

Geometry and Symmetry in Short-and-Sparse Deconvolution

Han-Wen Kuo^{1,2}, Yenson Lau^{1,2}, Yuqian Zhang³, John Wright^{1,2,4}

¹Department of Electrical Engineering, Columbia University

²Data Science Institute, Columbia University

³Department of Computer Science, Cornell University

⁴Department of Applied Physics and Applied Mathematics, Columbia University

January 3, 2019 Revised January 6, 2019

Abstract

We study the *Short-and-Sparse (SaS) deconvolution* problem of recovering a short signal \mathbf{a}_0 and a sparse signal \mathbf{x}_0 from their convolution. We propose a method based on nonconvex optimization, which under certain conditions recovers the target short and sparse signals, up to a signed shift symmetry which is intrinsic to this model. This symmetry plays a central role in shaping the optimization landscape for deconvolution. We give a *regional analysis*, which characterizes this landscape geometrically, on a union of subspaces. Our geometric characterization holds when the length- p_0 short signal \mathbf{a}_0 has shift coherence μ , and \mathbf{x}_0 follows a random sparsity model with sparsity rate $\theta \in \left[\frac{c_1}{p_0}, \frac{c_2}{p_0 \sqrt{\mu} + \sqrt{p_0}} \right] \cdot \frac{1}{\log^2 p_0}$. Based on this geometry, we give a provable method that successfully solves SaS deconvolution with high probability.

1 Introduction

Datasets in a wide range of areas, including neuroscience [Lew98], microscopy [CLC⁺17] and astronomy [Sah07], can be modeled as superpositions of translations of a basic motif. Data of this nature can be modeled mathematically as a convolution $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$, between a *short* signal \mathbf{a}_0 (the motif) and a longer *sparse* signal \mathbf{x}_0 , whose nonzero entries indicate where in the sample the motif is present. A very similar structure arises in image deblurring [CW98], where \mathbf{y} is a blurry image, \mathbf{a}_0 the blur kernel, and \mathbf{x}_0 the (edge map) of the target sharp image.

Motivated by these and related problems in imaging and scientific data analysis, we study the *Short-and-Sparse (SaS) Deconvolution* problem of recovering a short signal $\mathbf{a}_0 \in \mathbb{R}^{p_0}$ and a sparse signal $\mathbf{x}_0 \in \mathbb{R}^n$ ($n \gg p_0$) from their length- n cyclic convolution $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 \in \mathbb{R}^n$. This SaS model exhibits a basic *scaled shift symmetry*: for any nonzero scalar α and cyclic shift $s_\ell[\cdot]$,

$$\left(\alpha s_\ell[\mathbf{a}_0] \right) * \left(\frac{1}{\alpha} s_{-\ell}[\mathbf{x}_0] \right) = \mathbf{y}. \quad (1.1)$$

Because of this symmetry, we only expect to recover \mathbf{a}_0 and \mathbf{x}_0 up to a signed shift (see Figure 1). Our problem of interest can be stated more formally as:

Problem 1.1 (Short-and-Sparse Deconvolution). *Given the cyclic convolution $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 \in \mathbb{R}^n$ of $\mathbf{a}_0 \in \mathbb{R}^{p_0}$ short ($p_0 \ll n$), and $\mathbf{x}_0 \in \mathbb{R}^n$ sparse, recover \mathbf{a}_0 and \mathbf{x}_0 , up to a scaled shift.*

Despite a long history and many applications, until recently very little algorithmic theory was available for SaS deconvolution. Much of this difficulty can be attributed to the scale-shift symmetry: natural convex relaxations fail, and nonconvex formulations exhibit a complicated optimization landscape, with many

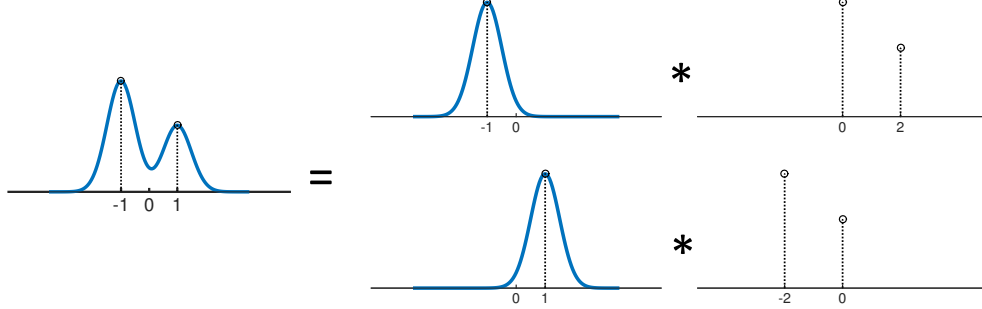


Figure 1: Shift symmetry in Short-and-Sparse deconvolution. An observation y (left) which is a convolution of a short signal a_0 and a sparse signal x_0 (top right) can be equivalently expressed as a convolution of $s_\ell[a_0]$ and $s_{-\ell}[x_0]$, where $s_\ell[\cdot]$ denotes a shift ℓ samples. The ground truth signals a_0 and x_0 can only be identified up to a scaled shift.

equivalent global minimizers (scaled shifts of the ground truth) and additional local minimizers (scaled shift truncations of the ground truth), and a variety of critical points [ZLK⁺17, ZKW18]. Currently available theory guarantees approximate recovery of a truncation¹ of a shift $s_\ell[a_0]$, rather than guaranteeing recovery of a_0 as a whole, and requires certain (complicated) conditions on the convolution matrix associated with a_0 [ZKW18].

In this paper, describe an algorithm which, under simpler conditions, *exactly* recovers a scaled shift of the pair (a_0, x_0) . Our algorithm is based on a formulation first introduced in [ZLK⁺17], which casts the deconvolution problem as (nonconvex) optimization over the sphere. We characterize the geometry of this objective function, and show that near a certain union of subspaces, every local minimizer is very close to a signed shift of a_0 . Based on this geometric analysis, we give provable methods for SaS deconvolution that exactly recover a scaled shift of (a_0, x_0) whenever a_0 is *shift-incoherent* and x_0 is a sufficiently sparse random vector. Our geometric analysis highlights the role of symmetry in shaping the objective landscape for SaS deconvolution.

Organization of this paper. The remainder of this paper is organized as follows. Section 2 introduces our optimization approach and modeling assumptions. Section 3 introduces our main results — both geometric and algorithmic — and compares them to the literature. Section 4-5 describes the main ideas of our analysis. Finally, Section 7 discusses two main limitations of our analysis and describes directions for future work.

2 Formulation and Assumptions

2.1 Nonconvex SaS over the Sphere

Bilinear Lasso. Our starting point is the (natural) formulation

$$\min_{\mathbf{a}, \mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2}_{\text{Data Fidelity}} + \lambda \underbrace{\|\mathbf{x}\|_1}_{\text{Sparsity}} \quad \text{s.t.} \quad \|\mathbf{a}\|_2 = 1. \quad (2.1)$$

We term this optimization problem the *Bilinear Lasso*, for its resemblance to the Lasso estimator in statistics. Indeed, letting

$$\varphi_{\text{lasso}}(\mathbf{a}) \equiv \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\} \quad (2.2)$$

denote the optimal Lasso cost, we see that (2.1) simply optimizes φ_{lasso} with respect to \mathbf{a} :

$$\min_{\mathbf{a}} \varphi_{\text{lasso}}(\mathbf{a}) \quad \text{s.t.} \quad \|\mathbf{a}\|_2 = 1. \quad (2.3)$$

¹I.e., the portion of the shifted signal $s_\ell[a_0]$ that falls in the window $\{0, \dots, p_0 - 1\}$.

In (2.1)-(2.3), we constrain \mathbf{a} to have unit ℓ^2 norm. This constraint breaks the scale ambiguity between \mathbf{a} and \mathbf{x} . Moreover, the choice of constraint manifold has surprisingly strong implications for computation: if \mathbf{a} is instead constrained to the simplex, the problem admits trivial global minimizers. In contrast, local minima of the sphere-constrained formulation often correspond to shifts (or shift truncations [ZLK⁺17]) of the ground truth \mathbf{a}_0 .

Simplifications and approximations. The problem (2.3) is defined in terms of the optimal Lasso cost. This function is challenging to analyze, especially far away from \mathbf{a}_0 . [ZLK⁺17] analyzes the local minima of a simplification of (2.3), obtained by approximating² the data fidelity term as

$$\begin{aligned} \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 &= \frac{1}{2} \|\mathbf{a} * \mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2, \\ &\approx \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2. \end{aligned} \quad (2.4)$$

This yields a simpler objective function

$$\varphi_{\ell^1}(\mathbf{a}) = \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}. \quad (2.5)$$

We make one further simplification to this problem, replacing the nondifferentiable penalty $\|\cdot\|_1$ with a smooth approximation $\rho(\mathbf{x})$.³ Our analysis allows for a variety of smooth sparsity surrogates $\rho(\mathbf{x})$; for concreteness, we state our main results for the particular penalty⁴

$$\rho(\mathbf{x}) = \sum_i (\mathbf{x}_i^2 + \delta^2)^{1/2}. \quad (2.6)$$

For $\delta > 0$, this is a smooth function of \mathbf{x} ; as $\delta \searrow 0$ it approaches $\|\mathbf{x}\|_1$. Replacing $\|\cdot\|_1$ with $\rho(\cdot)$, we obtain the objective function which will be our main object of study,

$$\varphi_{\rho}(\mathbf{a}) = \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 + \lambda \rho(\mathbf{x}) \right\}. \quad (2.7)$$

Core optimization problem. As in [ZLK⁺17], we optimize $\varphi_{\rho}(\mathbf{a})$ over the sphere \mathbb{S}^{p-1} :

$$\boxed{\min_{\mathbf{a}} \varphi_{\rho}(\mathbf{a}) \quad \text{s.t.} \quad \mathbf{a} \in \mathbb{S}^{p-1}.} \quad (2.8)$$

Here, we set $p = 3p_0 - 2$. As we will see, optimizing over this slightly higher dimensional sphere enables us to recover a (full) shift of \mathbf{a}_0 , rather than a *truncated* shift. Our approach will leverage the following fact: if we view $\mathbf{a} \in \mathbb{S}^{p-1}$ as indexed by coordinates $W = \{-p_0 + 1, \dots, 2p_0 - 1\}$, then for any shifts $\ell \in \{-p_0 + 1, \dots, p_0 - 1\}$, the support of ℓ -shifted short signal $s_{\ell}[\mathbf{a}_0]$ is entirely contained in interval W . We will give a provable method which recovers a scaled version of one of these canonical shifts.

2.2 Analysis Setting and Assumptions

For convenience, we assume that \mathbf{a}_0 has unit ℓ^2 norm, i.e., $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$.⁵ Our analysis makes two main assumptions, on the short motif \mathbf{a}_0 and the sparse map \mathbf{x}_0 , respectively:

Shift incoherence of \mathbf{a}_0 . The first is that distinct shifts \mathbf{a}_0 have small inner product. We define the *shift coherence* of $\mu(\mathbf{a}_0)$ to be the largest inner product between distinct shifts:

$$\mu(\mathbf{a}_0) = \max_{\ell \neq 0} |\langle \mathbf{a}_0, s_{\ell}[\mathbf{a}_0] \rangle| \quad (2.9)$$

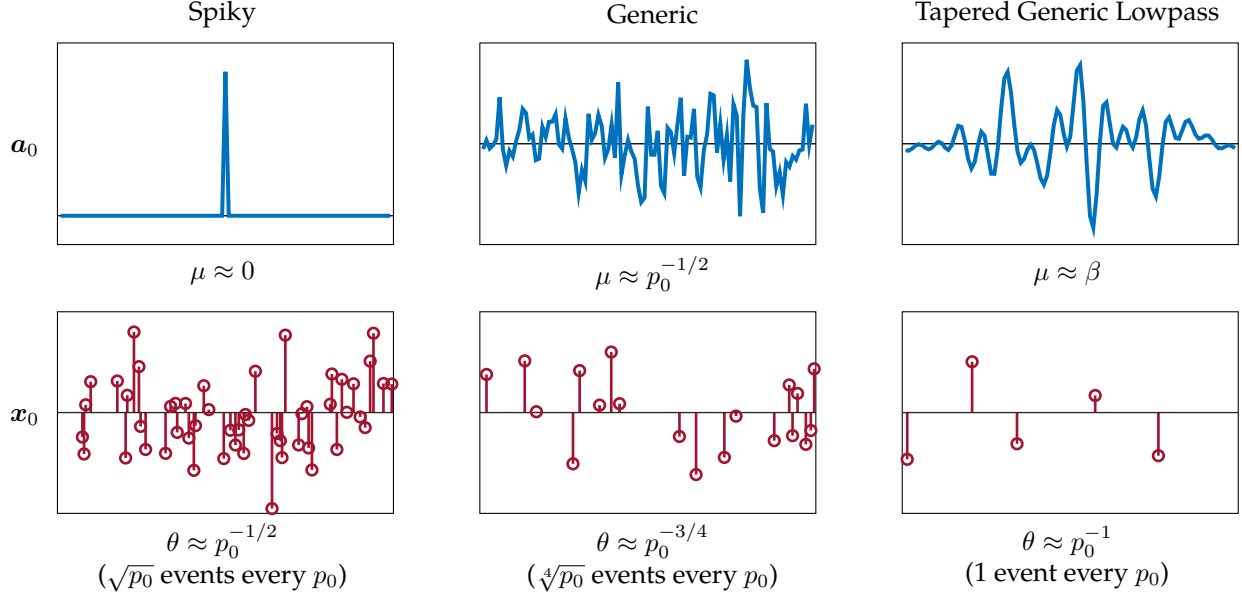


Figure 2: Sparsity-coherence tradeoff: Top: three families of motifs \mathbf{a}_0 with varying coherence μ . Bottom: maximum allowable sparsity θ and number of copies θp_0 within each length- p_0 window. Here, we suppress constants and logarithmic factors. When the target motif has smaller shift-coherence μ , our result allows larger θ , and vice versa. This sparsity-coherence tradeoff is made precise in our main result [Theorem 3.1](#), which, loosely speaking, asserts that when $\theta \lesssim 1/(p_0\sqrt{\mu} + \sqrt{p_0})$, our method succeeds.

The quantity $\mu(\mathbf{a}_0)$ is bounded between 0 and 1. Our theory allows any μ smaller than some numerical constant. [Figure 2](#) shows three examples of families of \mathbf{a}_0 that satisfy this assumption:

- *Spiky.* When \mathbf{a}_0 is close to the Dirac delta δ_0 , the shift coherence $\mu(\mathbf{a}_0) \approx 0$.⁶ Here, the observed signal \mathbf{y} consists of a superposition of sharp pulses. This is arguably the easiest instance of SaS deconvolution.
- *Generic.* If \mathbf{a}_0 is chosen uniformly at random from the sphere \mathbb{S}^{p_0-1} , its coherence is bounded as $\mu(\mathbf{a}_0) \lesssim \sqrt{1/p_0}$ with high probability.
- *Tapered Generic Lowpass.* Here, \mathbf{a}_0 is generated by taking a random conjugate symmetric superposition of the first L length- p_0 Discrete Fourier Transform (DFT) basis signals, windowing (e.g., with a Hamming window) and normalizing to unit ℓ^2 norm. When $L = p_0\sqrt{1-\beta}$, with high probability $\mu(\mathbf{a}_0) \lesssim \beta$. In this model, μ does not have to diminish as p_0 grows – it can be a fixed constant.⁷

Intuitively speaking, problems with smaller μ are easier to solve, a claim which will be made precise in our technical results.

²For a generic \mathbf{a} , we have $\langle s_i[\mathbf{a}], s_j[\mathbf{a}] \rangle \approx 0$ and hence $\|\mathbf{a} * \mathbf{x}\|_2^2 = \mathbf{x}^* \mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}} \mathbf{x} \approx \mathbf{x}^* \mathbf{I} \mathbf{x} = \|\mathbf{x}\|_2^2$.

³The objective φ_{ℓ^1} is not twice differentiable everywhere, and hence cannot be minimized using conventional second order methods.

⁴This particular surrogate is sometimes being named as the pseudo-Huber function.

⁵This is purely a technical convenience. Our theory guarantees recovery of a signed shift $(\pm s_{\ell}[\mathbf{a}_0], \pm s_{-\ell}[\mathbf{x}_0])$ of the truth. If \mathbf{a}_0 does not have unit norm, identical reasoning implies that our method recovers a scaled shift $(\alpha s_{\ell}[\mathbf{a}_0], \alpha^{-1} s_{-\ell}[\mathbf{x}_0])$ with $\alpha = \pm \frac{1}{\|\mathbf{a}_0\|_2}$.

⁶The use of “ \approx ” here suppresses constant and logarithmic factors.

⁷The upper right panel of [Figure 2](#) is generated using random DFT components with frequencies smaller than one-third Nyquist. Such a kernel is incoherent, with high probability. Many commonly occurring low-pass kernels have $\mu(\mathbf{a}_0)$ larger – very close to one. One of the most important limitations of our results is that they do not provide guarantees in this highly coherent situation.

Random sparsity model on x_0 . We assume that x_0 is a sparse random vector. More precisely, we assume that x_0 is Bernoulli-Gaussian, with rate θ :

$$x_{0i} = \omega_i g_i, \quad (2.10)$$

where $\omega_i \sim \text{Ber}(\theta)$, $g_i \sim \mathcal{N}(0, 1)$ and all random variables are jointly independent. We write this as

$$x_0 \sim_{\text{i.i.d.}} \text{BG}(\theta). \quad (2.11)$$

Here, θ is the probability that a given entry x_{0i} is nonzero. Problems with smaller θ are easier to solve. In the extreme case, when $\theta \ll 1/p_0$, the observation y contains many isolated copies of the motif a_0 , and a_0 can be determined by direct inspection. Our analysis will focus on the nontrivial scenario, when $\theta \gtrsim 1/p_0$.

Sparsity-Coherence tradeoffs. Our technical results will articulate *sparsity-coherence* tradeoffs, in which smaller coherence μ enables larger θ , and vice-versa. More specifically, in our main theorem, the sparsity-coherence relationship is captured in the form

$$\theta \lesssim 1/(p_0\sqrt{\mu} + \sqrt{p_0}). \quad (2.12)$$

When the target a_0 is highly shift-incoherent ($\mu \approx 0$), our method succeeds when each length- p_0 window contains about $\sqrt{p_0}$ copies of a_0 . When μ is larger (as in the generic lowpass model), our method succeeds as long as relatively few copies of a_0 overlap in the observed signal. In Figure 2, we illustrate these tradeoffs for the three models described above.

3 Main Results: Geometry and Algorithms

In this section, we introduce our main results – on the geometry of φ_ρ (Section 3.1) and its algorithmic implications (Section 3.2). Finally, in Section 3.3, we compare these results with the literature on deconvolution.

3.1 Geometry of the Objective φ_ρ

The goal in SaS deconvolution is to recover a_0 (and x_0) up to a signed shift — i.e., we wish to recover some $\pm s_\ell[a_0]$. The shifts $\pm s_\ell[a_0]$ play a key role in shaping the landscape of φ_ρ . In particular, we will argue that over a certain subset of the sphere, *every local minimum of φ_ρ is close to some $\pm s_\ell[a_0]$* .

Geometry near a single shift. To gain intuition into the properties of φ_ρ , we first visualize this function in the vicinity of a single shift $s_\ell[a_0]$ of the ground truth a_0 . In Figure 3, we plot the function value of φ_ρ over

$$\mathcal{B}_{\ell^2, r}(s_\ell[a_0]) \cap \mathbb{S}^{p-1},$$

where $\mathcal{B}_{\ell^2, r}(a)$ is a ball of radius r around a . We make two observations:

- The objective function φ_ρ is strongly convex on this neighborhood of $s_\ell[a_0]$.
- There is a local minimizer very close to $s_\ell[a_0]$.

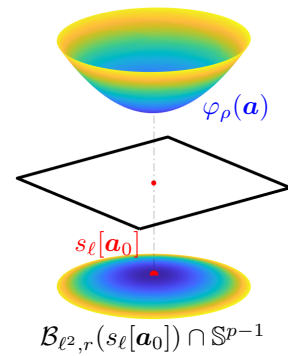


Figure 3: Geometry of φ_ρ near a shift of a_0 . Bottom: a portion of the sphere \mathbb{S}^{p-1} , colored according to φ_ρ . Top: φ_ρ visualized as height. φ_ρ is strongly convex in this region, and it has a minimizer very close to $s_\ell[a_0]$.

Geometry near the span of two shifts. We next visualize the objective function φ_ρ near the linear span of *two* different shifts $s_{\ell_1}[\mathbf{a}_0]$ and $s_{\ell_2}[\mathbf{a}_0]$. More precisely, we plot φ_ρ near the intersection (Figure 4, left) of the sphere \mathbb{S}^{p-1} and the linear subspace

$$\mathcal{S}_{\{\ell_1, \ell_2\}} = \{ \alpha_1 s_{\ell_1}[\mathbf{a}_0] + \alpha_2 s_{\ell_2}[\mathbf{a}_0] \mid \alpha_1, \alpha_2 \in \mathbb{R} \}.$$

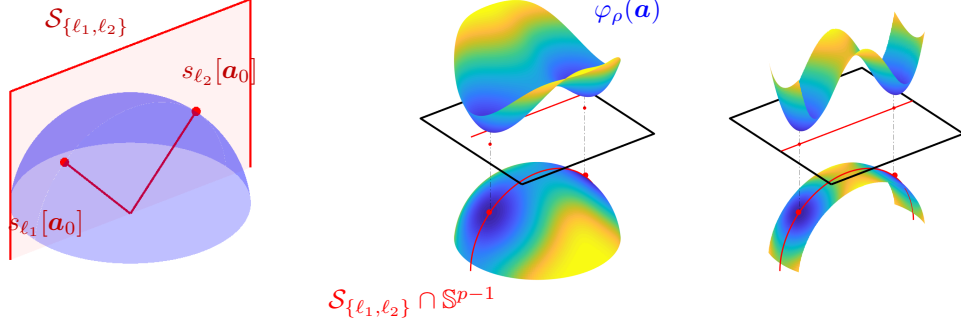


Figure 4: Geometry of φ_ρ near the span $\mathcal{S}_{\{\ell_1, \ell_2\}}$ of two shifts of \mathbf{a}_0 . Left: each pair of shifts $s_{\ell_1}[\mathbf{a}_0]$, $s_{\ell_2}[\mathbf{a}_0]$ defines a linear subspace $\mathcal{S}_{\{\ell_1, \ell_2\}}$ of \mathbb{R}^p . Center/right: every local minimum of φ_ρ near $\mathcal{S}_{\{\ell_1, \ell_2\}}$ (red line) is close to either $s_{\ell_1}[\mathbf{a}_0]$ or $s_{\ell_2}[\mathbf{a}_0]$; there is a negative curvature in the middle of $s_{\ell_1}[\mathbf{a}_0]$, $s_{\ell_2}[\mathbf{a}_0]$, and φ_ρ is convex in direction away from $\mathcal{S}_{\ell_1, \ell_2}$.

We make three observations:

- Again, there is a local minimizer near each shift $s_\ell[\mathbf{a}_0]$.
- These are the *only* local minimizers in the vicinity of $\mathcal{S}_{\{\ell_1, \ell_2\}}$. In particular, the objective function φ exhibits *negative curvature* along $\mathcal{S}_{\{\ell_1, \ell_2\}}$ at any superposition $\alpha_1 s_{\ell_1}[\mathbf{a}_0] + \alpha_2 s_{\ell_2}[\mathbf{a}_0]$ whose weights α_1 and α_2 are balanced, i.e., $|\alpha_1| \approx |\alpha_2|$.
- Furthermore, the function φ_ρ exhibits *positive curvature* in directions away from the subspace $\mathcal{S}_{\ell_1, \ell_2}$.

Geometry in the span of multiple shifts. Finally, we visualize φ_ρ over the intersection (Figure 5, left) of the sphere \mathbb{S}^{p-1} with the linear span of three shifts $s_{\ell_1}[\mathbf{a}_0]$, $s_{\ell_2}[\mathbf{a}_0]$, $s_{\ell_3}[\mathbf{a}_0]$ of the true kernel \mathbf{a}_0 :

$$\mathcal{S}_{\{\ell_1, \ell_2, \ell_3\}} = \{ \alpha_1 s_{\ell_1}[\mathbf{a}_0] + \alpha_2 s_{\ell_2}[\mathbf{a}_0] + \alpha_3 s_{\ell_3}[\mathbf{a}_0] \mid \alpha_1, \alpha_2, \alpha_3 \in \mathbb{R} \}$$

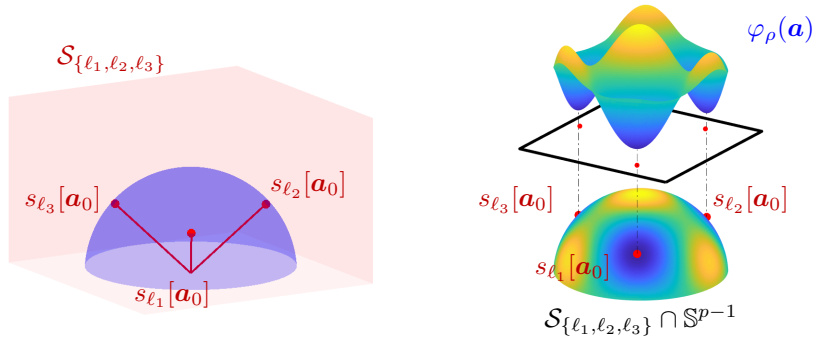


Figure 5: Geometry of φ_ρ over the span $\mathcal{S}_{\{\ell_1, \ell_2, \ell_3\}}$ of three shifts of \mathbf{a}_0 . The subspace $\mathcal{S}_{\{\ell_1, \ell_2, \ell_3\}}$ is three-dimensional; its intersection with the sphere \mathbb{S}^{p-1} is isomorphic to a two-dimensional sphere. On this set, φ_ρ has local minimizers near each of the $s_{\ell_i}[\mathbf{a}_0]$, and are the only minimizers near $\mathcal{S}_{\ell_1, \ell_2, \ell_3}$.

Again, *there is a local minimizer near each signed shift*. At roughly balanced superpositions of shifts, the objective function exhibits negative curvature. As a result, again, the *only* local minimizers are close to signed shifts.

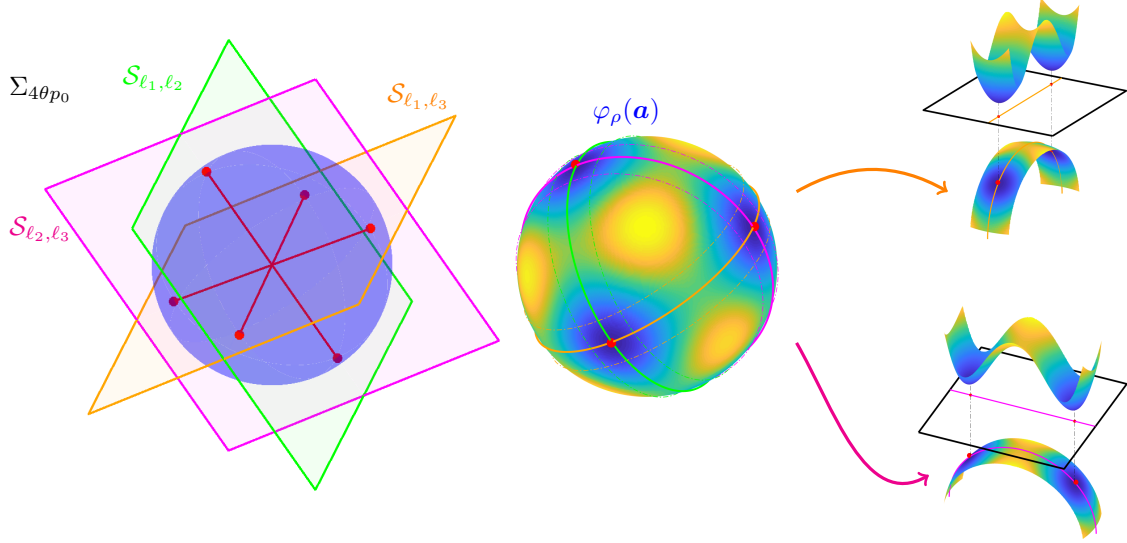


Figure 6: Geometry of φ_ρ over the union of subspaces $\Sigma_{4\theta p_0}$. Left: schematic representation of the union of subspaces $\Sigma_{4\theta p_0}$. For each set τ of at most $4\theta p_0$ shifts, we have a subspace \mathcal{S}_τ . Right: φ_ρ has good geometry near this union of subspaces.

Geometry of φ_ρ over a union of subspaces. Our main geometric result will show that these properties obtain on *every* subspace spanned by a few shifts of \mathbf{a}_0 . Indeed, for each subset

$$\tau \subseteq \{-p_0 + 1, \dots, p_0 - 1\}, \quad (3.1)$$

define a linear subspace

$$\mathcal{S}_\tau = \left\{ \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] \mid \alpha_{-p_0+1}, \dots, \alpha_{p_0-1} \in \mathbb{R} \right\}. \quad (3.2)$$

The subspace \mathcal{S}_τ is the linear span of the shifts $s_\ell[\mathbf{a}_0]$ indexed by ℓ in the set τ . Our geometric theory will show that with high probability the function φ_ρ has no spurious local minimizers near any \mathcal{S}_τ for which τ is not too large – say, $|\tau| \leq 4\theta p_0$. Combining all of these subspaces into a single geometric object, define the union of subspaces

$$\Sigma_{4\theta p_0} = \bigcup_{|\tau| \leq 4\theta p_0} \mathcal{S}_\tau. \quad (3.3)$$

Figure 6 (left) gives a schematic representation of this set. We claim:

- In the neighborhood of $\Sigma_{4\theta p_0}$, all local minimizers are near signed shifts.
- The value of φ_ρ grows in any direction away from $\Sigma_{4\theta p_0}$.

Main Geometric Result. Our main result formalizes the above observations, under two key assumptions: first, that the sparsity rate θ is sufficiently small (relative to the shift coherence μ of p_0), and, second, the signal length n is sufficiently large:

Theorem 3.1 (Main Geometric Theorem). *Let $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ with $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$ μ -shift coherent and $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$ with sparsity rate*

$$\theta \in \left[\frac{c_1}{p_0}, \frac{c_2}{p_0 \sqrt{\mu} + \sqrt{p_0}} \right] \cdot \frac{1}{\log^2 p_0}. \quad (3.4)$$

Choose $\rho(x) = \sqrt{x^2 + \delta^2}$ and set $\lambda = 0.1/\sqrt{p_0\theta}$ in φ_ρ . Then there exists $\delta > 0$ and numerical constant c such that if $n \geq \text{poly}(p_0)$, with high probability, every local minimizer $\bar{\mathbf{a}}$ of φ_ρ over $\Sigma_{4\theta p_0}$ satisfies $\|\bar{\mathbf{a}} - \sigma_{s_\ell}[\mathbf{a}_0]\|_2 \leq c \max\{\mu, p_0^{-1}\}$ for some signed shift $\sigma_{s_\ell}[\mathbf{a}_0]$ of the true kernel. Above, $c_1, c_2 > 0$ are positive numerical constants.

Proof. This follows from [Theorem 4.1](#). ■

The upper bound on θ in (3.4) yields the tradeoff between coherence and sparsity described in [Figure 2](#). Simply put, when \mathbf{a}_0 is better conditioned (as a kernel), its coherence μ is smaller and \mathbf{x}_0 can be denser.

At a technical level, our proof of [Theorem 3.1](#) shows that (i) $\varphi_\rho(\mathbf{a})$ is strongly convex in the vicinity of each signed shift, and that at every other point \mathbf{a} near $\Sigma_{4\theta p_0}$, there is either (ii) a nonzero gradient or (iii) a direction of strict negative curvature; furthermore (iv) the function φ_ρ grows away from $\Sigma_{4\theta p_0}$. Points (ii)-(iii) imply that near $\Sigma_{4\theta p_0}$ there are no “flat” saddles: every saddle point has a direction of strict negative curvature. We will leverage these properties to propose an efficient algorithm for finding a local minimizer near $\Sigma_{4\theta p_0}$. Moreover, this minimizer is close enough to a shift (here, $\|\bar{\mathbf{a}} - s_\ell[\mathbf{a}_0]\|_2 \lesssim \mu$) for us to exactly recover $s_\ell[\mathbf{a}_0]$: we will give a refinement algorithm that produces $(\pm s_\ell[\mathbf{a}_0], \pm s_{-\ell}[\mathbf{x}_0])$.

3.2 Provable Algorithm for SaS Deconvolution

The objective function φ_ρ has good geometric properties on (and near!) the union of subspaces $\Sigma_{4\theta p_0}$. In this section, we show how to use give an efficient method that exactly recovers \mathbf{a}_0 and \mathbf{x}_0 , up to shift symmetry. Although our geometric analysis only controls φ_ρ near $\Sigma_{4\theta p_0}$, we will give a descent method which, with appropriate initialization $\mathbf{a}^{(0)}$, produces iterates $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(k)}, \dots$ that remain close to $\Sigma_{4\theta p_0}$ for all k . In short, it is easy to *start* near $\Sigma_{4\theta p_0}$ and easy to *stay* near $\Sigma_{4\theta p_0}$. After finding a local minimizer $\bar{\mathbf{a}}$, we refine it to produce a signed shift of $(\mathbf{a}_0, \mathbf{x}_0)$ using alternating minimization.

The next two paragraphs give the main ideas behind the main steps of the algorithm. We then describe its components in more detail ([Algorithm 1](#)) and state our main algorithmic result ([Theorem 3.2](#)), which asserts that under appropriate conditions this method produces a signed shift of $(\mathbf{a}_0, \mathbf{x}_0)$.

Minimization: Starting and staying near $\Sigma_{4\theta p_0}$. Our algorithm starts with a initialization scheme which generates $\mathbf{a}^{(0)}$ near the union of subspaces $\Sigma_{4\theta p_0}$, which consists of linear combinations of just a few shifts of \mathbf{a}_0 . How can we find a point near this union? Notice that the data \mathbf{y} also consists of a linear combination of just a few shifts of \mathbf{a}_0 . Indeed:

$$\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 = \sum_{\ell \in \text{supp}(\mathbf{x}_0)} \mathbf{x}_{0\ell} s_\ell[\mathbf{a}_0]. \quad (3.5)$$

A length- p_0 segment of data $\mathbf{y}_{0,\dots,p_0-1} = [\mathbf{y}_0, \dots, \mathbf{y}_{p_0-1}]^*$ captures portions of roughly $2\theta p_0 \ll 4\theta p_0$ shifts $s_\ell[\mathbf{a}_0]$.

Many of these copies of \mathbf{a}_0 are truncated by the restriction to $\{0, \dots, p_0 - 1\}$. A relatively simple remedy is as follows: first, we zero-pad $\mathbf{y}_{0,\dots,p_0-1}$ to length $p = 3p_0 - 2$, giving

$$[\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}]. \quad (3.6)$$

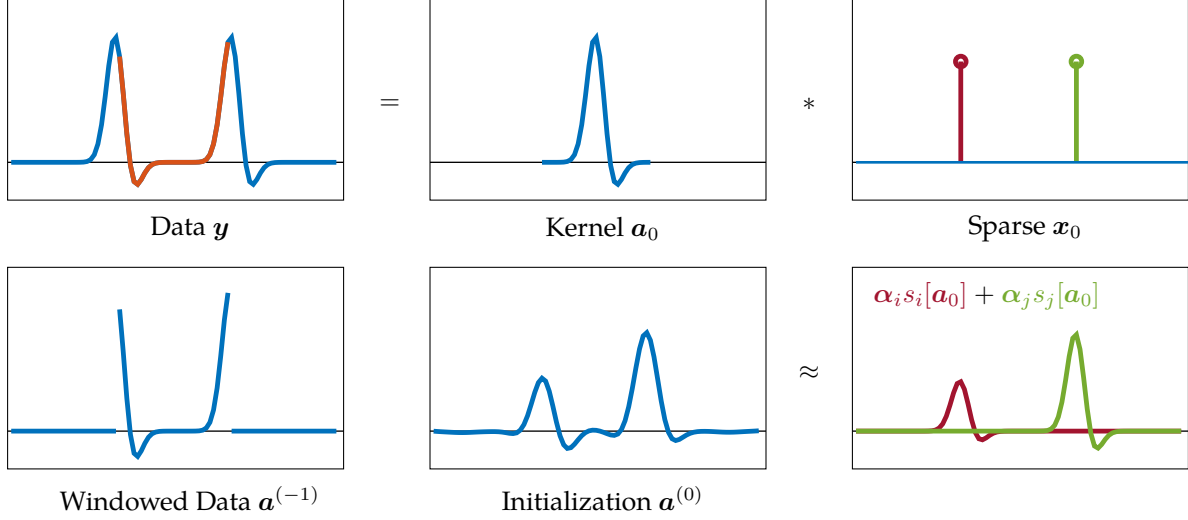


Figure 7: Data-driven initialization: using a piece of the observed data \mathbf{y} to generate an initial point $\mathbf{a}^{(0)}$ that is close to a superposition of shifts $s_\ell[\mathbf{a}_0]$ of the ground truth. Top: data $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ is a superposition of shifts of the true kernel \mathbf{a}_0 . Bottom: a length- p_0 window contains pieces of just a few shifts. Bottom middle: one step of the generalized power method approximately fills in the missing pieces, yielding a near superposition of shifts of \mathbf{a}_0 (right).

Zero padding provides enough space to accommodate any shift $s_\ell[\mathbf{a}_0]$ with $\ell \in \tau$. We then perform one step of the generalized power method⁸, writing

$$\mathbf{a}^{(0)} = -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\ell^1} \left(\mathbf{P}_{\mathbb{S}^{p-1}} \left[\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1} \right] \right), \quad (3.7)$$

where $\mathbf{P}_{\mathbb{S}^{p-1}}$ projects onto the sphere. The reasoning behind this construction may seem obscure. We will explain it at a more technical level in Section 5 after interpreting the gradient $\nabla \varphi_\rho$ in terms of its action on the shifts $s_\ell[\mathbf{a}_0]$ in Section 4. For now, we note that this operation has the effect of (approximately) filling in the missing pieces of the truncated shifts $s_\ell[\mathbf{a}_0]$ – see Figure 7 for an example. We will prove that with high probability $\mathbf{a}^{(0)}$ is indeed close to $\Sigma_{4\theta p_0}$.

The next key observation is that the function φ_ρ grows as we move away from the subspace \mathcal{S}_τ – see Figure 8. Because of this, a small-stepping descent method will not move far away from $\Sigma_{4\theta p_0}$. For concreteness, we will analyze a variant of the curvilinear search method [Gol80, GMWZ17], which moves in a linear combination of the negative gradient direction $-\mathbf{g}$ and a negative curvature direction $-\mathbf{v}$. At the k -th iteration, the algorithm updates $\mathbf{a}^{(k+1)}$ as

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)} \right] \quad (3.8)$$

with appropriately chosen step size t . The inclusion of a negative curvature direction allows the method to avoid stagnation near saddle points. Indeed, we will prove that starting from initialization $\mathbf{a}^{(0)}$, this method produces a sequence $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \dots$ which efficiently converges to a local minimizer $\tilde{\mathbf{a}}$ that is near some signed shift $\pm s_\ell[\mathbf{a}_0]$ of the ground truth.

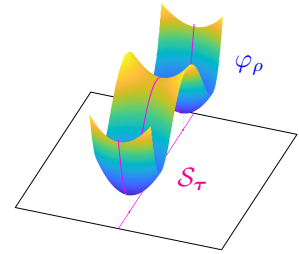


Figure 8: Growth of φ_ρ away from \mathcal{S}_τ . Because φ_ρ grows away from \mathcal{S}_τ , small-stepping descent methods stay near \mathcal{S}_τ .

⁸The power method for minimizing a quadratic form $\xi(\mathbf{a}) = \frac{1}{2} \mathbf{a}^* \mathbf{M} \mathbf{a}$ over the sphere consists of the iteration $\mathbf{a} \mapsto -\mathbf{P}_{\mathbb{S}^{p-1}} \mathbf{M} \mathbf{a}$. Notice that in this mapping, $-\mathbf{M} \mathbf{a} = -\nabla \xi(\mathbf{a})$. The generalized power method, for minimizing a function φ over the sphere consists of repeatedly projecting $-\nabla \varphi$ onto the sphere, giving the iteration $\mathbf{a} \mapsto -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi(\mathbf{a})$. (3.7) can be interpreted as one step of the generalized power method for the objective function φ_ρ .

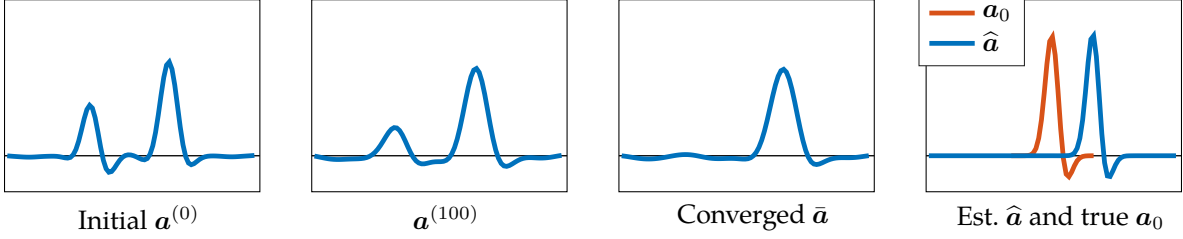


Figure 9: Local minimization and refinement. Left: data-driven initialization $\mathbf{a}^{(0)}$ consisting of a near-superposition of two shifts. Middle: minimizing φ_ρ produces a near shift of \mathbf{a}_0 . Right: rounded solution $\hat{\mathbf{a}}$ using the Lasso. $\hat{\mathbf{a}}$ is very close to a shift of \mathbf{a}_0 .

Refinement: Rounding a near-solution with homotopy alternating minimization. The second step of our algorithm *rounds* the local minimizer $\bar{\mathbf{a}} \approx \sigma_{s_\ell}[\mathbf{a}_0]$ to produce an exact solution $\hat{\mathbf{a}} = \sigma_{s_\ell}[\mathbf{a}_0]$. As a byproduct, it also exactly recovers the corresponding signed shift of the true sparse signal, $\hat{\mathbf{x}} = \sigma_{s_\ell}[\mathbf{x}_0]$.

Our rounding algorithm is an alternating minimization scheme, which alternates between minimizing the Lasso cost over \mathbf{a} with \mathbf{x} fixed, and minimizing the Lasso cost over \mathbf{x} with \mathbf{a} fixed. We make two modifications to this basic idea, both of which are important for obtaining exact recovery. First, unlike the standard Lasso cost, which penalizes all of the entries of \mathbf{x} , we maintain a running estimate $I^{(k)}$ of the support of \mathbf{x}_0 , and only penalize those entries that are not in $I^{(k)}$:

$$\frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i \notin I^{(k)}} |\mathbf{x}_i|. \quad (3.9)$$

This can be viewed as an extreme form of *reweighting* [CWB08]. Second, our algorithm gradually decreases penalty variable λ to 0, so that eventually

$$\hat{\mathbf{a}} * \hat{\mathbf{x}} \approx \mathbf{y}. \quad (3.10)$$

This can be viewed as a *homotopy* or *continuation* method [OPT00, EHJ⁺04]. For concreteness, at k -th iteration the algorithm reads:

$$\text{Update } \mathbf{x}: \quad \mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|, \quad (3.11)$$

$$\text{Update } \mathbf{a}: \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\underset{\mathbf{a}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right], \quad (3.12)$$

$$\text{Update } \lambda \text{ and } I: \quad \lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \operatorname{supp}(\mathbf{x}^{(k+1)}). \quad (3.13)$$

We prove that the iterates produced by this sequence of operations converge to the ground truth at a linear rate, as long as the initializer $\bar{\mathbf{a}}$ is sufficiently nearby.

