}) - 2\eta \underbrace{\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top (\mathbf{x} - \mathbf{x}^{(l)})}_{\textcircled{1}} \\
&\quad - \eta \underbrace{\int_0^1 \left(\nabla^2 g(\mathbf{x}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top \right) \right) (\mathbf{x} - \mathbf{x}^{(l)}) d\tau}_{\textcircled{2}} \\
&\quad - \eta \underbrace{\left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{x}^{(t,l)} \mathbf{x}^{(t,l)\top}) - \mathcal{P}_l(\mathbf{x}^{(t,l)} \mathbf{x}^{(t,l)\top}) \right) \mathbf{x}^{(t,l)}}_{\textcircled{3}} \\
&\quad + \eta \underbrace{\lambda^{(l)} \mathbf{u}^{(l)} \mathbf{u}^{(l)\top} (\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)})}_{\textcircled{4}} + \eta \underbrace{\left(\mathbf{M}^\circ - \lambda^{(l)} \mathbf{u}^{(l)} \mathbf{u}^{(l)\top} \right) (\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)})}_{\textcircled{5}} \\
&\quad + \eta \underbrace{\left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^\star) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^\star) \right) \mathbf{x}^{(t,l)}}_{\textcircled{6}} + \eta \underbrace{\left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) - \mathbf{E}^{(l)} \right) \mathbf{x}^{(t,l)}}_{\textcircled{7}},
\end{aligned}$$

where $\mathbf{x}^{(l)}(\tau) = \mathbf{x}^{(l)} + \tau(\mathbf{x} - \mathbf{x}^{(l)})$. $\textcircled{1}$ is easily bounded by

$$\|\textcircled{1}\|_2 \leq \|\tilde{\mathbf{x}}\|_2^2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \beta_0^2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \lambda^\star \sqrt{\frac{1}{\mu^5 \log^{26} n}} \sqrt{\frac{np}{n}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2.$$

From Lemma G.10 and (C.4), for all $0 \leq \tau \leq 1$, we have

$$\left\| \nabla^2 g(\mathbf{x}^{(l)}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top \right) \right\| \lesssim n \|\mathbf{x}(\tau) - \tilde{\mathbf{x}}\|_\infty (\|\mathbf{x}(\tau)\|_\infty + \|\tilde{\mathbf{x}}\|_\infty) + \sqrt{\frac{\log^3 n}{np}} \beta_0^2. \quad (\text{D.12})$$

From the definition of $\mathbf{x}(\tau)$, we have

$$\|\mathbf{x}^{(l)}(\tau) - \tilde{\mathbf{x}}\|_\infty \leq (1 - \tau) \|\mathbf{x}^{(l)} - \mathbf{x}\|_2 + \|\mathbf{x} - \tilde{\mathbf{x}}\|_\infty \lesssim \sqrt{\frac{\mu^3 \log^6 n}{np}} \frac{\beta_0}{\sqrt{n}},$$

where the last inequality is from the induction hypotheses (D.6), (D.8), and the fact that $t \leq T_1 \lesssim \log n$. Inserting this bound back to (D.12), we get

$$\left\| \nabla^2 g(\mathbf{x}^{(l)}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top \right) \right\| \lesssim \sqrt{\frac{\mu^3 \log^7 n}{np}} \beta_0^2, \quad (\text{D.13})$$

which also implies

$$\|\textcircled{2}\|_2 \lesssim \sqrt{\frac{\mu^3 \log^7 n}{np}} \beta_0^2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \lambda^\star \sqrt{\frac{1}{\mu^2 n \log^{19} n}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2.$$

It is implied from Lemma G.11 that

$$\|\textcircled{3}\|_2 \leq \|\mathbf{x}^{(l)}\|_\infty^2 \sqrt{\frac{\log n}{p}} \|\mathbf{x}^{(l)}\|_2 \lesssim \sqrt{\frac{\log^4 n}{np}} \frac{\beta_0^3}{\sqrt{n}} \lesssim \lambda^\star \sqrt{\frac{1}{\mu^5 n \log^{22} n}} \frac{\beta_0}{\sqrt{n}}.$$

A bound on $\textcircled{4}$ follows from the induction hypothesis (D.9) and the spectral bound (A.2).

$$\|\textcircled{4}\|_2 \leq \lambda^{(l)} \left| \mathbf{u}^{(l)\top} (\mathbf{x}^{(t)} - \mathbf{x}^{(t,l)}) \right| \lesssim \lambda^\star \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{n} (1 + \eta \lambda^\star)^t \lesssim \lambda^\star \frac{\sqrt{\mu^7 \log^{23} n}}{np} \frac{\beta_0}{\sqrt{n}}.$$

The second largest eigenvalue of $\mathbf{M}^{(l)}$ is at most $\|\mathbf{M}^{(l)} - \mathbf{M}^\star\|$ by Weyl's Theorem, and from Lemmas A.1 and A.2, we have

$$\|\mathbf{M}^{(l)} - \mathbf{M}^\star\| \leq \|\mathbf{M}^{(l)} - \mathbf{M}^\circ\| + \|\mathbf{M}^\circ - \mathbf{M}^\star\| \lesssim \lambda^\star \mu \sqrt{\frac{\log n}{np}}.$$

Hence, we get

$$\begin{aligned}\|\textcircled{5}\|_2 &\leq \left(\|\mathbf{M}^\circ - \mathbf{M}^{(l)}\| + \|\mathbf{M}^{(l)} - \lambda^{(l)} \mathbf{u}^{(l)} \mathbf{u}^{(l)\top}\| \right) \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \\ &\lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2.\end{aligned}$$

Lastly, we apply Lemmas G.11 and G.13 to get

$$\|\textcircled{6}\|_2 \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \frac{\beta_0}{\sqrt{n}}, \quad \|\textcircled{7}\|_2 \lesssim \lambda^* \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}},$$

There exists a universal constant $c_4 > 0$ such that

$$\eta \left(2\|\textcircled{1}\|_2 + \|\textcircled{2}\|_2 + \|\textcircled{5}\|_2 \right) \leq \frac{c_4}{\log^2 n} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2$$

if $n^2 p \gtrsim \mu^2 n \log^5 n$, and there exists a universal constant $c_5 > 0$ such that

$$\eta \left(\|\textcircled{3}\|_2 + \|\textcircled{4}\|_2 + \|\textcircled{6}\|_2 + \|\textcircled{7}\|_2 \right) \leq \frac{1}{2} c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}}.$$

if $n^2 p \gtrsim \mu^5 n \log^{20} n$. Combining all, we have

$$\begin{aligned}\|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s+1,l)}\|_2 &\leq \left(1 - \eta \|\tilde{\mathbf{x}}^{(s)}\|_2^2 + \frac{c_4}{\log^2 n} \right) \|\mathbf{x}^{(s)} - \mathbf{x}^{(s,l)}\|_2 + \frac{1}{2} c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} \\ &\leq \left(1 + \frac{c_4}{\log^2 n} \right) \|\mathbf{x}^{(s)} - \mathbf{x}^{(s,l)}\|_2 + \frac{1}{2} c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}}.\end{aligned}$$

for all $s \leq t$. An analysis on the recursive equation

$$x_{s+1} = \left(1 + \frac{c_4}{\log^2 n} \right) x_s + \frac{1}{2} c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}}, \quad x_0 = 0$$

gives the desired bound

$$\begin{aligned}\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)}\|_2 &\leq \frac{\log^2 n}{c_4} \left(\left(1 + \frac{c_4}{\log^2 n} \right)^{t+1} - 1 \right) \frac{1}{2} c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} \\ &\leq c_5 \mu \sqrt{\frac{\log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} (t+1),\end{aligned}$$

where we used basic inequalities $(1+x)^a \leq e^{ax}$ and $e^{ax} - 1 \leq 2ax$ which hold if x is small and ax is small, respectively.

