

---

# ReSync: Riemannian Subgradient-based Robust Rotation Synchronization

---

**Huikang Liu**

School of Information Management and Engineering  
Shanghai University of Finance and Economics  
liuhuikang@shufe.edu.cn

**Xiao Li\***

School of Data Science  
The Chinese University of Hong Kong, Shenzhen  
lixiao@cuhk.edu.cn

**Anthony Man-Cho So**

Department of Systems Engineering and Engineering Management  
The Chinese University of Hong Kong  
manchoso@se.cuhk.edu.hk

## Abstract

This work presents ReSync, a Riemannian subgradient-based algorithm for solving the robust rotation synchronization problem, which arises in various engineering applications. ReSync solves a least-squared minimization formulation over the rotation group, which is nonsmooth and nonconvex, and aims at recovering the underlying rotations directly. We provide strong theoretical guarantees for ReSync under the random corruption setting. Specifically, we first show that the initialization procedure of ReSync yields a proper initial point that lies in a local region around the ground-truth rotations. We next establish the weak sharpness property of the aforementioned formulation and then utilize this property to derive the local linear convergence of ReSync to the ground-truth rotations. By combining these guarantees, we conclude that ReSync converges linearly to the ground-truth rotations under appropriate conditions. Experiment results demonstrate the effectiveness of ReSync.

## 1 Introduction

Rotation synchronization (RS) is a fundamental problem in many engineering applications. For instance, RS (also known as “rotation averaging”) is an important subproblem of structure from motion (SfM) and simultaneous localization and mapping (SLAM) in computer vision [20, 22, 17], where the goal is to compute the absolute orientations of objects from relative rotations between pairs of objects. RS has also been applied to sensor network localization [41, 12], signal recovery from phaseless observations [2], digital communications [36], and cryo-EM imaging [35, 33].

Practical measurements of relative rotations are often *incomplete* and *corrupted*, leading to the problem of robust rotation synchronization (RRS) [28, 21, 22, 39, 9, 31]. The goal of RRS is to reconstruct a set of ground-truth rotations  $\mathbf{X}_1^*, \dots, \mathbf{X}_i^*, \dots, \mathbf{X}_n^* \in \text{SO}(d)$  from measurements of

---

\*Corresponding Author

relative rotations represented as

$$\mathbf{Y}_{ij} = \begin{cases} \mathbf{X}_i^* \mathbf{X}_j^{*\top}, & (i, j) \in \mathcal{A}, \\ \mathbf{O}_{ij}, & (i, j) \in \mathcal{E} \setminus \mathcal{A}, \\ \mathbf{0}, & (i, j) \in \mathcal{E}^c, \end{cases} \quad \text{with } (i, j) \in \begin{cases} \mathcal{A}, & \text{with ratio } pq, \\ \mathcal{E} \setminus \mathcal{A}, & \text{with ratio } (1-p)q, \\ \mathcal{E}^c, & \text{otherwise.} \end{cases} \quad (1)$$

Here,  $\text{SO}(d) := \{\mathbf{R} \in \mathbb{R}^{d \times d} : \mathbf{R}^\top \mathbf{R} = \mathbf{I}, \det(\mathbf{R}) = 1\}$  denotes the rotation group (also known as the special orthogonal group),  $\mathcal{E}$  represents the indices of all available observations,  $\mathcal{A}$  denotes the indices of true observations,  $\mathcal{A}^c := \mathcal{E} \setminus \mathcal{A}$  is the indices of *outliers*,  $\mathbf{O}_{ij} \in \text{SO}(d)$  is an outlying observation, and the missing observations are set to be  $\mathbf{0}$  by convention; see, e.g., [23, section 2.1]. We use  $q \in (0, 1)$  to denote the observation ratio and  $p \in (0, 1)$  to denote the ratio of true observations.

**Related works.** Due to the vast amount of research in this field, our overview will necessarily focus on theoretical investigations of RS. In the case where no outliers exist in the measurement model (1), i.e.,  $p = 1$ , a natural formulation is to minimize a smooth least-squares function  $\sum_{(i,j) \in \mathcal{E}} \|\mathbf{X}_i \mathbf{X}_j^\top - \mathbf{Y}_{ij}\|_F^2$  over  $\mathbf{X}_i \in \text{SO}(d)$ ,  $1 \leq i \leq n$ . Spectral relaxation and semidefinite relaxation (SDR) are typical approaches for addressing this problem [34, 3, 6, 5, 4, 30], where they provide strong recovery guarantees. However, these results cannot be directly applied to the corrupted model (1) due to the existence of outliers (i.e.,  $p < 1$ ) and the sensitivity of the least-squares solution to outlying observations.

Theoretical understanding of RRS is still rather limited. One typical setting for theoretical analysis of RRS is the random corruption model (RCM); see Section 2.2. The work [39] introduces a least-unsquared formulation and applies the SDR method to tackle it. Under the RCM and in the full observation case where  $q = 1$ , it is shown that the minimizer of the SDR reformulation exactly recovers the underlying Gram matrix (hence the ground-truth rotations) under the conditions that the true observation ratio  $p \geq 0.46$  for  $\text{SO}(2)$  (and  $p \geq 0.49$  for  $\text{SO}(3)$ ) and  $n \rightarrow \infty$ . In [23], the authors established the relationship between cycle-consistency and exact recovery and introduced a message-passing algorithm. Their method is tailored to find the corruption level in the graph, rather than recovering the ground-truth rotations directly. They provided linear convergence guarantees for their algorithm once the ratios satisfy  $p^8 q^2 = \Omega(\log n/n)$  under the RCM. However, it is unclear how this message-passing algorithm is related to other optimization procedures for solving the problem. Let us mention that they also provided guarantees for other compact groups and corruption settings. Following partly the framework established in [23], the work [32] presents an interesting nonconvex quadratic programming formulation of RRS. It is shown that the global minimizer of the nonconvex formulation recovers the true corruption level (still not the ground-true rotations directly) when  $p^2 q^2 = \Omega(\log n/n)$  under the RCM. Unfortunately, the work does not provide a concrete algorithm that provably finds a global minimizer of the nonconvex formulation. In [29], the authors introduced and analyzed a depth descent algorithm for recovering the underlying rotation matrices. In the context of the RCM, they showed asymptotic convergence of their algorithm to the underlying rotations without providing a specific rate. The result is achieved under the conditions that the algorithm is initialized near  $\mathbf{X}^*$ ,  $q \geq \mathcal{O}(\log n/n)$ , and  $p \geq 1 - 1/(d(d-1) + 2)$ . The latter requirement translates to  $p \geq 3/4$  for  $\text{SO}(2)$  and  $p \geq 7/8$  for  $\text{SO}(3)$ . It is important to note, however, that the primary focus of their research lies in the adversarial corruption setup rather than the RCM.

**Main contributions.** Towards tackling the RRS problem under the measurement model (1), we consider the following least-unsquared formulation, which was introduced in [39] as the initial step for applying the SDR method:

$$\begin{aligned} & \underset{\mathbf{X} \in \mathbb{R}^{n \times d \times d}}{\text{minimize}} \quad f(\mathbf{X}) := \sum_{(i,j) \in \mathcal{E}} \|\mathbf{X}_i \mathbf{X}_j^\top - \mathbf{Y}_{ij}\|_F \\ & \text{subject to} \quad \mathbf{X}_i \in \text{SO}(d), \quad 1 \leq i \leq n. \end{aligned} \quad (2)$$

Note that this problem is *nonsmooth* and *nonconvex* due to the unsquared Frobenius-norm loss and the rotation group constraint, respectively. We design a **Riemannian Subgradient synchronization** algorithm (ReSync) for addressing problem (2); see Algorithm 1. ReSync will first call an initialization procedure named SpectrIn (see Algorithm 2), which is a spectral relaxation method. Then, it implements an iterative Riemannian subgradient procedure. ReSync targets at directly recovering the ground-truth rotations  $\mathbf{X}^* \in \text{SO}(d)^n$  rather than the Gram matrix or the corruption level. Under the RCM (see Section 2.2), we provide the following strong theoretical guarantees for ReSync:

- (S.1) *Initialization.* The first step of ReSync is to call SpectrIn for computing the initial point  $\mathbf{X}^0$ . Theoretically, we establish that  $\mathbf{X}^0$  can be relatively close to  $\mathbf{X}^*$  depending on  $p$  and  $q$ ; see [Theorem 2](#).
- (S.2) *Weak sharpness.* We then establish a problem-intrinsic property of the formulation (2) called weak sharpness; see [Theorem 3](#). This property characterizes the geometry of problem (2) and is of independent interest.
- (S.3) *Convergence analysis.* Finally, we derive the local linear rate of convergence for ReSync based on the established weak sharpness property; see [Theorem 4](#).

The main idea is that the weak sharpness property in (S.2) helps to show *linear convergence* of ReSync to  $\mathbf{X}^*$  in (S.3). However, this result only holds *locally*. Thus, we need the initialization guarantee in (S.1) to initialize our algorithm in this local region and then argue that it will not leave this region once initialized. We refer to [Sections 3.1 to 3.3](#) for more technical challenges and our proof ideas. Combining the above theoretical results yields our overall guarantee: ReSync *converges linearly to the ground-truth rotations  $\mathbf{X}^*$  when  $p^7q^2 = \Omega(\log n/n)$* ; see [Theorem 1](#).

**Notation.** Our notation is mostly standard. We use  $\mathbb{R}^{nd \times d} \ni \mathbf{X} = (\mathbf{X}_1; \dots; \mathbf{X}_n) \in \text{SO}(d)^n$  to represent the Cartesian product of all the variables  $\mathbf{X}_i \in \text{SO}(d)$ ,  $1 \leq i \leq n$ . The same applies to the ground-truth rotations  $\mathbf{X}^* = (\mathbf{X}_1^*; \dots; \mathbf{X}_n^*)$ . Let  $\mathcal{E}_i = \{j \mid (i, j) \in \mathcal{E}\}$ ,  $\mathcal{A}_i = \{j \mid (i, j) \in \mathcal{A}\}$ , and  $\mathcal{A}_i^c = \mathcal{E}_i \setminus \mathcal{A}_i$ . We also define  $\mathcal{A}_{ij} = \mathcal{A}_i \cap \mathcal{A}_j$  for simplicity. For a set  $S$ , we use  $|S|$  to denote its cardinality. For any matrix  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{nd \times d}$ , we define the following distance up to a global rotation:

$$\text{dist}(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\mathbf{R}^*\|_F, \text{ where } \mathbf{R}^* = \arg \min_{\mathbf{R} \in \text{SO}(d)} \|\mathbf{X}\mathbf{R} - \mathbf{Y}\|_F^2 = \mathcal{P}_{\text{SO}(d)}(\mathbf{X}^\top \mathbf{Y}).$$

Besides, we introduce the following distances up to the global rotation  $\mathbf{R}^*$  defined above:

$$\text{dist}_1(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n \|\mathbf{X}_i - \mathbf{Y}_i \mathbf{R}^*\|_F, \quad \text{dist}_\infty(\mathbf{X}, \mathbf{Y}) = \max_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{Y}_i \mathbf{R}^*\|_F.$$

## 2 Algorithm and Setup

### 2.1 ReSync: Algorithm Development

In this subsection, we present ReSync for tackling the nonsmooth nonconvex formulation (2); see [Algorithm 1](#). Our algorithm has two main parts, i.e., initialization and an iterative Riemannian subgradient procedure.

**Initialization.** ReSync first calls a procedure SpectrIn (see [Algorithm 2](#)) for initialization. SpectrIn is a spectral relaxation-based initialization technique. SpectrIn computes the first  $d$  leading unit eigenvectors of the data matrix to form  $\tilde{\Phi} \in \mathbb{R}^{nd \times d}$ . We multiply  $\sqrt{n}$  to those eigenvectors to ensure that its norm matches that of  $\text{SO}(d)^n$ . We also construct  $\tilde{\Psi}$ , which reverses the sign of the last column of  $\tilde{\Phi}$  so that the determinants of  $\tilde{\Phi}$  and  $\tilde{\Psi}$  differ by a sign. Then, we compute the projection of  $\tilde{\Phi}$  and  $\tilde{\Psi}$  onto  $\text{SO}(d)^n$ . The projection is computed in a block-wise manner, namely

$$\tilde{\Phi}_i = \mathcal{P}_{\text{SO}(d)}(\tilde{\Phi}_i), \quad 1 \leq i \leq n,$$

where  $\tilde{\Phi}_i, \tilde{\Phi}_i \in \mathbb{R}^{d \times d}$  are the  $i$ -th block of  $\tilde{\Phi}$  and  $\tilde{\Phi}$ , respectively. The projection can be explicitly evaluated as

$$\tilde{\Phi}_i = \begin{cases} \mathbf{P}_i \mathbf{Q}_i^\top, & \text{if } \det(\tilde{\Phi}_i) > 0, \\ \hat{\mathbf{P}}_i \mathbf{Q}_i^\top, & \text{otherwise,} \end{cases} \quad 1 \leq i \leq n.$$

Here,  $\mathbf{P}_i, \mathbf{Q}_i \in \mathbb{R}^{d \times d}$  are the left and right singular vectors of  $\tilde{\Phi}_i$  (with descending order of singular values), respectively, and  $\hat{\mathbf{P}}_i$  is obtained by reversing the sign of the last column of  $\mathbf{P}_i$ . The initial point  $\mathbf{X}^0$  is chosen as  $\tilde{\Phi}$  or  $\tilde{\Psi}$ , depending on which is closer to  $\text{SO}(d)^n$ .

Let us mention that the computation of  $\tilde{\Psi}$  and Steps 5 - 9 in SpectrIn can practically improve the approximation error  $\text{dist}(\mathbf{X}^0, \mathbf{X}^*)$ . We demonstrate such a phenomenon in [Figure 1](#), in which ‘‘Naive SpectrIn’’ refers to outputting  $\mathbf{X}^0 = \tilde{\Phi}$  directly in [Algorithm 2](#).

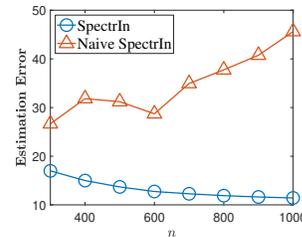


Figure 1: The average under 100 simulations of the initial distance  $\text{dist}(\mathbf{X}^0, \mathbf{X}^*)$  computed by [Algorithm 2](#) versus naive spectral initialization (i.e., outputting  $\mathbf{X}^0 = \tilde{\Phi}$  directly) with  $p = 0.2$ ,  $q = 0.2$  and  $d = 3$ .