Algorithm and Main Algorithmic Result. Our overall algorithm is summarized as [Algorithm 1](#). [Figure 9](#) illustrates the main steps of this algorithm. Our main algorithmic result states that under closely related hypotheses as above, [Algorithm 1](#) produces a signed shift of the ground truth $(\mathbf{a}_0, \mathbf{x}_0)$:

Algorithm 1 Short and Sparse Deconvolution

Input: Observation \mathbf{y} , motif length p_0 , sparsity θ , shift-coherence μ , and curvature threshold $-\eta_v$.

Minimization:

Set $\mathbf{a}^{(0)} \leftarrow -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi_\rho (\mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}])$.

Set $\lambda = 0.1/\sqrt{p_0\theta}$ ⁹ and $\delta > 0$ in φ_ρ . For $k = 1, 2, \dots, K_1$, let

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)}] \quad (3.14)$$

where $\mathbf{g}^{(k)}$ is the Riemannian gradient; $\mathbf{v}^{(k)}$ is the eigenvector of smallest Riemannian Hessian eigenvalue if less than $-\eta_v$ with $\langle \mathbf{v}^{(k)}, \mathbf{g}^{(k)} \rangle \geq 0$, otherwise let $\mathbf{v}^{(k)} = \mathbf{0}$; and $t \in (0, 0.1/n\theta]$ satisfies

$$\varphi_\rho(\mathbf{a}^{(k+1)}) < \varphi_\rho(\mathbf{a}^{(k)}) - \frac{1}{2}t\|\mathbf{g}^{(k)}\|_2^2 - \frac{1}{4}t^4\eta_v\|\mathbf{v}^{(k)}\|_2^2 \quad (3.15)$$

to obtain a near local minimizer $\bar{\mathbf{a}} \leftarrow \mathbf{a}^{(K_1)}$.

Refinement:

Set $\mathbf{a}^{(0)} \leftarrow \bar{\mathbf{a}}$, $\lambda^{(0)} \leftarrow 10(p\theta + \log n)(\mu + 1/p)$ and $I^{(0)} \leftarrow \mathcal{S}_{\lambda^{(0)}}[\text{supp}(\tilde{\mathbf{y}} * \bar{\mathbf{a}})]$. For $k = 1, 2, \dots, K_2$, let

$$\mathbf{x}^{(k+1)} \leftarrow \arg\min_{\mathbf{x}} \frac{1}{2}\|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|, \quad (3.16)$$

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\arg\min_{\mathbf{a}} \frac{1}{2}\|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2], \quad (3.17)$$

$$\lambda^{(k+1)} \leftarrow \lambda^{(k)}/2, \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}), \quad (3.18)$$

to obtain $(\hat{\mathbf{a}}, \hat{\mathbf{x}}) \leftarrow (\mathbf{a}^{(K_2)}, \mathbf{x}^{(K_2)})$.

Output: Return $(\hat{\mathbf{a}}, \hat{\mathbf{x}})$.

Theorem 3.2 (Main Algorithmic Theorem). Suppose $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ where $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$ is μ -truncated shift coherent such that $\max_{i \neq j} |\langle \boldsymbol{\iota}_{p_0}^* s_i[\mathbf{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu$ and $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$ with θ, μ satisfying

$$\theta \in \left[\frac{c_1}{p_0}, \frac{c_2}{(p_0\sqrt{\mu} + \sqrt{p_0}) \log^2 p_0} \right], \quad \mu \leq \frac{c_3}{\log^2 n} \quad (3.19)$$

for some constant $c_1, c_2, c_3 > 0$. If the signal lengths n, p_0 satisfy $n > \text{poly}(p_0)$ and $p_0 > \text{polylog}(n)$, then there exist $\delta, \eta_v > 0$ such that with high probability, [Algorithm 1](#) produces $(\hat{\mathbf{a}}, \hat{\mathbf{x}})$ that are equal to the ground truth up to signed shift symmetry:

$$\|(\hat{\mathbf{a}}, \hat{\mathbf{x}}) - \sigma(s_\ell[\mathbf{a}_0], s_{-\ell}[\mathbf{x}_0])\|_2 \leq \varepsilon \quad (3.20)$$

for some $\sigma \in \{-1, 1\}$ and $\ell \in \{-p_0 + 1, \dots, p_0 - 1\}$ if $K_1 > \text{poly}(n, p_0)$ and $K_2 > \text{polylog}(n, p_0, \varepsilon^{-1})$.

Proof. See [Theorem 5.1](#) and [Theorem 5.2](#). ■

3.3 Relationship to the Literature

Blind deconvolution is a classical problem in signal processing [[SCI75](#), [Can76](#)], and has been studied under a variety of hypotheses. In this section, we first discuss the relationship between our results and the existing literature on the short-and-sparse version of this problem, and then briefly discuss other deconvolution variants in the theoretical literature.

Applications of SaS Deconvolution. The short-and-sparse model arises in a number of applications. One class of applications involves finding basic motifs (repeated patterns) in datasets. This *motif discovery* problem arises in extracellular spike sorting [[Lew98](#), [ETS11](#)] and calcium imaging [[PSG⁺16](#)], where the observed signal exhibits repetitive *short* neuron excitation patterns occurring *sparsely* across time and/or space. Similarly, electron microscopy images [[CLC⁺17](#)] arising in study of nanomaterials often exhibit repeated motifs.

⁹In practice, we suggest setting $\lambda = c_\lambda/\sqrt{p_0\theta}$ with $c_\lambda \in [0.5, 0.8]$.

Another significant application of SaS deconvolution is *image deblurring*. Typically, the blur kernel is small relative to the image size (*short*) [AD88, YK96, Car01, LFDF07, LWDF11]. In natural image deblurring, the target image is often assumed to have relatively few sharp edges [FSH⁺06, JSK08, LWDF11], and hence have *sparse* derivatives. In scientific image deblurring, e.g., in astronomy [Lan92, HHSS09, BDH⁺13] and geophysics [KT98], the target image is often sparse, either in the spatial or wavelet domains, again leading to variants of the SaS model. The literature on blind image deconvolution is large; see, e.g., [KH96, CE16] for surveys.

Variants of the SaS deconvolution problem arise in many other areas of engineering as well. Examples include *blind equalization* in communications [Sat75, SW90, JSE⁺98], *dereverberation* in sound engineering [MK88, NG10] and image *super-resolution* [BK02, SGG⁺09, YWHM10].

Algorithmic theory for SaS deconvolution. These applications have motivated a great deal of algorithmic work on variants of the SaS problem [LB87, BPSW95, BS95, KH96, MC99, CE16, WJPH17]. In contrast, relatively little theory is available to explain when and why algorithms succeed. Our algorithm minimizes φ_ρ as an approximation to the Lasso cost over the sphere. Our formulation and results have strong precedent in the literature. Lasso-like objective functions have been widely used in image deblurring [YK96, CW98, FSH⁺06, LFDF07, SJA08, XJ10, DZSW11, KTF11, LWDF11, WZ14, PF14, ZLK⁺17]. A number of insights have been obtained into the geometry of sparse deconvolution – in particular, into the effect of various constraints on \mathbf{a} on the presence or absence of spurious local minimizers. In image deblurring, a simplex constraint ($\mathbf{a} \geq \mathbf{0}$ and $\|\mathbf{a}\|_1 = 1$) arises naturally from the physical structure of the problem [YK96, CW98]. Perhaps surprisingly, simplex-constrained deconvolution admits trivial global minimizers, at which the recovered kernel \mathbf{a} is a spike, rather than the target blur kernel [LWDF11, BVG13].

[WZ14] imposes the ℓ^2 regularization on \mathbf{a} and observes that this alternative constraint gives more reliable algorithm. [ZLK⁺17] studies the geometry of the simplified objective φ_{ℓ^1} over the sphere, and proves that in the dilute limit in which \mathbf{x}_0 has one nonzero entry, all strict local minima of φ_{ℓ^1} are close to signed shifts truncations of \mathbf{a}_0 . By adopting a different objective function (based on ℓ^4 maximization) over the sphere, [ZKW18] proves that on a certain region of the sphere every local minimum is near a *truncated* signed shift of \mathbf{a}_0 , i.e., the restriction of $s_\ell[\mathbf{a}_0]$ to the window $\{0, \dots, p_0 - 1\}$. The analysis of [ZKW18] allows the sparse sequence \mathbf{x}_0 to be denser ($\theta \sim p_0^{-2/3}$ for a generic kernel \mathbf{a}_0 , as opposed to $\theta \lesssim p_0^{-3/4}$ in our result). Both [ZLK⁺17] and [ZKW18] guarantee *approximate* recovery of a portion of $s_\ell[\mathbf{a}_0]$, under complicated conditions on the kernel \mathbf{a}_0 . Our core optimization problem is very similar to [ZLK⁺17]. However, we obtain *exact* recovery of both \mathbf{a}_0 and relatively dense \mathbf{x}_0 , under the much simpler assumption of shift incoherence.

Identifiability in SaS deconvolution. Other aspects of the SaS problem have been studied theoretically. One basic question is under what circumstances the problem is identifiable, up to the scaled shift ambiguity. [CM15] shows that the problem is ill-posed for worst case $(\mathbf{a}_0, \mathbf{x}_0)$ – in particular, for certain support patterns in which \mathbf{x}_0 does not have any isolated nonzero entries. This demonstrates that *some* modeling assumptions on the support of the sparse term are needed. At the same time, this worst case structure is unlikely to occur, either under the Bernoulli model, or in practical deconvolution problems.

Other low dimensional deconvolution models. Motivated by a variety of applications, many low-dimensional deconvolution models have been studied in the theoretical literature. In communication applications, the signals \mathbf{a}_0 and \mathbf{x}_0 either live in known low-dimensional subspaces, or are sparse in some known dictionary [ARR14, LLB16, Chi16, LS15, LLB17, LS17, KK17]. These theoretical works assume that the subspace / dictionary are chosen at random. This low-dimensional deconvolution model does not exhibit the signed shift ambiguity; nonconvex formulations for this model exhibit a different structure from that studied here. In fact, the variant in which both signals belong to known subspaces can be solved by convex relaxation [ARR14]. The SaS model does not appear to be amenable to convexification, and exhibits a more complicated nonconvex geometry, due to the shift ambiguity. The main motivation for tackling this model lies in the aforementioned applications in imaging and data analysis.

[WC16, LB18] study the related *multi-instance* sparse blind deconvolution problem (MISBD), where there are K observations $\mathbf{y}_i = \mathbf{a}_0 * \mathbf{x}_i$ consisting of multiple convolutions $i = 1, \dots, K$ of a kernel \mathbf{a}_0 and different sparse vectors \mathbf{x}_i . Both works develop provable algorithms. There are several key differences with our work. First, both the proposed algorithms and their analysis require the kernel to be invertible. Second, despite the apparent similarity between the SaS model and MISBD, these problems are not equivalent. It might seem possible to reduce SaS to MISBD by dividing the single observation \mathbf{y} into K pieces; this apparent reduction fails due to boundary effects.

3.4 Notations

Vectors and indices. All vectors/matrices are written in bold font \mathbf{a}/\mathbf{A} ; indexed values are written as $\mathbf{a}_i, \mathbf{A}_{ij}$. Zeros or ones vectors are defined as $\mathbf{0}$ or $\mathbf{1}$, and i -th canonical basis vector defined as \mathbf{e}_i . The indices for vectors/matrices all start from 0 and is taking modulo- n , thus a vector of length n should has its indices labeled as $\{0, 1, \dots, n-1\}$. We write $[n] = \{0, \dots, n-1\}$. We often use capital italic symbols I, J for subsets of $[n]$. We abuse notation slightly and write $[-p] = \{n-p+1, \dots, n-1, 0\}$ and $[\pm p] = \{n-p+1, \dots, n-1, 0, 1, \dots, p-1\}$. Index sets can be labels for vectors; $\mathbf{a}_I \in \mathbb{R}^{|I|}$ denotes the restriction of the vector \mathbf{a} to coordinates I . Also, we use check symbol for reversal operator on index set $\check{I} = -I$ and vectors $\check{\mathbf{a}}_i = \mathbf{a}_{-i}$.

Operators. We let P_C denote the projection operator associated with a compact set C . The zero-filling operator $\iota_I : \mathbb{R}^{|I|} \rightarrow \mathbb{R}^n$ injects the input vector to higher dimensional Euclidean space, via $(\iota_I \mathbf{x})_i = \mathbf{x}_{I^{-1}(i)}$ for $i \in I$ and 0 otherwise. Its adjoint operator ι_I^* can be understood as subset selection operator which picks up entries of coordinates I . A common zero-filling operator through out this paper ι is abbreviation of $\iota_{[p]}$, which is often being addressed as zero-padding operator and its adjoint ι^* as truncation operator.

Convolution The convolution operator are all circular with modulo- n : $(\mathbf{a} * \mathbf{x})_i = \sum_{j \in [n]} \mathbf{a}_j \mathbf{x}_{i-j}$, also, the convolution operator works on index set: $I * J = \text{supp}(\mathbf{1}_I * \mathbf{1}_J)$. Similarly, the shift operator $s_\ell[\cdot] : \mathbb{R}^p \rightarrow \mathbb{R}^n$ is circular with modulo- n without specification: $(s_\ell[\mathbf{a}])_j = (\iota_{[p]} \mathbf{a})_{j-\ell}$. Notice that here \mathbf{a} can be shorter $p \leq n$. Let $\mathbf{C}_\mathbf{a} \in \mathbb{R}^{n \times n}$ denote a circulant matrix (with modulo- n) for vector \mathbf{a} , whose j -th column is the cyclic shift of \mathbf{a} by j : $\mathbf{C}_\mathbf{a} \mathbf{e}_j = s_j[\mathbf{a}]$. It satisfies for any $\mathbf{b} \in \mathbb{R}^n$,

$$\mathbf{C}_\mathbf{a} \mathbf{b} = \mathbf{a} * \mathbf{b}. \quad (3.21)$$

The correlation between \mathbf{a} and \mathbf{b} can be also written in similar form of convolution operator which reverse one vector before convolution. Define two correlation matrices $\mathbf{C}_\mathbf{a}^*$ and $\check{\mathbf{C}}_\mathbf{a}$ as $\mathbf{C}_\mathbf{a}^* \mathbf{e}_j = s_j[\check{\mathbf{a}}]$ and $\check{\mathbf{C}}_\mathbf{a} \mathbf{e}_j = s_{-j}[\mathbf{a}]$. The two operators will satisfy

$$\mathbf{C}_\mathbf{a}^* \mathbf{b} = \check{\mathbf{a}} * \mathbf{b}, \quad \check{\mathbf{C}}_\mathbf{a} \mathbf{b} = \mathbf{a} * \check{\mathbf{b}}. \quad (3.22)$$

4 Geometry of φ_ρ in Shift Space

Underlying our main geometric and algorithmic results is a relationship between the geometry of the function φ_ρ and the symmetries of the deconvolution problem. In this section, we describe this relationship at a more technical level, by interpreting the gradient and hessian of the function φ_ρ in terms of the shifts $s_\ell[\mathbf{a}_0]$ and stating a key lemma which asserts that a certain neighborhood of the union of subspaces $\Sigma_{4\theta p_0}$ can be decomposed into regions of negative curvature, strong gradient, and strong convexity near the target solutions $\pm s_\ell[\mathbf{a}_0]$.

4.1 Shifts and Correlations

The set $\Sigma_{4\theta p_0}$ is a union of subspaces. Any point \mathbf{a} in one of these subspaces \mathcal{S}_τ is a superposition of shifts of \mathbf{a}_0 :

$$\mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0]. \quad (4.1)$$

This representation can be extended to a general point $\mathbf{a} \in \mathbb{S}^{p-1}$ by writing

$$\mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] + \sum_{\ell \notin \tau} \alpha_\ell s_\ell[\mathbf{a}_0]. \quad (4.2)$$

The vector α can be viewed as the coefficients of a decomposition of \mathbf{a} into different shifts of \mathbf{a}_0 . This representation is not unique. For \mathbf{a} close to \mathcal{S}_τ , we can choose a particular α for which α_{τ^c} is small, a notion that we will formalize below.

For convenience, we introduce a closely related vector $\beta \in \mathbb{R}^n$, whose entries are the inner products between \mathbf{a} and the shifts of \mathbf{a}_0 : $\beta_\ell = \langle \mathbf{a}, s_\ell[\mathbf{a}_0] \rangle$. Since the columns of C_{a_0} are the shifts of \mathbf{a}_0 , we can write

$$\beta = C_{a_0}^* \iota \mathbf{a} \quad (4.3)$$

$$= C_{a_0}^* \iota^* C_{a_0} \alpha =: M \alpha. \quad (4.4)$$

The matrix M is the Gram matrix of the truncated shifts $\iota^* s_\ell[\mathbf{a}_0]$: $M_{ij} = \langle \iota^* s_i[\mathbf{a}_0], \iota^* s_j[\mathbf{a}_0] \rangle$. When μ is small, the off-diagonal elements of M are small. In particular, on \mathcal{S}_τ we may take $\alpha_{\tau^c} = \mathbf{0}$, and $\beta \approx \alpha$, in the sense that $\beta_\tau \approx \alpha_\tau$ and the entries of β_{τ^c} are small. For detailed elaboration, see [Appendix B](#).

4.2 Shifts and the Calculus of φ_{ℓ^1}

Our main geometric claims pertain to the function φ_ρ , which is based on a smooth sparsity surrogate $\rho(\cdot) \approx \|\cdot\|_1$. In this section, we sketch the main ideas of the proof as if $\rho(\cdot) = \|\cdot\|_1$, by relating the geometry of the function φ_{ℓ^1} to the vectors α, β introduced above. Working with φ_{ℓ^1} simplifies the exposition; it is also faithful to the structure of our proof, which relates the derivatives of the smooth function φ_ρ to similar quantities associated with the nonsmooth function φ_{ℓ^1} .

The function φ_{ℓ^1} has a relatively simple closed form:

$$\varphi_{\ell^1}(\mathbf{a}) = -\frac{1}{2} \|\mathcal{S}_\lambda[\check{\mathbf{y}} * \mathbf{a}]\|_2^2. \quad (4.5)$$

Here, \mathcal{S}_λ is the *soft thresholding operator*, which is defined for scalars t as $\mathcal{S}_\lambda[t] = \text{sign}(t) \max\{|t| - \lambda, 0\}$, and is extended to vectors by applying it elementwise. The operator $\mathcal{S}_\lambda[\mathbf{x}]$ shrinks the elements of \mathbf{x} towards zero. Small elements become identically zero, resulting in a sparse vector.

Gradient: Sparsifying the Correlations β

Gradient over Euclidean space. Our goal is to understand the local minimizers of the function φ_{ℓ^1} over the sphere. The function φ_{ℓ^1} is differentiable. Clearly, any point \mathbf{a} at which its gradient (over the sphere) is nonzero cannot be a local minimizer. We first give an expression for the gradient of φ_{ℓ^1} over Euclidean space \mathbb{R}^p , and then extend it to the sphere \mathbb{S}^{p-1} . Using $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ and calculus gives

$$\begin{aligned} \nabla \varphi_{\ell^1}(\mathbf{a}) &= -\iota^* C_{a_0} \check{C}_{x_0} \mathcal{S}_\lambda \left[\check{C}_{x_0} C_{a_0}^* \iota \mathbf{a} \right] \\ &= -\iota^* C_{a_0} \check{C}_{x_0} \mathcal{S}_\lambda \left[\check{C}_{x_0} \beta \right] \\ &= -\iota^* C_{a_0} \chi[\beta], \end{aligned} \quad (4.6)$$

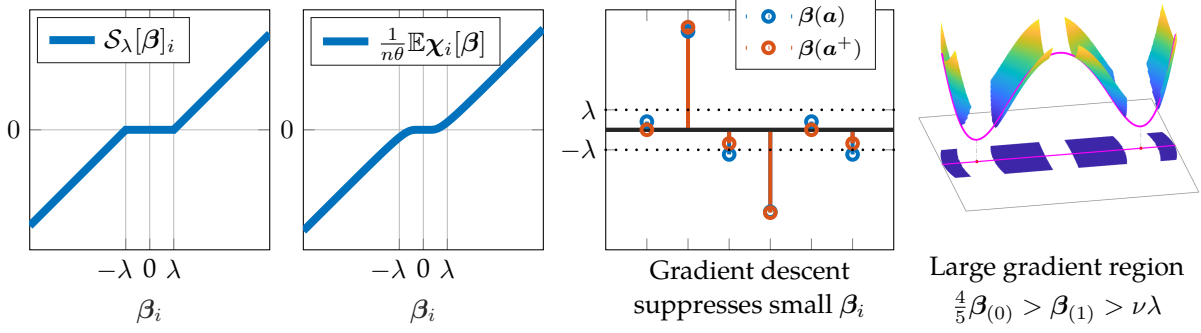


Figure 10: Gradient Sparsifies Correlations. Left: the soft thresholding operator $\mathcal{S}_\lambda[\beta]$ shrinks the entries of β towards zero, making it sparser. Middle left: the negative gradient $-\nabla\varphi_{\ell^1}$ is a superposition of shifts $s_\ell[\mathbf{a}_0]$, with coefficients $\chi_\ell[\beta] \approx \mathcal{S}_\lambda[\beta]_\ell$. Because of this, gradient descent sparsifies β . Middle right: $\beta(\mathbf{a})$ before, and $\beta(\mathbf{a}^+)$ after, one projected gradient step $\mathbf{a}^+ = \mathbf{P}_{\mathbb{S}^{p-1}}[\mathbf{a} - t \cdot \text{grad}[\varphi_{\ell^1}](\mathbf{a})]$. Notice that the small entries of β are shrunk towards zero. Right: the gradient $\text{grad}[\varphi_{\ell^1}](\mathbf{a})$ is large whenever it is easy to sparsify β ; in particular, when the largest entry $\beta_{(0)} \gg \beta_{(1)} \gg 0$.

where we have simplified the notation by introducing an operator $\chi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ as $\chi[\beta] = \widetilde{\mathcal{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda[\widetilde{\mathcal{C}}_{\mathbf{x}_0} \beta]$. This representation exhibits the (negative) gradient as a superposition of shifts of \mathbf{a}_0 with coefficients given by the entries of $\chi[\beta]$:

$$-\nabla\varphi_{\ell^1}(\mathbf{a}) = \sum_{\ell} \chi[\beta]_{\ell} s_{\ell}[\mathbf{a}_0]. \quad (4.7)$$

The operator χ appears complicated. However, its effect is relatively simple: *when \mathbf{x}_0 is a long random vector, $\chi[\beta]$ acts like a soft thresholding operator on the vector β . That is,*

$$\frac{1}{n\theta} \cdot \chi[\beta]_{\ell} \approx \begin{cases} \beta_{\ell} - \lambda, & \beta_{\ell} > \lambda \\ \beta_{\ell} + \lambda, & \beta_{\ell} < -\lambda \\ 0, & \text{otherwise} \end{cases}. \quad (4.8)$$

We show this rigorously below, in the proof of our main theorems. Here, we support this claim pictorially, by plotting the ℓ -th entry $\chi[\beta]_{\ell}$ as β_{ℓ} varies – see Figure 10 (middle left) and compare to Figure 10 (left). Because $\chi[\beta]$ suppresses small entries of β , the strongest contributions to $-\nabla\varphi_{\ell^1}$ in (4.7) will come from shifts $s_{\ell}[\mathbf{a}_0]$ with large β_{ℓ} . *In particular, the Euclidean gradient is large whenever there is a single preferred shift $s_{\ell}[\mathbf{a}_0]$, i.e., the largest entry of β is significantly larger than the second largest entry.*

Gradient over Sphere. The (Euclidean) gradient $\nabla\varphi_{\ell^1}$ measures the slope of φ_{ℓ^1} over \mathbb{R}^n . We are interested in the slope of φ_{ℓ^1} over the sphere \mathbb{S}^{p-1} , which is measured by the Riemannian gradient

$$\begin{aligned} \text{grad}[\varphi_{\ell^1}](\mathbf{a}) &= \mathbf{P}_{\mathbf{a}^\perp} \nabla\varphi_{\ell^1}(\mathbf{a}) \\ &= -\mathbf{P}_{\mathbf{a}^\perp} \sum_{\ell} \chi[\beta]_{\ell} s_{\ell}[\mathbf{a}_0]. \end{aligned} \quad (4.9)$$

The Riemannian gradient simply projects the Euclidean gradient onto the tangent space \mathbf{a}^\perp to \mathbb{S}^{p-1} at \mathbf{a} . The Riemannian gradient is large whenever

- (i) **Negative gradient points to one particular shift:** there is a single preferred shift $s_{\ell}[\mathbf{a}_0]$ so that the Euclidean gradient is large *and*

- (ii) **\mathbf{a} is not too close to any shift**: it is possible to move in the tangent space in the direction of this shift.¹⁰ Since the tangent space consists of those vectors orthogonal to \mathbf{a} , this is possible whenever $s_\ell[\mathbf{a}_0]$ is not too aligned with \mathbf{a} , i.e., \mathbf{a} is not too close to $s_\ell[\mathbf{a}_0]$.

Our technical lemma quantifies this situation in terms of the ordered entries of β . Write $|\beta_{(0)}| \geq |\beta_{(1)}| \geq \dots$, with corresponding shifts $s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0], \dots$. There is a strong gradient whenever $|\beta_{(0)}|$ is significantly larger than $|\beta_{(1)}|$ and $|\beta_{(1)}|$ is not too small compared to λ : in particular, when $\frac{4}{5}|\beta_{(0)}| > |\beta_{(1)}| > \frac{\lambda}{4 \log^2 \theta^{-1}}$. In this situation, gradient descent drives \mathbf{a} toward $s_{(0)}[\mathbf{a}_0]$, reducing $|\beta_{(1)}|, \dots$, and making the vector β sparser. We establish the technical claim that the (Euclidean) gradient of φ_{ℓ^1} sparsifies vectors in shift space in [Appendix C](#).

Hessian: Negative Curvature Breaks Symmetry

When there is no single preferred shift, i.e., when $|\beta_{(1)}|$ is close to $|\beta_{(0)}|$, the gradient can be small. Similarly, when \mathbf{a} is very close to $\pm s_{(0)}[\mathbf{a}_0]$, the gradient can be small. In either of these situations, we need to study the curvature of the function φ to determine whether there are local minimizers.

Nonsmoothness. Strictly speaking, the function φ_{ℓ^1} is not twice differentiable, due to the nonsmoothness of the soft thresholding operator $S_\lambda[t]$ at $t = \pm\lambda$. Indeed, φ_{ℓ^1} is nonsmooth at any point \mathbf{a} for which some entry of $\tilde{\mathbf{y}} * \mathbf{a}$ has magnitude λ . At other points \mathbf{a} , φ_{ℓ^1} is twice differentiable, and its Hessian is given by

$$\tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) = -\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \tilde{\mathbf{C}}_{\mathbf{x}_0} \mathbf{P}_I \tilde{\mathbf{C}}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}, \quad (4.10)$$

with $I = \text{supp}(\mathcal{S}_\lambda[\tilde{\mathbf{C}}_{\mathbf{y}} \mathbf{a}])$. We (formally) extend this expression to *every* $\mathbf{a} \in \mathbb{R}^n$, terming $\tilde{\nabla}^2 \varphi_{\ell^1}$ the *pseudo-Hessian* of φ_{ℓ^1} . For appropriately chosen smooth sparsity surrogate ρ , we will see that the (true) Hessian of the smooth function $\nabla^2 \varphi_\rho$ is close to $\tilde{\nabla}^2 \varphi_{\ell^1}$, and so $\tilde{\nabla}^2 \varphi_{\ell^1}$ yields useful information about the curvature of φ_ρ .

Curvature over Euclidean Space. As with the gradient, the Hessian is complicated, but becomes simpler when the sample size is large. The following approximation

$$\tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \approx - \sum_{\ell} s_\ell[\mathbf{a}_0] s_\ell[\mathbf{a}_0]^* \left(\frac{\partial}{\partial \beta_\ell} \chi_\ell[\beta] \right) \quad (4.11)$$

can be obtained from (4.7) noting that $\frac{\partial}{\partial \mathbf{a}} \chi_\ell[\beta] = \sum_j s_j[\mathbf{a}_0] \frac{\partial}{\partial \beta_j} \chi_\ell[\beta]$, that $\frac{\partial}{\partial \beta_j} \chi_\ell[\beta] \approx 0$ for $j \neq \ell$, and that

$$\frac{1}{n\theta} \cdot \frac{\partial \chi_\ell[\beta]}{\partial \beta_\ell} \approx \begin{cases} 0 & |\beta_\ell| \ll \lambda \\ 1 & |\beta_\ell| \gg \lambda \end{cases} \quad (4.12)$$

Again, we corroborate this approximation pictorially – see [Figure 11](#).

From this approximation, we can see that the quadratic form $\mathbf{v}^* \tilde{\nabla}^2 \varphi_{\ell^1} \mathbf{v}$ takes on a large negative value whenever \mathbf{v} is a shift $s_\ell[\mathbf{a}_0]$ corresponding to some $|\beta_\ell| \geq \lambda$, or whenever \mathbf{v} is a linear combination of such shifts. *In particular, if for some j , $|\beta_{(0)}|, |\beta_{(1)}|, \dots, |\beta_{(j)}| \gg \lambda$, then φ_{ℓ^1} will exhibit negative curvature in any direction $\mathbf{v} \in \text{span}(s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0], \dots, s_{(j)}[\mathbf{a}_0])$.*

Curvature over the Sphere. The (Euclidean) Hessian measures the curvature of the function φ_{ℓ^1} over \mathbb{R}^n . The Riemannian Hessian

$$\widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) = \mathbf{P}_{\mathbf{a}^\perp} \left(\underbrace{\tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a})}_{\text{Curvature of } \varphi_{\ell^1}} + \underbrace{\langle -\nabla \varphi_{\ell^1}(\mathbf{a}), \mathbf{a} \rangle \cdot \mathbf{I}}_{\text{Curvature of the sphere}} \right) \mathbf{P}_{\mathbf{a}^\perp}. \quad (4.13)$$

¹⁰...so the projection of the Euclidean gradient onto the tangent space does not vanish.

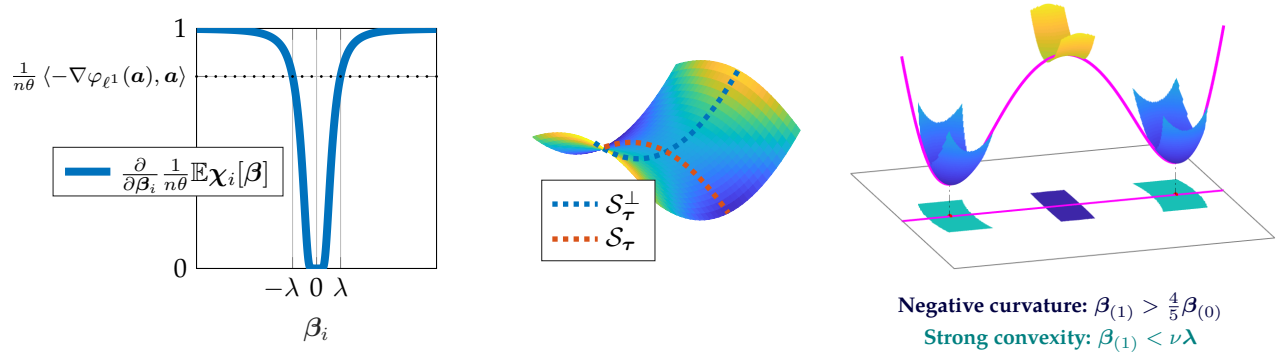


Figure 11: Hessian Breaks Symmetry. Left: contribution of $-s_i[\mathbf{a}_0]s_i[\mathbf{a}_0]^*$ to the Euclidean hessian. If $|\beta_i| \gg \lambda$ the Euclidean hessian exhibits a strong negative component in the $s_i[\mathbf{a}_0]$ direction. The Riemmanian hessian exhibits negative curvature in directions spanned by $s_i[\mathbf{a}_0]$ with corresponding $|\beta_i| \gg \lambda$ and positive curvature in directions spanned by $s_i[\mathbf{a}_0]$ with $|\beta_i| \ll \lambda$. Middle: this creates negative curvature along the subspace \mathcal{S}_τ and positive curvature orthogonal to this subspace. Right: our analysis shows that there is always a direction of negative curvature when $\beta_{(1)} > \frac{4}{5}\beta_{(0)}$; conversely when $\beta_{(1)} \ll \lambda$ there is positive curvature in every feasible direction and the function is strongly convex.

measures the curvature of φ_{ℓ^1} over the sphere. The projection $P_{\mathbf{a}^\perp}$ restricts its action to directions $\mathbf{v} \perp \mathbf{a}$ that are tangent to the sphere. The additional term $\langle -\nabla \varphi_{\ell^1}(\mathbf{a}), \mathbf{a} \rangle$ accounts for the curvature of the sphere. This term is always positive. The net effect is that directions of strong negative curvature of φ_{ℓ^1} over \mathbb{R}^n become directions of moderate negative curvature over the sphere. Directions of nearly zero curvature over \mathbb{R}^n become directions of positive curvature over the sphere. This has three implications for the geometry of φ_{ℓ^1} over the sphere:

- (i) **Negative curvature in symmetry breaking directions:** If $|\beta_{(0)}|, |\beta_{(1)}|, \dots, |\beta_{(j)}| \gg \lambda$, φ_{ℓ^1} will exhibit negative curvature in any tangent direction $\mathbf{v} \perp \mathbf{a}$ which is in the linear span

$$\text{span}(s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0], \dots, s_{(j)}[\mathbf{a}_0])$$

of the corresponding shifts of \mathbf{a}_0 .

- (ii) **Positive curvature in directions away from \mathcal{S}_τ :** The Euclidean Hessian quadratic form $\mathbf{v}^* \tilde{\nabla}^2 \varphi_{\ell^1} \mathbf{v}$ takes on relatively small values in directions orthogonal to the subspace \mathcal{S}_τ . The Riemannian Hessian is positive in these directions, creating positive curvature orthogonal to the subspace \mathcal{S}_τ .
- (iii) **Strong convexity around minimizers:** Around a minimizer $s_\ell[\mathbf{a}_0]$, only a single entry β_ℓ is large. Any tangent direction $\mathbf{v} \perp \mathbf{a}$ is nearly orthogonal to the subspace $\text{span}(s_\ell[\mathbf{a}_0])$, and hence is a direction of positive (Riemmanian) curvature. The objective function φ_ρ is strongly convex around the target solutions $\pm s_\ell[\mathbf{a}_0]$.

Figure 11 visualizes these regions of negative and positive curvature, and the technical claim of positivity/negativity of curvature in shift space is presented in detail in [Appendix D](#).

4.3 Any Local Minimizer is a Near Shift

We close this section by stating a key theorem, which makes the above discussion precise. We will show that a certain neighborhood of any subspace \mathcal{S}_τ can be covered by regions of *negative curvature*, *large gradient*, and regions of *strong convexity* containing target solutions $\pm s_\ell[\mathbf{a}_0]$. Furthermore, at the boundary of this neighborhood, the negative gradient points back—*retracts*—toward the subspace \mathcal{S}_τ , due to the (directional) convexity of φ_ρ away from the subspace.

Widened subspace region. To formally state the result, we need a way of measuring how close \mathbf{a} is to the subspace \mathcal{S}_τ . For technical reasons, it turns out to be convenient to do this in terms of the coefficients α in the representation

$$\mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] + \sum_{\ell' \in \tau^c} \alpha_{\ell'} s_{\ell'}[\mathbf{a}_0]. \quad (4.14)$$

If $\mathbf{a} \in \mathcal{S}_\tau$, we can take α with $\alpha_{\tau^c} = \mathbf{0}$. We can view the energy $\|\alpha_{\tau^c}\|_2$ as a measure of the distance from \mathbf{a} to \mathcal{S}_τ . A technical wrinkle arises, because the representation (4.14) is not unique. We resolve this issue by choosing the α that minimizes $\|\alpha_{\tau^c}\|_2$, writing:

$$d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = \inf \{ \|\alpha_{\tau^c}\|_2 : \sum_\ell \alpha_\ell s_\ell[\mathbf{a}_0] = \mathbf{a} \}. \quad (4.15)$$

The distance $d_\alpha(\mathbf{a}, \mathcal{S}_\tau)$ is zero for $\mathbf{a} \in \mathcal{S}_\tau$. Our analysis controls the geometric properties of φ_ρ over the set of \mathbf{a} for which $d_\alpha(\mathbf{a}, \mathcal{S}_\tau)$ is not too large. Similar to (3.3), we define an object which contains all points that are close to some \mathcal{S}_τ , in the above sense:

$$\Sigma_{4\theta p_0}^\gamma := \bigcup_{|\tau| \leq 4\theta p_0} \{ \mathbf{a} : d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma \}. \quad (4.16)$$

The aforementioned geometric properties hold over this set:

Theorem 4.1 (Three subregions). *Suppose that $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ where $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$ is μ -shift coherent and $\mathbf{x}_0 \sim \text{i.i.d.}$ $\text{BG}(\theta) \in \mathbb{R}^n$ satisfying*

$$\theta \in \left[\frac{c'}{p_0}, \frac{c}{p_0 \sqrt{\mu} + \sqrt{p_0}} \right] \cdot \frac{1}{\log^2 p_0} \quad (4.17)$$

for some constants $c', c > 0$. Set $\lambda = 0.1/\sqrt{p_0 \theta}$ in φ_ρ where $\rho(x) = \sqrt{x^2 + \delta^2}$. There exist numerical constants $C, c'', c''', c_1, c_4 > 0$ such that if $\delta \leq \frac{c'' \lambda \theta^8}{p^2 \log^2 n}$ and $n > C p_0^5 \theta^{-2} \log p_0$, then with probability at least $1 - c'''/n$, for every $\mathbf{a} \in \Sigma_{4\theta p_0}^\gamma$, we have:

- (Negative curvature): If $|\beta_{(1)}| \geq \nu_1 |\beta_{(0)}|$, then

$$\lambda_{\min}(\text{Hess}[\varphi_\rho](\mathbf{a})) \leq -c_1 n \theta \lambda; \quad (4.18)$$

- (Large gradient): If $\nu_1 |\beta_{(0)}| \geq |\beta_{(1)}| \geq \nu_2(\theta) \lambda$, then

$$\|\text{grad}[\varphi_\rho](\mathbf{a})\|_2 \geq c_2 n \theta \frac{\lambda^2}{\log^2 \theta^{-1}}; \quad (4.19)$$

- (Convex near shifts): If $\nu_2(\theta) \lambda \geq |\beta_{(1)}|$, then

$$\text{Hess}[\varphi_\rho](\mathbf{a}) \succ c_3 n \theta \mathbf{P}_{\mathbf{a}^\perp}; \quad (4.20)$$

- (Retraction to subspace): If $\frac{\gamma}{2} \leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma$, then for every α satisfying $\mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha$, there exists ζ satisfying $\text{grad}[\varphi_\rho](\mathbf{a}) = \iota^* C_{\mathbf{a}_0} \zeta$, such that

$$\langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq c_4 \|\zeta_{\tau^c}\|_2 \|\alpha_{\tau^c}\|_2; \quad (4.21)$$

- (Local minimizers): If \mathbf{a} is a local minimizer,

$$\min_{\substack{\ell \in [\pm p] \\ \sigma \in \{\pm 1\}}} \|\mathbf{a} - \sigma s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max \{ \mu, p_0^{-1} \}, \quad (4.22)$$

where $\nu_1 = \frac{4}{5}$, $\nu_2(\theta) = \frac{1}{4 \log^2 \theta^{-1}}$ and $\gamma = \frac{c \cdot \text{poly}(\sqrt{1/\theta}, \sqrt{1/\mu})}{\log^2 \theta^{-1}} \cdot \frac{1}{\sqrt{p_0}}$.

Proof. See [Appendix F.5](#). ■

The retraction property elaborated in (4.21) implies that the negative gradient at \mathbf{a} points in a direction that decreases $d_\alpha(\mathbf{a}, \mathcal{S}_\tau)$. This is a consequence of positive curvature away from \mathcal{S}_τ . It essentially implies that the gradient is monotone in α_{τ^c} space: choose any $\underline{\mathbf{a}} \in \mathcal{S}_\tau \cap \mathbb{S}^{p-1}$, write $\underline{\alpha}$ to be its coefficient, and let $\underline{\zeta}$ be the coefficient of $\text{grad}[\varphi_\rho](\underline{\mathbf{a}})$. Then $\underline{\alpha}_{\tau^c} = \mathbf{0}$, $\underline{\zeta}_{\tau^c} \approx \mathbf{0}$ and

$$\langle \underline{\zeta}_{\tau^c} - \underline{\zeta}_{\tau^c}, \alpha_{\tau^c} - \underline{\alpha}_{\tau^c} \rangle \approx \langle \underline{\zeta}_{\tau^c} - \mathbf{0}, \alpha_{\tau^c} - \mathbf{0} \rangle = \langle \underline{\zeta}_{\tau^c}, \alpha_{\tau^c} \rangle > 0.$$

Our main geometric claim in [Theorem 3.1](#) is a direct consequence of [Theorem 4.1](#). Moreover, it suggests that as long as we can minimize φ_ρ within the region $\Sigma_{4\theta p_0}^\gamma$, we will solve the SaS deconvolution problem.

5 Provable Algorithm

In light of [Theorem 4.1](#), in this section we introduce a two-part algorithm [Algorithm 1](#), which first applies the curvilinear descent method to find a local minimum of φ_ρ within $\Sigma_{4\theta p_0}^\gamma$, followed by refinement algorithm that uses alternating minimization to exactly recover the ground truth. This algorithm exactly solves SaS deconvolution problem.

5.1 Minimization

There are three major issues in finding a local minimizer within $\Sigma_{4\theta p_0}^\gamma$. We want ...

- (i) **Initialization.** the initializer $\mathbf{a}^{(0)}$ to reside within $\Sigma_{4\theta p_0}^\gamma$,
- (ii) **Negative curvature.** the method to avoid stagnating near the saddle points of φ_ρ ,
- (iii) **No exit.** the descent method to remain inside $\Sigma_{4\theta p_0}^\gamma$.

In the following paragraphs, we describe how our proposed algorithm achieves the above desiderata.