(D.9) at $(t+1)$ We decompose $\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)}$ as

$$\begin{aligned}
& \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \\
&= (\mathbf{x} - \eta \nabla f(\mathbf{x})) - (\mathbf{x}^{(l)} - \eta \nabla f^{(l)}(\mathbf{x}^{(l)})) \\
&= (\mathbf{x} - \eta \nabla g(\mathbf{x})) - (\mathbf{x}^{(l)} - \eta \nabla g^{(l)}(\mathbf{x}^{(l)})) + \eta (\mathbf{M}^\circ \mathbf{x} - \mathbf{M}^{(l)} \mathbf{x}^{(l)}) \\
&= (\mathbf{x} - \eta \nabla g(\mathbf{x})) - (\mathbf{x}^{(l)} - \eta \nabla g(\mathbf{x}^{(l)})) - \eta (\nabla g(\mathbf{x}^{(l)}) - \nabla g^{(l)}(\mathbf{x}^{(l)})) \\
&\quad + \eta (\mathbf{M}^\circ - \mathbf{M}^{(l)}) \mathbf{x} + \eta \mathbf{M}^{(l)} (\mathbf{x} - \mathbf{x}^{(l)}) \\
&= \int_0^1 (\mathbf{I} - \eta \nabla^2 g(\mathbf{x}^{(l)}(\tau)))(\mathbf{x} - \mathbf{x}^{(l)}) d\tau - \eta \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) - \mathcal{P}_l(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) \right) \mathbf{x}^{(l)} \\
&\quad + \eta \mathbf{M}^{(l)} (\mathbf{x} - \mathbf{x}^{(l)}) + \eta (\mathbf{M}^\circ - \mathbf{M}^{(l)}) (\mathbf{x} - \mathbf{x}^{(l)}) + \eta (\mathbf{M}^\circ - \mathbf{M}^{(l)}) \mathbf{x}^{(l)} \\
&= (1 - \eta \|\tilde{\mathbf{x}}\|_2^2)(\mathbf{x} - \mathbf{x}^{(l)}) - 2\eta \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top (\mathbf{x} - \mathbf{x}^{(l)}) \\
&\quad - \eta \int_0^1 \left(\nabla^2 g(\mathbf{x}^{(l)}(\tau)) - (\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top) \right) (\mathbf{x} - \mathbf{x}^{(l)}) d\tau \\
&\quad - \eta \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) - \mathcal{P}_l(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) \right) \mathbf{x}^{(l)} \\
&\quad + \eta \mathbf{M}^{(l)} (\mathbf{x} - \mathbf{x}^{(l)}) + \eta (\mathbf{M}^\circ - \mathbf{M}^{(l)}) (\mathbf{x} - \mathbf{x}^{(l)}) + \eta (\mathbf{M}^\circ - \mathbf{M}^{(l)}) \mathbf{x}^{(l)},
\end{aligned}$$

where $\mathbf{x}^{(l)}(\tau) = \mathbf{x}^{(l)} + \tau(\mathbf{x} - \mathbf{x}^{(l)})$. Then, we take inner product with $\mathbf{u}^{(l)}$ on both sides.

$$\begin{aligned}
\mathbf{u}^{(l)\top} (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)}) &= \underbrace{(1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{u}^{(l)\top} (\mathbf{x} - \mathbf{x}^{(l)}) - 2\eta \mathbf{u}^{(l)\top} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top (\mathbf{x} - \mathbf{x}^{(l)})}_{\textcircled{1}} \\
&\quad - \underbrace{\eta \int_0^1 \mathbf{u}^{(l)\top} \left(\nabla^2 g(\mathbf{x}(\tau)) - (\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top) \right) (\mathbf{x} - \mathbf{x}^{(l)}) d\tau}_{\textcircled{2}} \\
&\quad - \underbrace{\eta \mathbf{u}^{(l)\top} \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) - \mathcal{P}_l(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) \right) \mathbf{x}^{(l)}}_{\textcircled{3}} \\
&\quad + \underbrace{\eta \lambda^{(l)} \mathbf{u}^{(l)\top} (\mathbf{x} - \mathbf{x}^{(l)}) + \eta \mathbf{u}^{(l)\top} (\mathbf{M}^\circ - \mathbf{M}^{(l)}) (\mathbf{x} - \mathbf{x}^{(l)})}_{\textcircled{4}} \\
&\quad + \underbrace{\eta \mathbf{u}^{(l)\top} (\mathbf{M}^\circ - \mathbf{M}^{(l)}) \mathbf{x}^{(l)}}_{\textcircled{5}}
\end{aligned}$$

For the term $\textcircled{1}$, we have

$$\left| \mathbf{u}^{(l)\top} \tilde{\mathbf{x}} \right| \leq \left| \mathbf{u}^{(l)\top} \mathbf{x}^{(0)} \right| + (1 + \eta \lambda^*)^t \left| \mathbf{u}^{*\top} \mathbf{x}^{(0)} \right| \left| \mathbf{u}^{(l)\top} \mathbf{u}^* \right| \lesssim \sqrt{\frac{\log n}{n}} (1 + \eta \lambda^*)^t \beta_0$$

by Lemma A.5, and thus,

$$\begin{aligned}
|\textcircled{1}| &\leq \left| \mathbf{u}^{(l)\top} \tilde{\mathbf{x}} \right| \|\tilde{\mathbf{x}}\|_2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \sqrt{\frac{\log n}{n}} (1 + \eta \lambda^*)^t \beta_0^2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \\
&\lesssim \mu \sqrt{\frac{\log^3 n}{np}} \frac{\beta_0^3}{n} (1 + \eta \lambda^*)^t t \\
&\lesssim \lambda^* \sqrt{\frac{\mu}{np}} \frac{\beta_0}{n}.
\end{aligned}$$

The definition of Phase I was used to bound $(1 + \eta\lambda^*)^t t$ in deriving the last line. We use (D.13) to get

$$\begin{aligned} |\textcircled{2}| &\lesssim \int_0^1 \left\| \nabla^2 g(\mathbf{x}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\top \right) \right\| \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \|\mathbf{u}^{(l)}\|_2 d\tau \\ &\lesssim \sqrt{\frac{\mu^3 \log^7 n}{np}} \beta_0^2 \cdot \mu \sqrt{\frac{\log^2 n}{np} \frac{\beta_0}{\sqrt{n}}} t \lesssim \lambda^* \sqrt{\frac{\mu}{np \log^{15} n} \frac{\beta_0}{n}}. \end{aligned}$$

We apply Lemma G.12 to $\textcircled{3}$ to yield

$$\begin{aligned} |\textcircled{3}| &\lesssim \|\mathbf{x}^{(l)}\|_\infty^2 \sqrt{\frac{\log n}{p}} \left(\|\mathbf{u}^{(l)}\|_2 \|\mathbf{x}^{(l)}\|_\infty + \|\mathbf{x}^{(l)}\|_2 \|\mathbf{u}^{(l)}\|_\infty \right) \\ &\lesssim \sqrt{\frac{\mu \log^4 n}{np} \frac{\beta_0^3}{n}} \lesssim \lambda^* \sqrt{\frac{1}{\mu^3 n \log^{22} n} \frac{\beta_0}{n}}. \end{aligned}$$

We divide $\textcircled{4}$ into two terms that are related to sampling and noise, respectively.

$$\textcircled{4} = \mathbf{u}^{(l)\top} \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right) (\mathbf{x} - \mathbf{x}^{(l)}) + \mathbf{u}^{(l)\top} \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) - \mathbf{E}^{(l)} \right) (\mathbf{x} - \mathbf{x}^{(l)})$$

Then, Cauchy-Schwartz inequality is applied to yield

$$|\textcircled{4}| \leq \left\| \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right) \mathbf{u}^{(l)} \right\|_2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 + \left\| \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) - \mathbf{E}^{(l)} \right) \mathbf{u}^{(l)} \right\|_2 \|\mathbf{x} - \mathbf{x}^{(l)}\|_2.$$

Applying Lemmas G.11 and G.13 to the two terms, respectively, we get

$$|\textcircled{4}| \lesssim \lambda^* \sqrt{\frac{\mu^3 \log^2 n}{np} \frac{1}{\sqrt{n}}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \lambda^* \frac{\sqrt{\mu^5 \log^5 n} \beta_0}{np \frac{n}}{n}.$$

For the term $\textcircled{5}$, we decompose it into two terms as for $\textcircled{4}$.

$$|\textcircled{5}| \leq \left| \mathbf{u}^{(l)\top} \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right) \mathbf{x}^{(l)} \right| + \left| \mathbf{u}^{(l)\top} \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) - \mathbf{E}^{(l)} \right) \mathbf{x}^{(l)} \right|$$

Then, we apply Lemmas G.12 and G.14 to each term to obtain

$$|\textcircled{5}| \lesssim \lambda^* \sqrt{\frac{\mu^3 \log^3 n}{np} \frac{\beta_0}{n}}.$$

Combining all, there exists a universal constant $c_6 > 0$ such that

$$\eta (2|\textcircled{1}| + |\textcircled{2}| + |\textcircled{3}| + |\textcircled{4}| + |\textcircled{5}|) \leq \frac{\eta\lambda^*}{2} c_6 \sqrt{\frac{\mu^3 \log^3 n}{np} \frac{\beta_0}{n}}$$

if $n^2 p \gtrsim \mu^2 n \log^3 n$, and there exists a universal constant $c_7 > 0$ such that

$$\eta\lambda^{(l)} \leq \eta\lambda^* + \frac{c_7}{\log^2 n}$$

by (A.2) if $n^2 p \gtrsim \mu^2 n \log^5 n$. Finally, we have

$$\begin{aligned} \left| \mathbf{u}^{(l)\top} \left(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right) \right| &\leq \left(1 - \eta \|\tilde{\mathbf{x}}\|_2^2 + \eta\lambda^{(l)} \right) \left| \mathbf{u}^{(l)\top} (\mathbf{x} - \mathbf{x}^{(l)}) \right| + \frac{\eta\lambda^*}{2} c_6 \sqrt{\frac{\mu^3 \log^3 n}{np} \frac{\beta_0}{n}} \\ &\leq \left(1 + \eta\lambda^* + \frac{c_7}{\log^2 n} \right) \left| \mathbf{u}^{(l)\top} (\mathbf{x} - \mathbf{x}^{(l)}) \right| + \frac{\eta\lambda^*}{2} c_6 \sqrt{\frac{\mu^3 \log^3 n}{np} \frac{\beta_0}{n}}. \end{aligned}$$