**Riemannian subgradient update.** ReSync then implements an iterative Riemannian subgradient procedure after obtaining the initial point  $\mathbf{X}^0$ . The key is to compute the search direction (Riemannian subgradient)  $\tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i^k)$  and the retraction  $\text{Retr}_{\mathbf{X}_i^k}(\cdot)$  onto  $\text{SO}(d)$  for  $1 \leq i \leq n$ . Towards providing concrete formulas for the Riemannian subgradient update, let us impose the Euclidean inner product  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{A}^\top \mathbf{B})$  as the inherent Riemannian metric. Consequently, the tangent space to  $\text{SO}(d)$  at  $\mathbf{R} \in \text{SO}(d)$  is given by  $\text{T}_{\mathbf{R}} := \{\mathbf{R}\mathbf{S} : \mathbf{S} \in \mathbb{R}^{d \times d}, \mathbf{S} + \mathbf{S}^\top = 0\}$ . The Riemannian subgradient  $\tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i)$  can be computed as [40, Theorem 5.1]

$$\tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i) = \mathcal{P}_{\text{T}_{\mathbf{X}_i}}(\tilde{\nabla}f(\mathbf{X}_i)), \quad 1 \leq i \leq n, \quad (3)$$

where the projection can be computed as  $\mathcal{P}_{\text{T}_{\mathbf{X}_i}}(\mathbf{B}) = \mathbf{X}_i(\mathbf{X}_i^\top \mathbf{B} - \mathbf{B}^\top \mathbf{X}_i)/2$  for any  $\mathbf{B} \in \mathbb{R}^{d \times d}$  and  $\tilde{\nabla}f(\mathbf{X}_i)$  is the Euclidean subgradient of  $f$  with respect to the  $i$ -th block variable  $\mathbf{X}_i$ . Let us define  $f_{i,j}(\mathbf{X}) := \|\mathbf{X}_i \mathbf{X}_j^\top - \mathbf{Y}_{ij}\|_F$ . The Euclidean subdifferential  $\partial f(\mathbf{X}_i)$  with respect to the block variable  $\mathbf{X}_i$  is given by

$$\partial f(\mathbf{X}_i) = 2 \sum_{j:(i,j) \in \mathcal{E}} \partial f_{i,j}(\mathbf{X}_i), \quad \text{with} \quad \partial f_{i,j}(\mathbf{X}_i) = \begin{cases} \frac{\mathbf{X}_i - \mathbf{Y}_{ij} \mathbf{X}_j}{\|\mathbf{X}_i \mathbf{X}_j^\top - \mathbf{Y}_{ij}\|_F}, & \text{if } \|\mathbf{X}_i \mathbf{X}_j^\top - \mathbf{Y}_{ij}\|_F \neq 0, \\ \mathbf{V} \in \mathbb{R}^{d \times d}, \|\mathbf{V}\|_F \leq 1, & \text{otherwise.} \end{cases}$$

---

**Algorithm 1** ReSync: Riemannian Subgradient Synchronization

---

**Require:** Initialize  $\mathbf{X}^0 = \text{SpectrIn}(\mathbf{Y})$  (Algorithm 2), where  $\mathbf{Y} \in \mathbb{R}^{nd \times nd}$  and its  $(i, j)$ -th block is  $\mathbf{Y}_{i,j} \in \mathbb{R}^{d \times d}$ ;  
 1: Set iteration count  $k = 0$ ;  
 2: **while** stopping criterion not met **do**  
 3:   Update the step size  $\mu_k$ ;  
 4:   Riemannian subgradient update:

$$\mathbf{X}_i^{k+1} = \text{Retr}_{\mathbf{X}_i^k} \left( -\mu_k \tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i^k) \right)$$

for  $1 \leq i \leq n$ ;

5:   Update iteration count  $k = k + 1$ ;  
 6: **end while**

---



---

**Algorithm 2** SpectrIn: Spectral Initialization

---

1: **Input:**  $\mathbf{Y} \in \mathbb{R}^{nd \times nd}$ ;  
 2: Compute the  $d$  leading unit eigenvectors of  $\mathbf{Y}$ :  $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ ;  
 3: Set  $\tilde{\Phi} = \sqrt{n}[\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d] \in \mathbb{R}^{nd \times d}$  and  $\tilde{\Psi} = \sqrt{n}[\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{d-1}, -\mathbf{u}_d]$ ;  
 4: Compute  $\tilde{\Phi} = \mathcal{P}_{\text{SO}(d)^n}(\tilde{\Phi})$  and  $\tilde{\Psi} = \mathcal{P}_{\text{SO}(d)^n}(\tilde{\Psi})$ ;  
 5: **if**  $\|\tilde{\Phi} - \tilde{\Phi}\|_F \leq \|\tilde{\Psi} - \tilde{\Psi}\|_F$  **then**  
 6:    $\mathbf{X}^0 = \tilde{\Phi}$ ;  
 7: **else**  
 8:    $\mathbf{X}^0 = \tilde{\Psi}$ ;  
 9: **end if**  
 10: **Output:** Initial point  $\mathbf{X}^0$ .

---

Any element  $\tilde{\nabla}f(\mathbf{X}_i) \in \partial f(\mathbf{X}_i)$  is called a Euclidean subgradient. In ReSync, one can choose an arbitrary subgradient  $\tilde{\nabla}f(\mathbf{X}_i) \in \partial f(\mathbf{X}_i)$  at  $\mathbf{X}_i$ .

Mimicking the gradient method to update along the search direction  $\tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i)$  provides a point  $\mathbf{X}_i^+ = \mathbf{X}_i - \mu \tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i)$  on the tangent space  $\text{T}_{\mathbf{X}_i}$  at  $\mathbf{X}_i$ , which may violate the manifold constraint " $\mathbf{X}_i^+ \in \text{SO}(d)$ ". One common approach in Riemannian optimization is to employ a retraction operator to address the feasibility issue. For  $\text{SO}(d)$ , we can use a QR decomposition-based retraction and implement the Riemannian subgradient step as

$$\mathbf{X}_i^+ = \text{Retr}_{\mathbf{X}_i} \left( -\mu \tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i) \right) = \text{Qr} \left( \mathbf{X}_i - \mu \tilde{\nabla}_{\mathcal{R}}f(\mathbf{X}_i) \right), \quad 1 \leq i \leq n. \quad (4)$$

Here,  $\text{Qr}(\mathbf{B})$  returns the Q-factor in the thin QR decomposition of  $\mathbf{B}$ , while the diagonal entries of the R-factor are restricted to be positive [7].

Finally, setting  $\mathbf{X}_i = \mathbf{X}_i^k$ ,  $\mathbf{X}_j = \mathbf{X}_j^k$  for all  $j$  such that  $(i, j) \in \mathcal{E}$ ,  $\mu = \mu_k$  in (3) and (4) yields a concrete implementation of Step 4 in ReSync and leads to  $\text{SO}(d) \ni \mathbf{X}_i^{k+1} = \mathbf{X}_i^+$  for  $1 \leq i \leq n$ . This completes the description of one full iteration of ReSync. Note that the per-iteration complexity of the Riemannian subgradient procedure is  $\mathcal{O}(n^2q)$ , and Algorithm 2 has computational cost  $\mathcal{O}(n^3)$ .

## 2.2 RCM Setup for Theoretical Analysis

We develop our theoretical analysis of ReSync by adopting the random corruption model (RCM). The RCM was previously used in many works to analyze the performance of various synchronization algorithms; see, e.g., [39, 19, 23, 32]. Specifically, we can represent our measurement model (1) on a graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of  $n$  nodes representing  $\{\mathbf{X}_1^*, \dots, \mathbf{X}_n^*\}$  and  $\mathcal{E}$  is a set of edges

containing all the available measurements  $\{\mathbf{Y}_{i,j}, (i,j) \in \mathcal{E}\}$ . We assume that the graph  $\mathcal{G}$  follows the well-known Erdős-Rényi model  $G(n, q)$ , which implies that each edge  $(i, j) \in \mathcal{E}$  is observed with probability  $q$ , independently from every other edge. Each edge  $(i, j) \in \mathcal{E}$  is a true observation (i.e.,  $(i, j) \in \mathcal{A}$ ) with probability  $p$  and an outlier (i.e.,  $(i, j) \in \mathcal{A}^c$ ) with probability  $1 - p$ . Furthermore, the outliers  $\{\mathbf{O}_{i,j}\}_{(i,j) \in \mathcal{A}^c}$  are assumed to be independently and uniformly distributed on  $\text{SO}(d)$ .

### 3 Main Results

In this section, we present our theoretical results for ReSync. Our main results are summarized in the following theorem, which states that our proposed algorithm can converge at a linear rate to the underlying rotations  $\mathbf{X}^*$ . Our standing assumption in this section is stated below.

All our theoretical results in this section are based on the RCM; see [Section 2.2](#).

**Theorem 1** (overall). *Suppose that the ratios  $p$  and  $q$  satisfy*

$$p^7 q^2 = \Omega\left(\frac{\log n}{n}\right).$$

*With probability at least  $1 - \mathcal{O}(1/n)$ , ReSync with  $\mu_k = \mu_0 \gamma^k$ , where  $\mu_0 = \Theta(p^2/n)$  and  $\gamma = 1 - \frac{pq}{16}$ , converges linearly to the ground-truth rotations  $\mathbf{X}^*$  (up to a global rotation), i.e.,*

$$\text{dist}(\mathbf{X}^k, \mathbf{X}^*) \leq \xi_0 \gamma^k, \quad \text{dist}_\infty(\mathbf{X}^k, \mathbf{X}^*) \leq \delta_0 \gamma^k, \quad \forall k \geq 0.$$

*Here,  $\xi_0 = \Theta(\sqrt{np^5q})$  and  $\delta_0 = \Theta(p^2)$ .*

The basic idea of the proof is to establish the problem-intrinsic property of weak sharpness and then use it to derive a linear convergence result. However, the result only holds locally. Thus, we develop a procedure to initialize the algorithm in this local region and argue that ReSync will not leave this region afterwards. In the remaining parts of this section, we implement the above ideas and highlight the challenges and approaches to overcoming them.

#### 3.1 Analysis of SpectrIn with Leave-One-Out Technique

**Theorem 2** (initialization). *Let  $\mathbf{X}^0$  be generated by SpectrIn (see [Algorithm 2](#)). Suppose that the ratios  $p$  and  $q$  satisfy*

$$p^2 q = \Omega\left(\frac{\log n}{n}\right).$$

*Then, with probability at least  $1 - \mathcal{O}(1/n)$ , we have*

$$\text{dist}(\mathbf{X}^0, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right) \quad \text{and} \quad \text{dist}_\infty(\mathbf{X}^0, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{nq}}\right). \quad (5)$$

The works [34] and [11] show that exact reconstruction of  $\mathbf{X}^*$  is information-theoretically possible if the condition  $p^2 q = \Omega(\log n/n)$  holds for the cases  $d = 2$  and  $d = 3$ , respectively. Though [Theorem 2](#) does not provide exact recovery, it achieves an optimal sample complexity for reconstructing an approximate solution in the infinity norm. Specifically, [Theorem 2](#) shows that, as long as  $p^2 q \geq C \log n/n$  for some constant  $C > 0$  large enough, the  $\ell_\infty$ -distance  $\text{dist}_\infty(\mathbf{X}^0, \mathbf{X}^*)$  (i.e.,  $\max_{1 \leq i \leq n} \text{dist}(\mathbf{X}_i, \mathbf{X}_i^*)$ ) can be made relatively small. However, the  $\ell_2$ -distance  $\text{dist}(\mathbf{X}^0, \mathbf{X}^*)$  is of the order  $\Omega(\sqrt{n})$  under such a sample complexity.

The work [26] considers orthogonal and permutation group synchronization and shows that spectral relaxation-based methods achieve near-optimal performance bounds. Our result differs from that of [26] in twofold: 1) Our approach follows the standard leave-one-out analysis based on the standard “Dist” (up to  $\mathcal{O}(d)$  invariance) defined above [Lemma 3](#) in the Appendix. Nonetheless, we have to transfer the results to “dist” due to the structure of  $\text{SO}(d)$  in [Lemma 5](#), which is a nontrivial step due to the specific structure of  $\text{SO}(d)$ . 2) Our result can handle incomplete observations (i.e.,  $q < 1$ ). In the case of incomplete observations, the construction in (17) in the Appendix becomes more intricate; it has the additional third column, rendering the analysis of our [Lemma 2](#) more involved.

We prove [Theorem 2](#) with some matrix concentration bounds and the leave-one-out technique. We provide the proof sketch below and refer to [Appendix A](#) for the full derivations.

**Proof outline of Theorem 2.** According to (1) and the fact  $\mathbb{E}(\mathbf{O}_{ij}) = \mathbf{0}$  since outliers are assumed to be independently and uniformly distributed on  $\text{SO}(d)$  in the RCM (see Appendix A), we know that  $\mathbb{E}(\mathbf{Y}_{ij}) = pq\mathbf{X}_i^*\mathbf{X}_j^{*\top}$  for all  $(i, j) \in [n] \times [n]$ . This motivates us to introduce the noise matrix  $\mathbf{W}_{ij} = \mathbf{Y}_{ij} - pq\mathbf{X}_i^*\mathbf{X}_j^{*\top}$ , i.e.,

$$\mathbf{Y} = pq\mathbf{X}^*\mathbf{X}^{*\top} + \mathbf{W}. \quad (6)$$

The condition  $p^2q = \Omega(\log n/n)$  in Theorem 2 ensures that the expectation  $pq\mathbf{X}^*\mathbf{X}^{*\top}$  will dominate the noise matrix  $\mathbf{W}$  in the decomposition (6).

We first discuss how to bound  $\text{dist}(\mathbf{X}^0, \mathbf{X}^*)$ . Notice that  $\mathbf{X}^0$  and  $\mathbf{X}^*$  are the  $d$  leading eigenvectors of  $\mathbf{Y}$  (after projection onto  $\text{SO}(d)^n$ ) and  $pq\mathbf{X}^*\mathbf{X}^{*\top}$ , respectively. We can then use the matrix perturbation theory (see Lemma 3) to bound  $\text{dist}(\mathbf{X}^0, \mathbf{X}^*)$ . Towards this end, we need to estimate the operator norm  $\|\mathbf{W}\|_2$ , which could be done by applying the standard matrix Bernstein concentration inequality [38] since the blocks  $\{\mathbf{W}_{ij}\}$  are i.i.d. white noise with bounded operator norms and variances; see Lemma 2.

We next turn to bound the initialization error in the infinity norm, i.e.,  $\text{dist}_\infty(\mathbf{X}^0, \mathbf{X}^*)$ . Let us use  $(\mathbf{W}\mathbf{X}^0)_m \in \mathbb{R}^{d \times d}$  to denote the  $m$ -th block of  $\mathbf{W}\mathbf{X}^0 \in \mathbb{R}^{nd \times d}$  for  $1 \leq m \leq n$ . The main technical challenge lies in deriving a sharp bound for the term  $\max_{1 \leq m \leq n} \|(\mathbf{W}\mathbf{X}^0)_m\|_F$ , as it involves two *dependent* random quantities, i.e., the noise matrix  $\mathbf{W}$  and the initial  $\mathbf{X}^0$  that is obtained by projecting the first  $d$  leading eigenvectors of  $\mathbf{Y}$  onto  $\text{SO}(d)^n$ . To overcome such a statistical dependence, we utilize the leave-one-out technique. This technique was utilized in [42] to analyze the phase synchronization problem and was later applied to many other synchronization problems [1, 10, 15, 18, 26]. Let us define

$$\mathbf{Y}^{(m)} = pq\mathbf{X}^*\mathbf{X}^{*\top} + \mathbf{W}^{(m)} \quad \text{with} \quad \mathbf{W}_{kl}^{(m)} = \mathbf{W}_{kl} \cdot \mathbf{1}_{\{k \neq m\}} \cdot \mathbf{1}_{\{l \neq m\}}. \quad (7)$$

That is, we construct  $\mathbf{W}^{(m)} \in \mathbb{R}^{nd \times nd}$  by setting the  $m$ -th block-wise row and column of  $\mathbf{W}$  to be  $\mathbf{0}$ . Then, it is easy to see that  $\mathbf{Y}^{(m)}$  is statistically independent of  $\mathbf{W}_m^\top \in \mathbb{R}^{d \times nd}$ , where the latter denotes the  $m$ -th block-wise row of  $\mathbf{W}$ . Let  $\mathbf{X}^{(m)}$  be the  $d$  leading eigenvectors of  $\mathbf{Y}^{(m)}$ . Consequently,  $\mathbf{X}^{(m)}$  is also independent of  $\mathbf{W}_m^\top$ . Based on the above discussions, we can bound each  $\|(\mathbf{W}\mathbf{X}^0)_m\|_F$  in the following way:

$$\|(\mathbf{W}\mathbf{X}^0)_m\|_F = \|\mathbf{W}_m^\top \mathbf{X}^0\|_F \leq \|\mathbf{W}_m^\top \mathbf{X}^{(m)}\|_F + \|\mathbf{W}_m^\top (\mathbf{X}^0 - \mathbf{X}^{(m)})\|_F. \quad (8)$$

The first term  $\|\mathbf{W}_m^\top \mathbf{X}^{(m)}\|_F$  can be bounded using an appropriate concentration inequality due to the statistical independence between  $\mathbf{W}_m^\top$  and  $\mathbf{X}^{(m)}$ . The second term can be bounded as

$$\|\mathbf{W}_m^\top (\mathbf{X}^0 - \mathbf{X}^{(m)})\|_F \leq \|\mathbf{W}_m\|_2 \cdot \|\mathbf{X}^0 - \mathbf{X}^{(m)}\|_F,$$

in which  $\|\mathbf{W}_m\|_2$  can be further bounded by matrix concentration inequality (see Lemma 2) and  $\|\mathbf{X}^0 - \mathbf{X}^{(m)}\|_F$  can be bounded using standard matrix perturbation theory (see Lemma 4).

### 3.2 Weak Sharpness and Exact Recovery

We next present a property that is intrinsic to problem (2) in the following theorem.