Initialization within $\Sigma_{4\theta p_0}^\gamma$. Our data-driven initialization scheme produces $\mathbf{a}^{(0)}$, where

$$\begin{aligned} \mathbf{a}^{(0)} &= -P_{\mathbb{S}^{p-1}} \nabla \varphi_\rho (P_{\mathbb{S}^{p-1}} [\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}]) \\ &= -P_{\mathbb{S}^{p-1}} \nabla \varphi_\rho P_{\mathbb{S}^{p-1}} [P_{[p_0]}(\mathbf{a}_0 * \mathbf{x}_0)] , \\ &\approx -P_{\mathbb{S}^{p-1}} \nabla \varphi_\rho [P_{[p_0]}(\mathbf{a}_0 * \tilde{\mathbf{x}}_0)] , \end{aligned}$$

is the normalized gradient vector from a chunk of data $\mathbf{a}^{(-1)} := P_{[p_0]}(\mathbf{a}_0 * \tilde{\mathbf{x}}_0)$ with $\tilde{\mathbf{x}}_0$ a normalized Bernoulli-Gaussian random vector of length $2p_0 - 1$. Since $\nabla \varphi_\rho \approx \nabla \varphi_{\ell^1}$, expand the gradient $\nabla \varphi_{\ell^1}$ and rewrite the gradient $\nabla_{\ell^1}(\mathbf{a}^{(-1)})$ in shift space, we get

$$\begin{aligned} -\nabla \varphi_{\rho^1}(\mathbf{a}^{(-1)}) &\approx \boldsymbol{\iota}^* C_{\mathbf{a}_0} \check{C}_{\mathbf{x}_0} \mathcal{S}_\lambda [\check{C}_{\mathbf{x}_0} C_{\mathbf{a}_0}^* P_{[p_0]}(\mathbf{a}_0 * \tilde{\mathbf{x}}_0)] \\ &= \boldsymbol{\iota}^* C_{\mathbf{a}_0} \chi [C_{\mathbf{a}_0}^* P_{[p_0]} C_{\mathbf{a}_0} \tilde{\mathbf{x}}_0] \\ &\approx \boldsymbol{\iota}^* C_{\mathbf{a}_0} \chi [\tilde{\mathbf{x}}_0] \\ &\approx n\theta \cdot \boldsymbol{\iota}^* C_{\mathbf{a}_0} \mathcal{S}_\lambda [\tilde{\mathbf{x}}_0] , \end{aligned}$$

where the approximation in the third equation is accurate if the truncated shifts are incoherent

$$\max_{i \neq j} |\langle \boldsymbol{\iota}_{p_0}^* s_i[\mathbf{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu \ll 1. \quad (5.1)$$

With this simple approximation, it comes clear that the coefficients (in shift space) of initializer $\mathbf{a}^{(0)}$,

$$\mathbf{a}^{(0)} \approx \mathbf{P}_{\mathbb{S}^{p-1}} \mathbf{t}^* \mathbf{C}_{\mathbf{a}_0} \mathcal{S}_\lambda [\tilde{\mathbf{x}}_0], \quad (5.2)$$

approximate $\mathcal{S}_\lambda [\tilde{\mathbf{x}}_0]$, which resides near the subspace \mathcal{S}_τ , in which τ contains the nonzero entries of $\tilde{\mathbf{x}}_0$ on $\{-p_0 + 1, \dots, p_0 - 1\}$. With high probability, the number of non-zero entries is $|\tau| \lesssim 4\theta p_0$, we therefore conclude that our initializer $\mathbf{a}^{(0)}$ satisfies

$$\mathbf{a}^{(0)} \in \Sigma_{4\theta p_0}^\gamma. \quad (5.3)$$

Furthermore, since $\tilde{\mathbf{x}}_0$ is normalized, the largest magnitude for entries of $|\tilde{\mathbf{x}}_0|$ is likely to be around $1/\sqrt{2p_0\theta}$. To ensure that $\mathcal{S}_\lambda [\tilde{\mathbf{x}}_0]$ does not annihilate all nonzero entries of $\tilde{\mathbf{x}}_0$ (otherwise our initializer $\mathbf{a}^{(0)}$ will become $\mathbf{0}$), the ideal λ should be slightly less than the largest magnitude of $|\tilde{\mathbf{x}}_0|$. We suggest setting λ in φ_ρ as

$$\lambda = \frac{c}{\sqrt{p_0\theta}}. \quad (5.4)$$

for some $c \in (0, 1)$.

Minimize φ_ρ within $\Sigma_{4\theta p_0}^\gamma$. Many methods have been proposed to optimize functions whose saddle points exhibit strict negative curvature, including the noisy gradient method [GHJY15], trust region methods [AMS09, SQW17] and curvilinear search [WY13]. Any of the above methods can be adapted to minimize φ_ρ . In this paper, we use *curvilinear method with restricted stepsize* to demonstrate how to analyze an optimization problem using the geometric properties of φ_ρ over $\Sigma_{4\theta p_0}^\gamma$ – in particular, negative curvature in symmetry-breaking directions and positive curvature away from \mathcal{S}_τ .

Curvilinear search uses an update strategy that combines the gradient \mathbf{g} and a direction of negative curvature \mathbf{v} , which here we choose as an eigenvector of the hessian \mathbf{H} with smallest eigenvalue, scaled such that $\mathbf{v}^* \mathbf{g} \geq 0$. In particular, we set

$$\mathbf{a}^+ \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}] \quad (5.5)$$

For small t ,

$$\varphi(\mathbf{a}^+) \approx \varphi(\mathbf{a}) + \langle \mathbf{g}, \boldsymbol{\xi} \rangle + \frac{1}{2} \boldsymbol{\xi}^* \mathbf{H} \boldsymbol{\xi}. \quad (5.6)$$

Since $\boldsymbol{\xi}$ converges to $\mathbf{0}$ only if \mathbf{a} converges to the local minimizer (otherwise either gradient \mathbf{g} is nonzero or there is a negative curvature direction \mathbf{v}), this iteration produces a local minimizer for φ_ρ , whose saddle points near any \mathcal{S}_τ has negative curvature, we just need to ensure all iterates stays near some such subspace. We prove this by showing:

- When $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma$, curvilinear steps move a small distance away from the subspace:

$$|d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) - d_\alpha(\mathbf{a}, \mathcal{S}_\tau)| \leq \frac{\gamma}{2}. \quad (5.7)$$

- When $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \in [\frac{\gamma}{2}, \gamma]$, curvilinear steps retract toward subspace:

$$d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) \leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau). \quad (5.8)$$

Together, we can prove that the iterates $\mathbf{a}^{(k)}$ converge to a minimizer, and

$$\forall k = 1, 2, \dots, \quad \mathbf{a}^{(k)} \in \Sigma_{4\theta p_0}^\gamma. \quad (5.9)$$

We conclude this section with the following theorem:

Theorem 5.1 (Convergence of retractive curvilinear search). *Suppose signals $\mathbf{a}_0, \mathbf{x}_0$ satisfy the conditions of Theorem 4.1, $\theta > 10^3 c/p_0$ ($c > 1$), and \mathbf{a}_0 is μ -truncated shift coherent $\max_{i \neq j} |\langle \mathbf{t}_{p_0}^* s_i[\mathbf{a}_0], \mathbf{t}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu$. Write $\mathbf{g} = \text{grad}[\varphi_\rho](\mathbf{a})$ and $\mathbf{H} = \text{Hess}[\varphi_\rho](\mathbf{a})$. When the smallest eigenvalue of \mathbf{H} is strictly smaller than $-\eta_v$, let \mathbf{v} be the unit eigenvector of smallest eigenvalue, scaled so $\mathbf{v}^* \mathbf{g} \geq 0$; otherwise let $\mathbf{v} = \mathbf{0}$. Define a sequence $\{\mathbf{a}^{(k)}\}_{k \in \mathbb{N}}$ where $\mathbf{a}^{(0)}$ equals (3.7) and for $k = 1, 2, \dots, K_1$:*

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\mathbf{a}^{(k)} - t \mathbf{g}^{(k)} - t^2 \mathbf{v}^{(k)} \right] \quad (5.10)$$

with largest $t \in (0, \frac{0.1}{n\theta}]$ satisfying Armijo steplength:

$$\varphi_\rho(\mathbf{a}^{(k+1)}) < \varphi_\rho(\mathbf{a}^{(k)}) - \frac{1}{2} \left(t \|\mathbf{g}^{(k)}\|_2^2 + \frac{1}{2} t^4 \eta_v \|\mathbf{v}^{(k)}\|_2^2 \right), \quad (5.11)$$

then with probability at least $1 - 1/c$, there exists some signed shift $\bar{\mathbf{a}} = \pm s_i[\mathbf{a}_0]$ where $i \in [\pm p_0]$ such that $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \leq \mu + 1/p$ for all $k \geq K_1 = \text{poly}(n, p)$. Here, $\eta_v = c' n \theta \lambda$ for some $c' < c_1$ in Theorem 4.1.

Proof. See Appendix G.2. ■

5.2 Local Refinement

In this section, we describe and analyze an algorithm which refines an estimate $\bar{\mathbf{a}} \approx \mathbf{a}_0$ of the kernel to exactly recover $(\mathbf{a}_0, \mathbf{x}_0)$. Set

$$\mathbf{a}^{(0)} \leftarrow \bar{\mathbf{a}}, \quad \lambda^{(0)} \leftarrow C(p\theta + \log n)(\mu + 1/p), \quad I^{(0)} \leftarrow \text{supp}(\mathcal{S}_\lambda[C_{\bar{\mathbf{a}}}^* \mathbf{y}]). \quad (5.12)$$

We alternatively minimize the Lasso objective with respect to \mathbf{a} and \mathbf{x} :

$$\mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|, \quad (5.13)$$

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\underset{\mathbf{a}}{\text{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right], \quad (5.14)$$

$$\lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}). \quad (5.15)$$

One departure from standard alternating minimization procedures is our use of a continuation method, which (i) decreases λ and (ii) maintains a running estimate $I^{(k)}$ of the support set. Our analysis will show that $\mathbf{a}^{(k)}$ converges to one of the signed shifts of \mathbf{a}_0 at a linear rate, in the sense that

$$\min_{\sigma \in \pm 1, \ell \in [\pm p_0]} \|\mathbf{a}^{(k)} - \sigma \cdot s_\ell[\mathbf{a}_0]\|_2 \leq C' 2^{-k}. \quad (5.16)$$

Modified coherence and support density assumptions It should be clear that exact recovery is unlikely if \mathbf{x}_0 contains many consecutive nonzero entries: in fact in this situation, even *non-blind* deconvolution fails. Therefore to obtain exact recovery it is necessary to put an upper bound on signal dimension n . Here, we introduce the notation κ_I as an upper bound for number of nonzero entries of \mathbf{x}_0 in a length- p window:

$$\kappa_I := 6 \max\{\theta p, \log n\}, \quad (5.17)$$

where the indexing and addition should be interpreted modulo n . We will denote the support sets of true sparse vector \mathbf{x}_0 and recovered $\mathbf{x}^{(k)}$ in the intermediate k -th steps as

$$I = \text{supp}(\mathbf{x}_0), \quad I^{(k)} = \text{supp}(\mathbf{x}^{(k)}), \quad (5.18)$$

then in the Bernoulli-Gaussian model, with high probability,

$$\max_{\ell} |I \cap ([p] + \ell)| \leq \kappa_I. \quad (5.19)$$

The $\log n$ term reflects the fact that as n becomes enormous (exponential in p) eventually it becomes likely that some length- p window of \mathbf{x}_0 is densely occupied. In our main theorem statement, we preclude this possibility by putting an upper bound on signal length n with respect to window length p and shift coherence μ . We will assume

$$(\mu + 1/p) \cdot \kappa_I^2 < c \quad (5.20)$$

for some numerical constant $c \in (0, 1)$.

Alternating minimization produces \mathbf{a} that contracts toward \mathbf{a}_0 . Recall that (B.1) in Theorem 4.1 provides that

$$\|\bar{\mathbf{a}} - \mathbf{a}_0\|_2 \leq (\mu + 1/p), \quad (5.21)$$

which is sufficiently close to \mathbf{a}_0 as long as (5.19) holds true. Here, we will elaborate this by showing a single iteration of alternating minimization algorithm (5.13)-(5.15) is a contraction mapping for \mathbf{a} toward \mathbf{a}_0 .

To this end, at k -th iteration, write $T = I^{(k)}$, $J = I^{(k+1)}$ and $\boldsymbol{\sigma}^{(k)} = \text{sign}(\mathbf{x}^{(k)})$, then first observe that the solution to the reweighted Lasso problem (5.13) can be written as

$$\mathbf{x}^{(k+1)} = \boldsymbol{\iota}_J \left(\boldsymbol{\iota}_J^* \mathbf{C}_{\mathbf{a}^{(k)}}^* \mathbf{C}_{\mathbf{a}^{(k)}} \boldsymbol{\iota}_J \right)^{-1} \boldsymbol{\iota}_J^* \left(\mathbf{C}_{\mathbf{a}^{(k)}}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \lambda^{(k)} \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}^{(k+1)} \right), \quad (5.22)$$

and the solution to least squares problem (5.14) will be

$$\mathbf{a}^{(k+1)} = \left(\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}^{(k+1)}} \boldsymbol{\iota} \right)^{-1} \left(\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 \right). \quad (5.23)$$

Here, we are going to illustrate the relationship between $\mathbf{a}^{(k+1)} - \mathbf{a}_0$ and $\mathbf{a}^{(k)} - \mathbf{a}_0$ using simple approximations. First, let us assume that $\mathbf{a}^{(k)} \approx \mathbf{a}_0$, $\mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \approx \mathbf{I}$, and $\mathbf{I} \approx \mathbf{J} \approx \mathbf{T}$. Then (5.22) gives

$$\mathbf{x}^{(k+1)} \approx \mathbf{x}_0, \quad (5.24)$$

$$\begin{aligned} (\mathbf{x}^{(k+1)} - \mathbf{x}_0) &\approx \mathbf{P}_I \left(\mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}^{(k)}} \mathbf{x}_0 \right) \\ &\approx \mathbf{P}_I \left[\mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}_0 - \mathbf{a}^{(k)}) \right], \end{aligned} \quad (5.25)$$

which implies, while assuming $\mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{x}_0} \approx n\theta \mathbf{I}$, that from (5.23):

$$\begin{aligned} (\mathbf{a}^{(k+1)} - \mathbf{a}_0) &\approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}^{(k+1)}} \boldsymbol{\iota} \mathbf{a}_0 \\ &\approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{x}_0 - \mathbf{x}^{(k+1)}) \\ &\approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}^{(k)} - \mathbf{a}_0). \end{aligned} \quad (5.26)$$

Now since $\mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \approx n\theta \mathbf{e}_0 \mathbf{e}_0^*$, this suggests that $(n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}$ approximates a contraction mapping with fixed point \mathbf{a}_0 , as follows:

$$\begin{aligned} (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} &\approx \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{e}_0 \mathbf{e}_0^* \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \\ &\approx \mathbf{a}_0 \mathbf{a}_0^*. \end{aligned} \quad (5.27)$$

Hence, if we can ensure all above approximation is sufficiently and increasingly accurate as the iterate proceeds, the alternating minimization essentially is a power method which finds the leading eigenvector of matrix $\mathbf{a}_0 \mathbf{a}_0^*$ —and the solution to this algorithm is apparently \mathbf{a}_0 . Indeed, we prove that the iterates produced by this sequence of operations converge to the ground truth at a linear rate, as long as it is initialized sufficiently nearby:

Theorem 5.2 (Linear rate convergence of alternating minimization). *Suppose $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ where \mathbf{a}_0 is μ -shift coherent and $\mathbf{x}_0 \sim \text{BG}(\theta)$, then there exists some constants C, c, c_μ such that if $(\mu + 1/p) \kappa_I^2 < c_\mu$ and $n > C\theta^{-2}p^2 \log n$, then with probability at least $1 - c/n$, for any starting point $\mathbf{a}^{(0)}$ and $\lambda^{(0)}, I^{(0)}$ such that*

$$\|\mathbf{a}^{(0)} - \mathbf{a}_0\|_2 \leq \mu + 1/p, \quad \lambda^{(0)} = 5\kappa_I(\mu + 1/p), \quad I^{(0)} = \text{supp}(\mathbf{C}_{\mathbf{a}^{(0)}}^* \mathbf{y}), \quad (5.28)$$

and for $k = 1, 2, \dots$:

$$\mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|, \quad (5.29)$$

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\underset{\mathbf{a}}{\text{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right], \quad (5.30)$$

$$\lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}) \quad (5.31)$$

then

$$\|\mathbf{a}^{(k+1)} - \mathbf{a}_0\|_2 \leq (\mu + 1/p) 2^{-k} \quad (5.32)$$

for every $k = 0, 1, 2, \dots$.

Proof. See [Appendix H.3](#). ■

Remark 5.3. The estimates $\mathbf{x}^{(k)}$ also converges to the ground truth \mathbf{x}_0 at a linear rate.

6 Experiments

We demonstrate that the tradeoffs between the motif length p_0 and sparsity rate θ produce a transition region for successful SaS deconvolution under generic choices of \mathbf{a}_0 and \mathbf{x}_0 . For fixed values of $\theta \in [10^{-3}, 10^{-2}]$ and $p_0 \in [10^3, 10^4]$, we draw 50 instances of synthetic data by choosing $\mathbf{a}_0 \sim \text{Unif}(\mathbb{S}^{p_0-1})$ and $\mathbf{x}_0 \in \mathbb{R}^n$ with $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ where $n = 5 \times 10^5$. Note that choosing \mathbf{a}_0 this way implies $\mu(\mathbf{a}_0) \approx \frac{1}{\sqrt{p_0}}$.

For each instance, we recover \mathbf{a}_0 and \mathbf{x}_0 from $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ by minimizing problem (2.5). For ease of computation, we modify [Algorithm 1](#) by replacing curvilinear search with *accelerated Riemannian gradient descent* method ([Algorithm 2](#)), which is an adaptation of accelerated gradient descent [\[BT09\]](#) to the sphere. In particular, we apply momentum and increment by the Riemannian gradient via the exponential and logarithmic operators

$$\text{Exp}_{\mathbf{a}}(\mathbf{u}) := \cos(\|\mathbf{u}\|_2) \cdot \mathbf{a} + \sin(\|\mathbf{u}\|_2) \cdot \frac{\mathbf{u}}{\|\mathbf{u}\|_2}, \quad (6.1)$$

$$\text{Log}_{\mathbf{a}}(\mathbf{b}) := \arccos(\langle \mathbf{a}, \mathbf{b} \rangle) \cdot \frac{\mathbf{P}_{\mathbf{a}^\perp}(\mathbf{b} - \mathbf{a})}{\|\mathbf{P}_{\mathbf{a}^\perp}(\mathbf{b} - \mathbf{a})\|_2}, \quad (6.2)$$

derived from [\[AMS09\]](#). Here $\text{Exp}_{\mathbf{a}} : \mathbf{a}^\perp \rightarrow \mathbb{S}^{p-1}$ takes a tangent vector of \mathbf{a} and produces a new point on the sphere, whereas $\text{Log}_{\mathbf{a}} : \mathbb{S}^{p-1} \rightarrow \mathbf{a}^\perp$ takes a point $\mathbf{b} \in \mathbb{S}^{p-1}$ and returns the tangent vector which points from \mathbf{a} to \mathbf{b} .

For each recovery instance, we say the local minimizer \mathbf{a}_{\min} generated from [Algorithm 2](#) is sufficiently close to a solution of SaS deconvolution problem, if

$$\text{success}(\mathbf{a}_{\min}, \mathbf{a}_0) := \{ \max_{\ell} |\langle \mathbf{s}_\ell[\mathbf{a}_0], \mathbf{a}_{\min} \rangle| > 0.95 \}. \quad (6.3)$$

The result is shown in [Figure 12](#). Our source code can be accessed via the following address:

https://github.com/sbdsphere/sbd_experiments.git

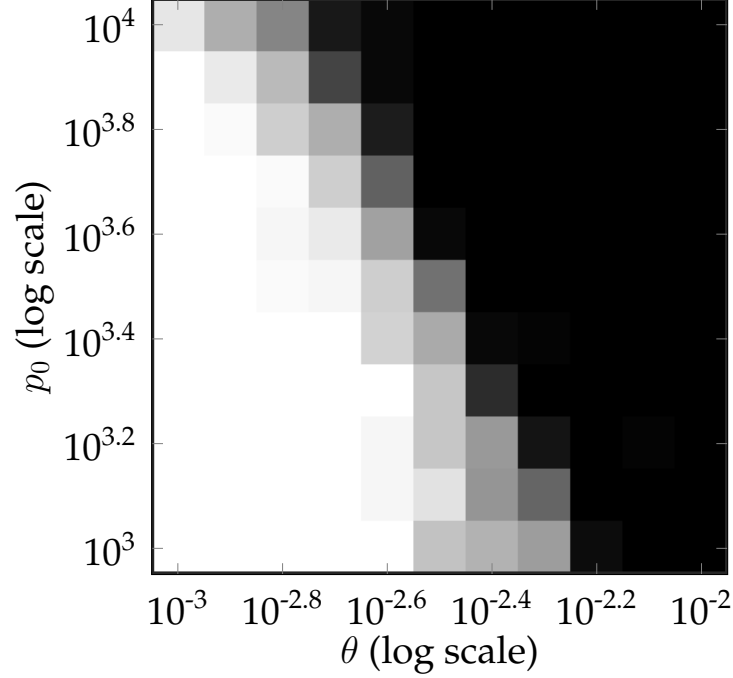


Figure 12: Success probability of SaS deconvolution under generic a_0, x_0 with varying kernel length p_0 , and sparsity rate θ . When sparsity rate decreases sufficiently with respect to kernel length, successful recovery becomes very likely (brighter), and vice versa (darker). A transition line is shown with slope $\frac{\log p_0}{\log \theta} \approx -2$, implying Algorithm 2 works with high probability when $\theta \lesssim \frac{1}{\sqrt{p_0}}$ in generic case.

Algorithm 2 SaS deconvolution with Accelerated Riemannian gradient descent

Input: Observation y , sparsity penalty $\lambda = 0.5/\sqrt{p_0\theta}$, momentum parameter $\eta \in [0, 1)$.

Initialize $a^{(0)} \leftarrow -P_{\mathbb{S}^{p-1}} \nabla \varphi_\rho (P_{\mathbb{S}^{p-1}} [0^{p_0-1}; [y_0, \dots, y_{p_0-1}]; 0^{p_0-1}])$,

for $k = 1, 2, \dots, K$ **do**

 Get momentum: $w \leftarrow \text{Exp}_{a^{(k)}}(\eta \cdot \text{Log}_{a^{(k-1)}}(a^{(k)}))$.

 Get negative gradient direction: $g \leftarrow -\text{grad}[\varphi_\rho](w)$.

 Armijo step $a^{(k+1)} \leftarrow \text{Exp}_w(tg)$, choosing $t \in (0, 1)$ s.t. $\varphi_\rho(a^{(k+1)}) - \varphi_\rho(w) < -t \|g\|_2^2$.

end for

Output: Return $a^{(K)}$.

7 Discussion

In this section, we close by discussing several of the most important limitations of our results, and highlighting corresponding directions for future work.

Minimizing φ_ρ does not accurately recover coherent kernels. The main drawback of our proposed method is that it does not succeed when the target motif a_0 has shift coherence very close to 1. For instance, a common scenario in image blind deconvolution involves deblurring an image with a smooth, low-pass point spread function (e.g., Gaussian blur). Both our analysis and numerical experiments show that in this situation minimizing φ_ρ does not find the generating signal pairs (a_0, x_0) consistently—the minimizer of φ_ρ is often spurious and is not close to any particular shift of a_0 . We do not suggest minimizing φ_ρ in this situation. On the other hand, minimizing the bilinear lasso objective φ_{lasso} over the sphere often succeeds even if the true signal pair (a_0, x_0) is coherent and dense.

Relation of φ_ρ to Bilinear Lasso. In light of the above observations, we view the analysis of the bilinear lasso as the most important direction for future theoretical work on SaS deconvolution. The drop quadratic formulation studied here has commonalities with the bilinear lasso: both exhibit local minima at signed shifts, and both exhibit negative curvature in symmetry breaking directions. A major difference (and hence, major challenge) is that gradient methods for bilinear lasso do not retract to a union of subspaces – they retract to a more complicated, nonlinear set.

Suboptimality in the analysis. Finally, there are several directions in which our analysis could be improved. Our lower bounds on the length n of the random vector x_0 required for success are clearly suboptimal. We also suspect our sparsity-coherence tradeoff between μ, θ (roughly, $\theta \lesssim 1/(\sqrt{\mu}p_0)$) is suboptimal, even for the φ_ρ objective. Articulating optimal sparsity-coherence tradeoffs for is another interesting direction for future work.

Acknowledgement

The authors gratefully acknowledge support from NSF 1343282, NSF CCF 1527809, and NSF IIS 1546411.

References

- [AD88] GR Ayers and J Christopher Dainty. Iterative blind deconvolution method and its applications. *Optics letters*, 13(7):547–549, 1988.
- [AMS09] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [ARR14] Ali Ahmed, Benjamin Recht, and Justin Romberg. Blind deconvolution using convex programming. *IEEE Transactions on Information Theory*, 60(3):1711–1732, 2014.
- [BC11] Heinz H. Bauschke and Patrick L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer Publishing Company, Incorporated, 1st edition, 2011.
- [BDH⁺13] David Briers, Donald D Duncan, Evan R Hirst, Sean J Kirkpatrick, Marcus Larsson, Wiendelt Steenbergen, Tomas Stromberg, and Oliver B Thompson. Laser speckle contrast imaging: theoretical and practical limitations. *Journal of biomedical optics*, 18(6):066018, 2013.
- [BK02] Simon Baker and Takeo Kanade. Limits on super-resolution and how to break them. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(9):1167–1183, 2002.
- [BPSW95] PJ Bones, CR Parker, BL Satherley, and RW Watson. Deconvolution and phase retrieval with use of zero sheets. *JOSA A*, 12(9):1842–1857, 1995.
- [BS95] Anthony J Bell and Terrence J Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural computation*, 7(6):1129–1159, 1995.
- [BT09] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [BVG13] Alexis Benichoux, Emmanuel Vincent, and Rémi Gribonval. A fundamental pitfall in blind deconvolution with sparse and shift-invariant priors. In *ICASSP-38th International Conference on Acoustics, Speech, and Signal Processing-2013*, 2013.
- [Can76] Michael Cannon. Blind deconvolution of spatially invariant image blurs with phase. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 24(1):58–63, 1976.
- [Car01] Alfred S Carasso. Direct blind deconvolution. *SIAM Journal on Applied Mathematics*, 61(6):1980–2007, 2001.
- [CE16] Patrizio Campisi and Karen Egiazarian. *Blind image deconvolution: theory and applications*. CRC press, 2016.
- [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.

- [CM15] Sunav Choudhary and Urbashi Mitra. Fundamental limits of blind deconvolution part ii: Sparsity-ambiguity trade-offs. *arXiv preprint arXiv:1503.03184*, 2015.
- [CW98] Tony F Chan and Chiu-Kwong Wong. Total variation blind deconvolution. *IEEE transactions on Image Processing*, 7(3):370–375, 1998.
- [CWB08] Emmanuel J Candes, Michael B Wakin, and Stephen P Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. *Journal of Fourier analysis and applications*, 14(5-6):877–905, 2008.
- [DZSW11] Weisheng Dong, Lei Zhang, Guangming Shi, and Xiaolin Wu. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*, 20(7):1838–1857, 2011.
- [EHJ⁺04] Bradley Efron, Trevor Hastie, Iain Johnstone, Robert Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [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.
- [FSH⁺06] Rob Fergus, Barun Singh, Aaron Hertzmann, Sam T Roweis, and William T Freeman. Removing camera shake from a single photograph. In *ACM transactions on graphics (TOG)*, volume 25, pages 787–794. ACM, 2006.
- [GHJY15] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.
- [GMWZ17] Donald Goldfarb, Cun Mu, John Wright, and Chaoxu Zhou. Using negative curvature in solving nonlinear programs. *Computational Optimization and Applications*, 68(3):479–502, 2017.
- [Gol80] Donald Goldfarb. Curvilinear path steplength algorithms for minimization which use directions of negative curvature. *Mathematical programming*, 18(1):31–40, 1980.
- [HHSS09] Stefan Harmeling, Michael Hirsch, Suvrit Sra, and Bernhard Scholkopf. Online blind deconvolution for astronomical imaging. In *2009 IEEE International Conference on Computational Photography (ICCP 2009)*, pages 1–7. IEEE, 2009.
- [JSE⁺98] Richard Johnson, Philip Schniter, Thomas J Endres, James D Behm, Donald R Brown, and Raúl A Casas. Blind equalization using the constant modulus criterion: A review. *Proceedings of the IEEE*, 86(10):1927–1950, 1998.
- [JSK08] Neel Joshi, Richard Szeliski, and David J Kriegman. Psf estimation using sharp edge prediction. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- [KH96] Deepa Kundur and Dimitrios Hatzinakos. Blind image deconvolution. *IEEE signal processing magazine*, 13(3):43–64, 1996.
- [KK17] Michael Kech and Felix Krahmer. Optimal injectivity conditions for bilinear inverse problems with applications to identifiability of deconvolution problems. *SIAM Journal on Applied Algebra and Geometry*, 1(1):20–37, 2017.
- [KT98] Kjetil F Kaaresen and Tofinn Taxt. Multichannel blind deconvolution of seismic signals. *Geophysics*, 63(6):2093–2107, 1998.
- [KTF11] Dilip Krishnan, Terence Tay, and Rob Fergus. Blind deconvolution using a normalized sparsity measure. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 233–240. IEEE, 2011.
- [Lan92] Richard G Lane. Blind deconvolution of speckle images. *JOSA A*, 9(9):1508–1514, 1992.
- [LB87] RG Lane and RHT Bates. Automatic multidimensional deconvolution. *JOSA A*, 4(1):180–188, 1987.
- [LB18] Yanjun Li and Yoram Bresler. Global geometry of multichannel sparse blind deconvolution on the sphere. *arXiv preprint arXiv:1404.4104*, 2018.
- [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):R53–R78, 1998.
- [LFDF07] Anat Levin, Rob Fergus, Fredo Durand, and William T Freeman. Deconvolution using natural image priors. *Massachusetts Institute of Technology, Computer Science and Artificial Intelligence Laboratory*, 3, 2007.

- [LLB16] Yanjun Li, Kiryung Lee, and Yoram Bresler. Identifiability in blind deconvolution with subspace or sparsity constraints. *IEEE Transactions on Information Theory*, 62(7):4266–4275, 2016.
- [LLB17] Yanjun Li, Kiryung Lee, and Yoram Bresler. Identifiability and stability in blind deconvolution under minimal assumptions. *IEEE Transaction of Information Theory*, 2017.
- [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, 2011.
- [MC99] Joanne Markham and José-Angel Conchello. Parametric blind deconvolution: a robust method for the simultaneous estimation of image and blur. *JOSA A*, 16(10):2377–2391, 1999.
- [MK88] Masato Miyoshi and Yutaka Kaneda. Inverse filtering of room acoustics. *IEEE Transactions on acoustics, speech, and signal processing*, 36(2):145–152, 1988.
- [NG10] Patrick A Naylor and Nikolay D Gaubitch. *Speech dereverberation*. Springer Science & Business Media, 2010.
- [OPT00] Michael R Osborne, Brett Presnell, and Berwin A Turlach. A new approach to variable selection in least squares problems. *IMA journal of numerical analysis*, 20(3):389–403, 2000.
- [PF14] Daniele Perrone and Paolo Favaro. Total variation blind deconvolution: The devil is in the details. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2909–2916, 2014.
- [PSG⁺16] Eftychios A Pnevmatikakis, Daniel Soudry, Yuanjun Gao, Timothy A Machado, Josh Merel, David Pfau, Thomas Reardon, Yu Mu, Clay Lacefield, Weijian Yang, et al. Simultaneous denoising, deconvolution, and demixing of calcium imaging data. *Neuron*, 89(2):285–299, 2016.
- [RV⁺13] Mark Rudelson, Roman Vershynin, et al. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18, 2013.
- [Sah07] Swapan K Saha. *Diffraction-limited imaging with large and moderate telescopes*. World Scientific, 2007.
- [Sat75] Yoichi Sato. A method of self-recovering equalization for multilevel amplitude-modulation systems. *IEEE Transactions on communications*, 23(6):679–682, 1975.
- [SCI75] Thomas G Stockham, Thomas M Cannon, and Robert B Ingebreetsen. Blind deconvolution through digital signal processing. *Proceedings of the IEEE*, 63(4):678–692, 1975.
- [SGG⁺09] Gleb Shtengel, James A Galbraith, Catherine G Galbraith, Jennifer Lippincott-Schwartz, Jennifer M Gillette, Suliana Manley, Rachid Sougrat, Clare M Waterman, Pakorn Kanchanawong, Michael W Davidson, et al. Interferometric fluorescent super-resolution microscopy resolves 3d cellular ultrastructure. *Proceedings of the National Academy of Sciences*, 106(9):3125–3130, 2009.
- [SJA08] Qi Shan, Jiaya Jia, and Aseem Agarwala. High-quality motion deblurring from a single image. In *Acm transactions on graphics (tog)*, volume 27, page 73. ACM, 2008.
- [SQW17] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere ii: Recovery by riemannian trust-region method. *IEEE Transactions on Information Theory*, 63(2):885–914, 2017.
- [SW90] Ofir Shalvi and Ehud Weinstein. New criteria for blind deconvolution of nonminimum phase systems (channels). *IEEE Transactions on information theory*, 36(2):312–321, 1990.
- [Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [WC16] Liming Wang and Yuejie Chi. Blind deconvolution from multiple sparse inputs. *IEEE Signal Processing Letters*, 23(10):1384–1388, 2016.
- [WJPH17] Philipp Walk, Peter Jung, Götz E Pfander, and Babak Hassibi. Blind deconvolution with additional autocorrelations via convex programs. *arXiv preprint arXiv:1701.04890*, 2017.
- [WY13] Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality constraints. *Mathematical Programming*, 142(1-2):397–434, 2013.
- [WZ14] David Wipf and Haichao Zhang. Revisiting bayesian blind deconvolution. *The Journal of Machine Learning Research*, 15(1):3595–3634, 2014.

- [XJ10] Li Xu and Jiaya Jia. Two-phase kernel estimation for robust motion deblurring. In *European conference on computer vision*, pages 157–170. Springer, 2010.
- [YK96] Yu-Li You and Mostafa Kaveh. Anisotropic blind image restoration. In *Image Processing, 1996. Proceedings., International Conference on*, volume 2, pages 461–464. IEEE, 1996.
- [YWHM10] Jianchao Yang, John Wright, Thomas S Huang, and Yi Ma. Image super-resolution via sparse representation. *IEEE transactions on image processing*, 19(11):2861–2873, 2010.
- [ZKW18] Yuqian Zhang, Han-Wen Kuo, and John Wright. Structured local optima in sparse blind deconvolution. *arXiv preprint arXiv:1806.00338*, 2018.
- [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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4894–4902, 2017.

A Basic bounds for Bernoulli-Gaussian vectors

In this section, we prove several lemmas pertaining to the sparse random vector $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$.

Lemma A.1 (Support of \mathbf{x}_0). *Let $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ and $I_0 = \text{supp}(\mathbf{x}_0) \subseteq [n]$. Suppose $n > 10\theta^{-1}$, then for any $\varepsilon \in (0, \frac{1}{10})$, with probability at least $1 - \varepsilon$ we have*

$$||I_0| - n\theta| \leq 2\sqrt{n\theta} \log \varepsilon^{-1}. \quad (\text{A.1})$$

And suppose $n \geq C\theta^{-2} \log p$ and θ , then with probability at least $1 - 2/n$, we have

$$\forall t \in [2p] \setminus \{0\}, \quad \frac{1}{2}n\theta^2 \leq |I_0 \cap (I_0 + t)| \leq 2n\theta^2 \quad (\text{A.2})$$

where C is a numerical constant.

Proof. Let $\mathbf{x}_0 = \boldsymbol{\omega} \cdot \mathbf{g} \sim_{\text{i.i.d.}} \text{BG}(\theta)$, notice that the support of the Bernoulli-Gaussian vector \mathbf{x}_0 is almost surely equal to the support of the Bernoulli vector $\boldsymbol{\omega}$. Applying Bernstein inequality [Lemma J.4](#) with $(\sigma^2, R) = (1, 1)$, then if $n\theta > 10$ we have

$$\mathbb{P} \left[\left| \sum_{k \in [n]} \omega_k - n\theta \right| > 2\sqrt{n\theta} \log \varepsilon^{-1} \right] \leq 2 \exp \left(\frac{-4n\theta \log^2 \varepsilon^{-1}}{2n\theta + 4\sqrt{n\theta} \log \varepsilon^{-1}} \right) \leq \varepsilon.$$

For [\(A.2\)](#), let $J_t := I_0 \cap (I_0 + t)$. The cardinality of J_t is an inner product between shifts of $\boldsymbol{\omega}$:

$$|J_t| = \sum_{k \in [n]} \omega_k \omega_{k-t}, \quad (\text{A.3})$$

and define two subset $J_{t1} \uplus J_{t2} = J_t$, as follows:

$$\begin{cases} J_{t1} = J_t \cap \mathcal{K}_1, & \mathcal{K}_1 := [n] \cap \{0, \dots, t-1, 2t, \dots, 3t-1, \dots\} \\ J_{t2} = J_t \cap \mathcal{K}_2, & \mathcal{K}_2 := [n] \cap \{t, \dots, 2t-1, 3t, \dots, 4t-1, \dots\} \end{cases}. \quad (\text{A.4})$$

Here, the size of sets $\mathcal{K}_1, \mathcal{K}_2$ has two-side bounds $0.4n \leq (n - 2p)/2 \leq |\mathcal{K}_2| \leq |\mathcal{K}_1| \leq (n + 2p)/2 \leq 0.6n$, thus the size of sets J_{t1}, J_{t2} can be derived using Bernstein inequality [Lemma J.4](#) with $n > C\theta^{-2} \log p$ as

$$\begin{aligned} \mathbb{P} \left[\max_{t \in [2p] \setminus \{0\}} |J_{t1}| \geq n\theta^2 \right] &= \mathbb{P} \left[\max_{t \in [2p] \setminus \{0\}} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k-t} \geq n\theta^2 \right] \leq 2p \cdot \mathbb{P} \left[\sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \geq n\theta^2 \right] \\ &\leq 2p \cdot \mathbb{P} \left[\sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} - \mathbb{E} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \geq n\theta^2 - 0.6n\theta^2 \right] \\ &\leq 4p \cdot \exp \left(\frac{-(0.4n\theta^2)^2}{2 \cdot 0.6n\theta^2 + 2 \cdot 0.4n\theta^2} \right) = \exp(\log(4p) - 0.08n\theta^2) \leq 1/n, \end{aligned} \quad (\text{A.5})$$

where the last two inequalities hold with $C > 10^5$. The lower bound can also derived as follows

$$\begin{aligned} \mathbb{P} \left[\min_{t \in [2p] \setminus \{0\}} |J_{t1}| \leq n\theta^2/4 \right] &= \mathbb{P} \left[\min_{t \in [2p] \setminus \{0\}} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k-t} \leq n\theta^2/4 \right] \leq 2p \cdot \mathbb{P} \left[\sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \leq n\theta^2/4 \right] \\ &\leq 2p \cdot \mathbb{P} \left[\sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} - \mathbb{E} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \leq n\theta^2/4 - 0.4n\theta^2 \right] \\ &\leq 4p \cdot \exp \left(\frac{-(0.15n\theta^2)^2}{2 \cdot 0.6n\theta^2 + 2 \cdot 0.15n\theta^2} \right) = \exp(\log(4p) - 0.0015n\theta^2) \leq 1/n. \end{aligned} \quad (\text{A.6})$$

The bound for $|J_2|$ can derived similarly to [\(A.5\)](#)-[\(A.6\)](#). ■

Lemma A.2 (Norms of \mathbf{x}_0). *Let $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$. If $n \geq 10\theta^{-1}$, then for any $\varepsilon \in (0, \frac{1}{10})$, with probability at least $1 - \varepsilon$,*

$$\left| \|\mathbf{x}_0\|_1 - \sqrt{2/\pi} n\theta \right| \leq 2\sqrt{n\theta} \log \varepsilon^{-1}, \quad \left| \|\mathbf{x}_0\|_2^2 - n\theta \right| \leq 3\sqrt{n\theta} \log \varepsilon^{-1} \quad (\text{A.7})$$

Proof. To bound $\|\mathbf{x}_0\|_1$, using Bernstein inequality with $(\sigma^2, R) = (\theta, 1)$ and with $n\theta \geq 10$ we have

$$\mathbb{P} \left[\left| \|\mathbf{x}_0\|_1 - \sqrt{\frac{2}{\pi}} n\theta \right| \geq 2\sqrt{n\theta} \log \varepsilon^{-1} \right] \leq 2 \exp \left(\frac{-4n\theta \log^2 \varepsilon^{-1}}{2n\theta + 4\sqrt{n\theta} \log \varepsilon^{-1}} \right) \leq \varepsilon$$

Similarly for $\|\mathbf{x}_0\|_2^2$, from Gaussian moments [Lemma J.2](#), we know the 2-norm $\sum_{i \in [n]} \mathbb{E} |x_{0i}|^4 = 3n\theta$ and q -norm $\sum_{i \in [n]} \mathbb{E} |x_{0i}|^{2q} \leq (n\theta)(2q-1)!! \leq \frac{1}{2}(3n\theta)2^{q-2}q!$ for $q \geq 3$. Let $(\sigma^2, R) = (3\theta, 2)$ in Bernstein inequality form [Lemma J.4](#), $n\theta \geq 10$ we have

$$\mathbb{P} \left[\left| \|\mathbf{x}_0\|_2^2 - n\theta \right| \geq 3\sqrt{n\theta} \log \varepsilon^{-1} \right] \leq 2 \exp \left(\frac{-9n\theta \log^2 \varepsilon^{-1}}{2(3n\theta) + 12\sqrt{n\theta} \log \varepsilon^{-1}} \right) \leq \varepsilon,$$

completing the proof. \blacksquare

Lemma A.3 (Norms of \mathbf{x}_0 subvectors). *Let $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$ and $n > 10$, then with probability at least $1 - 3/n$, we have*

$$\max_{\substack{U=[2p]+j \\ j \in [n]}} \|\mathbf{P}_U \mathbf{x}_0\|_2^2 \leq 2p\theta + 6 \left(\sqrt{p\theta} + \log n \right) \quad (\text{A.8})$$

and if \mathbf{a}_0 is μ -shift coherent and there exists a constance c_μ such that both $\theta^2 p < c_\mu$ and $\mu p^2 \theta < c_\mu$, then

$$\max_{\substack{U=[p]+j \\ j \in [n]}} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 \leq p\theta + \log n. \quad (\text{A.9})$$

Proof. Use Bernstein inequality with $(\sigma^2, R) = (3\theta, 2)$ and $t = \max \{ \sqrt{p\theta}, \log n \}$, with union bound we obtain:

$$\begin{aligned} \mathbb{P} \left[\max_{\substack{U=[2p]+j \\ j \in [n]}} \|\mathbf{P}_U \mathbf{x}_0\|_2^2 \geq 2p\theta + 6 \left(\sqrt{p\theta} + \log n \right) \right] &\leq 2n \exp \left(-\frac{36 \left(\sqrt{p\theta} + \log n \right)^2}{6p\theta + 12 \left(\sqrt{p\theta} + \log n \right)} \right) \\ &\leq 2 \exp \left(\log n - \frac{36t^2}{6t^2 + 12t} \right) \leq \frac{2}{n}. \end{aligned} \quad (\text{A.10})$$

For the second inequality, first we know calculate the expectation

$$\begin{aligned} \mathbb{E} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 &= \mathbb{E} [\mathbf{x}_0^* \mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0] \\ &= \theta \cdot \text{tr} (\mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0}) \|\mathbf{a}_0\|_2^2 + \theta \cdot \sum_{i=1}^{p-1} \|\boldsymbol{\iota}^* s_i[\mathbf{a}_0]\|_2^2 \\ &= p\theta. \end{aligned} \quad (\text{A.11})$$