An analysis on the recursive equation

$$x_{t+1} = \left(1 + \eta\lambda^* + \frac{c_7}{\log^2 n} \right) x_t + \frac{\eta\lambda^*}{2} c_6 \sqrt{\frac{\mu^3 \log^3 n}{np} \frac{\beta_0}{n}}, \quad x_0 = 0$$

gives the bound

$$\left| \mathbf{u}^{(l)\top} \left(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right) \right| \leq c_6 \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta\lambda^*)^{t+1} \frac{\beta_0}{n}.$$

(D.10) at $(t+1)$ We decompose $(\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l$ as

$$(\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l = \left(1 - \eta \|\tilde{\mathbf{x}}\|_2^2\right) (\mathbf{x}^{(l)} - \tilde{\mathbf{x}})_l + \eta \lambda^* \mathbf{u}^{*\top} (\mathbf{x}^{(l)} - \tilde{\mathbf{x}}) \mathbf{u}_l^* + \eta \left(\|\tilde{\mathbf{x}}\|_2^2 - \|\mathbf{x}^{(l)}\|_2^2\right) x_l^{(l)},$$

and this implies

$$\begin{aligned} \left|(\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l\right| &\leq \left(1 - \eta \|\tilde{\mathbf{x}}\|_2^2\right) \left|(\mathbf{x}^{(l)} - \tilde{\mathbf{x}})_l\right| \\ &\quad + \eta \lambda^* \|\mathbf{x}^{(l)} - \tilde{\mathbf{x}}\|_2 \|\mathbf{u}^*\|_\infty + \eta \|\tilde{\mathbf{x}}\|_2 \|\mathbf{x}^{(l)} - \tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty. \end{aligned}$$

From (D.5) and (D.8), we have

$$\|\mathbf{x}^{(l)} - \tilde{\mathbf{x}}\|_2 \leq \|\mathbf{x}^{(l)} - \mathbf{x}\|_2 + \|\mathbf{x} - \tilde{\mathbf{x}}\|_2 \leq 3c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 t.$$

If n is sufficiently large, we have

$$\eta \lambda^* \|\mathbf{u}^*\|_\infty + \eta \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty \leq \sqrt{\frac{\mu \log n}{n}}$$

for all $t \leq T_1$ because

$$\eta \lambda^* \|\mathbf{u}^*\|_\infty + \eta \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty \lesssim \eta \lambda^* \sqrt{\frac{\mu}{n}} + \eta \sqrt{\frac{\log^2 n}{n}} \beta_0^2 \lesssim \sqrt{\frac{\mu}{n}}.$$

Hence, we have

$$\left|(\mathbf{x}^{(s+1,l)} - \tilde{\mathbf{x}}^{(s+1)})_l\right| \leq \left|(\mathbf{x}^{(s,l)} - \tilde{\mathbf{x}}^{(s)})_l\right| + 3c_0 \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}}$$

for all $s \leq t$. Finally, we have

$$\begin{aligned} \left|(\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l\right| &\leq 3c_0 \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} \sum_{s=1}^t s \\ &\leq 3c_0 \sqrt{\frac{\mu^3 \log^2 n}{np}} \frac{\beta_0}{\sqrt{n}} (t+1)^2. \end{aligned}$$

(D.5) at $(t+1)$ We can obtain this through the combination of (D.7) and Lemma 5.2.

$$\begin{aligned} \left\|\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}\right\|_2 &\leq \left\|\mathbf{x}^{(t+1)} - \hat{\mathbf{x}}^{(t+1)}\right\|_2 + \left\|\hat{\mathbf{x}}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}\right\|_2 \\ &\leq c_2(3c_0 + c_5 + 1) \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta \lambda^*)^{t+1} \beta_0^3 + c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 (t+1) \\ &\leq c_2(3c_0 + c_5 + 1) \frac{1}{\lambda^*} \sqrt{\frac{\mu^3 \log^3 n}{np}} (1 + \eta \lambda^*)^{T_1} \beta_0^3 + c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 (t+1) \\ &\leq c_1 c_2 (3c_0 + c_5 + 1) \mu \sqrt{\frac{1}{np \log^2 n}} \beta_0 + c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 (t+1) \\ &\leq 2c_0 \mu \sqrt{\frac{\log n}{np}} \beta_0 (t+1) \end{aligned}$$

(D.6) at $(t + 1)$ The l th component of $\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}$ is bounded by

$$\begin{aligned}
\left| (\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| &\leq \left\| \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right\|_\infty + \left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| \\
&\leq \left\| \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right\|_2 + \left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| \\
&\leq c_5 \mu \sqrt{\frac{\log^2 n}{np} \frac{\beta_0}{\sqrt{n}}} (t+1) + 3c_0 \sqrt{\frac{\mu^3 \log^2 n}{np} \frac{\beta_0}{\sqrt{n}}} (t+1)^2 \\
&\leq (3c_0 + c_5) \sqrt{\frac{\mu^3 \log^2 n}{np} \frac{\beta_0}{\sqrt{n}}} (t+1)^2.
\end{aligned}$$

At $t = T_1$ Because $T_1 \lesssim \log n$, it is implied from Lemma D.2 that at $t = T_1$, there exists a constant $c_7 > 0$ such that

$$\begin{aligned}
\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 &\leq c_7 \mu \sqrt{\frac{\log^3 n}{np}} \beta_0, \\
\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty &\leq c_7 \sqrt{\frac{\mu^3 \log^8 n}{np} \frac{\beta_0}{\sqrt{n}}}, \\
\left\| \mathbf{x}^{(t)} - \mathbf{x}^{(t,l)} \right\|_2 &\leq c_7 \mu \sqrt{\frac{\log^4 n}{np} \frac{\beta_0}{\sqrt{n}}}, \\
\left| (\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l \right| &\leq c_7 \sqrt{\frac{\mu^3 \log^8 n}{np} \frac{\beta_0}{\sqrt{n}}}.
\end{aligned} \tag{D.14}$$

These bounds serve as a base case for the induction of the next part.

E Phase II

This section is mostly devoted to the proof of Lemma E.1 which is a formal version of Lemma 6.1.

Lemma E.1. *Suppose that (D.14) holds at $t = T_1$ and the initialization point $\mathbf{x}^{(0)}$ satisfies (B.1) to (B.6). Then, for all $T_1 < t \leq T_2$, we have*

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \leq 2c_7 \mu \sqrt{\frac{\log^3 n}{np}} \beta_0 (1 + \eta \lambda^*)^{t-T_1}, \tag{E.1}$$

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \leq c_{13} \sqrt{\frac{\mu^3 \log^8 n}{np} \frac{\beta_0}{\sqrt{n}}} (1 + \eta \lambda^*)^{t-T_1}, \tag{E.2}$$

$$\left\| \mathbf{x}^{(t)} - \mathbf{x}^{(t,l)} \right\|_2 \leq 2c_7 \mu \sqrt{\frac{\log^5 n}{np} \frac{\beta_0}{\sqrt{n}}} (1 + \eta \lambda^*)^{t-T_1}, \tag{E.3}$$

$$\left| (\mathbf{x}^{(t,l)} - \tilde{\mathbf{x}}^{(t)})_l \right| \leq 3c_7 \sqrt{\frac{\mu^3 \log^8 n}{np} \frac{\beta_0}{\sqrt{n}}} (1 + \eta \lambda^*)^{t-T_1}, \tag{E.4}$$

with high probability, where T_2 is the largest t such that $\tilde{\beta}_t^2 \leq \lambda^* \left(1 - \frac{1}{\log n}\right)$, and $c_{13} > 0$ is a constant.

