

Structured Local Optima in Sparse Blind Deconvolution

Yuqian Zhang, Han-Wen Kuo, John Wright
Department of Electrical Engineering and Data Science Institute
Columbia University

June 1, 2018

Abstract

Blind deconvolution is a ubiquitous problem of recovering two unknown signals from their convolution. Unfortunately, this is an ill-posed problem in general. This paper focuses on the *short and sparse* blind deconvolution problem, where the one unknown signal is short and the other one is sparsely and randomly supported. This variant captures the structure of the unknown signals in several important applications. We assume the short signal to have unit ℓ^2 norm and cast the blind deconvolution problem as a nonconvex optimization problem over the sphere. We demonstrate that (i) in a certain region of the sphere, every local optimum is close to some shift truncation of the ground truth, and (ii) for a generic short signal of length k , when the sparsity of activation signal $\theta \lesssim k^{-2/3}$ and number of measurements $m \gtrsim \text{poly}(k)$, a simple initialization method together with a descent algorithm which escapes strict saddle points recovers a near shift truncation of the ground truth kernel.

1 Introduction

Blind deconvolution is the problem of recovering two unknown signals \mathbf{a}_0 and \mathbf{x}_0 from their convolution $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$. This fundamental problem recurs across several fields, including astronomy, microscopy data processing [CLC⁺17], neural spike sorting [Lew98], computer vision [KH96], etc. However, this problem is ill-posed without further priors on the unknown signals, as there are infinitely many pairs of signals (\mathbf{a}, \mathbf{x}) whose convolution equals a given observation \mathbf{y} . Fortunately, in practice, the target signals (\mathbf{a}, \mathbf{x}) are often structured. In particular, a number of practical applications exhibit a common *short-and-sparse* structure:

In *Neural spike sorting*: Neurons in the brain fire brief voltage spikes when stimulated. The signatures of the spikes encode critical features of the neuron and the occurrence of such spikes are usually sparse and random in time [Lew98, ETS11].

In *Microscopy data analysis*: The nanoscale materials of interests are contaminated by randomly and sparsely distributed “defects”, which can dramatically change the electronic structure of the material [CLC⁺17].

In *Image deblurring*: Blurred images due to camera shake can be modeled as a convolution of the latent sharp image and a kernel capturing the motion of the camera. Although natural images are not sparse, they typically have (approximately) sparse gradients [CW98, LWDF11].

In the above applications, the observation signal $\mathbf{y} \in \mathbb{R}^m$ is generated via the convolution of a *short* kernel $\mathbf{a}_0 \in \mathbb{R}^k$ with $k \ll m$ and a *sparse* activation coefficient $\mathbf{x}_0 \in \mathbb{R}^m$ with $\|\mathbf{x}_0\|_0 \ll m$. Without loss of generality, we let \mathbf{y} denote the circular convolution of \mathbf{a}_0 and \mathbf{x}_0

$$\mathbf{y} = \mathbf{a}_0 \circledast \mathbf{x}_0 = \widetilde{\mathbf{a}}_0 \circledast \mathbf{x}_0, \quad (1.1)$$

with $\widetilde{\mathbf{a}}_0 \in \mathbb{R}^m$ denoting the zero padded m -length version of \mathbf{a}_0 , which can be expressed as $\widetilde{\mathbf{a}}_0 = \boldsymbol{\iota}_k \mathbf{a}_0$. Here, $\boldsymbol{\iota}_k : \mathbb{R}^k \rightarrow \mathbb{R}^m$ is a zero padding operator. Its adjoint $\boldsymbol{\iota}_k^* : \mathbb{R}^m \rightarrow \mathbb{R}^k$ acts as a projection onto the lower dimensional space by keeping the first k components.

The short-and-sparse blind deconvolution problem exhibits a *scaled-shift ambiguity*, which derives from the basic properties of a convolution operator. Namely, for any observation signal \mathbf{y} , and any nonzero scalar α and integer shift τ , the following equality always holds

$$\mathbf{y} = (\pm \alpha s_\tau[\widetilde{\mathbf{a}_0}]) \circledast (\pm \alpha^{-1} s_{-\tau}[\mathbf{x}_0]). \quad (1.2)$$

Here, $s_{-\tau}[\mathbf{v}]$ denotes the cyclic shift of the vector \mathbf{v} by τ entries:

$$s_\tau[\mathbf{v}](i) = \mathbf{v}([i - \tau - 1]_m + 1), \quad \forall i \in \{1, \dots, m\}. \quad (1.3)$$

Clearly, both scaling and cyclic shifts preserve the short-and-sparse structure of $(\mathbf{a}_0, \mathbf{x}_0)$. This *scaled-shift symmetry* raises nontrivial challenges for computation, making straightforward convexification approaches ineffective, and leading to complicated nonconvex optimization landscape. [ZLK⁺17] considers a natural nonconvex formulation of sparse blind deconvolution, in which the kernel $\mathbf{a} \in \mathbb{R}^k$ is constrained to have unit Frobenius norm. [ZLK⁺17] argues that under certain idealized conditions, this problem has well-structured local optima, in the sense that *every local optimum is close to some shift truncation of the ground truth*. The presence of these local optima can be viewed as a result of the shift symmetry associated to the convolution operator: the shifted and truncated kernel $\iota_k^* s_\tau[\widetilde{\mathbf{a}_0}]$ can be convolved with the sparse signal $s_{-\tau}[\mathbf{x}_0]$ (shifted in the other direction) to produce a near approximation to \mathbf{y} that

$$(\iota_k^* s_\tau[\widetilde{\mathbf{a}_0}]) \circledast s_{-\tau}[\mathbf{x}_0] \approx \mathbf{y}. \quad (1.4)$$

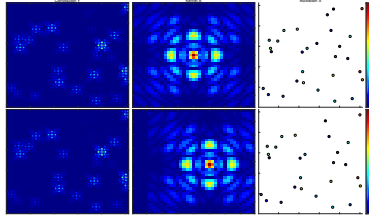


Figure 1: Local Minimum. Top: observation $\mathbf{y} = \mathbf{a}_0 \circledast \mathbf{x}_0$, and ground truth \mathbf{a}_0 , and \mathbf{x}_0 ; Bottom: recovered $\mathbf{a} \circledast \mathbf{x}$, \mathbf{a} , and \mathbf{x} at one local minimum of a natural formulation in [ZLK⁺17].

In [ZLK⁺17], the geometric insight about local minima is corroborated with a lot of experiments, but rigorous proof is only available under rather restrictive conditions. In this paper, we adopt the unit Frobenius norm constraint as in [ZLK⁺17] but consider a different objective function over the kernel sphere \mathbb{S}^{k-1} . We formulate the sparse blind deconvolution problem as the following optimization problem over the sphere:

$$\min -\|\check{\mathbf{y}} \circledast \mathbf{r}_{\mathbf{y}}(\mathbf{q})\|_4^4 \quad \text{s. t.} \quad \|\mathbf{q}\|_F = 1 \quad (1.5)$$

Here, $\check{\mathbf{y}}$ denotes the reversal¹ of \mathbf{y} and $\mathbf{r}_{\mathbf{y}}(\mathbf{q})$ is a preconditioner which we will discuss in detail later. Convolution $\check{\mathbf{y}} \circledast \mathbf{r}_{\mathbf{y}}(\mathbf{q})$ approximates the reversed underlying activation signal \mathbf{x}_0 , and $-\|\cdot\|_4^4$ serves as the sparsity penalty.

This paper studies the function landscape of the short-and-sparse blind deconvolution problem assuming the short k -length convolutional kernel lives on a unit Frobenius norm sphere, denoted as \mathbb{S}^{k-1} . We demonstrate that even when \mathbf{x}_0 is relatively dense, a shift truncation $\iota_k^* s_\tau[\widetilde{\mathbf{a}_0}]$ of the ground truth still can be obtained as one local minimum in certain region of the kernel sphere. Such benign region contains the sub-level set of small objective value, and an initial point with small objective value can be easily found. Specifically, for a generic kernel on the sphere² $\mathbf{a}_0 \in \mathbb{S}^{k-1}$, if the sparsity rate $\theta \lesssim k^{-2/3}$ and the number of measurement $m \gtrsim \text{poly}(k)$, initializing with some k consecutive entries of \mathbf{y} and applying any optimization method which (i) is a descent method, and (ii) converges to a local minimizer under a strict saddle hypothesis [JGN⁺17, XRKM17], produces a near shift-truncation of the ground truth.

¹Denote $\mathbf{y} = [y_1, y_2, \dots, y_{m-1}, y_m]^T$, then its reversal $\check{\mathbf{y}} = [y_1, y_m, y_{m-1}, \dots, y_2]^T$.

²Here, we refer a kernel sampled following a uniform distribution over the sphere as a generic kernel on the sphere.

1.1 Related Works

Even after accounting for the scale ambiguity, the general blind deconvolution problem remains ill-posed. Different types of prior knowledge about the unknown signals have been introduced and to make the blind deconvolution problem well posed. For example, if the signals a_0 and x_0 live on known linear subspaces, the blind deconvolution problem can be cast as a low-rank recovery problem, and solved via semidefinite programming. [ARR12] proves that if one of the subspaces is random and the other satisfies a spectral flatness condition, this approach recovers the pair (a_0, x_0) up to scale. [LLSW16] provides a more efficient nonconvex algorithm for blind deconvolution under this subspace model. [LS15] consider a more complicated model in which one of the signals is sparse in some known dictionary. [LLJB17] considers the case where both convolutional signals are sparse in some known dictionaries. These known dictionaries are assumed to be random (e.g., Gaussian or partial Fourier). Identifiability of these blind deconvolution problems is investigated in [LLB16, LLB17]. [LS17] further addresses a simultaneous demixing and deconvolution problem, where the observation is the superposition of multiple convolutions.

The above results offer efficient and guaranteed algorithms for blind deconvolution problems in which the signals of interest are sparse in a random dictionary. However, in the short-and-sparse blind deconvolution problem in microscopy image analysis or neural spike sorting, the sparse signal is *sparse with respect to the standard basis* rather than a random dictionary. Any cyclic shift of a standard basis is another standard basis, therefore the short-and-sparse blind deconvolution problem is only identifiable up to shifts. This is in contrast to the aforementioned random models, which only exhibit a scale ambiguity. When casting the short-and-sparse blind deconvolution problem as an optimization problem, this shift ambiguity creates a large group of equivalent global solutions (convolutional pairs of opposite shifts $s_\tau[\widetilde{a}_0]$ and $s_{-\tau}[x_0]$) and therefore much more complicated optimization landscape.

For sparsity in the standard basis, [CM14, CM15] show that sparsity alone is not sufficient for unique recovery, by demonstrating the existence of manifolds (a, x) of signals that are not identifiable from the convolution $y = a * x$. This construction requires both the support and magnitudes of the two signals to be regular: the support of x needs to have the form $J \cup s_1(J)$ for some set J , and the nonzero entries of x to take on specific values. When x is either Bernoulli or Bernoulli-Gaussian, with probability one, the pair (a, x) does not fall in this non-identifiable set. [Chi16] proposes a convex relaxation for a variant of the sparse blind deconvolution problem in which a lies in a random subspace and x is a superposition of spikes with continuous-valued locations. A strong point of this method is that it avoids discretization. Because of the random subspace model on a , the results of [Chi16] are not directly comparable to ours. However, if the rates from this work were adapted to the short-and-sparse setting, they would require x to be sparse enough that the observation y contains many isolated (non-overlapping) copies of a . This seems to reflect a fundamental limitation of convexification approaches in handling signals with multiple structures [OJF⁺15]. [WC16] studies another variant where multiple independent observations of circulant convolutions are available, motivated by multi-channel blind deconvolution. Although the convolution kernel is short compared to the total measurements, each independent "short" measurement is self contained. While in the short-and-sparse blind deconvolution problem, only one measurement is available and any "short" measurement heavily depends on adjacent measurements. This nuance leads to much more complicated optimization geometry.

Although the theory of short-and-sparse blind deconvolution remains completely open, many nonconvex algorithms have been developed and practiced in computer vision, where the convolution kernel captures the image blurring process due to camera shake [LWDF11]. Motivated by this physical model, people assume the convolutional kernel to be entry-wise nonnegative and sums up to 1, and then minimize the objective function of following form

$$\min_{a \geq 0, \|a\|_1=1} \min_x \frac{1}{2} \|y - a \circledast x\|_2^2 + \lambda \|x\|_* . \quad (1.6)$$

In the image deblurring application, x represents the gradient of a natural image and $\|\cdot\|_*$ penalizes the sparsity of x . However, such formulation always admits one local minimum obtained at the convolutional pair $(a, x) = (\delta, y)$ [BVG13, PF14]. In contrast, [WZ13, ZWZ13] carefully compare the difference in MAP and VB approaches, and propose to instead constrain a to have unit Frobenius norm – i.e., to reside on a high-dimensional sphere. [ZLK⁺17] studies the optimization landscape of the sphere constrained sparse

blind deconvolution and firstly identifies the structure of the local solutions. In particular, [ZLK⁺17] casts the short-and-sparse blind deconvolution problem as an optimization problem over the sphere:

$$\min_{\mathbf{a} \in \mathbb{S}^{k-1}} \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{a} \circledast \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (1.7)$$

and presents empirical evidence that *local minima* $\bar{\mathbf{a}}$ are close to certain shift truncations of \mathbf{a}_0 . [ZLK⁺17] further proves that a “linearized” version of (1.7), which neglects quadratic interactions in \mathbf{a} , satisfies this property, in the “dilute limit” in which the sparse signal \mathbf{x}_0 is a single spike. In this paper, we demonstrate that for a different objective function, this claim holds under much broader conditions than what is proved in [ZLK⁺17]. In particular, our results allow the sparse signal \mathbf{x}_0 to be much denser.

1.2 Assumptions and Notations

We assume that $\mathbf{x}_0 \in \mathbb{R}^m$ follows the Bernoulli-Gaussian (BG) model with sparsity level θ : $\mathbf{x}_0(i) = \omega_i g_i$ with $\omega_i \sim \text{Ber}(\theta)$ and $g_i \sim \mathcal{N}(0, 1)$, where all the different random variables are jointly independent. For simplicity, we write $\mathbf{x}_0 \sim_{i.i.d.} \text{BG}(\theta)$.

Throughout this paper, vectors $\mathbf{v} \in \mathbb{R}^k$ are indexed as $\mathbf{v} = [v_1, v_2, \dots, v_k]$, and $[\cdot]_m$ denotes the modulo operator of m . We use $\|\cdot\|_2$ to denote the operator norm, $\|\cdot\|_F$ to denote the Frobenius norm, and $\|\cdot\|_p$ to denote the entry wise ℓ^p norm. $(\cdot)_I$ denotes the projection onto subset with index I and $\mathcal{P}_\mathbb{S}[\cdot] = \frac{\cdot}{\|\cdot\|_F}$ denotes the projection onto the Frobenius sphere. $(\cdot)^{\circ p}$ is the entry wise p -th order exponent operator. We use C, c to denote positive constants, and their value change across the paper.

2 Problem Formulation and Main Results

In the short-and-sparse blind deconvolution problem, any k consecutive entries in \mathbf{y} only depend on $2k - 1$ consecutive entries in \mathbf{x}_0 :

$$\mathbf{y}_i = [y_i, \dots, y_{1+[i+k-1]_m}]^T = \sum_{\tau=-(k-1)}^{k-1} x_{1+[i+\tau-1]_m} \cdot \mathbf{t}_k^* s_\tau[\widetilde{\mathbf{a}}_0] \quad (2.1)$$

$$= \underbrace{\begin{bmatrix} a_k & a_{k-1} & \cdots & a_1 & \cdots & 0 & 0 \\ 0 & a_k & \cdots & a_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & a_{k-1} & \cdots & a_1 & 0 \\ 0 & 0 & \cdots & a_k & \cdots & a_2 & a_1 \end{bmatrix}}_{\mathbf{A}_0 \in \mathbb{R}^{k \times (2k-1)}} \underbrace{\begin{bmatrix} x_{1+[i-k]_m} \\ \vdots \\ x_i \\ \vdots \\ x_{1+[i+k-2]_m} \end{bmatrix}}_{\mathbf{x}_i \in \mathbb{R}^{(2k-1) \times 1}}. \quad (2.2)$$

Write $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m] \in \mathbb{R}^{k \times m}$ and $\mathbf{X}_0 = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{(2k-1) \times m}$. Using the above expression, we have that

$$\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0. \quad (2.3)$$

Each column \mathbf{x}_i of \mathbf{X}_0 only contains some $2k - 1$ entries of \mathbf{x}_0 . The *rows* of \mathbf{X}_0 are cyclic shifts of the reversal of \mathbf{x}_0 :

$$\mathbf{X}_0 = \begin{bmatrix} s_0[\tilde{\mathbf{x}}_0] \\ \vdots \\ s_{2k-2}[\tilde{\mathbf{x}}_0] \end{bmatrix}. \quad (2.4)$$

The shifts of $\tilde{\mathbf{x}}_0$ are *sparse vectors* in the linear subspace $\text{row}(\mathbf{X}_0)$. Note that if we could recover some shift $s_\tau[\mathbf{x}_0]$, we could subsequently determine $s_{-\tau}[\mathbf{a}_0]$ by solving a linear system of equations, and hence solve the deconvolution problem, up to the shift ambiguity.

2.1 Finding a Shifted Sparse Signal

In light of the above observations, a natural computational approach to sparse blind deconvolution is to attempt to find \mathbf{x}_0 by searching for a sparse vector in the linear subspace $\text{row}(\mathbf{X}_0)$, e.g., by solving an optimization problem

$$\min \|\mathbf{v}\|_* \quad \text{s. t.} \quad \mathbf{v} \in \text{row}(\mathbf{X}_0), \|\mathbf{v}\|_2 = 1, \quad (2.5)$$

where $\|\cdot\|_*$ is chosen to encourage sparsity of the target signal [SWW12, SQW15, QSW16, HSSS16].

In sparse blind deconvolution, we do not have access to the row space of \mathbf{X}_0 . Instead, we only observe the subspace $\text{row}(\mathbf{Y}) \subset \text{row}(\mathbf{X}_0)$. The subspace $\text{row}(\mathbf{Y})$ does not necessarily contain the desired sparse vector $\mathbf{e}_i^T \mathbf{X}_0$, but it *does* contain some approximately sparse vectors. In particular, consider following vector in $\text{row}(\mathbf{Y})$,

$$\mathbf{v} = \mathbf{Y}^T \mathbf{a}_0 = \underbrace{\check{\mathbf{x}}_0}_{\text{sparse}} + \underbrace{\sum_{i \neq 0} \langle \mathbf{a}_0, s_i[\mathbf{a}_0] \rangle s_i[\check{\mathbf{x}}_0]}_{\text{"noise"} \mathbf{z}}. \quad (2.6)$$

The vector \mathbf{v} is a superposition of a sparse signal $\check{\mathbf{x}}_0$ and its scaled shifts $\langle \mathbf{a}_0, s_i[\mathbf{a}_0] \rangle s_i[\check{\mathbf{x}}_0]$. If the shift-coherence $|\langle \mathbf{a}_0, s_\tau[\mathbf{a}_0] \rangle|$ is small³ and \mathbf{x}_0 is sparse enough, \mathbf{z} can be viewed as small noise.⁴ The vector \mathbf{v} is not sparse, but it is *spiky*: a few of its entries are much larger than the rest. We deploy a milder sparsity penalty $-\|\cdot\|_4^4$ to recover such a spiky vector, as $\|\cdot\|_4^4$ is very flat around 0 and insensitive to small noise in the signal.⁵ This gives

$$\min -\frac{1}{4} \|\mathbf{v}\|_4^4 \quad \text{s. t.} \quad \mathbf{v} \in \text{row}(\mathbf{Y}), \|\mathbf{v}\|_2 = 1. \quad (2.7)$$

We can express a generic unit vector $\mathbf{v} \in \text{row}(\mathbf{Y})$ as $\mathbf{v} = \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q}$, with $\|\mathbf{v}\|_2 = \|\mathbf{q}\|_2$. This leads to the following equivalent optimization problem over the sphere

$$\min \psi(\mathbf{q}) \doteq -\frac{1}{4m} \left\| \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_4^4 \quad \text{s. t.} \quad \|\mathbf{q}\|_2 = 1. \quad (2.8)$$

Interpretation: preconditioned shifts. This objective $\psi(\mathbf{q})$ can be rewritten as

$$\psi(\mathbf{q}) = -\frac{1}{4m} \left\| \check{\mathbf{y}} \circledast (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_4^4 = -\frac{1}{4m} \left\| \check{\mathbf{x}}_0 \circledast \mathbf{A}_0^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_4^4 \sim \|\check{\mathbf{x}}_0 \circledast \boldsymbol{\zeta}\|_4^4, \quad (2.9)$$

where $\boldsymbol{\zeta} = \mathbf{A}_0^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q}$. This approximation becomes accurate as m grows.⁶ This objective encourages the convolution of $\check{\mathbf{x}}_0$ and $\boldsymbol{\zeta}$ to be as spiky as possible. Reasoning analogous to (2.6) suggests that $\check{\mathbf{x}}_0 \circledast \boldsymbol{\zeta}$ will be spiky if

$$\boldsymbol{\zeta} = \mathbf{A}_0^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \approx \mathbf{e}_l, \quad l \in \{1, \dots, 2k-1\}. \quad (2.10)$$

For simplicity, we define the preconditioned convolution matrix

$$\mathbf{A} \doteq (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{A}_0 = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_{2k-1}], \quad (2.11)$$

with column coherence (preconditioned shift coherence) $\mu \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$. Then $\boldsymbol{\zeta}$ can also be interpreted as measuring the inner products of \mathbf{q} with columns of \mathbf{A} . Making this intuition rigorous, we will show that minimizing this objective over a certain region of the sphere yields a preconditioned shift truncate \mathbf{a}_l , from which we can recover a shift truncate of the original signal \mathbf{a}_0 .

2.2 Structured Local Minima

³For a generic kernel \mathbf{a}_0 , the shift-coherence is bounded by $|\langle \mathbf{a}_0, s_\tau[\mathbf{a}_0] \rangle| \approx 1/\sqrt{k}$ for any shift τ .

⁴In particular, under a Bernoulli-Gaussian model, for each j , $\mathbb{E}[\mathbf{z}_j^2] = \theta \sum_{i \neq 0} \langle \mathbf{a}_0, s_i[\mathbf{a}_0] \rangle^2$.

⁵In comparison, the classical choice $\|\cdot\|_* = \|\cdot\|_1$ is a strict sparsity penalty that essentially encourages all small entries to be 0.

⁶As $\mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)}[\mathbf{Y} \mathbf{Y}^T] = \mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)}[\mathbf{A}_0 \mathbf{X}_0 \mathbf{X}_0^T \mathbf{A}_0^T] = \theta m \mathbf{A}_0 \mathbf{A}_0^T$.

We will show that in a certain region $\mathcal{R}_{C_*} \subset \mathbb{S}^{k-1}$, the preconditioned shift truncations \mathbf{a}_l are the *only* local minimizers. Moreover, the other critical points in \mathcal{R}_{C_*} can be interpreted as resulting from competition between several of these local minima (Figure 2). At any saddle point, there exists strict negative curvature in the direction of a nearby local minimizer which breaks the balance in favor of some particular \mathbf{a}_l . The region \mathcal{R}_{C_*} is defined as follows:

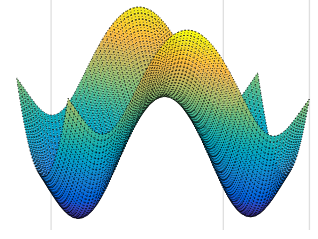


Figure 2: Saddles points are approximately superpositions of local minima.

Definition 2.1. For fixed $C_* > 0$, letting κ denote the condition number of \mathbf{A}_0 , and $\mu \doteq \max_{i \neq j} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$ the column coherence of \mathbf{A} , we define two regions \mathcal{R}_{C_*} , $\hat{\mathcal{R}}_{C_*} \subset \mathbb{S}^{k-1}$, as

$$\mathcal{R}_{C_*} \doteq \left\{ \mathbf{q} \in \mathbb{S}^{k-1} \mid \|\mathbf{A}^T \mathbf{q}\|_4^6 \geq C_* \mu \kappa^2 \|\mathbf{A}^T \mathbf{q}\|_3^3 \right\}. \quad (2.12)$$

$$\hat{\mathcal{R}}_{C_*} \doteq \left\{ \mathbf{q} \in \mathbb{S}^{k-1} \mid \|\mathbf{A}^T \mathbf{q}\|_4^6 \geq C_* \mu \kappa^2 \right\} \subseteq \mathcal{R}_{C_*}. \quad (2.13)$$

A simpler and smaller region $\hat{\mathcal{R}}_{C_*}$ is also introduced in Definition (2.1). This region $\hat{\mathcal{R}}_{C_*}$ can be viewed as a sub-level set for $-\|\mathbf{A}^T \mathbf{q}\|_4^4$, which is proportional to the objective value $\psi(\mathbf{q})$ assuming m is sufficiently large⁷. Therefore, once initialized within $\hat{\mathcal{R}}_{C_*}$, the iterates produced by a descent algorithm will stay in $\hat{\mathcal{R}}_{C_*}$.

In particular, at any stationary point $\mathbf{q} \in \mathcal{R}_{10}$, the local optimization landscape can be characterized in terms of the number of spikes (entries with nontrivial magnitude⁸) in ζ . If there is only one spike in ζ , then such stationary point \mathbf{q} is a local minimum that is close to one local minimizer; if there are more than two spikes in ζ , then such stationary point \mathbf{q} is saddle point. Based on the above characterizations of stationary points in \mathcal{R}_{C_*} with $C_* \geq 10$, we can deduce that any local minimum is close to some \mathbf{a}_l , a preconditioned shift truncation of the ground truth \mathbf{a}_0 .

Theorem 2.2 (Main Result). Assuming observation $\mathbf{y} \in \mathbb{R}^m$ is the circulant convolution of $\mathbf{a}_0 \in \mathbb{R}^k$ and $\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta) \in \mathbb{R}^m$, where the convolutional matrix \mathbf{A}_0 has minimum singular value $\sigma_{\min} > 0$ and condition number $\kappa \geq 1$, and \mathbf{A} has column incoherence $0 \leq \mu < 1$. There exists a positive constant C such that whenever the number of measurements

$$m \geq C \frac{\min \{ \mu^{-4/3}, \kappa^2 k^2 \}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \quad (2.14)$$

and $\theta \geq \log k/k$, then with high probability, any local optima $\bar{\mathbf{q}} \in \hat{\mathcal{R}}_{2C_*}$ satisfies

$$|\langle \bar{\mathbf{q}}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_l] \rangle| \geq 1 - c_* \kappa^{-2} \quad (2.15)$$

for some integer $1 \leq l \leq 2k-1$. Here, $C_* \geq 10$ and $c_* = 1/C_*$.

This theorem says that any local minimum in $\hat{\mathcal{R}}_{2C_*}$ is close to some normalized column of \mathbf{A} given polynomially many observation. The parameters σ_{\min} , κ and μ effectively measure the spectrum flatness of the ground truth kernel \mathbf{a}_0 and characterize how broad the results hold. A random like kernel usually has big σ_{\min} , small κ and μ , which equivalently implies the result holds in a large sub-level set $\hat{\mathcal{R}}_{2C_*}$, even with fewer observations.

Hence, once assuring the algorithm finds a local minimum in $\hat{\mathcal{R}}_{2C_*}$, then some shifted truncation of the ground truth kernel \mathbf{a}_0 can be recovered. In other words, if we can find an initialization point with small objective value, then a descent algorithm minimizing the objective function guarantees that \mathbf{q} always stays in $\hat{\mathcal{R}}_{2C_*}$ in proceeding iterations. Therefore, any descent algorithm that escapes a strict saddle point can be applied to find some \mathbf{a}_l , or some shift truncation of \mathbf{a}_0 .

⁷Please refer to Section 3 for more arguments.

⁸We call any ζ_l with magnitude no smaller than $2\mu \|\zeta\|_3^3 / \|\zeta\|_4^4$ to be nontrivial and defer technical reasonings to later sections.

2.3 Initialization with a Random Sample

Recall that $\mathbf{y}_i = \mathbf{A}_0 \mathbf{x}_i$, which is a sparse superposition of about $2\theta k$ columns of \mathbf{A}_0 . Intuitively speaking, such \mathbf{q}_{init} already encodes certain preferences towards a few preconditioned shift truncations of the ground truth. Therefore, we randomly choose an index i and set the initialization point as

$$\mathbf{q}_{\text{init}} = \mathcal{P}_{\mathbb{S}} \left[(\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{y}_i \right]. \quad (2.16)$$

Using $\mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)} [\mathbf{Y} \mathbf{Y}^T] = \theta m \mathbf{A}_0 \mathbf{A}_0^T$ again, we have

$$\zeta_{\text{init}} = \mathbf{A}^T \mathbf{q}_{\text{init}} \approx \mathcal{P}_{\mathbb{S}} [\mathbf{A}^T \mathbf{A} \mathbf{x}_i]. \quad (2.17)$$

For a generic kernel $\mathbf{a}_0 \in \mathbb{S}^{k-1}$, $\mathbf{A}^T \mathbf{A}$ is close to a diagonal matrix, as the magnitudes of off-diagonal entries are bounded by column incoherence μ . Hence, the sparse property of \mathbf{x}_i can be approximately preserved, that $\mathcal{P}_{\mathbb{S}} [\mathbf{A}^T \mathbf{A} \mathbf{x}_i]$ is spiky vector with small $\|\cdot\|_4^4$. By leveraging the sparsity level θ , one can make sure such initialization point \mathbf{q}_{init} falls in $\hat{\mathcal{R}}_{2C_*}$. Therefore, we propose Algorithm 1 for solving sparse blind deconvolution with its working conditions stated in Corollary 2.3. For the choice of descent algorithms which escape strict saddle points, there are several such algorithms specially tailored for sphere constrained optimization problems [ABG07, GWY09].

Algorithm 1 Short and Sparse Blind Deconvolution

Input: Observations $\mathbf{y} \in \mathbb{R}^m$ and kernel size k .

Output: Recovered Kernel $\bar{\mathbf{a}}$.

- 1: Generate random index $i \in [1, m]$ and set

$$\mathbf{q}_{\text{init}} = \mathcal{P}_{\mathbb{S}} \left[(\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{y}_i \right].$$

- 2: Solve following nonconvex optimization problem with a descent algorithm that escapes saddle point and find a local minimizer

$$\bar{\mathbf{q}} = \arg \min_{\mathbf{q} \in \mathbb{S}^{k-1}} \varphi(\mathbf{q}).$$

- 3: Set $\bar{\mathbf{a}} = \mathcal{P}_{\mathbb{S}} \left[(\mathbf{Y} \mathbf{Y}^T)^{1/2} \bar{\mathbf{q}} \right]$.
-

Corollary 2.3. Suppose the ground truth \mathbf{a}_0 kernel has preconditioned shift coherence $0 \leq \mu \leq \frac{1}{8 \times 48} \log^{-3/2}(k)$ and sparse coefficient $\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta) \in \mathbb{R}^m$. There exist positive constants $C \geq 2560^4$ and C' such that whenever the sparsity level

$$64k^{-1} \log k \leq \theta \leq \min \left\{ \frac{1}{48^2} \mu^{-2} k^{-1} \log^{-2} k, \left(\frac{1}{4} - \frac{640}{C^{1/4}} \right) (3C_* \mu \kappa^2)^{-2/3} k^{-1} (1 + 36\mu^2 k \log k)^{-2} \right\},$$

and signal length

$$m \geq \max \left\{ C \theta^2 \sigma_{\min}^{-2} \kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log(\kappa k), C' (1 - \theta)^{-2} \sigma_{\min}^{-2} \min \{ \mu^{-1}, \kappa^2 k^2 \} \kappa^8 k^4 \log^3(\kappa k) \right\},$$

then with high probability, Algorithm 1 recovers $\bar{\mathbf{a}}$ such that

$$\|\bar{\mathbf{a}} \pm \mathcal{P}_{\mathbb{S}} [\boldsymbol{\nu}_k s_\tau [\bar{\mathbf{a}}_0]]\|_2 \leq 4\sqrt{c_*} + ck^{-1} \quad (2.18)$$

for some integer shift $-(k-1) \leq \tau \leq k-1$.