Then apply Henson Wright inequality [Lemma J.6](#) with $\|\mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0}\|_F^2 = \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}\|_F^2 \leq p(1 + \mu p)$ and also $\|\mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0}\|_2 = \|\mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}\|_2^2 = 1 + \mu p$, we can derive

$$\mathbb{P} \left[\max_{\substack{U=[p]+j \\ j \in [n]}} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 \geq p\theta + \log n \right] \leq n \exp \left(-\min \left\{ \frac{\log^2 n}{64\theta^2 p(1 + \mu p)}, \frac{\log n}{8\sqrt{2}\theta(1 + \mu p)} \right\} \right)$$

$$\leq \exp \left(\log n - \min \left\{ \frac{\log^2 n}{128c_\mu}, \frac{\log n}{32c_\mu} \right\} \right) \leq \frac{1}{n} \quad (\text{A.12})$$

when $c_\mu < \frac{1}{300}$. ■

Lemma A.4 (Inner product between shifted \mathbf{x}_0). *Let $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$. There exists a numerical constant C such that if $n > C\theta^{-2} \log p$ and $p\theta \log^2 \theta^{-1} > 1$, with probability at least $1 - 4/n$, the following two statements hold simultaneously:*

$$\max_{i \neq j \in [2p]} \langle s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle \leq 6\sqrt{n\theta^2 \log n}; \quad (\text{A.13})$$

and for $\mathbf{x}_i = |\mathbf{x}_{0,i}| \in \mathbb{R}_+^n$ the vector of magnitudes of \mathbf{x}_0 ,

$$\max_{i \neq j \in [2p]} \langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle \leq 4n\theta^2. \quad (\text{A.14})$$

Proof. We will start from proving (A.14). Write $\mathbf{x} = |\mathbf{g}| \circ \boldsymbol{\omega}$ where $\mathbf{g} / \boldsymbol{\omega}$ are Gaussian/Bernoulli random vectors respectively. Let I_0 denote the support of $\boldsymbol{\omega}$ and $t = |j - i|$ with $0 < t < p$. Then (A.14) can be written as summation of Gaussian r.v.s. on intersection of support set between shifts:

$$\langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle = \sum_{k \in I_0 \cap (I_0 + t)} |\mathbf{g}_k| |\mathbf{g}_{k-t}| \quad (\text{A.15})$$

Define $J_t := I_0 \cap (I_0 + t) = J_{t1} \uplus J_{t2}$ same as (A.4). Notice that both $\sum_{k \in J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k-t}|$ and $\sum_{k \in J_{t2}} |\mathbf{g}_k| |\mathbf{g}_{k-t}|$ are sum of independent r.v.s.. We are left to consider the upper bound of $\sum_{j \in J_{t1}} |\mathbf{g}_j| |\mathbf{g}'_j|$ where \mathbf{g}, \mathbf{g}' are independent Gaussian vectors. We condition on the following event

$$\mathcal{E}_J := \{ \forall t \in [2p] \setminus \{0\}, n\theta^2/4 \leq |J_{t1}|, |J_{t2}| \leq n\theta^2 \}, \quad (\text{A.16})$$

which holds w.p. at least $1 - 2/n$ from Lemma A.1. Since $\sum_{j \in J_{t1}} |\mathbf{g}_j| |\mathbf{g}'_j| \leq \|\mathbf{g}_{J_{t1}}\|_2 \|\mathbf{g}'_{J_{t1}}\|_2$, we use Gaussian concentration Lemma J.3 and union bound to obtain

$$\begin{aligned} \mathbb{P} \left[\max_{t \in [2p] \setminus \{0\}} \sum_{j \in J_{t1}} |\mathbf{g}_j \mathbf{g}'_j| > 2|J_{t1}| \right] &\leq 2p \cdot \mathbb{P} [\|\mathbf{g}_{J_{t1}}\|_2 \|\mathbf{g}'_{J_{t1}}\|_2 - \mathbb{E} \|\mathbf{g}_{J_{t1}}\|_2 \|\mathbf{g}'_{J_{t1}}\|_2 > |J_{t1}|] \\ &\leq 4p \cdot \mathbb{P} [\|\mathbf{g}_{J_{t1}}\|_2 - \mathbb{E} \|\mathbf{g}_{J_{t1}}\|_2 > \sqrt{|J_{t1}|}/3] \\ &\leq 4p \exp(-(|J_{t1}|/9)/2) \leq 4p \exp(-n\theta^2/72) \leq 1/n \end{aligned} \quad (\text{A.17})$$

where the last inequality is derived simply via assuming $n = C\theta^{-2} \log p$ for some $C > 10^4$, such that

$$\begin{aligned} C > 400 * (4C)^{1/5} &\implies C \log p > 400 \log((4C)^{1/5} p) \implies C \log p > 72 \log(4Cp^5) > 72 \log(4Cp^2 \log^3 p) \\ &\implies n\theta^2 > 72 \log(p \cdot 4C\theta^{-2} \log p) = 72 \log(4np). \end{aligned}$$

Likewise for sum on set J_{t2} , we collect all above result and conclude for every $i \neq j \in [2p]$,

$$\langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle = \sum_{k \in J_{t1}} |\mathbf{g}_k| |\mathbf{g}'_{k-t}| + \sum_{k \in J_{t2}} |\mathbf{g}_k| |\mathbf{g}'_{k-t}| \leq 2(|J_{t1}| + |J_{t2}|) \leq 4n\theta^2. \quad (\text{A.18})$$

For (A.13) similarly condition on event \mathcal{E}_J , using Bernstein inequality Lemma J.4 with $(\sigma^2, R) = (1, 1)$:

$$\mathbb{P} \left[\max_{t \in [2p] \setminus \{0\}} \left| \sum_{j \in J_{t1}} \mathbf{g}_j \mathbf{g}'_j \right| > 3\sqrt{n\theta^2 \log n} \right] \leq p \cdot \exp \left(\frac{-9n\theta^2 \log n}{2|J_{t1}| + 6\sqrt{n\theta^2 \log n}} \right) \leq p \cdot \exp \left(\frac{-9n\theta^2 \log n}{3n\theta^2} \right) \leq \frac{1}{n} \quad (\text{A.19})$$

thus for every $i \neq j \in [2p]$,

$$|\langle s_i[\mathbf{x}_0], s_j[\mathbf{s}_0] \rangle| \leq \left| \sum_{k \in J_{t1}} \mathbf{g}_k \mathbf{g}'_{k-t} \right| + \left| \sum_{k \in J_{t2}} \mathbf{g}_k \mathbf{g}'_{k-t} \right| \leq 6\sqrt{n\theta^2 \log n}. \quad (\text{A.20})$$

Finally, both (A.18),(A.20) holds simultaneously with probability at least

$$1 - 2/n - 1/n - 1/n = 1 - 4/n \quad (\text{A.21})$$

■

Lemma A.5 (Convolution of \mathbf{x}_0). *Given $\mathbf{y} = \mathbf{x}_0 * \mathbf{a}_0$ where $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$ and $\mathbf{a}_0 \in \mathbb{R}^{p_0}$ is μ -shift coherent. Suppose $n \geq C\theta^{-2} \log p$ for some numerical constant $C > 0$, with probability at least $1 - 7/n$, we have the following two statement simultaneously hold:*

$$\|\mathbf{C}_{\mathbf{y}} \boldsymbol{\iota}\|_2^2 \leq 3(1 + \mu p)n\theta \quad (\text{A.22})$$

and for all $J \subseteq [n]$,

$$\|\mathbf{P}_J \mathbf{C}_{\mathbf{y}} \boldsymbol{\iota}\|_2^2 \leq 14|J|(1 + \mu p)(p\theta + \log n) \quad (\text{A.23})$$

Proof. Given any $\mathbf{a} \in \mathbb{S}^{p-1}$, write $\boldsymbol{\beta} = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{a}$ where $|\boldsymbol{\beta}| \leq 2p$. Apply $\|\mathbf{x}_0\|_2^2 \leq 2n\theta$ from Lemma A.2 by choosing $\varepsilon = 1/n$, also $|\langle s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle| \leq 6\sqrt{n\theta^2 \log n}$ from Lemma A.4 we get:

$$\begin{aligned} \|\mathbf{C}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}\|_2^2 &= \|\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\beta}\|_2^2 \leq \|\boldsymbol{\beta}\|_2^2 \|\mathbf{x}_0\|_2^2 + \sum_{i \neq j \in [\pm p]} |\beta_i \beta_j \langle s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle| \\ &\leq \|\boldsymbol{\beta}\|_2^2 \|\mathbf{x}_0\|_2^2 + \|\boldsymbol{\beta}\|_1^2 \max_{i \neq j \in [\pm p]} |\langle s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle| \\ &\leq \|\boldsymbol{\beta}\|_2^2 \cdot 2n\theta + p \|\boldsymbol{\beta}\|_2^2 \cdot 6\sqrt{n\theta^2 \log n} \leq 3 \|\boldsymbol{\beta}\|_2^2 n\theta \end{aligned}$$

where $n = C\theta^{-2} \log p$ with $C \geq 10^4$, and the statement holds with probability at least $1 - 5/n$.

For the bound of $\|\mathbf{P}_J \mathbf{C}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}\|_2^2$. Simply apply Lemma A.3 and utilize norm bound of $\|\boldsymbol{\beta}\|_2^2$, with probability at least $1 - 2/n$ we have:

$$\|\mathbf{P}_J \mathbf{C}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}\|_2^2 = \sum_{i \in J} |\langle s_i[\mathbf{x}_0], \boldsymbol{\beta} \rangle|^2 \leq |J| \max_{\substack{U=[2p]+j \\ j \in [n]}} \|\mathbf{P}_U \mathbf{x}_0\|_2^2 \|\boldsymbol{\beta}\|_2^2 \leq |J| \cdot 14(p\theta + \log n) \cdot \|\boldsymbol{\beta}\|_2^2$$

Finally apply Lemma B.4 and Gershgorin disc theorem obtain

$$\|\boldsymbol{\beta}\|_2^2 = \|\mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{a}\|_2^2 \leq \|\mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}\|_2^2 = \sigma_{\max}(\mathbf{M}) \leq 1 + \mu p. \quad (\text{A.24})$$

Remark A.6. When \mathbf{a}_0 is a basis vector \mathbf{e}_0 , the result of Lemma A.5 gives upper bound of $\|\mathbf{C}_{\mathbf{x}_0}\|_2 < 3n\theta$, whose lower bound can be derived similarly with $\|\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}\|_2 \geq \frac{2}{3}n\theta$

■

B Vectors in shift space

In this section, we will establish a number of properties of the coefficient vectors α and correlation vector β . Generally speaking, when \mathbf{a} is close to the subspace \mathcal{S}_τ , then both vectors α, β have most of their energy concentrated on the entries τ . In this section, we derive upper bounds on α_{τ^c} and β_{τ^c} under various assumptions.

In particular, we will introduce a relationship between the sparsity rate θ , coherence μ and size $|\tau|$, which we term the sparsity-coherence condition. In Lemma B.2 we prove that measuring the distance from \mathbf{a} to subspace \mathcal{S}_τ in terms of $\|\alpha_{\tau^c}\|_2$ gives a seminorm. We then use this distance to characterize a region $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ around the subspace \mathcal{S}_τ . Later, in Lemma B.4 we illustrate the relationship between α and β , where $\beta = C_{a_0}^* \iota^* C_{a_0} \alpha$. Finally in Lemma B.5 and Corollary B.6, controls the magnitude of α_{τ^c} and β_{τ^c} near \mathcal{S}_τ .

Definition B.1 (Sparsity-coherence condition). *Let $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$ with shift coherence μ . We say that $(\mathbf{a}_0, \theta, |\tau|)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$ with constant c_μ , if*

$$\theta \in \left[\frac{1}{p}, \frac{c_\mu}{4 \max\{|\tau|, \sqrt{p}\}} \right] \cdot \frac{1}{\log^2 \theta^{-1}}, \quad \mu \cdot \max\{|\tau|^2, p^2 \theta^2\} \cdot \log^2 \theta^{-1} \leq \frac{c_\mu}{4}, \quad (\text{B.1})$$

where $p = 3p_0 - 2$.

Lemma B.2 (d_α is a seminorm). *For every solution subspace \mathcal{S}_τ , the function $d_\alpha(\cdot, \mathcal{S}_\tau) : \mathbb{R}^p \rightarrow \mathbb{R}_+$ defined as*

$$d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{a_0} \alpha \}. \quad (\text{B.2})$$

is a seminorm, and for all $\mathbf{a} \in \mathcal{S}_\tau$, $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = 0$.

Proof. It is immediate from definition that $d_\alpha(\cdot, \mathcal{S}_\tau)$ is nonnegative and $\mathcal{S}_\tau \subseteq \{\mathbf{a} : d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = 0\}$. Subadditivity can be shown from simple norm inequalities and our definition of d_α , for all $\mathbf{a}_1, \mathbf{a}_2$ we have

$$\begin{aligned} d_\alpha(\mathbf{a}_1 + \mathbf{a}_2, \mathcal{S}_\tau) &= \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a}_1 + \mathbf{a}_2 = \iota^* C_{a_0} \alpha \} \\ &= \inf \{ \|\alpha_{1\tau^c} + \alpha_{2\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{a_0} \alpha_1, \quad \mathbf{a}_2 = \iota^* C_{a_0} \alpha_2 \} \\ &\leq \inf \{ \|\alpha_{1\tau^c}\|_2 + \|\alpha_{2\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{a_0} \alpha_1, \quad \mathbf{a}_2 = \iota^* C_{a_0} \alpha_2 \} \\ &= \inf \{ \|\alpha_{1\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{a_0} \alpha_1 \} + \inf \{ \|\alpha_{2\tau^c}\|_2 \mid \mathbf{a}_2 = \iota^* C_{a_0} \alpha_2 \} \\ &= d_\alpha(\mathbf{a}_1, \mathcal{S}_\tau) + d_\alpha(\mathbf{a}_2, \mathcal{S}_\tau). \end{aligned}$$

Similarly the absolute homogeneity, for any $c \in \mathbb{R}$:

$$\begin{aligned} d_\alpha(c \cdot \mathbf{a}, \mathcal{S}_\tau) &= \inf \{ \|\alpha'_{\tau^c}\|_2 \mid c \cdot \mathbf{a} = \iota^* C_{a_0} \alpha' \} = \inf \{ \|c \cdot \alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{a_0} \alpha \} \\ &= |c| \cdot \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{a_0} \alpha \} = |c| \cdot d_\alpha(\mathbf{a}, \mathcal{S}_\tau), \end{aligned}$$

which completes the proof that d_α is a seminorm. ■

Definition B.3 (Widened subspace). *For subspace \mathcal{S}_τ let*

$$\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu)) := \{ \mathbf{a} \in \mathbb{S}^{p-1} \mid d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma \} \quad (\text{B.3})$$

denote its widening by γ , in the seminorm d_α .

Our analysis works with a specific choice of width $\gamma(c_\mu)$, which depends on the problem parameters $\mathbf{a}_0, \theta, |\tau|$ and a constant c_μ , via

$$\gamma(c_\mu) = \frac{c_\mu}{4 \log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\tau|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\tau|} \right\}. \quad (\text{B.4})$$

Lemma B.4 (Properties of $C_{a_0}^* \iota^* C_{a_0}$). Let $M = C_{a_0}^* \iota^* C_{a_0}$, with $a_0 \in \mathbb{S}^{p_0-1}$ μ -shift coherent. The diagonal entries of M satisfy

$$\begin{cases} M_{ii} = 1 & i \in [-p_0 + 1, p_0 - 1] = [\pm p_0], \\ 0 \leq M_{ii} \leq 1 & i \in [-2p_0 + 2, -p_0] \cup [p_0, 2p_0 - 2], \\ M_{ii} = 0 & \text{otherwise,} \end{cases} \quad (\text{B.5})$$

and the off-diagonal entries satisfy

$$\begin{cases} |M_{ij}| \leq \mu & 0 < |i - j| < p_0, \{i \in [-p_0 + 1, p_0 - 1]\} \cup \{j \in [-p_0 + 1, p_0 - 1]\}, \\ |M_{ij}| < 1 & \{i, j \in [-2p_0 + 2, -p_0]\} \cup \{i, j \in [p_0, 2p_0 - 2]\}, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{B.6})$$

Furthermore, let $\tau \subset [\pm p_0]$, and $\tau^c = [\pm 2p_0 - 1] \setminus \tau$. The singular values of submatrix $\iota_\tau^* M \iota_\tau$ can be bounded as:

$$\begin{cases} 1 - \mu |\tau| \leq \sigma_{\min}(\iota_\tau^* M \iota_\tau) \leq \sigma_{\max}(\iota_\tau^* M \iota_\tau) \leq 1 + \mu |\tau|, \\ \sigma_{\max}(\iota_{\tau^c}^* M \iota_{\tau^c}) \leq \mu \sqrt{p |\tau|}, \\ \sigma_{\max}(\iota_{\tau^c}^* M \iota_{\tau^c}) \leq 1 + \mu p. \end{cases} \quad (\text{B.7})$$

Proof. Recall the definition of ι , which selects the entries $\{-p_0 + 1, \dots, 2p_0 - 2\}$. The entrywise properties of M can be derived by carefully counting the entries of the shifted support. The submatrix M on support $\{-2p_0 + 2, \dots, 2p_0 - 2\}$ has an upper bound to be characterized as follows:

$$\left| \iota_{[\pm 2p_0 - 1]}^* M \iota_{[\pm 2p_0 - 1]} \right| \leq \begin{bmatrix} \mathbf{J} & \mu \cdot \mathbf{1} & \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} & \mathbf{0} & \mathbf{0} \\ \mu \cdot \mathbf{1} & \mathbf{I} + \mu \cdot \mathbf{1}_o & \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix} & \mu \cdot \mathbf{1} & \mathbf{0} \\ [0 \cdots 0] & [\mu \cdots \mu] & 1 & [\mu \cdots \mu] & [0 \cdots 0] \\ \mathbf{0} & \mu \cdot \mathbf{1} & \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix} & \mathbf{I} + \mu \cdot \mathbf{1}_o & \mu \cdot \mathbf{1} \\ \mathbf{0} & \mathbf{0} & \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} & \mu \cdot \mathbf{1} & \mathbf{J} \end{bmatrix}. \quad (\text{B.8})$$

Here, the center row/column vector is indexed at 0, the matrices \mathbf{J} , \mathbf{I} , $\mathbf{1}$ and $\mathbf{1}_o$ are square and of size $(p_0 - 1)^2$. Among which, \mathbf{I} is the identity matrix, $\mathbf{1}$ is the ones matrix whereas $\mathbf{1}_o$ has all off diagonal entries equal 1. Also $|\mathbf{J}|$ has property $|J_{ij}| < 1$ for all i, j .

As for the singular values, notice that the first and second inequalities consider submatrix not containing \mathbf{J} since $\tau \subseteq [\pm p_0]$; thus the first inequality can be derived with Gershgorin disc theorem directly, and the second inequality with the upper bound with its Frobenius norm:

$$\sigma_{\max}(\iota_{\tau^c}^* M \iota_{\tau^c}) \leq \mu \sqrt{(2p_0 - 1) |\tau|} < \mu \sqrt{p |\tau|}. \quad (\text{B.9})$$

Finally by recalling $p = 3p_0 - 2 > 2p_0 - 1$. The last inequality is direct from bound of $\iota^* C_{a_0}$:

$$\sigma_{\max}(\iota_{\tau^c}^* M \iota_{\tau^c}) \leq \|C_{a_0}^* \iota^* C_{a_0}\|_2 = \|\iota^* C_{a_0} C_{a_0}^* \iota\|_2 = \|\iota^* C_{a_0}^* C_{a_0} \iota\|_2 \leq 1 + \mu p \quad (\text{B.10})$$

where the third equality is derived via commutativity of convolution. ■

Lemma B.5 (Shift space vectors in widened subspace). *Let $(\mathbf{a}_0, \theta, |\boldsymbol{\tau}|)$ satisfy the sparsity-coherence condition $\text{SCC}(c_\mu)$. Then for every $\mathbf{a} \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, every $\boldsymbol{\alpha}$ satisfying $\mathbf{a} = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}$ and $\|\boldsymbol{\alpha}_{\tau^c}\|_2 \leq \gamma(c_\mu)$ has*

$$|\|\boldsymbol{\alpha}_\tau\|_2 - 1| \leq c_\mu; \quad (\text{B.11})$$

moreover, $\boldsymbol{\beta} = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{a}$ satisfies

$$1 - 3c_\mu \leq \|\boldsymbol{\beta}_\tau\|_2^2 \leq 1 + \frac{c_\mu}{|\boldsymbol{\tau}| \log^2 \theta^{-1}}, \quad \|\boldsymbol{\beta}_{\tau^c}\|_\infty \leq \frac{c_\mu}{\sqrt{|\boldsymbol{\tau}|} \log^2 \theta^{-1}}, \quad \|\boldsymbol{\beta}_{\tau^c}\|_2 \leq \frac{c_\mu}{|\boldsymbol{\tau}| \theta \log \theta^{-1}} \min \left\{ \sqrt{\theta}, \gamma(c_\mu) \right\}. \quad (\text{B.12})$$

Proof. Write $-1/\log \theta = \theta_{\log}$ and $\gamma = \gamma(c_\mu)$ for convenience. First, by using bounds on γ in (B.4) and $\mu |\boldsymbol{\tau}| < 1$ we obtain:

$$\begin{cases} \gamma \cdot \sqrt{1 + \mu p} \leq \gamma(1 + \sqrt{\mu p}) \leq c_\mu \theta_{\log}^2 / 2 \\ \gamma \cdot \sqrt{1 + \mu^2 p} \leq \gamma(1 + \sqrt{\mu^2 p}) \leq \frac{c_\mu \theta_{\log}^2}{4} \left(\frac{1}{\sqrt{|\boldsymbol{\tau}|}} + \sqrt{\mu} \right) \leq \frac{c_\mu \theta_{\log}^2}{2\sqrt{|\boldsymbol{\tau}|}} \\ \gamma \cdot \mu \sqrt{p |\boldsymbol{\tau}|} \leq \gamma \cdot \sqrt{\mu p} \cdot \sqrt{\mu |\boldsymbol{\tau}|} \leq c_\mu \theta_{\log}^2 / 4 \end{cases} \quad (\text{B.13})$$

Let $\mathbf{a} = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}$ with $\|\boldsymbol{\alpha}_{\tau^c}\|_2 < \gamma$. Utilize properties of $\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}$ from Lemma B.4 and $\mu |\boldsymbol{\tau}| < c_\mu/4$ and (B.13), we have:

$$\begin{aligned} \|\boldsymbol{\alpha}_\tau\|_2 &\geq \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_\tau\|_2^{-1} (\|\mathbf{a}\|_2 - \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}_{\tau^c}\|_2) \geq \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_\tau\|_2^{-1} (1 - \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}\|_2 \|\boldsymbol{\alpha}_{\tau^c}\|_2) \\ &\geq \frac{1}{\sqrt{1 + \mu |\boldsymbol{\tau}|}} (1 - \gamma \cdot \sqrt{1 + \mu p}) \geq \frac{1 - c_\mu/2}{\sqrt{1 + c_\mu/4}} \geq 1 - c_\mu, \end{aligned} \quad (\text{B.14})$$

and similarly, the upper bound can be derived as:

$$\begin{aligned} \|\boldsymbol{\alpha}_\tau\|_2 &\leq \sigma_{\min}^{-1}(\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_\tau) (\|\mathbf{a}\|_2 + \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}_{\tau^c}\|_2) \leq \sigma_{\min}^{-1}(\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_\tau) (1 + \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}\|_2 \|\boldsymbol{\alpha}_{\tau^c}\|_2) \\ &\leq \frac{1}{\sqrt{1 - \mu |\boldsymbol{\tau}|}} (1 + \gamma \cdot \sqrt{1 + \mu p}) \leq \frac{1 + c_\mu/2}{\sqrt{1 - c_\mu/4}} \leq 1 + c_\mu. \end{aligned} \quad (\text{B.15})$$

The bound of $\|\boldsymbol{\beta}_\tau\|_2^2$ can be simply obtained using $\mu |\boldsymbol{\tau}| < c_\mu/4$ and γ bound from (B.13) as:

$$\|\boldsymbol{\beta}_\tau\|_2^2 \leq \sigma_{\max}^2(\boldsymbol{\iota}_\tau^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_\tau) \leq 1 + \mu |\boldsymbol{\tau}| \leq 1 + \frac{c_\mu \theta_{\log}^2}{|\boldsymbol{\tau}|} \quad (\text{B.16})$$

$$\begin{aligned} \|\boldsymbol{\beta}_\tau\|_2^2 &\geq (\sigma_{\min}(\boldsymbol{\iota}_\tau^* \mathbf{M} \boldsymbol{\iota}_\tau) \|\boldsymbol{\alpha}_\tau\|_2 - \sigma_{\max}(\boldsymbol{\iota}_\tau^* \mathbf{M} \boldsymbol{\iota}_{\tau^c}) \|\boldsymbol{\alpha}_{\tau^c}\|_2)^2 \\ &\geq \left((1 - \mu |\boldsymbol{\tau}|) (1 - c_\mu) - \mu \sqrt{p |\boldsymbol{\tau}|} \cdot \gamma \right)^2 \geq 1 - 3c_\mu. \end{aligned} \quad (\text{B.17})$$

As for the upper bound of $\|\boldsymbol{\beta}_{\tau^c}\|_\infty$, follow from (B.13), we have:

$$\begin{aligned} \|\boldsymbol{\beta}_{\tau^c}\|_\infty &\leq \|\boldsymbol{\iota}_{\tau^c}^* \mathbf{M} \boldsymbol{\alpha}_\tau\|_\infty + \|\boldsymbol{\iota}_{\tau^c}^* \mathbf{M} \boldsymbol{\alpha}_{\tau^c}\|_\infty \leq \mu \sqrt{|\boldsymbol{\tau}|} \|\boldsymbol{\alpha}_\tau\|_2 + \sqrt{1 + \mu^2 p} \|\boldsymbol{\alpha}_{\tau^c}\|_2 \\ &\leq \frac{c_\mu \theta_{\log}^2 (1 + c_\mu)}{4 |\boldsymbol{\tau}|} + \gamma \cdot \sqrt{1 + \mu^2 p} \leq \frac{c_\mu \theta_{\log}^2}{\sqrt{|\boldsymbol{\tau}|}}; \end{aligned} \quad (\text{B.18})$$

the bound for $\|\boldsymbol{\beta}_{\tau^c}\|_2$ requires two inequalities, we know

$$\|\boldsymbol{\beta}_{\tau^c}\|_2 \leq \|\boldsymbol{\iota}_{\tau^c}^* \mathbf{M} \boldsymbol{\alpha}_\tau\|_2 + \|\boldsymbol{\iota}_{\tau^c}^* \mathbf{M} \boldsymbol{\alpha}_{\tau^c}\|_2 \leq \mu \sqrt{p |\boldsymbol{\tau}|} \|\boldsymbol{\alpha}_\tau\|_2 + (1 + \mu p) \|\boldsymbol{\alpha}_{\tau^c}\|_2, \quad (\text{B.19})$$

for the first inequality, use $(\mu|\tau|^2)^{3/4}(\mu p^2\theta^2)^{1/4} = \mu\sqrt{p\theta}|\tau|^{3/2} < c_\mu\theta_{\log}^2/4$, definition of γ and $\theta|\tau| \leq c_\mu\theta_{\log}^2/4$ we have:

$$\begin{aligned} \text{(B.19)} &\leq \frac{\mu\sqrt{p\theta}|\tau|^{3/2}}{\sqrt{\theta}|\tau|}(1+c_\mu) + \frac{\sqrt{\theta}|\tau| \cdot \sqrt{|\tau|}\gamma}{\sqrt{\theta}|\tau|} + \frac{\mu p\sqrt{\theta}|\tau|\gamma}{\sqrt{\theta}|\tau|} \\ &\leq \frac{2c_\mu\theta_{\log}^2 + c_\mu\theta_{\log}^3 + c_\mu\theta_{\log}^2}{4\sqrt{\theta}|\tau|} \leq \frac{c_\mu\theta_{\log}^2}{\sqrt{\theta}|\tau|}, \end{aligned} \quad \text{(B.20)}$$

and similarly for the second inequality, use both conditions of μ , we have:

$$\begin{aligned} \text{(B.19)} &\leq \frac{\gamma}{\theta|\tau|} \cdot \frac{\mu\sqrt{p\theta}|\tau|^{3/2}}{\gamma}(1+c_\mu) + \gamma + \mu p\gamma \\ &\leq \frac{\gamma}{\theta|\tau|} \cdot \frac{4\mu\sqrt{p\theta}|\tau|^{3/2}}{c_\mu\theta_{\log}^2} \cdot \max\left\{\sqrt{|\tau|}, \sqrt{\mu p}, \mu p\sqrt{\theta}|\tau|\right\} + \frac{\gamma}{\theta|\tau|} \cdot \theta|\tau| + \frac{\gamma}{\theta|\tau|} \cdot \mu p\theta|\tau| \\ &\leq \frac{\gamma}{\theta|\tau|} \cdot \left(\frac{4}{c_\mu\theta_{\log}^2} \cdot \max\left\{\mu|\tau|^2 \cdot \sqrt{p\theta}, \mu(p\theta)|\tau| \cdot \sqrt{\mu|\tau|}, \mu\sqrt{p\theta}|\tau|^{3/2} \cdot \mu p\theta|\tau|\right\} + \frac{c_\mu\theta_{\log}^2}{4} + \frac{c_\mu\theta_{\log}^2}{4} \right) \\ &\leq \frac{\gamma}{\theta|\tau|} \left(\frac{c_\mu\theta_{\log}}{4} + \frac{c_\mu\theta_{\log}^2}{4} + \frac{c_\mu\theta_{\log}^2}{4} \right) \leq \frac{c_\mu\theta_{\log}\gamma}{\theta|\tau|}, \end{aligned} \quad \text{(B.21)}$$

which completes the proof. \blacksquare

Corollary B.6 ($|\langle \beta_{\tau^c}, \mathbf{x}_{0,\tau^c} \rangle|$ is small). *Given $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n and $|\tau|, c_\mu$ such that $(\mathbf{a}_0, \theta, |\tau|)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Write $\lambda = c_\lambda/\sqrt{|\tau|}$ with some $c_\lambda \geq 1/5$, then if $c_\mu \leq \frac{c_\lambda}{25}$,*

$$\mathbb{P}\left[\left|\sum_{i \in \tau^c} \beta_i \mathbf{x}_{0i}\right| > \frac{\lambda}{10}\right] \leq 2\theta, \quad \mathbb{P}\left[\left|\sum_i \beta_i \mathbf{x}_{0i}\right| > \frac{\lambda}{10}\right] \leq \theta|\tau| + 2\theta. \quad \text{(B.22)}$$

Proof. We bound tail probability of the first result with Gaussian moments [Lemma J.2](#) and Bernstein inequality [Lemma J.4](#). Via Hölder's inequality, $\sum_{i \in \tau^c} \mathbb{E}(\beta_i x_i)^q = \mathbb{E}x_0^q \|\beta_{\tau^c}\|_q^q = \theta(q-1)!! \|\beta_{\tau^c}\|_2^2 \|\beta_{\tau^c}\|_\infty^{q-2}$, thus

$$\mathbb{P}\left[\left|\sum_{i \in \tau^c} \beta_i \mathbf{x}_{0i}\right| > \lambda/10\right] \leq 2 \exp\left(\frac{-(\lambda/10)^2}{2\theta \|\beta_{\tau^c}\|_2^2 + 2(\lambda/10) \|\beta_{\tau^c}\|_\infty}\right) \quad \text{(B.23)}$$

Write $\theta_{\log} = -\frac{1}{\log \theta}$, [Lemma B.5](#) implies when $c_\mu \leq \frac{c_\lambda}{25}$, we have $\theta \|\beta_{\tau^c}\|_2^2 \leq \frac{c_\mu^2 \theta_{\log}^2}{|\tau|^2} \leq \frac{\theta_{\log} \lambda^2}{625}$ and $\|\beta_{\tau^c}\|_\infty \leq \frac{c_\mu \theta_{\log}}{\sqrt{|\tau|}} \leq \frac{\theta_{\log} \lambda}{25}$, therefore,

$$\text{(B.23)} \leq 2 \exp\left(\frac{-\lambda^2/100}{2\theta_{\log} \lambda^2/625 + 2(\theta_{\log} \lambda/25) \cdot (\lambda/10)}\right) \leq 2 \exp(\log \theta) \leq 2\theta \quad \text{(B.24)}$$

The second tail bound is straight forward from the first tail bound as follows:

$$\begin{aligned} \mathbb{P}\left[\left|\sum_i \beta_i \mathbf{x}_{0i}\right| > \frac{\lambda}{10}\right] &\leq \mathbb{P}[|\beta_{\tau}^* \mathbf{x}_{\tau}| + |\beta_{\tau^c}^* \mathbf{x}_{\tau^c}| > \lambda/10] \\ &\leq \mathbb{P}[\mathbf{x}_{\tau} \neq \mathbf{0}] + \mathbb{P}[\mathbf{x}_{\tau} = \mathbf{0}] \cdot \mathbb{P}[|\beta_{\tau^c}^* \mathbf{x}_{\tau^c}| > \lambda/10] \\ &\leq \theta|\tau| + 2\theta. \end{aligned} \quad \text{(B.25)}$$

\blacksquare

Corollary B.7 ($|\langle \beta_{\tau \setminus (0)}, x_{0, \tau \setminus (0)} \rangle|$ is small near shifts). Suppose that $x_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and $|\tau|, c_\mu$ such that $(\alpha_0, \theta, |\tau|)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$, then if $c_\mu \leq \frac{1}{10}$, for any α such that $|\beta_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}$, we have

$$\mathbb{P} \left[\left| \sum_{i \in \tau \setminus (0)} \beta_i x_{0i} \right| > \frac{2\lambda}{5} \right] \leq 2\theta \quad (\text{B.26})$$

Proof. For the last tail bound, write $x = \omega \circ g$. Wlog define β_0 be the largest correlation $\beta_{(0)}$, define random variables $s' = \langle \beta_{\tau \setminus \{0\}}, x_{\tau \setminus \{0\}} \rangle$. Firstly most of the entries of x_τ would be zero since via Bernstein inequality with $\theta |\tau| < 0.1$:

$$\mathbb{P} \left[\sum_{i \in \tau} \omega_i > \log \theta^{-1} \right] \leq \mathbb{P} \left[\sum_{i \in \tau} \omega_i > \theta |\tau| + 0.9 \log \theta^{-1} \right] \leq \exp \left(\frac{-0.9^2 \log^2 \theta^{-1}}{2(\theta |\tau| + 0.9 \log \theta^{-1}/3)} \right) \leq \theta \quad (\text{B.27})$$

thus with probability at least $1 - \theta$, we can write s' as a Gaussian r.v. with variation bounded as $\mathbb{E} s'^2 \leq \mathbb{E} \left[\sum_{i=1}^{\log \theta^{-1}} \beta_i g_i \right]^2 = \log \theta^{-1} \beta_{(1)}^2$, then via Gaussian tail bound [Lemma J.1](#):

$$\mathbb{P} [|s'| > 0.4\lambda] \leq \mathbb{P} \left[|g| > \frac{0.4\lambda}{\sqrt{\log \theta^{-1}} |\beta_{(1)}|} \right] + \mathbb{P} \left[\sum_{i \in \tau} \omega_i > \log \theta^{-1} \right] \leq \frac{2}{\sqrt{2\pi}} \exp(-1.2 \log \theta^{-1}) + \theta \leq 2\theta, \quad (\text{B.28})$$

■

C Euclidean gradient as soft-thresholding in shift space

In this section, we will study the Euclidean gradient (4.6), by deriving bounds showing that the χ operator approximates a soft-thresholding function in shift space (Lemma C.2 and Corollary C.4). Furthermore, we will show the operator $\chi[\beta_i]$ is monotone in $|\beta_i|$ from Lemma C.3. A figure of visualized χ operator is shown in Figure 13.

Expectation of χ operator. To understand the χ operator, we shall first consider a simple case—when \mathbf{x}_0 is highly sparse. By definition of β from (4.3) we can see that β has a short support of size at most $2p - 1$, when \mathbf{x}_0 has support entries separated by at least $2p$, the entries of vector $\chi[\beta]_i$ become sum of independent random variables as:

$$\chi[\beta]_i = \left\langle s_{-i}[\mathbf{x}_0], \mathcal{S}_\lambda \left[\mathbf{x}_0 * \check{\beta} \right] \right\rangle \underset{\mathbf{x}_0 \text{ sep.}}{=} \left\langle s_{-i}[\mathbf{x}_0], \mathcal{S}_\lambda [\beta_i s_{-i}[\mathbf{x}_0]] \right\rangle = \sum_{j \in \text{supp}(\mathbf{x}_0)} \mathbf{g}_j \cdot \mathcal{S}_\lambda [\mathbf{g}_j \cdot \beta_i]$$

where $(\mathbf{g}_j)_{j \in [n]}$ are standard Gaussian r.v.s.

The following lemma describes the behavior of the summands in the above expression:

Lemma C.1 (Gaussian smoothed soft-thresholding). *Let $g \sim \mathcal{N}(0, 1)$. Then for every $b, s \in \mathbb{R}$ and $\lambda > 0$,*

$$\mathbb{E}_g \left[g \mathcal{S}_\lambda [b \cdot g + s] \right] = b (1 - \text{erf}_b(\lambda, s)), \quad (\text{C.1})$$

where

$$\text{erf}_b(\lambda, s) = \frac{1}{2} \text{erf} \left(\frac{\lambda + s}{\sqrt{2}|b|} \right) + \frac{1}{2} \text{erf} \left(\frac{\lambda - s}{\sqrt{2}|b|} \right). \quad (\text{C.2})$$

Furthermore, for $s = 0$, $b \in [-1, 1]$ and $\varepsilon \in (0, 1/4)$, letting $\sigma = \text{sign}(b)$ we have

$$\sigma \mathcal{S}_{\nu'_2 \lambda} [b] \leq \sigma \mathbb{E}_g \left[g \mathcal{S}_\lambda [b \cdot g] \right] \leq \sigma \mathcal{S}_{\nu'_1(\varepsilon) \lambda} [b] + \varepsilon \quad (\text{C.3})$$

where $\nu'_1(\varepsilon) = 1/(2\sqrt{-\log \varepsilon})$ and $\nu'_2 = \sqrt{2/\pi}$.

Proof. Wlog assume $b > 0$. Write f as the pdf of standard Gaussian distribution. With integral by parts:

$$\int_{-\infty}^t t' f(t') dt' = -f(t), \quad \int_{-\infty}^t t'^2 f(t') dt' = \frac{1}{2} \text{erf} \left(\frac{t}{\sqrt{2}} \right) - t f(t)$$

Integrating, we obtain

$$\mathbb{E} \left[g \mathcal{S}_\lambda [b \cdot g + s] \right] = \int_{t \geq \frac{\lambda - s}{b}} (bt^2 - (\lambda - s)t) f(t) dt + \int_{t \leq -\frac{\lambda + s}{b}} (bt^2 + (\lambda + s)t) f(t) dt,$$

by writing $L = \lambda - s$, the integral of first summand

$$\int_{t \geq \frac{L}{b}} (bt^2 - Lt) f(t) dt = b \left[\frac{1}{2} - \frac{1}{2} \text{erf} \left(\frac{L}{\sqrt{2}b} \right) + \frac{L}{b} f \left(\frac{L}{b} \right) \right] - L f \left(\frac{L}{b} \right) = \frac{b}{2} - \frac{b}{2} \text{erf} \left(\frac{L}{\sqrt{2}b} \right),$$

and similarly for the second summand, which gives

$$\mathbb{E} \left[g \mathcal{S}_\lambda [b \cdot g + s] \right] = \frac{b}{2} - \frac{b}{2} \text{erf} \left(\frac{\lambda - s}{\sqrt{2}b} \right) + \frac{b}{2} - \frac{b}{2} \text{erf} \left(\frac{\lambda + s}{\sqrt{2}b} \right) = b (1 - \text{erf}_b(\lambda, s))$$

For $b < 0$, alternatively we have

$$\mathbb{E}[gS_\lambda[-|b| \cdot g + s]] = -\mathbb{E}[gS_\lambda[|b| \cdot g - s]] = -|b|(1 - \text{erf}_b(\lambda, -s)) = b(1 - \text{erf}_b(\lambda, s)),$$

To show (C.3), via definition of error function, for $x > 0$, we know:

$$\min \left\{ 1 - \varepsilon, \frac{1 - \varepsilon}{\sqrt{\log(1/\varepsilon)}} x \right\} \leq \text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \leq \frac{2x}{\sqrt{\pi}} \quad (\text{C.4})$$

where the lower bound is derived by first knowing erf is increasing thus for all $x > \sqrt{\log(1/\varepsilon)}$,

$$\text{erf}(x) \geq 1 - e^{-x^2} \geq 1 - e^{\log \varepsilon} = 1 - \varepsilon$$

and from concavity of erf we have for $0 < x < \sqrt{\log(1/\varepsilon)} = T$,

$$\text{erf}(x) \geq \frac{\text{erf}(T) - \text{erf}(0)}{T - 0} x + \text{erf}(0) \geq \frac{1 - \varepsilon}{\sqrt{\log(1/\varepsilon)}} x.$$

Lastly plug (C.4) into (C.1) and apply condition $|b| \leq 1$ and $\varepsilon < 1/4$ we have

$$|b| - \sqrt{\frac{2}{\pi}} \lambda \leq |b| - |b| \text{erf} \left(\frac{\lambda}{\sqrt{2}|b|} \right) \leq \max \left\{ |b| \varepsilon, |b| - \frac{\lambda(1 - \varepsilon)}{\sqrt{2 \log(1/\varepsilon)}} \right\} \leq \max \left\{ \varepsilon, |b| - \frac{\lambda}{2\sqrt{\log(1/\varepsilon)}} \right\},$$

which completes the proof. \blacksquare

This lemma establishes when \mathbf{x}_0 is separated, then χ is soft thresholding operator on β with threshold about $\lambda/2$. This phenomenon extends beyond the separated case, as long as when \mathbf{x}_0 is sufficiently sparse (when Definition B.1 holds). Recall that $\chi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is defined as

$$\chi[\beta] = \widetilde{\mathcal{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda [\widetilde{\mathcal{C}}_{\mathbf{x}_0} \beta]. \quad (\text{C.5})$$

The following lemma bounds its expectation:

Lemma C.2 (Expectation of $\chi(\beta)$). *Let $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ and $\lambda > 0$, then for every $\mathbf{a} \in \mathbb{S}^{p-1}$ and every $i \in [n]$, define the operator χ as in (C.5), then*

$$n^{-1} \mathbb{E} \chi[\beta]_i = \theta \beta_i (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) \quad (\text{C.6})$$

where $\mathbf{s}_i = \sum_{\ell \neq i} \beta_\ell \mathbf{x}_{0\ell}$. Suppose $(\mathbf{a}_0, \theta, |\tau|)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$ and $\lambda = c_\lambda / \sqrt{|\tau|}$ for some $c_\lambda > 1/5$ and $\sigma_i = \text{sign}(\beta_i)$, then there exists some numerical constant \bar{c} such that if $c_\mu \leq \bar{c}$ then for every $\mathbf{a} \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and every $i \in [n]$, (C.6) has upper bound

$$\sigma_i n^{-1} \mathbb{E} \chi[\beta]_i \leq \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta]_i} := \begin{cases} 4\theta^2 |\tau| |\beta_i| & |\beta_i| < \nu_1 \lambda, \\ \theta (|\beta_i| - \nu_1 \lambda / 2) & |\beta_i| \geq \nu_1 \lambda, \end{cases} \quad (\text{C.7})$$

and lower bound

$$\sigma_i n^{-1} \mathbb{E} \chi[\beta]_i \geq \sigma_i n^{-1} \underline{\mathbb{E} \chi[\beta]_i} =: \theta \mathcal{S}_{\nu_2 \lambda} [|\beta_i|], \quad (\text{C.8})$$

where $\nu_1 = 1 / (2\sqrt{\log \theta^{-1}})$, $\nu_2 = \sqrt{2/\pi}$.