Proof of Theorems 3.1 and 3.3 We first explain how Theorems 3.1 and 3.3 are derived from Lemmas D.2 and E.1. We first focus on $\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2$. For $t \leq T_1$, from Lemma D.2 and (C.3), we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \mu \sqrt{\frac{\log^3 n}{np}} \beta_0 \lesssim \frac{1}{\log n} \beta_0 \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2$$

provided that $n^2 p \gtrsim \mu^2 n \log^7 n$. For $T_1 < t \leq T_2$, from the definition of T_1 and Lemma E.1, we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \sqrt{\frac{1}{\log^{18} n} \frac{\beta_0}{\sqrt{n}}} (1 + \eta \lambda^*)^t.$$

From the lower bound of Lemma C.1, for all $t \leq T'_2$, we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \sqrt{\frac{1}{\log^{18} n} \left(1 + \frac{(1 + \eta \lambda^*)^t}{\sqrt{n}} \right)} \beta_0 \lesssim \sqrt{\frac{1}{\log^{17} n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2.$$

Now, for $T'_2 < t \leq T_2$, we have

$$\begin{aligned} \left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 &\lesssim \sqrt{\frac{1}{\log^{18} n} \left(1 + \frac{(1 + \eta \lambda^*)^{T'_2}}{\sqrt{n}} \right)} \beta_0 (1 + \eta \lambda^*)^{t - T'_2} \\ &\lesssim \sqrt{\frac{1}{\log^{17} n}} \left\| \tilde{\mathbf{x}}^{(T'_2)} \right\|_2 (1 + \eta \lambda^*)^{t - T'_2}. \end{aligned}$$

For any $T'_2 < t \leq T_2$, it is Lemma C.2 implied from Lemma C.2 that $(1 + \eta \lambda^*)^{t - T'_2} \leq \log^6 n$, and we have $\left\| \tilde{\mathbf{x}}^{(T'_2)} \right\|_2 \leq \left\| \tilde{\mathbf{x}} \right\|_2$. Hence, we get

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \sqrt{\frac{1}{\log^{17} n}} \sqrt{\log^{12} n} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2 \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2, \quad (\text{E.5})$$

and the proof for (11) of Theorem 3.3 is completed. If we combine this with (C.7), we are able to prove (4) of Theorem 3.1.

We move on to $\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty$. For $t \leq T_1$, from Lemma D.2 and (C.4), we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \sqrt{\frac{\mu^3 \log^8 n}{np} \frac{\beta_0}{\sqrt{n}}} \lesssim \frac{1}{\log n} \frac{\beta_0}{\sqrt{n}} \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty$$

provided that $n^2 p \gtrsim \mu^3 n \log^{10} n$. For $T_1 < t \leq T'_2$, from the definition of T_1 and Lemma E.1, we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \sqrt{\frac{1}{\log^{13} n} \frac{\beta_0}{\sqrt{n}}} (1 + \eta \lambda^*)^t.$$

From the lower bound of Lemma C.1, for all $t \leq T'_2$, we have

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \sqrt{\frac{1}{\log^{13} n} \left(1 + \frac{(1 + \eta \lambda^*)^t}{\sqrt{n}} \right)} \frac{\beta_0}{\sqrt{n}} \lesssim \sqrt{\frac{1}{\log^{12} n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty.$$

Now, for $T'_2 < t \leq T_2$, if we do the same as before, we get

$$\left\| \mathbf{x}^{(t)} - \tilde{\mathbf{x}}^{(t)} \right\|_\infty \lesssim \sqrt{\frac{1}{\log^{13} n}} \sqrt{\log^{12} n} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_2 \frac{1}{\sqrt{n}} \lesssim \frac{1}{\sqrt{\log n}} \left\| \tilde{\mathbf{x}}^{(t)} \right\|_\infty,$$

and the proof for (12) of Theorem 3.3 is completed. If we combine this with (C.8), we are able to prove (5) of Theorem 3.1.

Going through a similar way with (22) and (24), we can complete the proof of Theorems 3.1 and 3.3.

Proof of Lemma E.1 Before we start the proof, we define a function G as

$$G(\mathbf{x}) = \frac{1}{4} \left\| \mathbf{x} \mathbf{x}^\top \right\|_{\text{F}}^2.$$

The gradient of G satisfies

$$\nabla F(\mathbf{x}) = \nabla G(\mathbf{x}) - \mathbf{M}^* \mathbf{x}.$$

Now, we assume that the hypotheses hold up to the t th iteration and show that they hold at the $(t + 1)$ st iteration. For brevity, we drop the superscript (t) from $\mathbf{x}^{(t)}$, $\mathbf{x}^{(t,l)}$, $\tilde{\mathbf{x}}^{(t)}$ and denote them as \mathbf{x} , $\mathbf{x}^{(l)}$, $\tilde{\mathbf{x}}$, respectively.

(E.1) at $(t+1)$ We decompose $\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}$ as

$$\begin{aligned}
& \mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)} \\
&= (\mathbf{x} - \eta \nabla f(\mathbf{x})) - (\tilde{\mathbf{x}} - \eta \nabla F(\tilde{\mathbf{x}})) \\
&= (\mathbf{x} - \eta \nabla g(\mathbf{x})) - (\tilde{\mathbf{x}} - \eta \nabla G(\tilde{\mathbf{x}})) + \eta (\mathbf{M}^\circ \mathbf{x} - \mathbf{M}^* \tilde{\mathbf{x}}) \\
&= (\mathbf{x} - \eta \nabla g(\mathbf{x})) - (\tilde{\mathbf{x}} - \eta \nabla g(\tilde{\mathbf{x}})) - \eta (\nabla g(\tilde{\mathbf{x}}) - \nabla G(\tilde{\mathbf{x}})) + \eta \mathbf{M}^* (\mathbf{x} - \tilde{\mathbf{x}}) + \eta (\mathbf{M}^\circ - \mathbf{M}^*) \mathbf{x} \\
&= \int_0^1 (\mathbf{I} - \eta \nabla^2 g(\mathbf{x}(\tau))) (\mathbf{x} - \tilde{\mathbf{x}}) d\tau - \eta \|\tilde{\mathbf{x}}\|_2^2 (\mathbf{I}_{\tilde{\mathbf{x}}} - \mathbf{I}) \tilde{\mathbf{x}} + \eta \mathbf{M}^* (\mathbf{x} - \tilde{\mathbf{x}}) + \eta (\mathbf{M}^\circ - \mathbf{M}^*) \mathbf{x} \\
&= \underbrace{\left((1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{I} - 2\eta \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top + \eta \mathbf{M}^* \right)}_{\textcircled{1}} (\mathbf{x} - \tilde{\mathbf{x}}) - \underbrace{\eta \int_0^1 \left(\nabla^2 g(\mathbf{x}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \right) \right)}_{\textcircled{2}} (\mathbf{x} - \tilde{\mathbf{x}}) d\tau \\
&\quad - \underbrace{\eta \|\tilde{\mathbf{x}}\|_2^2 (\mathbf{I}_{\tilde{\mathbf{x}}} - \mathbf{I}) \tilde{\mathbf{x}}}_{\textcircled{3}} + \underbrace{\eta (\mathbf{M}^\circ - \mathbf{M}^*) \mathbf{x}}_{\textcircled{4}},
\end{aligned}$$

where $\mathbf{x}(\tau) = \tilde{\mathbf{x}} + \tau(\mathbf{x} - \tilde{\mathbf{x}})$. For the term $\textcircled{1}$, we require a bound on

$$\left\| (1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{I} - 2\eta \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top + \eta \mathbf{M}^* \right\|.$$

If we write $\tilde{\mathbf{x}} = \tilde{\alpha}_t \mathbf{u}^* + \tilde{\mathbf{x}}_\perp$, we have $\tilde{\alpha}_t^2 \leq \lambda^*$ and $\|\tilde{\mathbf{x}}_\perp\|_2 \lesssim \beta_0$. Then, we have

$$\begin{aligned}
& \left\| (1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{I} - 2\eta \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top + \eta \mathbf{M}^* \right\| \\
&= \left\| (1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{I} + \eta (\lambda^* - 2\tilde{\alpha}_t^2) \mathbf{u}^* \mathbf{u}^{*\top} - 2\eta (\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top - \tilde{\alpha}_t^2 \mathbf{u}^* \mathbf{u}^{*\top}) \right\| \\
&\leq \left\| (1 - \eta \|\tilde{\mathbf{x}}\|_2^2) \mathbf{I} + \eta (\lambda^* - 2\tilde{\alpha}_t^2) \mathbf{u}^* \mathbf{u}^{*\top} \right\| + 2\eta \|\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top - \tilde{\alpha}_t^2 \mathbf{u}^* \mathbf{u}^{*\top}\| \\
&\leq (1 + \eta \lambda^*) + 2\eta (2\alpha_t \|\tilde{\mathbf{x}}_\perp\|_2 + \|\tilde{\mathbf{x}}_\perp\|_2^2) \\
&\leq 1 + \eta \lambda^* + \frac{c_8}{\log^2 n}
\end{aligned}$$

for some universal constant $c_8 > 0$. This implies the desired bound

$$\|\textcircled{1}\|_2 \leq \left(1 + \eta \lambda^* + \frac{c_8}{\log^2 n} \right) \|\mathbf{x} - \tilde{\mathbf{x}}\|_2.$$

For all $0 \leq \tau \leq 1$, we have $\|\mathbf{x}(\tau) - \tilde{\mathbf{x}}\|_\infty \leq \|\mathbf{x} - \tilde{\mathbf{x}}\|_\infty$, and the induction hypothesis (E.2) gives

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_\infty \lesssim \sqrt{\frac{1}{\mu \log^{13} n}} (1 + \eta \lambda^*)^{T_2} \frac{\beta_0}{n}.$$