For a generic $\mathbf{a}_0 \in \mathbb{S}^{k-1}$, plugging in the numerical estimation of the parameters σ_{\min} , κ and μ (Figure 3), accurate recovery can be obtained with $m \gtrsim \theta^2 k^6 \text{poly log}(k)$ measurements and sparsity level $\theta \lesssim k^{-2/3} \text{poly log}(k)$. For bandpass kernels \mathbf{a}_0 , σ_{\min} is smaller and κ, μ are larger, and so our results require \mathbf{x}_0 to be longer and sparser.

3 Asymptotic Function Landscape

In the next two sections, we discuss some key elements of our analysis. In this section, we first investigate the stationary points of the “population” objective $\mathbb{E}_{\mathbf{x}_0}[\psi(\mathbf{q})]$. We demonstrate that any local minimizer in \mathcal{R}_{C_*} is close to a signed column of \mathbf{A} , a preconditioned shift truncation of \mathbf{a}_0 . In the next section, we then demonstrate that when m is sufficiently large, the “finite sample” objective $\psi(\mathbf{q})$ satisfies the same property.

In Section 3.1, we show how to accurately estimate the vector $\boldsymbol{\zeta} = \mathbf{A}^T \mathbf{q}$ at any stationary point $\mathbf{q} \in \mathcal{R}_{C_*}$. In Section 3.2, we show how the number of spikes in $\boldsymbol{\zeta}$ determines the geometry around a stationary point.

- For any stationary point $\mathbf{q} \in \mathcal{R}_{C_*}$, its preconditioned cross-correlation $\boldsymbol{\zeta}$ has *at least* one large entry (Section 3.2.1). This implies that any stationary point \mathbf{q} must be close some local minimizer.
- If $\boldsymbol{\zeta}$ has *only* one large entry, then \mathbf{q} is a local minimizer. (Section 3.2.2)
- If $\boldsymbol{\zeta}$ has *more than* one large entry, then \mathbf{q} is a strict saddle point. (Section 3.2.3)

With above three characterizations, we can deduce that any local minimizer in \mathcal{R}_{C_*} is close to some column of \mathbf{A} , a preconditioned shift truncation of \mathbf{a}_0 .

3.1 Stationary Points

Using $\mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)}[\mathbf{Y}\mathbf{Y}^T] = \theta m \mathbf{A}_0 \mathbf{A}_0^T$ again, the expectation of the objective function $\psi(\mathbf{q})$ can be approximated⁹ as

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)}[\psi(\mathbf{q})] &\approx \mathbb{E}_{\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta)} \left[-\frac{1}{m} \left\| \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_4^4 \right] \\ &= -\frac{1}{\theta^2 m^2} \left[3\theta (1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^4 + 3\theta^2 \|\mathbf{A}^T \mathbf{q}\|_2^4 \right] \\ &= -\frac{3(1-\theta)}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4 - \frac{3}{m^2}. \end{aligned} \quad (3.1)$$

In the next section, we will argue that the critical points of the finite sample objective $\psi(\mathbf{q})$ are close to those of the asymptotic approximation ϕ . We can therefore study the critical points of ψ by studying the simpler problem

$$\min_{\mathbf{q} \in \mathbb{R}^{k-1}} \varphi(\mathbf{q}) \doteq -\frac{1}{4} \|\mathbf{A}^T \mathbf{q}\|_4^4 = -\frac{1}{4} \|\boldsymbol{\zeta}\|_4^4. \quad (3.2)$$

The Euclidean gradient and Hessian for $\varphi(\mathbf{q})$ can be calculated as

$$\nabla \varphi(\mathbf{q}) = -\mathbf{A} \boldsymbol{\zeta}^{\circ 3}, \quad (3.3)$$

$$\nabla^2 \varphi(\mathbf{q}) = -3\mathbf{A} \text{diag}(\boldsymbol{\zeta}^{\circ 2}) \mathbf{A}^T. \quad (3.4)$$

We can study the critical points of φ over the sphere using the *Riemannian* gradient and hessian [AMS07]

$$\text{grad } \varphi(\mathbf{q}) = \mathbf{P}_{\mathbf{q}^\perp} [\nabla \varphi(\mathbf{q})] \quad (3.5)$$

$$= -\mathbf{A} \boldsymbol{\zeta}^{\circ 3} + \mathbf{q} \|\boldsymbol{\zeta}\|_4^4, \quad (3.6)$$

$$\text{Hess } \varphi(\mathbf{q}) = \mathbf{P}_{\mathbf{q}^\perp} [\nabla^2 \varphi(\mathbf{q}) - \langle \nabla \varphi(\mathbf{q}), \mathbf{q} \rangle \mathbf{I}] \mathbf{P}_{\mathbf{q}^\perp} \quad (3.7)$$

⁹In Lemma A.1 of the appendix, we give a detailed derivation of this approximation.

$$= -\mathbf{P}_{\mathbf{q}^\perp} \left[3\mathbf{A} \operatorname{diag}(\zeta^{\circ 2}) \mathbf{A}^T - \|\zeta\|_4^4 \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp}. \quad (3.8)$$

Here, $\mathbf{P}_{\mathbf{q}^\perp} = \mathbf{I} - \mathbf{q}\mathbf{q}^T$ denotes the projection onto the tangent space of the Frobenius sphere at point $\mathbf{q} \in \mathbb{S}^{k-1}$.

As in the Euclidean space, a stationary point on the sphere satisfies $\operatorname{grad} [\varphi] (\mathbf{q}) = \mathbf{0}$. Using (3.6), at any stationary point of φ ,

$$\mathbf{A}\zeta^{\circ 3} - \mathbf{q} \|\zeta\|_4^4 = \mathbf{0}. \quad (3.9)$$

Left-multiplying both sides of the equation by \mathbf{A}^T , we have

$$\mathbf{A}^T \mathbf{A} \zeta^{\circ 3} - \mathbf{A}^T \mathbf{q} \|\zeta\|_4^4 = \mathbf{0}. \quad (3.10)$$

For the i -th entry, following equality always holds

$$0 = \|\mathbf{a}_i\|_2^2 \zeta_i^3 + \sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3 - \zeta_i \|\zeta\|_4^4 \quad (3.11)$$

$$\Rightarrow 0 = \zeta_i^3 - \underbrace{\zeta_i \frac{\|\zeta\|_4^4}{\|\mathbf{a}_i\|_2^2}}_{\alpha_i} + \underbrace{\frac{\sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3}{\|\mathbf{a}_i\|_2^2}}_{\beta_i}. \quad (3.12)$$

For simplicity, we deploy the following notations

$$\alpha_i = \frac{\|\zeta\|_4^4}{\|\mathbf{a}_i\|_2^2}, \quad \beta_i = \frac{\sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3}{\|\mathbf{a}_i\|_2^2}. \quad (3.13)$$

If $\alpha_i \gg \beta_i$, Proposition 3.1 shows that ζ_i is very close to one of three values: 0, or $\pm\sqrt{\alpha_i}$.

Proposition 3.1. *Let $\mathbf{q} \in \mathbb{S}^{k-1}$ be a stationary point satisfying $\|\mathbf{A}^T \mathbf{q}\|_4^6 \geq 4\mu \|\mathbf{A}^T \mathbf{q}\|_3^3$, then the i -th entry of $\zeta = \mathbf{A}^T \mathbf{q}$ falls in the range*

$$\{0, \pm\sqrt{\alpha_i}\} \pm \frac{2\beta_i}{\alpha_i}, \quad (3.14)$$

with

$$\alpha_i = \frac{\|\zeta\|_4^4}{\|\mathbf{a}_i\|_2^2}, \quad \beta_i = \frac{\sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3}{\|\mathbf{a}_i\|_2^2}. \quad (3.15)$$

Proof Since $\|\zeta\|_4^6 \geq 4\mu \|\zeta\|_3^3$, in this case, $\beta_i \leq \frac{1}{4}\alpha_i^{3/2}$ is satisfied, as

$$\|\zeta\|_4^6 \geq 4\mu \|\zeta\|_3^3 \geq 4\|\mathbf{a}_i\|_2 \sum_{i \neq j} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3. \quad (3.16)$$

The roots can be estimated by applying Lemma A.2 with

$$\sqrt{\alpha_i} = \frac{\|\zeta\|_4^2}{\|\mathbf{a}_i\|_2}, \quad (3.17)$$

$$\frac{2\beta_i}{\alpha_i} = \frac{2 \sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3}{\|\zeta\|_4^4} \leq \frac{2\mu \|\zeta\|_3^3}{\|\zeta\|_4^4}. \quad (3.18)$$

■

This implies that either $|\langle \mathbf{a}_i, \mathbf{q} \rangle|$ is large ($\approx \sqrt{\alpha_i}$) or it is very close to zero.

3.2 Function Landscape on \mathcal{R}_{C_\star}

In this section, we study the optimization landscape around a stationary point \mathbf{q} by bounding the eigenvalues of the Riemannian Hessian $\text{Hess}[\varphi](\mathbf{q})$: if $\text{Hess}[\varphi](\mathbf{q})$ is positive semidefinite, then the φ is convex in a neighborhood of \mathbf{q} and hence \mathbf{q} is a local minimum; if $\text{Hess}[\varphi](\mathbf{q})$ has a negative eigenvalue, then there exists a direction along which the objective value decreases and hence \mathbf{q} is a saddle point.

Note that the Riemannian Hessian $\text{Hess}[\varphi](\mathbf{q})$ at stationary point \mathbf{q} is a function of ζ which can be accurately estimated when constrained in \mathcal{R}_{C_\star} with $C_\star \geq 10$. By plugging the estimation of ζ in the Riemannian Hessian, we can bound the eigenvalues of $\text{Hess}[\varphi](\mathbf{q})$, and hence we can characterize the optimization landscape around a stationary point \mathbf{q} .

3.2.1 Nontrivial Preference of a Stationary Point

First, we demonstrate that for any stationary point $\mathbf{q} \in \mathcal{R}_{C_\star}$ with $C_\star \geq 10$, ζ must have at least one large entry.

Lemma 3.2. *For any stationary point $\mathbf{q} \in \mathcal{R}_{C_\star}$ with $C_\star \geq 10$,*

$$\|\zeta\|_\infty \geq \frac{2\mu \|\zeta\|_3^3}{\|\zeta\|_4^4}. \quad (3.19)$$

Proof We give a proof by contradiction. Suppose that $\mathbf{q} \in \mathcal{R}_{C_\star}$ with $C_\star \geq 10$, and every entry of ζ has small magnitude such that $\|\zeta\|_\infty < 2\mu \|\zeta\|_3^3 / \|\zeta\|_4^4$, then

$$\|\zeta\|_4^4 \leq \|\zeta\|_\infty^2 \leq \left(\frac{2\beta_i}{\alpha_i} \right)^2 \leq \frac{4\mu^2 \|\zeta\|_3^6}{\|\zeta\|_4^8}, \quad (3.20)$$

which indicates $\|\zeta\|_4^6 \leq 2\mu \|\zeta\|_3^3$ and contradicts the assumption $\|\zeta\|_4^6 > C_\star \mu \kappa^2 \|\zeta\|_3^3$. Therefore, at least one entry of ζ has large enough magnitude. ■

Geometrically, the nontrivial entry ζ_i indicates the preference to corresponding column \mathbf{a}_i , as $\zeta_i = \langle \mathbf{a}_i, \mathbf{q} \rangle$. Therefore, Lemma 3.2 implies that any stationary point \mathbf{q} in \mathcal{R}_{C_\star} should be close to at least one column of \mathbf{A} .

3.2.2 Local Minima

Suppose $\mathbf{q} \in \mathcal{R}_{C_\star}$ ($C_\star \geq 10$) is a stationary point and vector ζ only has one nontrivial entry ζ_l , then we can demonstrate that the Riemannian Hessian $\text{Hess} \varphi(\mathbf{q})$ is positive definite, and hence \mathbf{q} is a local minimizer near \mathbf{a}_l .

Lemma 3.3. *Suppose \mathbf{q} is a stationary point in \mathcal{R}_{C_\star} with $C_\star \geq 10$, and $\zeta = \mathbf{A}^T \mathbf{q}$ has only one entry ζ_l of magnitude no smaller than $2\mu \|\zeta\|_3^3 / \|\zeta\|_4^4$. Then \mathbf{q} is a local minimum near \mathbf{a}_l and $|\langle \mathbf{q}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_l] \rangle| > 1 - 2c_\star \kappa^{-2}$ with $c_\star = 1/C_\star$.*

Proof Suppose ζ has only one big entry ζ_l , and other entries are bounded by $2\beta_l/\alpha_l$

$$\|\zeta\|_4^4 = \zeta_l^4 + \sum_{j \neq l} \zeta_j^4 \quad (3.21)$$

$$\leq \zeta_l^4 + \max_{j \neq l} \zeta_j^2 \cdot \sum_{j \neq l} \zeta_j^2 \quad (3.22)$$

$$\leq \zeta_l^4 + \frac{4\mu^2 \|\zeta\|_3^6}{\|\zeta\|_4^8}, \quad (3.23)$$

with $\|\zeta\|_4^6 \geq C_\star \mu \kappa^2 \|\zeta\|_3^3$, and for simplicity let $c_\star = 1/C_\star$, we have

$$\zeta_l^4 \geq \|\zeta\|_4^4 - \frac{4\mu^2 \|\zeta\|_3^6}{\|\zeta\|_4^8} \geq (1 - 4c_\star^2 \kappa^{-4}) \|\zeta\|_4^4. \quad (3.24)$$

On the other hand, we also have

$$\zeta_l^2 \leq \left(\sqrt{\alpha_l} + \frac{2\beta_l}{\alpha_l} \right)^2 \quad (3.25)$$

$$\leq \frac{\|\zeta\|_4^4}{\|\mathbf{a}_i\|_2^2} + \frac{4\mu \|\zeta\|_3^3}{\|\mathbf{a}_i\|_2 \|\zeta\|_4^2} + \frac{4\mu^2 \|\zeta\|_3^6}{\|\zeta\|_4^8} \quad (3.26)$$

$$\leq \frac{\|\zeta\|_4^4}{\|\mathbf{a}_i\|_2^2} (1 + 4c_\star \kappa^{-2} + 4c_\star^2 \kappa^{-4}). \quad (3.27)$$

Combining above two inequalities, we have

$$\zeta_l^2 \leq \frac{1 + 4c_\star \kappa^{-2} + 4c_\star^2 \kappa^{-4}}{1 - 4c_\star^2 \kappa^{-4}} \frac{\zeta_l^4}{\|\mathbf{a}_i\|_2^2}, \quad (3.28)$$

thus the local minimum \mathbf{q} is close to \mathbf{a}_l :

$$\frac{|\langle \mathbf{q}, \mathbf{a}_l \rangle|}{\|\mathbf{a}_l\|_2} \geq \frac{\sqrt{1 - 4c_\star^2 \kappa^{-4}}}{1 + 2c_\star \kappa^{-2}} \geq 1 - 2c_\star \kappa^{-2}. \quad (3.29)$$

Next, we need to verify that the Riemannian Hessian at $\bar{\mathbf{q}}$ is definite positive, recall that

$$\text{Hess } \varphi(\mathbf{q}) = -\mathbf{P}_{\mathbf{q}^\perp} \left[3\mathbf{A} \text{diag}(\zeta^{\circ 2}) \mathbf{A}^T - \|\zeta\|_4^4 \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp}. \quad (3.30)$$

Let \mathbf{v} be a unit vector such that $\mathbf{v} \perp \mathbf{q}$, then

$$\mathbf{v}^T \text{Hess } \varphi(\mathbf{q}) \mathbf{v} \quad (3.31)$$

$$= -\mathbf{v}^T \left(3\mathbf{A} \text{diag}(\zeta^{\circ 2}) \mathbf{A}^T - \|\zeta\|_4^4 \mathbf{I} \right) \mathbf{v} \quad (3.32)$$

$$= \|\zeta\|_4^4 - 3\mathbf{v}^T \mathbf{A} \text{diag}(\zeta^{\circ 2}) \mathbf{A}^T \mathbf{v} \quad (3.33)$$

$$= \|\zeta\|_4^4 - 3 \langle \mathbf{a}_l, \mathbf{v} \rangle^2 \zeta_l^2 - 3 \sum_{i \neq l} \langle \mathbf{a}_i, \mathbf{v} \rangle^2 \zeta_i^2 \quad (3.34)$$

$$\geq \|\zeta\|_4^4 - 3 \langle \mathbf{a}_l, \mathbf{v} \rangle^2 \zeta_l^2 - 3 \max_{i \neq l} \zeta_i^2. \quad (3.35)$$

The last inequality is due to $\sum_{i \neq l} \langle \mathbf{a}_i, \mathbf{v} \rangle^2 \leq \|\mathbf{A}^T \mathbf{v}\|_2^2 = 1$. Since $\mathbf{v} \perp \bar{\mathbf{q}}$ and ζ_l is the only entry with nontrivial magnitude, then derive from (3.29):

$$\langle \mathbf{a}_l, \mathbf{v} \rangle^2 \zeta_l^2 \leq 2c_\star \|\mathbf{a}_l\|_2^2 \left(\sqrt{\alpha_l} + \frac{2\beta_l}{\alpha_l} \right)^2 \quad (3.36)$$

$$\leq 2c_\star \|\mathbf{a}_l\|_2^2 \cdot (1 + 2c_\star)^2 \alpha_l \quad (3.37)$$

$$\leq 2c_\star (1 + 2c_\star^2)^2 \|\zeta\|_4^4, \quad (3.38)$$

and

$$\max_{i \neq l} \zeta_i^2 \leq \frac{4\beta^2}{\alpha^2} \leq \frac{4\mu^2 \|\zeta\|_3^6}{\|\zeta\|_4^8} \leq \frac{4c_\star^2 \|\zeta\|_4^{12}}{\|\zeta\|_4^8} \leq 4c_\star^2 \|\zeta\|_4^4. \quad (3.39)$$

Hence, the inequality $\mathbf{v}^T \text{Hess } \varphi(\mathbf{q}) \mathbf{v} \geq (1 - 6c_\star - 36c_\star^2 - 24c_\star^3) \|\zeta\|_4^4$ holds for any \mathbf{v} satisfying $\mathbf{v} \perp \mathbf{q}$, thus implies positive curvature along any tangent direction at such stationary point \mathbf{q} when $C_\star \geq 10$. ■

The lemma says if \mathbf{q} is a stationary point in \mathbf{R}_{C_\star} and \mathbf{q} is only close to one column \mathbf{a}_l , then \mathbf{q} is a local minimizer and satisfies $|\langle \mathbf{q}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_l] \rangle| > 1 - 2c_\star \kappa^{-2}$ with $c_\star = 1/C_\star$.

3.2.3 Saddle Points

At last, if $\mathbf{q} \in \mathcal{R}_{C_\star}$ ($C_\star \geq 10$) is a stationary point and vector $\boldsymbol{\zeta}$ has more than one nontrivial entry. Denote any two nontrivial entries of $\boldsymbol{\zeta}$ with ζ_l and $\zeta_{l'}$, then we can prove that the Riemannian Hessian $\text{Hess } \varphi(\mathbf{q})$ has negative curvature in the span of \mathbf{a}_l and $\mathbf{a}_{l'}$, hence \mathbf{q} is a saddle point.

Lemma 3.4. *Suppose \mathbf{q} is a stationary point in \mathcal{R}_{C_\star} with $C_\star \geq 10$, and $\boldsymbol{\zeta} = \mathbf{A}^T \mathbf{q}$ has at least two entries ζ_l and $\zeta_{l'}$ with magnitude magnitude $\geq 2\mu \|\boldsymbol{\zeta}\|_3^3 / \|\boldsymbol{\zeta}\|_4^4$, then the Riemannian Hessian at \mathbf{q} has at least one negative eigenvalue and \mathbf{q} is a saddle point.*

Proof Suppose $\boldsymbol{\zeta}$ has at least two big entries ζ_l and $\zeta_{l'}$ satisfying

$$\zeta_l^2 \geq \left(\sqrt{\alpha_l} - \frac{2\beta_l}{\alpha_l} \right)^2 \quad (3.40)$$

$$\geq \frac{\|\boldsymbol{\zeta}\|_4^4}{\|\mathbf{a}_l\|_2^2} - \frac{4\mu \|\boldsymbol{\zeta}\|_3^3}{\|\boldsymbol{\zeta}\|_4^2 \|\mathbf{a}_l\|_2} + \frac{4\mu^2 \|\boldsymbol{\zeta}\|_3^6}{\|\boldsymbol{\zeta}\|_4^8} \quad (3.41)$$

$$> \frac{\|\boldsymbol{\zeta}\|_4^4}{\|\mathbf{a}_l\|_2^2} - \frac{4\mu \|\boldsymbol{\zeta}\|_3^3}{\|\boldsymbol{\zeta}\|_4^2 \|\mathbf{a}_l\|_2}, \quad (3.42)$$

and $\zeta_{l'}$ likewise. Since the nontrivial entry $\zeta_l = \langle \mathbf{a}_l, \mathbf{q} \rangle$, and again let $c_\star = 1/C_\star$, it is easy to show that the norm of \mathbf{a}_l is sufficiently large:

$$\|\mathbf{a}_l\|_2^2 \geq \zeta_l^2 \geq \left(\sqrt{\alpha_l} - \frac{2\beta_l}{\alpha_l} \right)^2 \quad (3.43)$$

$$\geq (1 - 2c_\star)^2 \frac{\|\boldsymbol{\zeta}\|_4^4}{\|\mathbf{a}_l\|_2^2} \quad (3.44)$$

$$\geq (1 - c_\star)^2 C_\star^{2/3} \frac{\mu^{2/3} \|\boldsymbol{\zeta}\|_3^2}{\|\mathbf{a}_l\|_2^2}, \quad (3.45)$$

or

$$\|\mathbf{a}_l\|_2 \geq (1 - c_\star)^{1/2} C_\star^{1/6} \mu^{1/6} \|\boldsymbol{\zeta}\|_3^{1/2}. \quad (3.46)$$

Similar result holds for $\|\mathbf{a}_{l'}\|_2$, therefore

$$\frac{\mu}{\|\mathbf{a}_l\|_2 \|\mathbf{a}_{l'}\|_2} \leq \frac{\mu^{2/3}}{C_\star^{1/3} \|\boldsymbol{\zeta}\|_3} \leq \frac{C_\star^{-2/3} \|\boldsymbol{\zeta}\|_4^4}{C_\star^{1/3} \|\boldsymbol{\zeta}\|_3^3} \leq c_\star. \quad (3.47)$$

Now we are ready to show there exists a unit vector \mathbf{v} such that $\mathbf{v} \in \text{span}(\mathbf{a}_l, \mathbf{a}_{l'})$ and $\mathbf{v} \perp \mathbf{q}$, and the Hessian has negative curvature along such \mathbf{v} :

$$\begin{aligned} & \mathbf{v}^T \text{Hess } \varphi(\mathbf{q}) \mathbf{v} \\ &= -3\mathbf{v}^T \mathbf{A} \text{diag}(\boldsymbol{\zeta}^2) \mathbf{A}^T \mathbf{v} + \|\boldsymbol{\zeta}\|_4^4 \end{aligned} \quad (3.48)$$

$$\leq -3\mathbf{v}^T (\mathbf{a}_l \zeta_l^2 \mathbf{a}_l^T + \mathbf{a}_{l'} \zeta_{l'}^2 \mathbf{a}_{l'}^T) \mathbf{v} + \|\boldsymbol{\zeta}\|_4^4 \quad (3.49)$$

$$\begin{aligned} & < -3 \left(\left| \left\langle \frac{\mathbf{a}_l}{\|\mathbf{a}_l\|_2}, \mathbf{v} \right\rangle \right|^2 + \left| \left\langle \frac{\mathbf{a}_{l'}}{\|\mathbf{a}_{l'}\|_2}, \mathbf{v} \right\rangle \right|^2 \right) \|\boldsymbol{\zeta}\|_4^4 \\ & \quad + \frac{4\mu \|\boldsymbol{\zeta}\|_3^3}{\|\boldsymbol{\zeta}\|_4^2} (\|\mathbf{a}_l\|_2 + \|\mathbf{a}_{l'}\|_2) + \|\boldsymbol{\zeta}\|_4^4 \\ & < -3 \left(1 - \frac{\mu}{\|\mathbf{a}_l\|_2 \|\mathbf{a}_{l'}\|_2} \right) \|\boldsymbol{\zeta}\|_4^4 \end{aligned} \quad (3.50)$$

$$+ \frac{4\mu \|\zeta\|_3^3}{\|\zeta\|_4^2} (\|\mathbf{a}_l\|_2 + \|\mathbf{a}_{l'}\|_2) + \|\zeta\|_4^4 \quad (3.51)$$

$$\leq (-2 + 11c_\star) \|\zeta\|_4^4. \quad (3.52)$$

The third inequality is implied by [Lemma A.3](#) and is negative when $C_\star \geq 10$. ■

This lemma says if the stationary point \mathbf{q} has large inner product with any two columns \mathbf{a}_l and $\mathbf{a}_{l'}$, then this \mathbf{q} is a saddle point and the objective value decreases along the direction that breaks symmetry between \mathbf{a}_l and $\mathbf{a}_{l'}$. The saddle point \mathbf{q} can be seen as resulting from the competition between the two target solutions \mathbf{a}_l and $\mathbf{a}_{l'}$.

4 Large Sample Concentration

In this section, we argue that the geometric characteristics of $\psi(\mathbf{q})$ are similar to those of $\varphi(\mathbf{q})$, by demonstrating that the critical points of the finite sample objective function $\psi(\mathbf{q})$ are similar to those of the asymptotic objective function $\varphi(\mathbf{q})$:

- **Critical points are close.** The Riemannian gradient ([Lemma 4.2](#)) and Hessian ([Lemma 4.3](#)) concentrate, such that there is a bijection between critical points \mathbf{q}_φ of φ and critical points \mathbf{q}_ψ of ψ , with $\|\mathbf{q}_\varphi - \mathbf{q}_\psi\|_2$ small.
- **Curvature is preserved.** The Riemannian Hessian ([Lemma 4.3](#)) concentrates, such that $\text{Hess}[\psi](\mathbf{q}_{\text{fs}})$ has a negative eigenvalue if and only if $\text{Hess}[\varphi](\mathbf{q}_{\text{pop}})$ has a negative eigenvalue, and $\text{Hess}[\psi](\mathbf{q}_{\text{fs}})$ is positive definite if and only if $\text{Hess}[\varphi](\mathbf{q}_{\text{pop}})$ is positive definite.

This implies that every local minimizer of the finite sample objective function is close to a preconditioned shift-truncation ([Lemma 4.1](#)).

Lemma 4.1. *If the following inequalities hold*

$$\begin{aligned} & \left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq \frac{3c_\star}{2\kappa^2} \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^6, \end{aligned} \quad (4.1)$$

$$\begin{aligned} & \left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq 3(1 - 6c_\star - 36c_\star^2 - 24c_\star^3) \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4. \end{aligned} \quad (4.2)$$

for all $\mathbf{q} \in \mathcal{R}_{2C_\star}$ with $C_\star \geq 10$ and $c_\star = 1/C_\star$, then any local minimum $\bar{\mathbf{q}}$ of $\psi(\mathbf{q})$ in \mathcal{R}_{2C_\star} satisfies $|\langle \bar{\mathbf{q}}, \mathcal{P}_\mathbb{S}[\mathbf{a}_l] \rangle| \geq 1 - 2c_\star \kappa^{-2}$ for some index l .

Proof Please refer to [Appendix B](#). ■

The Riemannian gradient and Hessian of the finite sample objective function $\psi(\mathbf{q})$ have similar expressions as those of the asymptotic objective function $\varphi(\mathbf{q})$. Let $\boldsymbol{\eta} = \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \in \mathbb{S}^{m-1}$. Then

$$\psi(\mathbf{q}) = -\frac{1}{4m} \left\| \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_4^4 = -\frac{1}{4m} \|\boldsymbol{\eta}\|_4^4, \quad (4.3)$$

we calculate the Euclidean gradient and Hessian of the objective function

$$\nabla \psi(\mathbf{q}) = -\frac{1}{m} (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{Y} \boldsymbol{\eta}^{\circ 3}, \quad (4.4)$$

$$\nabla^2 \psi(\mathbf{q}) = -\frac{3}{m} (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \text{diag}(\boldsymbol{\eta}^{\circ 2}) \mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2}. \quad (4.5)$$

Similarly, the Riemannian gradient and Hessian have the form

$$\text{grad}[\psi](\mathbf{q}) = \mathbf{P}_{\mathbf{q}^\perp} [\nabla \psi(\mathbf{q})] \quad (4.6)$$

$$= -\frac{1}{m} (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \boldsymbol{\eta}^{\circ 3} + \frac{1}{m} \mathbf{q} \|\boldsymbol{\eta}\|_4^4, \quad (4.7)$$

$$\text{Hess}[\psi](\mathbf{q}) = \mathbf{P}_{\mathbf{q}^\perp} [\nabla^2 \psi(\mathbf{q}) - \langle \nabla \psi(\mathbf{q}), \mathbf{q} \rangle \mathbf{I}] \mathbf{P}_{\mathbf{q}^\perp} \quad (4.8)$$

$$= \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{3}{m} (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \text{diag}(\boldsymbol{\eta}^{\circ 2}) \mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2} + \frac{1}{m} \|\boldsymbol{\eta}\|_4^4 \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp}. \quad (4.9)$$

Since $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$, we can see that the Riemannian gradient and Hessian are (complicated) functions of the random circulant matrix \mathbf{X}_0 . Although the entries of the vector \mathbf{x}_0 are probabilistically independent, the entries of \mathbf{X}_0 are dependent random variables. To remove the dependence within the random circulant matrix \mathbf{X}_0 , we break \mathbf{X}_0 into submatrices $\mathbf{X}_1, \dots, \mathbf{X}_{2k-1}$ that

$$\mathbf{X}_i = [\mathbf{x}_i, \mathbf{x}_{i+(2k-1)}, \dots, \mathbf{x}_{i+(m-2k-1)}]. \quad (4.10)$$

Each of which is (marginally) distributed as a $(2k-1) \times \frac{m}{2k-1}$ i.i.d. BG(θ) random matrix. Indeed, there exists a permutation $\boldsymbol{\Pi}$ such that

$$\mathbf{X}_0 \boldsymbol{\Pi} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{2k-1}]. \quad (4.11)$$

A detailed analysis of (4.7)-(4.9) (see [Appendix E](#) and [Appendix F](#) in the Appendix) allows us to control the finite sample fluctuations of the gradient and Hessian in terms of analogous quantities for each \mathbf{X}_i . Because the \mathbf{X}_i are i.i.d., they are amenable to standard tools from measure concentration. Taking a union bound over i , we show that the gradient (Lemma [Lemma 4.2](#)) and hessian (Lemma [Lemma 4.3](#)) concentrate as desired:

Lemma 4.2. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$. There exists positive constant C that whenever

$$m \geq C \frac{\min \left\{ (2C_\star \mu)^{-1}, \kappa^2 k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3(\kappa k), \quad (4.12)$$

and $\theta \geq 1/k$, then with probability no smaller than $1 - \exp(-k) - \theta^2 (1-\theta)^2 k^{-4} - 2 \exp(-\theta k) - 48k^{-7} - 48m^{-5} - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right)$,

$$\left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \leq c \frac{1-\theta}{\theta m^2} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}, \quad (4.13)$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}$ with $c \leq 3/(2C_\star) \leq \frac{3}{20}$.