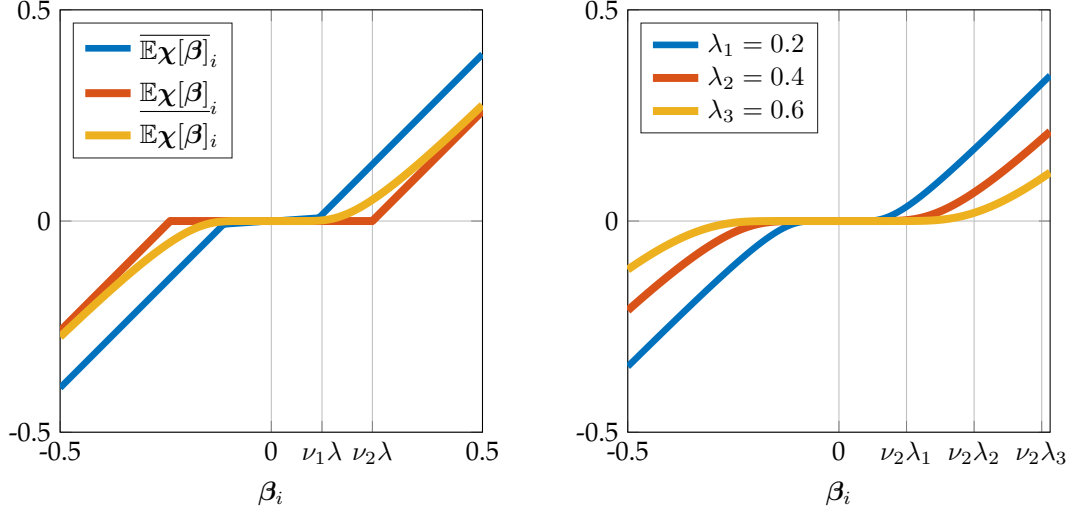


Figure 13: A numerical example of $\mathbb{E}\chi[\beta]_i$. We provide figures for the expectation of χ when entries of \mathbf{x}_0 are $2p$ -separated. Left: the yellow line is the function $\beta_i \rightarrow \beta_i (1 - \text{erf}_{\beta_i}(\lambda, 0))$ derived from (C.1), and the blue/red lines are its upper/lower bound (C.3) utilized in the analysis respectively. Right: functions of $\beta_i \rightarrow \beta_i (1 - \text{erf}_{\beta_i}(\lambda, 0))$ with different λ , the section of function of $\beta_i > \nu_2\lambda$ are close to linear.

This lemma shows the expectation of $\chi[\beta]_i$ acts like a shrinkage operation on $|\beta_i|$: for large $|\beta_i|$, it acts like a soft thresholding operation, and for small $|\beta_i|$, it reduces $|\beta_i|$ by multiplying a very small number $4\theta|\tau| \ll 1$. We rigorously prove this segmentation of χ operator as follows:

Proof. First, since $s_i[\mathbf{x}_0] \equiv_d s_j[\mathbf{x}_0]$,

$$\chi[\beta]_i = e_i^* \tilde{\mathcal{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda [\tilde{\mathcal{C}}_{\mathbf{x}_0} \beta] = \left\langle s_{-i}[\mathbf{x}_0], \mathcal{S}_\lambda [\mathbf{x}_0 * \tilde{\beta}] \right\rangle \equiv_d \left\langle s_{-j}[\mathbf{x}_0], \mathcal{S}_\lambda [s_{i-j}[\mathbf{x}_0] * \tilde{\beta}] \right\rangle = \chi[s_{j-i}[\beta]]_j$$

Thus wlog let us consider $i = 0$ and write \mathbf{x} as \mathbf{x}_0 . The random variable $\chi[\beta]_0$ can be written sum of random variables as:

$$\chi[\beta]_0 = \left\langle \mathbf{x}, \mathcal{S}_\lambda \left[\beta_0 \mathbf{x}_0 + \sum_{\ell \neq 0} \beta_\ell s_{-\ell}[\mathbf{x}] \right] \right\rangle = \sum_{j \in [n]} \mathbf{x}_j \mathcal{S}_\lambda \left[\beta_0 \mathbf{x}_j + \sum_{\ell \neq 0} \beta_\ell \mathbf{x}_{j+\ell} \right],$$

and a random variable $Z_j(\beta)$ is defined as

$$Z_j(\beta) = \mathbf{x}_j \mathcal{S}_\lambda \left[\beta_0 \mathbf{x}_j + \sum_{\ell \in [\pm p] \setminus 0} \beta_\ell \mathbf{x}_{j+\ell} \right], \quad (\text{C.9})$$

gives $\chi[\beta]_0 = \sum_{j \in [n]} Z_j(\beta)$ as sum of r.v.s. of same distribution and thus $n^{-1} \mathbb{E} \chi[\beta]_0 = \mathbb{E} Z_0(\beta)$. Define a random variable $\mathbf{s}_0 = \sum_{\ell \neq 0} \beta_\ell \mathbf{x}_\ell$, which is independent of \mathbf{x}_0 . From Lemma C.1, we can conclude

$$n^{-1} \mathbb{E} \chi[\beta]_0 = \mathbb{E}_{\mathbf{x}_0, \mathbf{s}_0} \mathbf{x}_0 \mathcal{S}_\lambda [\beta_0 \mathbf{x}_0 + \mathbf{s}_0] = \theta \beta_0 (1 - \mathbb{E}_{\mathbf{s}_0} \text{erf}_{\beta_0}(\lambda, \mathbf{s}_0)) \quad (\text{C.10})$$

so that (C.6) holds for $i = 0$, and hence for all i .

1. (Upper bound of $\mathbb{E}Z$) Wlog assume $\beta_0 \geq 0$ and write $Z = Z_0$. We derive the upper bound on $\mathbb{E}Z$ in two pieces.

(1). First, since $\mathbb{E}x_0\mathcal{S}_\lambda[0 \cdot x_0 + s_0] = 0$, we have

$$\begin{aligned}\mathbb{E}Z(\beta) &\leq \beta_0 \sup_{\beta \in [0, \beta_0]} \frac{d}{d\beta} \mathbb{E}_{x_0, s_0} x_0 \mathcal{S}_\lambda[\beta x_0 + s_0] = \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \frac{d}{d\beta} \int_{|\beta g + s_0| > \lambda} g(\beta g + s_0 - \text{sign}(\beta g + s_0) \cdot \lambda) d\mu(g) d\mu(s_0) \\ &= \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \mathbb{E}_{g, s_0} [g^2 \mathbf{1}_{\{|\beta g + s_0| > \lambda\}}] \leq \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \mathbb{E}_{g, s_0} \left[g^2 \left(\mathbf{1}_{\{|\beta g| > \frac{9\lambda}{10}\}} + \mathbf{1}_{\{|s_0| > \frac{\lambda}{10}\}} \right) \right] \\ &\leq \theta \beta_0 \left((\mathbb{E}g^6)^{1/3} \mathbb{P}[|\beta_0 g| > (9\lambda/10)]^{2/3} + \mathbb{P}[|s_0| > \lambda/10] \right)\end{aligned}\quad (\text{C.11})$$

We bound the tail probability of s_0 using [Corollary B.6](#) where

$$\mathbb{P}[|s_0| > \lambda/10] \leq \mathbb{P}[|\sum_i \beta_i x_i| > \lambda/10] \leq \theta |\tau| + 2\theta \leq 3\theta |\tau|. \quad (\text{C.12})$$

On the other hand, the first term in (C.11) can be derived by pdf of Gaussian r.v. [Lemma J.1](#) as:

$$(\mathbb{E}g^6)^{1/3} \mathbb{P}[|\beta_0 g| > (9\lambda/10)]^{2/3} \leq \sqrt[3]{15} \left(\frac{10\beta_0}{9\lambda\sqrt{2\pi}} \right)^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right) \leq \frac{3}{2} \left(\frac{\beta_0}{\lambda} \right)^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right). \quad (\text{C.13})$$

Combine (B.24), (C.13), when $\beta_0 < \nu_1 \lambda$, we know $e^{-\frac{\lambda^2}{4\beta_0^2}} \leq e^{\log \theta} \leq \theta |\tau|$. The first type of upper bound $\mathbb{E}Z$ is derived as

$$\forall \beta_0 \in [0, \nu_1 \lambda], \quad \mathbb{E}Z(\beta) \leq \theta \beta_0 \left(\frac{3}{2} \nu_1^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right) + 3\theta |\tau| \right) \leq 4\theta^2 |\tau| \beta_0. \quad (\text{C.14})$$

(2). The second type of upper bound can be derived directly from [Lemma C.1](#):

$$\begin{aligned}\mathbb{E}Z(\beta) &\leq \mathbb{E}_{x_0} \mathbb{E}_{s_0} x_0 \mathcal{S}_\lambda[\beta_0 x_0 + s_0] \leq \mathbb{E}_{x_0} x_0 \mathcal{S}_\lambda[\beta_0 x_0] + \mathbb{E}_{x_0} |x_0| \mathbb{E}_{s_0} |s_0| \\ &\leq \theta \cdot \left(\mathcal{S}_{\nu_1 \lambda}[\beta_0] + \varepsilon + \sqrt{2/\pi} \cdot \mathbb{E}|s_0| \right),\end{aligned}\quad (\text{C.15})$$

where $\mathbb{E}|s|$ can be bounded with $\|\beta\|_2$ and $\theta |\tau| < c_\mu \theta_{\log}$ from [Lemma B.5](#). When $c_\mu < \frac{1}{10}$, observe that

$$\mathbb{E}|s| \leq \sqrt{\sum_{\ell} \mathbb{E}x_{\ell}^2 \beta_{\ell}^2} \leq \sqrt{\theta} (\|\beta_{\tau}\|_2 + \|\beta_{\tau^c}\|_2) \leq \sqrt{\theta} (1 + c_\mu) + \frac{c_\mu \theta_{\log}}{|\tau|} \leq \frac{2c_\mu \theta_{\log}}{\sqrt{|\tau|}}. \quad (\text{C.16})$$

Now choose $\varepsilon = \theta \leq \frac{c_\mu \theta_{\log}}{|\tau|}$, so that $\nu'_1 = \nu_1 = \frac{\sqrt{\theta_{\log}}}{2}$ in (C.15). Since $c_\mu < \frac{c_\lambda}{25}$ we gain

$$\begin{aligned}\mathbb{E}Z(\beta) &\leq \theta \left(\mathcal{S}_{\nu_1 \lambda}[\beta_0] + \frac{c_\mu \theta_{\log}}{|\tau|} + \sqrt{\frac{2}{\pi}} \cdot \frac{2c_\mu \theta_{\log}}{\sqrt{|\tau|}} \right) \leq \theta \left(\mathcal{S}_{\nu_1 \lambda}[\beta_0] + \frac{3c_\mu \theta_{\log}}{\sqrt{|\tau|}} \right) \\ &\leq \theta \left(\mathcal{S}_{\nu_1 \lambda}[\beta_0] + \frac{\sqrt{\theta_{\log}}}{5} \lambda \right) \leq \theta \left(\mathcal{S}_{\nu_1 \lambda}[\beta_0] + \frac{1}{2} \nu_1 \lambda \right)\end{aligned}\quad (\text{C.17})$$

(3). Combine both (C.14) and (C.17), we can thus conclude that

$$\mathbb{E}Z(\beta) := \overline{\mathbb{E}Z(\beta)} \leq \begin{cases} 4\theta^2 |\tau| \beta_0 & \beta_0 \leq \nu_1 \lambda \\ \theta (\beta_0 - \frac{\nu_1}{2} \lambda) & \beta_0 > \nu_1 \lambda \end{cases}. \quad (\text{C.18})$$

2. (Lower bound of $\mathbb{E}Z$) On the other hand, for the lower bound for $\mathbb{E}Z$, use the fact that $\text{erf}_\beta(\lambda, s)$ is concave in s_0 , we have

$$\mathbb{E}Z(\beta) = \mathbb{E}_{s_0} \mathbb{E}_{x_0} x_0 \mathcal{S}_\lambda[\beta_0 x_0 + s_0] = \theta \cdot \mathbb{E}_{s_0} \left[\beta_0 - \frac{\beta_0}{2} \cdot \text{erf}\left(\frac{\lambda - s_0}{\sqrt{2}|\beta_0|}\right) - \frac{\beta_0}{2} \cdot \text{erf}\left(\frac{\lambda + s_0}{\sqrt{2}|\beta_0|}\right) \right]$$

$$\geq \theta \left(\beta_0 - \beta_0 \cdot \operatorname{erf} \left(\frac{\lambda}{\sqrt{2} |\beta_0|} \right) \right) \geq \theta \cdot \mathcal{S}_{\nu'_2 \lambda} [\beta_0] =: \underline{\mathbb{E}Z(\beta)}. \quad (\text{C.19})$$

The proof of $\beta_0 < 0$ is in the same vein. For cases of $i \neq 0$, since $\chi[\beta]_i \equiv_d \chi[s_{-i}[\beta]]_0$, replace β_0 with β_i we obtain the desired result. ■

Monotonicity of χ . Another convenient fact of $\mathbb{E}\chi[\beta]_i$ is that it is monotone increasing w.r.t. $|\beta_i|$. The monotonicity is clear in Figure 13; it is demonstrated rigorously with the following lemma:

Lemma C.3 (Monotonicity of $\mathbb{E}\chi(\beta)$). *Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and $|\tau|, c_\mu$ such that $(\mathbf{a}_0, \theta, |\tau|)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{|\tau|}$ in φ_{ℓ^1} where $c_\lambda \in [0, \frac{1}{4}]$, then there exists some numerical constant $\bar{c} > 0$, such that if $c_\mu < \bar{c}$, the expectation $|\mathbb{E}[\chi[\beta]]_i|$ is monotone increasing in $|\beta_i|$. In other words, if $|\beta_i| > |\beta_j|$ then*

$$\sigma_i \mathbb{E}\chi[\beta]_i \geq \sigma_j \mathbb{E}\chi[\beta]_j \quad (\text{C.20})$$

where $\sigma_i = \text{sign}(\beta_i)$.

The proof first operate simple calculus and then followed by studying cases of $|\beta_i| - |\beta_j|$ when either it is smaller are larger then λ .

Proof. 1. (Monotonicity by gradient negativity) Wlog assume $\beta_i > \beta_j > 0$, and from Lemma C.2 we can write $(n\theta)^{-1} \mathbb{E}\chi[\beta]_i = \beta_i (1 - \mathbb{E}_{\mathbf{s}_i} \operatorname{erf}_{\beta_i}(\lambda, \mathbf{s}_i))$. Consider $t \in [0, 1]$ and define $\ell(t) = t\beta_i - t\beta_j$. Write the random variable $\mathbf{s}_{ij} = \sum_{\ell \neq i, j} \beta_\ell \mathbf{x}_\ell$. Define h as a function of t such that

$$\begin{aligned} h(t) &= \mathbb{E}_{\mathbf{s}_{ij}} \left[((1-t)\beta_i + t\beta_j) (1 - \operatorname{erf}_{(1-t)\beta_i + t\beta_j}(\lambda, ((1-t)\beta_j + t\beta_i)x + \mathbf{s}_{ij})) \right] \\ &= \mathbb{E}_{\mathbf{s}_{ij}} \left[(\beta_i - \ell(t)) (1 - \operatorname{erf}_{\beta_i - \ell(t)}(\lambda, x \cdot (\beta_j + \ell(t)) + \mathbf{s}_{ij})) \right]. \end{aligned} \quad (\text{C.21})$$

Notice that $\mathbb{E}\chi[\beta]_i = h(0)$ and $\mathbb{E}\chi[\beta]_j = h(1)$ respectively, thus it suffices to prove $h'(t) < 0$ for all $t \in [0, 1]$. Write f as pdf of standard Gaussian r.v. where

$$\operatorname{erf}_\beta(\lambda, \mathbf{s}_{ij}) = \int_0^{\frac{\lambda + \mathbf{s}_{ij}}{\beta}} f(z) dz + \int_0^{\frac{\lambda - \mathbf{s}_{ij}}{\beta}} f(z) dz,$$

and use chain rule:

$$\begin{aligned} h'(t) &= \mathbb{E}_{\mathbf{s}_{ij}} \left[(\beta_j - \beta_i) (1 - \operatorname{erf}_{\beta_i - \ell(t)}(\lambda, x \cdot (\beta_j + \ell(t)) + \mathbf{s}_{ij})) \right. \\ &\quad - (\beta_i - \ell(t)) \cdot \frac{d}{dt} \left(\frac{\lambda + x \cdot (\beta_j + \ell(t)) + \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \cdot f \left(\frac{\lambda + x \cdot (\beta_j + \ell(t)) + \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \\ &\quad \left. - (\beta_i - \ell(t)) \cdot \frac{d}{dt} \left(\frac{\lambda - x \cdot (\beta_j + \ell(t)) - \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \cdot f \left(\frac{\lambda - x \cdot (\beta_j + \ell(t)) - \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \right] \\ &= (\beta_j - \beta_i) \mathbb{E}_{\mathbf{s}_{ij}} \left[1 - \operatorname{erf}_{\beta_i - \ell(t)}(\lambda, x \cdot (\beta_j + \ell(t)) + \mathbf{s}_{ij}) \right. \\ &\quad + \underbrace{\left(\frac{\lambda + x(\beta_j + \ell(t)) + \mathbf{s}_{ij}}{\beta_i - \ell(t)} + x \right)}_{z_{\lambda+}} \cdot f \left(\frac{\lambda + x(\beta_j + \ell(t)) + \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \\ &\quad \left. + \underbrace{\left(\frac{\lambda - x(\beta_j + \ell(t)) - \mathbf{s}_{ij}}{\beta_i - \ell(t)} - x \right)}_{z_{\lambda-}} \cdot f \left(\frac{\lambda - x(\beta_j + \ell(t)) - \mathbf{s}_{ij}}{\beta_i - \ell(t)} \right) \right] \\ &= (\beta_j - \beta_i) \mathbb{E}_{\mathbf{s}_{ij}} \left[1 - \int_0^{z_{\lambda+}} f(z) dz - \int_0^{z_{\lambda-}} f(z) dz + (z_{\lambda+} + x)f(z_{\lambda+}) + (z_{\lambda-} - x)f(z_{\lambda-}) \right]. \end{aligned} \quad (\text{C.22})$$

Consider the term only related to $z_{\lambda+}$, condition on cases that it is either positive or negative, observe that

$$\begin{cases} \mu_{+-} := \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda+} \leq 0} \left[\int_0^{z_{\lambda+}} f(z) dz - z_{\lambda+} f(z_{\lambda+}) \right] = \mathbb{E}_{x, \mathbf{s} | z_{\lambda+} \leq 0} \left[- \int_0^{-z_{\lambda+}} f(z) dz - z_{\lambda+} f(z_{\lambda+}) \right] \leq 0 \\ \mu_{++} := \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda+} > 0} \left[\int_0^{z_{\lambda+}} f(z) dz - z_{\lambda+} f(z_{\lambda+}) \right] \leq \min \left\{ \frac{1}{2}, \frac{1}{\sqrt{2\pi}} \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda+} > 0} z_{\lambda+} \right\} \end{cases},$$

where the negativity of the first equation can be observed by writing $v = -z_{\lambda+}$ and take derivative:

$$\begin{cases} - \int_0^v f(z) dz + v \cdot f(v) = 0 & v = 0 \\ \frac{d}{dv} \left\{ - \int_0^v f(z) dz + v \cdot f(v) \right\} = -f(v) + f(v) + v \cdot f'(v) < 0 & v > 0 \end{cases};$$

and similarly for $z_{\lambda-}$:

$$\begin{cases} \mu_{--} := \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda-} \leq 0} \left[\int_0^{z_{\lambda-}} f(z) dz - z_{\lambda-} f(z_{\lambda-}) \right] \leq 0 \\ \mu_{-+} := \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda-} > 0} \left[\int_0^{z_{\lambda-}} f(z) dz - z_{\lambda-} f(z_{\lambda-}) \right] \leq \min \left\{ \frac{1}{2}, \frac{1}{\sqrt{2\pi}} \mathbb{E}_{x, \mathbf{s}_{ij} | z_{\lambda-} > 0} z_{\lambda-} \right\} \end{cases},$$

then combine every term to (C.22) using tower property and from assumption $\beta_j - \beta_i < 0$ we obtain

$$\begin{aligned} (C.22) &\leq (\beta_j - \beta_i) (1 - \mathbb{P}[z_{\lambda+} > 0] \cdot \mu_{++} - \mathbb{P}[z_{\lambda-} > 0] \cdot \mu_{-+} + \mathbb{E}_{x, \mathbf{s}_{ij}} [x(f(z_{\lambda+}) - f(z_{\lambda-}))]) \\ &\leq (\beta_j - \beta_i) \left(1 - \min \left\{ \frac{\mathbb{P}[z_{\lambda+} > 0]}{2}, \frac{\mathbb{E}[z_{\lambda+}]}{\sqrt{2\pi}} \right\} - \min \left\{ \frac{\mathbb{P}[z_{\lambda-} > 0]}{2}, \frac{\mathbb{E}[z_{\lambda-}]}{\sqrt{2\pi}} \right\} - \frac{\theta}{\sqrt{2\pi}} \cdot \mathbb{E}[g] \right), \end{aligned} \quad (C.23)$$

where g is standard Gaussian r.v..

2. (Cases of varying β_i, β_j) Let $c_\lambda < \frac{1}{4}$. Suppose $\beta_i - \ell(t) \leq \frac{1}{4\sqrt{|\tau|}}$. Recall that $\|\beta_\tau\|_2^2 \geq 1 - 3c_\mu$. We are going to show there is at least one of the entry $\beta_* \in \{\beta_r\}_{r \in \tau \neq i, j} \cup \{\beta_j + \ell(t)\}$ is greater than $\frac{0.85}{\sqrt{|\tau|}}$. First, if both $i, j \notin \tau$, the lower bound is immediate since $\beta_*^2 = \|\beta_\tau\|_\infty^2 > \frac{1-3c_\mu}{|\tau|}$. On the other hand if at least one of i, j is in τ and all other β_τ entries are small where $\|\beta_{\tau \setminus \{i, j\}}\|_\infty^2 < \frac{1-3c_\mu}{|\tau|}$, then we know via norm inequalities,

$$(\beta_i + \beta_j)^2 > \beta_i^2 + \beta_j^2 > \|\beta_\tau\|_2^2 - (|\tau| - 1) \|\beta_{\tau \setminus \{i, j\}}\|_\infty^2 \geq \frac{1 - 3c_\mu}{|\tau|}, \quad (C.24)$$

which implies if $c_\mu < \frac{1}{100}$,

$$\beta_* = \beta_j + \ell(t) = (\beta_i + \beta_j) - (\beta_i - \ell(t)) \geq \frac{\sqrt{1-3c_\mu}}{\sqrt{|\tau|}} - \frac{1}{4\sqrt{|\tau|}} \geq \frac{0.72}{\sqrt{|\tau|}}. \quad (C.25)$$

In this case, adopt result from Corollary B.6 such that $\mathbb{P}[\sum \beta_\ell x_\ell > \lambda/10] \leq 3\theta |\tau| \leq .01$, we have

$$\begin{aligned} \mathbb{P}[z_{\lambda-} > 0] &= \mathbb{P}[z_{\lambda+} > 0] = 1 - \mathbb{P}[x(\beta_j + \ell(t)) + \mathbf{s}_{ij} < -\lambda] \\ &\leq 1 - \mathbb{P}[\mathbf{x}_* \beta_* < -11\lambda/10] \cdot \mathbb{P}[x(\beta_j + \ell(t)) + \mathbf{s}_{ij} - \mathbf{x}_* \beta_* < \lambda/10] \\ &\leq 1 - \theta \cdot \mathbb{P} \left[\mathbf{g}_* \cdot \frac{0.72}{\sqrt{|\tau|}} < \frac{-11c_\lambda}{10\sqrt{|\tau|}} \right] \cdot \left(1 - \mathbb{P} \left[\sum \beta_\ell x_\ell > \frac{\lambda}{10} \right] \right) \\ &\leq 1 - \theta \cdot \mathbb{P}[0.72 \cdot \mathbf{g}_* \leq -1.1 \cdot 0.25] \cdot (1 - 3c_\mu) \\ &\leq 1 - 0.35\theta. \end{aligned} \quad (C.26)$$

On the other hand, when $\beta_i - \ell(t) \geq \frac{1}{4\sqrt{|\tau|}}$, both $z_{\lambda+}, z_{\lambda-}$ are upper bounded via $|\tau| \theta \leq \frac{1}{800}$ such as:

$$\mathbb{E}_{x, \mathbf{s}_{ij}} |z_{\lambda-}| = \mathbb{E}_{x, \mathbf{s}_{ij}} |z_{\lambda+}| \leq \mathbb{E}_{x, \mathbf{s}_{ij}} \frac{\lambda + |x(\beta_j + \ell(t)) - \mathbf{s}_{ij}|}{\beta_i - \ell(t)} \leq 1 + 4\sqrt{|\tau|} \cdot \left(\mathbb{E}_{x, \mathbf{s}_{ij}} |x(\beta_j + \ell(t)) - \mathbf{s}_{ij}|^2 \right)^{1/2}$$

$$\leq 1 + 4\sqrt{|\boldsymbol{\tau}|\theta} \|\boldsymbol{\beta}\|_2 \leq 1 + 4\sqrt{|\boldsymbol{\tau}|\theta} \left(1 + c_\mu + \frac{c_\mu}{\sqrt{\theta}|\boldsymbol{\tau}|}\right) \leq 1.2. \quad (\text{C.27})$$

Combine (C.23), (C.26) we have

$$h'(t) \leq (\beta_j - \beta_i) \left(1 - 2 \cdot \frac{(1 - 0.35\theta)}{2} - \frac{\theta}{\sqrt{2\pi}} \cdot \sqrt{\frac{2}{\pi}}\right) \leq 0.03\theta(\beta_j - \beta_i) < 0, \quad (\text{C.28})$$

and combine (C.23), (C.27) and $\theta < c_\mu$ we have

$$h'(t) \leq (\beta_j - \beta_i) \left(1 - 2 \cdot \frac{1.2}{\sqrt{2\pi}} - \frac{\theta}{\sqrt{2\pi}} \cdot \sqrt{\frac{2}{\pi}}\right) \leq 0.03(\beta_j - \beta_i) < 0, \quad (\text{C.29})$$

which proves the monotonicity. ■

Finite sample deviation of χ . When the signal length of \mathbf{y} is sufficiently large, operator χ will be enough close to its expected value.

Corollary C.4 (Finite sample deviation of $\chi(\boldsymbol{\beta})$). *Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda/\sqrt{k}$ in φ_{ℓ^1} for some $c_\lambda > 1/5$, then there exists some numerical constants $C, c, \bar{c} > 0$, such that if $n \geq Cp^5\theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - 3/n$, for every $\mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and every $i \in [n]$, we have:*

$$|n^{-1}\chi[\boldsymbol{\beta}]_i - n^{-1}\mathbb{E}\chi[\boldsymbol{\beta}]_i| \leq c\theta/p^{3/2}, \quad (\text{C.30})$$

Proof. See [Appendix I.1](#) ■

D Euclidean Hessian as logic function in shift space

We can express the (pseudo) curvature (4.10) in direction $v \in \mathbb{S}^{p-1}$ in terms of the correlation $\gamma = C_{a_0}^* \iota v$ between v and a_0 , giving

$$v^* \tilde{\nabla}^2 \varphi_{\ell^1}(a) v = -\gamma^* \tilde{C}_{x_0} P_I \tilde{C}_{x_0} \gamma,$$

where

$$I(a) = \text{supp} \left(\mathcal{S}_\lambda \left[\tilde{C}_{x_0} C_{a_0}^* \iota a \right] \right) = \left\{ i \in [n] \mid \left| x_0 * \tilde{\beta} \right|_i > \lambda \right\}. \quad (\text{D.1})$$

The i -th diagonal entry of $\tilde{C}_{x_0} P_{I(a)} \tilde{C}_{x_0}$ is

$$-e_i^* \tilde{C}_{x_0} P_{I(a)} \tilde{C}_{x_0} e_i = -\left\| P_{I(a)} \tilde{C}_{x_0} e_i \right\|_2^2 = -\left\| P_{I(a)} s_{-i}[x_0] \right\|_2^2, \quad (\text{D.2})$$

which is the core component for us to study the curvature of objective φ_{ℓ^1} . We illustrate the expectation of diagonal term of Hessian in Lemma D.2 and Corollary D.3, whose figure of visualized $\left\| P_{I(a)} s_{-i}[x_0] \right\|_2$ is shown in Figure 13. Lastly, we also prove the off-diagonal terms $e_i^* \tilde{C}_{x_0} P_{I(a)} \tilde{C}_{x_0} e_j$ of Hessian is likely inconsequential in calculation of curvature in Lemma D.4.

Expectation of Hessian diagonals. We expect the Hessian to have stronger negative component in the $s_i[a_0]$ direction as $\left\| P_{I(a)} s_{-i}[x_0] \right\|_2^2$ becomes larger. This term can be tremendously simplified when x_0 is very sparse: suppose all entries of its support I_0 are separated by at least $2p - 1$ samples, then by implementing the definition of support from (D.1), we can derive

$$-\left\| P_{I(a)} s_{-i}[x_0] \right\|_2^2 = -\sum_{j \in I_0} x_{0j}^2 \mathbf{1}_{\left\{ \left| \sum_{\ell} \beta_\ell x_{0(\ell+j-i)} \right| > \lambda \right\}} \underbrace{\quad}_{\text{sep.}} - \sum_{j \in I_0} g_j^2 \mathbf{1}_{\left\{ |\beta_i g_j| > \lambda \right\}}, \quad (\text{D.3})$$

where $\mathbf{1}$ is the indicator function and g_j are independent standard Gaussian r.v.s.. In expectation, the summands in (D.3) acts like a smoothed logic function on entry β_i :

Lemma D.1 (Gaussian smoothed indicator). *Let $g \sim \mathcal{N}(0, 1)$, then for any $b, s \in \mathbb{R}$ and $\lambda > 0$.*

$$\mathbb{E}_g \left[g^2 \mathbf{1}_{\{|b \cdot g + s| > \lambda\}} \right] = 1 - \text{erf}_b(\lambda, s) + f_b(\lambda, s), \quad (\text{D.4})$$

where

$$f_b(\lambda, s) = \frac{1}{\sqrt{2\pi}} \left[\left(\frac{\lambda + s}{|b|} \right) e^{-\frac{(\lambda+s)^2}{2b^2}} + \left(\frac{\lambda - s}{|b|} \right) e^{-\frac{(\lambda-s)^2}{2b^2}} \right]. \quad (\text{D.5})$$

Proof. The proof can be derived via same calculation of integrals in Lemma C.1. \blacksquare

Although the definition (D.4) seems incomprehensible at first glance, we can actually interpret it as a smoothed indicator function which compares $|b|$ to the threshold $\sqrt{2/\pi} \lambda$. Once we assign $s = 0$, then we can see that $\mathbb{E} g^2 \mathbf{1}_{\{|b \cdot g| > \lambda\}}$ is an increasing function of $|b|$. Moreover by assigning different values for $|b|$ we obtain:

$$\mathbb{E} g^2 \mathbf{1}_{\{|b \cdot g| > \lambda\}} \approx \begin{cases} 1, & |b| \approx 1 \\ 1/2, & |b| \approx \sqrt{2/\pi} \lambda \\ 0, & |b| \approx 0 \end{cases}. \quad (\text{D.6})$$

Relate (D.6) to (D.3), when $|\beta_i|$ is close to 1 then we expect $-\frac{1}{n\theta} \left\| P_{I(a)} s_{-i}[x_0] \right\|_2^2$ to be close to -1 , and it increases to 0 as $|\beta_i|$ decreases, suggests that the Euclidean Hessian at point a has stronger negative component at $s_i[a_0]$ direction if $|\langle a, s_i[a_0] \rangle|$ is larger. See Figure 14 for a numerical example. This phenomenon can be extend beyond the idealistic separating case as follows:

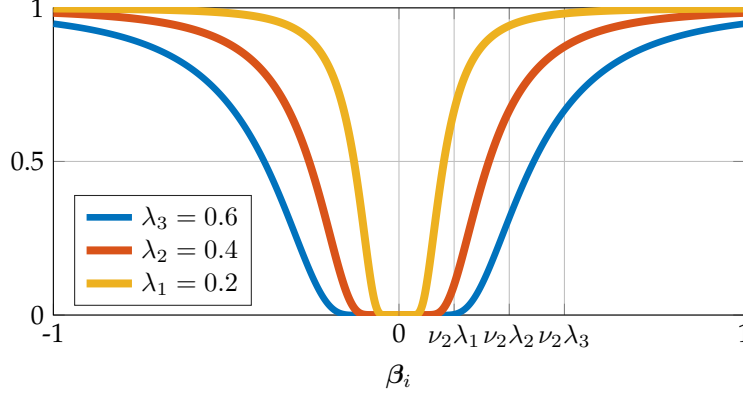


Figure 14: A numerical example for $\mathbb{E} \|P_{I(a)s_i}[x_0]\|_2^2$. We provide a figure to illustrate the expectation of $-\frac{1}{n\theta} \|P_{I(a)s_i}[x_0]\|_2^2$ when entries of x_0 are $2p$ -separated, as a function plot of $\beta_i \rightarrow 1 - \text{erf}_{\beta_i}(\lambda, 0) + f_{\beta_i}(\lambda, 0)$ from (D.4) with different λ . When $|\beta_i| \approx \nu_2 \lambda$ where $\nu_2 = \sqrt{2/\pi}$, then the its function value is close to 0.5. If $|\beta_i|$ is much larger then λ its value grow to 1, implies there is a negative curvature at $s_i[a_0]$ direction. Similarly if $|\beta_i|$ is much smaller then λ the function value is 0 thus the curvature is positive in $s_i[a_0]$ direction.

Lemma D.2 (Expected Hessian diagonals). *Let $x_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ and $\lambda > 0$, define the set $I(a)$ in (D.1), write $s_i = \sum_{\ell \neq i} \beta_\ell x_{0\ell}$, then for every $a \in \mathbb{S}^{p-1}$ and $i \in [n]$:*

$$n^{-1} \mathbb{E} \|P_{I(a)s-i}[x_0]\|_2^2 = \theta [1 - \mathbb{E}_{s_i} \text{erf}_{\beta_i}(\lambda, s_i) + \mathbb{E}_{s_i} f_{\beta_i}(\lambda, s_i)] \quad (\text{D.7})$$

Proof. Write x_0 as x . Observe that $y * \tilde{a} = x_0 * \tilde{\beta} = \sum_{\ell} \beta_\ell s_{-\ell}[x_0]$. Thus for any $j \in [n]$ and $i \in [\pm p]$:

$$(y * \tilde{a})_{j-i} = \left(\beta_i s_{-i}[x] + \sum_{\ell \neq i} \beta_\ell s_{-\ell}[x] \right)_{j-i} = \beta_i x_j + \sum_{\ell \neq i} \beta_\ell x_{j+\ell-i} =: \beta_i x_j + s_j, \quad (\text{D.8})$$

where x_j is independent of s_j , and both x_j, s_j are symmetric and identically distributed for all $j \in [n]$. Rewrite the random variable using (D.1) as

$$\|P_{I(a)s-i}[x_0]\|_2^2 = \left\| P_{I(a)} \sum_{j \in [n]} (x_{0j} e_{j-i}) \right\|_2^2 = \sum_{j \in [n]} x_{0j}^2 \mathbf{1}_{\{|y * \tilde{a}|_{j-i} > \lambda\}} = \sum_{j \in [n]} x_{0j}^2 \mathbf{1}_{\{|\beta_i x_{0j} + s_j| > \lambda\}}$$

Write $x = g \circ \omega$ as composition of Gaussian/Bernoulli r.v.s., the expectation has a simple form:

$$\mathbb{E} \|P_{I(a)s-i}[x_0]\|_2^2 = n\theta \cdot \mathbb{E} g_0^2 \mathbf{1}_{\{|\beta_i g_0 + s_0| > \lambda\}} = n\theta \cdot \mathbb{E} (1 - \text{erf}_{\beta_i}(\lambda, s_i) + f_{\beta_i}(\lambda, s_i))$$

where $s_i = \sum_{\ell \neq i} x_{0\ell} \beta_\ell$ with $x_{0i} \sim_{\text{i.i.d.}} \text{BG}(\theta)$, yielding the claimed expression. \blacksquare

Finite sample deviation of Hessian diagonals. When the signal length of y is sufficiently large, then i -th diagonal term for Hessian $\|P_{I(a)s-i}[x_0]\|_2^2$ will be close enough to its expected value.

Corollary D.3 (Large sample deviation of curvature). *Suppose $x_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that (a_0, θ, k) satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_{ℓ^1} for some $c_\lambda > 1/5$, then there exists some numerical constant $C, c, \bar{c} > 0$, such that if $n \geq Cp^4 \theta^{-1} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - 3/n$, for every $a \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and every $i \in [n]$, we have:*

$$\left| n^{-1} \|P_{I(a)s-i}[x_0]\|_2^2 - n^{-1} \mathbb{E} \|P_{I(a)s-i}[x_0]\|_2^2 \right| \leq c\theta/p \quad (\text{D.9})$$

Proof. See Appendix I.2. \blacksquare

Hessian off-diagonal terms near solution. The off-diagonal entries of Hessian in general are much smaller than the diagonal entries; however, it affects the region near sign shifts of \mathbf{a}_0 the most where we need to show strong convexity in the region. We provide an upper bound for off-diagonal entries in the vicinity of signed shifts. In these regions, only one entry of the correlations $|\beta_{(0)}|$ is large and the rest is small.

Lemma D.4 (Hessian off-diagonal term near solution). *Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Let $\lambda = c_\lambda/\sqrt{k}$ with $c_\lambda > 1/5$, then there exists some numerical constant $C, \bar{c} > 0$ such that if $n \geq C\theta^{-4} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - 4/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, where $|\beta_{(1)}| \leq \frac{1}{4 \log \theta^{-1}} \lambda$ and every $i \neq j \in [\pm p] \setminus \{(0)\}$, we have*

$$|s_i[\mathbf{x}_0]^*| \mathbf{P}_{I(\mathbf{a})} |s_j[\mathbf{x}_0]| < 8n\theta^3 \quad (\text{D.10})$$

Proof. Write $\theta_{\log} = -1/\log \theta$ and \mathbf{x}_0 as $\mathbf{x} = \omega \circ \mathbf{g}$. Wlog let β_0 be the largest correlation $\beta_{(0)}$. Define random variables $s' = \langle \beta_{\tau \setminus \{0, i, j\}}, \mathbf{x}_{\tau \setminus \{0, i, j\}} \rangle$. Firstly via [Corollary B.7](#) we have $\mathbb{P}[|s'| > 0.4\lambda] \leq 2\theta$; also define $s = \langle \beta_{\tau^c \setminus \{0, i, j\}}, \mathbf{x}_{\tau^c \setminus \{0, i, j\}} \rangle$, and base on [Corollary B.6](#) we have $\mathbb{P}[|s| > \lambda/10] \leq 2\theta$. Expand the $(-i, -j)$ -th cross term with $\theta < 0.1$ we have:

$$\begin{aligned} \mathbb{E} |s_{-i}[\mathbf{x}]^*| \mathbf{P}_{I(\mathbf{a})} |s_{-j}[\mathbf{x}]| &= \mathbb{E} \sum_{k \in [n]} |\mathbf{x}_{k+i} \mathbf{x}_{k+j}| \mathbf{1}_{\{|\beta_0 \mathbf{x}_k + \beta_i \mathbf{x}_{k+i} + \beta_j \mathbf{x}_{k+j} + s + s'| > \lambda\}} \\ &= n\theta^2 \cdot \mathbb{E} |\mathbf{g}_i \mathbf{g}_j| \mathbf{1}_{\{|\beta_0 \mathbf{x}_0 + \beta_i \mathbf{g}_i + \beta_j \mathbf{g}_j + s + s'| > \lambda\}} \\ &\leq n\theta^2 \cdot \mathbb{E} [|\mathbf{g}_i \mathbf{g}_j| (2\mathbf{1}_{\{|\beta_i \mathbf{g}_i| > \lambda/4\}} + \mathbb{P}[\mathbf{x}_0 \neq 0] + \mathbb{P}[|s| > 0.1\lambda] + \mathbb{P}[|s'| > 0.4\lambda])] \\ &\leq n\theta^2 \cdot (\exp(-\log^2 \theta^{-1}) + \theta + 2\theta + 2\theta) \\ &\leq 6n\theta^3. \end{aligned} \quad (\text{D.11})$$

Write (D.10) as two summation of independent random variables with $t = j - i$ by separating sum into two sets J_{t1}, J_{t2} defined in (A.4) with both $|J_{t1}|, |J_{t2}| < n\theta^2$ with probability at least $1 - 2/n$ from [Lemma A.1](#)

$$\mathbb{E} |s_{-i}[\mathbf{x}]^*| \mathbf{P}_{I(\mathbf{a})} |s_{-j}[\mathbf{x}]| = \sum_{(k-i) \in I(\mathbf{a})} |\mathbf{x}_k| |\mathbf{x}_{k+t}| = \sum_{(k-i) \in I(\mathbf{a}) \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| + \sum_{(k-i) \in I(\mathbf{a}) \cap J_{t2}} |\mathbf{g}_k| |\mathbf{g}_{k+t}|,$$

whose first summands can be upper bounded w.h.p. via Bernstein inequality [Lemma J.4](#) with $(\sigma^2, R) = (1, 1)$ and writes $\mathcal{C} := \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu)) \cap \left\{ \mathbf{a} \mid |\beta_{(1)}| \leq \frac{1}{4 \log \theta^{-1}} \lambda \right\}$, then we have

$$\begin{aligned} &\mathbb{P} \left[\max_{\substack{i \neq j \in [\pm p] \setminus \{0\} \\ \mathbf{a} \in \mathcal{C}}} \left(\sum_{(k-i) \in I(\mathbf{a}) \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| - \mathbb{E} \sum_{(k-i) \in I(\mathbf{a}) \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| \right) \geq n\theta^3 \right] \\ &\mathbb{P} \left[\max_{\substack{i \neq j \in [\pm p] \setminus \{0\}}} \left(\sum_{(k-i) \in J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| - \mathbb{E} \sum_{(k-i) \in J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| \right) \geq n\theta^3 \right] \\ &\leq 4p^2 \cdot \exp \left(\frac{-n^2\theta^6}{2|J_{t1}| + 2n\theta^3} \right) \leq \exp \left(8 \log p - \frac{n^2\theta^6}{3n\theta^2} \right) \leq \exp \left(-\frac{n\theta^4}{10} \right) \leq \frac{1}{n} \end{aligned} \quad (\text{D.12})$$

when $n = C\theta^{-4} \log p$ with $C > 10^4$ and $\theta \log^2 \theta^{-1} \geq 1/p$. Thus for all $i \neq j \in [\pm p] \setminus \{0\}$ and \mathbf{a} satisfies our condition of lemma, from (D.11) and (D.12) we can conclude :

$$|s_{-i}[\mathbf{x}]^*| \mathbf{P}_{I(\mathbf{a})} |s_{-j}[\mathbf{x}]| \leq \sum_{I(\mathbf{a}) \cap J_{t1}} \mathbb{E} |\mathbf{g}_k| |\mathbf{g}_{k+t}| + \sum_{I(\mathbf{a}) \cap J_{t2}} \mathbb{E} |\mathbf{g}_k| |\mathbf{g}_{k+t}| + 2n\theta^3 \leq 8n\theta^3$$

which holds with probability at least $1 - 2/n - 2 \cdot 1/n = 1 - 4/n$ base on [Lemma A.1](#) and (D.12). \blacksquare

E Geometric relation between ρ and ℓ^1 -norm

In this section, we discuss how to ensure that the smooth sparsity surrogate ρ approximates $\|\cdot\|_1$ accurately enough that guarantees φ_ρ inherits the good properties of φ_{ℓ^1} . We prove several lemmas which allow us to transfer properties of φ_{ℓ^1} to φ_ρ . Our result does not pertain to the suggested pseudo-Huber surrogate $\rho(x)_i = \sqrt{x_i^2 + \delta^2}$ in the main script, and is general for a class of function class defined in [Definition E.2](#) that is smooth and well approximates ℓ^1 when the proper smoothing parameter δ is chosen from the result of [Lemma E.6](#). In particular we ask the regularizer $\rho_\delta(x)$ to be uniformly bounded to $|x|$ by $\delta/2$:

$$\forall x \in \mathbb{R}, \quad |\rho_\delta(x) - |x|| \leq \delta/2 \quad (\text{E.1})$$

then if $\delta \rightarrow 0$ we have for every \mathbf{a} near subspace,

$$\|\text{prox}_{\lambda\ell^1}[\tilde{\mathbf{a}} * \mathbf{y}] - \text{prox}_{\lambda\rho_\delta}[\tilde{\mathbf{a}} * \mathbf{y}]\|_2 \rightarrow 0, \quad (\text{E.2})$$

$$\|\nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_{\rho_\delta}(\mathbf{a})\|_2 \rightarrow 0, \quad (\text{E.3})$$

$$\|\tilde{\nabla}^2\varphi_{\ell^1}(\mathbf{a}) - \nabla^2\varphi_{\rho_\delta}(\mathbf{a})\|_2 \rightarrow 0. \quad (\text{E.4})$$

An example choices of eligible smooth sparse surrogate is demonstrated in [Table 1](#).