Hence, by Lemma G.10, we have

$$\begin{aligned}
\left\| \nabla^2 g(\mathbf{x}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \right) \right\| &\lesssim n \|\mathbf{x} - \tilde{\mathbf{x}}\|_\infty (\|\mathbf{x}\|_\infty + \|\tilde{\mathbf{x}}\|_\infty) + \sqrt{\frac{n \log n}{p}} \|\tilde{\mathbf{x}}\|_\infty^2 \\
&\lesssim \left(\sqrt{\frac{1}{\log^{12} n}} + \sqrt{\frac{\mu^2 \log^3 n}{np}} \right) (1 + \eta \lambda^*)^{2T_2} \frac{\beta_0^2}{n} \\
&\lesssim \lambda^* \sqrt{\frac{1}{\log^{12} n}}
\end{aligned}$$

if $n^2 p \gtrsim \mu^2 n \log^{15} n$ because $\|\tilde{\mathbf{x}}\|_\infty \lesssim \sqrt{\mu \log n} (1 + \eta \lambda^*)^{T_2} \frac{\beta_0}{n}$ and $(1 + \eta \lambda^*)^{T_2} \lesssim \sqrt{\lambda^*} \frac{\sqrt{n}}{\beta_0}$. This gives

$$\|\textcircled{2}\|_2 \lesssim \lambda^* \sqrt{\frac{1}{\log^{12} n}} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2. \quad (\text{E.6})$$

For the term ③, we use Lemma G.8 to obtain

$$\|\textcircled{3}\|_2 \lesssim \lambda^* \sqrt{\frac{\mu \log n}{np}} \|\tilde{\mathbf{x}}\|_2.$$

Lastly, the term ④ is bounded with

$$\|\textcircled{4}\|_2 \lesssim \|\mathbf{M}^\circ - \mathbf{M}^*\| \|\mathbf{x}\|_2 \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \|\tilde{\mathbf{x}}\|_2.$$

Combining all, there exists a universal constant $c_9 > 0$ such that

$$\begin{aligned} \|\textcircled{1}\|_2 + \|\textcircled{2}\|_2 &\leq \left(1 + \eta\lambda^* + \frac{c_9}{\log^2 n}\right) \|\mathbf{x} - \tilde{\mathbf{x}}\|_2, \\ \eta (\|\textcircled{3}\|_2 + \|\textcircled{4}\|_2) &\leq c_9 \mu \sqrt{\frac{\log n}{np}} \|\tilde{\mathbf{x}}\|_2. \end{aligned}$$

Because $\|\mathbf{x}^{(T_1)}\|_2 \lesssim \sqrt{\log n} \beta_0$ by (C.3) and $\|\tilde{\mathbf{x}}^{(t)}\|_2$ can grow at a rate at most $(1 + \eta\lambda^*)$, there exists a universal constant $c_{10} > 0$ such that

$$c_8 \mu \sqrt{\frac{\log n}{np}} \|\tilde{\mathbf{x}}^{(t)}\|_2 \leq c_{10} \mu \sqrt{\frac{\log^2 n}{np}} (1 + \eta\lambda^*)^{t-T_1} \beta_0. \quad (\text{E.7})$$

Hence, for all $T_1 \leq s \leq t$, we have

$$\|\mathbf{x}^{(s+1)} - \tilde{\mathbf{x}}^{(s+1)}\|_2 \leq \left(1 + \eta\lambda^* + \frac{c_9}{\log^2 n}\right) \|\mathbf{x}^{(s)} - \tilde{\mathbf{x}}^{(s)}\|_2 + c_{10} \mu \sqrt{\frac{\log^2 n}{np}} (1 + \eta\lambda^*)^{t-T_1} \beta_0.$$

An analysis on the recursive equation

$$x_{s+1} = \left(1 + \eta\lambda^* + \frac{c_9}{\log^2 n}\right) x_s + c_{10} \mu \sqrt{\frac{\log^2 n}{np}} (1 + \eta\lambda^*)^{t-T_1} \beta_0, \quad x_{T_1} = c_7 \mu \sqrt{\frac{\log^3 n}{np}} \beta_0$$

proves that

$$\|\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}\|_2 \leq 2c_7 \mu \sqrt{\frac{\log^3 n}{np}} \beta_0 (1 + \eta\lambda^*)^{t+1-T_1}.$$

(E.3) at $(t+1)$ Similar to the proof of (D.8), we have the decomposition

$$\begin{aligned} \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} &= \underbrace{(1 - \eta \|\tilde{\mathbf{x}}\|_2^2 - 2\eta \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top)(\mathbf{x} - \mathbf{x}^{(l)})}_{\textcircled{1}} \\ &\quad - \underbrace{\eta \int_0^1 \left(\nabla^2 g(\mathbf{x}^{(l)}(\tau)) - \left(\|\tilde{\mathbf{x}}\|_2^2 \mathbf{I} + 2\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \right) \right) (\mathbf{x} - \mathbf{x}^{(l)}) d\tau}_{\textcircled{2}} \\ &\quad - \underbrace{\eta \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) - \mathcal{P}_l(\mathbf{x}^{(l)} \mathbf{x}^{(l)\top}) \right) \mathbf{x}^{(l)}}_{\textcircled{3}} \\ &\quad + \underbrace{\eta \mathbf{M}^*(\mathbf{x} - \mathbf{x}^{(l)}) + \eta (\mathbf{M}^\circ - \mathbf{M}^*)(\mathbf{x} - \mathbf{x}^{(l)})}_{\textcircled{4}} \\ &\quad + \underbrace{\eta \left(\frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{M}^*) - \mathcal{P}_{\Omega}^{(l)}(\mathbf{M}^*) \right) \mathbf{x}^{(l)}}_{\textcircled{5}} + \underbrace{\eta \left(\frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) - \mathbf{E}^{(l)} \right) \mathbf{x}^{(l)}}_{\textcircled{6}}, \end{aligned}$$

where $\mathbf{x}^{(l)}(\tau) = \mathbf{x}^{(l)} + \tau(\mathbf{x} - \mathbf{x}^{(l)})$. Both of the terms ① and ② can be bounded similar to ① and ② of $\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}$ as

$$\begin{aligned}\|\textcircled{1}\|_2 &\leq \left(1 + \eta\lambda^* + \frac{c_{11}}{\log^2 n}\right) \|\mathbf{x} - \mathbf{x}^{(l)}\|_2, \\ \|\textcircled{2}\|_2 &\lesssim \lambda^* \sqrt{\frac{1}{\log^{12} n}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2\end{aligned}$$

for some universal constant $c_{11} > 0$. For the terms ③ and ⑤, we use Lemma G.11 to obtain

$$\begin{aligned}\|\textcircled{3}\|_2 &\lesssim \sqrt{\frac{\log n}{p}} \|\mathbf{x}^{(l)}\|_2 \|\mathbf{x}^{(l)}\|_\infty^2 \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \frac{1}{\sqrt{n}} \|\tilde{\mathbf{x}}\|_2, \\ \|\textcircled{5}\|_2 &\lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \frac{1}{\sqrt{n}} \|\tilde{\mathbf{x}}\|_2,\end{aligned}$$

and use Lemma G.13 to obtain

$$\|\textcircled{6}\|_2 \lesssim \lambda^* \mu \sqrt{\frac{\log^2 n}{np}} \frac{1}{\sqrt{n}} \|\tilde{\mathbf{x}}\|_2.$$

From Lemmas A.1 and A.3, the term ④ is bounded as

$$\|\textcircled{4}\|_2 \leq \|\mathbf{M}^\circ - \mathbf{M}^*\| \|\mathbf{x} - \mathbf{x}^{(l)}\|_2 \lesssim \lambda^* \mu \sqrt{\frac{\log n}{np}} \|\mathbf{x} - \mathbf{x}^{(l)}\|_2.$$

Combining all with (E.7), there exists a universal constant $c_{12} > 0$ such that

$$\begin{aligned}\|\textcircled{1}\|_2 + \|\textcircled{2}\|_2 + \|\textcircled{4}\|_2 &\leq \left(1 + \eta\lambda^* + \frac{c_{12}}{\log^2 n}\right) \|\mathbf{x} - \tilde{\mathbf{x}}\|_2, \\ \eta (\|\textcircled{3}\|_2 + \|\textcircled{5}\|_2 + \|\textcircled{6}\|_2) &\leq c_{12} \mu \sqrt{\frac{\log^3 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}.\end{aligned}$$

Hence, we have

$$\|\mathbf{x}^{(s+1)} - \mathbf{x}^{(s+1,l)}\|_2 \leq \left(1 + \eta\lambda^* + \frac{c_{12}}{\log^2 n}\right) \|\mathbf{x}^{(s)} - \mathbf{x}^{(s,l)}\|_2 + c_{12} \mu \sqrt{\frac{\log^3 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1},$$

for all $T_1 \leq s \leq t$. An analysis on the recursive equation

$$x_{s+1} = \left(1 + \eta\lambda^* + \frac{c_{12}}{\log^2 n}\right) x_s + c_{12} \mu \sqrt{\frac{\log^3 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}, \quad x_{T_1} = c_7 \mu \sqrt{\frac{\log^4 n}{np}} \frac{\beta_0}{\sqrt{n}}$$

proves that

$$\begin{aligned}\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)}\|_2 &\leq 2 \left(c_{12} \mu \sqrt{\frac{\log^3 n}{np}} (t+1-T_1) + c_7 \mu \sqrt{\frac{\log^4 n}{np}} \right) \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1} \\ &\leq c_{13} \mu \sqrt{\frac{\log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1}\end{aligned}$$

holds for some universal constant $c_{13} > 0$ because $T_2 - T_1 \lesssim \log n$.