**Theorem 3** (weak sharpness). *Suppose that the ratios  $p$  and  $q$  satisfy*

$$p^2q^2 = \Omega\left(\frac{\log n}{n}\right).$$

*Then, with probability at least  $1 - \mathcal{O}(1/n)$ , for any  $\mathbf{X} \in \text{SO}(d)^n$  satisfying  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ , we have*

$$f(\mathbf{X}) - f(\mathbf{X}^*) \geq \frac{npq}{8} \text{dist}_1(\mathbf{X}, \mathbf{X}^*).$$

Some remarks on Theorem 3 are in order. This theorem shows that problem (2) possesses the *weak sharpness* property [8], which is intrinsic to the problem and independent of the algorithm used to solve it. It is known that with this property, various subgradient-type methods can achieve linear convergence [14, 25]. We will establish a similar linear convergence result for ReSync in the next subsection based on Theorem 3.

The weak sharpness property shown in Theorem 3 is of independent interest, as it could be helpful when analyzing other optimization algorithms (not just ReSync) for solving problem (2). Currently,

only a few applications are known to produce sharp optimization problems, such as robust low-rank matrix recovery [25], robust phase retrieval [16], and robust subspace recovery [24]. Furthermore, sharp instances of manifold optimization problems are especially scarce. Hence, [Theorem 3](#) extends the list of optimization problems that possess the weak sharpness property and contributes to the growing literature on the geometry of structured nonsmooth nonconvex optimization problems.

It is worth noting that [Theorem 3](#) also establishes the *exact recovery* property of the formulation (2). Specifically, up to a global rotation, the ground-truth  $\mathbf{X}^*$  is guaranteed to be the unique global minimizer of  $f$  over the region  $\text{SO}(d)^n \cap \{\mathbf{X} : \text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)\}$ . Consequently, recovering the underlying  $\mathbf{X}^*$  reduces to finding the global minimizer of  $f$  over the aforementioned region. As we will show in the next subsection, ReSync will converge linearly to the global minimizer  $\mathbf{X}^*$  when initialized in this region. However, the initialization requirement is subject to the stronger condition  $p^4q = \Omega(\log n/n)$  on the ratios  $p$  and  $q$ , which is ensured by [Theorem 2](#).

We list our main ideas for proving [Theorem 3](#) below. The full proof can be found in [Appendix B](#).

**Proof outline of [Theorem 3](#).** Note that the objective function  $f$  can be decomposed into two parts:

$$f(\mathbf{X}) = \overbrace{\sum_{(i,j) \in \mathcal{A}} \|\mathbf{X}_i^\top \mathbf{X}_j - \mathbf{X}_i^{*\top} \mathbf{X}_j^*\|_F}^{g(\mathbf{X})} + \overbrace{\sum_{(i,j) \in \mathcal{A}^c} \|\mathbf{X}_i^\top \mathbf{X}_j - \mathbf{O}_{ij}\|_F}^{h(\mathbf{X})}. \quad (9)$$

It is easy to see that  $g(\mathbf{X}^*) = 0$  and  $g(\mathbf{X}) \geq 0$ . Based on the fact that the true observation is uniformly distributed in all the indices, we have  $\mathbb{E}(g(\mathbf{X})) = pq \sum_{1 \leq i, j \leq n} \|\mathbf{X}_i^\top \mathbf{X}_j - \mathbf{X}_i^{*\top} \mathbf{X}_j^*\|_F \geq \frac{npq}{2} \text{dist}_1(\mathbf{X}, \mathbf{X}^*)$ ; see [Appendix B.2](#) for the last inequality. A traditional way to lower bounding  $g(\mathbf{X})$  using  $\mathbb{E}(g(\mathbf{X}))$  for all  $\mathbf{X} \in \text{SO}(d)$  is to apply concentration inequality and an epsilon-net covering argument. Unfortunately, the sample complexity condition  $p^2q^2 = \Omega(\log n/n)$  does not lead to a high probability result in this way. Instead, our approach is to apply the concentration theory on the cardinalities of index sets rather than on  $\mathbf{X}$  directly; see the following lemma.

**Lemma 1** (concentration of cardinalities of index sets). *Given any  $\epsilon = \Omega\left(\frac{\sqrt{\log n}}{\sqrt{npq}}\right)$ , with probability at least  $1 - \mathcal{O}(1/n)$ , we have*

$$\begin{aligned} (1 - \epsilon)npq &\leq |\mathcal{E}_i| \leq (1 + \epsilon)npq, & (1 - \epsilon)npq &\leq |\mathcal{A}_i| \leq (1 + \epsilon)npq, \\ (1 - \epsilon)npq^2 &\leq |\mathcal{E}_i \cap \mathcal{A}_j| \leq (1 + \epsilon)npq^2, & (1 - \epsilon)np^2q^2 &\leq |\mathcal{A}_{ij}| \leq (1 + \epsilon)np^2q^2 \end{aligned}$$

for any  $1 \leq i, j \leq n$ . See [Section 1](#) for the notation.

We then provide a sharp lower bound on  $g(\mathbf{X})$  based on [Lemma 1](#).

**Proposition 1.** *Under the conditions of [Theorem 3](#), with probability at least  $1 - \mathcal{O}(1/n)$ , we have*

$$g(\mathbf{X}) \geq \frac{3npq}{16} \text{dist}_1(\mathbf{X}, \mathbf{X}^*), \quad \forall \mathbf{X} \in \text{SO}(d)^n. \quad (10)$$

Next, to lower bound  $h(\mathbf{X}) - h(\mathbf{X}^*) = \sum_{(i,j) \in \mathcal{A}^c} (\|\mathbf{X}_i^\top \mathbf{X}_j - \mathbf{O}_{ij}\|_F - \|\mathbf{X}_i^{*\top} \mathbf{X}_j^* - \mathbf{O}_{ij}\|_F)$  we first bound

$$h(\mathbf{X}) - h(\mathbf{X}^*) \geq \sum_{(i,j) \in \mathcal{A}^c} \left\langle \frac{\mathbf{X}_i^{*\top} \mathbf{X}_j^* - \mathbf{O}_{ij}}{\|\mathbf{X}_i^{*\top} \mathbf{X}_j^* - \mathbf{O}_{ij}\|_F}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{X}_i^{*\top} \mathbf{X}_j^* \right\rangle,$$

where the inequality comes from the convexity of the norm function  $U \mapsto \|U\|_F$  whenever  $\mathbf{X}_i^{*\top} \mathbf{X}_j^* - \mathbf{O}_{ij} \neq \mathbf{0}$ . Then, using the orthogonality of each block of  $\mathbf{X}^*$ , we further have

$$h(\mathbf{X}) - h(\mathbf{X}^*) \geq \sum_{(i,j) \in \mathcal{A}^c} \left\langle \frac{\mathbf{I} - \mathbf{X}_i^* \mathbf{O}_{ij} \mathbf{X}_j^{*\top}}{\|\mathbf{I} - \mathbf{X}_i^* \mathbf{O}_{ij} \mathbf{X}_j^{*\top}\|_F}, \mathbf{X}_i^* \mathbf{X}_i^\top \mathbf{X}_j \mathbf{X}_j^\top - \mathbf{I} \right\rangle. \quad (11)$$

Recall that since the outliers  $\{\mathbf{O}_{i,j}\}_{(i,j) \in \mathcal{A}^c}$  are independently and uniformly distributed on  $\text{SO}(d)$ , so are  $\{\mathbf{X}_i^* \mathbf{O}_{ij} \mathbf{X}_j^{*\top}\}_{(i,j) \in \mathcal{A}^c}$ . This observation indicates that  $\{\mathbf{I} - \mathbf{X}_i^* \mathbf{O}_{ij} \mathbf{X}_j^{*\top} / \|\mathbf{I} - \mathbf{X}_i^* \mathbf{O}_{ij} \mathbf{X}_j^{*\top}\|_F\}_{(i,j) \in \mathcal{A}^c}$  are i.i.d. random matrices. Hence, by invoking concentration results that utilize the randomness of the outliers  $\{\mathbf{O}_{i,j}\}_{(i,j) \in \mathcal{A}^c}$  and the cardinalities  $(i, j) \in \mathcal{A}^c$ , we obtain the following result.

**Proposition 2.** Under the conditions of [Theorem 3](#), with probability at least  $1 - \mathcal{O}(1/n)$ , we have

$$h(\mathbf{X}) - h(\mathbf{X}^*) \geq -\frac{npq}{16} \text{dist}_1(\mathbf{X}, \mathbf{X}^*) \quad (12)$$

for all  $\mathbf{X} \in \text{SO}(d)^n$  satisfying  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ .

Combining [Proposition 1](#) and [Proposition 2](#) gives [Theorem 3](#).

### 3.3 Convergence Analysis and Proof of [Theorem 1](#)

Let us now turn to utilize the weak sharpness property shown in [Theorem 3](#) to establish the local linear convergence of ReSync. As a quick corollary of [Theorem 3](#), we have the following result.

**Corollary 1.** Under the conditions of [Theorem 3](#), with probability at least  $1 - \mathcal{O}(1/n)$ , for any  $\mathbf{X} \in \text{SO}(d)^n$  satisfying  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ , we have

$$\left\langle \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}), \mathbf{X}^* - \mathbf{X} \right\rangle \leq -\frac{npq}{16} \text{dist}_1(\mathbf{X}, \mathbf{X}^*), \quad \forall \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}) \in \partial_{\mathcal{R}} f(\mathbf{X}). \quad (13)$$

This condition indicates that any Riemannian subgradient  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X})$  provides a descent direction pointing towards  $\mathbf{X}^*$ . However, it only holds for  $\mathbf{X} \in \text{SO}(d)^n$  satisfying  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ . Our key idea for establishing local convergence is to show that the Riemannian subgradient update in ReSync is a contraction operator in both the Euclidean and infinity norm-induced distances using [Corollary 1](#), i.e., if  $\mathbf{X}^k$  lies in the local region, then  $\mathbf{X}^{k+1}$  also lies in the region. This idea motivates us to define two sequences of neighborhoods as follows:

$$\mathcal{N}_F^k = \{\mathbf{X} \mid \text{dist}(\mathbf{X}, \mathbf{X}^*) \leq \xi_k\} \quad \text{and} \quad \mathcal{N}_\infty^k = \{\mathbf{X} \mid \text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) \leq \delta_k\}. \quad (14)$$

Here,  $\xi_k = \xi_0 \gamma^k$ ,  $\delta_k = \delta_0 \gamma^k$ , where  $\xi_0$ ,  $\delta_0$ , and  $\gamma \in (0, 1)$  will be specified later. Thus, these two sequences of sets  $\{\mathcal{N}_F^k\}$  and  $\{\mathcal{N}_\infty^k\}$  will linearly shrink to the ground-truth. It remains to show that if  $\mathbf{X}^k \in \mathcal{N}_F^k \cap \mathcal{N}_\infty^k$ , then  $\mathbf{X}^{k+1} \in \mathcal{N}_F^{k+1} \cap \mathcal{N}_\infty^{k+1}$ , which is summarized in the following theorem.

**Theorem 4** (convergence analysis). Suppose that  $\delta_0 = \mathcal{O}(p^2)$  and  $\xi_0 = \mathcal{O}(\sqrt{npq}\delta_0)$ . Set  $\gamma = 1 - \frac{pq}{16}$  and  $\mu_k = \delta_k/n$  in ReSync. If  $\mathbf{X}^k \in \mathcal{N}_F^k \cap \mathcal{N}_\infty^k$  for any  $k \geq 0$ , then with probability at least  $1 - \mathcal{O}(1/n)$ , we have

$$\mathbf{X}^{k+1} \in \mathcal{N}_F^{k+1} \cap \mathcal{N}_\infty^{k+1}.$$

**Proof outline of [Theorem 4](#).** The proof consisted of two parts. On the one hand, we need to show that  $\mathbf{X}^{k+1} \in \mathcal{N}_F^{k+1}$ , which can be achieved by applying [Corollary 1](#). On the other hand, in order to show that  $\mathbf{X}^{k+1} \in \mathcal{N}_\infty^{k+1}$ , we need a good estimate of each block of  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X})$ . See [Appendix C](#).

Having developed the necessary tools, we are now ready to prove [Theorem 1](#).

**Proof of [Theorem 1](#).** Based on [Theorem 2](#), we know that  $\mathbf{X}^0 \in \mathcal{N}_F^0 \cap \mathcal{N}_\infty^0$  if  $\xi_0$  and  $\delta_0$  satisfy

$$\xi_0 = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right) \quad \text{and} \quad \delta_0 = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{npq}}\right). \quad (15)$$

According to [Theorem 4](#), by choosing  $\delta_0 = \Theta(p^2)$  and  $\xi_0 = \Theta(\sqrt{np^5q})$ , condition (15) holds when  $p^7q^2 = \Omega(\log n/n)$ . This completes the proof of [Theorem 1](#).

## 4 Experiments

In this section, we conduct experiments on ReSync for solving the RRS problem on both synthetic and real data, providing empirical support for our theoretical findings. Our experiments are conducted on a personal computer with a 2.90GHz 8-core CPU and 32GB memory. All our experiment results are averaged over 20 independent trials. Our code is available at <https://github.com/Huikang2019/ReSync>.

### 4.1 Synthetic Data

We consider the rotation group  $\text{SO}(3)$  in all our experiments. We generate  $\mathbf{X}_1^*, \dots, \mathbf{X}_n^*$  by first generating matrices of the same dimension with i.i.d. standard Gaussian entries and then projecting

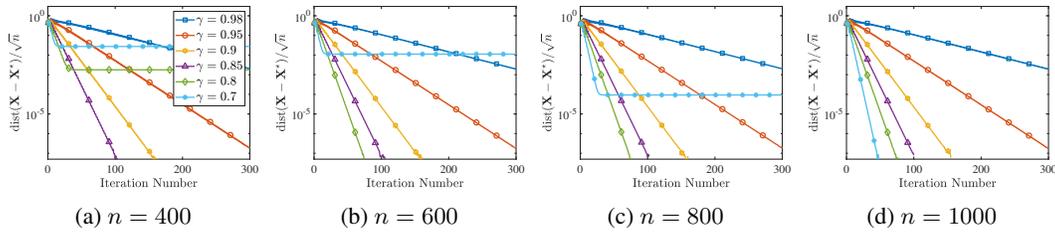


Figure 2: Convergence of ReSync with  $p = q = (\log n/n)^{1/3}$ .

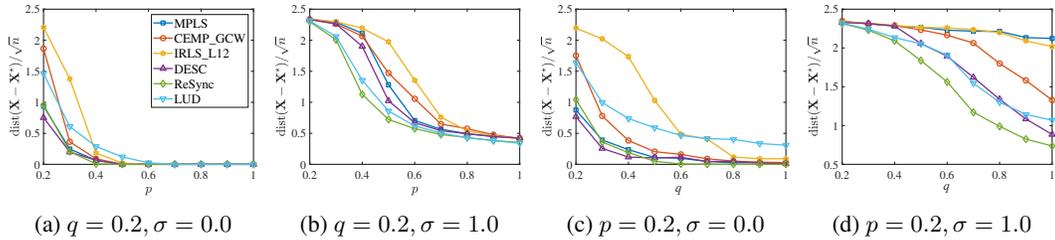


Figure 3: Comparison with state-of-the-art synchronization algorithms.

each of them onto  $SO(3)$ . The underlying graph, outliers, and relative rotations in the measurement model (1) are generated according to the RCM as described in Section 2.2. In our experiments, we also consider the case where the true observations are contaminated by additive noise, namely,  $\{\mathbf{Y}_{i,j}\}_{(i,j) \in \mathcal{A}}$  in (1) is generated using the formula

$$\mathbf{Y}_{i,j} = \mathcal{P}_{SO(3)}(\mathbf{X}_i^{*\top} \mathbf{X}_j^* + \sigma \mathbf{G}_{i,j}) \quad \text{for } (i,j) \in \mathcal{A}, \quad (16)$$

where  $\mathbf{G}_{i,j}$  consists of i.i.d. entries following the standard Gaussian distribution and  $\sigma \geq 0$  controls the variance level of the noise.

**Convergence verification of ReSync.** We evaluate the convergence performance of ReSync with the noise level  $\sigma = 0$  in (16). We set  $p = q = (\log n/n)^{1/3}$  in the measurement model (1), which satisfies  $p^2q = \log n/n$ . We use the initial step size  $\mu_0 = 1/npq$  and the decaying factor  $\gamma \in \{0.7, 0.8, 0.85, 0.90, 0.95, 0.98\}$  in ReSync. We test the performance for various  $n$  selected from  $\{400, 600, 800, 1000\}$ . Figure 2 displays the experiment results. It can be observed that (i) ReSync converges linearly to ground-truth rotations for a wide range of  $\gamma$  and (ii) a smaller  $\gamma$  often leads to faster convergence speed. These corroborate our theoretical findings. However, it is worth noting that excessively small  $\gamma$  values may result in an early stopping phenomenon (e.g.,  $\gamma \leq 0.8$  when  $n = 400$ ). In addition, ReSync performs better with a larger  $n$ , as it allows for a smaller  $\gamma$  (e.g.,  $\gamma = 0.7$  when  $n = 1000$ ) and hence converges to the ground-truth rotations faster.