Proof Please refer to section [E](#). ■

Lemma 4.3. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$. There exists positive constant C that whenever

$$m \geq C \frac{\min \left\{ (2C_\star \mu \kappa^2)^{-4/3}, k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^6 k^4 \log^3(\kappa k), \quad (4.14)$$

and $\theta \geq 1/k$, then with probability no smaller than $1 - \exp(-k) - \theta^2 (1-\theta)^2 k^{-4} - 2 \exp(-\theta k) - 48k^{-7} - 48m^{-5} - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right)$,

$$\left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \leq c \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4, \quad (4.15)$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}$ with positive constant $c \leq 0.048 \leq 3(1 - 6c_\star - 36c_\star^2 - 24c_\star^3)$.

Proof Please refer to section F. ■

5 Experiments

5.1 Properties of a Random Kernel

Our results are stated in terms of several parameters, including the condition number κ of \mathbf{A}_0 and the column coherence of \mathbf{A} . In Figure 3, we demonstrate the typical values of σ_0 , κ , and μ for generic unit-norm kernels of varying dimension $k = 10, 20, \dots, 1000$.

From this figure, for a generic unit-norm kernel, we have following estimates:

$$\sigma_0 \approx \log^{-1}(k), \quad (5.1)$$

$$\kappa \approx \log^{4/3}(k), \quad (5.2)$$

$$\mu \approx \sqrt{\log(k)/k}. \quad (5.3)$$

On the other hand, if the kernel \mathbf{a}_0 is bandpass, then both κ and μ are larger. In this situation, our results

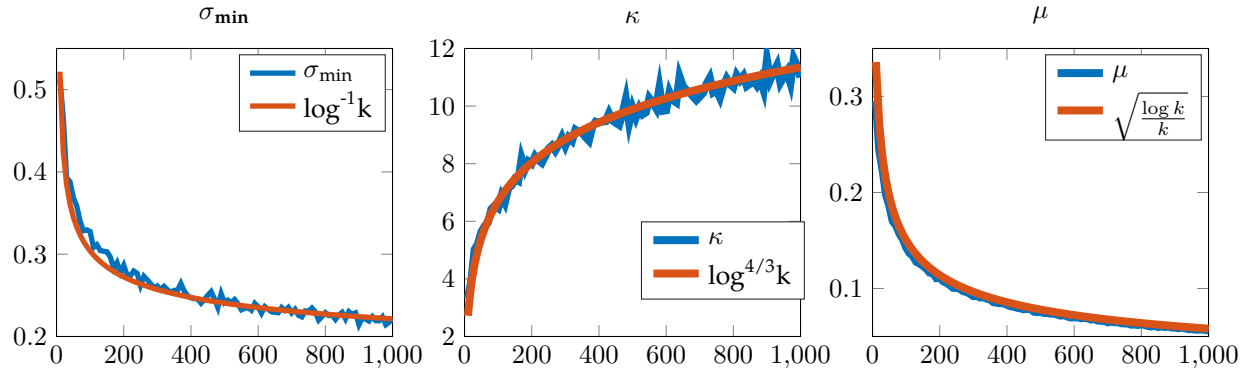


Figure 3: Average of Parameters σ_{\min} , κ , and μ of a random unit norm kernel \mathbf{a}_0 over 50 independent trials, as a function of dimension k .

require more observations m and smaller sparsity rate θ .

5.2 Recovery Accuracy of Local Minima

We next investigate the performance of Algorithm 1 under varying settings. We define the recover error as $\text{err} = 1 - \max_{\tau} |\langle \bar{\mathbf{a}}, \mathcal{P}_{\mathbb{S}}[\mathcal{L}_{k,s_{\tau}}^*[\bar{\mathbf{a}}_0]] \rangle|$, and calculate the average error from 50 independent experiments. In Figure 4, the left figure plots the average error when we fix the kernel size $k = 50$, and vary the dimension m and the sparsity θ of \mathbf{x}_0 .¹⁰ The right figure plots the average error when we vary the dimensions k, m of both convolution signals, and set the sparsity as $\theta = k^{-2/3}$.

This figure agrees with the theory developed in this paper: when the activation coefficient \mathbf{x}_0 is long and sparse (large m and small θ), the algorithm obtains a closer estimate of a shift-truncation of the ground truth.

5.3 Recovery Accuracy of the Ground Truth Kernel

In this section, we provide experiment results for the recovery of the ground truth kernel obtained by the annealing algorithm proposed in [ZLK⁺17]. The annealing algorithm recovers the ground truth kernel by minimizing the Lasso cost in (1.7), initialized at the zero-padded shift truncated kernel rendered from Algorithm

¹⁰Note that the x -axis is indexed with overlapping ratio $k \cdot \theta$, which indicates how many times the kernel \mathbf{a}_0 present in a k -length window of \mathbf{y} on average.

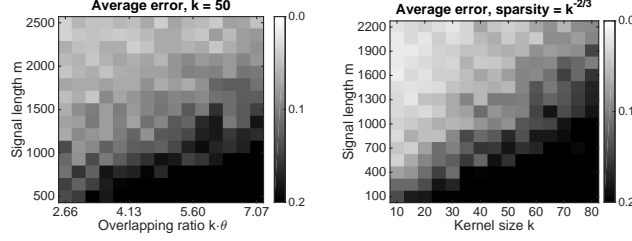


Figure 4: Recovery Error of the Shift Truncated Kernel of Algorithm 1.

Algorithm 1. The recovery accuracy presented in Figure Figure 6 is measured as $\text{err} = \min_{\tau} \|\bar{\mathbf{a}}^{(+)} \pm s_{\tau}[\widetilde{\mathbf{a}}_0]\|_2$. Here, $\bar{\mathbf{a}}^{(+)}$ denote the local minimum in the lifted optimization space.

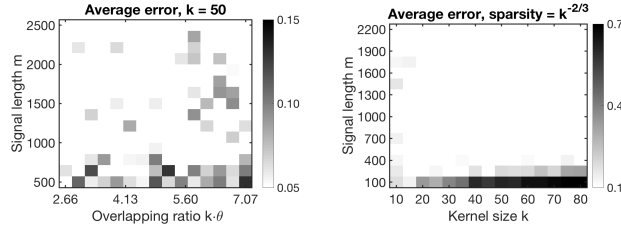


Figure 5: Recovery Error of the Ground Truth Kernel with Algorithm Algorithm 1 finding a shift truncated kernel and the annealing Lasso problem recovering the ground truth kernel.

For comparison, we also present experiment results of the algorithm proposed by [ZLK⁺17], which is composed of solving two Lasso minimization problems over the original kernel sphere and lifted kernel sphere respectively.

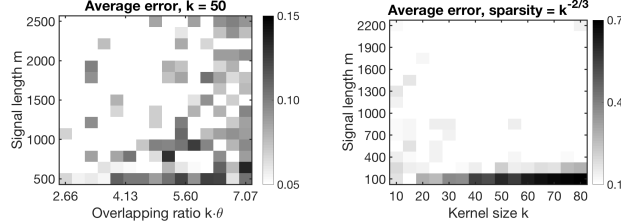


Figure 6: Recovery Error of the Ground Truth Kernel by minimizing the Lasso objective function recovering both the shift truncated kernel as well as the ground truth kernel.

In terms of the recovery accuracy of the ground truth kernel, Algorithm 1 proposed in this paper achieves better recovery for sparser and longer observations, while the [ZLK⁺17] manifests slight advantages when the observations is limited. As the optimization landscape studied in [ZLK⁺17] varies with different choice of sparsity parameter λ , it is possible that experiment results for [ZLK⁺17] could be improved. On the other hand, only empirical knowledge about the choice of λ is available while there is little disciplined understanding. In contrast, Algorithm 1 does not depend on any parameter tuning and guarantees recovery once the working conditions are met.

6 Discussions

Finally, we provide some comments about the results and proof strategy presented in this paper, and discuss directions for future research.

This paper casts the sparse blind deconvolution problem as finding a *spiky* vector in a subspace and studies its optimization landscape. We prove that the geometric property that *any local solution is close to a shift-truncation of the ground truth kernel* holds on a sub-level set of the sphere. This holds even when the observation contains densely overlapping copies of the true kernel. In addition, we propose a simple initialization scheme such that any descent algorithm that escapes strict saddles can recover the local minimum, which is a near shift-truncation of the ground truth kernel.

Sample Complexity. The sample complexity shown in this paper $m \sim k^6$ is suboptimal. Our proofs relies heavily on “worst case” tools such as the triangle inequality, multiplication of operator norm, and union bound. In particular, we believe that the sample complexity can be improved by replacing the sample splitting argument in Section [Appendix E](#) and [Appendix F](#) in the Appendix with more sophisticated arguments based on decoupling (see also [\[QZEW17\]](#)).

Global Geometry. The theoretical results presented in this paper demonstrate that “all local optima are benign” in the sub-level set \mathcal{R}_{C_*} . Our empirical results suggest that this is a property holds over the whole sphere. Proving this could be challenging, as our characterization of the saddle points only applies when $\|\zeta\|_4^4$ is large. It would be exciting to see if further research investigating other techniques for nonconvex optimization problems could be motivated by our current work.

Convolutional Dictionary Learning. This is a natural and practical extension of blind deconvolution, where the observation is the superposition of several convolutions. The empirical observations and algorithm proposed in [\[ZLK⁺17\]](#) hold in this more challenging situation. It would be interesting to develop efficient and provable algorithms for convolutional dictionary learning based on the ℓ^4 formulation.

Acknowledgement

The authors gratefully acknowledge support from NSF 1343282, NSF CCF 1527809, and NSF IIS 1546411. It is a great pleasure to acknowledge conversations with Yenson Lau, Sky Cheung, and Abhay Pasupathy.

References

- [ABG07] P.-A. Absil, C.G. Baker, and K.A. Gallivan. Trust-region methods on riemannian manifolds. *Foundations of Computational Mathematics*, 7(3):303–330, Jul 2007.
- [AMS07] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, Princeton, NJ, USA, 2007.
- [ARR12] Ali Ahmed, Benjamin Recht, and Justin Romberg. Blind deconvolution using convex programing. *arXiv preprint:1211.5608*, 2012.
- [Bha97] Rajendra Bhatia. *Matrix analysis*, 1997.
- [BVG13] Alexis Benichoux, Emmanuel Vincent, and Remi Gribonval. A fundamental pitfall in blind deconvolution with sparse and shift-invariant priors. *38th International Conference on Acoustics, Speech, and Signal Processing*, May 2013.
- [Chi16] Yuejie Chi. Guaranteed blind sparse spikes deconvolution via lifting and convex optimization. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):782–794, June 2016.
- [CLC⁺17] Sky Cheung, Yenson Lau, Zhengyu Chen, Ju Sun, Yuqian Zhang, John Wright, and Abhay Pasupathy. Beyond the fourier transform: A nonconvex optimization approach to microscopy analysis. *Submitted*, 2017.

- [CM14] Sunav Choudhary and Urbashi Mitra. Fundamental limits of blind deconvolution part I: Ambiguity kernel. *ArXiv e-prints*, abs/1411.3810, November 2014.
- [CM15] Sunav Choudhary and Urbashi Mitra. Fundamental limits of blind deconvolution part II: Sparsity-ambiguity trade-offs. *ArXiv e-prints*, abs/1503.03184, March 2015.
- [CW98] Tony F. Chan and Chiu-Kwong Wong. Total variation blind deconvolution. *IEEE Transactions on Image Processing*, 7(3):370–375, Mar 1998.
- [DIPG99] Victor De la Pena and Evarist Giné. *Decoupling: from dependence to independence*. Springer, 1999.
- [ETS11] Chaitanya Ekanadham, Daniel Tranchina, and Eero P. Simoncelli. A blind sparse deconvolution method for neural spike identification. In *Advances in Neural Information Processing Systems 24*, pages 1440–1448. 2011.
- [FR13] Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*. Springer, 2013.
- [GWY09] Donald Goldfarb, Zaiwen Wen, and Wotao Yin. A curvilinear search method for p-harmonic flows on spheres. *SIAM J. Imaging Sciences*, 2(1):84–109, 2009.
- [HSSS16] Samuel B. Hopkins, Tselil Schramm, Jonathan Shi, and David Steurer. Fast spectral algorithms from sum-of-squares proofs: Tensor decomposition and planted sparse vectors. In *Proceedings of the Forty-eighth Annual ACM Symposium on Theory of Computing*, STOC ’16, pages 178–191, 2016.
- [JGN⁺17] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. *arXiv preprint arXiv:1703.00887*, 2017.
- [KH96] Deepa Kundur and Dimitrios Hatzinakos. Blind image deconvolution. *Signal Processing Magazine, IEEE*, 13(3):43–64, May 1996.
- [Lew98] Michael S. Lewicki. A review of methods for spike sorting: the detection and classification of neural action potentials. *Network: Computation in Neural Systems*, 9(4):53–78, 1998.
- [LLB16] Yanjun Li, Kiryung Lee, , and Yoram Bresler. Identifiability in blind deconvolution with subspace or sparsity constraints. *IEEE Transaction of Information Theory*, 62(7):4266 – 4275, July 2016.
- [LLB17] Yanjun Li, Kiryung Lee, and Yoram Bresler. Identifiability and stability in blind deconvolution under minimal assumptions. *IEEE Transaction of Information Theory*, 2017.
- [LLJB17] Kiryung Lee, Yanjun Li, Marius Junge, and Yoram Bresler. Blind recovery of sparse signals from subsampled convolution. *IEEE Transaction of Information Theory*, 63(2):802–821, February 2017.
- [LLSW16] Xiaodong Li, Shuyang Ling, Thomas Strohmer, and Ke Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. preprint, 2016.
- [LS15] Shuyang Ling and Thomas Strohmer. Self-calibration and biconvex compressive sensing. *Inverse Problems*, 31(11):115002, 2015.
- [LS17] Shuyang Ling and Thomas Strohmer. Blind deconvolution meets blind demixing: Algorithms and performance bounds. *IEEE Transactions on Information Theory*, 63(7):4497–4520, 2017.
- [LWDF11] Anat Levin, Yair Weiss, Fredo Durand, and William T. Freeman. Understanding blind deconvolution algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(12):2354–2367, Dec 2011.
- [OJF⁺15] S. Oymak, A. Jalali, M. Fazel, Y. C. Eldar, and B. Hassibi. Simultaneously structured models with application to sparse and low-rank matrices. *IEEE Transactions on Information Theory*, 61(5):2886–2908, May 2015.
- [PF14] Daniele Perrone and Paolo Favaro. Total variation blind deconvolution: The devil is in the details. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [QSW16] Qing Qu, Ju Sun, and John Wright. Finding a sparse vector in a subspace: linear sparsity using alternating directions. *IEEE Transactions on Information Theory*, 2016.
- [QZEW17] Qing Qu, Yuqian Zhang, Yonina C. Eldar, and John Wright. Convolutional phase retrieval via gradient descent. preprint, 2017.
- [SQW15] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere. preprint, 2015.
- [SWW12] Daniel Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. preprint, 2012.
- [Tro12] Joel A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012.

- [WC16] L. Wang and Y. Chi. Blind deconvolution from multiple sparse inputs. *IEEE Signal Processing Letters*, 23(10):1384–1388, Oct 2016.
- [WZ13] David Wipf and Haichao Zhang. Revisiting bayesian blind deconvolution. *arXiv preprint:1305.2362*, 2013.
- [XRKM17] Peng Xu, Farbod Roosta-Khorasani, and Michael W. Mahoney. Second-order optimization for non-convex machine learning: An empirical study. *arXiv preprint arXiv:1708.07827*, 2017.
- [ZLK⁺17] Yuqian Zhang, Yenson Lau, Han-Wen Kuo, Sky Cheung, Abhay Pasupathy, and John Wright. On the global geometry of sphere-constrained sparse blind deconvolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [ZWZ13] Haichao Zhang, David Wipf, and Yanning Zhang. Multi-image blind deblurring using a coupled adaptive sparse prior. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, January 2013.

Appendix

[Appendix A](#) contains some basic lemmas for quantities used repeatedly; [Appendix B](#) presents the proofs of the main theorem and corollary of this paper. [Appendix C](#) and [Appendix D](#) proves for lemmas around the initialization point \mathbf{q}_{init} and the preconditioning term $\mathbf{Y}^T \mathbf{Y}$ (or $\mathbf{A}_0^T \mathbf{A}_0$) respectively. Finite sample concentration for the Riemannian gradient and Hessian are presented in [Appendix E](#) and [Appendix F](#) respectively.

A Basics

Lemma A.1 (Expectation of the Approximate Objective Function). *Assuming $\mathbf{x}_0 \sim \text{i.i.d.}$ $\text{BG}(\theta) \in \mathbb{R}^m$, then*

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_0} \left[\frac{1}{m} \left\| \mathbf{Y}^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_4^4 \right] \\ = 3\theta (1 - \theta) \left\| \mathbf{A}^T \mathbf{q} \right\|_4^4 + 3\theta^2 \left\| \mathbf{A}^T \mathbf{q} \right\|_2^4. \end{aligned} \quad (\text{A.1})$$

Proof Let $\mathbf{g} \in \mathbb{R}^{2k-1}$ be a standard random Gaussian vector and \mathbf{P}_I be the projection operator onto Bernoulli vector $I \sim \text{Ber}(\theta)$. Then any column $\mathbf{x}_i \in \mathbb{R}^{2k-1}$ of \mathbf{X}_0 is equal in distribution to $\mathbf{x}_i = \mathbf{P}_I \mathbf{g}$ with $\mathbf{g} \sim \text{i.i.d.}$ $\mathcal{N}(0, 1)$.

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_0} \left[\frac{1}{m} \left\| \mathbf{Y}^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_4^4 \right] \\ = \frac{1}{m} \mathbb{E}_I \mathbb{E}_{\mathbf{g}} \left\| \mathbf{q}^T \mathbf{A} \mathbf{X}_0 \right\|_4^4 \end{aligned} \quad (\text{A.2})$$

$$= \mathbb{E}_I \mathbb{E}_{\mathbf{g}} \left\| \mathbf{q}^T \mathbf{A} \mathbf{x}_i \right\|_4^4 \quad (\text{A.3})$$

$$= \mathbb{E}_I \mathbb{E}_{\mathbf{g}} (\mathbf{q}^T \mathbf{A} \mathbf{P}_I \mathbf{g})^4 \quad (\text{A.4})$$

$$= 3 \mathbb{E}_I (\mathbf{q}^T \mathbf{A} \mathbf{P}_I \mathbf{A}^T \mathbf{q})^2 \quad (\text{A.5})$$

$$= 3 \mathbb{E}_I \left(\sum_{i \in I} \langle \mathbf{a}_i, \mathbf{q} \rangle^4 + \sum_{\{i \neq j\} \in I} \langle \mathbf{a}_i, \mathbf{q} \rangle^2 \langle \mathbf{a}_j, \mathbf{q} \rangle^2 \right) \quad (\text{A.6})$$

$$= 3\theta (1 - \theta) \left\| \mathbf{A}^T \mathbf{q} \right\|_4^4 + 3\theta^2 \left\| \mathbf{A}^T \mathbf{q} \right\|_2^4 \quad (\text{A.7})$$

■

Lemma A.2 (Root Estimation for Cubic Gradient Function). *Consider an equation of the form*

$$f(x) = x(\alpha - x^2) - \beta = 0, \quad (\text{A.8})$$

with $\alpha > 0$. Suppose that $\beta < \frac{1}{4}\alpha^{3/2}$. Then $f(x) = 0$ has three solutions, x_1, x_2, x_3 satisfying

$$\max \{ |x_1 - \sqrt{\alpha}|, |x_2 + \sqrt{\alpha}|, |x_3| \} \leq \frac{2\beta}{\alpha}. \quad (\text{A.9})$$

Proof Suppose first that $\beta > 0$. Then $f(0) < 0$. Moreover,

$$f\left(\frac{2\beta}{\alpha}\right) = 2\beta - 8\beta^3/\alpha^3 - \beta \quad (\text{A.10})$$

$$= \beta(1 - 8\beta^2/\alpha^3) \quad (\text{A.11})$$

$$> 0. \quad (\text{A.12})$$

Hence, f has at least one root in the interval $\left[0, \frac{2\beta}{\alpha}\right]$. Similarly, notice that $f(\sqrt{\alpha}) < 0$ and that

$$\begin{aligned} f\left(\sqrt{\alpha} - \frac{2\beta}{\alpha}\right) &= \alpha^{3/2} - 2\beta - \left(\sqrt{\alpha} - \frac{2\beta}{\alpha}\right)^3 - \beta \\ &= \alpha^{3/2} - 3\beta - \alpha^{3/2} + 6\beta - 12\beta^2/\alpha^{3/2} + 8\beta^3/\alpha^3 \end{aligned} \quad (\text{A.13})$$

$$= \alpha^{3/2} - 3\beta - \alpha^{3/2} + 6\beta - 12\beta^2/\alpha^{3/2} + 8\beta^3/\alpha^3 \quad (\text{A.14})$$

$$= \beta \left(3 - \frac{12\beta}{\alpha^{3/2}} + \frac{8\beta^2}{\alpha^3}\right) \quad (\text{A.15})$$

$$> 0. \quad (\text{A.16})$$

Thus, there is at least one root in the interval $\left[\sqrt{\alpha} - \frac{2\beta}{\alpha}, \sqrt{\alpha}\right]$. Finally, note that $f(-\sqrt{\alpha}) < 0$, $\frac{df}{dx}(-\sqrt{\alpha}) = -2\alpha$, and $\frac{d^2f}{dx^2}(x') = -3x'$ is positive for $x' \leq -\sqrt{\alpha}$. Hence, convexity gives that

$$\begin{aligned} f\left(-\sqrt{\alpha} - \frac{2\beta}{\alpha}\right) &\geq f(-\sqrt{\alpha}) + \frac{df}{dx}(-\sqrt{\alpha}) \times (-2\beta/\alpha) \\ &= -\beta + (-2\alpha) \times (-2\beta/\alpha) \end{aligned} \quad (\text{A.17})$$

$$= -\beta + (-2\alpha) \times (-2\beta/\alpha) \quad (\text{A.18})$$

$$= 3\beta > 0. \quad (\text{A.19})$$

Under this condition, there is at least one root in the interval, $[-\sqrt{\alpha} - 2\beta/\alpha, -\sqrt{\alpha}]$. These three intervals do not overlap, as long as $\frac{4\beta}{\alpha} < \sqrt{\alpha}$, or $\beta < \frac{1}{4}\alpha^{3/2}$.

In the case that $\beta \leq 0$, a symmetric argument applies. Thus there are exactly three solutions to equation (A.8) in the specified intervals. ■

Lemma A.3. Let \mathbf{a}_l and $\mathbf{a}_{l'}$ be two nonzero vectors with inner product $\mu_{l,l'} \doteq \langle \mathbf{a}_l, \mathbf{a}_{l'} \rangle$. Then for any unit vector $\mathbf{v} \in \text{span}(\mathbf{a}_l, \mathbf{a}_{l'})$,

$$\left| \left\langle \frac{\mathbf{a}_l}{\|\mathbf{a}_l\|_2}, \mathbf{v} \right\rangle \right|^2 + \left| \left\langle \frac{\mathbf{a}_{l'}}{\|\mathbf{a}_{l'}\|_2}, \mathbf{v} \right\rangle \right|^2 \geq 1 - \frac{|\mu_{l,l'}|}{\|\mathbf{a}_l\|_2 \|\mathbf{a}_{l'}\|_2}. \quad (\text{A.20})$$

Proof Let \mathbf{u} and \mathbf{u}^\perp be two orthogonal unit vectors, such that

$$\mathbf{a}_l = \|\mathbf{a}_l\|_2 \mathbf{u}, \quad (\text{A.21})$$

$$\mathbf{a}_{l'} = \frac{\mu_{l,l'}}{\|\mathbf{a}_l\|_2} \mathbf{u} + \sqrt{\|\mathbf{a}_{l'}\|_2^2 - \frac{\mu_{l,l'}^2}{\|\mathbf{a}_l\|_2^2}} \mathbf{u}^\perp. \quad (\text{A.22})$$

Suppose $\mathbf{v} = a\mathbf{u} + b\mathbf{u}^\perp$ with $a^2 + b^2 = 1$. Let $\mu_{\text{rel}} = \frac{\mu_{l,l'}}{\|\mathbf{a}_l\|_2 \|\mathbf{a}_{l'}\|_2}$, then we can expand the quantity of interests as

$$\begin{aligned} &\left| \left\langle \frac{\mathbf{a}_l}{\|\mathbf{a}_l\|_2}, \mathbf{v} \right\rangle \right|^2 + \left| \left\langle \frac{\mathbf{a}_{l'}}{\|\mathbf{a}_{l'}\|_2}, \mathbf{v} \right\rangle \right|^2 \\ &= \left| \langle \mathbf{u}, a\mathbf{u} + b\mathbf{u}^\perp \rangle \right|^2 \\ &\quad + \left| \left\langle \mu_{\text{rel}}\mathbf{u} + \sqrt{1 - \mu_{\text{rel}}^2} \mathbf{u}^\perp, a\mathbf{u} + b\mathbf{u}^\perp \right\rangle \right|^2 \end{aligned} \quad (\text{A.23})$$

$$= a^2 + \left(a\mu_{\text{rel}} + b\sqrt{1 - \mu_{\text{rel}}^2} \right)^2 \quad (\text{A.24})$$

$$= a^2 + b^2 + (a^2 - b^2)\mu_{\text{rel}}^2 + 2ab\mu_{\text{rel}}\sqrt{1 - \mu_{\text{rel}}^2} \quad (\text{A.25})$$

$$= 1 + [a^2 - b^2, 2ab] \left[\mu_{\text{rel}}^2, \mu_{\text{rel}} \sqrt{1 - \mu_{\text{rel}}^2} \right]^T \quad (\text{A.26})$$

Since $[a^2 - b^2, 2ab]$ is a unit vector, then above equation is lower bounded by

$$\begin{aligned} 1 - \left\| \left[\mu_{\text{rel}}^2, \mu_{\text{rel}} \sqrt{1 - \mu_{\text{rel}}^2} \right] \right\|_2 \\ = 1 - |\mu_{\text{rel}}| \end{aligned} \quad (\text{A.27})$$

$$= 1 - \frac{|\mu_{l,l'}|}{\|\mathbf{a}_l\|_2 \|\mathbf{a}_{l'}\|_2} \quad (\text{A.28})$$

as claimed. \blacksquare

Lemma A.4 (Nonzeros in a Bernoulli Vector). *Let $\mathbf{v} \sim_{\text{i.i.d.}} \text{Ber}(\theta) \in \mathbb{R}^n$, then*

$$\mathbb{P}[\|\mathbf{v}\|_0 \geq (1+t)\theta n] \leq 2 \exp\left(-\frac{3t^2}{2t+6}\theta n\right). \quad (\text{A.29})$$

Proof As $\|\mathbf{v}\|_0 = v_0 + \dots + v_{n-1}$, and

$$|v_i - \theta| \leq 1, \quad \mathbb{E}[(v_i - \theta)^2] = \theta(1 - \theta) \leq \theta \quad (\text{A.30})$$

with Bernstein's inequality, we obtain that

$$\begin{aligned} \mathbb{P}[\|\mathbf{v}\|_0 \geq (1+t)\theta n] \\ \leq 2 \exp\left(-\frac{t^2 \theta^2 n^2}{2(\theta - \theta^2)n + \frac{2}{3}t\theta n}\right) \end{aligned} \quad (\text{A.31})$$

$$\leq 2 \exp\left(-\frac{3t^2}{2t+6}\theta n\right), \quad (\text{A.32})$$

as claimed. \blacksquare

Lemma A.5 (Entry-wise Truncation of a Bernoulli Gaussian Vector). *Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$, then*

$$\mathbb{P}[\|\mathbf{x}_0\|_\infty > t] \leq 2\theta m e^{-t^2/2}. \quad (\text{A.33})$$

Proof A Bernoulli-Gaussian variable $x = \omega \cdot g$ satisfies

$$\mathbb{P}[|x| \geq t] = \theta \cdot \mathbb{P}[|g| \geq t] \leq 2\theta e^{-t^2/2}, \quad (\text{A.34})$$

Taking a union bound over the m entries of \mathbf{x}_0 , we obtain

$$\mathbb{P}[\|\mathbf{x}_0\|_\infty > t] \leq m \mathbb{P}[|x| > t] \quad (\text{A.35})$$

$$\leq 2\theta m e^{-t^2/2}, \quad (\text{A.36})$$

as claimed. \blacksquare

Lemma A.6 (Operator Norm of a Bernoulli Gaussian Circulant Matrix). *Let $\mathbf{C}_{\mathbf{x}_0} \in \mathbb{R}^{m \times m}$ be the circulant matrix generated from $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$, then*

$$\mathbb{P}[\|\mathbf{C}_{\mathbf{x}_0}\|_2 \geq t] \leq 2m \exp\left(-\frac{t^2}{2\theta m + 2t}\right). \quad (\text{A.37})$$

Proof The operator norm of a circulant matrix is

$$\|\mathbf{C}_{\mathbf{x}_0}\|_2 = \max_l |\langle \mathbf{x}_0, \mathbf{w}_l \rangle|, \quad (\text{A.38})$$

where \mathbf{w}_l is the l -th (discrete) Fourier basis vector

$$\mathbf{w}_l = \left[1, e^{l \frac{2\pi j}{m}}, \dots, e^{l(m-1) \frac{2\pi j}{m}} \right]^T, \quad l = 0, \dots, m-1, \quad (\text{A.39})$$

and j is the imaginary unit. With moment control Bernstein inequality, we obtain

$$\begin{aligned} \mathbb{P} [|\langle \mathbf{x}_0, \mathbf{w}_l \rangle| \geq t] &\leq 2 \exp \left(-\frac{t^2}{2\theta \|\mathbf{w}_l\|_2^2 + 2 \|\mathbf{w}_l\|_\infty t} \right) \\ &\leq 2 \exp \left(-\frac{t^2}{2\theta m + 2t} \right) \end{aligned} \quad (\text{A.40})$$

together with the union bound,

$$\mathbb{P} [\|\mathbf{C}_{\mathbf{x}_0}\|_2 \geq t] \leq m \mathbb{P} [|\langle \mathbf{x}_0, \mathbf{w}_l \rangle| \geq t] \quad (\text{A.41})$$

$$\leq 2m \exp \left(-\frac{t^2}{2\theta m + 2t} \right), \quad (\text{A.42})$$

as claimed. \blacksquare

Lemma A.7 (Norms of $\boldsymbol{\eta}$ and $\bar{\boldsymbol{\eta}}$). Suppose $\delta = \left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \leq 1/(2\kappa^2)$, then vectors $\boldsymbol{\eta} = \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q}$ and $\bar{\boldsymbol{\eta}} = \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q}$ satisfy

$$\|\boldsymbol{\eta}\|_\infty \leq \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right) \left(\frac{2k}{\theta m} \right)^{1/2} \|\mathbf{x}_0\|_\infty, \quad (\text{A.43})$$

$$\|\bar{\boldsymbol{\eta}}\|_\infty \leq \left(\frac{2k}{\theta m} \right)^{1/2} \|\mathbf{x}_0\|_\infty, \quad (\text{A.44})$$

$$\|\boldsymbol{\eta}\|_6^6 \leq \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right)^4 \frac{4k^2}{\theta^2 m^2} \|\mathbf{x}_0\|_\infty^4, \quad (\text{A.45})$$

$$\|\bar{\boldsymbol{\eta}}\|_2 \leq 1 + \delta/2, \quad (\text{A.46})$$

$$\|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_\infty \leq \frac{4\kappa^3 \delta}{\sigma_{\min}} \left(\frac{2k}{\theta m} \right)^{1/2} \|\mathbf{x}_0\|_\infty, \quad (\text{A.47})$$

$$\|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \leq (1 + \delta/2) \frac{4\kappa^3 \delta}{\sigma_{\min}}. \quad (\text{A.48})$$