(E.4) at $(t+1)$ We use the same bound

$$\left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| \leq \left(1 - \eta \|\tilde{\mathbf{x}}\|_2^2\right) \left| (\mathbf{x}^{(l)} - \tilde{\mathbf{x}})_l \right| + \eta (\lambda^* \|\mathbf{u}^*\|_\infty + \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty) \|\mathbf{x}^{(l)} - \tilde{\mathbf{x}}\|_2.$$

that was used in the proof of (D.10). From (E.1) and (E.3), we have

$$\left\| \mathbf{x}^{(l)} - \tilde{\mathbf{x}} \right\|_2 \leq \left\| \mathbf{x}^{(l)} - \mathbf{x} \right\|_2 + \left\| \mathbf{x} - \tilde{\mathbf{x}} \right\|_2 \leq 3c_7\mu \sqrt{\frac{\log^5 n}{np}} \beta_0 (1 + \eta\lambda^*)^{t-T_1}.$$

Combined with the fact that $\lambda^* \|\mathbf{u}^*\|_\infty + \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty \leq 3\lambda^* \sqrt{\frac{\mu}{n}}$, there exists a universal constant $c_{14} > 0$ such that

$$\eta(\lambda^* \|\mathbf{u}^*\|_\infty + \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{x}}\|_\infty) \left\| \mathbf{x}^{(l)} - \tilde{\mathbf{x}} \right\|_2 \leq c_{14} \sqrt{\frac{\mu^3 \log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1}.$$

Hence, for all $T_1 \leq s \leq t$, we have

$$\left| (\mathbf{x}^{(s+1,l)} - \tilde{\mathbf{x}}^{(s+1)})_l \right| \leq \left| (\mathbf{x}^{(s,l)} - \tilde{\mathbf{x}}^{(s)})_l \right| + c_{14} \sqrt{\frac{\mu^3 \log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t-T_1},$$

and this implies

$$\begin{aligned} \left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| &\leq c_{14} \sqrt{\frac{\mu^3 \log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} \sum_{s=T_1}^t (1 + \eta\lambda^*)^{s-T_1} + c_7 \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} \\ &\leq 2c_7 \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1}. \end{aligned}$$

(E.2) at $(t+1)$ The l th component of $\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)}$ is bounded by

$$\begin{aligned} \left| (\mathbf{x}^{(t+1)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| &\leq \left\| \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right\|_\infty + \left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| \\ &\leq \left\| \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1,l)} \right\|_2 + \left| (\mathbf{x}^{(t+1,l)} - \tilde{\mathbf{x}}^{(t+1)})_l \right| \\ &\leq c_{13}\mu \sqrt{\frac{\log^5 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1} + 2c_7 \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1} \\ &\leq 3c_7 \sqrt{\frac{\mu^3 \log^8 n}{np}} \frac{\beta_0}{\sqrt{n}} (1 + \eta\lambda^*)^{t+1-T_1}. \end{aligned}$$

F Fixed Initialization Size

In Section 3, we claimed that the estimation error is improved to $\frac{1}{\sqrt{np}} + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}}$ if the initialization size is fixed to $n^{-1/4}$ regardless of the sample complexity. We briefly discuss how the proofs should change in such a case. For clear presentation, μ and $\log n$ factors are ignored in this section.

For every bound of Phase I (Lemmas 5.1 to 5.5), $\frac{1}{\sqrt{np}}$ is changed to $\frac{1}{\sqrt{np}} + \frac{\sigma}{\lambda^*} \sqrt{\frac{n}{p}}$, while allowing σ to be as large as $\frac{\lambda^* \mu}{n} \sqrt{np}$. More importantly, the definition of Phase I is changed to be the largest t such that $(1 + \eta\lambda^*)^t \leq \sqrt{n}$, so it is lengthened by \sqrt{np} times than before. In the original proof, the estimation error of $\frac{1}{\sqrt{np}}$ obtained at the end of Phase I was increased to $\frac{1}{\text{poly}(\log n)}$ during the first part of Phase II. However, if $(1 + \eta\lambda^*)^t$ equals \sqrt{n} at the end of Phase I, we do not have such a part in Phase II, and the estimation error obtained at the end of Phase I is maintained through Phase II.

G Technical Lemmas

We introduce some technical lemmas in this section. Most of them are the results of classical concentration inequalities.

Theorem G.1 (Matrix Bernstein Inequality). *Let $\{\mathbf{X}_i\}$ be $n \times n$ independent symmetric random matrices. Assume that each random matrix satisfies $\mathbb{E} \mathbf{X}_i = \mathbf{0}$ and $\|\mathbf{X}_i\| \leq L$ almost surely. Then, for all $\tau \geq 0$, we have*

$$\mathbb{P} \left[\left\| \sum_i \mathbf{X}_i \right\| \geq \tau \right] \leq n \exp \left(\frac{-\tau^2/2}{V + L\tau/3} \right),$$

where $V = \|\sum_i \mathbb{E}(\mathbf{X}_i^2)\|$.

Corollary G.2 (Matrix Bernstein Inequality). *Let $\{\mathbf{X}_i\}$ be $n \times n$ independent symmetric random matrices. Assume that each random matrix satisfies $\mathbb{E} \mathbf{X}_i = \mathbf{0}$ and $\|\mathbf{X}_i\| \leq L$ almost surely. Then, with high probability, we have*

$$\left\| \sum_i \mathbf{X}_i \right\| \lesssim \sqrt{V \log n} + L \log n,$$

where $V = \|\sum_i \mathbb{E}(\mathbf{X}_i^2)\|$.

Lemma G.3. *For any fixed matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M} \right\| \lesssim \sqrt{\frac{n \log n}{p}} \|\mathbf{M}\|_\infty + \frac{\log n}{p} \|\mathbf{M}\|_\infty$$

with high probability.

Proof. We decompose the matrix into the sum of independent symmetric matrices.

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}) - \mathbf{M} = \sum_{i < j} \left(\frac{\delta_{ij}}{p} - 1 \right) M_{ij} (\mathbf{e}_i \mathbf{e}_j^\top + \mathbf{e}_j \mathbf{e}_i^\top) + \sum_i \left(\frac{\delta_{ii}}{p} - 1 \right) M_{ii} \mathbf{e}_i \mathbf{e}_i^\top$$

We calculate L and V of Corollary G.2. We have $L \leq \frac{1}{p} \|\mathbf{M}\|_\infty$ because

$$\begin{aligned} \left\| \left(\frac{\delta_{ij}}{p} - 1 \right) M_{ij} (\mathbf{e}_i \mathbf{e}_j^\top + \mathbf{e}_j \mathbf{e}_i^\top) \right\| &\leq \frac{1}{p} \|\mathbf{M}\|_\infty, \\ \left\| \left(\frac{\delta_{ii}}{p} - 1 \right) M_{ii} \mathbf{e}_i \mathbf{e}_i^\top \right\| &\leq \frac{1}{p} \|\mathbf{M}\|_\infty. \end{aligned}$$

We also have the following bound on V .

$$V = \frac{1-p}{p} \left\| \sum_{i,j} M_{ij}^* \mathbf{e}_i \mathbf{e}_i^\top \right\| \leq \frac{n}{p} \|\mathbf{M}\|_\infty^2$$

Hence, Corollary G.2 implies the desired result. \square

We can prove Lemma A.1 by applying Lemma G.3 to \mathbf{M}^* and using $\|\mathbf{M}^*\|_\infty = \lambda^* \frac{\mu}{n}$.

We introduce classical Bernstein inequality and the results obtained from it.

Theorem G.4 (Bernstein Inequality). *Let $\{X_i\}$ be independent random variables. Assume that each random variable satisfies $\mathbb{E} X_i = 0$ and $|X_i| \leq L$ almost surely. Then, for all $\tau \geq 0$, we have*

$$\mathbb{P} \left[\left| \sum_i X_i \right| \geq \tau \right] \leq 2 \exp \left(\frac{-\tau^2/2}{V + L\tau/3} \right),$$

where $V = \sum_i \mathbb{E}[X_i^2]$.