**Comparison with the state-of-the-arts.** We next compare ReSync with state-of-the-art synchronization algorithms, including IRLS\_L12 [9], MPLS [31], CEMP\_GCW [23, 32], DESC [32], and LUD [39]. We obtain the implementation of the first four algorithms from <https://github.com/CoLeWyeth/DESC>, while LUD's implementation is obtained through private communication with its authors. In our comparisons, we use their default parameter settings. For ReSync, we set the initial step size to  $\mu_0 = 1/npq$  and the decaying factor to  $\gamma = 0.95$ , as suggested by the previous experiment. We fix  $n = 200$  and vary the true observation ratio  $p$  (or the observation ratio  $q$ ) while keeping  $q = 0.2$  (or  $p = 0.2$ ) fixed. We display the experiment results for  $\sigma = 0$  and  $\sigma = 1$  in Figures 3a and 3b, respectively, where  $p$  is selected from  $\{0.2, 0.3, 0.4, \dots, 1\}$ . When  $\sigma = 0$ , ReSync achieves competitive performance compared to other robust synchronization algorithms. When the additive noise level is  $\sigma = 1$ , ReSync outperforms other algorithms. In Figures 3c and 3d, we present the results with varying  $q$  chosen from  $\{0.2, 0.3, 0.4, \dots, 1\}$  for noise-free ( $\sigma = 0$ ) and noisy ( $\sigma = 1$ ) cases, respectively. In the noise-free case, DESC performs best when  $q < 0.5$ , while ReSync slightly outperforms others when  $q \geq 0.5$ . In the noisy case, it is clear that ReSync achieves the best performance for a large range of  $q$ .

## 4.2 Real Data

We consider the global alignment problem of three-dimensional scans from the Lucy dataset, which is a down-sampled version of the dataset containing 368 scans with a total number of 3.5 million triangles. We refer to [39] for more details about the experiment setting. We apply three algorithms LUD [39], DESC [32] and our ReSync on this dataset since they have the best performance on noisy synthetic data. As Figure 4 shows, ReSync outperforms the other two methods.

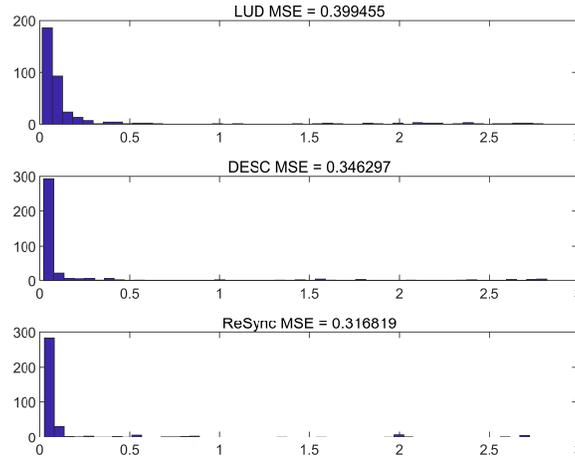


Figure 4: Histogram of the unsquared residuals of LUD, DESC, and ReSync for the Lucy dataset.

## 5 Conclusion and Discussions on Limitations

In this work, we introduced ReSync, a Riemannian subgradient-based algorithm with spectral initialization for solving RRS. We established strong theoretical results for ReSync under the RCM. In particular, we first presented an initialization guarantee for SpectrIn, which is a procedure embedded in ReSync for initialization. Then, we established a problem-intrinsic property called weak sharpness for our nonsmooth nonconvex formulation, which is of independent interest. Based on the established weak sharpness property, we derived linear convergence of ReSync to the underlying rotations once it is initialized in a local region. Combining these theoretical results demonstrates that ReSync converges linearly to the ground-truth rotations under the RCM.

**Limitations.** Our overall guarantee in Theorem 1 requires the sample complexity of  $p^7 q^2 = \Omega(\log n/n)$ , which does not match the currently best known lower bound  $p^2 q = \Omega(\log n/n)$  for exact recovery [34, 11]. We showed in Theorem 2 that approximate recovery with an optimal sample complexity is possible. Moreover, we showed in Theorem 3 that exact recovery with  $p^2 q^2 = \Omega(\log n/n)$  is possible if we have a global minimizer of the objective function of problem (2) within a certain local region. However, due to the nonconvexity of problem (2), it is non-trivial to obtain the said minimizer. We circumvented this difficulty by establishing the linear convergence of ReSync to a desired minimizer in Theorem 4. Nevertheless, a strong requirement on initialization is needed, which translates to the weaker final complexity result of  $p^7 q^2 = \Omega(\log n/n)$ .

Although our theory allows for  $p \rightarrow 0$  as  $n \rightarrow \infty$ , our argument relies heavily on the randomness of the outliers  $\{\mathcal{O}_{i,j}\}$  and the absence of additive noise. In practice, adversarial outliers that arbitrarily corrupt a measurement and additive noise contamination are prevalent. It remains unknown how well ReSync performs in such scenarios.

The above challenges are significant areas for future research and improvements.

## Acknowledgments and Disclosure of Funding

The authors thank Dr. Zengde Deng (Cainiao Network) and Dr. Shixiang Chen (University of Science and Technology of China) for providing helpful advice. They also thank the reviewers for their insightful comments, which have helped greatly to improve the quality and presentation of the manuscript.

Huikang Liu is supported in part by the National Natural Science Foundation of China (NSFC) Grant 72192832. Xiao Li is supported in part by the National Natural Science Foundation of China (NSFC) under grants 12201534 and 72150002, and in part by the Shenzhen Science and Technology Program under grants RCBS20210609103708017 and RYX20221008093033010. Anthony Man-Cho So is supported in part by the Hong Kong Research Grants Council (RGC) General Research Fund (GRF) project CUHK 14205421.

## References

- [1] Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, and Yiqiao Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. *Annals of Statistics*, 48(3):1452–1474, 2020.
- [2] Boris Alexeev, Afonso S Bandeira, Matthew Fickus, and Dustin G Mixon. Phase retrieval with polarization. *SIAM Journal on Imaging Sciences*, 7(1):35–66, 2014.
- [3] Mica Arie-Nachimson, Shahar Z Kovalsky, Ira Kemelmacher-Shlizerman, Amit Singer, and Ronen Basri. Global motion estimation from point matches. In *2012 2nd International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*, pages 81–88. IEEE, 2012.
- [4] Afonso S Bandeira. Random Laplacian matrices and convex relaxations. *Foundations of Computational Mathematics*, 18(2):345–379, 2018.
- [5] Afonso S Bandeira, Nicolas Boumal, and Amit Singer. Tightness of the maximum likelihood semidefinite relaxation for angular synchronization. *Mathematical Programming*, 163(1–2):145–167, 2017.
- [6] Afonso S Bandeira, Amit Singer, and Daniel A Spielman. A Cheeger inequality for the graph connection Laplacian. *SIAM Journal on Matrix Analysis and Applications*, 34(4):1611–1630, 2013.
- [7] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15(1):1455–1459, 2014.
- [8] James V Burke and Michael C Ferris. Weak sharp minima in mathematical programming. *SIAM Journal on Control and Optimization*, 31(5):1340–1359, 1993.
- [9] Avishek Chatterjee and Venu Madhav Govindu. Robust relative rotation averaging. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):958–972, 2017.
- [10] Yuxin Chen, Jianqing Fan, Cong Ma, and Kaizheng Wang. Spectral method and regularized MLE are both optimal for top- $K$  ranking. *Annals of Statistics*, 47(4):2204–2235, 2019.
- [11] Yuxin Chen and Andrea J Goldsmith. Information recovery from pairwise measurements. In *2014 IEEE International Symposium on Information Theory*, pages 2012–2016. IEEE, 2014.
- [12] Mihai Cucuringu, Yaron Lipman, and Amit Singer. Sensor network localization by eigenvector synchronization over the Euclidean group. *ACM Transactions on Sensor Networks*, 8(3):1–42, 2012.
- [13] Chandler Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- [14] Damek Davis, Dmitriy Drusvyatskiy, Kellie J MacPhee, and Courtney Paquette. Subgradient methods for sharp weakly convex functions. *Journal of Optimization Theory and Applications*, 179(3):962–982, 2018.

- [15] Shaofeng Deng, Shuyang Ling, and Thomas Strohmer. Strong consistency, graph Laplacians, and the stochastic block model. *Journal of Machine Learning Research*, 22(1):5210–5253, 2021.
- [16] John C Duchi and Feng Ruan. Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *Information and Inference: A Journal of the IMA*, 8(3):471–529, 2019.
- [17] Anders Eriksson, Carl Olsson, Fredrik Kahl, and Tat-Jun Chin. Rotation averaging with the chordal distance: Global minimizers and strong duality. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):256–268, 2019.
- [18] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An  $l_\infty$  eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018.
- [19] Tingran Gao and Zhizhen Zhao. Multi-frequency phase synchronization. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2132–2141. PMLR, 2019.
- [20] Venu Madhav Govindu. Lie-algebraic averaging for globally consistent motion estimation. In *CVPR 2004*, volume 1, pages I–I. IEEE, 2004.
- [21] Richard Hartley, Khurram Aftab, and Jochen Trunpf. L1 rotation averaging using the Weiszfeld algorithm. In *CVPR 2011*, pages 3041–3048. IEEE, 2011.
- [22] Richard Hartley, Jochen Trunpf, Yuchao Dai, and Hongdong Li. Rotation averaging. *International Journal of Computer Vision*, 103(3):267–305, 2013.
- [23] Gilad Lerman and Yunpeng Shi. Robust group synchronization via cycle-edge message passing. *Foundations of Computational Mathematics*, 22(6):1665–1741, 2022.
- [24] Xiao Li, Shixiang Chen, Zengde Deng, Qing Qu, Zhihui Zhu, and Anthony Man-Cho So. Weakly convex optimization over Stiefel manifold using Riemannian subgradient-type methods. *SIAM Journal on Optimization*, 31(3):1605–1634, 2021.
- [25] Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and René Vidal. Nonconvex robust low-rank matrix recovery. *SIAM Journal on Optimization*, 30(1):660–686, 2020.
- [26] Shuyang Ling. Near-optimal performance bounds for orthogonal and permutation group synchronization via spectral methods. *Applied and Computational Harmonic Analysis*, 60:20–52, 2022.
- [27] Huikang Liu, Man-Chung Yue, and Anthony Man-Cho So. A unified approach to synchronization problems over subgroups of the orthogonal group. *Applied and Computational Harmonic Analysis*, 66:320–372, 2023.
- [28] Daniel Martinec and Tomas Pajdla. Robust rotation and translation estimation in multiview reconstruction. In *CVPR 2007*, pages 1–8. IEEE, 2007.
- [29] Tyler Maunu and Gilad Lerman. Depth descent synchronization in  $SO(D)$ . *International Journal of Computer Vision*, 131(4):968–986, 2023.
- [30] David M Rosen, Luca Carlone, Afonso S Bandeira, and John J Leonard. SE-Sync: A certifiably correct algorithm for synchronization over the special euclidean group. *International Journal of Robotics Research*, 38(2–3):95–125, 2019.
- [31] Yunpeng Shi and Gilad Lerman. Message passing least squares framework and its application to rotation synchronization. In *Proceedings of the 37th International Conference on Machine Learning*, pages 8796–8806. PMLR, 2020.
- [32] Yunpeng Shi, Cole M Wyeth, and Gilad Lerman. Robust group synchronization via quadratic programming. In *Proceedings of the 39th International Conference on Machine Learning*, pages 20095–20105. PMLR, 2022.

- [33] Yoel Shkolnisky and Amit Singer. Viewing direction estimation in cryo-EM using synchronization. *SIAM Journal on Imaging Sciences*, 5(3):1088–1110, 2012.
- [34] Amit Singer. Angular synchronization by eigenvectors and semidefinite programming. *Applied and Computational Harmonic Analysis*, 30(1):20–36, 2011.
- [35] Amit Singer, Zhizhen Zhao, Yoel Shkolnisky, and Ronny Hadani. Viewing angle classification of cryo-electron microscopy images using eigenvectors. *SIAM Journal on Imaging Sciences*, 4(2):723–759, 2011.
- [36] Anthony Man-Cho So. Probabilistic analysis of the semidefinite relaxation detector in digital communications. In *Proceedings of the 21st annual ACM-SIAM Symposium on Discrete Algorithms*, pages 698–711. SIAM, 2010.
- [37] Gilbert W Stewart and Ji-guang Sun. *Matrix Perturbation Theory*. Academic Press, 1990.
- [38] Joel A Tropp. An Introduction to Matrix Concentration Inequalities. *Foundations and Trends® in Machine Learning*, 8(1–2):1–230, 2015.
- [39] Lanhui Wang and Amit Singer. Exact and stable recovery of rotations for robust synchronization. *Information and Inference: A Journal of the IMA*, 2(2):145–193, 2013.
- [40] Wei Hong Yang, Lei-Hong Zhang, and Ruyi Song. Optimality conditions for the nonlinear programming problems on Riemannian manifolds. *Pacific Journal of Optimization*, 10(2):415–434, 2014.
- [41] Stella Yu. Angular embedding: A robust quadratic criterion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(1):158–173, 2011.
- [42] Yiqiao Zhong and Nicolas Boumal. Near-optimal bounds for phase synchronization. *SIAM Journal on Optimization*, 28(2):989–1016, 2018.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Algorithm and Setup</b>	<b>3</b>
2.1	ReSync: Algorithm Development . . . . .	3
2.2	RCM Setup for Theoretical Analysis . . . . .	4
<b>3</b>	<b>Main Results</b>	<b>5</b>
3.1	Analysis of SpectrIn with Leave-One-Out Technique . . . . .	5
3.2	Weak Sharpness and Exact Recovery . . . . .	6
3.3	Convergence Analysis and Proof of Theorem 1 . . . . .	8
<b>4</b>	<b>Experiments</b>	<b>8</b>
4.1	Synthetic Data . . . . .	8
4.2	Real Data . . . . .	10
<b>5</b>	<b>Conclusion and Discussions on Limitations</b>	<b>10</b>
<b>A</b>	<b>Full Proof of Theorem 2</b>	<b>15</b>
A.1	Initialization Error in Euclidean Norm . . . . .	16
A.2	Initialization Error in Infinity Norm . . . . .	17
<b>B</b>	<b>Full Proof of Theorem 3</b>	<b>18</b>
B.1	Proof of Lemma 1 . . . . .	18
B.2	Proof of Proposition 1 . . . . .	18
B.3	Proof of Proposition 2 . . . . .	19
<b>C</b>	<b>Full Proof of Theorem 4</b>	<b>21</b>
C.1	Proof of Corollary 1 . . . . .	21
C.2	Proof of Contraction . . . . .	22

## A Full Proof of Theorem 2

In this section, we present the full proof of Theorem 2. We will use the notation defined in the proof outline of Theorem 2 in Section 3.1.

Firstly, we know that

$$\mathbf{W}_{ij} = \begin{cases} (1-pq)\mathbf{X}_i^* \mathbf{X}_j^{*\top}, & \text{with probability } pq, \\ \mathbf{O}_{ij} - pq\mathbf{X}_i^* \mathbf{X}_j^{*\top}, & \text{with probability } (1-p)q, \\ -pq\mathbf{X}_i^* \mathbf{X}_j^{*\top}, & \text{otherwise.} \end{cases} \quad (17)$$

Since  $\mathbf{O}_{ij}$  is assumed to be uniformly distributed on  $\text{SO}(d)$  in the RCM, given any matrix  $\mathbf{Q} \in \text{SO}(d)$ , it is easy to see that  $\mathbf{O}_{ij}\mathbf{Q}$  is also uniformly distributed on  $\text{SO}(d)$ , so we have

$$\mathbb{E}(\mathbf{O}_{ij}\mathbf{Q}) = \mathbb{E}(\mathbf{O}_{ij}) \cdot \mathbf{Q}, \quad \forall \mathbf{Q} \in \text{SO}(d).$$

Let  $\mathbf{E}_{kl} \in \text{SO}(d)$ ,  $k \neq l$  denote the diagonal matrix whose  $k$ -th and  $l$ -th diagonal entries are  $-1$  and others are  $1$ , then we have  $\mathbb{E}(\mathbf{O}_{ij}) = \mathbb{E}(\mathbf{O}_{ij}) \cdot \mathbf{E}_{kl}$ , which implies

$$d\mathbb{E}(\mathbf{O}_{ij}) = \mathbb{E}(\mathbf{O}_{ij}) \cdot (\mathbf{E}_{12} + \mathbf{E}_{23} + \cdots + \mathbf{E}_{d1}) = (d-4)\mathbb{E}(\mathbf{O}_{ij}).$$

Thus, we have  $\mathbb{E}(\mathbf{O}_{ij}) = \mathbf{0}$ . Then, it is easy to see that  $\mathbb{E}(\mathbf{W}_{ij}) = \mathbf{0}$  and

$$\text{Var}(\mathbf{W}_{ij}) = (1-pq)^2\mathbf{I} + (1-q)p^2q^2\mathbf{I} + (1-p)q(1+p^2q^2)\mathbf{I} = q(1-p^2q)\mathbf{I}. \quad (18)$$

Based on the above calculations, we can derive the following Lemma, which is a direct result of the matrix Bernstein inequality [38]. Similar results can also be found in [42, Lemma 9], [27, Proposition 3], and [26, Eq. (5.12)].