Proof Since $\delta = \left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2$, then

$$\|\mathbf{X}_0\|_2 \leq (\theta m)^{1/2} \sqrt{1 + \delta} \quad (\text{A.49})$$

$$\leq (\theta m)^{1/2} (1 + \delta/2). \quad (\text{A.50})$$

As $\boldsymbol{\eta} = \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} = \mathbf{X}_0^T \mathbf{A}_0^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q}$, together with [Lemma D.3](#):

$$\begin{aligned} &\left\| \mathbf{A}_0^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_\infty \\ &\leq \left\| \mathbf{A}_0^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_2 \end{aligned} \quad (\text{A.51})$$

$$\begin{aligned} &\leq \left\| \mathbf{A}_0^T \left((\mathbf{Y}\mathbf{Y}^T)^{-1/2} - (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right) \mathbf{q} \right\|_2 \\ &\quad + \left\| \mathbf{A}_0^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_2 \end{aligned} \quad (\text{A.52})$$

$$\leq (\theta m)^{-1/2} \frac{4\kappa^3 \delta}{\sigma_{\min}} \|\mathbf{q}\|_2 + (\theta m)^{-1/2} \|\mathbf{A}^T \mathbf{q}\|_2 \quad (\text{A.53})$$

$$\leq (\theta m)^{-1/2} \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right) \quad (\text{A.54})$$

Norms of $\boldsymbol{\eta}$. Since $\|\mathbf{X}_0 \mathbf{e}_l\|_2 \leq \sqrt{2k-1} \|\mathbf{X}_0 \mathbf{e}_l\|_\infty$, we have

$$\|\boldsymbol{\eta}\|_\infty = \max_{l \in [1, \dots, m]} \left\langle \mathbf{X}_0 \mathbf{e}_l, \mathbf{A}_0^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{q} \right\rangle \quad (\text{A.55})$$

$$\leq \max_l \|\mathbf{X}_0 \mathbf{e}_l\|_2 \left\| \mathbf{A}_0^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_2 \quad (\text{A.56})$$

$$\leq \sqrt{2k} \|\mathbf{x}_0\|_\infty \cdot (\theta m)^{-1/2} \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right). \quad (\text{A.57})$$

At the same time, plugging in $\|\boldsymbol{\eta}\|_2 = 1$, we have

$$\|\boldsymbol{\eta}\|_6^6 \leq \|\boldsymbol{\eta}\|_2^2 \|\boldsymbol{\eta}\|_\infty^4 \leq \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right)^4 \frac{4k^2}{\theta^2 m^2} \|\mathbf{x}_0\|_\infty^4. \quad (\text{A.58})$$

Norms of $\bar{\boldsymbol{\eta}}$. Here, $\bar{\boldsymbol{\eta}} = \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} = \mathbf{X}_0^T \mathbf{A}_0^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q}$ with

$$\begin{aligned} &\left\| \mathbf{A}_0^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_\infty \\ &\leq \left\| \mathbf{A}_0^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_2 \end{aligned} \quad (\text{A.59})$$

$$= (\theta m)^{-1/2}, \quad (\text{A.60})$$

therefore

$$\begin{aligned} \|\bar{\boldsymbol{\eta}}\|_\infty &\leq \max_l \|\mathbf{X}_0 \mathbf{e}_l\|_2 \left\| \mathbf{A}_0 (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_2 \\ &\leq \left(\frac{2k}{\theta m} \right)^{1/2} \|\mathbf{x}_0\|_\infty, \end{aligned} \quad (\text{A.61})$$

$$\begin{aligned} \|\bar{\boldsymbol{\eta}}\|_2 &\leq \|\mathbf{X}_0^T\|_2 \left\| \mathbf{A}_0 (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_2 \\ &\leq 1 + \delta/2. \end{aligned} \quad (\text{A.62})$$

Norms of $\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}$. With similar reasoning, we can obtain

$$\begin{aligned} &\|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_\infty \\ &= \left\| \mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{q} - \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right\|_\infty \\ &\leq \max_{l \in [1, \dots, m]} \|\mathbf{X}_0 \mathbf{e}_l\|_2 (\theta m)^{-1/2} \times \\ &\quad \left\| \mathbf{A}_0^T \left(\frac{1}{\theta m} \mathbf{Y}\mathbf{Y}^T \right)^{-1/2} - \mathbf{A}_0^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right\|_2 \end{aligned} \quad (\text{A.63})$$

$$\leq \frac{4\kappa^3 \delta}{\sigma_{\min}} \left(\frac{2k}{\theta m} \right)^{1/2} \|\mathbf{x}_0\|_\infty, \quad (\text{A.64})$$

and

$$\begin{aligned} & \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \\ & \leq \|\mathbf{X}_0\|_2 (\theta m)^{-1/2} \|\mathbf{q}\|_2 \times \\ & \quad \left\| \mathbf{A}_0^T \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} - \mathbf{A}_0^T (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right\|_2 \end{aligned} \quad (\text{A.65})$$

$$\leq (\theta m)^{-1/2} \frac{4\kappa^3 \delta}{\sigma_{\min}} \|\mathbf{X}_0\|_2 \quad (\text{A.66})$$

$$\leq (1 + \delta/2) \frac{4\kappa^3 \delta}{\sigma_{\min}}, \quad (\text{A.67})$$

completing the proof. \blacksquare

B Proof of the Main Theorem and Corollary

B.1 Proof of the Main Theorem

Lemma B.1. *If following inequalities hold*

$$\begin{aligned} & \left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq \frac{3c_\star}{2\kappa^2} \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^6, \end{aligned} \quad (\text{B.1})$$

$$\begin{aligned} & \left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq 3(1 - 6c_\star - 36c_\star^2 - 24c_\star^3) \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4. \end{aligned} \quad (\text{B.2})$$

for all $\mathbf{q} \in \mathcal{R}_{2C_\star}$ with $C_\star \geq 10$ and $c_\star = 1/C_\star$, then any local minimum $\bar{\mathbf{q}}$ of $\psi(\mathbf{q})$ in \mathcal{R}_{2C_\star} satisfies $|\langle \bar{\mathbf{q}}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_l] \rangle| \geq 1 - 2c_\star \kappa^{-2}$ for some index l .

Proof Let

$$\boldsymbol{\delta}_{\text{grad}} = \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}), \quad (\text{B.3})$$

and let

$$\bar{\boldsymbol{\delta}}_{\text{grad}} = \frac{\theta m^2}{3(1-\theta)} \boldsymbol{\delta}_{\text{grad}}. \quad (\text{B.4})$$

Then at any stationary point of $\psi(\mathbf{q})$, we have

$$\mathbf{0} = \mathbf{A}^T \text{grad}[\psi](\mathbf{q}) \quad (\text{B.5})$$

$$= \frac{3(1-\theta)}{\theta m^2} \mathbf{A}^T \text{grad}[\varphi](\mathbf{q}) + \mathbf{A}^T \boldsymbol{\delta}_{\text{grad}}. \quad (\text{B.6})$$

Hence for any index i , following equality always holds

$$\begin{aligned} 0 &= \|\mathbf{a}_i\|_2^2 \zeta_i^3 + \sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3 - \zeta_i \|\boldsymbol{\zeta}\|_4^4 + \langle \mathbf{a}_i, \bar{\boldsymbol{\delta}}_{\text{grad}} \rangle \\ &= \underbrace{\zeta_i^3 - \zeta_i \frac{\|\boldsymbol{\zeta}\|_4^4}{\|\mathbf{a}_i\|_2^2}}_{\alpha_i} + \underbrace{\sum_{j \neq i} \langle \mathbf{a}_i, \mathbf{a}_j \rangle \zeta_j^3 + \langle \mathbf{a}_i, \bar{\boldsymbol{\delta}}_{\text{grad}} \rangle}_{\beta'_i} \end{aligned} \quad (\text{B.7})$$

with $\zeta = \mathbf{A}^T \mathbf{q}$. Under the assumption that

$$\left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \leq \frac{3c_\star}{2\kappa^2} \frac{1-\theta}{\theta m^2} \|\zeta\|_4^6, \quad (\text{B.8})$$

the perturbed part can be bounded via

$$|\langle \mathbf{a}_i, \bar{\delta}_{\text{grad}} \rangle| \leq \|\mathbf{a}_i\|_2 \|\bar{\delta}_{\text{grad}}\|_2 \leq \frac{c_\star}{2\kappa^2} \|\mathbf{a}_i\|_2 \|\zeta\|_4^6, \quad (\text{B.9})$$

and also

$$\frac{\beta'_i}{\alpha_i^{3/2}} \leq \frac{\mu \|\zeta\|_3^3 + \frac{1}{2} c_\star \kappa^{-2} \|\zeta\|_4^6}{\|\zeta\|_4^6} \leq c_\star \kappa^{-2} \leq \frac{1}{4}. \quad (\text{B.10})$$

Then by [Lemma A.2](#), at every stationary point $\bar{\mathbf{q}}$, the i -th entry of ζ resides in the set $\bigcup_{x \in \{0, \pm\sqrt{\alpha_i}\}} [x - \frac{2\beta'_i}{\alpha_i}, x + \frac{2\beta'_i}{\alpha_i}]$ – i.e., ζ is nearly a trinary vector.

Moreover, we can characterize the curvature of critical points in terms of the number of large entries of ζ . Indeed, whenever ζ has at least two entries in

$$\bigcup_{x \in \{\pm\sqrt{\alpha_i}\}} \left[x - \frac{2\beta'_i}{\alpha_i}, x + \frac{2\beta'_i}{\alpha_i} \right],$$

using (3.52), there exists a direction of strict negative curvature, provided

$$\begin{aligned} \text{Hess}[\psi](\mathbf{q}) &\prec \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \\ &\quad + 3(2 - 11c_\star) \frac{1-\theta}{\theta m^2} \|\zeta\|_4^4 \mathbf{I}. \end{aligned} \quad (\text{B.11})$$

Similarly, whenever ζ has only one entry in

$$\bigcup_{x \in \{\pm\sqrt{\alpha_i}\}} \left[x - \frac{2\beta'_i}{\alpha_i}, x + \frac{2\beta'_i}{\alpha_i} \right],$$

using (3.35), we have that $\text{Hess}[\psi](\mathbf{q}) \succ \mathbf{0}$, provided

$$\begin{aligned} \text{Hess}[\psi](\mathbf{q}) &\succ \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \\ &\quad - 3(1 - 6c_\star - 36c_\star^2 - 24c_\star^3) \frac{1-\theta}{\theta m^2} \|\zeta\|_4^4 \mathbf{I}. \end{aligned} \quad (\text{B.12})$$

When $C_\star \geq 10$ and $c_\star \leq 0.1$, we have $2 - 11c_\star > 1 - 6c_\star - 36c_\star^2 - 24c_\star^3 \geq 0.016$, and so above characterization obtains. \blacksquare

Theorem B.2 (Main Result). Assume the observation $\mathbf{y} \in \mathbb{R}^m$ is the cyclic convolution of $\mathbf{a}_0 \in \mathbb{R}^k$ and $\mathbf{x}_0 \sim \text{i.i.d. BG}(\theta) \in \mathbb{R}^m$, where the convolution matrix $\mathbf{A}_0 \in \mathbb{R}^{k \times (2k-1)}$ has minimum singular value $\sigma_{\min} > 0$ and condition number $\kappa \geq 1$, and \mathbf{A} has column incoherence μ . If

$$m \geq C \frac{\min \left\{ (2C_\star \mu)^{-1}, \kappa^2 k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3(\kappa k) \quad (\text{B.13})$$

and $\theta \geq \log k/k$, then with probability no smaller than $1 - \exp(-k) - \theta^2(1-\theta)^2 k^{-4} - 2 \exp(-\theta k) - 48k^{-7} - 48m^{-5} - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right)$, any local minimum $\bar{\mathbf{q}}$ of ψ in $\hat{\mathcal{R}}_{2C_\star}$ satisfies $|\langle \bar{\mathbf{q}}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_\tau] \rangle| \geq 1 - c_\star \kappa^{-2}$ for some integer τ .

Proof From the concentration analysis for the Riemannian gradient (Lemma 4.2) and Hessian (Lemma 4.3), if

$$m \geq C \frac{\min \left\{ (2C_*\mu)^{-1}, \kappa^2 k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3(\kappa k), \quad (\text{B.14})$$

then with probability no smaller than $1 - \exp(-k) - \theta^2(1-\theta)^2 k^{-4} - 2 \exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right) - 48k^{-7} - 48m^{-5}$,

$$\begin{aligned} & \left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq \frac{3c_*}{2\kappa^2} \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^6, \end{aligned} \quad (\text{B.15})$$

$$\begin{aligned} & \left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \\ & \leq 3(1-6c_* - 36c_*^2 - 24c_*^3) \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4. \end{aligned} \quad (\text{B.16})$$

hold for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$ with $C_* \geq 10$ and $c_* = 1/C_*$. Therefore, by Lemma 4.1 any local minimum $\bar{\mathbf{q}}$ of $\psi(\mathbf{q})$ in \mathcal{R}_{2C_*} satisfies $|\langle \bar{\mathbf{q}}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_l] \rangle| \geq 1 - 2c_*\kappa^{-2}$ for some index l . ■

B.2 Proof of the Main Corollary

Corollary B.3. Suppose the ground truth kernel \mathbf{a}_0 has induces coherence $0 \leq \mu \leq \frac{1}{8 \times 48} \log^{-3/2}(k)$ and sparse coefficient $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$. there exist positive constants $C \geq 2560^4$ and C' such that whenever the sparsity level

$$\begin{aligned} 64k^{-1} \log k \leq \theta \leq \min \left\{ \frac{1}{48^2} \mu^{-2} k^{-1} \log^{-2} k, \right. \\ \left. \left(\frac{1}{4} - \frac{640}{C^{1/4}} \right) (3C_*\mu\kappa^2)^{-2/3} k^{-1} (1 + 36\mu^2 k \log k)^{-2} \right\}, \end{aligned} \quad (\text{B.17})$$

and signal length

$$\begin{aligned} m \geq \max \left\{ C\theta^2 \sigma_{\min}^{-2} \kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log(\kappa k), \right. \\ \left. C' (1-\theta)^{-2} \sigma_{\min}^{-2} \min \left\{ \mu^{-1}, \kappa^2 k^2 \right\} \kappa^8 k^4 \log^3(\kappa k) \right\}, \end{aligned} \quad (\text{B.18})$$

then Algorithm 1 recovers $\bar{\mathbf{a}}$ such that

$$\|\bar{\mathbf{a}} \pm \mathcal{P}_{\mathbb{S}}[\mathbf{u}_k s_\tau [\widetilde{\mathbf{a}}_0]]\|_2 \leq 4\sqrt{c_*} + ck^{-1} \quad (\text{B.19})$$

for some integer shift $\tau \in [-(k-1), k-1]$ with probability no smaller than $1 - k^{-1} - 8k^{-2} - \exp(-k) - \theta^2(1-\theta)^2 k^{-4} - 2 \exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right) - 48k^{-7} - 48m^{-5}$.

Proof From the concentration results for the Riemannian gradient, at every point $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$, the objective value of $\psi(\mathbf{q})$ satisfies

$$\begin{aligned} & \left| \psi(\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \varphi(\mathbf{q}) + \frac{3}{4m^2} \right| \\ & \leq \left| \frac{\left\| \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q} \right\|_4^4}{4m} - \frac{3(1-\theta) \|\boldsymbol{\zeta}\|_4^4}{4\theta m^2} - \frac{3}{4m^2} \right| \end{aligned} \quad (\text{B.20})$$

$$\leq \left\langle \mathbf{q}, \frac{(\mathbf{Y}\mathbf{Y}^T)^{-1/2}\mathbf{Y}\boldsymbol{\eta}^{\circ 3}}{4m} - \frac{3(1-\theta)}{4\theta m^2}\mathbf{A}\boldsymbol{\zeta}^{\circ 3} - \frac{3}{4m^2}\mathbf{q} \right\rangle \quad (\text{B.21})$$

$$\leq \left\| \frac{(\mathbf{Y}\mathbf{Y}^T)^{-1/2}\mathbf{Y}\boldsymbol{\eta}^{\circ 3}}{4m} - \frac{3(1-\theta)}{4\theta m^2}\mathbf{A}\boldsymbol{\zeta}^{\circ 3} - \frac{3}{4m^2}\mathbf{q} \right\|_2 \quad (\text{B.22})$$

$$\begin{aligned} &\leq \frac{1}{4m} \left\| (\mathbf{Y}\mathbf{Y}^T)^{-1/2}\mathbf{Y}\boldsymbol{\eta}^{\circ 3} - (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2}\mathbf{Y}\boldsymbol{\eta}^{\circ 3} \right\|_2 \\ &\quad + \frac{1}{4\theta^{1/2}m^{3/2}} \left\| (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2}\mathbf{Y}(\boldsymbol{\eta}^{\circ 3} - \bar{\boldsymbol{\eta}}^{\circ 3}) \right\|_2 \\ &\quad + \left\| \frac{1}{4\theta^{1/2}m^{3/2}} (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2}\mathbf{Y}\bar{\boldsymbol{\eta}}^{\circ 3} - \frac{3(1-\theta)}{4\theta m^2}\mathbf{A}\boldsymbol{\zeta}^{\circ 3} - \frac{3}{4m^2}\mathbf{q} \right\|_2 \end{aligned} \quad (\text{B.23})$$

$$\leq \frac{3c_\star}{8\kappa^2} \frac{1-\theta}{\theta m^2} \min_{\mathbf{q} \in \mathcal{R}_{2C_\star}} \|\mathbf{A}^T \mathbf{q}\|_4^6 \quad (\text{B.24})$$

with probability no smaller than $1 - 2 \exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) - 48k^{-7} - 48m^{-5}$. The last inequality is derived with similar arguments in [Lemma 4.2](#), for simplicity, we do not present them here. Moreover, with [Lemma C.1](#), we can obtain an initialization point \mathbf{q}_{init} such that

$$\|\mathbf{A}^T \mathbf{q}_{\text{init}}\|_4^4 \geq (3C_\star \mu \kappa^2)^{2/3} \quad (\text{B.25})$$

$$\geq (2C_\star \mu \kappa^2)^{2/3} + \mu/2. \quad (\text{B.26})$$

Consider any descent method for ψ , which generates a sequence of iterates $\mathbf{q}^{(0)} = \mathbf{q}_{\text{init}}, \mathbf{q}^{(1)}, \dots, \mathbf{q}^{(k)}, \dots$ such that $\psi(\mathbf{q}^{(k)})$ is non-increasing with k . Then

$$\psi(\mathbf{q}^{(k)}) \leq \psi(\mathbf{q}_{\text{init}}) \quad (\text{B.27})$$

$$\begin{aligned} &\leq \frac{3(1-\theta)}{\theta m^2} \varphi(\mathbf{q}_{\text{init}}) + \frac{3}{4m^2} \\ &\quad + \frac{3c_\star}{8\kappa^2} \frac{1-\theta}{\theta m^2} \min_{\mathbf{q} \in \mathcal{R}_{2C_\star}} \|\mathbf{A}^T \mathbf{q}\|_4^6. \end{aligned} \quad (\text{B.28})$$

On the other hand, the finite sample objective function value ψ is close to that of $\frac{3(1-\theta)}{\theta m^2} \varphi(\mathbf{q}) - \frac{3}{4m^2}$,

$$\begin{aligned} &\frac{3(1-\theta)}{\theta m^2} \varphi(\mathbf{q}^{(k)}) \\ &\leq \psi(\mathbf{q}^{(k)}) + \frac{3}{4m^2} + \frac{3c_\star}{8\kappa^2} \frac{1-\theta}{\theta m^2} \min_{\mathbf{q} \in \mathcal{R}_{2C_\star}} \|\mathbf{A}^T \mathbf{q}\|_4^6 \end{aligned} \quad (\text{B.29})$$

$$\leq \frac{3(1-\theta)}{\theta m^2} \varphi(\mathbf{q}_{\text{init}}) + \frac{3c_\star}{4\kappa^2} \frac{1-\theta}{\theta m^2} \min_{\mathbf{q} \in \mathcal{R}_{2C_\star}} \|\mathbf{A}^T \mathbf{q}\|_4^6, \quad (\text{B.30})$$

Therefore, we obtain that

$$\varphi(\mathbf{q}^{(k)}) \leq \varphi(\mathbf{q}_{\text{init}}) + \frac{\mu}{2} \quad (\text{B.31})$$

$$\leq \varphi(\mathbf{q}_{\text{init}}) + \frac{c_\star}{4\kappa^2} \min_{\mathbf{q} \in \mathcal{R}_{2C_\star}} \|\mathbf{A}^T \mathbf{q}\|_4^6, \quad (\text{B.32})$$

which implies that $\mathbf{q}^{(k)} \in \hat{\mathcal{R}}_{2C_*}$ always holds. At last, [Theorem B.2](#) says that any local minimum $\bar{\mathbf{q}}$ is close to $\pm \mathbf{a}_i$ for some i , in the sense that

$$|\langle \bar{\mathbf{q}}, \mathcal{P}_{\mathbb{S}}[\mathbf{a}_i] \rangle| \geq 1 - c_* \kappa^{-2}. \quad (\text{B.33})$$

Write $\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T = \mathbf{A}_0 (\mathbf{I} + \Delta) \mathbf{A}_0^T$ with $\|\Delta\|_2 \leq \delta$, and let

$$\bar{\mathbf{q}} = \pm \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} + \sqrt{2 \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} \delta, \quad (\text{B.34})$$

with $\|\delta\|_2 = 1$. Since

$$\mathbf{a}_i = (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \boldsymbol{\iota}_k^* s_{-(k-i)}[\widetilde{\mathbf{a}}_0], \quad (\text{B.35})$$

we have

$$\begin{aligned} & \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \bar{\mathbf{q}} \\ &= \pm \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \left[\frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} + \sqrt{2 \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} \delta \right] \end{aligned} \quad (\text{B.36})$$

$$\begin{aligned} &= \pm \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \frac{\boldsymbol{\iota}_k^* s_{-(k-i)}[\widetilde{\mathbf{a}}_0]}{\|\mathbf{a}_i\|_2} \\ &\quad + \sqrt{2 \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \delta \end{aligned} \quad (\text{B.37})$$

therefore the error can be bounded as

$$\begin{aligned} & \left\| \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \bar{\mathbf{q}} \pm \frac{\boldsymbol{\iota}_k^* s_{-(k-i)}[\widetilde{\mathbf{a}}_0]}{\|\mathbf{a}_i\|_2} \right\|_2 \\ &\leq \left\| \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{I} \right\|_2 \left\| \frac{\boldsymbol{\iota}_k^* s_{-(k-i)}[\widetilde{\mathbf{a}}_0]}{\|\mathbf{a}_i\|_2} \right\|_2 \\ &\quad + \sqrt{2 \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} \left\| \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \right\|_2. \end{aligned} \quad (\text{B.38})$$

Finally, using the fact that for any nonzero vectors \mathbf{u} and \mathbf{v} that $\langle \mathbf{u}, \mathbf{v} \rangle \geq 0$,

$$\left\| \frac{\mathbf{u}}{\|\mathbf{u}\|_2} - \frac{\mathbf{v}}{\|\mathbf{v}\|_2} \right\|_2 \leq \frac{\sqrt{2}}{\|\mathbf{v}\|_2} \|\mathbf{u} - \mathbf{v}\|_2 \quad (\text{B.39})$$

always holds. Therefore,

$$\begin{aligned} & \|\bar{\mathbf{a}} \pm \mathcal{P}_{\mathbb{S}}[\boldsymbol{\iota}_k s_i[\widetilde{\mathbf{a}}_0]]\|_2 \\ &= \left\| \mathcal{P}_{\mathbb{S}} \left[\left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \bar{\mathbf{q}} \right] \pm \mathcal{P}_{\mathbb{S}}[\boldsymbol{\iota}_k s_i[\widetilde{\mathbf{a}}_0]] \right\|_2 \end{aligned} \quad (\text{B.40})$$

$$\leq \frac{\sqrt{2} \|\mathbf{a}_i\|_2}{\|\boldsymbol{\iota}_k^* s_{i-k}[\widetilde{\mathbf{a}}_0]\|_2} \left\| \left(\frac{\mathbf{Y} \mathbf{Y}^T}{\theta m} \right)^{1/2} \bar{\mathbf{q}} \pm \frac{\boldsymbol{\iota}_k^* s_{i-k}[\widetilde{\mathbf{a}}_0]}{\|\mathbf{a}_i\|_2} \right\|_2 \quad (\text{B.41})$$

$$\leq \kappa \sqrt{2(1+\delta) \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} \quad (\text{B.42})$$

$$\begin{aligned}
& + \sqrt{2}\kappa \left\| \left(\frac{\mathbf{Y}\mathbf{Y}^T}{\theta m} \right)^{1/2} (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} - \mathbf{I} \right\|_2 \\
& \leq 2\kappa \sqrt{2 \left(1 - \left| \left\langle \bar{\mathbf{q}}, \frac{\mathbf{a}_i}{\|\mathbf{a}_i\|_2} \right\rangle \right| \right)} + \sqrt{2}\kappa^3 \delta / \sigma_{\min}
\end{aligned} \tag{B.43}$$

(Lemma D.2)

$$\leq 4\sqrt{c_\star} + 10\sqrt{2}\kappa^3 \sigma_{\min}^{-1} \sqrt{k \log m / m} \tag{B.44}$$

$$\leq 4\sqrt{c_\star} + ck^{-1}, \tag{B.45}$$

completing the proof. \blacksquare

C Initialization

Lemma C.1. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$. There exists a positive constant $C > 2560^4$ such that whenever

$$m \geq C \theta^2 \sigma_{\min}^{-2} \kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log(\kappa k / \sigma_{\min}) \tag{C.1}$$

and the sparsity rate

$$\begin{aligned}
64k^{-1} \log k \leq \theta \leq \min \left\{ \frac{1}{48^2} \mu^{-2} k^{-1} \log^{-2} k, \right. \\
\left. \left(\frac{1}{4} - \frac{640}{C^{1/4}} \right) (3C_\star \mu \kappa^2)^{-2/3} k^{-1} (1 + 36\mu^2 k \log k)^{-2} \right\},
\end{aligned} \tag{C.2}$$

Then the initialization $\mathbf{q}_{\text{init}} = \mathcal{P}_{\mathbb{S}} \left[(\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{y}_i \right]$ satisfies

$$\|\mathbf{A}^T \mathbf{q}_{\text{init}}\|_4^6 \geq 3C_\star \mu \kappa^2, \tag{C.3}$$

namely $\mathbf{q}_{\text{init}} \in \hat{\mathcal{R}}_{3C_\star}$, with probability no smaller than $1 - k^{-1} - 8k^{-2} - 2\exp(-\theta k) - 48k^{-7} - 48m^{-5} - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right)$.

Proof Since

$$m \geq C \frac{\theta^2}{\sigma_{\min}^2} \kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log(\kappa k / \sigma_{\min}) \tag{C.4}$$

with $C \geq 2560^4$, then from Lemma D.1, then with probability no smaller than $1 - 2\exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right) - 48k^{-7} - 48m^{-5}$, we have

$$\delta \doteq \left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \tag{C.5}$$

$$\leq 10\sqrt{k \log m / m} \tag{C.6}$$

$$\begin{aligned}
& \leq \frac{10\sigma_{\min}}{\theta \sigma_{\min}^{-1} \kappa^3 k (1 + 36\mu^2 k \log k)^2} \times \\
& \sqrt{\frac{\log\left(\frac{C \kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log\left(\frac{\kappa k}{\sigma_{\min}}\right)}{\sigma_{\min}^2}\right)}{C \log(\kappa k / \sigma_{\min})}}
\end{aligned} \tag{C.7}$$

$$\leq \frac{20\sigma_{\min}}{C^{1/4} \theta \kappa^3 k (1 + 36\mu^2 k \log k)^2} \tag{C.8}$$

obtains, and the last inequality holds when $C \geq 1000$ that

$$\log(37^4 C) \leq \log 2\sqrt{C}. \quad (\text{C.9})$$

Therefore

$$\begin{aligned} & C\sigma_{\min}^{-2}\kappa^6 k^3 (1 + 36\mu^2 k \log k)^4 \log(\kappa k/\sigma_{\min}) \\ & \leq 37^4 C (\kappa k/\sigma_{\min})^7 \log^5(\kappa k/\sigma_{\min}) \end{aligned} \quad (\text{C.10})$$

$$\leq 37^4 C (\kappa k/\sigma_{\min})^{12} \quad (\text{C.11})$$

or

$$\begin{aligned} & \sqrt{\frac{\log\left(\frac{C\kappa^6 k^3 (1+36\mu^2 k \log k)^4}{\sigma_{\min}^2} \log\left(\frac{\kappa k}{\sigma_{\min}}\right)\right)}{C \log(\kappa k/\sigma_{\min})}} \\ & \leq \sqrt{\frac{\log(37^4 C) + 12 \log(\kappa k/\sigma_{\min})}{C \log(\kappa k/\sigma_{\min})}} \end{aligned} \quad (\text{C.12})$$

$$\leq \sqrt{\frac{\log 2}{\sqrt{C} \log(\kappa k/\sigma_{\min})}} + \frac{12}{C} \quad (\text{C.13})$$

$$\leq \frac{2}{C^{1/4}} \quad (k \geq 2, C \geq 16) \quad (\text{C.14})$$

Moreover, $\kappa^2 \delta \leq 1/2$ always holds provided

$$C \geq \left(\frac{40}{\theta k (1 + 36\mu^2 k \log k)^2} \right)^4. \quad (\text{C.15})$$

Notice that because θ is lower bounded by $c \log k/k$, the right hand side is indeed bounded by an absolute constant.

Set $\zeta_{\text{init}} = \mathbf{A}^T \mathbf{q}_{\text{init}}$ and $\hat{\zeta}_{\text{init}} = \mathcal{P}_{\mathbb{S}}[\mathbf{A}^T \mathbf{A} \mathbf{x}_i]$. Then using for any nonzero vectors \mathbf{u} and \mathbf{v} ,

$$\left\| \frac{\mathbf{u}}{\|\mathbf{u}\|_2} - \frac{\mathbf{v}}{\|\mathbf{v}\|_2} \right\|_2 \leq \frac{2}{\|\mathbf{v}\|_2} \|\mathbf{u} - \mathbf{v}\|_2, \quad (\text{C.16})$$

we have that

$$\begin{aligned} & \left\| \zeta_{\text{init}} - \hat{\zeta}_{\text{init}} \right\|_2 \\ & = \left\| \mathbf{A}^T \mathcal{P}_{\mathbb{S}} \left[(\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{A}_0 \mathbf{x}_i \right] - \mathcal{P}_{\mathbb{S}}[\mathbf{A}^T \mathbf{A} \mathbf{x}_i] \right\|_2 \end{aligned} \quad (\text{C.17})$$

$$= \left\| \frac{\mathbf{A}^T \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 \mathbf{x}_i}{\left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 \mathbf{x}_i \right\|_2} - \frac{\mathbf{A}^T \mathbf{A} \mathbf{x}_i}{\|\mathbf{A}^T \mathbf{A} \mathbf{x}_i\|_2} \right\|_2 \quad (\text{C.18})$$

$$\leq \frac{2}{\|\mathbf{A} \mathbf{x}_i\|_2} \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 \mathbf{x}_i - \mathbf{A} \mathbf{x}_i \right\|_2 \quad (\text{C.19})$$

$$\leq 2 \|\mathbf{A}_0\|_2 \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right\|_2 \quad (\text{C.20})$$

$$\leq \frac{8\kappa^3 \delta}{\sigma_{\min}}, \quad (\text{C.21})$$

where we have used [Lemma D.3](#) in the final bound.