Corollary G.5 (Bernstein Inequality). *Let $\{X_i\}$ be independent random variables. Assume that each random variable satisfies $\mathbb{E} X_i = 0$ and $|X_i| \leq L$ almost surely. Then, with high probability, we have*

$$\left| \sum_i X_i \right| \lesssim \sqrt{V \log n} + L \log n,$$

where $V = \sum_i \mathbb{E}[X_i^2]$.

Lemma G.6. Let and $\{X_i\}$ be independent Bernoulli random variables with expectation p . Then, for any fixed vector \mathbf{a} , we have

$$\left| \sum_i \left(\frac{X_i}{p} - 1 \right) a_i \right| \lesssim \sqrt{\frac{\log n}{p}} \|\mathbf{a}\|_2 + \frac{\log n}{p} \|\mathbf{a}\|_\infty$$

with high probability.

Proof. We can apply Corollary G.5 with $L = \frac{1}{p} \|\mathbf{a}\|_\infty$ and $V = \frac{1-p}{p} \|\mathbf{a}\|_2^2$. \square

Lemma G.7. If $n^2 p \gtrsim \mu n \log n$, we have

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) \right\|_{2,\infty} \lesssim \lambda^* \sqrt{\frac{\mu}{np}}$$

with high probability.

Proof. Let us consider ℓ_2 -norm of the i th row of \mathbf{M}° .

$$\begin{aligned} \left\| \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) \right)_{i*} \right\|_2^2 &= \lambda^{*2} u_i^{*2} \sum_j \frac{1}{p^2} \delta_{ij} u_j^{*2} \\ &\leq \frac{1}{p} \lambda^{*2} \|\mathbf{u}^*\|_\infty^2 \left(\|\mathbf{u}^*\|_2^2 + \left(\sum_j \frac{1}{p} \delta_{ij} u_j^{*2} - \|\mathbf{u}^*\|_2^2 \right) \right) \\ &\lesssim \frac{\lambda^{*2} \mu}{np} \left(1 + \sqrt{\frac{\log n}{np}} \right) \lesssim \frac{\lambda^{*2} \mu}{np} \end{aligned}$$

The third line follows from Lemma G.6. \square

Proof of Lemma A.2. The spectral norm of a symmetric matrix that has nonzero entries only on the l th row/column is bounded by twice of the norm of its l th row. Hence,

$$\begin{aligned} \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right\| &\leq 2 \left\| \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathcal{P}_\Omega^{(l)}(\mathbf{M}^*) \right)_{l*} \right\|_2 = 2 \left\| \left(\frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) - \mathbf{M}^* \right)_{l*} \right\|_2 \\ &\lesssim \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{M}^*) \right\|_{2,\infty} + \|\mathbf{M}^*\|_{2,\infty} \lesssim \lambda^* \sqrt{\frac{\mu}{np}}, \end{aligned}$$

where the last inequality follows from Lemma G.7. \square

Lemma G.8. Let \mathbf{y} be a vector that is independent from the sampling. Then, if $n^2 p \gtrsim n \log n$, we have

$$\max_{i \in [n]} \left| \|\mathbf{x}\|_{2,i}^2 - \|\mathbf{y}\|_2^2 \right| \lesssim n \|\mathbf{x} - \mathbf{y}\|_\infty (\|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty) + \sqrt{\frac{\log n}{p}} \|\mathbf{y}\|_2 \|\mathbf{y}\|_\infty + \frac{\log n}{p} \|\mathbf{y}\|_\infty^2$$

with very high probability.

Proof. Let us fix i and decompose the difference as

$$\|\mathbf{x}\|_{2,i}^2 - \|\mathbf{y}\|_2^2 = \frac{1}{p} \sum_{j=1}^n \delta_{ij} (x_j^2 - y_j^2) + \sum_{j=1}^n \left(\frac{\delta_{ij}}{p} - 1 \right) y_j^2.$$

The first term is bounded as

$$\left| \frac{1}{p} \sum_{j=1}^n \delta_{ij} (x_j^2 - y_j^2) \right| \leq \|\mathbf{x} - \mathbf{y}\|_\infty \|\mathbf{x} + \mathbf{y}\|_\infty \frac{1}{p} \sum_{j=1}^n \delta_{ij} \lesssim n \|\mathbf{x} - \mathbf{y}\|_\infty \|\mathbf{x} + \mathbf{y}\|_\infty,$$

and the second term is bounded as

$$\left| \sum_{j=1}^n \left(\frac{\delta_{ij}}{p} - 1 \right) y_j^2 \right| \lesssim \sqrt{\frac{\log n}{p}} \|\mathbf{y}\|_2 \|\mathbf{y}\|_\infty + \frac{\log n}{p} \|\mathbf{y}\|_\infty^2$$

by Lemma G.6. \square

Lemma G.9. *Let \mathbf{y} be a vector that is independent from the sampling. Then, if $n^2 p \gtrsim n \log n$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{x}\mathbf{x}^\top) - \mathbf{y}\mathbf{y}^\top \right\| \lesssim n \|\mathbf{x} - \mathbf{y}\|_\infty (\|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty) + \sqrt{\frac{n \log n}{p}} \|\mathbf{y}\|_\infty^2$$

with very high probability.

Proof. We have the following sequence of inequalities

$$\begin{aligned} \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{x}\mathbf{x}^\top) - \mathbf{y}\mathbf{y}^\top \right\| &\leq \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{x}\mathbf{x}^\top) - \frac{1}{p} \mathcal{P}_\Omega(\mathbf{y}\mathbf{y}^\top) \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{y}\mathbf{y}^\top) - \mathbf{y}\mathbf{y}^\top \right\| \\ &\leq \|\mathbf{x}\mathbf{x}^\top - \mathbf{y}\mathbf{y}^\top\|_\infty \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{1}\mathbf{1}^\top) \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{y}\mathbf{y}^\top) - \mathbf{y}\mathbf{y}^\top \right\| \\ &\lesssim \|\mathbf{x} - \mathbf{y}\|_\infty (\|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty) \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{1}\mathbf{1}^\top) \right\| + \left\| \frac{1}{p} \mathcal{P}_\Omega(\mathbf{y}\mathbf{y}^\top) - \mathbf{y}\mathbf{y}^\top \right\| \\ &\lesssim n \|\mathbf{x} - \mathbf{y}\|_\infty (\|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty) + \sqrt{\frac{n \log n}{p}} \|\mathbf{y}\|_\infty^2, \end{aligned}$$

where the second line is derived from a basic inequality $\|\mathbf{A}\| \leq \|\mathbf{A}\|_\infty$ that holds for any matrix \mathbf{A} , and the last line follows by applying Lemma G.3 to $\mathbf{1}\mathbf{1}^\top$ and $\mathbf{y}\mathbf{y}^\top$. \square

Lemma G.10. *Let \mathbf{y} be a vector that is independent from the sampling. Then, if $n^2 p \gtrsim n \log n$, we have*

$$\begin{aligned} \left\| \nabla^2 g(\mathbf{x}) - \left(\|\mathbf{y}\|_2^2 \mathbf{I} + 2\mathbf{y}\mathbf{y}^\top \right) \right\| &\lesssim n \|\mathbf{x} - \mathbf{y}\|_\infty (\|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty) \\ &\quad + \sqrt{\frac{\log n}{p}} \|\mathbf{y}\|_2 \|\mathbf{y}\|_\infty + \frac{\log n}{p} \|\mathbf{y}\|_\infty^2 + \sqrt{\frac{n \log n}{p}} \|\mathbf{y}\|_\infty^2 \end{aligned}$$

Proof. This follows directly from Lemmas G.8 and G.9. \square

Let us define an operator \mathcal{P}_{Ω_l} such that an entry of $\mathcal{P}_{\Omega_l}(\mathbf{X})$ is equal to that of \mathbf{X} if it is contained both in the l th row/column and Ω , and otherwise 0. We also define an operator \mathcal{P}_l that makes the entries outside the l th row/column zero. Then, we have

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{X}) - \mathcal{P}_\Omega^{(l)}(\mathbf{X}) = \frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{X}) - \mathcal{P}_l(\mathbf{X}).$$

Also, note that

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) - \mathbf{E}^{(l)} = \frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{E}).$$

The following lemma was also introduced in [14], but we include the proof for completeness.

Lemma G.11. *Suppose that a matrix \mathbf{M} and a vector \mathbf{v} are independent from sampling of the l th row/column. If $n^2 p \gtrsim n \log n$, we have*

$$\left\| \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{M}) - \mathcal{P}_l(\mathbf{M}) \right) \mathbf{v} \right\|_2 \lesssim \|\mathbf{M}\|_\infty \left(\sqrt{\frac{\log n}{p}} \|\mathbf{v}\|_2 + \frac{\log n}{p} \|\mathbf{v}\|_\infty + \sqrt{\frac{n}{p}} \|\mathbf{v}\|_\infty \right)$$

with high probability.