**Lemma 2.** *With probability at least  $1 - \mathcal{O}(1/n)$ , the following holds for any  $m \in [n]$ :*

$$\|\mathbf{W}\|_2 = \mathcal{O}\left(\sqrt{nq \log n}\right), \quad \|\mathbf{W}^{(m)}\|_2 = \mathcal{O}\left(\sqrt{nq \log n}\right), \quad \|\mathbf{W}_m\|_2 = \mathcal{O}\left(\sqrt{nq \log n}\right).$$

*Proof.* Note that  $\mathbf{W} = \sum_{i < j} \mathbf{W}^{(ij)}$ , where  $\mathbf{W}^{(ij)} \in \mathbb{R}^{nd \times nd}$  denotes the matrix with the  $(i, j)$  and  $(j, i)$ -th block equal to  $\mathbf{W}_{ij}$  and  $\mathbf{W}_{ji}$  and others equal to 0. So  $\{\mathbf{W}^{(ij)}\}$  are i.i.d. centered and bounded random matrices. Besides, we have

$$\left\| \mathbb{E}(\mathbf{W}\mathbf{W}^\top) \right\|_2 = \left\| \sum_{i < j} \mathbb{E}(\mathbf{W}^{(ij)}(\mathbf{W}^{(ij)})^\top) \right\|_2 = 2(n-1) \|\text{Var}(\mathbf{W}_{ij})\|_2 = \mathcal{O}(nq).$$

According to the matrix Bernstein inequality [38], we have that  $\|\mathbf{W}\|_2 = \mathcal{O}\left(\sqrt{nq \log n}\right)$  holds with probability at least  $1 - \mathcal{O}(1/n^2)$ . The above argument also holds for each  $\mathbf{W}^{(m)}$ , then taking a union bound over the choice of  $m \in [n]$  yields the second result. For  $\mathbf{W}_m$ , we have

$$\|\mathbf{W}\|_2 = \max_{\|\mathbf{u}\|_2=1} \|\mathbf{W}\mathbf{u}\|_2 \geq \max_{\|\mathbf{u}\|_2=1} \|\mathbf{W}_m\mathbf{u}\|_2 = \|\mathbf{W}_m\|_2,$$

which gives the last result.  $\square$

The following lemma follows from [13]; see also [26, Theorem A.2]. [42, Lemma 11] is a special case where  $d = 1$ . Before that, we need to introduce the distance up to a global orthogonal matrix:

$$\text{Dist}(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\mathbf{R}^*\|_F, \quad \text{where } \mathbf{R}^* = \arg \min_{\mathbf{R} \in \mathcal{O}(d)} \|\mathbf{X}\mathbf{R} - \mathbf{Y}\|_F^2 = \mathcal{P}_{\mathcal{O}(d)}(\mathbf{X}^\top \mathbf{Y}).$$

**Lemma 3** (Davis-Kahan sin  $\Theta$  Theorem). *Suppose that  $\mathbf{A}, \mathbf{E} \in \mathbb{C}^{n \times n}$  are Hermitian matrices and  $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{E}$ . Let  $\delta = \lambda_d(\mathbf{A}) - \lambda_{d+1}(\mathbf{A}) > 0$  be the gap between the  $d$ -th eigenvalue and  $d+1$ -th eigenvalue of  $\mathbf{A}$  for some  $1 \leq d \leq n-1$ . Furthermore, let  $\mathbf{U}, \tilde{\mathbf{U}}$  be the  $d$ -leading eigenvectors of  $\mathbf{A}$  and  $\hat{\mathbf{A}}$ , respectively, which are normalized such that  $\|\mathbf{U}\|_F = \|\tilde{\mathbf{U}}\|_F = \sqrt{nd}$ . Then, we have*

$$\text{Dist}(\mathbf{U}, \tilde{\mathbf{U}}) \leq \frac{\sqrt{2}\|\mathbf{E}\mathbf{U}\|_F}{\delta - \|\mathbf{E}\|_2}. \quad (19)$$

### A.1 Initialization Error in Euclidean Norm

Based on Lemma 2 and Lemma 3, we have the following bound on Dist. Note that the notations  $\Phi, \Psi, \tilde{\Phi}, \tilde{\Psi}$  used in the following analysis are defined in SpectrIn (i.e., Algorithm 2).

**Lemma 4.** Let  $\Phi$  and  $\Phi^{(m)}$  be the  $d$ -leading eigenvectors of  $\mathbf{Y}$  and  $\mathbf{Y}^{(m)}$ , respectively, which are normalized such that  $\|\Phi\|_F = \|\Phi^{(m)}\|_F = \sqrt{nd}$ . Then, we have

$$\text{Dist}(\Phi, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right) \quad \text{and} \quad \text{Dist}(\Phi, \Phi^{(m)}) = \mathcal{O}(1) \quad (20)$$

hold with probability at least  $1 - \mathcal{O}(1/n)$ .

*Proof.* Let us Choose  $\mathbf{A} = \mathbf{Y}^{(m)}$  and  $\mathbf{E} = \Delta\mathbf{W}^{(m)} = \mathbf{W} - \mathbf{W}^{(m)}$  in Lemma 3, then  $\tilde{\mathbf{A}} = \mathbf{Y}$ ,  $\mathbf{U} = \Phi^{(m)}$  and  $\tilde{\mathbf{U}} = \Phi$ . Since  $\Phi^{(m)}$  is independent of  $\Delta\mathbf{W}^{(m)}$ , similar to Lemma 2, we apply the matrix Bernstein inequality [38] to obtain, with probability at least  $1 - \mathcal{O}(1/n^2)$ , that

$$\|\mathbf{E}\mathbf{U}\|_F = \|\Delta\mathbf{W}^{(m)}\Phi^{(m)}\|_F = \mathcal{O}(\sqrt{nq \log n}).$$

In addition,  $\mathbf{Y}^{(m)} = pq\mathbf{X}^*\mathbf{X}^{*\top} + \mathbf{W}^{(m)}$  implies that

$$\delta = \lambda_d(\mathbf{Y}^{(m)}) - \lambda_{d+1}(\mathbf{Y}^{(m)}) \geq \lambda_d(pq\mathbf{X}^*\mathbf{X}^{*\top}) - \|\mathbf{W}^{(m)}\|_2 \geq npq - \mathcal{O}(\sqrt{nq \log n}).$$

where the second inequality holds due to  $\lambda_d(\mathbf{X}^*\mathbf{X}^{*\top}) = n$  and the last inequality is from  $\lambda_d(\mathbf{X}^*\mathbf{X}^{*\top}) = n$  and Lemma 2. Based on the condition that  $p^2q = \Omega\left(\frac{\log n}{n}\right)$ , as long as  $p^2q \geq \frac{C \log n}{n}$  for some large enough constant  $C$ , we have

$$\delta - \|\mathbf{E}\|_2 \geq npq - \mathcal{O}(\sqrt{nq \log n}) - \|\mathbf{E}\|_2 \geq npq - \mathcal{O}(\sqrt{nq \log n}) \geq \frac{1}{2}npq,$$

where the second inequality holds because of  $\|\mathbf{E}\|_2 \leq \|\mathbf{W}\|_2 + \|\mathbf{W}^{(m)}\|_2 = \mathcal{O}(\sqrt{nq \log n})$ . Hence, by applying Lemma 3, we get

$$\text{Dist}(\Phi, \Phi^{(m)}) \leq \frac{\sqrt{2}\|\mathbf{E}\mathbf{U}\|_F}{\delta - \|\mathbf{E}\|_2} \leq \frac{\mathcal{O}(\sqrt{nq \log n})}{npq} = \mathcal{O}(1) \quad (21)$$

where the last inequality is because of  $p\sqrt{q} = \Omega(\sqrt{\log n/n})$ . Similarly, by choosing  $\mathbf{A} = pq\mathbf{X}^*\mathbf{X}^{*\top}$  and  $\mathbf{E} = \mathbf{W}$ , we can show that

$$\text{Dist}(\Phi, \mathbf{X}^*) \leq \frac{\sqrt{2}\|\mathbf{W}\mathbf{X}^*\|_F}{npq - \mathcal{O}(\sqrt{nq \log n})} \leq \frac{\sqrt{2}\|\mathbf{W}\|_2\|\mathbf{X}^*\|_F}{\frac{1}{2}npq} = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right).$$

Here, the last inequality holds because  $\|\mathbf{W}\|_2 = \mathcal{O}(\sqrt{nq \log n})$  (see Lemma 2) and the fact that  $\|\mathbf{X}^*\|_F = \sqrt{nd}$ .  $\square$

Following the same analysis as in Lemma 4, we are also able to show that  $\text{Dist}(\Psi, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right)$ , where  $\Psi$  reverses the sign of the last column of  $\Phi$  so that the determinants of  $\Phi$  and  $\Psi$  differ by a sign. However, Lemma 4 only provides the upper bound on ‘‘Dist’’, i.e., the distance up to an orthogonal matrix. The following lemma translates the result to that on ‘‘dist’’.

**Lemma 5.** Suppose that  $\|\tilde{\Phi} - \Phi\|_F \leq \|\tilde{\Psi} - \Psi\|_F$ , where  $\tilde{\Phi} = \mathcal{P}_{\text{SO}(d)^n}(\Phi)$  and  $\tilde{\Psi} = \mathcal{P}_{\text{SO}(d)^n}(\Psi)$ , then we have

$$\text{dist}(\Phi, \mathbf{X}^*) = \text{Dist}(\Phi, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right).$$

*Proof.* Based on the structure of  $O(d)$  and  $\text{SO}(d)$ , it is easy to see that

$$\text{Dist}(\Phi, \mathbf{X}^*) = \min\{\text{dist}(\Phi, \mathbf{X}^*), \text{dist}(\Psi, \mathbf{X}^*)\}.$$

Our remaining task is to prove  $\text{Dist}(\Phi, \mathbf{X}^*) = \text{dist}(\Phi, \mathbf{X}^*)$  based on the condition that  $\|\tilde{\Phi} - \Phi\|_F \leq \|\tilde{\Psi} - \Psi\|_F$ . If  $\text{Dist}(\Phi, \mathbf{X}^*) = \text{dist}(\Psi, \mathbf{X}^*)$ , then we have

$$\mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right) = \text{Dist}(\Phi, \mathbf{X}^*) = \text{dist}(\Psi, \mathbf{X}^*) \geq \|\tilde{\Psi} - \Psi\|_F \geq \|\tilde{\Phi} - \Phi\|_F$$

where the first inequality holds because  $\tilde{\Psi}$  is the projection of  $\Psi$  on  $\text{SO}(d)^n$ . It is easy to see that

$$\|\tilde{\Phi} - \Phi\|_F + \|\tilde{\Psi} - \Psi\|_F = \|\Phi - \mathcal{P}_{\text{SO}(d)^n}(\Phi)\|_F + \|\Phi - \mathcal{P}_{(O(d) \setminus \text{SO}(d))^n}(\Phi)\|_F \geq 2\sqrt{n}.$$

The equality holds because the mapping that reverses the sign of the last column of a matrix is a bijection between  $\text{SO}(d)$  and  $O(d) \setminus \text{SO}(d)$ , and the inequality holds since the minimum distance between  $\text{SO}(d)$  and  $O(d) \setminus \text{SO}(d)$  is 2. Under the condition that  $p^2q = \Omega\left(\frac{\log n}{n}\right)$ , as long as  $p^2q \geq \frac{C \log n}{n}$  for some large enough constant  $C$ , we could have  $\|\tilde{\Phi} - \Phi\|_F + \|\tilde{\Psi} - \Psi\|_F = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right) < 2\sqrt{n}$ , which contradict to the above inequality. Thus, we have

$$\text{dist}(\Phi, \mathbf{X}^*) = \text{Dist}(\Phi, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{q}}\right).$$

□

## A.2 Initialization Error in Infinity Norm

Next, based on the Davis-Kahan theorem, we can bound the distance between  $\tilde{\mathbf{X}}^0$  and  $\mathbf{X}^*$  in the infinity norm, which is stated in the following result.

**Lemma 6.** *Let  $\Phi$  be the  $d$ -leading eigenvectors of  $\mathbf{Y}$ . Then, we have*

$$\text{dist}_\infty(\Phi, \mathbf{X}^*) = \mathcal{O}\left(\frac{\sqrt{\log n}}{p\sqrt{nq}}\right) \quad (22)$$

holds with probability at least  $1 - \mathcal{O}(1/n)$ .

*Proof.* Let  $\Pi^m = \text{argmin}_{\Pi \in \text{SO}(d)} \|\Phi - \Phi^{(m)}\Pi\|_F$ . We can compute

$$\begin{aligned} \|(\mathbf{W}\Phi)_m\|_F &= \|\mathbf{W}_m^\top \Phi\|_F \leq \|\mathbf{W}_m^\top \Phi^{(m)}\Pi^m\|_F + \|\mathbf{W}_m^\top (\Phi - \Phi^{(m)})\Pi^m\|_F \\ &\leq \|\mathbf{W}_m^\top \Phi^{(m)}\|_F + \|\mathbf{W}_m\|_2 \|\Phi - \Phi^{(m)}\Pi^m\|_F. \end{aligned} \quad (23)$$

The fact that  $\Phi^{(m)}$  is independent from  $\mathbf{W}_m$  implies  $\|\mathbf{W}_m^\top \Phi^{(m)}\|_F = \mathcal{O}(\sqrt{nq \log n})$ , so we have

$$\|(\mathbf{W}\Phi)_m\|_F = \mathcal{O}(\sqrt{nq \log n}) (1 + \mathcal{O}(1)) = \mathcal{O}(\sqrt{nq \log n}). \quad (24)$$

Next, let  $\Pi^* = \text{argmin}_{\Pi \in \text{SO}(d)} \|\Phi - \mathbf{X}^*\Pi\|_F$ , then we have

$$(\mathbf{Y}\Phi)_m = pq\mathbf{X}_m^* \mathbf{X}^{*\top} \Phi + (\mathbf{W}\Phi)_m = npq\mathbf{X}_m^* \Pi^* + pq\mathbf{X}_m^* \mathbf{X}^{*\top} (\Phi - \mathbf{X}^* \Pi^*) + (\mathbf{W}\Phi)_m.$$

Since  $\Phi$  is the  $d$ -leading eigenvector of  $\mathbf{Y}$ , we have  $\mathbf{Y}\Phi = \Phi\Sigma$  with  $\Sigma \in \mathbb{R}^{d \times d}$  consisting of the  $d$ -leading eigenvalues of  $\mathbf{Y}$ . Applying the standard eigenvalue perturbation theory, e.g., [37, Theorem 4.11], we have

$$\|\Sigma - npq\mathbf{I}\|_2 = \mathcal{O}(\|\mathbf{W}\|_2) = \mathcal{O}(\sqrt{nq \log n}).$$

Based on the assumption that  $p^2q = \Omega\left(\frac{\log n}{n}\right)$ , we can show that  $\Sigma \geq \frac{3}{4}npq\mathbf{I}$ . Note that  $\mathbf{Y}\Phi = \Phi\Sigma$  implies  $\Phi_m = (\mathbf{Y}\Phi)_m \Sigma^{-1}$  for each  $m \in [n]$ , then we have

$$\|\Phi_m\|_F \leq \frac{4}{3npq} \|(\mathbf{Y}\Phi)_m\|_F \leq \frac{4}{3npq} \|pq\mathbf{X}_m^* \mathbf{X}^{*\top} \Phi_m\|_F + \frac{4}{3npq} \|(\mathbf{W}\Phi)_m\|_F \leq 2\sqrt{d},$$

where the last inequality holds because  $\frac{4}{3npq} \|pq\mathbf{X}_m^* \mathbf{X}^{*\top} \Phi_m\|_F = \frac{4}{3n} \|\mathbf{X}^{*\top} \Phi_m\|_F \leq \frac{4\sqrt{d}}{3}$  and (24). Therefore, for each  $m \in [n]$ ,

$$\begin{aligned} npq(\Phi_m - \mathbf{X}_m^* \Pi^*) &= \Phi_m(npq\mathbf{I} - \Sigma) + \Phi_m \Sigma - npq\mathbf{X}_m^* \Pi^* \\ &= \Phi_m(npq\mathbf{I} - \Sigma) + (\mathbf{Y}\Phi)_m - npq\mathbf{X}_m^* \Pi^*. \end{aligned}$$

This further implies

$$\begin{aligned} npq\|\Phi_m - \mathbf{X}_m^* \Pi^*\|_F &\leq \|\Sigma - npq\mathbf{I}\|_2 \|\Phi_m\|_F + pq\|\mathbf{X}_m^* \mathbf{X}^{*\top} (\Phi - \mathbf{X}^* \Pi^*)\|_F + \|(\mathbf{W}\Phi)_m\|_F \\ &\leq 2\sqrt{d}\|\mathbf{W}\|_2 + \sqrt{n}pq\|\Phi - \mathbf{X}^* \Pi^*\|_F + \mathcal{O}(\sqrt{nq \log n}) \\ &= \mathcal{O}(\sqrt{nq \log n}), \end{aligned}$$

where the last inequality holds because of Lemma 2 and Lemma 4. This completes the proof. □

Finally, since  $\mathbf{X}^0 = \mathcal{P}_{\text{SO}(d)^n}(\Phi)$ , according to Lemma 2 in [27], we have

$$\text{dist}(\mathbf{X}^0, \mathbf{X}^*) \leq 2 \text{dist}(\Phi, \mathbf{X}^*) \quad \text{and} \quad \text{dist}_\infty(\mathbf{X}^0, \mathbf{X}^*) \leq 2 \text{dist}_\infty(\Phi, \mathbf{X}^*),$$

which complete the proof of Theorem 2.