Since $\|\cdot\|_4^4$ is convex, $\|\zeta_{\text{init}}\|_4^4$ can be lower bounded via

$$\|\zeta_{\text{init}}\|_4^4 \geq \|\hat{\zeta}_{\text{init}}\|_4^4 + 4 \left\langle \hat{\zeta}_{\text{init}}^{\circ 3}, \zeta_{\text{init}} - \hat{\zeta}_{\text{init}} \right\rangle \quad (\text{C.22})$$

$$\geq \|\hat{\zeta}_{\text{init}}\|_4^4 - 4 \|\zeta_{\text{init}} - \hat{\zeta}_{\text{init}}\|_2 \quad (\text{C.23})$$

$$\geq \|\hat{\zeta}_{\text{init}}\|_4^4 - \frac{32\kappa^3\delta}{\sigma_{\min}}. \quad (\text{C.24})$$

Let $I = \text{supp}(\mathbf{x}_i)$, then the vector $\hat{\zeta}_{\text{init}} = \mathcal{P}_{\mathbb{S}}[\mathbf{A}^T \mathbf{A} \mathbf{x}_i]$ is composed of $|I|$ large components and small components on the off-support I^c of \mathbf{x}_i .

Dense Component of $\hat{\zeta}_{\text{init}}$. Note that $\|(\mathbf{A}^T \mathbf{A})_{I^c, I} \mathbf{x}_i\|_2 \leq \|\text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2$ with $\|\text{offdiag}(\mathbf{A}^T \mathbf{A})\|_{\infty} \leq \mu$. We have

$$\mathbb{E}[\text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i] = \mathbf{0} \quad (\text{C.25})$$

$$\begin{aligned} \mathbb{E}[|e_j^T \text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i|^2] &= \theta \|e_j^T \text{offdiag}(\mathbf{A}^T \mathbf{A})\|_2^2 \\ &\leq \mu^2 \theta k \end{aligned} \quad (\text{C.26})$$

With Bernstein's Inequality, the summation of moment-bounded independent random variables can be controlled via

$$\mathbb{P}[|e_j^T \text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i| \geq \mu t] \leq 2 \exp\left(-\frac{t^2}{2\theta k + 2t}\right) \quad (\text{C.27})$$

and via union bound

$$\mathbb{P}[\|\text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2 \geq 2k(\mu t)^2] \leq 4k \exp\left(-\frac{t^2}{2\theta k + 2t}\right) \quad (\text{C.28})$$

Therefore, setting $t^2 = 9\theta k \log k$, we obtain

$$\|\text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2 \leq 18\mu^2 \theta k^2 \log k \quad (\text{C.29})$$

with failure probability bounded by

$$\begin{aligned} &4k \exp\left(-\frac{9\theta k \log k}{2\theta k + 2\sqrt{9\theta k \log k}}\right) \\ &= 4k \exp\left(-\frac{9 \log k}{2 + 6\sqrt{(\theta k)^{-1} \log k}}\right) \end{aligned} \quad (\text{C.30})$$

$$\leq 4k^{-2} \quad (\text{C.31})$$

The last inequality is derived under the assumption $(\theta k)^{-1} \log k \leq \frac{1}{64}$.

Spiky Component of $\hat{\zeta}_{\text{init}}$. On the other hand,

$$\mathbb{E}[\|\text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2] = \theta \|\text{diag}(\mathbf{A}^T \mathbf{A})\|_F^2 \quad (\text{C.32})$$

$$= \theta k. \quad (\text{C.33})$$

For $\text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i$, applying the moment control Bernstein Inequality, we have

$$\mathbb{P} \left[\left| \|\text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2 - \mathbb{E}[\cdot] \right| \geq t \right] \leq 2 \exp \left(-\frac{t^2}{2\theta k + 2t} \right). \quad (\text{C.34})$$

By setting $t = 2\sqrt{\theta k \log k}$, we obtain that with probability no smaller than $1 - k^{-1}$,

$$\|\text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2 \geq \theta k - 2\sqrt{\theta k \log k}. \quad (\text{C.35})$$

Denote the following events for the entry-wise magnitude

$$\mathcal{E}_j = \{|e_j^T \text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i| \leq \mu t\}, \quad (\text{C.36})$$

and for the support size

$$\mathcal{E}_{\text{supp}} = \{\|\mathbf{x}_i\|_0 \leq 4\theta k\}. \quad (\text{C.37})$$

On their intersection $\mathcal{E}_{\text{supp}} \cap \bigcap_{j=1}^{2k} \mathcal{E}_j$, we have

$$\|\text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i\|_2^2 \leq 4\theta k(\mu t)^2. \quad (\text{C.38})$$

The the failure probability can be bounded from the union bound as

$$\begin{aligned} & \mathbb{P} \left[\|\text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i\|_2^2 \geq 4\theta k(\mu t)^2 \right] \\ & \leq \mathbb{P} \left[\left(\mathcal{E}_{\text{supp}} \cap \bigcap_j \mathcal{E}_j \right)^c \right] \end{aligned} \quad (\text{C.39})$$

$$= \mathbb{P} \left[\mathcal{E}_{\text{supp}}^c \cup \bigcup_j \mathcal{E}_j^c \right] \quad (\text{C.40})$$

$$\leq \mathbb{P}[\mathcal{E}_{\text{supp}}^c] + \sum_j \mathbb{P}[\mathcal{E}_j^c] \quad (\text{C.41})$$

$$\leq \exp(-\theta k) + 4k \exp \left(-\frac{t^2}{2\theta k + 2t} \right). \quad (\text{C.42})$$

Therefore, by setting $t^2 = 9\theta k \log k$, we obtain

$$\|\text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i\|_2^2 \leq 36\mu^2 \theta^2 k^2 \log k \quad (\text{C.43})$$

with probability no smaller than $1 - \exp(-\theta k) - 8k^{-2}$. Therefore, with probability no smaller than $1 - k^{-1} - 8k^{-2} - \exp(-\theta k)$,

$$\|\text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i\|_2^2 \geq \theta k - 2\sqrt{\theta k \log k} \quad (\text{C.44})$$

$$\|\text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i\|_2^2 \leq 36\mu^2 \theta^2 k^2 \log k \quad (\text{C.45})$$

and via Cauchy-Schwartz inequality, we obtain

$$\left\| (\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right\|_2^2 \quad (\text{C.46})$$

$$= \left\| \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i + \text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right\|_2^2 \quad (\text{C.47})$$

$$\begin{aligned}
&= \left\| \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i \right\|_2^2 + \left\| \text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right\|_2^2 \\
&\quad + 2 \left\langle \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i, \text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right\rangle
\end{aligned} \tag{C.48}$$

$$\begin{aligned}
&\geq \left\| \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i \right\|_2^2 \\
&\quad - 2 \left\| \text{diag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i \right\|_2 \left\| \text{offdiag}(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right\|_2
\end{aligned} \tag{C.49}$$

$$\geq \theta k \left(1 - 2\sqrt{(\theta k)^{-1} \log k} - 12\mu\sqrt{\theta k \log k} \right) \tag{C.50}$$

$$\geq \theta k/2. \tag{C.51}$$

The last equation is derived by plugging in

$$(\theta k)^{-1} \log k \leq \frac{1}{64}, \quad \mu^2 \theta k \log k \leq \frac{1}{48^2} \tag{C.52}$$

under the assumption

$$64k^{-1} \log k \leq \theta \leq \frac{1}{48^2} \mu^{-2} k^{-1} \log^{-1} k. \tag{C.53}$$

Lower Bound of $\|\cdot\|_4^4$. Since with probability no smaller than $1 - 4k^{-2}$, $\left\| \text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i \right\|_2^2 \leq 36\mu^2 \theta k^2 \log k$ obtains and the relative $\|\cdot\|_2^2$ norm between the flat entries to the spiky entries in $\mathbf{A}^T \mathbf{A} \mathbf{x}_i$ can be bounded as

$$\frac{\left\| (\mathbf{A}^T \mathbf{A})_{I^c, I} \mathbf{x}_i \right\|_2^2}{\left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i \right\|_2^2} \leq \frac{\left\| \text{offdiag}(\mathbf{A}^T \mathbf{A}) \mathbf{x}_i \right\|_2^2}{\left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i \right\|_2^2} \tag{C.54}$$

$$\leq 36\mu^2 k \log k \doteq r. \tag{C.55}$$

Since

$$\left\| \hat{\zeta}_{\text{init}} \right\|_4^4 = \left\| \mathcal{P}_{\mathbb{S}} [\mathbf{A}^T \mathbf{A} \mathbf{x}_i] \right\|_4^4 \tag{C.56}$$

$$\begin{aligned}
&= \frac{1}{\left\| \mathbf{A}^T \mathbf{A} \mathbf{x}_i \right\|_2^4} \left\| (\mathbf{A}^T \mathbf{A})_{I^c, I} \mathbf{x}_i \right\|_4^4 \\
&\quad + \frac{1}{\left\| \mathbf{A}^T \mathbf{A} \mathbf{x}_i \right\|_2^4} \left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i \right\|_4^4
\end{aligned} \tag{C.57}$$

$$\geq \frac{1}{\left\| \mathbf{A}^T \mathbf{A} \mathbf{x}_i \right\|_2^4} \left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i \right\|_4^4 \tag{C.58}$$

$$= \frac{\left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i \right\|_2^4 \left\| \mathcal{P}_{\mathbb{S}} [(\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i] \right\|_4^4}{\left\| (\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i + (\mathbf{A}^T \mathbf{A})_{I^c, I} \mathbf{x}_i \right\|_2^4} \tag{C.59}$$

$$\geq \frac{1}{(1+r)^2} \left\| \mathcal{P}_{\mathbb{S}} [(\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i] \right\|_4^4 \tag{C.60}$$

and with high probability $1 - \exp(-\theta k)$ according to [Lemma A.4](#), $\mathcal{P}_{\mathbb{S}} [(\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i]$ satisfies

$$\left\| \mathcal{P}_{\mathbb{S}} [(\mathbf{A}^T \mathbf{A})_{I, I} \mathbf{x}_i] \right\|_4^4 \geq \frac{1}{\left\| \mathbf{x}_i \right\|_0} \geq \frac{1}{2\theta(2k-1)}, \tag{C.61}$$

Together, we have

$$\left\| \zeta_{\text{init}} \right\|_4^4 \geq \left\| \hat{\zeta}_{\text{init}} \right\|_4^4 - \frac{32\kappa^3 \delta}{\sigma_{\min}} \tag{C.62}$$

$$\begin{aligned} &\geq \frac{1}{(1+r)^2} \left\| \mathcal{P}_{\mathbb{S}} \left[(\mathbf{A}^T \mathbf{A})_{I,I} \mathbf{x}_i \right] \right\|_4^4 \\ &\quad - \frac{640C^{-1/4}}{\theta k (1 + 36\mu^2 k \log k)^2} \end{aligned} \quad (\text{C.63})$$

$$\geq \left(\frac{1}{4} - \frac{640}{C^{1/4}} \right) \frac{1}{\theta k (1 + 36\mu^2 k \log k)^2} \quad (\text{C.64})$$

holds with probability no smaller than $1 - k^{-1} - 8k^{-2} - 2\exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) - 48k^{-7} - 48m^{-5}$. To make sure $\|\zeta_{\text{init}}\|_4^6 \geq 3C_*\mu\kappa^2$ as desired, we require the sparsity to satisfy

$$\theta \leq \left(\frac{1}{4} - \frac{640}{C^{1/4}} \right) (3C_*\mu\kappa^2)^{-2/3} k^{-1} (1 + 36\mu^2 k \log k)^{-2}, \quad (\text{C.65})$$

then the initialization $\mathbf{q}_{\text{init}} \in \hat{\mathcal{R}}_{3C_*}$ follows by [Definition 2.1](#). \blacksquare

D Preconditioning

Lemma D.1. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$, then following inequality holds

$$\left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \leq 10\sqrt{k \log m / m}, \quad (\text{D.1})$$

with probability no smaller than $1 - 2\exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) - 48k^{-7} - 48m^{-5}$.

Proof Since

$$\begin{aligned} &\left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \\ &\leq \left\| \text{diag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) - \mathbf{I} \right\|_2 \\ &\quad + \left\| \text{offdiag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) \right\|_2, \end{aligned} \quad (\text{D.2})$$

which is bounded by δ with probability no smaller than $1 - \varepsilon_d - \varepsilon_o$ whenever the probability that each of the terms is upper bounded by $\delta/2$ satisfies

$$\mathbb{P} \left[\left\| \text{diag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) - \mathbf{I} \right\|_2 \geq \delta/2 \right] \leq \varepsilon_d, \quad (\text{D.3})$$

$$\mathbb{P} \left[\left\| \text{offdiag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) - \mathbf{I} \right\|_2 \geq \delta/2 \right] \leq \varepsilon_o. \quad (\text{D.4})$$

Diagonal of $\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T$. Note that $\text{diag}(\mathbf{X}_0 \mathbf{X}_0^T) = \|\mathbf{x}_0\|_2^2 \mathbf{I}$, so

$$\left\| \text{diag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) - \mathbf{I} \right\|_2 = \left| \frac{1}{\theta m} \|\mathbf{x}_0\|_2^2 - 1 \right|. \quad (\text{D.5})$$

We calculate the moment for each summand of $\|\mathbf{x}_0\|_2^2$. The summands can be seen as a χ_1^2 random variable but populated with probability θ , whence

$$\mathbb{E}_{\mathbf{x}_i \sim \text{BG}(\theta)} \left[(x_i^2)^p \right] = \theta \mathbb{E}_{X_i \sim \chi_1^2} [X_i^p] \quad (\text{D.6})$$

$$= \theta \frac{\Gamma(p + \frac{1}{2})}{\Gamma(\frac{1}{2})} \quad (\text{D.7})$$

$$\leq \frac{\theta p! (2)^p}{2} \quad (\text{D.8})$$

$$= \frac{p!}{2} \sigma^2 R^{p-2}. \quad (\text{D.9})$$

Apply Bernstein's inequality for moment bounded random variables (G.4) with $R = 2, \sigma^2 = 4\theta$, then

$$\mathbb{P} \left[\left| \frac{1}{m} \|\mathbf{x}_0\|_2^2 - \theta \right| \geq t \right] \leq 2 \exp \left(-\frac{mt^2}{8\theta + 4t} \right). \quad (\text{D.10})$$

By taking $t = \frac{1}{2}\theta\delta$, we obtain

$$\begin{aligned} & \mathbb{P} \left[\left\| \text{diag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) - \mathbf{I} \right\|_2 \geq \delta/2 \right] \\ & \leq 2 \exp \left(-\frac{\theta m \delta^2}{32 + 8\delta} \right) \end{aligned} \quad (\text{D.11})$$

$$\leq 2 \exp \left(-\frac{100\theta k \log m}{32 + 80\sqrt{k \log m/m}} \right) \quad (\text{D.12})$$

$$\leq 2 \exp(-\theta k). \quad (\text{D.13})$$

Off-diagonal of $\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T$. Note that $\text{offdiag}(\mathbf{X}_0 \mathbf{X}_0^T)$ is a sub-circulant matrix generated by

$$\mathbf{r}_{\mathbf{x}_0} = [r_{\mathbf{x}_0}(2k-2), \dots, 0, \dots, r_{\mathbf{x}_0}(2k-2)]^T \quad (\text{D.14})$$

with $r_{\mathbf{x}_0}(\tau) = \langle \mathbf{x}_0, s_\tau[\mathbf{x}_0] \rangle$ for $\tau = 1, \dots, 2k-2$. Equivalently, we can write

$$\mathbf{r}_{\mathbf{x}_0} = \mathbf{R}_{\mathbf{x}_0}^T \mathbf{x}_0, \quad (\text{D.15})$$

with

$$\mathbf{R}_{\mathbf{x}_0} = [s_{2k-2}[\mathbf{x}_0], \dots, \mathbf{0}, \dots, s_{2k-2}[\mathbf{x}_0]] \in \mathbb{R}^{m \times (4k-3)}. \quad (\text{D.16})$$

Operator norm of a circulant matrix is defined as the following

$$\left\| \text{offdiag} \left(\frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T \right) \right\|_2 = \max_{l=0, \dots, 4k-4} \left| \left\langle \mathbf{v}_l, \frac{1}{\theta m} \mathbf{r}_{\mathbf{x}_0} \right\rangle \right|, \quad (\text{D.17})$$

where \mathbf{v}_l is the l -th (discrete) Fourier basis vector

$$\mathbf{v}_l = \left[1, e^{l \frac{2\pi j}{4k-3}}, \dots, e^{l(4k-4) \frac{2\pi j}{4k-3}} \right]^T, \quad (\text{D.18})$$

and j is the imaginary unit. Let $v_{l,\tau} = \mathbf{v}_l(2k-2-\tau) + \mathbf{v}_l(2k-2+\tau)$, then

$$\langle \mathbf{v}_l, \mathbf{r}_{\mathbf{x}_0} \rangle = \sum_{\tau=1}^{2k-2} v_{l,\tau} \langle \mathbf{x}_0, s_\tau[\mathbf{x}_0] \rangle \quad (\text{D.19})$$

$$= \sum_{\tau=1}^{2k-2} v_{l,\tau} \sum_{i=0}^{m-1} \mathbf{x}_0(i) \mathbf{x}_0([i+\tau]_m). \quad (\text{D.20})$$

By decoupling (Theorem 3.4.1 of [DIPG99]), the tail probability of the weighted autocorrelation $\langle \mathbf{v}_l, \mathbf{r}_{\mathbf{x}_0} \rangle$ can be upper bounded via

$$\mathbb{P} [|\langle \mathbf{v}_l, \mathbf{r}_{\mathbf{x}_0} \rangle| > t]$$

$$= \mathbb{P} \left[\left| \sum_{\tau=1}^{2k-2} v_{l,\tau} \langle \mathbf{x}_0, s_\tau[\mathbf{x}_0] \rangle \right| > t \right] \quad (\text{D.21})$$

$$\leq 6 \mathbb{P} \left[\left| \sum_{\tau=1}^{2k-2} v_{l,\tau} \langle \mathbf{x}_0, s_\tau[\mathbf{x}'_0] \rangle \right| > \frac{t}{6} \right], \quad (\text{D.22})$$

where $\mathbf{x}'_0 \sim \text{i.i.d. BG}(\theta)$ is an independent copy of the random vector \mathbf{x}_0 , we have Plugging in $\langle \mathbf{v}_l, \mathbf{r}_{\mathbf{x}_0} \rangle = \langle \mathbf{v}_l, \mathbf{R}_{\mathbf{x}_0}^T \mathbf{x}_0 \rangle = \langle \mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l, \mathbf{x}_0 \rangle$.

$$\mathbb{P} \left[\left| \left\langle \mathbf{v}_l, \frac{1}{\theta m} \mathbf{r}_{\mathbf{x}_0} \right\rangle \right| > t \right] \leq 6 \mathbb{P} \left[\left| \frac{1}{\theta m} \langle \mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l, \mathbf{x}_0 \rangle \right| > \frac{t}{6} \right]. \quad (\text{D.23})$$

Again with Bernstein's inequality for moment bounded random variable, we have

$$\begin{aligned} & \mathbb{P} \left[\left| \frac{1}{\theta m} \langle \mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l, \mathbf{x}_0 \rangle \right| \geq t \right] \\ & \leq 2 \exp \left(- \frac{\theta m^2 t^2}{2 \|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_2^2 + 2 \|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_\infty m t} \right) \end{aligned} \quad (\text{D.24})$$

Control $\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_2$.

$$\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_2^2 \leq \|\mathbf{R}_{\mathbf{x}_0}\|_2^2 \|\mathbf{v}_l\|_2^2 = k \|\mathbf{R}_{\mathbf{x}_0}\|_2^2 \quad (\text{D.25})$$

With tail bound of the operator norm of a circulant matrix in [Lemma A.6](#), we have

$$\mathbb{P} [\|\mathbf{R}_{\mathbf{x}_0}\|_2 \geq t] \leq 4m \exp \left(- \frac{t^2}{2\theta m + 2t} \right) \quad (\text{D.26})$$

Control $\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_\infty$. For a discrete Fourier basis \mathbf{v}_l as defined, we have

$$\|\mathbf{v}_l\|_2^2 = \|\mathbf{v}_l\|_0 = 4k - 3, \quad \|\mathbf{v}_l\|_\infty = 1 \quad (\text{D.27})$$

Note that

$$\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_\infty = \max_{\tau=1, \dots, 2k-2} |\langle s_\tau[\mathbf{x}_0], \mathbf{v}_l \rangle| \quad (\text{D.28})$$

and moment control Bernstein inequality implies that

$$\mathbb{P} [|\langle s_\tau[\mathbf{x}_0], \mathbf{v}_l \rangle| \geq t] \leq 2 \exp \left(- \frac{t^2}{2\theta \|\mathbf{v}_l\|_2^2 + 2 \|\mathbf{v}_l\|_\infty t} \right). \quad (\text{D.29})$$

with union bound, we obtain

$$\mathbb{P} [\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_\infty \geq t] \leq \sum_{\tau=1}^{2k-2} \mathbb{P} [|\langle s_\tau[\mathbf{x}_0], \mathbf{v}_l \rangle| \geq t] \quad (\text{D.30})$$

$$\leq 4k \exp \left(- \frac{t^2}{8\theta k + 2t} \right) \quad (\text{D.31})$$

Therefore, by plugging in

$$\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_\infty \leq t_1 = 10\sqrt{\theta k \log k}, \quad (\text{D.32})$$

$$\|\mathbf{R}_{\mathbf{x}_0} \mathbf{v}_l\|_2 \leq t_2 = 5\sqrt{\theta m \log m}, \quad (\text{D.33})$$

we obtain the following probabilities

$$\begin{aligned}\mathbb{P}\left[\|\mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l\|_\infty \geq t_1\right] &\leq 4k \exp\left(-\frac{t_1^2}{8\theta k + 2t_1}\right) \\ &\leq 4k^{-8},\end{aligned}\tag{D.34}$$

$$\begin{aligned}\mathbb{P}\left[\|\mathbf{R}_{\mathbf{x}'_0}\|_2 \geq t_2\right] &\leq 4m \exp\left(-\frac{t_2^2}{2\theta m + 2t_2}\right) \\ &\leq 4m^{-6}.\end{aligned}\tag{D.35}$$

Denoting event

$$\mathbf{E} = \left\{\|\mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l\|_\infty \leq t_1, \|\mathbf{R}_{\mathbf{x}'_0}\|_2 \leq t_2\right\},\tag{D.36}$$

and combining these bounds with (D.23), we obtain

$$\begin{aligned}\mathbb{P}\left[\left\|\text{offdiag}\left(\frac{1}{\theta m}\mathbf{X}_0\mathbf{X}_0^T\right)\right\|_2 \geq \delta/2\right] \\ \leq 6\mathbb{P}\left[\max_l \left|\frac{1}{\theta m}\langle \mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l, \mathbf{x}_0 \rangle\right| \geq \frac{\delta}{12}\right]\end{aligned}\tag{D.37}$$

$$\leq 12k\mathbb{P}\left[\left|\frac{1}{\theta m}\langle \mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l, \mathbf{x}_0 \rangle\right| \geq \frac{\delta}{12}\right]\tag{D.38}$$

$$\begin{aligned}&\leq 12k\mathbb{P}\left[\|\mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l\|_\infty > t_1\right] + 12k\mathbb{P}\left[\|\mathbf{R}_{\mathbf{x}'_0}\|_2 > t_2\right] \\ &\quad + 12k\mathbb{P}\left[\left|\frac{1}{\theta m}\langle \mathbf{R}_{\mathbf{x}'_0}\mathbf{v}_l, \mathbf{x}_0 \rangle\right| \geq \frac{\delta}{12} \mid \mathbf{E}\right]\end{aligned}\tag{D.39}$$

$$\begin{aligned}&\leq 24k \exp\left(-\frac{100\theta km \log m/144}{50\theta m \log m + \frac{200}{12}k\sqrt{\theta m \log k \log m}}\right) \\ &\quad + 12k(4k^{-8} + 4m^{-6})\end{aligned}\tag{D.40}$$

$$\begin{aligned}&\quad \left(t_1 = 10\sqrt{\theta k \log k}, t_2 = 5\sqrt{\theta m \log m}\right) \\ &\leq 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right) + 48k^{-7} + 48m^{-5}\end{aligned}\tag{D.41}$$

At last, by combining the control for both the diagonal and off-diagonal term, we obtain that with probability no smaller than $1 - 2\exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\left\{k, 3\sqrt{\theta m}\right\}\right) - 48k^{-7} - 48m^{-5}$,

$$\left\|\frac{1}{\theta m}\mathbf{X}_0\mathbf{X}_0^T - \mathbf{I}\right\|_2 \leq 10\sqrt{k \log m/m},\tag{D.42}$$

holds and completes the proof. ■

Lemma D.2. Suppose $\delta = \left\|\frac{1}{\theta m}\mathbf{X}_0\mathbf{X}_0^T - \mathbf{I}\right\|_2 \leq 1/(2\kappa^2)$, then

$$\left\|\left(\frac{1}{\theta m}\mathbf{Y}\mathbf{Y}^T\right)^{1/2}(\mathbf{A}_0\mathbf{A}_0^T)^{-1/2} - \mathbf{I}\right\|_2 \leq \kappa^2\delta/\sigma_{\min}.\tag{D.43}$$

Proof As in by [Bha97], we denote the directional derivative of f at direction Δ with

$$Df(\mathbf{M})(\Delta) = \left.\frac{d}{dt}\right|_{t=0} f(\mathbf{M} + t\Delta),\tag{D.44}$$

Denote symmetric matrix $M = A_0 A_0^T = U \Lambda U^T$, with λ_{\max} and λ_{\min} being its maximum and minimum eigenvalue. Then we have

$$\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T = M + \Delta, \quad \|\Delta\|_2 \leq \lambda_{\max} \delta. \quad (\text{D.45})$$

Then derivative of f with $Df(M)$. By differential calculus, we can obtain that

$$\begin{aligned} & \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{1/2} (A_0 A_0^T)^{-1/2} - I \right\|_2 \\ &= \left\| (A_0 A_0^T + \Delta)^{1/2} (A_0 A_0^T)^{-1/2} - I \right\|_2 \end{aligned} \quad (\text{D.46})$$

$$= \left\| (A_0 A_0^T)^{-1/2} \int_{t=0}^1 Df(A_0 A_0^T + t\Delta) (\Delta) dt \right\|_2 \quad (\text{D.47})$$

$$\leq \sup_{t \in [0,1]} \|Df(A_0 A_0^T + t\Delta)\|_2 \|\Delta\|_2 \left\| (A_0 A_0^T)^{-1/2} \right\|_2 \quad (\text{D.48})$$

$$\leq \sup_{t \in [0,1]} \|Df(A_0 A_0^T + t\Delta)\|_2 \lambda_{\max} \delta / \sigma_{\min} \quad (\text{D.49})$$

Moreover, we denote $f(t) = t^{1/2}$ and $g(t) = t^2$, then $f = g^{-1}$. The directional derivative of g has following form

$$Dg(M)(X) = MX + XM, \quad (\text{D.50})$$

and directional derivative $Z = Df(M)(X)$ satisfies

$$MZ + ZM = X. \quad (\text{D.51})$$

Denote $M = U \Lambda U^T$ with U orthogonal, without loss of generality,

$$\Lambda Z + Z \Lambda = X. \quad (\text{D.52})$$

Applying Theorem VII.2.3 of [Bha97], we have

$$\|Df(M)(X)\|_2 = \sup_{\|X\|_2 \leq 1} \|Z\|_2 \quad (\text{D.53})$$

$$\leq \int_{t=0}^{\infty} \|e^{-\Lambda t} X e^{-\Lambda t}\|_2 dt \quad (\text{D.54})$$

$$\leq \int_{t=0}^{\infty} e^{-2\lambda_{\min} t} \|X\|_2 dt \quad (\text{D.55})$$

and

$$\begin{aligned} & \sup_{t \in [0,1]} \|Df(A_0 A_0^T + t\Delta)\|_2 \\ & \leq \frac{\|X\|_2}{2(\lambda_{\min} - \lambda_{\max} \delta)} \end{aligned} \quad (\text{D.56})$$

$$\leq 1/\lambda_{\min}. \quad (\text{D.57})$$

Therefore,

$$\left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{1/2} (A_0 A_0^T)^{-1/2} - I \right\|_2 \leq \kappa^2 \delta / \sigma_{\min}. \quad (\text{D.58})$$

■

Lemma D.3. Suppose \mathbf{A}_0 has condition number κ and

$$\delta = \left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \leq 1 / (2\kappa^2) \quad (\text{D.59})$$

then

$$\left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right\|_2 \leq 4\kappa^2 \delta / \sigma_{\min}^2. \quad (\text{D.60})$$

Proof Denote symmetric matrix

$$\mathbf{M} = \mathbf{A}_0 \mathbf{A}_0^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T, \quad (\text{D.61})$$

with λ_{\max} and λ_{\min} being its maximum and minimum eigenvalue. Then we have

$$\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T = \mathbf{M} + \mathbf{\Delta}, \quad \|\mathbf{\Delta}\|_2 \leq \lambda_{\max} \delta. \quad (\text{D.62})$$

Then

$$\begin{aligned} & \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \right\|_2 \\ &= \left\| (\mathbf{M} + \mathbf{\Delta})^{-1/2} - \mathbf{M}^{-1/2} \right\|_2 \end{aligned} \quad (\text{D.63})$$

$$\leq \|\mathbf{\Delta}\|_2 \cdot \sup_{0 \leq t \leq 1} \|Df(\mathbf{M} + t\mathbf{\Delta})\|_2. \quad (\text{D.64})$$

Here, $f(t) = t^{-1/2}$ and Df is the derivative of function f . In addition, we define function $g(t) = t^{-2}$, $h(t) = t^{-1}$, $w(t) = t^2$, and following function compositions hold

$$f = g^{-1}, \quad g = h \circ w. \quad (\text{D.65})$$

For differential function g and if $Dg(f(\mathbf{M})) \neq 0$, we have

$$Df(\mathbf{M}) = [Dg(f(\mathbf{M}))]^{-1}. \quad (\text{D.66})$$

The derivative of function g satisfies the chain rule that

$$Dg(\mathbf{M}) = Dh(w(\mathbf{M}))(Dw(\mathbf{M})). \quad (\text{D.67})$$

Plug in

$$Dh(\mathbf{M})(\mathbf{X}) = -\mathbf{M}^{-1} \mathbf{X} \mathbf{M}^{-1}, \quad (\text{D.68})$$

$$Dw(\mathbf{M})(\mathbf{X}) = \mathbf{M} \mathbf{X} + \mathbf{X} \mathbf{M}, \quad (\text{D.69})$$

we obtain that

$$Dg(\mathbf{M})(\mathbf{X}) \quad (\text{D.70})$$

$$= Dh(w(\mathbf{M}))(Dw(\mathbf{M})(\mathbf{X})) \quad (\text{D.71})$$

$$= Dh(w(\mathbf{M}))[\mathbf{M} \mathbf{X} + \mathbf{X} \mathbf{M}] \quad (\text{D.72})$$

$$= Dh(\mathbf{M}^2)[\mathbf{M} \mathbf{X} + \mathbf{X} \mathbf{M}] \quad (\text{D.73})$$

$$= -\mathbf{M}^{-2}[\mathbf{M} \mathbf{X} + \mathbf{X} \mathbf{M}] \mathbf{M}^{-2} \quad (\text{D.74})$$

$$= -[\mathbf{M}^{-1} \mathbf{X} \mathbf{M}^{-2} + \mathbf{M}^{-2} \mathbf{X} \mathbf{M}^{-1}]. \quad (\text{D.75})$$