Calculus of φ_ρ . The marginal minimizer over \mathbf{x} in (2.7) can be expressed in terms of the proximal operator [BC11] of ρ at point $\tilde{\mathbf{a}} * \mathbf{y}$:

$$\text{prox}_{\lambda\rho}[\tilde{\mathbf{a}} * \mathbf{y}] = \underset{\mathbf{x} \in \mathbb{R}^n}{\text{argmin}} \left\{ \lambda\rho(\mathbf{x}) + \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle \right\}.$$

Plugging in, we obtain

$$\varphi_\rho(\mathbf{a}) = \lambda\rho(\text{prox}_{\lambda\rho}[\tilde{\mathbf{a}} * \mathbf{y}]) + \frac{1}{2} \|\tilde{\mathbf{a}} * \mathbf{y} - \text{prox}_{\lambda\rho}[\tilde{\mathbf{a}} * \mathbf{y}]\|_2^2 - \frac{1}{2} \|\tilde{\mathbf{a}} * \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2 \quad (\text{E.5})$$

The objective function $\varphi_\rho(\mathbf{a})$ is a differentiable function of \mathbf{a} . This can be seen, e.g., by noting that

$$\varphi_\rho(\mathbf{a}) = \epsilon(\lambda\rho)(\tilde{\mathbf{a}} * \mathbf{y}) - \frac{1}{2} \|\tilde{\mathbf{a}} * \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2, \quad (\text{E.6})$$

where $\epsilon(g)(\mathbf{z}) = g(\text{prox}_g(\mathbf{z})) + \frac{1}{2} \|\mathbf{z} - \text{prox}_g(\mathbf{z})\|_2^2$ is the *Moreau envelope* of a function g . The Moreau envelope is differentiable:

Fact E.1 (Derivative of Moreau envelope, [BC11], Prop.12.29). *Let f be a proper lower semicontinuous convex function and $\lambda > 0$ then the Moreau envelope $\epsilon(\lambda f)(\mathbf{z}) = \lambda f(\text{prox}_{\lambda f}[\mathbf{z}]) + \frac{1}{2} \|\mathbf{z} - \text{prox}_{\lambda f}[\mathbf{z}]\|_2^2$ is Fréchet differentiable with $\nabla\epsilon(\lambda f)(\mathbf{z}) = \mathbf{z} - \text{prox}_{\lambda\rho}[\mathbf{z}]$.*

Furthermore, φ_ρ is twice differentiable whenever $\text{prox}_{\lambda\rho}$ is differentiable. In this case, the (Euclidean) gradient and hessian of φ_ρ are given by

$$\nabla\varphi_\rho(\mathbf{a}) = -\boldsymbol{\iota}^* \tilde{\mathbf{C}}_{\mathbf{y}} \text{prox}_{\lambda\rho}[\tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}], \quad (\text{E.7})$$

$$\nabla^2\varphi_\rho(\mathbf{a}) = -\boldsymbol{\iota}^* \tilde{\mathbf{C}}_{\mathbf{y}} \nabla \text{prox}_{\lambda\rho}[\tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}] \tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota}. \quad (\text{E.8})$$

The Riemannian gradient and hessian over \mathbb{S}^{p-1} are

$$\text{grad}[\varphi_\rho](\mathbf{a}) = -\mathbf{P}_{\mathbf{a}^\perp} \boldsymbol{\iota}^* \tilde{\mathbf{C}}_{\mathbf{y}} \text{prox}_{\lambda\rho}[\tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}], \quad (\text{E.9})$$

$$\text{Hess}[\varphi_\rho](\mathbf{a}) = -\mathbf{P}_{\mathbf{a}^\perp} \left(\boldsymbol{\iota}^* \tilde{\mathbf{C}}_{\mathbf{y}} \nabla \text{prox}_{\lambda\rho}[\tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}] \tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} - \langle \nabla\varphi_\rho(\mathbf{a}), \mathbf{a} \rangle \mathbf{I} \right) \mathbf{P}_{\mathbf{a}^\perp}. \quad (\text{E.10})$$

Surrogate class	$\rho_i(x)$	$\nabla \rho_i(x)$	$\nabla^2 \rho_i(x)$
Log hyperbolic cosine	$\frac{\delta}{2} \log(e^{2x/\delta} + e^{-2x/\delta})$	$\frac{e^{4x/\delta} - 1}{e^{4x/\delta} + 1}$	$\frac{4e^{4x/\delta}}{\delta(e^{4x/\delta} + 1)^2}$
Pseudo Huber	$\sqrt{x^2 + \delta^2}$	$\frac{x}{\sqrt{x^2 + \delta^2}}$	$\frac{\delta^2}{(x^2 + \delta^2)^{3/2}}$
Gaussian convolution	$\int x - t f_\delta(t) dt$	$\text{erf}(x/\sqrt{2\delta})$	$2f_\delta(x)$

Table 1: Classes of smooth sparse surrogate ρ and how to set its parameter. Three common classes are listed with parameter δ to tune the smoothness. All the listed functions are greater than $|x|$ pointwise and has largest distance to $|x|$ at origin where $\rho(0) - |x| \leq \delta$, satisfies the condition (E.11). Also its second order derivatives $\nabla^2 \rho_i(x)$ are monotone decreasing w.r.t. $|x|$, hence are certified to be eligible δ -smoothed ℓ^1 surrogates.

Sparse regularizer ρ as smoothed ℓ^1 function. Our analysis accommodates any sufficiently accurate smooth approximation ρ to the ℓ^1 function. The requisite sense of approximation is captured in the following definition:

Definition E.2 (δ -smoothed ℓ^1 function). *We call an additively separable function $\rho(\mathbf{x}) = \sum_{i=1}^n \rho_i(\mathbf{x}_i) : \mathbb{R}^n \rightarrow \mathbb{R}$, a δ -smoothed ℓ^1 function with $\delta > 0$ if for each $i \in [n]$, ρ_i is even, convex, twice differentiable and $\nabla^2 \rho_i(x)$ being monotone decreasing w.r.t. $|x|$, where, there exists some constant c , such that for all $x \in \mathbb{R}$:*

$$|\rho_i(x) - |x|| + c \leq \delta/2 \quad (\text{E.11})$$

The proximal operator of the ℓ^1 norm is the entrywise soft thresholding function \mathcal{S}_λ ; the proximal operator associated to a smoothed ℓ^1 function turns out to be a differentiable approximation to \mathcal{S}_λ . In particular, we will show that it approximates \mathcal{S}_λ in the following sense:

Definition E.3 ($\sqrt{\delta}$ -smoothed soft threshold). *An odd function $\mathcal{S}_\lambda^\delta[\cdot] : \mathbb{R} \rightarrow \mathbb{R}$ is a $\sqrt{\delta}$ -smoothed soft thresholding function with parameter $\delta > 0$ if it is a strictly monotone odd function and is differentiable everywhere, whose function value satisfies*

$$0 \leq \text{sign}(z) (\mathcal{S}_\lambda^\delta[z] - \mathcal{S}_\lambda[z]) \leq \sqrt{\lambda\delta}, \quad \forall z \in \mathbb{R} \quad (\text{E.12})$$

and its derivative satisfies for any given $B \in (0, \lambda)$:

$$|\nabla \mathcal{S}_\lambda^\delta[z] - \nabla \mathcal{S}_\lambda[z]| \leq \sqrt{\lambda\delta}/B, \quad ||z| - \lambda| \geq B. \quad (\text{E.13})$$

If ρ is a δ -smooth ℓ^1 function, then for all $i \in [n]$, we have that $\text{prox}_{\lambda\rho}[\mathbf{z}]_i$ is a $\sqrt{\delta}$ -smoothed soft threshold function of \mathbf{z}_i . This can be proven with the following lemma:

Lemma E.4 (Proximal operator for smoothed ℓ^1). *Suppose ρ is a δ -smoothed ℓ^1 function, then $\mathbf{z}_i \mapsto \text{prox}_{\lambda\rho}[\mathbf{z}]_i$ is a $\sqrt{\delta}$ -smoothed soft threshold function.*

Proof. We know that

$$\mathbf{x}_z := \text{prox}_{\lambda\rho}[\mathbf{z}] = \underset{\mathbf{x} \in \mathbb{R}^n}{\text{argmin}} \lambda\rho(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2. \quad (\text{E.14})$$

This optimization problem is strongly convex, and so the minimizer \mathbf{x}_z is unique. Using the stationarity condition and since ρ is separable, for all $i \in [n]$, we have $\lambda \nabla \rho_i(\mathbf{x}_{zi}) + \mathbf{x}_{zi} - \mathbf{z}_i = 0$, implies

$$\mathbf{x}_{zi} = (\text{Id} + \lambda \nabla \rho_i)^{-1}(\mathbf{z}_i). \quad (\text{E.15})$$

Since ρ_i is convex and even, $\nabla \rho_i$ is monotone increasing and odd. By inverse function theorem, we know that strict monotonicity and differentiability of $\text{Id} + \lambda \nabla \rho_i$ implies its inverse is differentiable and is a strictly monotone increasing odd function. Furthermore, it implies $\nabla \mathbf{x}_{zi}$ has the form

$$\nabla \mathbf{x}_{zi} = \nabla_i (\text{Id} + \lambda \nabla \rho_i)^{-1}(\mathbf{z}_i) = \frac{1}{\lambda \nabla^2 \rho_i(\mathbf{x}_{zi}) + 1} < 1. \quad (\text{E.16})$$

Notice that since $\nabla^2 \rho_i(x)$ is monotone decreasing when $x \geq 0$, hence $\nabla \mathbf{x}_{zi}$ is monotone increasing in $\mathbf{z}_i \geq 0$.

Now we are left to show that (E.12) and (E.13) hold, and since $\text{prox}_{\lambda \rho}[\cdot]_i$ is an odd function it suffices to consider the case when the input vector \mathbf{z}_i is nonnegative. Firstly, via convexity and entrywise bounded difference $|\rho_i(x) - |x|| \leq \delta/2$ we are going to show

$$|\nabla \rho_i(x)| \leq 1 \quad \forall x \in \mathbb{R}, \quad \nabla \rho_i(x) \geq 1 - \sqrt{\delta/\lambda} \quad \forall x \geq \sqrt{\lambda \delta}. \quad (\text{E.17})$$

Consider a positive x with $\nabla \rho_i(x) > 1 + \varepsilon$ for some $\varepsilon > 0$, by convexity if $\tilde{x} > x$ then $\nabla \rho_i(\tilde{x}) > 1 + \varepsilon$, hence

$$\rho_i(x + \delta/\varepsilon) \geq \rho_i(x) + \nabla \rho_i(x) \cdot (\delta/\varepsilon) > x - \delta/2 + (1 + \varepsilon) \cdot (\delta/\varepsilon) = (x + \delta/\varepsilon) + \delta/2,$$

contradicts the boundedness condition. Secondly, use mean value theorem we know for all $x \geq \sqrt{\lambda \delta}$:

$$\nabla \rho_i(x) \geq \frac{\rho_i(\sqrt{\lambda \delta}) - \rho_i(0)}{\sqrt{\lambda \delta} - 0} \geq \frac{(\sqrt{\lambda \delta} - \delta/2) - (0 + \delta/2)}{\sqrt{\lambda \delta} - 0} \geq 1 - \sqrt{\frac{\delta}{\lambda}}.$$

To prove (E.12), when $0 \leq \mathbf{z}_i \leq \lambda$, then $\mathcal{S}_\lambda[\mathbf{z}_i] = 0$ and $\mathbf{x}_{zi} \leq \sqrt{\lambda \delta}$ since if $\mathbf{x}_{zi} > \sqrt{\lambda \delta}$, by (E.17):

$$\lambda \nabla \rho_i(\mathbf{x}_{zi}) + \mathbf{x}_{zi} > \lambda(1 - \sqrt{\delta/\lambda}) + \sqrt{\lambda \delta} = \lambda \geq \mathbf{z}_i$$

then \mathbf{x}_{zi} violate the stationary condition in (E.15), resulting $0 \leq \mathbf{x}_{zi} - \mathcal{S}_\lambda[\mathbf{z}_i] \leq \sqrt{\lambda \delta}$ whenever $0 \leq \mathbf{z}_i \leq \lambda$. Likewise in the case of $\mathbf{z}_i \geq \lambda$ where $\mathcal{S}_\lambda[\mathbf{z}_i] = \mathbf{z}_i - \lambda$, (E.17) provides:

$$\begin{cases} \forall \mathbf{x}_{zi} > \mathbf{z}_i - \lambda + \sqrt{\lambda \delta}, & \lambda \nabla \rho_i(\mathbf{x}_{zi}) + \mathbf{x}_{zi} > \lambda(1 - \sqrt{\delta/\lambda}) + \mathbf{z}_i - \lambda + \sqrt{\lambda \delta} = \mathbf{z}_i \\ \forall \mathbf{x}_{zi} < \mathbf{z}_i - \lambda, & \lambda \nabla \rho_i(\mathbf{x}_{zi}) + \mathbf{x}_{zi} < \lambda + \mathbf{z}_i - \lambda = \mathbf{z}_i \end{cases}$$

again violates (E.15) and therefore (E.12) holds for all $\mathbf{z}_i \in \mathbb{R}$.

Lastly (E.13) is a direct result of (E.12). For all $\mathbf{z}_i \leq \lambda - B$, recall that $\nabla \mathbf{x}_{zi}$ is monotone increasing in \mathbf{z}_i :

$$\nabla \mathbf{x}_{zi} \leq \min_{y \in [\lambda - B, \lambda]} \nabla \mathbf{x}_{yi} \leq \frac{\mathbf{x}_{\lambda i} - \mathbf{x}_{(\lambda - B)i}}{\lambda - (\lambda - B)} \leq \frac{(\sqrt{\lambda \delta} + \mathcal{S}_\lambda[\lambda]) - \mathcal{S}_\lambda[\lambda - B]}{B} = \frac{\sqrt{\lambda \delta}}{B};$$

and similarly for all $\mathbf{z}_i > \lambda + B$:

$$\nabla \mathbf{x}_{zi} \geq \max_{y \in [\lambda, \lambda + B]} \nabla \mathbf{x}_{yi} \geq \frac{\mathbf{x}_{(\lambda + B)i} - \mathbf{x}_{\lambda i}}{(\lambda + B) - \lambda} \geq \frac{\mathcal{S}_\lambda[\lambda + B] - (\mathcal{S}_\lambda[\lambda] + \sqrt{\lambda \delta})}{B} = 1 - \frac{\sqrt{\lambda \delta}}{B},$$

implies (E.13) holds. ■

Approximate geometry of φ_ρ using φ_{ℓ^1} Based on (E.9)-(E.10) and denote $\check{\mathbf{C}}_y \boldsymbol{\iota} \mathbf{a} = \check{\mathbf{a}} * \mathbf{y}$, the only differences of Riemannian gradient and Hessian between φ_ρ and φ_{ℓ^1} comes from the difference of $\text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}]$ and $\text{prox}_{\lambda \|\cdot\|_1}[\check{\mathbf{a}} * \mathbf{y}]$. Thus for the purpose of obtaining good geometric approximation of φ_ρ with that of objective φ_{ℓ^1} , we may apply both Definition E.3 and Lemma E.4, together suggest if ρ is a δ -smoothed ℓ^1 function, then the i -th entry of $\text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}]$ will be $\sqrt{\lambda \delta}$ -close to the authentic soft thresholding function $\mathcal{S}_\lambda[\check{\mathbf{a}} * \mathbf{y}]_i$, and its gradient $\nabla \text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}]$ is $\sqrt{\lambda \delta}/B$ -close to $\nabla \mathcal{S}_\lambda[\check{\mathbf{a}} * \mathbf{y}]$ as long as $(\check{\mathbf{a}} * \mathbf{y})_i$ is not close to $\pm \lambda$ by distance B .

Firstly, we will show by utilizing the random structure of \mathbf{y} , such that with high probability, only a fraction of entries of $\check{\mathbf{a}} * \mathbf{y}$ will be close to $\pm \lambda$.

Lemma E.5 (Gradients discontinuity entries). For each $\mathbf{a} \in \mathbb{S}^{p-1}$, let

$$J_B(\mathbf{a}) := \left\{ i \mid \left(\widetilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a} \right)_i \in [-\lambda - B, -\lambda + B] \cup [\lambda - B, \lambda + B] \right\}. \quad (\text{E.18})$$

Suppose the subspace dimension is at most k and signal \mathbf{y} satisfies [Definition B.1](#). Let $\lambda = c_\lambda / \sqrt{k}$ and $B \leq c' \lambda \theta^2 / p \log n$ for some $c_\lambda, c' \in (0, 1)$, then there is a numerical constant $C > 0$ such that if $n \geq Cp^5 \theta^{-2} \log p$, then with probability at least $1 - 3/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, we have

$$|J_B(\mathbf{a})| \leq \frac{24c'n\theta^2}{p \log n} \quad (\text{E.19})$$

Proof. See [Appendix I.3](#). ■

The geometric approximation between φ_{ℓ^1} and φ_ρ necessarily consists of three parts: the gradient, the Hessian, and the coefficients. Here we conclude the approximation result with the following lemma:

Lemma E.6 (φ_{ℓ^1} approximates φ_ρ). Suppose $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Let $\rho \in \mathbb{R}^n \rightarrow \mathbb{R}$ be a δ -smoothed ℓ^1 function with

$$\lambda = \frac{c_\lambda}{\sqrt{k}}, \quad \delta \leq \frac{c'^4 \theta^8}{p^2 \log^2 n} \lambda \quad (\text{E.20})$$

with some $c', c_\lambda \in (0, 1)$, then there is a numerical constant $C, \bar{c} > 0$ such that if $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - 10/n$, the following statements hold simultaneously for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$:

(1). The coefficients has norm difference

$$\left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \text{prox}_{\lambda \ell^1}[\widetilde{\mathbf{a}} * \mathbf{y}] - \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \text{prox}_{\lambda \rho}[\widetilde{\mathbf{a}} * \mathbf{y}] \right\|_2 \leq c' n \theta^4. \quad (\text{E.21})$$

(2). The gradient has norm difference

$$\|\nabla \varphi_{\ell^1}(\mathbf{a}) - \nabla \varphi_\rho(\mathbf{a})\|_2 \leq c' n \theta^4. \quad (\text{E.22})$$

(3). The (pesudo) Riemmannian curvature difference is bounded in all directions $\mathbf{v} \in \mathbb{S}^{p-1}$ via

$$\forall \mathbf{v} \in \mathbb{S}^{p-1}, \quad \left| \mathbf{v}^* \left(\widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) - \text{Hess}[\varphi_\rho](\mathbf{a}) \right) \mathbf{v} \right| \leq 200c' n \theta^2. \quad (\text{E.23})$$

Proof. 1. (Coefficients) From [Lemma E.4](#), the proximal δ -smoothed ℓ^1 function satisfies

$$|\mathcal{S}_\lambda[\widetilde{\mathbf{a}} * \mathbf{y}] - \mathcal{S}_\lambda^\delta[\widetilde{\mathbf{a}} * \mathbf{y}]|_j < \sqrt{\lambda \delta} \quad \forall j \in [n].$$

Since the support of coefficient vectors are contained in $[\pm p]$, using simple norm inequality:

$$\left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda[\widetilde{\mathbf{a}} * \mathbf{y}] - \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda^\delta[\widetilde{\mathbf{a}} * \mathbf{y}] \right\|_2 \leq \sqrt{\lambda \delta n} \cdot \left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \right\|_2. \quad (\text{E.24})$$

Apply [Lemma A.5](#) by replacing \mathbf{a}_0 with standard basis \mathbf{e}_0 and extend support of $\boldsymbol{\iota}$ to $\boldsymbol{\iota}_{[\pm p]}$, notice that in this case we have $\mu = 0$. Condition on the event

$$\left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \right\|_2 \leq \left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{e}_0}^* \right\|_2 \leq \sqrt{3(1 + 2\mu p)n\theta} \leq \sqrt{3n\theta},$$

and we gain

$$(\text{E.24}) \leq \sqrt{\lambda \delta n} \cdot \sqrt{3n\theta} \leq n\sqrt{3\lambda \theta \delta} \leq c' n \theta^4.$$

2. (Gradient) From definition of Riemannian gradient (E.9) and apply similar norm bound of (E.24), and condition on the following events of Lemma A.5 holds, obtain

$$\|\nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a})\|_2 \leq \sqrt{\lambda\delta n} \cdot \|\iota^* \check{\mathbf{C}}_{\mathbf{y}}\|_2 \leq n\sqrt{3\lambda\theta(1+\mu p)\delta} \leq c'n\theta^4. \quad (\text{E.25})$$

3. (Hessian) For every realization of $J_B(\mathbf{a})$ from $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, base on Lemma E.5, condition on the event such that

$$B \leq \frac{c'\lambda\theta^2}{p \log n}, \quad |J| \leq \frac{24c'n\theta^2}{p \log n}; \quad (\text{E.26})$$

and rewrite $J_B(\mathbf{a})$ as J . Also condition on the event using Lemma A.5 and $(1+\mu p)\theta \log \theta^{-1} < 1$

$$\|\iota^* \check{\mathbf{C}}_{\mathbf{y}}\|_2 \leq \sqrt{3n}, \quad \|\iota^* \check{\mathbf{C}}_{\mathbf{y}} \mathbf{P}_J\|_2 \leq \sqrt{8|J|p \log n}, \quad (\text{E.27})$$

then the difference of Hessian (E.10), in direction $\mathbf{v} \in \mathbb{S}^{p-1}$ can be bounded as

$$\begin{aligned} & \left| \mathbf{v}^* \left(\widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) - \text{Hess}[\varphi_\rho](\mathbf{a}) \right) \mathbf{v} \right| \\ & \leq \left| \mathbf{v}^* \iota^* \check{\mathbf{C}}_{\mathbf{y}} \left(\mathbf{P}_{I(\mathbf{a})} - \text{diag} \left[\nabla S_\lambda^\delta \left[\check{\mathbf{C}}_{\mathbf{y}} \iota \mathbf{a} \right] \right] \right) \check{\mathbf{C}}_{\mathbf{y}} \iota \mathbf{v} \right| + \|\nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a})\|_2 \end{aligned} \quad (\text{E.28})$$

where $I(\mathbf{a})$ is defined in (D.1). Let $\mathbf{D} = \mathbf{P}_{I(\mathbf{a})} - \text{diag} \left[\nabla S_\lambda^\delta \left[\check{\mathbf{C}}_{\mathbf{y}} \iota \mathbf{a} \right] \right]$ and notice that \mathbf{D} is a diagonal matrix, which suggests (E.28) can be decomposed using

$$(\mathbf{P}_J + \mathbf{P}_{J^c})\mathbf{D}(\mathbf{P}_J + \mathbf{P}_{J^c}) = \mathbf{P}_J \mathbf{D} \mathbf{P}_J + \mathbf{P}_{J^c} \mathbf{D} \mathbf{P}_{J^c},$$

where, from with property of $\sqrt{\delta}$ -smoothed ℓ^1 function Lemma E.4:

$$\max_j |\mathbf{P}_J \mathbf{D} \mathbf{P}_J|_{jj} \leq 1, \quad \max_j |\mathbf{P}_{J^c} \mathbf{D} \mathbf{P}_{J^c}|_{jj} \leq \sqrt{\lambda\delta}/B.$$

Finally, once again apply δ bound from (E.20) and bounds for $B, |J|, \mathbf{y}$ from (E.26)-(E.27), we gain

$$\begin{aligned} (\text{E.28}) & \leq \left\| \iota^* \check{\mathbf{C}}_{\mathbf{y}} \mathbf{P}_J \right\|_2^2 + \frac{\sqrt{\lambda\delta}}{B} \left\| \iota^* \check{\mathbf{C}}_{\mathbf{y}} \right\|_2^2 + \|\nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a})\|_2 \\ & \leq 8|J|p \log n + \frac{3n\sqrt{\lambda\delta}}{B} + c'n\theta^2 \\ & \leq 8 \cdot \frac{24c'n\theta^2}{p \log n} \cdot p \log n + \frac{3n(c'^4\lambda^2\theta^8/p^2 \log^2 n)^{1/2}}{c'\lambda\theta^2/p \log p} + c'n\theta^2 \\ & \leq 200c'n\theta^2, \end{aligned}$$

where all above result holds with probability at least $1 - 10/n$ from Lemma E.5 and Lemma A.5. ■

F Analysis of geometry

In this section we prove major geometrical result in [Theorem 4.1](#). This lemma consists of three parts of geometry of φ_ρ , including the negative curvature region [Corollary F.2](#), large gradient region [Corollary F.4](#), strong convexity region near shift [Corollary F.6](#), and retraction to subspace [Corollary F.8](#), which are respectively base on geometric properties of φ_{ℓ^1} in [Lemma F.1](#), [Lemma F.3](#), [Lemma F.5](#) and [Lemma F.7](#). We will handle each individual region in the following subsections. To shed light on the technical detail of the proof, we will begin with two figures for illustration of a toy example, which demonstrate the geometry near a two dimension solution subspace $\mathcal{S}_{\{i,j\}}$, as follows:

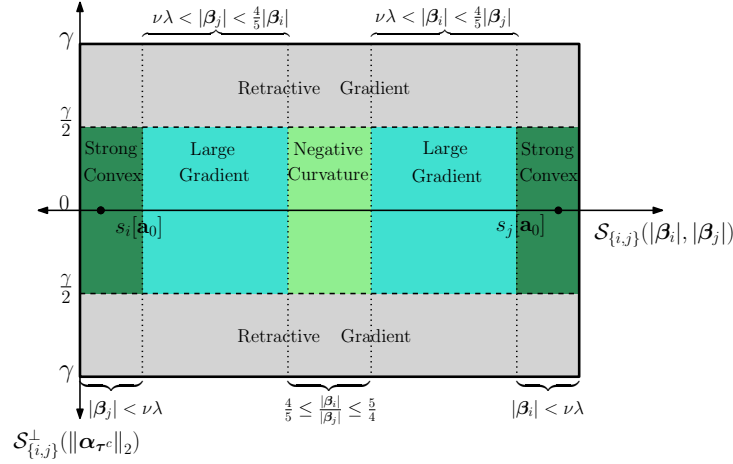


Figure 15: The top view of geometry over subspace $\mathcal{S}_{\{i,j\}}$. We display the geometric properties in the neighborhood of subspace $\mathcal{S}_{\{i,j\}}$ (horizontal axis) which contains the solutions $s_i[\mathbf{a}_0]$ and $s_j[\mathbf{a}_0]$. When \mathbf{a} lies near middle of two shifts (light green region) such that $|\beta_i| \approx |\beta_j|$, then there exists a negative curvature direction in subspace $\mathcal{S}_{\{i,j\}}$. When \mathbf{a} leans closer to one of the shifts $s_i[\mathbf{a}_0]$ (blue green region), its negative gradient direction points at that nearest shift. When \mathbf{a} is in the neighborhood of the shift $s_i[\mathbf{a}_0]$ (dark green region) such that $|\beta_i| \ll \lambda$, it will be strongly convex at \mathbf{a} , and the unique minimizer within the convex region will be close to $s_i[\mathbf{a}_0]$. Finally, the negative gradient will be pointing back toward the subspace $\mathcal{S}_{\{i,j\}}$ if near boundary (grey region).

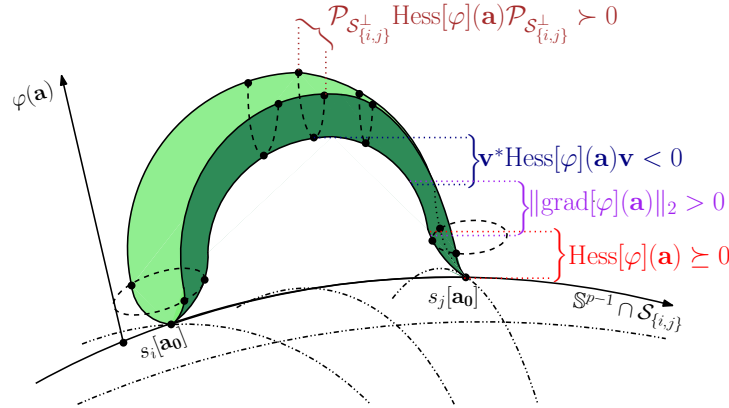


Figure 16: The side view of geometry of subspace $\mathcal{S}_{\{i,j\}}$ on sphere. We illustrate the geometry of $\mathcal{S}_{\{i,j\}}$ over the sphere, in which the properties of the three regions are denoted. In negative curvature region, there exists a direction \mathbf{v} such that $\mathbf{v}^* \text{Hess}[\varphi](\mathbf{a}) \mathbf{v}$ is negative. In large gradient region, the norm of Riemannian gradient $\|\text{grad}[\varphi](\mathbf{a})\|_2$ will be strictly greater than 0 and pointing at the nearest shift. Finally there is a convex region near all shifts such that $\text{Hess}[\varphi](\mathbf{a})$ is positive semidefinite.

F.1 Negative curvature

For any $\mathbf{a} \in \mathbb{S}^{p-1}$ near the subspace \mathcal{S}_τ such that the entries of leading correlation vector $\beta_{(0)}, \beta_{(1)}$ have balanced magnitude, the Hessian of $\varphi_\rho(\mathbf{a})$ exhibits negative curvature in the span of $s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0]$. We will first demonstrate the pseudo negative curvature of φ_{ℓ^1} in Lemma F.1, then show φ_ρ approximates φ_{ℓ^1} in terms of Hessian in Corollary F.2 when ρ is properly defined as in Appendix E.

Lemma F.1 (Negative curvature for φ_{ℓ^1}). *Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Set $\lambda = c_\lambda/\sqrt{k}$ in φ_{ℓ^1} with $c_\lambda \in [\frac{1}{5}, \frac{1}{4}]$. There exist numerical constants $C, c, c', \bar{c} > 0$ such that if $n > Cp^5\theta^{-2} \log p$, and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$ the following holds at every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $|\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}|$: for $\mathbf{v} \in \mathcal{S}_{\{(0), (1)\}} \cap \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$,*

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} \leq -cn\theta\lambda. \quad (\text{F.1})$$

Proof. First of all the regional condition $\left| \frac{\beta_{(0)}}{\beta_{(1)}} \right| \leq \frac{5}{4}$ provides a two side bound for the two leading β 's

$$0.79 \geq \frac{|\beta_{(0)}|}{\sqrt{\beta_{(0)}^2 + \beta_{(1)}^2}} \|\beta_\tau\|_2 \geq |\beta_{(0)}| \geq |\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}| \geq \frac{4}{5} \cdot \frac{\|\beta_\tau\|_2}{\sqrt{|\tau|}} \geq \frac{0.79}{\sqrt{|\tau|}} \quad (\text{F.2})$$

Set $J = \{(0), (1)\}$, choose $\mathbf{v} = \iota^* \mathbf{C}_{\mathbf{a}_0} \iota_J \gamma$ with $\|\mathbf{v}\|_2 = 1$ then $\left| \|\gamma\|_2^2 - 1 \right| \leq \mu$. There exists such \mathbf{v} satisfies condition above with $\mathbf{a} \perp \mathbf{v}$ by choosing γ as

$$\mathbf{a}^* \mathbf{v} = \mathbf{a}^* \iota^* \mathbf{C}_{\mathbf{a}_0} \iota_J \gamma = \gamma_{(0)} \beta_{(0)} + \gamma_{(1)} \beta_{(1)} = 0,$$

hence $\left| \frac{\gamma_{(1)}}{\gamma_{(0)}} \right| = \left| \frac{\beta_{(0)}}{\beta_{(1)}} \right| \leq \frac{5}{4}$. This implies $\gamma_{(0)}^2 \geq \frac{16}{25} \gamma_{(1)}^2 \geq \frac{16}{25} (1 - \mu - \gamma_{(0)}^2)$ where $\mu \leq \frac{c_\mu}{4} \leq \frac{1}{100}$, it gives the lower bound of $\gamma_{(0)}$ as

$$\gamma_{(0)}^2 \geq \frac{(1 - \mu) \cdot 16}{25 + 16} \geq 0.385 \quad (\text{F.3})$$

1. (Expand the Hessian) The (pseudo) curvature along direction \mathbf{v} is written as

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} = \mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} - \langle \nabla \varphi_{\ell^1}(\mathbf{a}), \mathbf{a} \rangle = -\gamma^* \iota_J^* \mathbf{M} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{M} \iota_J \gamma + \beta^* \chi[\beta] \quad (\text{F.4})$$

expand the first term of (F.4) we obtain

$$\begin{aligned} & -\gamma^* \iota_J^* \mathbf{M} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{M} \iota_J \gamma \\ &= -\gamma^* \iota_J^* \mathbf{M} (\mathbf{P}_{(0)} + \mathbf{P}_{(1)} + \mathbf{P}_{J^c}) \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} (\mathbf{P}_{(0)} + \mathbf{P}_{(1)} + \mathbf{P}_{J^c}) \mathbf{M} \iota_J \gamma \\ &\leq -\sum_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_i \right\|_2^2 (e_i^* \mathbf{M} \iota_J \gamma)^2 + 2 \sum_{\substack{(i,j) \in \{J, J^c\} \\ (i,j) = ((0), (1))}} \left| e_i^* \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| |(e_i^* \mathbf{M} \iota_J \gamma) (e_j^* \mathbf{M} \iota_J \gamma)| \\ &\leq -\sum_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_i \right\|_2^2 (|\gamma_i| - \mu)^2 \\ &\quad + 2 \max_{i \neq j \in [\pm p]} \left| e_i^* \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| (\|\iota_J^* \mathbf{M} \iota_J \gamma\|_1 \|\iota_{J^c}^* \mathbf{M} \iota_J \gamma\|_1 + (|\gamma_{(0)}| + \mu) (|\gamma_{(1)}| + \mu)) \end{aligned} \quad (\text{F.5})$$

Consider the following events

$$\begin{cases} \mathcal{E}_{\text{cross}} := \left\{ \forall \mathbf{a} \in \mathbb{S}^{p-1}, \max_{i \neq j \in [\pm p]} \left| e_i^* \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| < 4n\theta^2 \right\} \\ \mathcal{E}_{\text{ncurv}} := \left\{ \forall \mathbf{a} \in \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu)), \min_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}] \right\|_2^2 \geq n\theta (1 - \mathbb{E}_{\mathbf{s}_i}(\lambda, \mathbf{s}_i) + \mathbb{E}_{\mathbf{s}_i}(\lambda, \mathbf{s}_i)) - \frac{c_\mu n\theta}{p} \right\} \end{cases}, \quad (\text{F.6})$$

and from Lemma B.4 we know

$$\|\iota_J^* \mathbf{M} \iota_J \gamma\|_1 \leq \|\gamma\|_1 + 2\mu \leq 1.5, \quad \|\iota_{J^c}^* \mathbf{M} \iota_J \gamma\|_1 \leq \mu p \|\gamma\|_1 \leq 1.5\mu p,$$

on the event $\mathcal{E}_{\text{cross}} \cap \mathcal{E}_{\text{ncurv}}$, we have

$$\begin{aligned} & -\gamma^* \iota_J^* \mathbf{M} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(a)} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{M} \iota_J \gamma \\ & \leq \underbrace{-n\theta \cdot \sum_{i \in J} (|\gamma_i| - \mu)^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i) + \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i)) + (18\mu p + 8) n\theta^2 + \frac{2c_\mu n\theta}{\sqrt{|\tau|}}}_{g_1(\beta)} \end{aligned} \quad (\text{F.7})$$

Meanwhile, for the latter term of (F.4), consider the following event $\mathcal{E}_{\bar{\chi}}$ where we write $\sigma_i = \text{sign}(\beta_i)$ as:

$$\mathcal{E}_{\bar{\chi}} := \left\{ \sigma_i \chi[\beta]_i \leq \begin{cases} n\theta \cdot |\beta_i| (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) + \frac{c_\mu n\theta}{p}, & \forall i \in \tau \\ n\theta \cdot |\beta_i| 4\theta |\tau| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau^c \end{cases} \right\}, \quad (\text{F.8})$$

and use both $\|\beta\|_1 \leq \frac{c_\mu p}{\sqrt{|\tau|}}$, $\|\beta_{\tau^c}\|_2^2 \leq \frac{c_\mu}{\theta |\tau|^2}$. On this event we have

$$\begin{aligned} \beta^* \chi[\beta] & \leq n\theta \cdot \sum_{i \in \tau} \beta_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) + 4n\theta^2 |\tau| \|\beta_{\tau^c}\|_2^2 + \frac{c_\mu n\theta}{p} \|\beta\|_1 \\ & \leq n\theta \cdot \underbrace{\sum_{i \in \tau} \beta_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i))}_{g_2(\beta)} + \frac{5c_\mu n\theta}{\sqrt{|\tau|}}. \end{aligned} \quad (\text{F.9})$$

2. (Lower bound $\mathbb{E} f_{\beta_i}$) Combine the first term from each of the (F.7) and (F.9). Use $\mu \leq c_\mu \leq \frac{1}{300}$ and (F.3) to obtain $(|\gamma_{(0)}| - \mu)^2 > 0.38$, we have

$$\begin{aligned} \frac{1}{n\theta} (g_1(\beta) + g_2(\beta)) & \leq - \sum_{i \in J} \left[(|\gamma_i| - \mu)^2 - \beta_i^2 \right] (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) \\ & \quad + \sum_{i \in \tau \setminus J} \beta_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i), \end{aligned} \quad (\text{F.10})$$

now use Taylor expansion¹¹ for f_{β_i} and apply upper bound $\mathbb{E} \mathbf{s}_i^2 \leq \theta \|\beta\|_2^2 \leq \theta \left(1 + \frac{c_\mu}{\sqrt{|\tau|}} + \frac{c_\mu}{\theta |\tau|^2} \right) \leq \frac{3c_\mu}{|\tau|}$,

$$\mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \geq \mathbb{E}_{\mathbf{s}_i} \frac{1}{\sqrt{2\pi}} \cdot \left(\frac{2\lambda}{|\beta_i|} - \frac{\lambda^3}{|\beta_i|^3} \left(1 + \frac{3\mathbf{s}_i^2}{\lambda^2} \right) \right) \geq \frac{1}{\sqrt{2\pi}} \cdot \underbrace{\left(\frac{2\lambda}{|\beta_i|} - \frac{1}{|\beta_i|^3} \left(\lambda^3 + \frac{9c_\mu \lambda}{|\tau|} \right) \right)}_{f(\beta)},$$

where $f(\beta)$ is concave at stationary point since

$$\begin{cases} f'(\beta_*) = 0 \implies 2\lambda\beta_*^2 = 3\lambda \left(\lambda^2 + \frac{9c_\mu}{|\tau|} \right) \\ f''(\beta_*) = \frac{1}{|\beta_*|^3} \left(4\lambda - \frac{12\lambda}{\beta_*^2} \left(\lambda^2 + \frac{9c_\mu}{|\tau|} \right) \right) = \frac{1}{|\beta_*|^3} \left(4\lambda - \frac{12}{3/2} \lambda \right) < 0 \end{cases},$$

then combine with regional condition (F.2), and also apply assumption $c_\lambda \leq \frac{1}{3}$ and $c_\mu \leq \frac{1}{300}$, we gain

$$0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \geq 0.3 \min_{\beta = \frac{0.79}{\sqrt{|\tau|}}, 0.79} f(\beta)$$

¹¹ Apply $\exp[-x^2/2] > 1 - x^2/2$

$$\begin{aligned}
&\geq 0.3 \min \left\{ \frac{2c_\lambda}{0.79} - \frac{c_\lambda^3 + 9c_\mu c_\lambda}{0.79^3}, \lambda \left(\frac{2}{0.79} - \frac{c_\lambda^2 + 9c_\mu}{0.79^3} \right) \right\} \\
&\geq 0.3 \min \{2c_\lambda, 2\lambda\} \geq 0.6\lambda.
\end{aligned} \tag{F.11}$$

3. (Upper bound $\mathbb{E}\chi[\beta]_i$) When $\beta_{(0)}^2 = (|\gamma_{(0)}| - \mu)^2 - \eta$ for some $\eta > 0$. With monotonicity [Lemma C.3](#), which implies:

$$(1 - \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)})) \geq (1 - \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)})) \geq (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)), \tag{F.12}$$

then combine (F.11)-(F.12) and use $\mu \leq \frac{c_\mu}{4\sqrt{|\tau|}}$ from [Lemma B.5](#)

$$\begin{aligned}
\text{(F.10)} &\leq - \underbrace{\left((|\gamma_{(0)}|^2 - \mu)^2 - \beta_{(0)}^2 - \eta \right)}_{=0} (1 - \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)})) \\
&\quad + \left(\sum_{i \in \tau \setminus (0)} \beta_i^2 - (|\gamma_{(1)}| - \mu)^2 - \eta \right) \underbrace{(1 - \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}))}_{<1} - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \\
&\leq \left(\|\beta_\tau\|_2^2 - \|\gamma\|_2^2 + 2\mu \|\gamma\|_1 \right) - 0.6\lambda \\
&\leq \frac{2c_\mu}{\sqrt{|\tau|}} - 0.6\lambda.
\end{aligned} \tag{F.13}$$

On the other hand, when $\beta_{(0)}^2 \geq (|\gamma_{(0)}| - \mu)^2 > 0.38$, combining (F.11)-(F.12) gives:

$$\begin{aligned}
\text{(F.10)} &\leq \left(\|\beta_\tau\|_2^2 - \|\gamma\|_2^2 + 2\mu \|\gamma\|_1 \right) + \left((|\gamma_{(0)}| - \mu)^2 - \beta_{(0)}^2 \right) \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)}) \\
&\quad + \left((|\gamma_{(1)}| - \mu)^2 - \sum_{i \in \tau \setminus (0)} \beta_i^2 \right) \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \\
&\leq \left(\frac{c_\mu}{\sqrt{|\tau|}} + 4\mu \right) + \left(\gamma_{(1)}^2 - \|\beta_\tau\|_2^2 + \beta_{(0)}^2 \right) \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) - 0.6\lambda,
\end{aligned} \tag{F.14}$$

where [Lemma C.2](#) provides the upper bound for $\mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)})$ as

$$\begin{aligned}
\mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) &= 1 - \frac{1}{n\theta\beta_{(1)}} \mathbb{E}\chi[\beta]_{(1)} \leq 1 - \frac{\sigma_{(1)}}{n\theta|\beta_{(1)}|} \mathbb{E}\chi[\beta]_{(1)} = 1 - \frac{1}{|\beta_{(1)}|} \left(|\beta_{(1)}| - \sqrt{\frac{2}{\pi}} \lambda \right) \\
&\leq \sqrt{\frac{2}{\pi}} \cdot \frac{\lambda}{|\beta_{(1)}|},
\end{aligned} \tag{F.15}$$

then calculate the constant for the second term in (F.14) by writing $\kappa = \left| \frac{\gamma_{(1)}}{\gamma_{(0)}} \right| = \left| \frac{\beta_{(0)}}{\beta_{(1)}} \right| \leq \frac{5}{4}$, which provides

$\gamma_{(1)}^2 \leq \frac{(1+\mu)\kappa^2}{\kappa^2+1}$ and $\beta_{(0)}^2 \leq \frac{\|\beta_\tau\|_2^2 \kappa^2}{\kappa^2+1}$ where $\mu < \frac{c_\mu}{4}$, and by applying $|\beta_{(1)}| > \frac{4}{5} |\beta_{(0)}| \geq 0.3$, we have

$$\frac{(\gamma_{(1)}^2 - 1) + c_\mu + \beta_{(0)}^2}{|\beta_{(1)}|} \leq -\frac{\kappa}{(\kappa^2 + 1)|\beta_{(0)}|} + \kappa |\beta_{(0)}| + \frac{\mu + c_\mu}{0.3} \leq \frac{\kappa^2 - 1}{\sqrt{\kappa^2 + 1}} + \kappa \left(\|\beta_\tau\|_2^2 - 1 \right) + 4.2c_\mu \leq 0.36 + 6c_\mu, \tag{F.16}$$

and finally combine (F.15)-(F.16), follow from (F.14) and use $c_\lambda \leq \frac{1}{3}$:

$$\text{(F.10)} \leq \frac{2c_\mu}{\sqrt{|\tau|}} + \sqrt{\frac{2}{\pi}} \left(\gamma_{(1)}^2 - 1 + c_\mu + \beta_{(0)}^2 \right) \frac{\lambda}{|\beta_{(1)}|} - 0.6\lambda$$

$$\begin{aligned}
&\leq \frac{2c_\mu}{\sqrt{|\boldsymbol{\tau}|}} + \sqrt{\frac{2}{\pi}} \left(0.36\lambda + \frac{6c_\mu c_\lambda}{0.3} \right) - 0.6\lambda \\
&\leq \frac{4c_\mu}{\sqrt{|\boldsymbol{\tau}|}} - 0.3\lambda
\end{aligned} \tag{F.17}$$

3. (Collect all results) Combine the components of pseudo Hessian (F.7), (F.9) with bounds for $g_1 + g_2$ from (F.13) and (F.17), and use Lemma B.5 which provides both $\mu p \theta |\boldsymbol{\tau}| < \frac{c_\mu}{4}$ and $\theta |\boldsymbol{\tau}| < \frac{c_\mu}{4}$ where $c_\mu < \frac{1}{300}$ and $c_\lambda \geq \frac{1}{5}$, we can obtain:

$$\begin{aligned}
\mathbf{v}^* \widetilde{\text{Hess}}_{\varphi_{\ell^1}}[\mathbf{a}] \mathbf{v} &\leq g_1(\boldsymbol{\beta}) + g_2(\boldsymbol{\beta}) + \frac{7c_\mu n \theta}{\sqrt{|\boldsymbol{\tau}|}} + (18\mu p + 8) n \theta^2 \\
&\leq n \theta \cdot \left(\frac{4c_\mu}{\sqrt{|\boldsymbol{\tau}|}} - 0.3\lambda \right) + n \theta \cdot \frac{7c_\mu}{\sqrt{|\boldsymbol{\tau}|}} + n \theta \cdot \frac{6.5c_\mu}{|\boldsymbol{\tau}|} \\
&\leq \frac{n \theta}{\sqrt{|\boldsymbol{\tau}|}} (0.059 - 0.06) \leq -0.001 n \theta \lambda
\end{aligned} \tag{F.18}$$

Finally, the curvature is negative along \mathbf{v} direction with probability at least

$$1 - \underbrace{\mathbb{P}[\mathcal{E}_{\text{cross}}^c]}_{\text{Lemma A.4}} - \underbrace{\mathbb{P}[\mathcal{E}_{\text{ncurv}}^c]}_{\text{Corollary D.3}} - \underbrace{\mathbb{P}[\mathcal{E}_{\bar{\chi}}^c]}_{\text{Corollary C.4}}. \tag{F.19}$$

■

Similarly for objective φ_ρ , we have that

Corollary F.2 (Negative curvature for φ_ρ). *Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_ρ where $c_\lambda \in [\frac{1}{5}, \frac{1}{4}]$, then there exists some numerical constants $C, c, c', c'', \bar{c} > 0$ such that if ρ is δ -smoothed ℓ^1 function where $\delta \leq c'' \lambda \theta^8 / p^2 \log^2 n$, $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $\mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $|\boldsymbol{\beta}_{(1)}| \geq \frac{4}{5} |\boldsymbol{\beta}_{(0)}|$: for $\mathbf{v} \in \mathcal{S}_{\{(0), (1)\}} \cap \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$,*

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_\rho](\mathbf{a}) \mathbf{v} \leq -cn\theta\lambda \tag{F.20}$$

Proof. Choose $\mathbf{v} \in \mathbb{S}^{p-1}$ according to Lemma F.1 and (E.23) from Lemma E.6 with constant multiplier δ satisfies $c''^{1/4} < 10^{-3}c$, we gain

$$\mathbf{v}^* \text{Hess}[\varphi_\rho](\mathbf{a}) \mathbf{v} \leq -cn\theta\lambda + 200c'n\theta^2 \leq -cn\theta\lambda/2 \tag{F.21}$$

■

F.2 Large gradient

For any $\mathbf{a} \in \mathbb{S}^{p-1}$ near subspace and the second largest correlation $\boldsymbol{\beta}_{(1)}$ much smaller than the first correlation $\boldsymbol{\beta}_{(0)}$ while not being near 0, the negative gradient of $\varphi_\rho(\mathbf{a})$ will point at the largest shift. We show this in Lemma F.3, and the φ_ρ version in Corollary F.4 when ρ is properly defined as in Appendix E.