## B Full Proof of Theorem 3

In this section, we provide the full proof of Lemma 1, Proposition 1, and Proposition 2, which finishes the proof of Theorem 3. To simplify the notation and the theoretical derivations, we assume without loss of generality that  $\mathbf{X}_i^* = \mathbf{I}$  for all  $1 \leq i \leq n$  as one can separately rotate the space that each variable  $\mathbf{X}_i$  lies in such that the corresponding ground-truth  $\mathbf{X}_i^*$  is rotated to identity [39, Lemma 4.1]. Consequently, we have  $\mathbf{Y}_{ij} = \mathbf{I}$  for  $(i, j) \in \mathcal{A}$ .

### B.1 Proof of Lemma 1

For each  $1 \leq i \leq n$ ,  $|\mathcal{E}_i| = \sum_{1 \leq j \leq n} \mathbf{1}_{\mathcal{E}_i}(j)$ , where  $\mathbf{1}_{\mathcal{E}_i}(\cdot)$  denotes the indicator function w.r.t  $\mathcal{E}_i$ . Based on our model,  $\sum_{1 \leq j \leq n} \mathbf{1}_{\mathcal{E}_i}(j)$  follows the binomial distribution  $B(n, q)$ . According to the Bernstein inequality [38], for any constant  $\epsilon \in (0, 1)$ , we have

$$\begin{aligned} \Pr \left( \left| \sum_{1 \leq j \leq n} \mathbf{1}_{\mathcal{E}_i}(j) - nq \right| \geq \epsilon nq \right) &\leq 2 \exp \left( - \frac{\frac{1}{2} \epsilon^2 n^2 q^2}{\sum_{1 \leq j \leq n} E\{\mathbf{1}_{\mathcal{E}_i}^2(j)\} + \epsilon nq/3} \right) \\ &= 2 \exp \left( - \frac{\frac{1}{2} \epsilon^2 n^2 q^2}{nq + \epsilon nq/3} \right) \leq 2 \exp \left( - \frac{3}{8} \epsilon^2 nq \right). \end{aligned}$$

The last inequality holds because of  $\epsilon < 1$ . Therefore,

$$\Pr \left( \bigcup_{1 \leq i \leq n} \{|\mathcal{E}_i| - nq\} \geq \epsilon nq \right) \leq 2n \exp \left( - \frac{3}{8} \epsilon^2 nq \right) \leq \frac{2}{n^2}, \quad (25)$$

where the last inequality holds because we assume that  $\epsilon \geq \frac{\sqrt{8 \log n}}{\sqrt{npq}}$ . Similarly, we have

$$\Pr \left( \bigcup_{1 \leq i \leq n} \{|\mathcal{A}_i| - npq\} \geq \epsilon npq \right) \leq 2n \exp \left( - \frac{3}{8} \epsilon^2 npq \right) \leq \frac{2}{n^2}, \quad (26)$$

$$\Pr \left( \bigcup_{1 \leq i, j \leq n} \{|\mathcal{E}_i \cap \mathcal{A}_j| - npq^2\} \geq \epsilon npq^2 \right) \leq 2n^2 \exp \left( - \frac{3}{8} \epsilon^2 npq^2 \right) \leq \frac{2}{n}, \quad (27)$$

and

$$\Pr \left( \bigcup_{1 \leq i, j \leq n} \{|\mathcal{A}_{ij}| - np^2q^2\} \geq \epsilon np^2q^2 \right) \leq 2n^2 \exp \left( - \frac{3}{8} \epsilon^2 np^2q^2 \right) \leq \frac{2}{n}. \quad (28)$$

Hence, we complete the proof of Lemma 1 once  $n \geq 4$ .

### B.2 Proof of Proposition 1

We can first compute

$$\begin{aligned} \sum_{1 \leq i, j \leq n} \|\mathbf{X}_i - \mathbf{X}_j\|_F &\leq \sum_{1 \leq i, j \leq n} \frac{1}{|\mathcal{A}_{ij}|} \sum_{k \in \mathcal{A}_{ij}} (\|\mathbf{X}_i - \mathbf{X}_k\|_F + \|\mathbf{X}_j - \mathbf{X}_k\|_F) \\ &\leq \frac{1}{(1 - \epsilon)np^2q^2} \sum_{1 \leq i, j \leq n} \sum_{k \in \mathcal{A}_{ij}} (\|\mathbf{X}_i - \mathbf{X}_k\|_F + \|\mathbf{X}_j - \mathbf{X}_k\|_F). \end{aligned} \quad (29)$$

Here, the first inequality comes from the triangle inequality, while the second one follows from Lemma 1. Now, invoking Lemma 1, which tells  $|\mathcal{A}_k| \leq (1 + \epsilon)npq$ , gives

$$\begin{aligned} \sum_{1 \leq i, j \leq n} \sum_{k \in \mathcal{A}_{ij}} \|\mathbf{X}_i - \mathbf{X}_k\|_F &= \sum_{1 \leq i \leq n} \sum_{k \in \mathcal{A}_i} \sum_{j \in \mathcal{A}_k} \|\mathbf{X}_i - \mathbf{X}_k\|_F \\ &\leq (1 + \epsilon)npq \sum_{1 \leq i \leq n} \sum_{k \in \mathcal{A}_i} \|\mathbf{X}_i - \mathbf{X}_k\|_F \end{aligned}$$

$$= (1 + \epsilon)npq \sum_{(i,k) \in \mathcal{A}} \|\mathbf{X}_i - \mathbf{X}_k\|_F.$$

By symmetry, we conclude that

$$\sum_{1 \leq i, j \leq n} \sum_{k \in \mathcal{A}_{ij}} (\|\mathbf{X}_i - \mathbf{X}_k\|_F + \|\mathbf{X}_j - \mathbf{X}_k\|_F) \leq 2(1 + \epsilon)npq \sum_{(i,j) \in \mathcal{A}} \|\mathbf{X}_i - \mathbf{X}_j\|_F. \quad (30)$$

Furthermore, we claim the following bound for any  $\mathbf{X} \in \text{SO}(d)^n$ :

$$\sum_{1 \leq i, j \leq n} \|\mathbf{X}_i - \mathbf{X}_j\|_F \geq \frac{n}{2} \text{dist}_1(\mathbf{X}, \mathbf{X}^*). \quad (31)$$

Combining (29), (30), (31), and the fact  $g(\mathbf{X}) = \sum_{(i,j) \in \mathcal{A}} \|\mathbf{X}_i - \mathbf{X}_j\|_F$  establishes Proposition 1.

Hence, it remains to show (31). First of all, according to the triangle inequality, we have

$$\sum_{1 \leq i, j \leq n} \|\mathbf{X}_i - \mathbf{X}_j\|_F \geq \sum_{1 \leq i \leq n} \left\| n\mathbf{X}_i - \sum_{1 \leq j \leq n} \mathbf{X}_j \right\|_F = n \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \bar{\mathbf{X}}\|_F, \quad (32)$$

where  $\bar{\mathbf{X}} = \frac{1}{n} \sum_{1 \leq j \leq n} \mathbf{X}_j$  can be taken as a diagonal matrix (since the Frobenius norm is invariant up to a global rotation) with its first  $(d - 1)$  diagonal entries being positive. Finally, applying the following lemma to (32) provides (31).

**Lemma 7.** For any  $\mathbf{A} \in \text{SO}(d)$  and  $\mathbf{B} = \text{Diag}(b_1, \dots, b_d)$  satisfying  $b_1, \dots, b_{d-1} \in [0, 1]$  and  $b_d \in [-1, 1]$ , we have

$$\|\mathbf{A} - \mathbf{B}\|_F \geq \frac{1}{2} \|\mathbf{A} - \mathbf{I}\|_F.$$

*Proof.* It is equivalent to show that  $\|\mathbf{A} - \mathbf{B}\|_F^2 \geq \frac{1}{4} \|\mathbf{A} - \mathbf{I}\|_F^2$ . Since  $\mathbf{A} \in \text{SO}(d)$ , we simplify as

$$d + \langle \mathbf{A}, \mathbf{I} - 4\mathbf{B} \rangle + 2\|\mathbf{B}\|_F^2 = \sum_{1 \leq i \leq d} 1 + \mathbf{A}_{ii}(1 - 4b_i) + 2b_i^2 \geq 0. \quad (33)$$

To prove (33) holds for all  $\mathbf{A} \in \text{SO}(d)$ , we choose  $\bar{\mathbf{A}} = \text{argmin}_{\mathbf{A} \in \text{SO}(d)} \langle \mathbf{A}, \mathbf{I} - 4\mathbf{B} \rangle = \mathcal{P}_{\text{SO}(d)}(4\mathbf{B} - \mathbf{I})$ . It is easy to see that  $\bar{\mathbf{A}}$  is also a diagonal matrix. For any  $1 \leq i \leq d - 1$ , we have

$$1 + \mathbf{A}_{ii}(1 - 4b_i) + 2b_i^2 = \begin{cases} 2(b_i - 1)^2, & \mathbf{A}_{ii} = 1; \\ 2b_i^2 + 4b_i, & \mathbf{A}_{ii} = -1. \end{cases}$$

So it is always nonnegative since  $b_i \geq 0$  for any  $1 \leq i \leq d - 1$ . For the last summation in (33), on the one hand, if  $\bar{\mathbf{A}}_{dd} = 1$ , then  $1 + \bar{\mathbf{A}}_{dd}(1 - 4b_d) + 2b_d^2 = 2(b_d - 1)^2 \geq 0$ . On the other hand, if  $\bar{\mathbf{A}}_{dd} = -1$ , then there exist another  $1 \leq k \leq d - 1$  such that  $\bar{\mathbf{A}}_{kk} = -1$ . So we have  $(1 + \bar{\mathbf{A}}_{kk}(1 - 4b_k) + 2b_k^2) + (1 + \bar{\mathbf{A}}_{dd}(1 - 4b_d) + 2b_d^2) = 2b_k^2 + 2b_d^2 + 4(b_k + b_d) \geq 0$ . The last inequality holds because  $|b_d| \leq b_k$ . We complete the proof.  $\square$

### B.3 Proof of Proposition 2

Based on (11) and our simplification that  $\mathbf{X}_i^* = \mathbf{I}$ ,  $1 \leq i \leq n$ , we have

$$h(\mathbf{X}) - h(\mathbf{X}^*) \geq \sum_{(i,j) \in \mathcal{A}^c} \left\langle \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \right\rangle. \quad (34)$$

Let  $\mathbf{A} = \mathbb{E} \left\{ \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\}$ ,  $\mathbf{Z}_{ij} = \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \cdot \mathbf{1}_{\mathcal{A}^c}(ij) - (1 - p)q\mathbf{A}$ , and  $\mathbf{Z} \in \mathbb{R}^{nd \times nd}$  collects each  $\mathbf{Z}_{ij}$  in its  $(i, j)$ -th block. We can compute

$$\begin{aligned} & \left| \sum_{(i,j) \in \mathcal{A}^c} \left\langle \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \right\rangle \right| = \left| \sum_{1 \leq i, j \leq n} \langle \mathbf{Z}_{ij} + (1 - p)q\mathbf{A}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \rangle \right| \\ & \leq \left| \sum_{1 \leq i, j \leq n} \langle (1 - p)q\mathbf{A}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \rangle \right| \end{aligned} \quad (35)$$

$$+ \left| \sum_{1 \leq i, j \leq n} \langle \mathbf{Z}_{ij}, (\mathbf{X}_i - \mathbf{I})^\top (\mathbf{X}_j - \mathbf{I}) + (\mathbf{X}_i - \mathbf{I})^\top + (\mathbf{X}_j - \mathbf{I}) \rangle \right|.$$

To bound the first term in (35), we can proceed as

$$\begin{aligned} \left| \sum_{1 \leq i, j \leq n} \langle (1-p)q\mathbf{A}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \rangle \right| &= \left| \left\langle (1-p)q\mathbf{A}, \left( \sum_{i=1}^n \mathbf{X}_i \right)^\top \left( \sum_{i=1}^n \mathbf{X}_i \right) - n^2 \mathbf{I} \right\rangle \right| \\ &= \left| \left\langle (1-p)q\mathbf{A}, \left( \sum_{i=1}^n \mathbf{X}_i + n\mathbf{I} \right)^\top \left( \sum_{i=1}^n \mathbf{X}_i - n\mathbf{I} \right) \right\rangle \right| \\ &\leq \left\| (1-p)q\mathbf{A} \left( \sum_{i=1}^n \mathbf{X}_i + n\mathbf{I} \right) \right\|_F \left\| \sum_{i=1}^n \mathbf{X}_i - n\mathbf{I} \right\|_F \\ &\leq \|(1-p)q\mathbf{A}\|_F \left\| \sum_{i=1}^n \mathbf{X}_i + n\mathbf{I} \right\|_2 \left\| \sum_{i=1}^n \mathbf{X}_i - n\mathbf{I} \right\|_F \\ &\leq 2(1-p)qn \|\mathbf{A}\|_F \left\| \sum_{i=1}^n \mathbf{X}_i - n\mathbf{I} \right\|_F. \end{aligned} \tag{36}$$

Here, the second equality is true since  $\sum_{i=1}^n \mathbf{X}_i$  can be taken as a diagonal matrix. According to [39, Lemma A.1], we know that  $\|\mathbf{A}\|_F = \left\| \mathbb{E} \left\{ \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\} \right\|_F \leq \frac{1}{\sqrt{2}}$  for all  $d \geq 2$ . On the other hand, we know that

$$\text{dist}(\mathbf{X}, \mathbf{X}^*)^2 = \sum_{i=1}^n \|\mathbf{X}_i - \mathbf{I}\|_F^2 = 2 \text{trace} \left( n\mathbf{I} - \sum_{i=1}^n \mathbf{X}_i \right) \geq 2 \left\| \sum_{i=1}^n \mathbf{X}_i - n\mathbf{I} \right\|_F, \tag{37}$$

where the last inequality holds because  $n\mathbf{I} - \sum_{i=1}^n \mathbf{X}_i$  is a nonnegative diagonal matrix. Therefore, we obtain

$$\begin{aligned} \left| \sum_{1 \leq i, j \leq n} \langle (1-p)q\mathbf{A}, \mathbf{X}_i^\top \mathbf{X}_j - \mathbf{I} \rangle \right| &\leq \frac{(1-p)qn}{\sqrt{2}} \text{dist}(\mathbf{X}, \mathbf{X}^*)^2 \\ &\leq \frac{(1-p)qn}{\sqrt{2}} \max_i \|\mathbf{X}_i - \mathbf{I}\|_F \sum_{i=1}^n \|\mathbf{X}_i - \mathbf{I}\|_F. \end{aligned} \tag{38}$$

To further bound the second term in (35), we can proceed as

$$\begin{aligned} &\left| \sum_{1 \leq i, j \leq n} \langle \mathbf{Z}_{ij}, (\mathbf{X}_i - \mathbf{I})^\top (\mathbf{X}_j - \mathbf{I}) + (\mathbf{X}_i - \mathbf{I})^\top + (\mathbf{X}_j - \mathbf{I}) \rangle \right| \\ &\leq (\mathbf{X} - \mathbf{I}^n) \mathbf{Z} (\mathbf{X} - \mathbf{I}^n)^\top + 2 \left| \sum_{1 \leq i \leq n} \left\langle \sum_{1 \leq j \leq n} \mathbf{Z}_{ij}, \mathbf{X}_i - \mathbf{I} \right\rangle \right| \\ &\leq \|\mathbf{Z}\|_{\text{op}} \|\mathbf{X} - \mathbf{I}^n\|_F^2 + 2 \sum_{1 \leq i \leq n} \left\| \sum_{1 \leq j \leq n} \mathbf{Z}_{ij} \right\|_F \|\mathbf{X}_i - \mathbf{I}\|_F \\ &\leq \left( 2\sqrt{d} \|\mathbf{Z}\|_{\text{op}} + 2 \max_{1 \leq i \leq n} \left\| \sum_{1 \leq j \leq n} \mathbf{Z}_{ij} \right\|_F \right) \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{I}\|_F, \end{aligned} \tag{39}$$

where  $\mathbf{I}^n \in \mathbb{R}^{nd \times nd}$  collects  $n$  identity matrix together and the last inequality holds because  $\|\mathbf{X} - \mathbf{I}^n\|_F^2 = \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{I}\|_F^2 \leq \sum_{1 \leq i \leq n} (\|\mathbf{X}_i\|_F + \|\mathbf{I}\|_F) \|\mathbf{X}_i - \mathbf{I}\|_F = 2\sqrt{d} \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{I}\|_F$ .