Since the function g is differentiable and $Dg(\mathbf{M}) \neq \mathbf{0}$, then

$$Df(\mathbf{M}) = [Dg(f(\mathbf{M}))]^{-1} \quad (\text{D.76})$$

$$= \left[Dg\left(\mathbf{M}^{-1/2}\right) \right]^{-1}. \quad (\text{D.77})$$

Hence, directional derivative $\mathbf{Z} \doteq Df(\mathbf{M})(\mathbf{X})$ satisfies

$$\mathbf{M}^{1/2} \mathbf{Z} \mathbf{M} + \mathbf{M} \mathbf{Z} \mathbf{M}^{1/2} = -\mathbf{X}. \quad (\text{D.78})$$

Denote $\mathbf{M} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ with $\mathbf{\Lambda} \succ 0$ and \mathbf{U} orthogonal, without loss of generality

$$\mathbf{\Lambda} \mathbf{Z} \mathbf{\Lambda}^{1/2} + \mathbf{\Lambda}^{1/2} \mathbf{Z} \mathbf{\Lambda} = -\mathbf{X}. \quad (\text{D.79})$$

Above equation can be reformulated as a Sylvester equation as following

$$\mathbf{\Lambda}^{1/2} \mathbf{Z} - \mathbf{Z} \left(-\mathbf{\Lambda}^{1/2} \right) = -\mathbf{\Lambda}^{-1/2} \mathbf{X} \mathbf{\Lambda}^{-1/2}. \quad (\text{D.80})$$

From Theorem VII.2.3 of [Bha97], when there are no common eigenvalues of $\mathbf{\Lambda}^{1/2}$ and $-\mathbf{\Lambda}^{1/2}$, then there exists a closed form solution for matrix \mathbf{Z} that

$$\mathbf{Z} = \int_{t=0}^{\infty} e^{-\mathbf{\Lambda}^{1/2} t} \left(-\mathbf{\Lambda}^{-1/2} \mathbf{X} \mathbf{\Lambda}^{-1/2} \right) e^{-\mathbf{\Lambda}^{1/2} t} dt \quad (\text{D.81})$$

Therefore, the operator norm of $Df(\mathbf{M})$ can be obtained as

$$\|Df(\mathbf{M})(\mathbf{X})\|_2 = \sup_{\|\mathbf{X}\|_2 \leq 1} \|\mathbf{Z}\|_2 \quad (\text{D.82})$$

$$\leq \int_{t=0}^{\infty} \left\| e^{-\mathbf{\Lambda}^{1/2} t} \left(\mathbf{\Lambda}^{-1/2} \mathbf{X} \mathbf{\Lambda}^{-1/2} \right) e^{-\mathbf{\Lambda}^{1/2} t} \right\| dt \quad (\text{D.83})$$

$$\leq \int_{t=0}^{\infty} e^{-\lambda_{\min} t} \left\| \mathbf{\Lambda}^{-1/2} \mathbf{X} \mathbf{\Lambda}^{-1/2} \right\| dt \quad (\text{D.84})$$

$$\leq \frac{\|\mathbf{X}\|}{\lambda_{\min}^2}. \quad (\text{D.85})$$

Therefore

$$\begin{aligned} & \left\| (\mathbf{M} + \mathbf{\Delta})^{-1/2} - \mathbf{M}^{-1/2} \right\|_2 \\ & \leq \frac{\|\mathbf{\Delta}\|_2}{(\lambda_{\min} - \|\mathbf{\Delta}\|_2)^2} \end{aligned} \quad (\text{D.86})$$

$$\leq \frac{4\|\mathbf{\Delta}\|_2}{\lambda_{\min}^2} \quad (\delta \leq 1/(2\kappa^2)) \quad (\text{D.87})$$

$$\leq \frac{4\lambda_{\max}\delta}{\lambda_{\min}^2} \quad (\text{D.88})$$

$$= \frac{4\kappa^2\delta}{\sigma_{\min}^2} \quad (\text{D.89})$$

■

E Concentration for Gradient (Lemma 4.2)

Lemma E.1. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$. There exists a positive constant C such that whenever

$$m \geq C \frac{\min \left\{ (2C_*\mu)^{-1}, \kappa^2 k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \quad (\text{E.1})$$

and $\theta > \log k/k$, then with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4} - 2 \exp(-\theta k) - 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) - 48k^{-7} - 48m^{-5}$,

$$\left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \leq c \frac{1-\theta}{\theta m^2} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}, \quad (\text{E.2})$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$ with positive constant $c \leq 3/(2C_*)$.

Proof Denote $\boldsymbol{\eta} = \mathbf{Y}^T (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{q}$ and $\bar{\boldsymbol{\eta}} = \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} = (\theta m)^{-1/2} \mathbf{X}_0^T \boldsymbol{\zeta}$, then

$$\begin{aligned} & \left\| \text{grad}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad}[\varphi](\mathbf{q}) \right\|_2 \\ &= \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{1}{m} (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \boldsymbol{\eta}^{\circ 3} - \frac{3(1-\theta)}{\theta m^2} \mathbf{A} \boldsymbol{\zeta}^{\circ 3} \right] \right\|_2 \\ &\leq \underbrace{\frac{1}{m} \left\| (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \boldsymbol{\eta}^{\circ 3} - (\theta m)^{-1/2} \mathbf{A} \mathbf{X}_0 \boldsymbol{\eta}^{\circ 3} \right\|_2}_{\Delta_1^g} \\ &\quad + \underbrace{\frac{1}{\theta^{1/2} m^{3/2}} \left\| \mathbf{A} \mathbf{X}_0 \boldsymbol{\eta}^{\circ 3} - \mathbf{A} \mathbf{X}_0 \bar{\boldsymbol{\eta}}^{\circ 3} \right\|_2}_{\Delta_2^g} \\ &\quad + \underbrace{\left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{1}{\theta^{1/2} m^{3/2}} \mathbf{A} \mathbf{X}_0 \bar{\boldsymbol{\eta}}^{\circ 3} - \frac{3(1-\theta)}{\theta m^2} \mathbf{A} \boldsymbol{\zeta}^{\circ 3} \right] \right\|_2}_{\Delta_3^g}. \end{aligned}$$

First, let us note that

$$\begin{aligned} & C (1-\theta)^{-2} \sigma_{\min}^{-2} \kappa^{10} k^6 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \\ &\leq C \left(\frac{\kappa k}{\sigma_{\min} (1-\theta)} \right)^{10} \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \end{aligned} \quad (\text{E.3})$$

$$\leq C \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)^{13}, \quad (\text{E.4})$$

hence

$$\begin{aligned} & \frac{\log^3 \left(C (1-\theta)^{-2} \sigma_{\min}^{-2} \kappa^{10} k^6 \log^3 \left((1-\theta)^{-1} \sigma_{\min}^{-1} \kappa k \right) \right)}{C \log^3 \left((1-\theta)^{-1} \sigma_{\min}^{-1} \kappa k \right)} \\ &\leq \left(\frac{\log C + 13 \log \left((1-\theta)^{-1} \sigma_{\min}^{-1} \kappa k \right)}{C^{1/3} \log \left((1-\theta)^{-1} \sigma_{\min}^{-1} \kappa k \right)} \right)^3 \end{aligned} \quad (\text{E.5})$$

$$\leq \left(\frac{\log C}{C^{1/3} \log \left((1-\theta)^{-1} \sigma_{\min}^{-1} \kappa k \right)} + \frac{13}{C^{1/3}} \right)^3 \quad (\text{E.6})$$

$$\leq \left(\frac{1}{C^{1/6}} + \frac{1}{2} \frac{1}{C^{1/6}} \right)^3 \quad (C \geq 10^8) \quad (\text{E.7})$$

$$\leq \frac{4}{C^{1/2}}. \quad (\text{E.8})$$

Given

$$m \geq C \frac{\min \left\{ (2C_* \mu)^{-1}, \kappa^2 k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^8 k^4 \log^3 \left(\frac{\kappa k}{\sigma_{\min} (1-\theta)} \right), \quad (\text{E.9})$$

as the ratio $\log^3 m/m$ decreases with increasing m , then

$$\begin{aligned} \frac{\log^3 m}{m} &\leq \frac{\log^3 \left(\frac{C \kappa^{10} k^6}{(1-\theta)^2 \sigma_{\min}^2} \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \right)}{C \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)} \\ &\quad \times \frac{(1-\theta)^2 \sigma_{\min}^2}{\min \left\{ (2C_* \mu)^{-1}, \kappa^2 k^2 \right\} \kappa^8 k^4} \end{aligned} \quad (\text{E.10})$$

$$\leq \frac{4}{C^{1/2}} \frac{(1-\theta)^2 \sigma_{\min}^2}{\min \left\{ (2C_* \mu)^{-1}, \kappa^2 k^2 \right\} \kappa^8 k^4} \quad (\text{E.11})$$

According to [Lemma D.1](#), following inequality always holds

$$\left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \leq \delta \quad (\text{E.12})$$

$$\leq 10 \sqrt{k \log m / m} \quad (\text{E.13})$$

$$\leq \frac{20 (1-\theta) \sigma_{\min} \max \left\{ (2C_* \mu)^{1/2}, (\kappa k)^{-1} \right\}}{C^{1/4} \kappa^4 k^{3/2} \log m} \quad (\text{E.14})$$

$$\leq \frac{20 \sigma_{\min} (1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^6}{C^{1/4} \kappa^3 \kappa^2 k \log m}, \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*}. \quad (\text{E.15})$$

with probability no smaller than $1 - \varepsilon_0$ with $\varepsilon_0 = 2 \exp(-\theta k) + 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) + 48k^{-7} + 48m^{-5}$.

Moreover, $4\kappa^3 \delta / \sigma_{\min} \leq 1/2$ whenever

$$C \geq \left(\frac{160 (1-\theta)}{k \log m} \right)^4, \quad (\text{E.16})$$

whence $\delta \leq 1/(8\kappa^2)$, and [Lemma D.3](#) implies that

$$\begin{aligned} &\left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{A}_0 \right\|_2 \\ &\leq 4\kappa^3 \delta / \sigma_{\min} \end{aligned} \quad (\text{E.17})$$

$$\leq \frac{80 (1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^6}{C^{1/4} k \log m \kappa^2}, \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*}. \quad (\text{E.18})$$

At the same time,

$$\|\mathbf{X}_0\|_2 \leq (\theta m)^{1/2} \sqrt{1+\delta} \leq (\theta m)^{1/2} (1 + \delta/2). \quad (\text{E.19})$$

Moreover, [Lemma A.5](#) implies that with probability no smaller than $1 - \varepsilon_B$, we have

$$\|\mathbf{x}_0\|_\infty \leq \sqrt{2} \log^{1/2} \left(\frac{2\theta m}{\varepsilon_B} \right). \quad (\text{E.20})$$

Upper Bound for Δ_1^g . Using [Lemma A.7](#), on the an event of probability at least $1 - \varepsilon_0 - \varepsilon_B$,

$$\|\boldsymbol{\eta}^{\circ 3}\|_2 = \|\boldsymbol{\eta}\|_6^3 \quad (\text{E.21})$$

$$\leq \left(1 + \frac{4\kappa^3\delta}{\sigma_{\min}}\right)^2 \frac{2k}{\theta m} \|\mathbf{x}_0\|_\infty^2 \quad (\text{E.22})$$

$$\leq \frac{9k}{\theta m} \log \left(\frac{2\theta m}{\varepsilon_B} \right). \quad (\text{E.23})$$

Therefore, we can obtain following upper bound

$$\Delta_1^g = \frac{1}{m} \left\| (\mathbf{Y}\mathbf{Y}^T)^{-1/2} \mathbf{Y} \boldsymbol{\eta}^{\circ 3} - (\theta m)^{-1/2} \mathbf{A} \mathbf{X}_0 \boldsymbol{\eta}^{\circ 3} \right\|_2 \quad (\text{E.24})$$

$$\leq \frac{1}{\theta^{1/2} m^{3/2}} \|\mathbf{X}_0\|_2 \|\boldsymbol{\eta}^{\circ 3}\|_2 \times \left\| \left(\frac{1}{\theta m} \mathbf{Y}\mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 - \mathbf{A} \right\|_2 \quad (\text{E.25})$$

$$\leq \frac{5}{4m} \cdot \frac{4\kappa^3\delta}{\sigma_{\min}} \cdot \frac{9k}{\theta m} \log \left(\frac{2\theta m}{\varepsilon_B} \right) \quad (\text{E.26})$$

$$\leq \frac{900(1-\theta) \log(2\theta m/\varepsilon_B)}{C^{1/4} \theta m^2 \log m} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2} \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*}.$$

Upper Bound for Δ_2^g . Similarly, with probability no smaller than $1 - \varepsilon_0 - \varepsilon_B$, together with [Lemma A.7](#), following upper bound can be obtained

$$\|\boldsymbol{\eta}^{\circ 3} - \bar{\boldsymbol{\eta}}^{\circ 3}\|_2 = \|\boldsymbol{\eta}^{\circ 3} - \text{diag}(\boldsymbol{\eta}^{\circ 2}) \bar{\boldsymbol{\eta}} + \text{diag}(\boldsymbol{\eta}^{\circ 2}) \bar{\boldsymbol{\eta}} - \bar{\boldsymbol{\eta}}^{\circ 3}\|_2 \quad (\text{E.27})$$

$$\leq \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \|\text{diag}(\boldsymbol{\eta}^{\circ 2})\|_2 + \|\bar{\boldsymbol{\eta}}\|_2 \|\text{diag}(\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2})\|_2 \quad (\text{E.28})$$

$$= \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \|\boldsymbol{\eta}\|_\infty^2 + \|\bar{\boldsymbol{\eta}}\|_2 \|\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2}\|_\infty \quad (\text{E.29})$$

$$\leq \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \|\boldsymbol{\eta}\|_\infty^2 + \|\bar{\boldsymbol{\eta}}\|_2 \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_\infty \|\boldsymbol{\eta} + \bar{\boldsymbol{\eta}}\|_\infty \quad (\text{E.30})$$

$$\leq 4(1+\delta/2) \frac{4\kappa^3\delta}{\sigma_{\min}} \frac{k}{\theta m} \log(2\theta m/\varepsilon_B) \times \left[\left(1 + \frac{4\kappa^3\delta}{\sigma_{\min}}\right)^2 + \left(2 + \frac{4\kappa^3\delta}{\sigma_{\min}}\right) \right] \quad (\text{E.31})$$

$$\leq \frac{24k}{\theta m} \log(2\theta m/\varepsilon_B) \cdot \frac{4\kappa^3\delta}{\sigma_{\min}}. \quad (\text{E.32})$$

Therefore, we can obtain following upper bound

$$\Delta_2^g = \frac{1}{\theta^{1/2} m^{3/2}} \|\mathbf{A} \mathbf{X}_0^T \boldsymbol{\eta}^{\circ 3} - \mathbf{A} \mathbf{X}_0^T \bar{\boldsymbol{\eta}}^{\circ 3}\|_2 \quad (\text{E.33})$$

$$\leq \frac{1}{\theta^{1/2} m^{3/2}} \|\mathbf{A}\|_2 \|\mathbf{X}_0\|_2 \|\boldsymbol{\eta}^{\circ 3} - \bar{\boldsymbol{\eta}}^{\circ 3}\|_2 \quad (\text{E.34})$$

$$\leq \frac{5}{4m} \cdot \frac{24k}{\theta m} \log(2\theta m/\varepsilon_B) \cdot \frac{4\kappa^3 \delta}{\sigma_{\min}} \quad (\text{E.35})$$

$$\leq \frac{2400}{C^{1/4}} \frac{1-\theta}{\theta m^2} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2} \cdot \frac{\log(2\theta m/\varepsilon_B)}{\log m}. \quad (\text{E.36})$$

For both Δ_1^g and Δ_2^g to be bounded by $\frac{1}{2C_\star} \frac{1-\theta}{\theta m^2} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}$, we set

$$C \geq \left(4800 C_\star \frac{\log(2\theta m/\varepsilon_B)}{\log m} \right)^4. \quad (\text{E.37})$$

Notice that the right hand side is indeed bounded by a numerical constant for all m .

Tail Bound for Δ_3^g . Note that

$$\begin{aligned} & (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{Y} \bar{\boldsymbol{\eta}}^{\circ 3} \\ &= (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{A}_0 \mathbf{X}_0 \left(\mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} \right)^{\circ 3} \end{aligned} \quad (\text{E.38})$$

$$= (\theta m)^{-3/2} \mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3}, \quad (\text{E.39})$$

and its expectation with respect to \mathbf{x}_0

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{m} \mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} \right] \\ &= \mathbb{E} \left[\mathbf{A} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{A}^T \mathbf{q})^3 \right] \end{aligned} \quad (\text{E.40})$$

$$= 3\theta (1-\theta) \mathbf{A} \boldsymbol{\zeta}^{\circ 3} + 3\theta^2 \|\mathbf{A}^T \mathbf{q}\|_2^2 \mathbf{A} \mathbf{A}^T \mathbf{q} \quad (\text{E.41})$$

$$= 3\theta (1-\theta) \mathbf{A} \boldsymbol{\zeta}^{\circ 3} + 3\theta^2 \mathbf{q}, \quad (\text{E.42})$$

hence

$$\begin{aligned} & \mathbf{P}_{\mathbf{q}^\perp} \left[\mathbb{E} \left[\frac{1}{m} \mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} \right] \right] \\ &= \mathbf{P}_{\mathbf{q}^\perp} [3\theta (1-\theta) \mathbf{A} \boldsymbol{\zeta}^{\circ 3}]. \end{aligned} \quad (\text{E.43})$$

Therefore, the Δ_3^g term can be simplified as

$$\Delta_3^g = \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{1}{\theta^{1/2} m^{3/2}} \mathbf{A} \mathbf{X}_0^T \bar{\boldsymbol{\eta}}^{\circ 3} - \frac{3(1-\theta)}{\theta m^2} \mathbf{A} \boldsymbol{\zeta}^{\circ 3} \right] \right\|_2 \quad (\text{E.44})$$

$$= \frac{1}{\theta^2 m^2} \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{\mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 3}}{m} - 3\theta (1-\theta) \mathbf{A} \boldsymbol{\zeta}^{\circ 3} \right] \right\|_2 \quad (\text{E.45})$$

$$\begin{aligned} & \leq \frac{1}{\theta^2 m^2} \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{\mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 3}}{m} - \mathbb{E}[\cdot] \right] \right\|_2 \\ & \quad + \frac{1}{\theta^2 m^2} \|\mathbf{P}_{\mathbf{q}^\perp} [3\theta^2 \mathbf{q}]\|_2 \end{aligned} \quad (\text{E.46})$$

$$\leq \frac{1}{\theta^2 m^2} \left\| \frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E}[\cdot] \right\|_2. \quad (\text{E.47})$$

Under the assumption that

$$m \geq \frac{C}{(1-\theta)^2} \min \{ \mu^{-1}, \kappa^2 k^2 \} \kappa^2 k^4 \log^3 (\kappa k), \quad (\text{E.48})$$

applying [Lemma E.2](#), we have

$$\left\| \frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} - \mathbb{E} [\cdot] \right\|_2 \leq c \theta (1-\theta) \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}. \quad (\text{E.49})$$

with probability larger than $1 - c_2 \exp(-k) - c_2 k^{-4}$. At last, taking $\varepsilon_B = \theta^2 k^{-4}$, we obtain that

$$\begin{aligned} & \left\| \text{grad} [\psi] (\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{grad} [\varphi] (\mathbf{q}) \right\|_2 \\ & \leq c \frac{1-\theta}{\theta m^2} \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}, \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*} \end{aligned} \quad (\text{E.50})$$

with probability larger than $1 - c_2 \exp(-k) - c_2 k^{-4} - \varepsilon_B - \varepsilon_0$ as desired. \blacksquare

E.1 Proof of [Lemma E.2](#)

Lemma E.2. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^m$. There exist positive constant C such that whenever

$$m \geq \frac{C}{(1-\theta)^2} \min \{ (2C_* \mu)^{-1}, \kappa^2 k^2 \} \kappa^2 k^4 \log^3 (\kappa k) \quad (\text{E.51})$$

and $\theta k \geq 1$, then with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4}$,

$$\left\| \frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} - \mathbb{E} [\cdot] \right\|_2 \leq c \theta (1-\theta) \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2} \quad (\text{E.52})$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$ with positive constant $c \leq 1/(2C_*)$.

Proof Let $\bar{\mathbf{x}}_i \in \mathbb{R}^{2k-1}$ be generated via

$$\bar{\mathbf{x}}_i = \begin{cases} \mathbf{x}_i & \|\mathbf{x}_i\|_\infty \leq B \text{ and } \|\mathbf{x}_i\|_0 \leq 4\theta k \log m \\ \mathbf{0} & \text{else} \end{cases} \quad (\text{E.53})$$

Let $\bar{\mathbf{X}}_0 \in \mathbb{R}^{(2k-1) \times m}$ denote the circulant submatrix generated by $\bar{\mathbf{x}}_0$. Then $\bar{\mathbf{X}}_0 = \mathbf{X}_0$ obtains whenever

1. $\|\mathbf{x}_0\|_\infty \leq B$, which happens with probability no smaller than $1 - 2\theta m e^{-B^2/2}$ according to [Lemma A.5](#);
2. $\|\mathbf{x}_i\|_0 \leq 4\theta k \log m$ holds for any index i , applying [Lemma A.4](#) and Boole's inequality we have

$$\begin{aligned} & \mathbb{E} [\mathbf{1}_{\cup_i \|\mathbf{x}_i\|_0 > 4\theta k \log m}] \\ & \leq m \mathbb{P} [\|\mathbf{x}_i\|_0 > 4\theta k \log m] \end{aligned} \quad (\text{E.54})$$

$$\leq 2m \exp \left(-\frac{3}{4} \theta k \log m \right). \quad (\text{E.55})$$

Denote $\boldsymbol{\zeta} = \mathbf{A}^T \mathbf{q}$ and

$$\mathbf{g}_E = \mathbb{E} \left[\frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} \right], \quad (\text{E.56})$$

$$\bar{\mathbf{g}}_E = \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} \right], \quad (\text{E.57})$$

then,

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbf{g}_E \right\|_2 \geq c\theta(1-\theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2} \right] \\ & \leq \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbf{g}_E \right\|_2 \geq c\theta(1-\theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2} \right] \\ & \quad + 2\theta m e^{-B^2/2} + 2m \exp \left(-\frac{3}{4}\theta k \log m \right) \end{aligned} \quad (\text{E.58})$$

With triangle inequality, we have

$$\begin{aligned} & \left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbf{g}_E \right\|_2 \\ & \leq \left\| \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} \right] - \bar{\mathbf{g}}_E \right\|_2 + \|\bar{\mathbf{g}}_E - \mathbf{g}_E\|_2. \end{aligned} \quad (\text{E.59})$$

Hence, provided

$$\|\bar{\mathbf{g}}_E - \mathbf{g}_E\|_2 \leq \frac{c}{2}\theta(1-\theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2}, \quad (\text{E.60})$$

we have

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbf{g}_E \right\|_2 \geq c\theta(1-\theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2} \right] \\ & \leq \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \bar{\mathbf{g}}_E \right\|_2 \geq \frac{c}{2}\theta(1-\theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2} \right]. \end{aligned} \quad (\text{E.61})$$

Truncation Level Next, we choose a large enough entry-wise truncation level B such that the expectation of the gradient $\mathbb{E} \left[\frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 3} \right]$ is close to that of its truncation $\mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} \right]$.

Moreover, we introduce following events notation

$$\mathcal{E}_i \doteq \{ \|\mathbf{x}_i\|_\infty > B \cup \|\mathbf{x}_i\|_0 > 4\theta k \log m \}, \quad (\text{E.62})$$

then

$$\begin{aligned} & \|\bar{\mathbf{g}}_E - \mathbf{g}_E\|_2 \\ & = \left\| \mathbb{E} \left[\frac{1}{m} \sum_i \mathbf{x}_i \langle \mathbf{x}_i, \boldsymbol{\zeta} \rangle^3 \cdot \mathbf{1}_{\mathcal{E}_i} \right] \right\|_2 \end{aligned} \quad (\text{E.63})$$

$$\leq \frac{1}{m} \sum_i \left\| \mathbb{E} \left[\mathbf{x}_i \langle \mathbf{x}_i, \boldsymbol{\zeta} \rangle^3 \cdot \mathbf{1}_{\mathcal{E}_i} \right] \right\|_2 \quad (\text{E.64})$$

$$\leq \frac{1}{m} \sum_i \left(\mathbb{E} \left[\left\| \mathbf{x}_i (\mathbf{x}_i^T \boldsymbol{\zeta})^{\circ 3} \right\|_2^2 \right] \cdot \mathbb{E} [\mathbf{1}_{\mathcal{E}_i}] \right)^{1/2} \quad (\text{E.65})$$

$$\begin{aligned} & \leq \left(\mathbb{E} \left[\|\mathbf{x}_i\|_2^8 \right] \right)^{1/2} \times \\ & \quad \sqrt{\mathbb{E} [\mathbf{1}_{\|\mathbf{x}_i\|_\infty > B}] + \mathbb{E} [\mathbf{1}_{\|\mathbf{x}_i\|_0 > 4\theta k \log m}]} \end{aligned} \quad (\text{E.66})$$

$$\leq 50k^2 \sqrt{4\theta k e^{-B^2/2} + \exp(-\frac{3}{4}\theta k \log m)} \quad (\text{E.67})$$

By setting

$$B \geq C' \log^{1/2} \left(\frac{\kappa^4 k^8}{\theta (1-\theta)^2} \right), \quad (\text{E.68})$$

we have

$$\theta k e^{-B^2/2} \leq \frac{1}{2} \left(\frac{c}{100} \right)^2 \theta^2 (1-\theta)^2 \frac{\|\zeta\|_4^{12}}{\kappa^4 k^4}. \quad (\text{E.69})$$

In addition, whenever

$$\theta k \geq \frac{4}{3 \log m} \log \left(\frac{400^2 \kappa^4 k^4}{c^2 \theta^2 (1-\theta)^2 \|\zeta\|_4^{12}} \right), \quad (\text{E.70})$$

we have

$$\exp(-\frac{3}{4}\theta k \log m) \leq \frac{1}{2} \left(\frac{c}{100} \right)^2 \theta^2 (1-\theta)^2 \frac{\|\zeta\|_4^{12}}{\kappa^4 k^4}. \quad (\text{E.71})$$

Therefore,

$$\sqrt{4\theta k e^{-B^2/2} + \exp(-\frac{3}{4}\theta k \log m)} \leq \frac{c}{2} \theta (1-\theta) \frac{\|\zeta\|_4^6}{50 \kappa^2 k^2}. \quad (\text{E.72})$$

In addition,

$$\left(\mathbb{E} \left[\|\mathbf{x}_i\|_2^8 \right] \right)^{1/2} \leq (7!! \cdot 2^4 k^4)^{1/2} < 50k^2. \quad (\text{E.73})$$

Plugging in Eq (E.73) and (E.72) back to (E.67), we obtain that

$$\|\bar{\mathbf{g}}_E - \mathbf{g}_E\|_2 \leq \frac{c}{2} \theta (1-\theta) \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2}, \quad (\text{E.74})$$

and hence

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \zeta)^{\circ 3} - \mathbf{g}_E \right\|_2 \geq c \theta (1-\theta) \frac{\|\zeta\|_4^6}{\kappa^2} \right] \\ & \leq \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \zeta)^{\circ 3} - \bar{\mathbf{g}}_E \right\|_2 \geq \frac{c}{2} \theta (1-\theta) \frac{\|\zeta\|_4^6}{\kappa^2} \right]. \end{aligned} \quad (\text{E.75})$$

Independent Submatrices. To deal with the complicated dependence within the random circulant matrix \mathbf{X}_0 , we break \mathbf{X}_0 into submatrices $\mathbf{X}_1, \dots, \mathbf{X}_{2k-1}$, each of which is (marginally) distributed as a $(2k-1) \times \frac{m}{2k-1}$ i.i.d. BG(θ) random matrix. Indeed, there exists a permutation $\mathbf{\Pi}$ such that

$$\mathbf{X}_0 \mathbf{\Pi} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{2k-1}], \quad (\text{E.76})$$

with

$$\mathbf{X}_i = [\mathbf{x}_i, \mathbf{x}_{i+(2k-1)}, \dots, \mathbf{x}_{i+(m-2k-1)}]. \quad (\text{E.77})$$

We apply similar matrix breaking approach for the truncated matrix $\bar{\mathbf{X}}$. The summands within each term $\bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \zeta)^{\circ 3}$ are mutually independent and hence is amenable to classical concentration results.

$$\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \zeta)^{\circ 3} = \frac{1}{m} \sum_{l=1}^m \langle \bar{\mathbf{x}}_l, \zeta \rangle^3 \bar{\mathbf{x}}_l \quad (\text{E.78})$$

$$= \sum_{i=1}^{2k-1} \frac{1}{m} \left(\sum_{j=0}^{\frac{m}{2k-1}-1} \langle \bar{\mathbf{x}}_{i+(2k-1)j}, \zeta \rangle^3 \bar{\mathbf{x}}_{i+(2k-1)j} \right) \quad (\text{E.79})$$

$$= \sum_{i=1}^{2k-1} \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3}. \quad (\text{E.80})$$

We conservatively bound the quantity of interest, $\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3}$, by ensuring that for each k , $\bar{\mathbf{X}}_k (\bar{\mathbf{X}}_k^T \boldsymbol{\zeta})^{\circ 3}$ be close to its expectation.