Lemma F.3 (Large gradient for φ_{ℓ^1}). *Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_{ℓ^1} with some $c_\lambda \in [\frac{1}{5}, \frac{1}{4}]$, then there exists some numerical constants $C, c', c, \bar{c} > 0$, such that if $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $\mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $\frac{4}{5} |\boldsymbol{\beta}_{(0)}| > |\boldsymbol{\beta}_{(1)}| > \frac{1}{4 \log \theta^{-1}} \lambda$,*

$$\langle \boldsymbol{\sigma}_{(0)} \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_{\ell^1}](\mathbf{a}) \rangle \geq cn\theta (\log^{-2} \theta^{-1}) \lambda^2 \tag{F.22}$$

where $\boldsymbol{\sigma}_i = \text{sign}(\boldsymbol{\beta}_i)$.

Proof. 1. (Properties for α, β) Define $\theta_{\log} = \frac{1}{\log \theta^{-1}}$, we first derive upper bound on the dominant entry $|\beta_{(0)}|$ as follows. Write the geodesic distance between \mathbf{a} and $\iota^* s_i[\mathbf{a}_0]$ as a function of β_i as $d_{\mathbb{S}}(\mathbf{a}, \pm \iota^* s_i[\mathbf{a}_0]) = \cos^{-1}(\beta_i)$, then by triangle inequality we have:

$$\begin{aligned} d_{\mathbb{S}}(\mathbf{a}, \pm \iota^* s_{(0)}[\mathbf{a}_0]) &\geq d_{\mathbb{S}}(\pm \iota^* s_{(0)}[\mathbf{a}_0], \iota^* s_{(1)}[\mathbf{a}_0]) - d_{\mathbb{S}}(\mathbf{a}, \iota^* s_{(1)}[\mathbf{a}_0]) \\ \implies \cos^{-1} \pm \beta_{(0)} &\geq \cos^{-1} \mu - \cos^{-1} |\beta_{(1)}| \\ \implies \pm \beta_{(0)} &\leq \cos(\cos^{-1} \mu - \cos^{-1} |\beta_{(1)}|) = \mu |\beta_{(1)}| + \sqrt{(1 - \mu^2)(1 - \beta_{(1)}^2)} \leq 1 - \frac{1}{2} (|\beta_{(1)}| - \mu)^2. \end{aligned}$$

Use the regional condition $|\beta_{(1)}| \geq \frac{\theta_{\log}}{4} \lambda$ and since $\mu |\tau|^{3/2} < \frac{c_{\lambda}}{100} \theta_{\log}$ from [Definition B.1](#), implies

$$|\beta_{(0)}| \leq 1 - \frac{\beta_{(1)}^2}{2} \left(1 - \frac{4\mu\sqrt{|\tau|}}{\theta_{\log} c_{\lambda}}\right) \leq 1 - 0.49 \beta_{(1)}^2 =: \beta_{\text{ub}}. \quad (\text{F.23})$$

Meanwhile a lower bound for $\beta_{(0)}$ can be easily determined by the other side of regional condition:

$$|\beta_{(0)}| \geq \frac{5}{4} |\beta_{(1)}| =: \beta_{\text{lb}}. \quad (\text{F.24})$$

Also since $\beta = M\alpha$, based on properties of M from [Lemma B.4](#). When $\|\alpha_{\tau}\|_2 \leq 1 + c_{\mu}$ and $\|\alpha_{\tau^c}\|_2 \leq \gamma \leq \frac{c_{\mu} \theta_{\log}^2}{4\mu\sqrt{p}|\tau|}$, we gain:

$$\begin{aligned} \beta_{(0)} &= \alpha_{(0)} + e_{(0)}^* M \alpha_{\setminus(0)} \\ \implies |\alpha_{(0)} - \beta_{(0)}| &\leq \mu \sqrt{|\tau|} \|\alpha_{\tau}\|_2 + \mu \sqrt{p} \|\alpha_{\tau^c}\|_2 \leq \frac{c_{\mu} \theta_{\log}^2 (1 + c_{\mu})}{4|\tau|} + \mu \sqrt{p} \gamma \leq \frac{c_{\mu} \theta_{\log}^2}{|\tau|}. \end{aligned} \quad (\text{F.25})$$

and therefore $|\alpha_{(0)}| \leq |\beta_{(0)}| + \frac{c_{\mu} \theta_{\log}^2}{|\tau|} \leq 1 - .49 \left(\frac{\theta_{\log}}{4} \lambda\right)^2 + \frac{c_{\mu} \theta_{\log}^2}{|\tau|} < 1$.

2. (Upper bound of $\beta^* \chi[\beta]$) Define a piecewise smooth convex upper bound h for $\beta_i \chi[\beta]_i$ as:

$$h(\beta_i) := \begin{cases} \beta_i^2 - \frac{\nu_1 \lambda}{2} |\beta_i| & |\beta_i| \geq \nu_1 \lambda \\ \frac{1}{2} \beta_i^2 & |\beta_i| \leq \nu_1 \lambda \end{cases},$$

then [Lemma J.7](#) tells us since $\|\beta_{\tau \setminus (0)}\|_{\infty} \leq \beta_{(1)}$:

$$\begin{aligned} \sum_{i \in \tau \setminus (0)} h(\beta_i) &\leq \|\beta_{\tau \setminus (0)}\|_2^2 \left(1 - \frac{\nu_1 \lambda \beta_{(1)}}{2 \beta_{(1)}^2}\right) \leq \left(1 + \frac{c_{\mu} \theta_{\log}^2}{|\tau|} - \beta_{(0)}^2\right) \left(1 - \frac{\nu_1 \lambda}{2 \beta_{(1)}}\right) \\ &\leq \left(1 - \frac{\nu_1 \lambda}{2 \beta_{(1)}}\right) (1 - \beta_{(0)}^2) + \frac{c_{\mu} \theta_{\log}^2}{|\tau|}, \end{aligned}$$

then condition on the following event using [Corollary C.4](#),

$$\mathcal{E}_{\bar{\chi}} := \left\{ \beta_i \chi[\beta]_i \leq \begin{cases} n\theta \cdot h(\beta_i) + \frac{c_{\mu} \theta}{p^{3/2}} |\beta_i|, & \forall i \in \tau \setminus (0) \\ n\theta \cdot 4\beta_i^2 \theta |\tau| + \frac{c_{\mu} \theta}{p^{3/2}} |\beta_i|, & \forall i \in \tau^c \end{cases} \right\},$$

which provides the upper bound of $\beta^* \chi[\beta]$ by applying $5p > \log^{8/3}(p \log^2 p) > (\theta_{\log}^2)^{4/3}$ from lower bound of θ from [Definition B.1](#), $\|\beta_{\tau^c}\|_2 \leq \frac{c_{\mu} \theta_{\log}}{\sqrt{\theta} |\tau|}$ from [Lemma B.5](#), $|\tau| \leq \sqrt{p}$ from lemma assumption and let $c_{\mu} < \frac{1}{100}$:

$$\beta^* \chi[\beta] \leq \chi[\beta]_{(0)} \beta_{(0)} + \sum_{i \in \tau \setminus (0)} \beta_i \chi[\beta]_i + \langle \beta_{\tau^c}, \chi[\beta]_{\tau^c} \rangle$$

$$\begin{aligned}
&\leq \chi[\beta]_{(0)}\beta_{(0)} + n \left(\theta \sum_{i \in \tau \setminus (0)} h(\beta_i) + 4\theta^2 |\tau| \|\beta_{\tau^c}\|_2^2 + \frac{c_\mu \theta}{p^{3/2}} \left(\sqrt{|\tau|} \|\beta_\tau\|_2 + \sqrt{p} \|\beta_{\tau^c}\|_2 \right) \right) \\
&\leq \chi[\beta]_{(0)}\beta_{(0)} + n \left(\theta \cdot \eta(1 - \beta_{(0)}^2) + \theta \cdot \frac{c_\mu \theta_{\log}^2}{|\tau|} + \frac{4\theta^2 |\tau| c_\mu^2 \theta_{\log}^2}{\theta |\tau|^2} + c_\mu \theta \left(\frac{1 + c_\mu}{p^{3/4} |\tau|} + \frac{c_\mu \theta_{\log}}{p \sqrt{\theta} |\tau|} \right) \right) \\
&\leq \chi[\beta]_{(0)}\beta_{(0)} + n\theta \left(\eta(1 - \beta_{(0)}^2) + \frac{6c_\mu \theta_{\log}^2}{|\tau|} \right), \tag{F.26}
\end{aligned}$$

where $\eta = 1 - \frac{\nu_1 \lambda}{2\beta_{(1)}}$.

3. (Align the gradient with $\iota^* s_{(0)}[a_0]$) Base on the definition β , since $\beta_{(0)} = \langle a, \iota^* s_{(0)}[a_0] \rangle$, we can expect that the negative gradient is likely aligned with direction toward one of the candidate solution $\pm \iota^* s_{(0)}[a_0]$. Wlog assume that both $\beta_{(0)}, \beta_{(1)}$ are positive, then expand the gradient and use incoherent property for a_0 [Lemma B.4](#) we have:

$$\begin{aligned}
\langle \iota^* s_{(0)}[a_0], -\text{grad}_{\varphi_{t_1}}[a] \rangle &= \langle \iota^* s_{(0)}[a_0], \iota^* C_{a_0} (\chi[\beta] - \beta^* \chi[\beta] \alpha) \rangle \\
&\geq (\chi[\beta]_{(0)} - \beta^* \chi[\beta] \alpha_{(0)}) - \mu \|\chi[\beta]_{\setminus (0)} - \beta^* \chi[\beta] \alpha_{\setminus (0)}\|_1, \tag{F.27}
\end{aligned}$$

where $\setminus (0)$ is an abbreviation of the complement set $[\pm 2p_0] \setminus (0)$. The latter part of (F.27) has an upper bound using bounds of $\beta^* \chi[\beta] < \frac{3n\theta}{2}$, $\|\chi[\beta]_{\tau^c}\|_2 < \frac{n\theta\gamma_2}{20}$ from (F.62), and $\|\chi[\beta]_{\tau \setminus (0)}\|_2 \leq n\theta \|\beta_{\tau \setminus (0)}\|_2$ in event $\mathcal{E}_{\bar{\chi}}$, we obtain:

$$\begin{aligned}
&\mu \|\chi[\beta]_{\setminus (0)} - \beta^* \chi[\beta] \alpha_{\setminus (0)}\|_1 \\
&\leq \mu \left(\sqrt{|\tau|} \|\chi[\beta]_{\tau \setminus (0)}\|_2 + \beta^* \chi[\beta] \sqrt{|\tau|} \|\alpha_{\tau \setminus (0)}\|_2 + \sqrt{p} \|\chi[\beta]_{\tau^c}\|_2 + \beta^* \chi[\beta] \sqrt{p} \|\alpha_{\tau^c}\|_2 \right) \\
&\leq n\theta \cdot \left[\mu \sqrt{|\tau|} (\|\beta_\tau\|_2 - |\beta_{(0)}|) + \mu \sqrt{|\tau|} (\|\alpha_\tau\|_2 - |\alpha_{(0)}|) + \frac{1}{20} \mu \sqrt{p} \gamma_2 + \frac{3}{2} \mu \sqrt{p} \gamma_2 \right] \\
&\leq n\theta \cdot \frac{c_\mu \theta_{\log}^2}{4 |\tau|} \left[2(1 + c_\mu) - |\beta_{(0)}| - |\alpha_{(0)}| + \left(\frac{1}{20} + \frac{3}{2} \right) c_\mu \right] \\
&\leq n\theta \cdot \frac{c_\mu \theta_{\log}^2}{|\tau|} (0.5 + c_\mu - 0.5\beta_{(0)}). \tag{F.28}
\end{aligned}$$

On the other hand, the former term of (F.27) possesses a lower bound using (F.25)-(F.26), $\chi[\beta]_{(0)} > n\theta \left(\beta_{(0)} - \frac{\nu_1}{2} \lambda - \frac{c_\mu}{p} \right) \geq n\theta (\beta_{(0)} - 0.51\nu_1 \lambda)$ and $\alpha_{(0)} \leq 1$:

$$\begin{aligned}
&\chi[\beta]_{(0)} - \beta^* \chi[\beta] \alpha_{(0)} \\
&\geq (1 - \alpha_{(0)}\beta_{(0)}) \chi[\beta]_{(0)} - n\theta \cdot \left[\eta(1 - \beta_{(0)}^2) + \frac{6c_\mu \theta_{\log}^2}{|\tau|} \right] \alpha_{(0)} \\
&\geq n\theta \underbrace{\left(1 - \left(\beta_{(0)} + \frac{c_\mu \theta_{\log}^2}{|\tau|} \right) \beta_{(0)} \right)}_{(a)} (\beta_{(0)} - 0.51\nu_1 \lambda) - n\theta \underbrace{\left[\eta(1 - \beta_{(0)}^2) \left(\beta_{(0)} + \frac{c_\mu \theta_{\log}^2}{|\tau|} \right) + \frac{6c_\mu \theta_{\log}^2}{|\tau|} \alpha_{(0)} \right]}_{(b)} \\
&\geq n\theta \left[\underbrace{\left(1 - \beta_{(0)}^2 \right) (\beta_{(0)} - 0.51\nu_1 \lambda) - \frac{c_\mu \theta_{\log}^2 \beta_{(0)}^2}{|\tau|}}_{(a)} - \underbrace{\left(1 - \beta_{(0)}^2 \right) \eta \beta_{(0)} - \eta \frac{c_\mu \theta_{\log} (1 - \beta_{(0)}^2)}{|\tau|} - \frac{6c_\mu \theta_{\log}^2}{|\tau|}}_{(b)} \right]
\end{aligned}$$

$$\geq n\theta \left[\left(1 - \beta_{(0)}^2\right) \left((1 - \eta) \beta_{(0)} - 0.51\nu_1\lambda\right) - \frac{c_\mu\theta_{\log}^2}{|\tau|} \left((1 - \eta)\beta_{(0)}^2 + 7\right) \right], \quad (\text{F.29})$$

combine (F.27) with (F.28)-(F.29) and $\eta > 0$, we have

$$\begin{aligned} (\text{F.27}) &\geq n\theta \left[\left(1 - \beta_{(0)}^2\right) \left((1 - \eta) \beta_{(0)} - 0.51\nu_1\lambda\right) - \frac{c_\mu\theta_{\log}^2}{|\tau|} \left((1 - \eta)\beta_{(0)}^2 + 7\right) \right] - n\theta \cdot \frac{c_\mu\theta_{\log}^2}{|\tau|} (0.5 + c_\mu - 0.5\beta_{(0)}) \\ &\geq n\theta \left[\underbrace{\left(1 - \beta_{(0)}^2\right) \left(\frac{\nu_1\lambda}{2\beta_{(1)}}\beta_{(0)} - 0.51\nu_1\lambda\right)}_{f(\beta)} - \frac{8c_\mu\theta_{\log}^2}{|\tau|} \right]. \end{aligned} \quad (\text{F.30})$$

4. (Lower bound of $f(\beta)$) Given a fixed $\beta_{(1)}$, the cubic function $f(\beta_{(0)})$ has zeros set $\beta_{(0)} \in \{\pm 1, 1.02\beta_{(1)}\}$ and has negative leading coefficient. Combine with the condition of $\beta_{(0)} \in \{\beta_{\text{lb}}, \beta_{\text{ub}}\}$ from (F.23)-(F.24), we can observe that

$$\beta_{(0)} \in [\beta_{\text{lb}}, \beta_{\text{ub}}] = \left[\frac{5}{4}\beta_{(1)}, 1 - 0.49\beta_{(1)}^2\right] \subseteq [1.02\beta_{(1)}, 1],$$

therefore the cubic term is always positive and minimizer is either one of the boundary point. When $\beta_{(0)} = \beta_{\text{lb}}$, use $(1 + \frac{25}{16})\beta_{(1)}^2 < 1.01$, and use $\nu_1\lambda < \frac{\sqrt{\theta_{\log}}}{2\sqrt{|\tau|}} \leq \frac{1}{2\sqrt{2}}$, since $|\tau| \geq 2$, we have:

$$f(\beta_{\text{lb}}) \geq (1 - \beta_{\text{lb}}^2) \left(\frac{\nu_1\lambda}{2\beta_{(1)}}\beta_{\text{lb}} - 0.51\nu_1\lambda\right) \geq (1 - 0.616) \cdot \left(\frac{5}{8} - 0.51\right) \nu_1\lambda \geq \frac{1}{16\sqrt{2}}\nu_1\lambda \geq \frac{\theta_{\log}^2}{32}\lambda^2. \quad (\text{F.31})$$

On the other hand when $\beta_{(0)} = \beta_{\text{ub}}$:

$$f(\beta_{\text{ub}}) \geq (1 - \beta_{\text{ub}}^2) \left(\frac{\nu_1\lambda}{2\beta_{(1)}}\beta_{\text{ub}} - 0.51\nu_1\lambda\right) \geq 0.49\beta_{(1)}^2 \cdot \left(\frac{\nu_1\lambda}{2\beta_{(1)}}(1 - 0.49\beta_{(1)}^2) - 0.51\nu_1\lambda\right),$$

which is a cubic function of $\beta_{(1)}$ with negative leading coefficient, whose zeros set is $\{-0.73, 0, 2.81\}$. Thus it minimizes at the boundary points of $\beta_{(1)} \in \left[\frac{\lambda}{4\log\theta^{-1}}, 1\right] \subset [0, 2.81]$, thus assign $\beta_{(1)} = \frac{\lambda}{4\log\theta^{-1}}$, we have:

$$f(\beta_{\text{ub}}) \geq 0.49 \left(\frac{\lambda}{4\log\theta^{-1}}\right)^2 \cdot \left(\frac{1}{2} \left(1 - 0.49 \left(\frac{\lambda}{4\log\theta^{-1}}\right)^2\right) - 0.51\nu_1\lambda\right) \geq \frac{1}{6} \left(\frac{\lambda}{4\log\theta^{-1}}\right)^2 \geq \frac{\theta_{\log}^2}{96}\lambda^2. \quad (\text{F.32})$$

Finally combine (F.30) with the lower bound of cubic function (F.31)-(F.32) together with condition $c_\mu < \frac{c_\lambda^2}{800}$ and $\nu_1 = \frac{\sqrt{\theta_{\log}}}{2}$, obtain

$$\begin{aligned} \left\langle \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}_{\varphi_{\ell_1}}[\mathbf{a}] \right\rangle &\geq n\theta \cdot \left(\min\{f(\beta_{\text{ub}}), f(\beta_{\text{lb}})\} - \frac{8c_\mu\theta_{\log}^2}{|\tau|} \right) \\ &\geq n\theta \left(\frac{\theta_{\log}^2 c_\lambda^2}{96|\tau|} - \frac{8\theta_{\log}^2 c_\lambda^2}{800|\tau|} \right) \geq 6 \times 10^{-3} n\theta_{\log}^2 c_\lambda^2. \end{aligned} \quad (\text{F.33})$$

The proof for the case where $\beta_{(0)}$ negative can be derived in the same manner. ■

As a consequence, we have that

Corollary F.4 (Large gradient for φ_ρ). Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_ρ with $c_\lambda \in [\frac{1}{5}, \frac{1}{4}]$, then there exists some numerical constants $C, c, c', c'', \bar{c} > 0$ such that if ρ is δ -smoothed ℓ^1 function where $\delta \leq c'' \lambda \theta^8 / p^2 \log^2 n$ with $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $\frac{4}{5} |\beta_{(0)}| > |\beta_{(1)}| > \frac{1}{4 \log \theta^{-1} \lambda}$,

$$\langle \sigma_{(0)} \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq cn\theta (\log^{-2} \theta^{-1}) \lambda^2 \quad (\text{F.34})$$

where $\sigma_i = \text{sign}(\beta_i)$.

Proof. Choose $\iota^* s_{(0)}[\mathbf{a}_0]$ as in Lemma F.3, and apply (E.22) from Lemma E.6 with the constant multiplier of δ satisfies $c''^4 < c/4$, then utilize $\theta |\tau| \log^2 \theta^{-1} < c_\mu$ from Definition B.1 we have

$$\langle \sigma_{(0)} \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq cn\theta (\log^{-2} \theta^{-1}) \lambda - c'' n \theta^2 \geq cn\theta (\log^{-2} \theta^{-1}) \lambda / 2 \quad (\text{F.35})$$

F.3 Convex near solutions

For any $\mathbf{a} \in \mathbb{S}^{p-1}$ near subspace and the second largest correlation $\beta_{(1)}$ smaller than $\frac{1}{4 \log \theta^{-1}} \lambda$, then φ_ρ will be strongly convex at \mathbf{a} . We show this in Lemma F.5, and the φ_ρ version in Corollary F.6 when ρ is properly defined as in Appendix E.

Lemma F.5 (Strong convexity of φ_{ℓ^1} near shift). Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_{ℓ^1} with $c_\lambda \in [\frac{1}{4}, \frac{1}{5}]$, then there exists some numerical constants $C, c, c', \bar{c} > 0$ such that if $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $|\beta_{(1)}| < \frac{1}{4 \log \theta^{-1}} \lambda$: for all $\mathbf{v} \in \mathbb{S}^{p-1} \cap \mathbf{v}^\perp$,

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} > cn\theta; \quad (\text{F.36})$$

furthermore, there exists $\bar{\mathbf{a}}$ as an local minimizer such that

$$\min_{\ell} \|\bar{\mathbf{a}} - s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max \{\mu, p^{-1}\}. \quad (\text{F.37})$$

Proof. 1. (Expectation of χ near shifts) We will write \mathbf{x} as \mathbf{x}_0 through out this proof. When \mathbf{a} is near one of the shift, the χ operator shrinks all other smaller entries of correlation vector $\beta_{\setminus(0)}$ in an even larger shrinking ratio. Firstly we can show $|\langle \beta_{\setminus(0)}, \mathbf{x}_{\setminus(0)} \rangle|$ is no larger than $\lambda/2$ with probability at least $1 - 4\theta$, since

$$\mathbb{P} \left[|\langle \beta_{\setminus(0)}, \mathbf{x}_{\setminus(0)} \rangle| > \frac{\lambda}{2} \right] \leq \mathbb{P} \left[|\langle \beta_{\tau \setminus(0)}, \mathbf{x}_{\tau \setminus(0)} \rangle| > \frac{2\lambda}{5} \right] + \mathbb{P} \left[|\langle \beta_{\tau^c}, \mathbf{x}_{\tau^c} \rangle| > \frac{\lambda}{10} \right] \leq 4\theta \quad (\text{F.38})$$

via Corollary B.6 and Corollary B.7. Now recall from Lemma C.2 and the derivation of (C.10)-(C.11), we know for every $i \neq (0)$,

$$\begin{aligned} \sigma_i \mathbb{E} \chi[\beta]_i &= n\theta |\beta_i| \mathbb{E}_{\mathbf{s}_i} [1 - \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)] \\ &\leq n\theta |\beta_i| \mathbb{E}_{g, \mathbf{x}_{\setminus i}} \left[g^2 \mathbf{1}_{\{|\beta_i g + \beta_{(0)} \mathbf{x}_{(0)} + \beta_{\setminus \{(0), i\}}^* \mathbf{x}_{\setminus \{(0), i\}}| > \lambda\}} \right] \\ &\leq n\theta |\beta_i| (\mathbb{E} g^2 \mathbf{1}_{\{|\beta_i g| > \lambda/2\}} + \mathbb{P}[\mathbf{x}_{(0)} \neq 0] + \mathbb{P}[|\langle \beta_{\setminus \{(0), i\}}, \mathbf{x}_{\setminus \{(0), i\}} \rangle| > \lambda/2]) \\ &\leq n\theta |\beta_i| \left((\mathbb{E} g^2)^{1/2} \mathbb{P}[|\beta_{(1)} g| > \lambda/2]^{1/2} + \theta + 4\theta \right) \\ &\leq n\theta |\beta_i| (\exp(-\log^2 \theta^{-1}) + 5\theta) \\ &\leq 6n\theta^2 |\beta_i| \end{aligned} \quad (\text{F.39})$$

where the third inequality is derived using union bound; the the fourth inequality is the result of (F.38), and the fifth inequality is derived from Gaussian tail bound lemma J.1.

2. (Local strong convexity) Let $\gamma = C_{a_0}^* \iota v$, for any $\|v\|_2 = 1$ we have $\|\gamma\|_2^2 \leq 1 + \mu p$. Furthermore:

$$\begin{aligned} |\gamma_{(0)}| &= |\langle \iota^* s_{(0)}[a_0], v \rangle| = |\langle P_{a^\perp} \iota^* s_{(0)}[a_0], v \rangle| = |\langle \iota^* s_{(0)}[a_0] - \beta_{(0)} a, v \rangle| \\ &\leq \|\iota^* s_{(0)}[a_0] - \beta_{(0)} a\|_2 \leq \sqrt{1 - \beta_{(0)}^2}. \end{aligned} \quad (\text{F.40})$$

Consider any such v , the pseudo Hessian can be lower bounded as

$$\begin{aligned} v^* \tilde{\nabla}^2 \varphi_{\ell^1}(a) v &= -\gamma^* \tilde{C}_x P_{I(a)} \tilde{C}_x \gamma \\ &\geq -\gamma_{(0)}^2 \left\| P_{I(a)} \tilde{C}_x e_{(0)} \right\|_2^2 - \sum_{i \neq (0)} \left\| P_{I(a)} \tilde{C}_x e_i \right\|_2^2 \gamma_i^2 - 2 \sum_{i \neq j} \left| e_i^* \tilde{C}_x P_{I(a)} \tilde{C}_x e_j \right| |\gamma_i| |\gamma_j| \\ &\geq -\left(1 - \beta_{(0)}^2\right) \|x\|_2^2 - \max_{i \neq (0)} \left\| P_{I(a)} s_{-i}[x] \right\|_2^2 \|\gamma\|_2^2 - 2 \max_{i \neq j} \left| e_i^* \tilde{C}_x P_{I(a)} \tilde{C}_x e_j \right| \|\gamma\|_1^2, \end{aligned} \quad (\text{F.41})$$

where the second term is bounded by using its expectation derived in Lemma D.2, and utilize $\mathbb{P}[|s_i| > \lambda/2] < 4\theta$ from (F.38), $\mathbb{E}\chi$ from (F.39) and regional condition $|\beta_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}$ to acquire

$$\begin{aligned} \mathbb{E} \left\| P_{I(a)} s_{-i}[x] \right\|_2^2 &= n\theta [1 - \mathbb{E}_{s_i} \text{erf}_{\beta_i}(\lambda, s_i) + \mathbb{E}_{s_i} f_{\beta_i}(\lambda, s_i)] \\ &\leq \frac{|\mathbb{E}\chi[\beta]_i|}{|\beta_i|} + n\theta \cdot \left(\max_{|s_i| \leq \frac{\lambda}{2}} f_{\beta_i}(\lambda, s_i) + \mathbb{P}\left[|s_i| > \frac{\lambda}{2}\right] \right) \\ &\leq 6n\theta^2 + \frac{2n\theta}{\sqrt{2\pi}} \max_{|s_i| \leq \frac{\lambda}{2}} \left(\frac{\lambda + |s_i|}{|\beta_i|} \cdot \exp\left[-\frac{(\lambda - |s_i|)^2}{2\beta_i^2}\right] \right) + 4n\theta^2 \\ &\leq 10n\theta^2 + n\theta \cdot \log \theta^{-1} \exp(-2 \log^2 \theta^{-1}) \\ &\leq 11n\theta^2, \end{aligned} \quad (\text{F.42})$$

and define the events $\mathcal{E}_{\|x\|_2}$, $\mathcal{E}_{\text{cross}}$ and $\mathcal{E}_{\text{pcurv}}$ as follows:

$$\begin{cases} \mathcal{E}_{\text{pcurv}} := \left\{ \forall a \in \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu)), \left\| P_{I(a)} s_{-i}[x] \right\|_2^2 \leq 11n\theta^2 + \frac{c_\mu n\theta}{p} \right\} \\ \mathcal{E}_{\text{cross}} := \left\{ \forall a \in \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu)), |\beta_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}, \max_{i \neq j \in [\pm p]} \left| e_i^* \tilde{C}_x P_{I(a)} \tilde{C}_x e_j \right| \leq 8n\theta^3 \right\} \\ \mathcal{E}_{\|x\|_2} := \left\{ \|x\|_2^2 \leq n\theta + 3\sqrt{n\theta} \log n \right\} \end{cases} \quad (\text{F.43})$$

For the Hessian term, on the event $\mathcal{E}_{\text{pcurv}} \cap \mathcal{E}_{\text{cross}} \cap \mathcal{E}_{\|x\|_2}$, and use all $\mu p^2 \theta^2$, $\mu p \theta |\tau|$ and $\theta \sqrt{p}$ are all less than $\frac{c_\mu}{4 \log^2 \theta^{-1}}$, from Lemma B.5, and from lemma assumption with sufficiently large C we have $n > \theta^{-1} 36 \log^2 n$, thus $v^* \tilde{\nabla}^2 \varphi_{\ell^1}(a) v$ can be lower bounded from (F.41) as

$$\begin{aligned} v^* \tilde{\nabla}^2 \varphi_{\ell^1}(a) v &\geq -\left(1 - \beta_{(0)}^2\right) \left(n\theta + 3\sqrt{n\theta} \log n\right) - (1 + \mu p) \left(11n\theta^2 + \frac{c_\mu n\theta}{p}\right) - 8p(1 + \mu p) \cdot 8n\theta^3 \\ &\geq -\frac{1}{2}n\theta \cdot (1 - \beta_{(0)}^2) - n\theta \cdot \left(\frac{11c_\mu}{4} + c_\mu^2 + \frac{64c_\mu}{4} + \frac{64c_\mu}{4}\right) \\ &\geq -\frac{1}{2}n\theta \cdot \left(1 - \beta_{(0)}^2 + 20c_\mu\right). \end{aligned} \quad (\text{F.44})$$

The bounds of $\beta^* \chi[\beta]$ can be derive on the event whose expectation is drawn from Lemma C.2 and (F.39) as

$$\mathcal{E}_\chi := \left\{ \begin{cases} \sigma_i \chi[\beta]_i \geq n\theta \mathcal{S}_{\nu_2 \lambda}(|\beta_i|) - \frac{c_\mu n\theta}{p}, & \forall i \in [\pm p] \\ \sigma_i \chi[\beta]_i \leq 6n\theta^2 |\beta_i| + \frac{c_\mu n\theta}{p^{3/2}}, & \forall i \neq (0) \end{cases} \right\},$$

then use $\|\beta\|_1 \leq 1 + \frac{\lambda p}{4 \log \theta^{-1}} \leq \frac{\lambda p}{2}$, implies:

$$\begin{aligned} \beta^* \chi[\beta] &\geq n\theta \left(|\beta_{(0)}| (|\beta_{(0)}| - \nu_2 \lambda) - c_\mu \|\beta\|_1 \frac{n\theta}{p} \right) \\ &\geq n\theta \left(\beta_{(0)}^2 - \sqrt{\frac{2}{\pi}} \lambda - \frac{c_\mu}{2} \lambda \right) \\ &\geq n\theta \left(\beta_{(0)}^2 - \lambda \right). \end{aligned} \quad (\text{F.45})$$

Finally via the regional condition $|\beta_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}$, the absolute value of leading correlation

$$\beta_{(0)}^2 \geq \|\beta_\tau\|_2^2 - |\tau| \beta_{(1)}^2 \geq 1 - 2c_\mu - 0.1^2 > 0.9, \quad (\text{F.46})$$

then we collect all above results and obtain:

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} = \mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} - \beta^* \chi[\beta] \geq \left(1.5\beta_{(0)}^2 - 0.5 - \lambda - 20c_\mu \right) n\theta \geq 0.3n\theta, \quad (\text{F.47})$$

with probability at least

$$1 - \underbrace{\mathbb{P}[\mathcal{E}_{\text{cross}}^c]}_{\text{Lemma D.4}} - \underbrace{\mathbb{P}[\mathcal{E}_{\text{pcurv}}^c]}_{\text{Corollary D.3}} - \underbrace{\mathbb{P}[\mathcal{E}_{\|\mathbf{x}\|_2}^c]}_{\text{Lemma A.2}} - \underbrace{\mathbb{P}[\mathcal{E}_\chi^c]}_{\text{Corollary C.4}} \geq 1 - c'/n. \quad (\text{F.48})$$

3. (Identify local minima) Wlog let \mathbf{a}_* be a local minimum where its gradient is zero that is close to \mathbf{a}_0 . The strong convexity (F.47), provides the upper bound on $\|\mathbf{a}_* - \mathbf{a}_0\|_2^2$ via

$$\begin{aligned} \varphi_{\ell^1}(\mathbf{a}_*) &\geq \varphi_{\ell^1}(\mathbf{a}_0) + \langle \mathbf{a}_* - \mathbf{a}_0, \text{grad}[\varphi_{\ell^1}](\mathbf{a}_0) \rangle + \frac{0.3}{2} n\theta \|\mathbf{a}_* - \mathbf{a}_0\|_2^2 \\ \implies \|\text{grad}[\varphi_{\ell^1}](\mathbf{a}_0)\|_2 &\geq 0.15n\theta \|\mathbf{a}_* - \mathbf{a}_0\|_2 \end{aligned} \quad (\text{F.49})$$

Thus we only require to bound the gradient at \mathbf{a}_0 , whose coefficients $\alpha = e_0$ and correlation β has properties $\beta_0 = 1$ and $\|\beta_{\setminus 0}\|_\infty \leq \mu$ hence $\|\beta_{\setminus 0}\|_2 \leq \sqrt{2p}\mu$. Expand the gradient term and condition on \mathcal{E}_χ , since $\mu p^2 \theta^2 \leq \frac{c_\mu}{4}$ and $\theta < \frac{c_\mu}{4\sqrt{p}}$, we can upper bound the gradient at \mathbf{a}_0 as

$$\begin{aligned} \|\text{grad}[\varphi_{\ell^1}](\mathbf{a}_0)\|_2 &= \|\iota^* \mathbf{C}_{\mathbf{a}_0} (\chi[\beta] - \beta^* \chi[\beta] e_0)\|_2 \leq \|\iota^* \mathbf{C}_{\mathbf{a}_0}\|_2 \|\chi[\beta]_{\setminus 0}\|_2 \\ &\leq \sqrt{1 + \mu p} \left(6n\theta^2 \|\beta_{\setminus 0}\|_2 + n\theta \cdot \frac{c_\mu}{p^{3/2}} \cdot \sqrt{2p} \right) \\ &\leq n\theta \sqrt{1 + \mu p} \left(6\mu \sqrt{2p} \cdot \theta + \frac{2c_\mu}{p} \right) \\ &\leq n\theta \left(3c_\mu \mu + 6\mu \cdot \sqrt{2\mu} \cdot p\theta + \frac{2c_\mu}{p} + \frac{2c_\mu \sqrt{\mu}}{\sqrt{p}} \right) \\ &\leq 7\sqrt{c_\mu} n\theta \cdot \max \left\{ \mu, \frac{1}{p} \right\}. \end{aligned} \quad (\text{F.50})$$

Thus we conclude that with sufficiently small c_μ :

$$\|\mathbf{a}_* - \mathbf{a}_0\|_2 \leq 50\sqrt{c_\mu} \max \left\{ \mu, p^{-1} \right\} \leq \frac{1}{2} \max \left\{ \mu, p^{-1} \right\}. \quad (\text{F.51})$$

and we complete the proof by generalize this result from minima near \mathbf{a}_0 to any of its shifts $s_i[\mathbf{a}_0]$. ■

Similarly, for objective φ_ρ we have

Corollary F.6 (Strong convexity of φ_ρ of near shift). *Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_ρ with $c_\lambda \in [\frac{1}{5}, \frac{1}{4}]$, then there*

exists some numerical constant $C, c, c', c'', \bar{c} > 0$ such that if ρ is δ -smoothed ℓ^1 function where $\delta \leq c' \lambda \theta^8 / p^2 \log^2 n$ and $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c''/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ satisfying $|\beta_{(1)}| < \nu_1 \lambda$: for all $\mathbf{v} \in \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$,

$$\mathbf{v}^* \widetilde{\text{Hess}}[\varphi_\rho](\mathbf{a}) \mathbf{v} > cn\theta; \quad (\text{F.52})$$

furthermore, there exists $\bar{\mathbf{a}}$ as an local minimizer such that

$$\min_{\ell} \|\bar{\mathbf{a}} - s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max\{\mu, p^{-1}\} \quad (\text{F.53})$$

Proof. The strong convexity (F.52) is derived by combining (F.36) and (E.23) by letting constant multiplier of δ satisfies $c'^{1/4} < 10^{-3}c$. On the other hand the local minimizer near solution (F.53) is derived via combining (F.49), (E.21) and utilize both $\theta\sqrt{p} < c_\mu$ and $\mu p^2 \theta^2 < c_\mu$ such that:

$$\begin{aligned} \|\text{grad}[\varphi_\rho](\mathbf{a})\|_2 &\leq \|\iota^* C_{\mathbf{a}_0}\|_2 \left\| \chi[\beta] - \widetilde{C}_{x_0} S_\lambda^\delta [\widetilde{C}_y \iota \mathbf{a}] \right\|_2 + \|\iota^* C_{\mathbf{a}_0}\|_2 \|\chi[\beta]_{\setminus 0}\|_2 \\ &\leq \sqrt{1 + \mu p} \cdot n\theta^3 + 7\sqrt{c_\mu} n\theta \cdot \max\{\mu, p^{-1}\} \\ &\leq 8n\theta\sqrt{c_\mu} \cdot \max\{\mu, p^{-1}\} \end{aligned} \quad (\text{F.54})$$

■

F.4 Retraction toward subspace

As in Figure 16, the function value grows in direction away from subspace \mathcal{S}_τ , we will illustrate this phenomenon by proving the negative gradient direction $-\mathbf{g}$ will point toward the subspace \mathcal{S}_τ . To show this, we prove for every coefficients of \mathbf{a} as α , there exists coefficients of \mathbf{g} as ζ satisfies

$$\langle \alpha_{\tau^c}(\mathbf{g}), \alpha_{\tau^c}(\mathbf{a}) \rangle > c \|\alpha_{\tau^c}\|_2 \|\zeta_{\tau^c}\|_2 \quad (\text{F.55})$$

whenever $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \in [\frac{\gamma}{2}, \gamma]$. Apparently, the gradient will decrease $d_\alpha(\mathbf{a}, \mathcal{S}_\tau)$, hence being addressed as *retractive toward subspace* \mathcal{S}_τ . This retractive phenomenon is true for gradient of both φ_{ℓ^1} and φ_ρ .