Using the randomness of  $\mathbf{O}_{ij}$ , we claim that with probability at least  $1 - 4d/n$ , we have

$$\|\mathbf{Z}\|_{\text{op}} \leq \sqrt{8n(1-p)q \log n}, \quad \text{and} \quad \max_{1 \leq i \leq n} \left\| \sum_{1 \leq j \leq n} \mathbf{Z}_{ij} \right\|_F \leq \sqrt{8n(1-p)q \log n}. \quad (40)$$

Combining the above bounds gives

$$\begin{aligned} & h(\mathbf{X}) - h(\mathbf{X}^*) \\ & \geq - \left( 2(\sqrt{d} + 1)\sqrt{8n(1-p)q \log n} + \frac{(1-p)qn}{\sqrt{2}} \max_i \|\mathbf{X}_i - \mathbf{I}\|_F \right) \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{I}\|_F \\ & \geq - \frac{npq}{16} \cdot \text{dist}_1(\mathbf{X}, \mathbf{X}^*), \end{aligned} \quad (41)$$

where the last inequality holds because we assume  $p^2q^2 = \Omega\left(\frac{\log n}{n}\right)$ ,  $\max_i \|\mathbf{X}_i - \mathbf{I}\|_F = \text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ , and  $\sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{I}\|_F = \text{dist}_1(\mathbf{X}, \mathbf{X}^*)$ .

Finally, it remains to show that the two inequalities in (40) holds with probability at least  $1 - 4d/n$ . It is quick to verify that  $\mathbb{E}(\mathbf{Z}_{ij}) = \mathbf{0}$ ,  $\|\mathbf{Z}_{ij}\|_{\text{op}} \leq 1 + (1-p)q\|\mathbf{A}\|_F \leq 2$ , and

$$\begin{aligned} \mathbb{E}(\mathbf{Z}^2) &= \text{BlkDiag} \left( \sum_j \mathbb{E}(\mathbf{Z}_{1j} \mathbf{Z}_{1j}^\top), \dots, \sum_j \mathbb{E}(\mathbf{Z}_{nj} \mathbf{Z}_{nj}^\top) \right) \\ &= \sum_j \mathbb{E}(\mathbf{Z}_{1j} \mathbf{Z}_{1j}^\top) \otimes \mathbf{I}_n \\ &= n(1-p)q \mathbb{E} \left( \frac{(\mathbf{I} - \mathbf{O}_{ij})(\mathbf{I} - \mathbf{O}_{ij})^\top}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F^2} - (1-p)^2q^2 \mathbf{A} \mathbf{A}^\top \right) \otimes \mathbf{I}_n, \end{aligned}$$

where  $\text{BlkDiag}(\cdot)$  means the block diagonal matrix,  $\otimes$  means the Kronecker product and  $\mathbf{I}_n$  denotes the  $n$ -by- $n$  identity matrix. Thus, we have  $\|\mathbb{E}(\mathbf{Z}^2)\|_{\text{op}} \leq n(1-p)q$ . According to the Matrix Bernstein inequality [38], we have

$$\begin{aligned} \Pr \left( \|\mathbf{Z}\|_{\text{op}} \leq \sqrt{8n(1-p)q \log n} \right) &\geq 1 - 2nd \exp \left( \frac{-4n(1-p)q \log n}{n(1-p)q + 2\sqrt{8n(1-p)q \log n}/3} \right) \\ &\geq 1 - \frac{2d}{n}. \end{aligned} \quad (42)$$

Here, the last inequality holds because we assume  $2\sqrt{8n(1-p)q \log n}/3 \leq n(1-p)q$ , which is true as long as  $p = \Omega(\log n/n)$ .

For any fixed  $1 \leq i \leq n$ , a similar argument based on Matrix Bernstein inequality [38] shows that

$$\begin{aligned} \Pr \left( \left\| \sum_j \mathbf{Z}_{ij} \right\|_F \leq \sqrt{8n(1-p)q \log n} \right) &\geq 1 - 2d \exp \left( \frac{-4n(1-p)q \log n}{n(1-p)q + \sqrt{8n(1-p)q \log n}/3} \right) \\ &\geq 1 - \frac{2d}{n^2}, \end{aligned} \quad (43)$$

which implies  $\Pr \left( \max_i \left\| \sum_j \mathbf{Z}_{ij} \right\|_F \leq \sqrt{8n(1-p)q \log n} \right) \geq 1 - 2d/n$ .

## C Full Proof of Theorem 4

### C.1 Proof of Corollary 1

Combining the convexity of  $f$  and Theorem 3, we have

$$-\frac{npq}{8} \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{X}_i^*\|_F \geq f(\mathbf{X}^*) - f(\mathbf{X}) \geq \left\langle \tilde{\nabla} f(\mathbf{X}), \mathbf{X}^* - \mathbf{X} \right\rangle, \quad \forall \tilde{\nabla} f(\mathbf{X}) \in \partial f(\mathbf{X}). \quad (44)$$

for all  $\mathbf{X} \in \text{SO}(d)^n$  satisfying  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) = \mathcal{O}(p)$ . For any  $\tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}) \in \partial_{\mathcal{R}}^\perp f(\mathbf{X})$ , we can further compute

$$\begin{aligned} \left\langle \tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}), \mathbf{X} - \mathbf{X}^* \right\rangle &= \sum_{1 \leq i \leq n} \left\langle \tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}_i), \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}^\perp}(\mathbf{X}_i - \mathbf{X}_i^*) \right\rangle \\ &\leq \sum_{1 \leq i \leq n} \|\tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}_i)\|_F \cdot \|\mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}^\perp}(\mathbf{X}_i - \mathbf{X}_i^*)\|_F. \end{aligned}$$

On the one hand, notice that

$$\begin{aligned} \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}^\perp}(\mathbf{X}_i - \mathbf{X}_i^*) &= \mathbf{X}_i - \mathbf{X}_i^* - \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}}(\mathbf{X}_i - \mathbf{X}_i^*) \\ &= \mathbf{X}_i (\mathbf{X}_i^\top (\mathbf{X}_i - \mathbf{X}_i^*) + (\mathbf{X}_i - \mathbf{X}_i^*)^\top \mathbf{X}_i) / 2 \\ &= \mathbf{X}_i (\mathbf{X}_i - \mathbf{X}_i^*)^\top (\mathbf{X}_i - \mathbf{X}_i^*) / 2, \end{aligned}$$

which implies

$$\|\mathcal{P}_{\mathbf{T}_{\mathbf{X}}^\perp}(\mathbf{X}_i - \mathbf{X}_i^*)\|_F = \frac{1}{2} \|\mathbf{X}_i (\mathbf{X}_i - \mathbf{X}_i^*)^\top (\mathbf{X}_i - \mathbf{X}_i^*)\|_F \leq \frac{1}{2} \|\mathbf{X}_i - \mathbf{X}_i^*\|_F^2.$$

On the other hand, according to Lemma 1, with probability at least  $1 - \mathcal{O}(1/n)$ , for any  $i \in [n]$ ,

$$\|\tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}_i)\|_F \leq \|\tilde{\nabla} f(\mathbf{X}_i)\|_F = \left\| \sum_{j \in \mathcal{E}_i} \tilde{\nabla} f_{i,j}(\mathbf{X}_i) \right\|_F \leq \sum_{j \in \mathcal{E}_i} \|\tilde{\nabla} f_{i,j}(\mathbf{X}_i)\|_F \leq |\mathcal{E}_i| \leq 2nq,$$

where  $\tilde{\nabla} f_{i,j}(\mathbf{X}_i) \in \partial f_{i,j}(\mathbf{X}_i)$  satisfies  $\|\tilde{\nabla} f_{i,j}(\mathbf{X}_i)\|_F \leq 1$ . Hence, we have

$$\left\langle \tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}), \mathbf{X} - \mathbf{X}^* \right\rangle \leq nq \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{X}_i^*\|_F^2 \leq \frac{npq}{16} \sum_i \|\mathbf{X}_i - \mathbf{X}_i^*\|_F$$

for any  $\mathbf{X}$  such that  $\text{dist}_\infty(\mathbf{X}, \mathbf{X}^*) \leq \frac{p}{16}$ . Invoking the above bounds into (44) yields the desired result

$$\left\langle \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}), \mathbf{X} - \mathbf{X}^* \right\rangle = \left\langle \tilde{\nabla} f(\mathbf{X}) - \tilde{\nabla}_{\mathcal{R}}^\perp f(\mathbf{X}), \mathbf{X} - \mathbf{X}^* \right\rangle \geq \frac{npq}{16} \sum_{1 \leq i \leq n} \|\mathbf{X}_i - \mathbf{X}_i^*\|_F.$$

## C.2 Proof of Contraction

Let us first present some preliminary results, which will be used in our later derivations. By noticing that  $\tilde{\nabla} f(\mathbf{X}_i) = \sum_{j \in \mathcal{E}_i} \tilde{\nabla} f_{i,j}(\mathbf{X}_i)$  where  $\tilde{\nabla} f_{i,j}(\mathbf{X}_i) \in \partial f_{i,j}(\mathbf{X}_i)$ , we define

$$\tilde{\nabla} g(\mathbf{X}_i) = \sum_{j \in \mathcal{A}_i} \tilde{\nabla} f_{i,j}(\mathbf{X}_i), \quad \tilde{\nabla} h(\mathbf{X}_i) = \sum_{j \in \mathcal{E}_i \setminus \mathcal{A}_i} \tilde{\nabla} f_{i,j}(\mathbf{X}_i).$$

Recall that  $\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_i) = \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}}(\tilde{\nabla} g(\mathbf{X}_i))$  and  $\tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i) = \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}}(\tilde{\nabla} h(\mathbf{X}_i))$ . Similarly, we have  $\tilde{\nabla} f(\mathbf{X}_i) = \tilde{\nabla} g(\mathbf{X}_i) + \tilde{\nabla} h(\mathbf{X}_i)$  and  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i) = \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_i) + \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i)$ . Furthermore, the QR decomposition-based retraction satisfies the second-order boundedness property, i.e., there exists some  $M \geq 1$  such that

$$\begin{aligned} \|\mathbf{X}_i^{k+1} - \mathbf{X}_i^*\|_F &= \|\text{Retr}_{\mathbf{X}_i^k}(-\mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k)) - \mathbf{X}_i^*\|_F \\ &\leq \|\mathbf{X}_i^k - \mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) - \mathbf{X}_i^*\|_F + M \cdot \mu_k^2 \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k)\|_F^2. \end{aligned} \quad (45)$$

Recall that  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) = \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_i^k) + \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i^k)$ , we have

$$\begin{aligned} \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k)\|_F &\leq \|\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_i^k)\|_F + \|\tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i^k)\|_F \\ &\leq 5\sqrt{2\delta_0}nq + (1 + \epsilon)npq \leq (1 + 2\epsilon)npq, \end{aligned} \quad (46)$$

where the second inequality comes from

$$\|\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_i^k)\|_F = \left\| \sum_{j \in \mathcal{A}_i} \tilde{\nabla}_{\mathcal{R}} f_{i,j}(\mathbf{X}_i) \right\|_F \leq \sum_{j \in \mathcal{A}_i} \|\tilde{\nabla}_{\mathcal{R}} f_{i,j}(\mathbf{X}_i)\|_F \leq |\mathcal{A}_i| \leq (1 + \epsilon)npq,$$

and Lemma 10 and the last inequality is due to the choice  $\delta_0 \leq \epsilon^2 p^2 / 50$  (i.e.,  $\delta_0 = \mathcal{O}(p^2)$ ). Thus, by choosing  $\mu_k = \mathcal{O}(\frac{\delta_k}{n})$ , and  $\delta_0 = \mathcal{O}(p^2)$ , the second-order term  $M \cdot \mu_k^2 \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k)\|_F^2 = \mathcal{O}(p^6 q^2)$ , which is a very high-order error. In the following analysis, we will ignore this term to simplify our derivations.

Using the above preliminaries and Corollary 1, we are ready to establish two key lemmas, which show that if  $\mathbf{X}^k \in \mathcal{N}_F^k \cap \mathcal{N}_\infty^k$ , then  $\mathbf{X}^{k+1} \in \mathcal{N}_F^{k+1}$  (Lemma 8) and  $\mathbf{X}^{k+1} \in \mathcal{N}_\infty^{k+1}$  (Lemma 9), respectively. This completes the proof of Theorem 4.