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \bar{\mathbf{g}}_E \right\|_2 \geq \frac{c}{2} \theta (1 - \theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2} \right] \\ & \leq \sum_{i=1}^{2k-1} \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \frac{\bar{\mathbf{g}}_E}{2k-1} \right\|_2 \geq \frac{c}{2} \theta (1 - \theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2 (2k-1)} \right] \\ & = \sum_{i=1}^{2k-1} \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \bar{\mathbf{g}}_E \right\|_2 \geq \frac{c}{2} \theta (1 - \theta) \frac{\|\boldsymbol{\zeta}\|_4^6}{\kappa^2 (2k-1)} \right] \end{aligned}$$

Applying Bernstein inequality for matrix variables as in [Lemma G.7](#), with $d_1 = 2k - 1$, $d_2 = 1$, we can obtain that for independent random vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ with

$$\sigma^2 = \sum_{i=1}^n \mathbb{E}[\|\mathbf{v}_i\|_2^2] \quad (\text{E.81})$$

and ensuring that

$$\|\mathbf{v}_i\|_2 \leq R \quad a.s. \quad (\text{E.82})$$

we obtain that

$$\mathbb{P} \left[\left\| \sum_i \mathbf{v}_i - \mathbb{E}[\cdot] \right\| > t \right] \leq 4k \exp \left(\frac{-t^2/2}{\sigma^2 + 2Rt/3} \right) \quad (\text{E.83})$$

Here, we have used that

$$\left\| \sum_{i=1}^n \mathbb{E}[\mathbf{v}_i \mathbf{v}_i^*] \right\| \leq \text{tr} \sum_{i=1}^n \mathbb{E}[\mathbf{v}_i \mathbf{v}_i^*] \quad (\text{E.84})$$

$$= \sum_{i=1}^n \mathbb{E}[\|\mathbf{v}_i\|_2^2]. \quad (\text{E.85})$$

and

$$\mathbf{w}_i = \bar{\mathbf{x}}_i \langle \bar{\mathbf{x}}_i, \boldsymbol{\zeta} \rangle^3. \quad (\text{E.86})$$

Notice that

$$\|\mathbf{w}_i\|_2 \leq \|\bar{\mathbf{x}}_i\|_2^4 \quad (\text{E.87})$$

$$\leq (4B^2 \theta k \log m)^2 \quad (\text{E.88})$$

$$= 16B^4 \theta^2 k^2 \log m. \quad (\text{E.89})$$

Let us further note that

$$\begin{aligned} & \sum_{\substack{j_1, \\ j_2 \neq j_3 \neq j_4}} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \boldsymbol{\zeta}_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \boldsymbol{\zeta}_{j_3}^2 \bar{\mathbf{x}}_i(j_4)^2 \boldsymbol{\zeta}_{j_4}^2] \\ & = 3 \sum_{j_1 \neq j_2 \neq j_3} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^4 \boldsymbol{\zeta}_{j_1}^2 \bar{\mathbf{x}}_i(j_2)^2 \boldsymbol{\zeta}_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \boldsymbol{\zeta}_{j_3}^2] \end{aligned}$$

$$\begin{aligned}
& + \sum_{j_1=1}^{2k-1} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2] \times \\
& \quad \sum_{j_1 \neq j_2 \neq j_3 \neq j_4} \mathbb{E} [\bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \bar{\mathbf{x}}_i(j_4)^2 \zeta_{j_4}^2]
\end{aligned} \tag{E.90}$$

$$\leq 2\theta k \times \theta^3 \|\zeta\|_2^6 + 3 \times 3\theta^3 \|\zeta\|_2^6 \tag{E.91}$$

In similar vein, we can obtain that

$$\begin{aligned}
& \sum_{j_1, j_2 \neq j_3} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^4 \zeta_{j_3}^4] \\
& = \sum_{j_1} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2] \sum_{j_2 \neq j_3 \neq j_1} \mathbb{E} [\bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^4 \zeta_{j_3}^4] \\
& \quad + \sum_{j_1 \neq j_2} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^4 \zeta_{j_1}^2 \bar{\mathbf{x}}_i(j_2)^4 \zeta_{j_2}^4] \\
& \quad + \sum_{j_1 \neq j_2} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \zeta_{j_1}^2 \bar{\mathbf{x}}_i(j_2)^6 \zeta_{j_2}^4]
\end{aligned} \tag{E.92}$$

$$\leq 2\theta k \times 3\theta^2 \|\zeta\|_2^2 \|\zeta\|_4^4 + (9 + 15) \theta^2 \|\zeta\|_2^2 \|\zeta\|_4^4 \tag{E.93}$$

and

$$\begin{aligned}
& \sum_{j_1, j_2} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^6 \zeta_{j_2}^6] \\
& = \sum_{j_1} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2] \sum_{j_2 \neq j_1} \mathbb{E} [\bar{\mathbf{x}}_i(j_2)^6 \zeta_{j_2}^6] \\
& \quad + \sum_{j_1} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^8 \zeta_{j_1}^6]
\end{aligned} \tag{E.94}$$

$$\leq 2\theta k \times 15\theta \|\zeta\|_6^6 + 105\theta \|\zeta\|_6^6 \tag{E.95}$$

Now we calculate

$$\mathbb{E} [\|\mathbf{w}_i\|_2^2] = \mathbb{E} [\|\bar{\mathbf{x}}_i\|_2^2 \langle \bar{\mathbf{x}}_i, \zeta \rangle^6] \tag{E.96}$$

$$= \mathbb{E} \left[\sum_{j_1, \dots, j_7} \bar{\mathbf{x}}_i(j_1)^2 \prod_{\ell=2}^7 \bar{\mathbf{x}}_i(j_\ell) \zeta_{j_\ell} \right] \tag{E.97}$$

$$\begin{aligned}
& = 15 \sum_{\substack{j_1, \\ j_2 \neq j_3 \neq j_4}} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \bar{\mathbf{x}}_i(j_4)^2 \zeta_{j_4}^2] \\
& \quad + 15 \sum_{j_1, j_2 \neq j_3} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^4 \zeta_{j_3}^4] \\
& \quad + \sum_{j_1, j_2} \mathbb{E} [\bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^6 \zeta_{j_2}^6]
\end{aligned} \tag{E.98}$$

$$\begin{aligned}
& \leq 15\theta^3 \|\zeta\|_2^6 (2\theta k + 9) \\
& \quad + 15\theta^2 \|\zeta\|_4^4 (6 + 24) \\
& \quad + \theta \|\zeta\|_6^6 (30\theta k + 105)
\end{aligned} \tag{E.99}$$

$$\leq 150\theta^2 k + 600\theta \tag{E.100}$$

whence for $\theta > 1/k$,

$$\mathbb{E} [\|\mathbf{w}_i\|_2^2] \leq C\theta^2 k, \quad (\text{E.101})$$

and hence

$$\sigma^2 \leq C'\theta^2 m. \quad (\text{E.102})$$

Matrix Bernstein gives that

$$\begin{aligned} & \mathbb{P} [\|\bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E} [\cdot]\|_2 \geq t] \\ & \leq 4k \exp \left(\frac{-t^2/2}{C\theta^2 m + C'B^4\theta^2 k^2 \log^2 kt} \right). \end{aligned} \quad (\text{E.103})$$

Setting $t = \frac{c}{4} \frac{m\theta(1-\theta)\|\boldsymbol{\zeta}\|_4^6}{\kappa^2(2k-1)}$, we obtain that

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E} [\cdot] \right\|_2 \geq \frac{c}{4} \frac{\theta(1-\theta)\|\boldsymbol{\zeta}\|_4^6}{\kappa^2(2k-1)} \right] \\ & \leq 4k \exp \left(- \frac{c'' m (1-\theta)^2 \|\boldsymbol{\zeta}\|_4^{12}}{\kappa^4 k^2 + \theta(1-\theta) B^4 \kappa^2 k^3 \|\boldsymbol{\zeta}\|_4^6} \right) \end{aligned} \quad (\text{E.104})$$

ε -Net Covering To obtain a probability bound for all $\mathbf{q} \in \mathbb{S}^{k-1}$, we choose a set of $\boldsymbol{\zeta}_n = \mathbf{A}^T \mathbf{q}_n$ with $n = 1, \dots, N$. Suppose for any $\mathbf{q} \in \mathbb{S}^{k-1}$, there exists \mathbf{q}_n such that $\|\mathbf{q} - \mathbf{q}_n\|_2 \leq \varepsilon$, then

$$\left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} \right\|_2 \leq L \|\mathbf{q} - \mathbf{q}_n\|_2. \quad (\text{E.105})$$

For entry wise bounded $\bar{\mathbf{X}}_i \in \mathbb{R}^{(2k-1) \times \frac{m}{2k-1}}$, we have

$$\|\bar{\mathbf{X}}_i\|_2 \leq \sqrt{2\theta m} B, \quad \|\bar{\mathbf{X}}_i \mathbf{e}_j\|_2 \leq \sqrt{4\theta k} B, \quad (\text{E.106})$$

then the Lipschitz constant L can be bounded as

$$L \leq \frac{1}{m} \|\bar{\mathbf{X}}_i\|_2 \left\| \text{diag} (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \right\|_2 \|\bar{\mathbf{X}}_i^T \mathbf{A}^T\|_2 \quad (\text{E.107})$$

$$\leq 8\theta^2 k B^4. \quad (\text{E.108})$$

With triangle inequality, we have

$$\begin{aligned} & \left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} \right] \right\|_2 \\ & \leq \left\| \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} \right] - \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} \right] \right\|_2 \\ & \quad + \left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} - \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} \right] \right\|_2 \\ & \quad + \left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} \right\|_2 \end{aligned} \quad (\text{E.109})$$

$$\begin{aligned} & \leq \left\| \frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} - \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3} \right] \right\|_2 \\ & \quad + 2L\varepsilon. \end{aligned} \quad (\text{E.110})$$

Hence we need to choose the ε -net to cover the sphere of \mathbf{q} with

$$\varepsilon = \frac{c}{4} \frac{\theta(1-\theta)}{\kappa^2(2k-1)} \frac{1}{L} \min_{\mathbf{q} \in \mathbb{S}^{k-1}} \|\boldsymbol{\zeta}\|_4^6, \quad (\text{E.111})$$

plug in $L \leq 4\theta^2 k B^4$ and number of sample N suffice

$$N \leq \left(\frac{3}{\varepsilon}\right)^k \quad (\text{E.112})$$

$$\leq \exp\left(k \ln\left(\frac{3}{\varepsilon}\right)\right) \quad (\text{E.113})$$

$$\leq \exp\left[k \ln\left(C \frac{\theta^2 \kappa^2 k^4 B^4}{\theta(1-\theta)}\right)\right] \quad (\text{E.114})$$

For $n = 1, \dots, N$, denote

$$P_i(\mathbf{q}_n) = \mathbb{P}\left[\left\|\frac{\bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 3}}{m} - \mathbb{E}[\cdot]\right\|_2 \geq \frac{c\theta(1-\theta)\|\boldsymbol{\zeta}_n\|_4^6}{4\kappa^2(2k-1)}\right], \quad (\text{E.115})$$

then together with union bound over all \mathbf{q}_n , we obtain that,

$$\begin{aligned} & \mathbb{P}\left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \frac{\left\|\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E}[\cdot]\right\|_2}{\|\boldsymbol{\zeta}\|_4^6} \geq \frac{c}{2} \frac{\theta(1-\theta)}{\kappa^2(2k-1)}\right] \\ & \leq \sum_{\mathbf{q}_n \in \hat{\mathcal{R}}_{2C_\star}} P_i(\mathbf{q}_n) \end{aligned} \quad (\text{E.116})$$

$$\leq N \max_{\mathbf{q}_n \in \hat{\mathcal{R}}_{2C_\star}} P_i(\mathbf{q}_n) \quad (\text{E.117})$$

$$\begin{aligned} & \leq 4k \sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \exp\left(-\frac{cm(1-\theta)^2 \|\boldsymbol{\zeta}\|_4^{12}}{\kappa^4 k^2 + \theta(1-\theta) B^4 \kappa^2 k^3 \|\boldsymbol{\zeta}\|_4^6}\right) \times \\ & \exp\left(k \ln\left(\frac{3}{\varepsilon}\right)\right). \end{aligned} \quad (\text{E.118})$$

Hence

$$\begin{aligned} & \mathbb{P}\left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \frac{\left\|\frac{1}{m} \bar{\mathbf{X}}_0 (\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E}[\cdot]\right\|_2}{\|\boldsymbol{\zeta}\|_4^6} \geq \frac{c}{2} \frac{\theta(1-\theta)}{\kappa^2}\right] \\ & \leq \sum_i \mathbb{P}\left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \frac{\left\|\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E}[\cdot]\right\|_2}{\|\boldsymbol{\zeta}\|_4^6} \geq \frac{c\theta(1-\theta)}{2\kappa^2(2k-1)}\right] \end{aligned} \quad (\text{E.119})$$

$$\begin{aligned} & \leq (2k-1) \max_i \mathbb{P}\left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \frac{\left\|\frac{1}{m} \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 3} - \mathbb{E}[\cdot]\right\|_2}{\|\boldsymbol{\zeta}\|_4^6} \geq \frac{c\theta(1-\theta)}{2\kappa^2(2k-1)}\right] \end{aligned} \quad (\text{E.120})$$

$$\begin{aligned} & \leq 8k^2 \sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}} \exp\left(-\frac{cm(1-\theta)^2 \|\boldsymbol{\zeta}\|_4^{12}}{\kappa^4 k^2 + \theta(1-\theta) B^4 \kappa^2 k^3 \|\boldsymbol{\zeta}\|_4^6}\right) \times \\ & \exp\left(k \ln\left(\frac{3}{\varepsilon}\right)\right), \end{aligned} \quad (\text{E.121})$$

which is bounded by $\exp(-k)$ as long as

$$m \geq C \frac{\min\left\{(2C_\star \mu)^{-2}, \kappa^2 k^2\right\}}{(1-\theta)^2} \kappa^2 k^4 \log^3(\kappa k) \quad (\text{E.122})$$

$$\begin{aligned} &\geq C' k \log \left(\frac{\theta \kappa^2 k^2 B^4}{(1-\theta) \|\zeta\|_4^6} \right) \times \\ &\max \left\{ \frac{\kappa^4 k^2}{(1-\theta)^2 \|\zeta\|_4^{12}}, \frac{\theta B^4 \kappa^2 k^3}{(1-\theta) \|\zeta\|_4^6} \right\}. \end{aligned} \quad (\text{E.123})$$

To sum up, we obtain that for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}$, inequality

$$\left\| \frac{1}{m} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} - \mathbb{E}[\cdot] \right\|_2 \leq c\theta (1-\theta) \frac{\|\mathbf{A}^T \mathbf{q}\|_4^6}{\kappa^2} \quad (\text{E.124})$$

holds with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4} - c_3 \exp(-\theta k)$. \blacksquare

F Concentration for Hessian (Lemma 4.3)

Lemma F.1. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$. There exists positive constant C that whenever

$$m \geq C\theta \frac{\min \left\{ (2C_\star \mu \kappa^2)^{-4/3}, k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^6 k^4 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \quad (\text{F.1})$$

and $\theta \geq \log k/k$, then with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4} - 48k^{-7} - 48m^{-5} - 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right)$,

$$\left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \leq c \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4, \quad (\text{F.2})$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}$ with positive constant $c \leq 0.048 \leq 3(1 - 6c_\star - 36c_\star^2 - 24c_\star^3)$.

Proof Denote $\boldsymbol{\eta} = \mathbf{Y}^T (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{q}$ and $\bar{\boldsymbol{\eta}} = \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q} = (\theta m)^{-1/2} \mathbf{X}_0^T \zeta$, and

$$\mathbf{W} = \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2}, \quad (\text{F.3})$$

$$\hat{\mathbf{Y}} = (\mathbf{Y} \mathbf{Y}^T)^{-1/2} \mathbf{Y}. \quad (\text{F.4})$$

Then we have

$$\begin{aligned} &\left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \\ &= \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{3}{m} \hat{\mathbf{Y}} \text{diag}(\boldsymbol{\eta}^{\circ 2}) \hat{\mathbf{Y}}^T - \langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp} \right. \\ &\quad \left. - \frac{3(1-\theta)}{\theta m^2} \mathbf{P}_{\mathbf{q}^\perp} \left[3\mathbf{A} \text{diag}(\zeta^{\circ 2}) \mathbf{A}^T - \|\zeta\|_4^4 \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp} \right\|_2 \end{aligned} \quad (\text{F.5})$$

$$\begin{aligned} &\leq \left\| \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{3}{m} \hat{\mathbf{Y}} \text{diag}(\boldsymbol{\eta}^{\circ 2}) \hat{\mathbf{Y}}^T \right] \mathbf{P}_{\mathbf{q}^\perp} \right. \\ &\quad \left. - \mathbf{P}_{\mathbf{q}^\perp} \left[\frac{9(1-\theta)}{\theta m^2} \mathbf{A} \text{diag}(\zeta^{\circ 2}) \mathbf{A}^T - \frac{3}{m^2} \mathbf{I} \right] \mathbf{P}_{\mathbf{q}^\perp} \right\|_2 \\ &\quad + \left\| \left[\langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle - \frac{3(1-\theta)}{\theta m^2} \|\zeta\|_4^4 - \frac{3}{m^2} \right] \mathbf{P}_{\mathbf{q}^\perp} \right\|_2 \end{aligned} \quad (\text{F.6})$$

$$\begin{aligned}
& \leq \underbrace{\frac{3}{\theta m^2} \left\| \mathbf{W} \mathbf{Y} \text{diag}(\boldsymbol{\eta}^{\circ 2}) \mathbf{Y}^T \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \right\|_2}_{\Delta_1^H} \\
& + \underbrace{\frac{3}{\theta m^2} \left\| \mathbf{A} \mathbf{X}_0 \text{diag}(\boldsymbol{\eta}^{\circ 2}) \mathbf{Y}^T \mathbf{W} \right\|_2}_{\Delta_2^H} \\
& + \underbrace{\frac{3}{\theta m^2} \left\| \mathbf{A} \mathbf{X}_0 \text{diag}(\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T \mathbf{A}^T \right\|_2}_{\Delta_3^H} \\
& + \underbrace{\frac{3}{\theta m^2} \left\| \mathbf{P}_{q^\perp} [\mathbf{A} \mathbf{X}_0 \text{diag}(\bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T \mathbf{A}^T] \mathbf{P}_{q^\perp} \right.}_{\Delta_4^H} \\
& \quad \left. - \mathbf{P}_{q^\perp} [3(1-\theta) \mathbf{A} \text{diag}(\boldsymbol{\zeta}^{\circ 2}) \mathbf{A}^T + \theta \mathbf{I}] \mathbf{P}_{q^\perp} \right\|_2}_{\Delta_4^H} \\
& + \underbrace{\left\| \left[\langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle - \frac{3(1-\theta)}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4 - \frac{3}{m^2} \right] \mathbf{P}_{q^\perp} \right\|_2}_{\Delta_5^H}
\end{aligned} \tag{F.7}$$

In the rest of the proof, we prove that

$$\Delta_i^H \leq \frac{c}{9} \frac{1-\theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4, \quad i = 1, 2, 3. \tag{F.8}$$

and

$$\Delta_i^H \leq \frac{c}{3} \frac{1-\theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4, \quad i = 4, 5. \tag{F.9}$$

First, let us note that

$$C(1-\theta)^{-2} \sigma_{\min}^{-2} \kappa^6 k^5 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \tag{F.10}$$

$$\leq C \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)^6 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \tag{F.11}$$

$$\leq C \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)^9 \tag{F.12}$$

or

$$\begin{aligned}
& \frac{\log^3 \left(C(1-\theta)^{-2} \sigma_{\min}^{-2} \kappa^6 k^5 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right) \right)}{C \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)} \\
& \leq \left(\frac{\log C + 9 \log \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)}{C^{1/3} \log \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)} \right)^3
\end{aligned} \tag{F.13}$$

$$\leq \left(\frac{\log C}{C^{1/3} \log \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right)} + \frac{9}{C^{1/3}} \right)^3 \tag{F.14}$$

$$\leq \left(\frac{1}{C^{1/6}} + \frac{1}{2} \frac{1}{C^{1/6}} \right)^3 \quad (C \geq 10^8) \quad (\text{F.15})$$

$$\leq \frac{4}{C^{1/2}}. \quad (\text{F.16})$$

Since

$$m \geq C \frac{\min \left\{ (2C_* \mu \kappa^2)^{-4/3}, k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^6 k^4 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right), \quad (\text{F.17})$$

as the ratio $\log^3 m/m$ decreases with increasing m , then

$$\begin{aligned} & \frac{\log^3 m}{m} \\ & \leq \frac{\log^3 \left(C \frac{\kappa^6 k^5}{(1-\theta)^2 \sigma_{\min}^2} \log^3 \left(\frac{\kappa k}{\sigma_{\min}(1-\theta)} \right) \right)}{C \log^3 \left(\frac{\kappa k}{\sigma_{\min}(1-\theta)} \right)} \\ & \quad \times \frac{(1-\theta)^2 \sigma_{\min}^2}{\min \left\{ (2C_* \mu \kappa^2)^{-2/3}, k \right\} \kappa^6 k^4} \end{aligned} \quad (\text{F.18})$$

$$\leq \frac{4}{C^{1/2}} \frac{(1-\theta)^2 \sigma_{\min}^2}{\min \left\{ (2C_* \mu \kappa^2)^{-2/3}, k \right\} \kappa^6 k^4} \quad (\text{F.19})$$

According to [Lemma D.1](#), following inequality obtains

$$\left\| \frac{1}{\theta m} \mathbf{X}_0 \mathbf{X}_0^T - \mathbf{I} \right\|_2 \leq \delta \quad (\text{F.20})$$

$$\leq 10 \sqrt{k \log m/m} \quad (\text{F.21})$$

$$\leq \frac{20(1-\theta) \sigma_{\min} \max \left\{ (2C_* \mu \kappa^2)^{2/3}, k^{-1} \right\}}{C^{1/4} \kappa^3 k^{3/2} \log m} \quad (\text{F.22})$$

$$\leq \frac{20 \sigma_{\min}}{C^{1/4} \kappa^3} \cdot \frac{(1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^4}{k^{3/2} \log m}, \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*} \quad (\text{F.23})$$

with probability no smaller than $1 - \varepsilon_0$ with $\varepsilon_0 = 2 \exp(-\theta k) + 24k \exp\left(-\frac{1}{144} \min \left\{ k, 3\sqrt{\theta m} \right\}\right) + 48k^{-7} + 48m^{-5}$.

We have $4\kappa^3 \delta / \sigma_{\min} \leq 1/2$ whenever

$$C \geq \left(\frac{160(1-\theta)}{k^{3/2} \log m} \right)^4 \quad (\text{F.24})$$

whence $\delta \leq 1/(8\kappa^2)$, and [Lemma D.3](#) implies that

$$\begin{aligned} & \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 - (\mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{A}_0 \right\|_2 \\ & \leq 4\kappa^3 \delta / \sigma_{\min} \end{aligned} \quad (\text{F.25})$$

$$\leq \frac{80(1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^4}{C^{1/4} k^{3/2} \log m}, \quad \forall \mathbf{q} \in \hat{\mathcal{R}}_{2C_*}. \quad (\text{F.26})$$

Moreover,

$$\|\mathbf{X}_0\|_2 \leq (\theta m)^{1/2} \sqrt{1 + \delta} \quad (\text{F.27})$$

$$\leq (\theta m)^{1/2} (1 + \delta/2) \quad (\text{F.28})$$

$$\leq \frac{17}{16} (\theta m)^{1/2}. \quad (\text{F.29})$$

Finally, [Lemma A.5](#) implies that with probability no smaller than $1 - \varepsilon_B$, we have

$$\|\mathbf{x}_0\|_\infty \leq \sqrt{2} \log^{1/2} \left(\frac{2\theta m}{\varepsilon_B} \right). \quad (\text{F.30})$$

Upper Bound for Δ_1^H and Δ_2^H . With probability no smaller than $1 - \varepsilon_0 - \varepsilon_B$, the norms of $\boldsymbol{\eta}$ are upper bounded as in [Lemma A.7](#),

$$\begin{aligned} \Delta_1^H &\leq \frac{3}{\theta m^2} \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 - \mathbf{A} \right\|_2 \times \\ &\quad \|\mathbf{X}_0\|_2^2 \|\boldsymbol{\eta}\|_\infty^2 \left\| \mathbf{A}_0^T \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \right\|_2 \end{aligned} \quad (\text{F.31})$$

$$\begin{aligned} &\leq \frac{3}{\theta m^2} \cdot \frac{4\kappa^3 \delta}{\sigma_{\min}} \cdot (1 + \delta/2)^2 \theta m \cdot \\ &\quad \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right)^3 \frac{4k}{\theta m} \log(2\theta m / \varepsilon_B) \end{aligned} \quad (\text{F.32})$$

$$\leq \frac{3660}{C^{1/4}} \frac{1 - \theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4 \cdot \frac{\log(2\theta m / \varepsilon_B)}{k^{1/2} \log m}. \quad (\text{F.33})$$

A similar result holds for

$$\begin{aligned} \Delta_2^H &\leq \frac{3}{\theta m^2} \|\mathbf{X}_0\|_2^2 \|\text{diag}(\boldsymbol{\eta}^{\circ 2})\|_2 \times \\ &\quad \left\| \left(\frac{1}{\theta m} \mathbf{Y} \mathbf{Y}^T \right)^{-1/2} \mathbf{A}_0 - \mathbf{A} \right\|_2 \end{aligned} \quad (\text{F.34})$$

$$\leq \frac{2440}{C^{1/4}} \frac{1 - \theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4 \cdot \frac{\log(2\theta m / \varepsilon_B)}{k^{1/2} \log m}. \quad (\text{F.35})$$

To make $\Delta_1^H \leq \frac{c}{9} \frac{1 - \theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4$ and $\Delta_2^H \leq \frac{c}{9} \frac{1 - \theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4$, we require

$$C \geq \left(9 \times 3660 c^{-1} \frac{\log(2\theta m / \varepsilon_B)}{k^{1/2} \log m} \right)^4. \quad (\text{F.36})$$

The right hand side is bounded by an absolute constant for all m .

Upper Bound for Δ_3^H . With probability no smaller than $1 - \varepsilon_0 - \varepsilon_B$, the difference between $\bar{\boldsymbol{\eta}}^{\circ 2}$ and $\boldsymbol{\eta}^{\circ 2}$ is upper bounded as in [Lemma A.7](#),

$$\begin{aligned} &\|\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2}\|_\infty \\ &\leq \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_\infty \|\boldsymbol{\eta} + \bar{\boldsymbol{\eta}}\|_\infty \end{aligned} \quad (\text{F.37})$$

$$\leq \frac{4\kappa^3 \delta}{\sigma_{\min}} \left(2 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right) \frac{2k}{\theta m} \log(2\theta m / \varepsilon_B) \quad (\text{F.38})$$

$$\leq \frac{5k}{\theta m} \log(2\theta m/\varepsilon_B) \cdot \frac{4\kappa^3 \delta}{\sigma_{\min}}. \quad (\text{F.39})$$

Therefore

$$\Delta_3^H = \frac{3}{\theta m^2} \|\mathbf{A} \mathbf{X}_0 \text{diag}(\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T \mathbf{A}^T\|_2 \quad (\text{F.40})$$

$$\leq \frac{3}{\theta m^2} \|\mathbf{A}\|_2^2 \|\mathbf{X}_0\|_2^2 \|\text{diag}(\boldsymbol{\eta}^{\circ 2} - \bar{\boldsymbol{\eta}}^{\circ 2})\|_2 \quad (\text{F.41})$$

$$\leq \frac{15k}{\theta m^2} (1 + \delta/2)^2 \log(2\theta m/\varepsilon_B) \cdot \frac{4\kappa^3 \delta}{\sigma_{\min}} \quad (\text{F.42})$$

$$\leq \frac{1400(1-\theta) \log(2\theta m/\varepsilon_B)}{C^{1/4} \theta k^{1/2} m^2 \log m} \|\boldsymbol{\zeta}\|_4^4. \quad (\text{F.43})$$

Again, Δ_3^H is bounded by $\frac{c}{9} \frac{1-\theta}{\theta m^2} \|\boldsymbol{\zeta}\|_4^4$ whenever

$$C \geq \left(9 \times 1400 c^{-1} \frac{\log(2\theta m/\varepsilon_B)}{k^{1/2} \log m} \right)^4 \quad (\text{F.44})$$

Upper Bound for Δ_4^H . Recall that

$$\bar{\boldsymbol{\eta}} = \mathbf{Y}^T (\theta m \mathbf{A}_0 \mathbf{A}_0^T)^{-1/2} \mathbf{q}, \quad (\text{F.45})$$

then

$$\begin{aligned} & \mathbb{E} [\mathbf{X}_0 \text{diag}(\bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T] \\ &= \mathbb{E} \left[\frac{1}{\theta m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \mathbf{X}_0^T \right] \end{aligned} \quad (\text{F.46})$$

$$\begin{aligned} &= 3(1-\theta) \text{diag}(\mathbf{A}^T \mathbf{q})^{\circ 2} + 2\theta \mathbf{A}^T \mathbf{q} \mathbf{q}^T \mathbf{A} \\ &\quad + \theta \|\mathbf{A}^T \mathbf{q}\|_2^2 \mathbf{I}, \end{aligned} \quad (\text{F.47})$$

once including the projection $\mathbf{P}_{\mathbf{q}^\perp}$, we have

$$\begin{aligned} & \mathbf{P}_{\mathbf{q}^\perp} \mathbb{E} [\mathbf{A} \mathbf{X}_0 \text{diag}(\bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T \mathbf{A}^T] \mathbf{P}_{\mathbf{q}^\perp} \\ &= \mathbf{P}_{\mathbf{q}^\perp} [3(1-\theta) \mathbf{A} \text{diag}(\boldsymbol{\zeta}^{\circ 2}) \mathbf{A}^T + \theta \mathbf{I}] \mathbf{P}_{\mathbf{q}^\perp}. \end{aligned} \quad (\text{F.48})$$

Therefore

$$\begin{aligned} \Delta_4^H &= \frac{3}{\theta m^2} \left\| \mathbf{P}_{\mathbf{q}^\perp} [\mathbf{A} \mathbf{X}_0 \text{diag}(\bar{\boldsymbol{\eta}}^{\circ 2}) \mathbf{X}_0^T \mathbf{A}^T] \mathbf{P}_{\mathbf{q}^\perp} \right. \\ &\quad \left. - \mathbf{P}_{\mathbf{q}^\perp} [3(1-\theta) \mathbf{A} \text{diag}(\boldsymbol{\zeta}^{\circ 2}) \mathbf{A}^T + \theta \mathbf{I}] \mathbf{P}_{\mathbf{q}^\perp} \right\|_2 \end{aligned} \quad (\text{F.49})$$

$$\leq \frac{3}{\theta^2 m^2} \left\| \frac{1}{m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 2} \mathbf{X}_0^T - \mathbb{E}[\cdot] \right\|_2 \quad (\text{F.50})$$

Under the assumption for sample size that $m \geq C(1-\theta)^{-2} \kappa^4 \min\left\{(2C_\star \mu)^{-2/3}, k\right\} k^3 \log^5(\kappa k)$, applying [Lemma F.2](#), we have

$$\left\| \frac{1}{m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 2} \mathbf{X}_0^T - \mathbb{E}[\cdot] \right\|_2 \leq \frac{c}{9} \theta (1-\theta) \|\boldsymbol{\zeta}\|_4^4. \quad (\text{F.51})$$

simultaneously at every $\mathbf{q} \in \hat{\mathcal{R}}_{2C_\star}$ with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4}$.