Proof. If we consider the contribution of l th term and the other terms separately, we have

$$\begin{aligned} \left\| \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{M}) - \mathcal{P}_l(\mathbf{M}) \right) \mathbf{v} \right\|_2 &\leq \left| \sum_{j=1}^n \left(\frac{\delta_{lj}}{p} - 1 \right) M_{lj} v_j \right| + |v_l| \sqrt{\sum_{i=1}^n \left(\frac{\delta_{il}}{p} - 1 \right)^2 M_{il}^2} \\ &\leq \|\mathbf{M}\|_\infty \left(\left| \sum_{j=1}^n \left(\frac{\delta_{lj}}{p} - 1 \right) v_j \right| + \|\mathbf{v}\|_\infty \sqrt{\sum_{i=1}^n \left(\frac{\delta_{il}}{p} - 1 \right)^2} \right) \end{aligned}$$

From Lemma G.6, we have

$$\left| \sum_{j=1}^n \left(\frac{\delta_{lj}}{p} - 1 \right) v_j \right| \lesssim \sqrt{\frac{\log n}{p}} \|\mathbf{v}\|_2 + \frac{\log n}{p} \|\mathbf{v}\|_\infty$$

with high probability. Regarding the second term, notice that

$$\sum_{i=1}^n \left(\frac{\delta_{il}}{p} - 1 \right)^2 = n + \left(\frac{1}{p} - 2 \right) \sum_{i=1}^n \frac{\delta_{il}}{p}.$$

Lemma G.6 implies that $\sum_{i=1}^n \frac{\delta_{il}}{p} \asymp n$ with high probability if $n^2 p \gtrsim n \log n$. Hence, we have

$$\sum_{i=1}^n \left(\frac{\delta_{il}}{p} - 1 \right)^2 \lesssim \frac{n}{p},$$

and this finishes the proof. \square

Lemma G.12. *Let \mathbf{M} be a matrix and \mathbf{v}, \mathbf{w} be vectors that are independent from sampling of the l th row/column. Then, if $n^2 p \gtrsim n \log n$, we have*

$$\begin{aligned} &\left| \mathbf{w}^\top \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{M}) - \mathcal{P}_l(\mathbf{M}) \right) \mathbf{v} \right| \\ &\lesssim \|\mathbf{M}\|_\infty \left(\sqrt{\frac{\log n}{p}} (\|\mathbf{v}\|_2 \|\mathbf{w}\|_\infty + \|\mathbf{w}\|_2 \|\mathbf{v}\|_\infty) + \frac{\log n}{p} \|\mathbf{v}\|_\infty \|\mathbf{w}\|_\infty \right) \end{aligned}$$

Proof. We can consider the l th row and column separately by

$$\begin{aligned} &\left| \mathbf{w}^\top \left(\frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{M}) - \mathcal{P}_l(\mathbf{M}) \right) \mathbf{v} \right| \\ &\leq \left| v_l \sum_i \left(\frac{\delta_{il}}{p} - 1 \right) M_{il} w_i \right| + \left| w_l \sum_j \left(\frac{\delta_{lj}}{p} - 1 \right) M_{lj} v_j \right| + \left| \left(\frac{\delta_{ll}}{p} - 1 \right) M_{ll} v_l w_l \right| \\ &\leq \|\mathbf{M}\|_\infty \left(\|\mathbf{v}\|_\infty \left| \sum_i \left(\frac{\delta_{il}}{p} - 1 \right) w_i \right| + \|\mathbf{w}\|_\infty \left| \sum_j \left(\frac{\delta_{lj}}{p} - 1 \right) v_j \right| + \frac{1}{p} \|\mathbf{v}\|_\infty \|\mathbf{w}\|_\infty \right) \end{aligned}$$

If we apply Lemma G.6 to the summations, we get the desired result. \square

Lemma G.13. *Let \mathbf{E} be a symmetric matrix whose upper and on diagonal entries are drawn from Gaussian distribution $\mathcal{N}(0, \sigma^2)$ independently. Let \mathbf{v} be a vector that is independent from sampling of the l th row and column. Then, if $n^2 p \gtrsim n \log^2 n$, we have*

$$\left\| \frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{E}) \mathbf{v} \right\|_2 \lesssim \sigma \left(\sqrt{\frac{\log n}{p}} \|\mathbf{v}\|_2 + \frac{\sqrt{\log^3 n}}{p} \|\mathbf{v}\|_\infty + \sqrt{\frac{n}{p}} \|\mathbf{v}\|_\infty \right)$$

Proof. If we consider the contribution of l th term and the other terms separately, we have

$$\left\| \frac{1}{p} \mathcal{P}_{\Omega_l}(\mathbf{E}) \mathbf{v} \right\|_2 \leq \frac{1}{p} \left| \sum_{j=1}^n \delta_{lj} E_{lj} v_j \right| + \frac{1}{p} |v_l| \sqrt{\sum_{i=1}^n \delta_{il} E_{il}^2}$$

For the first term, we will calculate V and L of Corollary G.5. V is calculated as

$$V = \sum_{j=1}^n \mathbb{E}[(\delta_{lj} E_{lj} v_j)^2] = p\sigma^2 \|\mathbf{v}\|_2^2.$$

To find L , we first note that $\|\mathbf{E}_{l*}\|_\infty \lesssim \sigma\sqrt{\log n}$ with high probability, where \mathbf{E}_{l*} is the l th row of \mathbf{E} . Thus, for all $j \in [n]$, we have

$$|\delta_{lj} E_{lj} v_j| \lesssim \sigma\sqrt{\log n} \|\mathbf{v}\|_\infty.$$

Corollary G.5 implies that the first term is bounded as

$$\frac{1}{p} \left| \sum_{j=1}^n \delta_{lj} E_{lj} v_j \right| \lesssim \sigma \left(\sqrt{\frac{\log n}{p}} \|\mathbf{v}\|_2 + \frac{\sqrt{\log^3 n}}{p} \|\mathbf{v}\|_\infty \right). \quad (\text{G.1})$$

For the second term, it suffices to bound

$$\left| \sum_{i=1}^n \delta_{il} (E_{il}^2 - \sigma^2) \right|.$$

As before, we obtain V and L through

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}[(\delta_{il} (E_{il}^2 - \sigma^2))^2] &= p \sum_{i=1}^n \mathbb{E}[E_{il}^4 - 2\sigma^2 E_{il}^2 + \sigma^4] = 2\sigma^4 np, \\ |\delta_{il} (E_{il}^2 - \sigma^2)| &\lesssim \sigma^2 \log n. \end{aligned}$$

Corollary G.5 implies that

$$\left| \sum_{i=1}^n \delta_{il} (E_{il}^2 - \sigma^2) \right| \lesssim \sigma^2 \left(\sqrt{np \log n} + \log^2 n \right).$$

Because $\sum_{i=1}^n \delta_{il} \asymp np$, we have

$$\sum_{i=1}^n \delta_{il} E_{il}^2 \lesssim \sigma^2 \left(np + \sqrt{np \log n} + \log^2 n \right) \lesssim \sigma^2 np \quad (\text{G.2})$$

if $n^2 p \gtrsim n \log^2 n$. Combining (G.1) and (G.2), we get the desired bound. \square

Lemma G.14. *Let \mathbf{E} be a symmetric matrix whose upper and on diagonal entries are drawn from Gaussian distribution $\mathcal{N}(0, \sigma^2)$ independently. Let \mathbf{v}, \mathbf{w} be vectors that are independent from sampling of the l th row and column. Then, if $n^2 p \gtrsim n \log n$, we have*

$$\frac{1}{p} |\mathbf{w}^\top \mathcal{P}_{\Omega_l}(\mathbf{E}) \mathbf{v}| \lesssim \sigma \left(\sqrt{\frac{\log n}{p}} (\|\mathbf{v}\|_2 \|\mathbf{w}\|_\infty + \|\mathbf{w}\|_2 \|\mathbf{v}\|_\infty) + \frac{\sqrt{\log^3 n}}{p} \|\mathbf{v}\|_\infty \|\mathbf{w}\|_\infty \right).$$

Proof. We can consider the l th row and column separately by

$$\begin{aligned} \frac{1}{p} |\mathbf{w}^\top \mathcal{P}_{\Omega_l}(\mathbf{E}) \mathbf{v}| &\leq \frac{1}{p} \left| v_l \sum_i \delta_{il} E_{il} w_i \right| + \frac{1}{p} \left| w_l \sum_j \delta_{lj} E_{lj} v_j \right| + \frac{1}{p} |\delta_{ll} E_{ll} v_l w_l| \\ &\leq \frac{\|\mathbf{v}\|_\infty}{p} \left| \sum_i \delta_{il} E_{il} w_i \right| + \frac{\|\mathbf{w}\|_\infty}{p} \left| \sum_j \delta_{lj} E_{lj} v_j \right| + \frac{1}{p} \|\mathbf{v}\|_\infty \|\mathbf{w}\|_\infty |E_{ll}|. \end{aligned}$$

We bound the two summations similar to (G.1) and for the last term, we note that $|E_{ll}| \lesssim \sigma\sqrt{\log n}$ with high probability. \square