**Lemma 8.** *With high probability, suppose that  $\mathbf{X}^k \in \mathcal{N}_F^k \cap \mathcal{N}_\infty^k$ ,  $\mu_k = \mathcal{O}(\frac{\delta_k}{n})$ , and*

$$\delta_0 = \mathcal{O}(p^2), \quad \text{and} \quad \xi_0 = \Theta(\sqrt{npq}\delta_0), \quad (47)$$

*then  $\mathbf{X}^{k+1} \in \mathcal{N}_F^{k+1}$ .*

*Proof.* By ignoring the high-order error term in (45), in order to bound  $\|\mathbf{X}^{k+1} - \mathbf{X}^*\|_F^2$  we can first compute

$$\begin{aligned} \|\mathbf{X}^{k+1} - \mathbf{X}^*\|_F^2 &= \sum_{1 \leq i \leq n} \|\mathbf{X}_i^{k+1} - \mathbf{X}_i^*\|_F^2 \leq \sum_{1 \leq i \leq n} \|\mathbf{X}_i^k - \mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) - \mathbf{X}_i^*\|_F^2 \\ &= \|\mathbf{X}^k - \mathbf{X}^*\|_F^2 - 2\mu_k \left\langle \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}^k), \mathbf{X}^k - \mathbf{X}^* \right\rangle + \mu_k^2 \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}^k)\|_F^2. \end{aligned}$$

Then, according to Corollary 1, we have

$$\begin{aligned} \left\langle \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}^k), \mathbf{X}^k - \mathbf{X}^* \right\rangle &\geq \frac{npq}{16} \sum_{1 \leq i \leq n} \|\mathbf{X}_i^k - \mathbf{X}_i^*\|_F \geq \frac{npq}{16\delta_k} \sum_{1 \leq i \leq n} \|\mathbf{X}_i^k - \mathbf{X}_i^*\|_F^2 \\ &= \frac{npq}{16\delta_k} \|\mathbf{X}^k - \mathbf{X}^*\|_F^2, \end{aligned}$$

where the second inequality holds because  $\|\mathbf{X}_i^k - \mathbf{X}_i^*\|_F \leq \delta_k$  (i.e.,  $\mathbf{X}^k \in \mathcal{N}_\infty^k$ ). Combining the above two inequalities gives

$$\begin{aligned} \|\mathbf{X}^{k+1} - \mathbf{X}^*\|_F^2 &\leq \left(1 - \mu_k \cdot \frac{npq}{8\delta_k}\right) \|\mathbf{X}^k - \mathbf{X}^*\|_F^2 + \mu_k^2 \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}^k)\|_F^2 \\ &\leq \left(1 - \frac{pq}{8}\right) \|\mathbf{X}^k - \mathbf{X}^*\|_F^2 + (1 + 2\epsilon)^2 n^3 p^2 q^2 \mu_k^2, \end{aligned}$$

where the last inequality is due to (46). Since  $\mu_k = \mathcal{O}(\frac{\delta_k}{n})$  and  $\xi_0 = \Theta(\sqrt{npq}\delta_0)$  (i.e.,  $\xi_k = \Theta(\sqrt{npq}\delta_k)$ ), we have  $n^3 p^2 q^2 \mu_k^2 = \mathcal{O}(pq\xi_k^2)$ , which implies

$$\|\mathbf{X}^{k+1} - \mathbf{X}^*\|_F^2 \leq \left(1 - \frac{pq}{8}\right) \xi_k^2 + \frac{pq}{16} \xi_k^2 = \left(1 - \frac{pq}{16}\right) \xi_k^2 = \xi_{k+1}^2.$$

This completes the proof.  $\square$

**Lemma 9.** *With high probability, suppose that  $\mathbf{X}^k \in \mathcal{N}_F^k \cap \mathcal{N}_\infty^k$ ,  $\mu_k = \mathcal{O}(\frac{\delta_k}{n})$ , and*

$$\delta_0 = \mathcal{O}(p^2) \quad \text{and} \quad \xi_0 = \mathcal{O}(\sqrt{npq}\delta_0), \quad (48)$$

*then  $\mathbf{X}^{k+1} \in \mathcal{N}_\infty^{k+1}$ .*

*Proof.* As stated at the beginning of Appendix B, we can assume without loss of generality that  $\mathbf{R}^* = \mathbf{I}$  and  $\mathbf{X}_i^* = \mathbf{I}$  for all  $1 \leq i \leq n$ . We divide the index set  $[n]$  into three sets

$$\begin{aligned} \mathcal{I}_1 &= \left\{ i \mid \|\mathbf{X}_i^k - \mathbf{I}\|_F \leq \frac{\delta_k}{4} \right\}, \quad \mathcal{I}_2 = \left\{ i \mid \frac{\delta_k}{4} < \|\mathbf{X}_i^k - \mathbf{I}\|_F \leq \frac{3\delta_k}{4} \right\}, \\ \text{and } \mathcal{I}_3 &= \left\{ i \mid \frac{3\delta_k}{4} < \|\mathbf{X}_i^k - \mathbf{I}\|_F \leq \delta_k \right\}. \end{aligned}$$

For any  $i \in \mathcal{I}_1 \cup \mathcal{I}_2$ , we have

$$\|\mathbf{X}_i^k - \mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) - \mathbf{I}\|_F \leq \|\mathbf{X}_i^k - \mathbf{I}\|_F + \mu_k \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k)\|_F \leq \frac{3\delta_k}{4} + 2\mu_k npq \leq \delta_{k+1}, \quad (49)$$

where the last inequality holds because we choose  $\mu_k = \mathcal{O}(\frac{\delta_k}{n}) \leq \frac{\delta_k}{16n}$ .

It remains to consider the case  $i \in \mathcal{I}_3$ . Firstly, it is easy to see that

$$\text{dist}(\mathbf{X}^k, \mathbf{X}^*)^2 \geq \sum_{i \in \mathcal{I}_2 \cup \mathcal{I}_3} \|\mathbf{X}^k - \mathbf{X}^*(\delta_k)\|_F^2 \geq \frac{\delta_k^2}{16} |\mathcal{I}_2 \cup \mathcal{I}_3|.$$

Note that we have  $|\mathcal{I}_2 \cup \mathcal{I}_3| \leq \frac{\text{dist}(\mathbf{X}^k, \mathbf{X}^*)^2}{\delta_k^2/16} \leq \frac{16\xi_k^2}{\delta_k^2} = \mathcal{O}(npq)$  according to the assumption  $\xi_0 = \mathcal{O}(\sqrt{npq}\delta_0)$ . Hence, for any  $i \in \mathcal{I}_3$ , we have

$$\tilde{\nabla} f(\mathbf{X}_i) = \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\mathbf{X}_i - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} + \sum_{j \in \mathcal{A}_i \setminus \mathcal{I}_1} \frac{\mathbf{X}_i - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} + \tilde{\nabla} h(\mathbf{X}_i). \quad (50)$$

Since  $|\mathcal{I}_2 \cup \mathcal{I}_3| = \mathcal{O}(npq)$ , we can choose  $\xi_0$  properly such that  $|\mathcal{I}_2 \cup \mathcal{I}_3| \leq \epsilon npq$ . Then, we have

$$|\mathcal{A}_i \setminus \mathcal{I}_1| \leq |\mathcal{I}_2 \cup \mathcal{I}_3| \leq \epsilon npq \quad \text{and} \quad |\mathcal{A}_i \cap \mathcal{I}_1| \geq |\mathcal{A}_i| - |\mathcal{A}_i \setminus \mathcal{I}_1| \geq (1 - 2\epsilon)npq.$$

Let us use  $\tilde{\nabla} g^1(\mathbf{X}_i)$  to denote the first term on the RHS of (50). We have

$$\begin{aligned} \tilde{\nabla} g^1(\mathbf{X}_i) &= \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\mathbf{X}_i - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} = \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\mathbf{X}_i - \mathbf{I}}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} + \frac{\mathbf{I} - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} \\ &= \sigma(\mathbf{X}_i - \mathbf{I}) + \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\mathbf{I} - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F}, \end{aligned} \quad (51)$$

where  $\sigma = \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{1}{\|\mathbf{X}_i - \mathbf{X}_j\|_F}$ . For the last term in the above equation, we have

$$\left\| \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\mathbf{I} - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} \right\|_F \leq \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{\|\mathbf{I} - \mathbf{X}_j\|_F}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} \leq \frac{\delta_k}{4} \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{1}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} = \frac{\delta_k \sigma}{4}.$$

In addition, the projection of  $\mathbf{X}_i - \mathbf{I}$  onto the cotangent space can be bounded as

$$\|\mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}^\perp}(\mathbf{X}_i - \mathbf{I})\|_F = \|(\mathbf{X}_i - \mathbf{I})\|_F^2/2 \leq \delta_k^2/2,$$

which implies that  $\tilde{\nabla}_{\mathcal{R}} g^1(\mathbf{X}_i) = \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}}(\tilde{\nabla} g^1(\mathbf{X}_i))$  satisfies

$$\|\tilde{\nabla}_{\mathcal{R}} g^1(\mathbf{X}_i) - \sigma(\mathbf{X}_i - \mathbf{I})\|_F \leq \frac{\delta_k \sigma}{4} + \frac{\delta_k^2 \sigma}{2} \leq \frac{\delta_k \sigma}{2}.$$

Then, the fact  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i) = \tilde{\nabla}_{\mathcal{R}} g^1(\mathbf{X}_i) + \mathcal{P}_{\mathbf{T}_{\mathbf{X}_i}^\perp} \left( \sum_{j \in \mathcal{A}_i \setminus \mathcal{I}_1} \frac{\mathbf{X}_i - \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} \right) + \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i)$  implies

$$\begin{aligned} \|\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i) - \sigma(\mathbf{X}_i - \mathbf{I})\|_F &\leq \|\tilde{\nabla}_{\mathcal{R}} g^1(\mathbf{X}_i) - \sigma(\mathbf{X}_i - \mathbf{I})\|_F + |\mathcal{A}_i \setminus \mathcal{I}_1| + \|\tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i)\|_F \\ &\leq \frac{\delta_k \sigma}{2} + \epsilon npq + 5\sqrt{2\delta_0}nq. \end{aligned}$$

Next, motivated by the update of ReSync, we can construct

$$\mathbf{X}_i^k - \mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) - \mathbf{I} = (1 - \mu_k \sigma)(\mathbf{X}_i^k - \mathbf{I}) + \mu_k (\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i) - \sigma(\mathbf{X}_i - \mathbf{I})),$$

which implies

$$\begin{aligned} \|\mathbf{X}_i^k - \mu_k \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_i^k) - \mathbf{I}\|_F &\leq (1 - \mu_k \sigma) \|\mathbf{X}_i^k - \mathbf{I}\| + \mu_k \left( \frac{\delta_k \sigma}{4} + \frac{\delta_k^2 \sigma}{2} + \epsilon npq + 5\sqrt{2\delta_0}nq \right) \\ &\leq (1 - \mu_k \sigma) \delta_k + \mu_k \left( \frac{\delta_k \sigma}{2} + \epsilon npq + 5\sqrt{2\delta_0}nq \right) \\ &= \delta_k - \mu_k \left( \frac{\delta_k \sigma}{2} - \epsilon npq - 5\sqrt{2\delta_0}nq \right). \end{aligned} \quad (52)$$

In order to further upper bound the above inequality, we can compute

$$\begin{aligned} \sigma &= \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{1}{\|\mathbf{X}_i - \mathbf{X}_j\|_F} \geq \sum_{j \in \mathcal{A}_i \cap \mathcal{I}_1} \frac{1}{\|\mathbf{X}_i - \mathbf{I}\|_F + \|\mathbf{X}_j - \mathbf{I}\|_F} \\ &\geq \frac{4}{5\delta_k} |\mathcal{A}_i \cap \mathcal{I}_1| = \frac{4(1 - 2\epsilon)npq}{5\delta_k}. \end{aligned}$$

where the second inequality holds because  $\|\mathbf{X}_i - \mathbf{I}\|_F + \|\mathbf{X}_j - \mathbf{I}\|_F \leq 5\delta_k/4$  as  $j \in \mathcal{I}_1$ . It implies

$$\frac{\delta_k \sigma}{2} - \epsilon npq - 5\sqrt{2\delta_0}nq \geq \frac{2(1-2\epsilon)npq}{5} - \epsilon npq - 5\sqrt{2\delta_0}nq \geq \frac{npq}{4}.$$

By plugging the above bound into (52) and ignoring the high-order error term in (45), we complete the proof.  $\square$

Finally, we present Lemma 10 and its proof.

**Lemma 10.** *With high probability, the following holds for all  $\mathbf{X}^0 \in \mathcal{N}_\infty^0$ :*

$$\left\| \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i) \right\|_F \leq 5\sqrt{2\delta_0}nq \quad \forall i \in [n]. \quad (53)$$

*Proof of Lemma 10.* Firstly, define  $\tilde{\mathbf{O}}_{ij} = \mathbf{O}_{ij} \mathbf{X}_j \mathbf{X}_i^\top$ , then

$$\tilde{\nabla} h(\mathbf{X}_i) = \sum_{j \in \mathcal{E}_i/\mathcal{A}_i} \frac{\mathbf{X}_i - \mathbf{O}_{ij} \mathbf{X}_j}{\|\mathbf{X}_i - \mathbf{O}_{ij} \mathbf{X}_j\|_F} = \sum_{j \in \mathcal{E}_i/\mathcal{A}_i} \frac{\mathbf{I} - \mathbf{O}_{ij} \mathbf{X}_j \mathbf{X}_i^\top}{\|\mathbf{X}_i - \mathbf{O}_{ij} \mathbf{X}_j\|_F} \mathbf{X}_i = \sum_{j \in \mathcal{E}_i/\mathcal{A}_i} \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} \mathbf{X}_i$$

The fact that  $\mathbf{X} \in \mathcal{N}_\infty^0$  implies that  $\|\mathbf{X}_i - \mathbf{X}_j\|_F \leq 2\delta_0$ , i.e.,  $\|\mathbf{I} - \mathbf{X}_j \mathbf{X}_i^\top\|_F \leq 2\delta_0$ . For any  $\|\mathbf{O}_{ij} - \mathbf{I}\| \geq \sqrt{2\delta_0}$ , we have

$$\begin{aligned} \left\| \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F &\leq \left\| \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F \\ &\quad + \left\| \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} - \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F \leq 2\sqrt{2\delta_0} \end{aligned}$$

where the last inequality holds because

$$\begin{aligned} \left\| \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F &= \|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F \left| \frac{1}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - \frac{1}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right| \\ &= \left| 1 - \frac{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right| \\ &= \frac{|\|\mathbf{I} - \mathbf{O}_{ij}\|_F - \|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F|}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \leq \frac{\|\mathbf{O}_{ij} - \tilde{\mathbf{O}}_{ij}\|_F}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \\ &= \frac{\|\mathbf{I} - \mathbf{X}_j \mathbf{X}_i^\top\|_F}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \leq \sqrt{2\delta_0}, \end{aligned}$$

and

$$\left\| \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} - \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F = \frac{\|\mathbf{O}_{ij} - \tilde{\mathbf{O}}_{ij}\|_F}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} = \frac{\|\mathbf{I} - \mathbf{X}_j \mathbf{X}_i^\top\|_F}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \leq \sqrt{2\delta_0}.$$

Let  $\Phi_i = \{j \in \mathcal{E}_i/\mathcal{A}_i \mid \|\mathbf{O}_{ij} - \mathbf{I}\| \leq \sqrt{2\delta_0}\}$ . According to the fact that  $|\mathcal{E}_i/\mathcal{A}_i| \leq (1+\epsilon)nq$  and the randomness of  $\mathbf{O}_{ij}$ , it is easy to show that  $|\Phi_i| \leq \sqrt{2\delta_0}nq$  hold for all  $1 \leq i \leq n$  with high probability. Thus, by splitting the sum  $\sum_{j \in \mathcal{E}_i/\mathcal{A}_i}$  into two parts:  $j \in \Phi_i$  and  $j \notin \Phi_i$ , we have

$$\left\| \sum_{j \in \mathcal{E}_i/\mathcal{A}_i} \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - \sum_{j \in \mathcal{E}_i/\Omega_i} \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\|_F \leq 4\sqrt{2\delta_0}nq. \quad (54)$$

Besides, since  $\mathbf{O}_{ij}$  is uniformly distributed on  $\text{SO}(d)$ , according to Lemma A.1 in [39], we know that  $\mathbb{E} \left\{ \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} \right\} = c(d)\mathbf{I}$ . Then, the matrix Bernstein's inequality [38] tells us that, with high probability,

$$\left\| \sum_{j \in \mathcal{E}_i/\Omega_i} \frac{\mathbf{I} - \mathbf{O}_{ij}}{\|\mathbf{I} - \mathbf{O}_{ij}\|_F} - |\mathcal{E}_i/\Omega_i| \cdot c(d)\mathbf{I} \right\|_F \leq \sqrt{2\delta_0}nq. \quad (55)$$

This, together with (54), implies that

$$\left\| \sum_{j \in \mathcal{E}_i / \mathcal{A}_i} \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} - |\mathcal{E}_i / \Omega_i| \cdot c(d) \mathbf{I} \right\|_F \leq 5\sqrt{2\delta_0} nq.$$

The fact that  $\tilde{\nabla} h(\mathbf{X}_i) = \sum_{j \in \mathcal{E}_i / \mathcal{A}_i} \frac{\mathbf{I} - \tilde{\mathbf{O}}_{ij}}{\|\mathbf{I} - \tilde{\mathbf{O}}_{ij}\|_F} \mathbf{X}_i$  implies that

$$\left\| \tilde{\nabla} h(\mathbf{X}_i) - |\mathcal{E}_i / \Omega_i| \cdot c(d) \mathbf{X}_i \right\|_F \leq 5\sqrt{2\delta_0} nq. \quad (56)$$

We complete the proof by taking the projection operator  $\tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}_i) = \mathcal{P}_{\mathbb{T}_{\mathbf{X}_i}}(\tilde{\nabla} h(\mathbf{X}_i))$  and the fact that  $\mathcal{P}_{\mathbb{T}_{\mathbf{X}_i}}(\mathbf{X}_i) = 0$ .  $\square$