Upper Bound for Δ_5^H . Note that this term is essentially the difference between

$$\Delta_5^H = \left\| \left[\langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle - \frac{3(1-\theta)}{\theta m^2} \|\zeta\|_4^4 - \frac{3}{m^2} \right] \mathbf{P}_{\mathbf{q}^\perp} \right\|_2 \quad (\text{F.52})$$

$$\leq \left| \langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle - \frac{3(1-\theta)}{\theta m^2} \|\zeta\|_4^4 - \frac{3}{m^2} \right| \quad (\text{F.53})$$

$$\begin{aligned} &\leq \frac{1}{\theta^2 m^2} \left| \frac{1}{m} \|\mathbf{X}_0^T \zeta\|_4^4 - 3\theta(1-\theta) \|\zeta\|_4^4 - 3\theta^2 \right| \\ &\quad + \left| \langle \mathbf{q}, \nabla \psi(\mathbf{q}) \rangle - \frac{1}{\theta^2 m^2} \|\mathbf{X}_0^T \zeta\|_4^4 \right| \end{aligned} \quad (\text{F.54})$$

$$\begin{aligned} &\leq \frac{1}{\theta^2 m^2} \left\| \frac{\mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \zeta)^{\circ 3}}{m} - 3\theta(1-\theta) \mathbf{A}^T \zeta^{\circ 3} - 3\theta^2 \mathbf{q} \right\|_2 \\ &\quad + \frac{1}{m} \left| \|\boldsymbol{\eta}\|_4^4 - \|\bar{\boldsymbol{\eta}}\|_4^4 \right| \end{aligned} \quad (\text{F.55})$$

Recall that

$$\begin{aligned} &\mathbb{E} \left[\frac{1}{m} \mathbf{A} \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 3} \right] \\ &= \mathbb{E} \left[\mathbf{A} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{A}^T \mathbf{q})^3 \right] \end{aligned} \quad (\text{F.56})$$

$$= 3\theta(1-\theta) \mathbf{A} \zeta^{\circ 3} + 3\theta^2 \mathbf{q}, \quad (\text{F.57})$$

With similar argument as in [Lemma 4.2](#), we can show that this term can be bounded by $\frac{c}{6} \frac{1-\theta}{\theta m^2} \|\boldsymbol{\eta}\|_4^4$ whenever

$$m \geq C' \frac{\min \left\{ (\mu \kappa^2)^{-4/3}, k^2 \right\}}{(1-\theta)^2 \sigma_{\min}^2} \kappa^6 k^4 \log^3 \left(\frac{\kappa k}{(1-\theta) \sigma_{\min}} \right). \quad (\text{F.58})$$

Moreover, with probability $1 - \varepsilon_0 - \varepsilon_B$

$$\begin{aligned} &\frac{1}{m} \left| \|\boldsymbol{\eta}\|_4^4 - \|\bar{\boldsymbol{\eta}}\|_4^4 \right| \\ &\leq \frac{1}{m} |\langle \boldsymbol{\eta} - \bar{\boldsymbol{\eta}}, 4\boldsymbol{\eta}^{\circ 3} \rangle| \end{aligned} \quad (\text{F.59})$$

$$\leq \frac{4}{m} \|\boldsymbol{\eta} - \bar{\boldsymbol{\eta}}\|_2 \|\boldsymbol{\eta}\|_6^3 \quad (\text{F.60})$$

$$\leq \frac{16\kappa^3 \delta}{\sigma_{\min} m} (1 + \delta/2) \left(1 + \frac{4\kappa^3 \delta}{\sigma_{\min}} \right)^2 \frac{4k}{\theta m} \log(2\theta m / \varepsilon_B) \quad (\text{F.61})$$

$$\leq \frac{153k}{\theta m^2} \log(2\theta m / \varepsilon_B) \cdot \frac{\kappa^3 \delta}{\sigma_{\min}} \quad (\text{F.62})$$

$$\leq \frac{3060}{C^{1/4}} \frac{(1-\theta)}{\theta m^2} \|\zeta\|_4^4 \cdot \frac{\log(2\theta m / \varepsilon_B)}{k^{1/2} \log m}, \quad (\text{F.63})$$

which is bounded by $\frac{c}{6} \frac{1-\theta}{\theta m^2} \|\zeta\|_4^4$ whenever

$$C \geq \left(6 \times 3060 c^{-1} \frac{(1-\theta) \log(2\theta m / \varepsilon_B)}{k^{1/2} \log m} \right)^4. \quad (\text{F.64})$$

The right hand side is bounded by an absolute constant for all m .

Adding up failure probabilities, we have that with probability larger than $1 - c_2 \exp(-k) - c_2 k^{-4} - \varepsilon_0$,

$$\left\| \text{Hess}[\psi](\mathbf{q}) - \frac{3(1-\theta)}{\theta m^2} \text{Hess}[\varphi](\mathbf{q}) \right\|_2 \leq c \frac{1-\theta}{\theta m^2} \|\mathbf{A}^T \mathbf{q}\|_4^4 \quad (\text{F.65})$$

holds as desired for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$, where $\varepsilon_0 = 2 \exp(-\theta k) + 24k \exp\left(-\frac{1}{144} \min\{k, 3\sqrt{\theta m}\}\right) + 48k^{-7} + 48m^{-5}$. \blacksquare

F.1 Proof of Lemma F.2

Lemma F.2. Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$. There exist constants $C > 0$ that whenever

$$m \geq C \frac{\min\left\{(2C_* \mu \kappa^2)^{-4/3}, k^2\right\}}{(1-\theta)^2} k^4 \log^3(\kappa k), \quad (\text{F.66})$$

and $\theta k > 1$, then with probability no smaller than $1 - c_1 \exp(-k) - c_2 k^{-4}$,

$$\begin{aligned} & \left\| \frac{1}{m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \mathbf{X}_0^T - \mathbb{E}[\cdot] \right\|_2 \\ & \leq c\theta(1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^4, \end{aligned} \quad (\text{F.67})$$

holds for all $\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}$ with positive constant $c \leq 0.005 \leq (1 - 6c_* - 36c_*^2 - 24c_*^3)/3$.

Proof The proof strategy for the finite sample concentration of the Hessian is similar to that of the gradient as presented in Lemma E.2. For simplicity, we will only demonstrate some key steps here, please refer to Lemma E.2 for detailed arguments.

Again, from Lemma A.5, the coefficient satisfies $\|\mathbf{x}_0\|_\infty \leq B$ with probability no smaller than $1 - 2\theta m e^{-B^2/2}$. We write $\bar{\mathbf{x}}_0(i) = \mathbf{x}_0(i) \mathbb{1}_{|\mathbf{x}_0(i)| \leq B}$, and let $\bar{\mathbf{X}}_0$ denote the circulant matrix generated by the truncated vector $\bar{\mathbf{x}}_0$. Denote

$$\mathbf{H}_E = \mathbb{E} \left[\frac{1}{m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \mathbf{X}_0^T \right], \quad (\text{F.68})$$

$$\bar{\mathbf{H}}_E = \mathbb{E} \left[\frac{1}{m} \bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \bar{\mathbf{X}}_0^T \right], \quad (\text{F.69})$$

then

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{\mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \boldsymbol{\zeta})^{\circ 2} \mathbf{X}_0^T}{m} - \mathbf{H}_E \right\|_2 \geq c\theta(1-\theta) \|\boldsymbol{\zeta}\|_4^4 \right] \\ & \leq \mathbb{P} \left[\left\| \frac{\bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_0^T}{m} - \bar{\mathbf{H}}_E \right\|_2 \geq c\theta(1-\theta) \|\boldsymbol{\zeta}\|_4^4 \right] \\ & \quad + 2\theta m e^{-B^2/2} + 2m \exp\left(-\frac{3}{4} \theta k \log m\right) \end{aligned} \quad (\text{F.70})$$

while via triangle inequality,

$$\begin{aligned} & \left\| \frac{1}{m} \bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \bar{\mathbf{X}}_0^T - \mathbf{H}_E \right\|_2 \\ & \leq \left\| \frac{1}{m} \bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \bar{\mathbf{X}}_0^T - \bar{\mathbf{H}}_E \right\|_2 \end{aligned}$$

$$+ \|\bar{\mathbf{H}}_E - \mathbf{H}_E\|_2. \quad (\text{F.71})$$

Truncation Level. Next, we choose a large enough entry-wise truncation level B such that the expectation of the Hessian $\mathbb{E} \left[\mathbf{X}_0 \text{diag} (\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \mathbf{X}_0^T \right]$ is close to that of its truncation $\mathbb{E} \left[\bar{\mathbf{X}}_0 \text{diag} (\bar{\mathbf{X}}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \bar{\mathbf{X}}_0^T \right]$. Moreover, we introduce following events notation

$$\mathcal{E}_i \doteq \{\|\mathbf{x}_i\|_\infty > B \cup \|\mathbf{x}_i\|_0 > 4\theta k \log m\}, \quad (\text{F.72})$$

then

$$\begin{aligned} & \|\bar{\mathbf{H}}_E - \mathbf{H}_E\|_2 \\ &= \left\| \mathbb{E} \left[\frac{1}{m} \sum_i \langle \mathbf{x}_i, \boldsymbol{\zeta} \rangle^2 \mathbf{x}_i \mathbf{x}_i^T \cdot \mathbf{1}_{\mathbf{E}_i} \right] \right\|_F \end{aligned} \quad (\text{F.73})$$

$$\leq \frac{1}{m} \sum_i \left\| \mathbb{E} \left[\langle \mathbf{x}_i, \boldsymbol{\zeta} \rangle^2 \mathbf{x}_i \mathbf{x}_i^T \cdot \mathbf{1}_{\mathbf{E}_i} \right] \right\|_F \quad (\text{F.74})$$

$$\leq \frac{1}{m} \sum_i \left(\mathbb{E} \left[\left\| \langle \mathbf{x}_i, \boldsymbol{\zeta} \rangle^2 \mathbf{x}_i \mathbf{x}_i^T \right\|_F^2 \right] \cdot \mathbb{E} [\mathbf{1}_{\mathbf{E}_i}] \right)^{1/2} \quad (\text{F.75})$$

$$\leq \left(\mathbb{E} [\|\mathbf{x}_i\|_2^8] \right)^{1/2} \times \sqrt{\mathbb{E} [\mathbf{1}_{\|\mathbf{x}_i\|_\infty > B}] + \mathbb{E} [\mathbf{1}_{\|\mathbf{x}_i\|_0 > 4\theta k \log m}]} \quad (\text{F.76})$$

$$\leq 50k^2 \sqrt{4\theta k e^{-B^2/2} + \exp(-\frac{3}{4}\theta k \log m)} \quad (\text{F.77})$$

By setting

$$B \geq C' \log^{1/2} \left(\frac{k^7}{\theta (1-\theta)^2} \right) \quad (\text{F.78})$$

we have

$$\theta k e^{-B^2/2} \leq c' \theta^2 (1-\theta)^2 \frac{\|\boldsymbol{\zeta}\|_4^8}{k^4} \quad (\text{F.79})$$

In addition, whenever

$$\theta k \geq \frac{4}{3 \log m} \log \left(\frac{400^2 k^4}{c^2 \theta^2 (1-\theta)^2 \|\boldsymbol{\zeta}\|_4^8} \right), \quad (\text{F.80})$$

we have

$$\exp(-\frac{3}{4}\theta k \log m) \leq \frac{1}{2} \left(\frac{c\theta(1-\theta)}{100} \right)^2 \frac{\|\boldsymbol{\zeta}\|_4^8}{k^4}. \quad (\text{F.81})$$

Hence,

$$\sqrt{4\theta k e^{-B^2/2} + \exp(-\frac{3}{4}\theta k \log m)} \leq \frac{c\theta(1-\theta)}{100k^2} \|\boldsymbol{\zeta}\|_4^4. \quad (\text{F.82})$$

Therefore, we can obtain that

$$\|\bar{\mathbf{H}}_E - \mathbf{H}_E\|_2 \leq \frac{c}{2} \theta (1-\theta) \|\boldsymbol{\zeta}\|_4^4 \quad (\text{F.83})$$

always holds, hence

$$\begin{aligned} & \mathbb{P} \left[\left\| \frac{\bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_0^T}{m} - \mathbf{H}_E \right\|_2 \geq c\theta(1-\theta) \|\boldsymbol{\zeta}\|_4^4 \right] \\ & \leq \mathbb{P} \left[\left\| \frac{\bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_0^T}{m} - \bar{\mathbf{H}}_E \right\|_2 \geq \frac{c}{2} \theta(1-\theta) \|\boldsymbol{\zeta}\|_4^4 \right]. \end{aligned} \quad (\text{F.84})$$

Independent Sub-matrices. As we did in [Lemma E.2](#), we remove the dependence in \mathbf{X}_0 by sampling every $2k-1$ column such that

$$\mathbf{X}_0 \boldsymbol{\Pi} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_{2k-1}], \quad (\text{F.85})$$

where

$$\mathbf{X}_i = [\mathbf{x}_i, \mathbf{x}_{i+(2k-1)}, \dots, \mathbf{x}_{i+(m-2k-1)}], \quad (\text{F.86})$$

and $\boldsymbol{\Pi}$ is a certain permutation of the columns of \mathbf{X}_0 .

Applying Bernstein inequality for matrix variables as in [Lemma G.7](#), with $\mathbf{M}_i = \langle \bar{\mathbf{x}}_i, \mathbf{A}^T \mathbf{q} \rangle^2 \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \in \mathbb{R}^{(2k-1) \times (2k-1)}$. Since

$$\|\mathbf{M}_i\|_2 = \left\| \langle \bar{\mathbf{x}}_i, \mathbf{A}^T \mathbf{q} \rangle^2 \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right\|_2 \quad (\text{F.87})$$

$$\leq \|\bar{\mathbf{x}}_i\|_2^4 \quad (\text{F.88})$$

$$\leq 4B^4 k^2 \quad (\text{F.89})$$

and

$$\|\mathbb{E}[\mathbf{M}_i \mathbf{M}_i^*]\| = \|\mathbb{E}[\mathbf{M}_i^* \mathbf{M}_i]\| \quad (\text{F.90})$$

$$= \left\| \mathbb{E} \left[\langle \bar{\mathbf{x}}_i, \mathbf{A}^T \mathbf{q} \rangle^4 \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right] \right\| \quad (\text{F.91})$$

$$= \left\| \mathbb{E} \left[\langle \bar{\mathbf{x}}_i, \boldsymbol{\zeta} \rangle^4 \|\bar{\mathbf{x}}_i\|_2^2 \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \right] \right\| \quad (\text{F.92})$$

$$\leq \mathbb{E} \left[\langle \bar{\mathbf{x}}_i, \boldsymbol{\zeta} \rangle^4 \|\bar{\mathbf{x}}_i\|_2^4 \right], \quad (\text{F.93})$$

we obtain the following upper bound:

$$\begin{aligned} & \mathbb{E} \left[\langle \bar{\mathbf{x}}_i, \boldsymbol{\zeta} \rangle^4 \|\bar{\mathbf{x}}_i\|_2^4 \right] \\ & = \mathbb{E} \left[\sum_{j_1, j_2}^{2k-1} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \sum_{j_3, \dots, j_6} \prod_{\ell=3}^6 \bar{\mathbf{x}}_i(j_\ell) \zeta_{j_\ell} \right] \end{aligned} \quad (\text{F.94})$$

$$\begin{aligned} & = 3 \mathbb{E} \left[\sum_{j_1, j_2}^{2k-1} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \sum_{j_3 \neq j_4} \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \bar{\mathbf{x}}_i(j_4)^2 \zeta_{j_4}^2 \right] \\ & \quad + \mathbb{E} \left[\sum_{j_1, j_2}^{2k-1} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \cdot \sum_{j_3} \bar{\mathbf{x}}_i(j_3)^4 \zeta_{j_3}^4 \right] \end{aligned} \quad (\text{F.95})$$

$$= 3 \mathbb{E} \left[\sum_{\substack{j_1 \neq j_2 \\ \neq j_3 \neq j_4}} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \bar{\mathbf{x}}_i(j_4)^2 \zeta_{j_4}^2 \right]$$

$$\begin{aligned}
& + 3\mathbb{E} \left[\sum_{j_1 \neq j_2 \neq j_3} \bar{\mathbf{x}}_i(j_1)^4 \bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \right] \\
& + 6\mathbb{E} \left[\sum_{j_1 \neq j_2} \bar{\mathbf{x}}_i(j_1)^6 \zeta_{j_1}^2 \bar{\mathbf{x}}_i(j_2)^2 \zeta_{j_2}^2 \right] \\
& + 6\mathbb{E} \left[\sum_{j_1 \neq j_2 \neq j_3} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^4 \zeta_{j_2}^2 \bar{\mathbf{x}}_i(j_3)^2 \zeta_{j_3}^2 \right] \\
& + 6\mathbb{E} \left[\sum_{j_1 \neq j_2} \bar{\mathbf{x}}_i(j_1)^4 \zeta_{j_1}^2 \bar{\mathbf{x}}_i(j_2)^4 \zeta_{j_2}^2 \right] \\
& + 2\mathbb{E} \left[\sum_{j_1 \neq j_2} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^6 \zeta_{j_2}^4 \right] \\
& + \mathbb{E} \left[\sum_{j_1 \neq j_2 \neq j_3} \bar{\mathbf{x}}_i(j_1)^2 \bar{\mathbf{x}}_i(j_2)^2 \bar{\mathbf{x}}_i(j_3)^4 \zeta_{j_3}^4 \right] \\
& + \mathbb{E} \left[\sum_{j_1 \neq j_2} \bar{\mathbf{x}}_i(j_1)^4 \bar{\mathbf{x}}_i(j_2)^4 \zeta_{j_2}^4 \right] \\
& + \mathbb{E} \left[\sum_j \bar{\mathbf{x}}_i(j)^8 \zeta_j^4 \right] \tag{F.96} \\
& \leq (105\theta + 18\theta^2 k + 60\theta^2 k + 12\theta^3 k^2) \|\zeta\|_4^4 \\
& \quad + 3(21\theta^2 + 30\theta^2 + 4\theta^4 k^2 + 12\theta^2 k) \|\zeta\|_2^4 \tag{F.97} \\
& \leq C\theta^3 k^2 \tag{F.98}
\end{aligned}$$

Assuming $\theta m \geq 1$, hence

$$\sigma^2 = C\theta^3 km. \tag{F.99}$$

Setting $t = \frac{c}{2} \frac{\theta(1-\theta)m\|\zeta\|_4^4}{2k-1}$ in Matrix Bernstein gives that

$$\begin{aligned}
& \mathbb{P} \left[\left\| \bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \zeta)^{\circ 3} - \mathbb{E}[\cdot] \right\|_2 > t \right] \\
& \leq 8k \exp \left(\frac{-t^2/2}{C\theta^3 km + C'B^4\theta^2 k^2 t} \right), \tag{F.100}
\end{aligned}$$

we obtain that

$$\begin{aligned}
& \mathbb{P} \left[\left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \zeta)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2 > c \frac{\theta(1-\theta)\|\zeta\|_4^6}{2k-1} \right] \\
& \leq 8k \exp \left(-\frac{cm(1-\theta)^2 \|\zeta\|_4^8}{\theta k^3 + \theta(1-\theta)B^4 k^3 \|\zeta\|_4^4} \right). \tag{F.101}
\end{aligned}$$

ε -Net Covering To obtain a probability bound for all $\mathbf{q} \in \mathbb{S}^{k-1}$, we choose a set of $\zeta_n = \mathbf{A}^T \mathbf{q}_n$ with $n =$

$1, \dots, N$. Since for any $\mathbf{q}, \mathbf{q}' \in \mathbb{S}^{k-1}$ and $\boldsymbol{\zeta}' = \mathbf{A}^T \mathbf{q}'$, we have

$$\begin{aligned} & \left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}')^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right\|_2 \\ &= \frac{1}{m} \left\| \bar{\mathbf{X}}_i \text{diag} \left[(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} - (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}')^{\circ 2} \right] \bar{\mathbf{X}}_i^T \right\|_2 \end{aligned} \quad (\text{F.102})$$

$$\leq \frac{\|\bar{\mathbf{X}}_i\|_2^2}{m} \left\| \text{diag} \left[(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} - (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}')^{\circ 2} \right] \right\|_2 \quad (\text{F.103})$$

$$\leq \frac{\|\bar{\mathbf{X}}_i\|_2^2}{m} \|\bar{\mathbf{X}}_i^T \boldsymbol{\zeta} + \bar{\mathbf{X}}_i^T \boldsymbol{\zeta}'\|_\infty \|\bar{\mathbf{X}}_i^T \boldsymbol{\zeta} - \bar{\mathbf{X}}_i^T \boldsymbol{\zeta}'\|_\infty \quad (\text{F.104})$$

$$\leq L \|\mathbf{q} - \mathbf{q}'\|_2 \quad (\text{F.105})$$

Then the Lipschitz constant L is upper bounded by

$$L \leq \frac{\|\bar{\mathbf{X}}_i\|_2^2}{m} \|\bar{\mathbf{X}}_i^T \mathbf{A}^T\|_2 (\|\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}\|_\infty + \|\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}'\|_\infty) \quad (\text{F.106})$$

$$\leq \frac{2}{m} \|\bar{\mathbf{X}}_i\|_2^4 \quad (\text{F.107})$$

$$\leq 8\theta^2 m B^4. \quad (\text{F.108})$$

With triangle inequality, we have

$$\begin{aligned} & \left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2 \\ & \leq \left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right\|_2 \\ & \quad + \left\| \mathbb{E} \left[\frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right] - \mathbb{E} \left[\frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right] \right\|_2 \\ & \quad + \left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E} \left[\frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right] \right\|_2 \end{aligned} \quad (\text{F.109})$$

$$\begin{aligned} & \leq \left\| \frac{\bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E} \left[\frac{\bar{\mathbf{X}}_i (\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} \right] \right\|_2 \\ & \quad + 2L\varepsilon \end{aligned} \quad (\text{F.110})$$

Next, we are going to choose the ε -net to cover the sphere of \mathbf{q} with

$$\varepsilon = \frac{c}{4} \frac{\theta(1-\theta)}{(2k-1)L} \min_{\mathbf{q} \in \mathbb{S}^{k-1}} \|\boldsymbol{\zeta}\|_4^4, \quad (\text{F.111})$$

hence the number of samples N is bounded by

$$N = \left(\frac{3}{\varepsilon} \right)^k \quad (\text{F.112})$$

$$\leq \exp(-k \ln \varepsilon) \quad (\text{F.113})$$

$$\leq C \exp \left[k \log \left(\frac{\theta B^4 k^2 m}{1-\theta} \right) \right]. \quad (\text{F.114})$$

For $n = 1, \dots, N$, denote

$$P_i(\mathbf{q}_n) = \mathbb{P} \left[\left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta}_n)^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2 \geq \frac{c\theta(1-\theta)\|\boldsymbol{\zeta}_n\|_4^4}{4(2k-1)} \right] \quad (\text{F.115})$$

together with union bound over all \mathbf{q}_n , we obtain

$$\mathbb{P} \left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \frac{\left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2}{\|\boldsymbol{\zeta}\|_4^4} \geq \frac{c\theta(1-\theta)}{2(2k-1)} \right] \leq \sum_{\mathbf{q}_n \in \hat{\mathcal{R}}_{2C_*}} P_i(\mathbf{q}_n) \quad (\text{F.116})$$

$$\leq N \max_{\mathbf{q}_n \in \hat{\mathcal{R}}_{2C_*}} P_i(\mathbf{q}_n) \quad (\text{F.117})$$

$$\leq 8k \sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \exp \left(-\frac{cm(1-\theta)^2 \|\boldsymbol{\zeta}\|_4^8}{\theta k^3 + \theta(1-\theta)B^4 k^3 \|\boldsymbol{\zeta}\|_4^4} \right) \times \exp \left(k \ln \left(\frac{3}{\varepsilon} \right) \right). \quad (\text{F.118})$$

Hence

$$\mathbb{P} \left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \frac{\left\| \frac{\bar{\mathbf{X}}_0 \text{diag}(\bar{\mathbf{X}}_0^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_0^T}{m} - \mathbb{E}[\cdot] \right\|_2}{\|\boldsymbol{\zeta}\|_4^4} \geq \frac{c}{2}\theta(1-\theta) \right] \leq \sum_i \mathbb{P} \left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \frac{\left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2}{\|\boldsymbol{\zeta}\|_4^4} \geq \frac{c\theta(1-\theta)}{2(2k-1)} \right] \quad (\text{F.119})$$

$$\leq (2k-1) \max_i \mathbb{P} \left[\sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \frac{\left\| \frac{\bar{\mathbf{X}}_i \text{diag}(\bar{\mathbf{X}}_i^T \boldsymbol{\zeta})^{\circ 2} \bar{\mathbf{X}}_i^T}{m} - \mathbb{E}[\cdot] \right\|_2}{\|\boldsymbol{\zeta}\|_4^4} \geq \frac{c\theta(1-\theta)}{2(2k-1)} \right] \quad (\text{F.120})$$

$$\leq 16k^2 \sup_{\mathbf{q} \in \hat{\mathcal{R}}_{2C_*}} \exp \left(-\frac{c'm(1-\theta)^2 \|\boldsymbol{\zeta}\|_4^8}{\theta k^3 + \theta(1-\theta)B^4 k^3 \|\boldsymbol{\zeta}\|_4^4} \right) \times \exp \left(k \ln \left(\frac{3}{\varepsilon} \right) \right) \quad (\text{F.121})$$

Therefore, by taking

$$m \geq \frac{C\theta}{(1-\theta)^2} \min \left\{ (2C_* \mu \kappa^2)^{-4/3}, k^2 \right\} k^4 \log^3 k \quad (\text{F.122})$$

$$\geq C' \theta k \log \left(\frac{\theta k m B^4}{(1-\theta) \|\zeta\|_4^4} \right) \frac{k^3 + (1-\theta) B^4 k^3 \|\zeta\|_4^4}{(1-\theta)^2 \|\zeta\|_4^8} \quad (\text{F.123})$$

and adding up failure probability, we obtain

$$\begin{aligned} & \left\| \frac{1}{m} \mathbf{X}_0 \text{diag}(\mathbf{X}_0^T \mathbf{A}^T \mathbf{q})^{\circ 2} \mathbf{X}_0^T - \mathbb{E}[\cdot] \right\|_2 \\ & \leq c \theta (1-\theta) \|\mathbf{A}^T \mathbf{q}\|_4^4 \end{aligned} \quad (\text{F.124})$$

with probability no smaller than $1 - c_1 \exp(-k) - c_2 \theta (1-\theta)^2 k^{-4} - c_3 \exp(-\theta k)$. \blacksquare

G Tools

Lemma G.1 (Moments of the Gaussian Random Variables). *If $X \sim \mathcal{N}(0, \sigma^2)$, then it holds for all integer $p \geq 1$ that*

$$\mathbb{E}[|X|^p] = \sigma^p (p-1)!! \left[\sqrt{\frac{2}{\pi}} \mathbb{1}_{p \text{ odd}} + \mathbb{1}_{p \text{ even}} \right] \quad (\text{G.1})$$

$$\leq \sigma^p (p-1)!! \quad (\text{G.2})$$

Lemma G.2 (Moments of the χ^2 Random Variables). *If $X \sim \chi^2(n)$, then it holds for all integer $p \geq 1$,*

$$\mathbb{E}[X^p] = 2^p \frac{\Gamma(p + n/2)}{\Gamma(n/2)} \quad (\text{G.3})$$

$$= \prod_{k=1}^p (n + 2k - 2) \leq p! (2n)^p / 2 \quad (\text{G.4})$$

Lemma G.3 (Moments of the χ Random Variables). *If $X \sim \chi(n)$, then it holds for all integer $p \geq 1$,*

$$\mathbb{E}[X^p] = 2^{p/2} \frac{\Gamma(p/2 + n/2)}{\Gamma(n/2)} \leq p! n^{p/2}. \quad (\text{G.5})$$

Lemma G.4 (Moment-Control Bernstein's Inequality for Scalar RVs, Theorem 2.10 of [FR13]). *Let X_1, \dots, X_p be i.i.d. real-valued random variables. Suppose that there exist some positive number R and σ^2 such that*

$$\mathbb{E}[|X_k|^m] \leq \frac{m!}{2} \sigma^2 R^{m-2}, \quad \text{for all integers } m \geq 2.$$

Let $S \doteq \frac{1}{p} \sum_{k=1}^p X_k$, then for all $t > 0$, it holds that

$$\mathbb{P}[|S - \mathbb{E}[S]| \geq t] \leq 2 \exp \left(-\frac{pt^2}{2\sigma^2 + 2Rt} \right). \quad (\text{G.6})$$

Corollary G.5 (Moment-Control Bernstein's Inequality for Vector RVs, Corollary A.10 of [SQW15]). *Let $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^d$ be i.i.d. random vectors. Suppose there exist some positive number R and σ^2 such that*

$$\mathbb{E}[\|\mathbf{x}_k\|^m] \leq \frac{m!}{2} \sigma^2 R^{m-2}, \quad \text{for all integers } m \geq 2.$$

Let $\mathbf{s} = \frac{1}{p} \sum_{k=1}^p \mathbf{x}_k$, then for any $t > 0$, it holds that

$$\mathbb{P}[\|\mathbf{s} - \mathbb{E}[\mathbf{s}]\| \geq t] \leq 2(d+1) \exp \left(-\frac{pt^2}{2\sigma^2 + 2Rt} \right). \quad (\text{G.7})$$

Lemma G.6 (Moment-Control Bernstein's Inequality for Matrix RVs, Theorem 6.2 of [Tro12]). *Let $\mathbf{X}_1, \dots, \mathbf{X}_p \in \mathbb{R}^{d \times d}$ be i.i.d. random, symmetric matrices. Suppose there exist some positive number R and σ^2 such that*

$$\mathbb{E}[\mathbf{X}_k^m] \preceq \frac{m!}{2} \sigma^2 R^{m-2} \mathbf{I}, \quad (\text{G.8})$$

$$-\mathbb{E}[\mathbf{X}_k^m] \preceq \frac{m!}{2} \sigma^2 R^{m-2} \mathbf{I}. \quad (\text{G.9})$$

for all integers $m \geq 2$. Let $\mathbf{S} \doteq \frac{1}{p} \sum_{k=1}^p \mathbf{X}_k$, then for all $t > 0$, it holds that

$$\mathbb{P}[\|\mathbf{S} - \mathbb{E}[\mathbf{S}]\| \geq t] \leq 2d \exp\left(-\frac{pt^2}{2\sigma^2 + 2Rt}\right). \quad (\text{G.10})$$

Lemma G.7 (Bernstein's Inequality for Uncentered Matrix RVs). *The matrix Bernstein inequality states that for independent random matrices $\mathbf{M}_1, \dots, \mathbf{M}_n \in \mathbb{R}^{d_1 \times d_2}$, if*

$$\sigma^2 = \max\left\{\left\|\sum_{i=1}^n \mathbb{E}[\mathbf{M}_i \mathbf{M}_i^*]\right\|, \left\|\sum_{i=1}^n \mathbb{E}[\mathbf{M}_i^* \mathbf{M}_i]\right\|\right\}, \quad (\text{G.11})$$

and

$$\|\mathbf{M}_i\|_2 \leq R \quad \text{a.s.}, \quad (\text{G.12})$$

then

$$\mathbb{P}\left[\left\|\sum_i \mathbf{M}_i - \mathbb{E}[\cdot]\right\| > t\right] \leq (d_1 + d_2) \exp\left(\frac{-t^2/2}{\sigma^2 + 2Rt/3}\right). \quad (\text{G.13})$$

Proof For zero mean random matrices

$$\mathbf{M}_1 - \mathbb{E}\mathbf{M}_1, \dots, \mathbf{M}_n - \mathbb{E}\mathbf{M}_n \in \mathbb{R}^{d_1 \times d_2}, \quad (\text{G.14})$$

we have that

$$\|\mathbf{M}_i - \mathbb{E}\mathbf{M}_i\|_2 \leq 2R, \quad (\text{G.15})$$

and

$$\mathbf{0} \preceq \sum_{i=1}^n \mathbb{E}[(\mathbf{M}_i - \mathbb{E}\mathbf{M}_i)(\mathbf{M}_i - \mathbb{E}\mathbf{M}_i)^*] \quad (\text{G.16})$$

$$\preceq \sum_{i=1}^n \mathbb{E}[\mathbf{M}_i \mathbf{M}_i^*], \quad (\text{G.17})$$

$$\mathbf{0} \preceq \sum_{i=1}^n \mathbb{E}[(\mathbf{M}_i - \mathbb{E}\mathbf{M}_i)^*(\mathbf{M}_i - \mathbb{E}\mathbf{M}_i)] \quad (\text{G.18})$$

$$\preceq \sum_{i=1}^n \mathbb{E}[\mathbf{M}_i^* \mathbf{M}_i]. \quad (\text{G.19})$$

Plugging corresponding quantities back to Theorem 1.6 of [Tro12], we obtain that

$$\mathbb{P}\left[\left\|\sum_i \mathbf{M}_i - \mathbb{E}[\cdot]\right\| > t\right] \leq (d_1 + d_2) \exp\left(\frac{-t^2/2}{\sigma^2 + 2Rt/3}\right). \quad (\text{G.20})$$

■