Lemma F.7 (Retraction of φ_{ℓ^1} toward subspace). *Suppose that $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{k}$ in φ_{ℓ^1} with $c_\lambda \in (0, \frac{1}{3}]$, then there exists some numerical constants $C, c, \bar{c} > 0$ such that if $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ such that if*

$$d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \geq \gamma(c_\mu)/2 \quad (\text{F.56})$$

then for every α satisfying $\mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha$, there exists some ζ satisfying $\text{grad}[\varphi_{\ell^1}](\mathbf{a}) = \iota^* C_{\mathbf{a}_0} \zeta$ that

$$\langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{1}{4n\theta} \|\zeta_{\tau^c}\|_2^2. \quad (\text{F.57})$$

Proof. Write $\gamma = \gamma(c_\mu)$ Recall the gradient can be derived as

$$\text{grad}[\varphi_{\ell^1}](\mathbf{a}) = -P_{\mathbf{a}^\perp} \iota^* C_{\mathbf{a}_0} \chi[\beta] = (\mathbf{a} \mathbf{a}^* - \mathbf{I}) \iota^* C_{\mathbf{a}_0} \chi[\beta] = \iota^* C_{\mathbf{a}_0} (\beta^* \chi[\beta] \alpha - \chi[\beta]), \quad (\text{F.58})$$

for every α satisfies $\mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha$. Now via Corollary C.4, condition on the event:

$$\mathcal{E}_\chi := \left\{ \sigma_i \chi[\beta]_i \leq \begin{cases} n\theta \cdot |\beta_i| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau \\ n\theta \cdot |\beta_i| 4\theta |\tau| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau^c \end{cases}, \quad \sigma_i \chi[\beta]_i \geq n\theta \cdot \mathcal{S}_{\sqrt{2/\pi\lambda}}[|\beta_i|] \right\}, \quad (\text{F.59})$$

and on this event, utilize [Lemma B.5](#), bounds of $\beta^* \chi[\beta]$ and $\|\chi[\beta]_{\tau^c}\|_2$ can be derived with $c_\mu < \frac{1}{100}$ as:

$$\beta^* \chi[\beta] \leq n\theta \left(\|\beta_\tau\|_2^2 + 4\theta |\tau| \|\beta_{\tau^c}\|_2^2 + c_\mu \right) \geq n\theta (1 + c_\mu + 4c_\mu^2 + c_\mu) \leq \frac{3}{2}n\theta \quad (\text{F.60})$$

$$\beta^* \chi[\beta] \geq n\theta \left(\|\beta_\tau\|_2^2 - \sqrt{2/\pi\lambda} \|\beta_\tau\|_1 - c_\mu \right) \geq n\theta \left(1 - 4c_\mu - \sqrt{2/\pi\lambda} c_\mu \right) \geq \frac{1}{2}n\theta \quad (\text{F.61})$$

$$\|\chi[\beta]_{\tau^c}\|_2 \leq 4n\theta^2 |\tau| \|\beta_{\tau^c}\|_2 + \frac{c_\mu n\theta}{p} \sqrt{p} \leq n\theta (4c_\mu \gamma + c_\mu \gamma) \leq \frac{1}{20}n\theta \gamma. \quad (\text{F.62})$$

Let $\alpha(g) = \beta^* \chi[\beta] \alpha - \chi[\beta]$, derive

$$\begin{aligned} & \langle \alpha(g)_{\tau^c}, \alpha_{\tau^c} \rangle - \frac{1}{4n\theta} \|\alpha(g)_{\tau^c}\|_2^2 \\ &= \beta^* \chi[\beta] \|\alpha_{\tau^c}\|_2^2 - \langle \alpha_{\tau^c}, \chi[\beta]_{\tau^c} \rangle - \frac{1}{4n\theta} \|\beta^* \chi[\beta] \alpha_{\tau^c} - \chi[\beta]_{\tau^c}\|_2^2 \\ &\geq \beta^* \chi[\beta] \|\alpha_{\tau^c}\|_2^2 - \|\alpha_{\tau^c}\|_2 \|\chi[\beta]_{\tau^c}\|_2 - \frac{1}{2n\theta} |\beta^* \chi[\beta]|^2 \|\alpha_{\tau^c}\|_2^2 - \frac{1}{2n\theta} \|\chi[\beta]_{\tau^c}\|_2^2 \\ &\geq (\beta^* \chi[\beta] - \frac{1}{2n\theta} (\beta^* \chi[\beta])^2) \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20}n\theta \gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000}n\theta \gamma^2, \end{aligned} \quad (\text{F.63})$$

notice that this is a quadratic function of $\beta^* \chi[\beta]$ with negative leading coefficient and zeros at $\{0, 2n\theta\}$, hence [\(F.63\)](#) is minimized when $\beta^* \chi[\beta] = \frac{1}{2}n\theta$. Plugging in,

$$(\text{F.63}) \geq \frac{3}{8}n\theta \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20}n\theta \gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000}n\theta \gamma^2 \quad (\text{F.64})$$

then again this is a quadratic function of $\|\alpha_{\tau^c}\|_2$ with positive leading coefficient and zeros at $\{0, \frac{8}{60}\gamma\}$, thus [\(F.64\)](#) is minimized at $\|\alpha_{\tau^c}\|_2 = \frac{\gamma}{2}$. Plugging in again,

$$(\text{F.64}) \geq \frac{3}{8}n\theta \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20}n\theta \gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000}n\theta \gamma^2 \geq \left(\frac{3}{32} - \frac{1}{80} - \frac{1}{1000} \right) n\theta \gamma^2 > 0 \quad (\text{F.65})$$

which concludes our proof. \blacksquare

As a consequence, we have that

Corollary F.8 (Retraction of φ_ρ toward the subspace). *Suppose that $x_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$ in \mathbb{R}^n , and k, c_μ such that (a_0, θ, k) satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$. Define $\lambda = c_\lambda / \sqrt{|k|}$ in φ_ρ with $c_\lambda \in (0, \frac{1}{3}]$, then there exists some numerical constants $C, c, c', c'', \bar{c} > 0$ such that if ρ is δ -smoothed ℓ^1 function where $\delta \leq c'' \lambda \theta^8 / p^2 \log^2 n$ and $n > Cp^5 \theta^{-2} \log p$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - c'/n$, for every $a \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ such that if*

$$d_\alpha(a, \mathcal{S}_\tau) \geq \gamma(c_\mu)/2 \quad (\text{F.66})$$

then for every α satisfying $a = \iota^* C_{a_0} \alpha$, there exists some ζ satisfying $\text{grad}[\varphi_\rho](a) = \iota^* C_{a_0} \zeta$ that

$$\langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{1}{6n\theta} \|\zeta_{\tau^c}\|_2^2. \quad (\text{F.67})$$

Proof. Write $\gamma = \gamma(c_\mu)$. Define

$$\chi_{\ell^1}[\beta] = \check{C}_{x_0} \mathcal{S}_\lambda [\check{a} * y], \quad \chi_\rho[\beta] = \check{C}_{x_0} \mathcal{S}_\lambda^\delta [\check{a} * y],$$

which, and on event [\(F.59\)](#) and [Lemma E.6](#), we know

$$\beta^* \chi_{\ell^1}[\beta] \leq \frac{3}{2}n\theta, \quad (\text{F.68})$$

$$\|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 \leq \frac{1}{20}n\theta \gamma, \quad (\text{F.69})$$

$$\|\chi_{\ell^1}[\beta] - \chi_\rho[\beta]\|_2 \leq c_1 n\theta^4, \quad (\text{F.70})$$

for some constant $c_1 > 0$. Now given any α satisfies $a = \iota^* C_{a_0} \alpha$, the gradient of both objective can be derived as

$$\text{grad}[\varphi_{\ell^1}](a) = -P_{a^\perp} \iota^* C_{a_0} \text{prox}_{\lambda \|\cdot\|_1} [\check{a} * y] = (aa^* - I) \iota^* C_{a_0} \chi_{\ell^1}[\beta]$$

$$= \iota^* C_{a_0} (\beta^* \chi_{\ell^1} [\beta] \alpha - \chi_{\ell^1} [\beta]), \quad (\text{F.71})$$

$$\begin{aligned} \text{grad}[\varphi_\rho](\mathbf{a}) &= -\mathbf{P}_{a^\perp} \iota^* C_{a_0} \text{prox}_{\lambda\rho}[\tilde{\mathbf{a}} * \mathbf{y}] = (\mathbf{a}\mathbf{a}^* - \mathbf{I}) \iota^* C_{a_0} \chi_\rho[\beta] \\ &= \iota^* C_{a_0} (\beta^* \chi_\rho[\beta] \alpha - \chi_\rho[\beta]). \end{aligned} \quad (\text{F.72})$$

In the same spirit, define the coefficient of each gradient vector

$$\zeta_{\ell^1} = \beta^* \chi_{\ell^1} [\beta] \alpha - \chi_{\ell^1} [\beta], \quad (\text{F.73})$$

$$\zeta_\rho = \beta^* \chi_\rho [\beta] \alpha - \chi_\rho [\beta], \quad (\text{F.74})$$

which, by norm inequality from (F.68)-(F.70) and Lemma F.7, we can derive

$$\|\zeta_{\ell^1} - \zeta_\rho\|_2 \leq \|(\mathbf{I} - \alpha\beta^*)(\chi_\rho[\beta] - \chi_{\ell^1}[\beta])\|_2 \leq c_1 n \theta^4, \quad (\text{F.75})$$

$$\|(\zeta_{\ell^1})_{\tau^c}\|_2 \geq |\beta^* \chi_{\ell^1} [\beta]| \|\alpha_{\tau^c}\|_2 - \|\chi_{\ell^1} [\beta]_{\tau^c}\|_2 \geq \frac{1}{5} n \theta \gamma, \quad (\text{F.76})$$

$$\langle (\zeta_{\ell^1})_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{1}{4n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2, \quad (\text{F.77})$$

where the first inequality is derived by observing $(\mathbf{I} - \alpha\beta^*)$ is a projection operator, as such:

$$\begin{aligned} \beta^* \alpha &= \mathbf{a}^* \iota^* C_{a_0} \alpha = \mathbf{a}^* \mathbf{a} = 1, \\ (\mathbf{I} - \alpha\beta^*)^2 &= \mathbf{I} - 2\alpha\beta^* + \alpha(\beta^* \alpha)\beta^* = \mathbf{I} - \alpha\beta^*. \end{aligned}$$

Now we are ready to derive (F.67):

$$\begin{aligned} \langle (\zeta_\rho)_{\tau^c}, \alpha_{\tau^c} \rangle &\geq \langle (\zeta_{\ell^1})_{\tau^c}, \alpha_{\tau^c} \rangle - \|\alpha_{\tau^c}\|_2 \|\zeta_\rho - \zeta_{\ell^1}\|_2 \\ &\geq \frac{1}{4n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2 - c_1 n \theta^4 \gamma \\ &\geq \frac{1}{12n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2 \\ &\quad + \frac{1}{6n\theta} \left(\|(\zeta_\rho)_{\tau^c}\|_2^2 - 2 \|(\zeta_{\ell^1})_{\tau^c}\|_2 \|\zeta_{\ell^1} - \zeta_\rho\|_2 - \|\zeta_{\ell^1} - \zeta_\rho\|_2^2 \right) - c_1 n \theta^4 \gamma \\ &\geq \frac{1}{6n\theta} \|(\zeta_\rho)_{\tau^c}\|_2^2 + \frac{1}{12n\theta} \left(\frac{1}{5} n \theta \gamma \right)^2 - \frac{1}{3n\theta} \left(\frac{1}{5} n \theta \gamma \right) (c_1 n \theta^4) - \frac{1}{6n\theta} (c_1 n \theta^4)^2 - c_1 n \theta^4 \gamma \\ &\geq \frac{1}{6n\theta} \|(\zeta_\rho)_{\tau^c}\|_2^2. \end{aligned} \quad (\text{F.78})$$

where the last inequality is true since $\theta^3 \ll \gamma$. ■

F.5 Proof of Theorem 4.1

By collecting result from above, we are ready to prove the acclaimed geometric result in Theorem 4.1. It guarantees that for every \mathbf{a} near \mathcal{S}_τ , either one of the following is true

$$\lambda_{\min}(\text{Hess}[\varphi_\rho](\mathbf{a})) \leq -c_1 n \theta \lambda, \quad (\text{F.79})$$

$$\langle \sigma_{(0)} \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq c_2 n \theta (\log^{-2} \theta^{-1}) \lambda^2, \quad (\text{F.80})$$

$$\text{Hess}[\varphi_\rho](\mathbf{a}) \succ c_3 n \theta \cdot \mathbf{P}_{a^\perp}, \quad (\text{F.81})$$

all local minimizer $\bar{\mathbf{a}}$ satisfies for some $\mathbf{a}_* \in \{\pm \iota^* s_\ell[\mathbf{a}] \mid \ell \in [\pm p_0]\}$,

$$\|\bar{\mathbf{a}} - \mathbf{a}_*\|_2 \leq c_4 \sqrt{c_\mu} \max\{\mu, p_0^{-1}\}, \quad (\text{F.82})$$

and whenever $\frac{\gamma}{2} \leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma$, coefficient of \mathbf{a} and its gradient \mathbf{g} , α , written as ζ , satisfies

$$\langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{c_5}{n\theta} \|\zeta_{\tau^c}\|_2^2. \quad (\text{F.83})$$

To connect the geometric results introduced in Lemma F.1, Lemma F.3, Lemma F.5 and Lemma F.7, we are only required to prove the required signal condition claimed in Theorem 4.1 is necessary from Definition B.1.

In particular, when the subspace dimension $|\tau| \leq 4p_0\theta$. On top of that, we are also required to show the chosen smooth parameter δ in the pseudo-Huber penalty $\rho(x) = \sqrt{x^2 + \delta^2}$ approximate $|x|$ sufficiently well, hence results of [Corollary F.2](#), [Corollary F.4](#), [Corollary F.6](#) and [Corollary F.8](#) also holds.

Proof. Firstly we will show when largest solution subspace dimension $k = 4p_0\theta$, the signal condition of [Definition B.1](#) will be satisfied. Recall that the signal condition of [Theorem 4.1](#) requests

$$\frac{2}{p_0 \log^2 p_0} \leq \theta \leq \frac{c}{(p_0 \sqrt{\mu} + \sqrt{p_0}) \log^2 p_0}, \quad (\text{F.84})$$

since $p = 3p_0 - 2$, this implies the lower bounds for sparsity θ as:

$$\theta \geq \frac{1}{2p_0 \left(\frac{1}{2} \log p_0\right)^2} \geq \frac{1}{p \log^2 \theta^{-1}}; \quad (\text{F.85})$$

the upper bound of θ via $\theta \sqrt{p_0} \log^2 p_0 \leq c$:

$$\theta \leq \frac{9c}{\sqrt{p_0}(3 \log p_0)^2} \leq \frac{16c}{\sqrt{p} \log^2 \theta^{-1}}, \quad \theta \leq \frac{4c^2}{k \log^4 p_0} \leq \frac{36c^2}{k(3 \log p_0)^2} \leq \frac{36c^2}{k \log^2 \theta^{-1}}; \quad (\text{F.86})$$

and the upper bound for coherence μ as:

$$\mu \max \{k^2, (p\theta)^2\} \log^2 \theta^{-1} \leq \mu \max \{16(p_0\theta)^2, 9(p_0\theta)^2\} \log^2 \theta^{-1} \leq 16(\sqrt{\mu}p_0\theta)^2 \log^2 p_0 \leq 16c. \quad (\text{F.87})$$

Therefore [Definition B.1](#) holds if $\max \{16c, 36c^2\} \leq c_\mu/4$ via (F.85)-(F.87).

Furthermore, we know from lemma assumption all interested \mathbf{a} are near subspace \mathcal{S}_τ by

$$d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \frac{c}{\sqrt{p_0} \log^2 \theta^{-1}} \cdot \min \left\{ \frac{1}{\sqrt{\theta}}, \frac{1}{\sqrt{\mu}} \cdot \frac{1}{\mu (p_0\theta)^{3/2}} \right\} \leq \frac{c}{\log^2 \theta^{-1}} \min \left\{ \frac{2}{\sqrt{k}}, \frac{1}{\sqrt{p_0\mu}}, \frac{4}{\mu p_0 \sqrt{\theta} k} \right\} \leq \gamma \quad (\text{F.88})$$

where γ is defined in [Definition B.3](#) of widened subspace $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$.

Lastly, the pseudo-Huber function $\rho(x) = \sqrt{x^2 + \delta^2}$ is an ℓ^1 smoothed sparse surrogate defined in [Definition E.2](#), by observing that it is convex, smooth, even, whose second order derivative (according to [Table 1](#)) $\nabla^2 \rho(x) = \frac{\delta^2}{(x^2 + \delta^2)^{3/2}}$ is monotone decreasing in $|x|$. More importantly

$$\sup_{x \in \mathbb{R}} |\rho(x) - |x|| = |\rho(0) - |0|| = \delta. \quad (\text{F.89})$$

Hence, by choosing $\delta \leq \frac{c'^4 \theta^8}{p^2 \log^2 n} \lambda$, for some sufficiently small constant c' and letting $\lambda = 0.2\sqrt{k} = 0.1/\sqrt{p_0\theta}$ in φ_ρ . We obtain the geometrical results in [Corollary F.2](#) when $|\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}|$, [Corollary F.4](#) when $\frac{4}{5} |\beta_{(0)}| \geq |\beta_{(1)}| \geq \frac{\lambda}{4 \log^2 \theta^{-1}}$ and [Corollary F.6](#) when $\frac{\lambda}{4 \log^2 \theta^{-1}} \geq |\beta_{(1)}|$, and the retraction result in [Corollary F.8](#). ■

G Analysis of algorithm — minimization within widened subspace

In this section, we prove convergence of the first part of our algorithm—minimization of φ_ρ near \mathcal{S}_τ . We begin by proving the initialization method guarantees that $\mathbf{a}^{(0)}$ is near \mathcal{S}_τ , in the sense that

$$d_\alpha(\mathbf{a}^{(0)}, \mathcal{S}_\tau) \leq \gamma, \quad (\text{G.1})$$

where the distance d_α is defined in (4.15). We then demonstrate that small-stepping curvilinear search converges to a desired local minimum of φ_ρ at rate $O(1/k)$, where k is the iteration number. To do this, it is important to utilize (i) the *retractive* property to show that the iterates stay near \mathcal{S}_τ and (ii) the geometric properties of φ_ρ near \mathcal{S}_τ .

G.1 Initialization near subspace

The following lemma shows that the initialization $\mathbf{a}^{(0)} = \mathbf{P}_{\mathbb{S}^{p-1}} [\nabla \varphi_{\ell^1}(\mathbf{a}^{(-1)})]$, where

$$\mathbf{a}^{(-1)} = \mathbf{P}_{\mathbb{S}^{p-1}} \left[\sum_{\ell \in \tau} \mathbf{x}_{0\ell} \mathbf{l}_{p_0}^* s_\ell[\mathbf{a}_0] \right], \quad (\text{G.2})$$

and is very close to the subspace \mathcal{S}_τ :

Lemma G.1 (Initialization from a piece of data). *Let $\bar{\mathbf{x}} \in \mathbb{R}^{2p_0-1}$ indexed by $[\pm p_0]$, with $\bar{\mathbf{x}}_i \sim_{\text{i.i.d.}} \text{BG}(\theta)$. Define $\bar{\mathbf{y}} = \bar{\mathbf{x}} * \mathbf{a}_0$, and $\mathbf{a}^{(0)}$ as*

$$\mathbf{a}^{(0)} = -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\ell^1} \left(\mathbf{P}_{\mathbb{S}^{p-1}} \left[\mathbf{0}^{p_0-1}; [\bar{\mathbf{y}}_0; \dots; \bar{\mathbf{y}}_{p_0-1}]; \mathbf{0}^{p_0-1} \right] \right), \quad (\text{G.3})$$

with $\lambda = 0.2/\sqrt{p\theta}$ in φ_1 . Set $\tau = \text{supp}(\bar{\mathbf{x}})$. Suppose that $(\mathbf{a}_0, \theta, k)$ satisfies the sparsity-coherence condition $\text{SCC}(c_\mu)$ and \mathbf{a}_0 satisfies $\max_{i \neq j} |\langle \mathbf{l}_{p_0}^* s_i[\mathbf{a}_0], \mathbf{l}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu$. Then there exists some constant $c, \bar{c} > 0$ such that if $p_0\theta > 1000c$ and $c_\mu \leq \bar{c}$, then with probability at least $1 - 1/c$, we have

$$d_\alpha(\mathbf{a}^{(0)}, \mathcal{S}_\tau) \leq \frac{c_\mu}{4 \log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\tau|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\tau|} \right\}. \quad (\text{G.4})$$

Proof. 1. (Distance to \mathcal{S}_τ from $\mathbf{a}^{(0)}$) Let $\eta = \|\mathbf{l}_{p_0}^*(\mathbf{a}_0 * \mathbf{x})\|_2 = \|\mathbf{l}_{p_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}\|_2$ and $\gamma = \gamma(c_\mu)$, as in (G.4). Expand the expression of $\mathbf{a}^{(0)}$ from (G.3) we have

$$\mathbf{a}^{(0)} = \mathbf{P}_{\mathbb{S}^{p-1}} \mathbf{l}^* \widetilde{\mathbf{C}}_{\mathbf{y}} \mathcal{S}_\lambda \left[\widetilde{\mathbf{C}}_{\mathbf{y}} \mathbf{l}_{p_0} \mathbf{P}_{\mathbb{S}^{p_0-1}} \mathbf{l}_{p_0}^*(\mathbf{a}_0 * \mathbf{x}) \right] = \mathbf{P}_{\mathbb{S}^{p-1}} \mathbf{l}^* \mathbf{C}_{\mathbf{a}_0} \chi \left[\frac{1}{\eta} \mathbf{C}_{\mathbf{a}_0}^* \mathbf{l}_{p_0} \mathbf{l}_{p_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x} \right] \quad (\text{G.5})$$

To relate $\mathbf{a}^{(0)}$ to its coefficient, introduce the truncated autocorrelation matrix $\widetilde{\mathbf{M}} = \mathbf{C}_{\mathbf{a}_0}^* \mathbf{l}_{p_0} \mathbf{l}_{p_0}^* \mathbf{C}_{\mathbf{a}_0}$, define $\tilde{\alpha}, \tilde{\beta}$ as

$$\tilde{\beta} = \frac{1}{\eta} \widetilde{\mathbf{M}} \mathbf{x}, \quad \tilde{\alpha} = \chi \left[\frac{1}{\eta} \widetilde{\mathbf{M}} \mathbf{x} \right] = \chi[\tilde{\beta}] \quad (\text{G.6})$$

and note that $\widetilde{\mathbf{M}}$ is bounded entrywise as

$$|\widetilde{\mathbf{M}}_{ij}| \leq \begin{cases} 1 & i = j \in [-p_0 + 1, p_0 - 1] \\ \mu & i \neq j \in [-p_0 + 1, p_0 - 1], |i - j| < p_0 \\ 0 & \text{otherwise} \end{cases} \quad (\text{G.7})$$

From (G.5), we can write $\mathbf{a}^{(0)} = \mathbf{P}_{\mathbb{S}^{p-1}} \mathbf{l}^* \mathbf{C}_{\mathbf{a}_0} \tilde{\alpha}$, meaning that the normalized version of $\tilde{\alpha}$ is a valid coefficient vector for $\mathbf{a}^{(0)}$. Let $\tau^c = [\pm 2p_0] \setminus \tau$. The distance d_α to subspace \mathcal{S}_τ (4.15) is upper bounded as

$$d_\alpha(\mathbf{a}^{(0)}, \mathcal{S}_\tau) \leq \frac{\|\tilde{\alpha}_{\tau^c}\|_2}{\|\mathbf{l}^* \mathbf{C}_{\mathbf{a}_0} \tilde{\alpha}\|_2} \leq \frac{\|\tilde{\alpha}_{\tau^c}\|_2}{\|\mathbf{l}^* \mathbf{C}_{\mathbf{a}_0} \tilde{\alpha}_\tau\|_2 - \|\mathbf{l}^* \mathbf{C}_{\mathbf{a}_0} \tilde{\alpha}_{\tau^c}\|_2}$$

$$\leq \frac{\|\tilde{\alpha}_{\tau^c}\|_2}{\sqrt{1-\mu|\tau|}\|\tilde{\alpha}_{\tau}\|_2 - \sqrt{1+\mu p}\|\tilde{\alpha}_{\tau^c}\|_2}$$

where the last inequality is derived with [Lemma B.4](#). Therefore, it is sufficient to show

$$(1 + \gamma\sqrt{1+\mu p})\|\tilde{\alpha}_{\tau^c}\|_2 \leq \gamma\sqrt{1-\mu|\tau|}\|\tilde{\alpha}_{\tau}\|_2 \quad (\text{G.8})$$

to complete the proof that $d_{\alpha}(\mathbf{a}^{(0)}, \mathcal{S}_{\tau}) \leq \gamma$.

2. (Bound η) Condition on the following two events

$$\mathcal{E}_{\tau} := \{|\tau| < 4p_0\theta\}, \quad \mathcal{E}_{\|x\|_2} := \left\{\sqrt{p_0\theta} \leq \|x\|_2 \leq \sqrt{3p_0\theta}\right\} \quad (\text{G.9})$$

and utilize μ bound from [Lemma B.5](#) such that $\mu|\tau| < 0.1$. An upper bound on η can be obtained using properties of \tilde{M} of [\(G.7\)](#):

$$\eta = \|\iota_{p_0}^* \mathbf{C}_{a_0} \mathbf{x}\|_2 \leq \|\iota^* \mathbf{C}_{a_0} \mathbf{x}\|_2 \leq \sqrt{1+\mu|\tau|}\|\mathbf{x}\|_2 \leq 2\sqrt{p_0\theta} \quad (\text{G.10})$$

To lower bound η , use $\eta^2 = \mathbf{g}^* \mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau} \mathbf{g}$ where \mathbf{g} is the standard Gaussian vector. Observe the submatrix of \tilde{M} is diagonal dominant:

$$\begin{cases} \tilde{M}_{ii} = \|\iota_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 \in [0, 1] \\ \text{tr}(\tilde{M}) = \sum_{i \in [\pm p_0]} \|\iota_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 = \|\mathbf{a}_0\|_2^2 + \sum_{i=1}^{p_0-1} \left(\|\iota_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 + \|\iota_{p_0}^* s_{i-p_0}[\mathbf{a}_0]\|_2^2 \right) = p_0 \end{cases} \quad (\text{G.11})$$

Write $\mathbf{x} = \mathbf{g} \circ \mathbf{w}$ where \mathbf{w} and \mathbf{g} are Bernoulli and Gaussian vector respectively with $\text{supp}(\mathbf{w}) = \tau$, then the trace of $\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}$ can be written as sum of independent r.v.s as:

$$\text{tr}(\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}) = \sum_{i \in [\pm p_0]} w_i \|\iota_{p_0}^* s_i[\mathbf{a}_0]\|_2^2,$$

Bernstein inequality [Lemma J.4](#) and [\(G.11\)](#) gives

$$\mathbb{P}\left[\text{tr}(\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}) < \frac{3p_0\theta}{4}\right] \leq \mathbb{P}\left[\text{tr}(\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}) - p_0\theta \leq -\frac{p_0\theta}{4}\right] \leq 2 \exp\left(\frac{-(p_0\theta/4)^2}{2p_0\theta + p_0\theta/2}\right) \leq 2 \exp\left(\frac{-p_0\theta}{40}\right), \quad (\text{G.12})$$

thus condition on ω satisfies $\text{tr}(\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}) \geq 3p_0\theta/4$ and \mathcal{E}_{τ} , expectation η^2 has lower bound

$$\mathbb{E}_{\mathbf{g}|\mathbf{w}} \eta^2 = \mathbb{E}_{\mathbf{g}|\mathbf{w}} [\mathbf{g}^* \mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau} \mathbf{g}] = \text{tr}(\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau}) \geq \frac{3p_0\theta}{4}$$

then apply Bernstein inequality again by first writing svd of $\mathbf{P}_{\tau} \tilde{M} \mathbf{P}_{\tau} = \mathbf{U} \Sigma \mathbf{U}^*$ with Σ being rank $|\tau| < 4p_0\theta$ and square orthobasis \mathbf{U} . Let $\mathbf{g}' = \mathbf{U}^* \mathbf{g}$, then \mathbf{g}' is standard i.i.d. Gaussian vector, provides alternative expression $\eta^2 < \sum_{i=1}^{4p_0\theta} g_i'^2 \sigma_i$ where $\sigma_i \leq 1 + \mu|\tau| \leq 1.1$. We obtain probability of η^2 to be small as

$$\mathbb{P}_{\mathbf{g}|\mathbf{w}} \left[\eta^2 < \frac{p_0\theta}{2} \right] \leq \mathbb{P}_{\mathbf{g}|\mathbf{w}} \left[\eta^2 - \mathbb{E}_{\mathbf{g}|\mathbf{w}} \eta^2 < -\frac{p_0\theta}{4} \right] \leq 2 \exp\left(\frac{-(p_0\theta/4)^2}{2(1+\mu|\tau|)(12p_0\theta + p_0\theta/2)}\right) \leq 2 \exp\left(\frac{-p_0\theta}{440}\right) \quad (\text{G.13})$$

by applying moment bounds $(\sigma^2, R) = (12p_0\theta(1 + \mu|\tau|), 2(1 + \mu|\tau|))$. We thereby define event

$$\mathcal{E}_\eta = \left\{ \sqrt{p_0\theta/2} \leq \eta \leq 2\sqrt{p_0\theta} \right\}, \quad (\text{G.14})$$

which holds w.h.p. based on (G.9), (G.12) and (G.13).

3. (Bound $\tilde{\alpha}$) Condition on $\mathcal{E}_\eta \cap \mathcal{E}_{\|x\|_2} \cap \mathcal{E}_\tau$. Use definition $\tilde{\beta} = \frac{1}{\eta} \tilde{M}x$ from (G.6), and properties of \tilde{M} from (G.7) we can obtain:

$$\begin{cases} \|\tilde{\beta}_{\tau^c}\|_2 \leq \frac{1}{\eta} \left\| \iota_{\tau^c}^* \tilde{M} \iota_{\tau} \right\|_2 \|x\|_2 \leq \frac{\mu\sqrt{p_0|\tau|}}{\sqrt{p_0\theta/2}} \cdot \sqrt{3p_0\theta} \leq 3\mu\sqrt{p_0|\tau|} \\ \|\tilde{\beta}_\tau\|_2 \geq \frac{1}{\eta} \left\| \iota_\tau^* \tilde{M} \iota_{\tau^c} \right\|_2 \|x\|_2 \geq \frac{\sqrt{1-\mu}|\tau|}{2\sqrt{p_0\theta}} \cdot \sqrt{p_0\theta} \geq 0.45 \end{cases}. \quad (\text{G.15})$$

Use definition $\|\tilde{\alpha}\|_2 = \|\chi[\tilde{\beta}]\|_2$, condition on event

$$\mathcal{E}_\chi := \left\{ \begin{cases} \sigma_i \chi[\beta]_i \geq n\theta \mathcal{S}_{\nu_2\lambda} [|\beta_i|] - \frac{c_\mu n\theta}{p}, & \forall i \in \tau \\ \sigma_i \chi[\beta]_i \leq 4n\theta^2 |\tau| |\beta_i| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau^c \end{cases} \right\},$$

also from Definition B.1 we have $\mu(p\theta)^{1/2} |\tau|^{3/2} < \frac{c_\mu}{4\log^2 \theta^{-1}}$ and from lemma assumption $\lambda = \frac{1}{5\sqrt{p\theta}}$, provides bounds of $\|\tilde{\alpha}\|_2$ via triangle inequality as:

$$\begin{cases} \|\tilde{\alpha}_{\tau^c}\|_2 \leq 4n\theta^2 |\tau| \cdot \|\tilde{\beta}_{\tau^c}\|_2 + \frac{c_\mu n\theta}{p} \cdot \sqrt{2p_0} \leq 3c_\mu n\theta \left(\frac{\sqrt{\theta}}{\log^2 \theta^{-1}} + \frac{c_\mu}{p} \right) \\ \|\tilde{\alpha}_\tau\|_2 \geq n\theta \left(\|\tilde{\beta}_\tau\|_2 - \nu_2 \lambda \sqrt{|\tau|} - \frac{c_\mu}{p} \sqrt{|\tau|} \right) \geq n\theta \left(0.45 - \sqrt{\frac{2}{\pi}} \cdot \frac{1}{5} - c_\mu \right) \geq 0.2n\theta \end{cases}, \quad (\text{G.16})$$

since both $\theta|\tau|, \mu p\theta|\tau| < c_\mu$, we have

$$\begin{cases} \sqrt{1 + \mu p} \|\tilde{\alpha}_{\tau^c}\|_2 \leq 3c_\mu n\theta \sqrt{1 + \mu p} (\sqrt{\theta} + p^{-1}) \leq 6c_\mu n\theta \\ \|\tilde{\alpha}_{\tau^c}\|_2 \leq \frac{6c_\mu^{3/2} n\theta}{\log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\tau|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\tau|} \right\} \leq 24\sqrt{c_\mu} n\theta \gamma \end{cases},$$

which satisfies (G.8), since $\mu|\tau| < c_\mu < \frac{1}{1000}$,

$$(1 + \gamma\sqrt{1 + \mu p}) \|\tilde{\alpha}_{\tau^c}\|_2 \leq (24\sqrt{c_\mu} + 6c_\mu) n\theta \gamma \leq 0.1n\theta \gamma \leq \gamma\sqrt{1 - \mu|\tau|} \|\tilde{\alpha}_\tau\|_2. \quad (\text{G.17})$$

Finally, given $p_0\theta > 1000c$, this result holds with probability at least

$$1 - \underbrace{\mathbb{P}[\mathcal{E}_\tau^c]}_{\text{Lemma A.1}} - \underbrace{\mathbb{P}[\mathcal{E}_{\|x\|_2}^c]}_{\text{Lemma A.2}} - \underbrace{\mathbb{P}[\mathcal{E}_\eta^c]}_{(\text{G.14})} - \underbrace{\mathbb{P}[\mathcal{E}_\chi^c]}_{\text{Corollary C.4}} \geq 1 - \frac{2}{p_0\theta} - \frac{1}{n} - 4 \exp\left(\frac{-p_0\theta}{440}\right) \geq 1 - \frac{1}{c} \quad (\text{G.18})$$

■

G.2 Minimization near subspace (Proof of Theorem 5.1)

Before we start the proof of theorem, writing $g = \text{grad}[\varphi_\rho](a)$ and $H = \text{Hess}[\varphi_\rho](a)$, we will first restate the results of Theorem 4.1 in simplified terms. The theorem shows that for any $a \in \mathbb{S}^{p-1}$ whose distance to subspace $d_\alpha(a, \mathcal{S}_\tau) \leq \gamma$, then at least one of the the following statement hold:

$$\|g\|_2 \geq \eta_g \quad (\text{G.19})$$

$$\lambda_{\min}(\mathbf{H}) \leq -\eta_v \quad (\text{G.20})$$

$$\mathbf{H} \succ \eta_c \cdot \mathbf{P}_{\mathbf{a}^\perp}. \quad (\text{G.21})$$

Furthermore, φ_ρ is retractive near \mathcal{S}_τ : wherever $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \geq \frac{\gamma}{2}$, writing $\alpha(\mathbf{a}), \alpha(\mathbf{g})$ to be the coefficient of \mathbf{a}, \mathbf{g} , we have

$$\langle \alpha(\mathbf{a})_{\tau^c}, \alpha(\mathbf{g})_{\tau^c} \rangle \geq \eta_r \|\alpha(\mathbf{g})_{\tau^c}\|_2. \quad (\text{G.22})$$

Also, the the gradient, Hessian and the third order derivative are all bounded as follows:

Remark G.2. *With high probability, for every \mathbf{a} whose $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) < \gamma$, its $\max\{\|\mathbf{g}\|_2, \|\mathbf{H}\|_2, \|\nabla \mathbf{H}\|_2\} \leq \bar{\eta} = \text{poly}(n, p)$.*

We state [Remark G.2](#) without explicit proof since its derivation is similar to the proof in [Theorem 4.1](#).

We prove that if the negative curvature direction $-\mathbf{v}$ is chosen to be the least eigenvector with $\mathbf{v}^* \mathbf{H} \mathbf{v} < -\eta_v$ and $\mathbf{v}^* \mathbf{g}$ (if cannot, let $\mathbf{v} = \mathbf{0}$), then the iterates

$$\mathbf{a}^{(k+1)} = \mathbf{P}_{\mathbb{S}^{p-1}} \left[\mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)} \right] \quad (\text{G.23})$$

converges toward the minimizer $\bar{\mathbf{a}}$ in ℓ^2 -norm with rate $O(1/k)$. Notice that here all $\eta_g, \eta_v, \eta_c, \eta_r, \bar{\eta}$ are all greater than 0 and are rational functions of the dimension parameters n, p .

Finally, we should note that \mathbf{a}_0 being μ -truncated shift coherent implies that \mathbf{a}_0 is at at most 2μ -shift coherent. Hence we utilize the usual incoherence condition in the proof.

Proof. Notice that when \mathbf{a} is in the region near some signed shift $\bar{\mathbf{a}}$ of \mathbf{a}_0 , the function φ_ρ is strongly convex, and the iterates coincide with the Riemannian gradient method, which converges at a linear rate. Indeed, if for all k larger than some \bar{k} , $\mathbf{a}^{(k)}$ is in this region, then $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \leq (1 - t\eta_c)^{-(k-\bar{k})} \|\mathbf{a}^{(\bar{k})} - \bar{\mathbf{a}}\|_2$ [[AMS09](#)](Theorem 4.5.6) where the step size $t = \Omega(1/n\theta)$ hence $t\eta_c = \Omega(1)$. We will argue that the iterates $\mathbf{a}^{(k)}$ remain close to the subspace \mathcal{S}_τ and that after $\bar{k} = \text{poly}(n, p)$ iterations they indeed remain in the strongly convex region around some $\bar{\mathbf{a}}$.

1. (Existence of Armijo steplength). First, we show there exists a nontrivial step size t at every iteration, in the sense that for all $\mathbf{a} \in \mathbb{S}^{p-1}$, there exists $T > 0$ such that for all $t \in (0, T)$, the Armijo step condition (5.11) is satisfied. Note that since φ_ρ is a smooth function, $\mathbf{a} \rightarrow \varphi_\rho \circ \mathbf{P}_{\mathbb{S}^{p-1}}(\mathbf{a})$ admits a version of Taylor's theorem (see also [[AMS09](#)](Section 7.1.3)): for any $\boldsymbol{\xi} \perp \mathbf{a}$, writing $\mathbf{a}^+ = \mathbf{P}_{\mathbb{S}^{p-1}}[\mathbf{a} + \boldsymbol{\xi}]$,

$$|\varphi_\rho(\mathbf{a}^+) - (\varphi_\rho(\mathbf{a}) + \langle \text{grad}[\varphi_\rho](\mathbf{a}), \boldsymbol{\xi} \rangle + \frac{1}{2} \boldsymbol{\xi}^* \text{Hess}[\varphi_\rho](\mathbf{a}) \boldsymbol{\xi})| \leq \bar{\eta} \|\boldsymbol{\xi}\|_2^3, \quad (\text{G.24})$$

using $\|\nabla \mathbf{H}\|_2 \leq \bar{\eta}$. Now, let $\boldsymbol{\xi} = -t\mathbf{g} - t^2\mathbf{v}$ as in the iterates (5.10). Suppose the Armijo step condition (5.11) does not hold, so

$$\varphi_\rho(\mathbf{a}^+) > \varphi_\rho(\mathbf{a}) - \frac{1}{2} \left(t \|\mathbf{g}\|_2^2 + \frac{1}{2} t^4 \eta_v \|\mathbf{v}\|_2^2 \right). \quad (\text{G.25})$$

Since $\mathbf{g}^* \mathbf{v} \geq 0$ and $\mathbf{v}^* \mathbf{H} \mathbf{v} \leq -\eta_v \|\mathbf{v}\|_2^2$ or $\mathbf{v} = \mathbf{0}$, using $\|\mathbf{a} + \mathbf{b}\|_2^3 \leq 4 \|\mathbf{a}\|_2^3 + 4 \|\mathbf{b}\|_2^3$ (Hölder's inequality) and $\|\mathbf{H}\|_2 < \bar{\eta}$, we can derive

$$\begin{aligned} & \langle \mathbf{g}, -t\mathbf{g} - t^2\mathbf{v} \rangle + \frac{1}{2} (t\mathbf{g} + t^2\mathbf{v})^* \mathbf{H} (t\mathbf{g} + t^2\mathbf{v}) + c \|\mathbf{g} + t^2\mathbf{v}\|_2^3 > -\frac{1}{2} \left(t \|\mathbf{g}\|_2^2 + \frac{1}{2} t^4 \eta_v \|\mathbf{v}\|_2^2 \right) \\ \implies & -\frac{1}{2} t \|\mathbf{g}\|_2^2 + \frac{1}{2} t^2 \mathbf{g}^* \mathbf{H} \mathbf{g} + t^3 \mathbf{v}^* \mathbf{H} \mathbf{g} - \frac{1}{4} t^4 \eta_v \|\mathbf{v}\|_2^2 + 4\bar{\eta} t^3 \|\mathbf{g}\|_2^3 + 4\bar{\eta} t^6 \|\mathbf{v}\|_2^3 > 0 \\ \implies & -\frac{1}{2} t \|\mathbf{g}\|_2^2 + t^2 \left(\frac{1}{2} \bar{\eta} \|\mathbf{g}\|_2^2 + t\bar{\eta} \|\mathbf{v}\|_2 \|\mathbf{g}\|_2 + 4\bar{\eta} t \|\mathbf{g}\|_2^3 \right) - \frac{1}{4} t^4 \eta_v \|\mathbf{v}\|_2^2 + 4\bar{\eta} t^6 \|\mathbf{v}\|_2^3 > 0. \end{aligned} \quad (\text{G.26})$$

If

$$t < T = \min \left\{ \frac{\|\mathbf{g}\|_2}{\bar{\eta} \|\mathbf{g}\|_2 + 2\bar{\eta} t \|\mathbf{v}\|_2 + 8\bar{\eta} t \|\mathbf{g}\|_2^2}, \sqrt{\frac{\eta_v}{16\bar{\eta} \|\mathbf{v}\|_2}} \right\}, \quad (\text{G.27})$$

then (G.26) < 0 contradicting (G.25). Using our bounds on $\|g\|_2$, $\bar{\eta}$, η_v and $\|v\|$, it follows that T is lower bounded by a polynomial $\text{poly}(n^{-1}, p^{-1})$.

2. (Bounds on $d_\alpha(g, \mathcal{S}_\tau)$, $d_\alpha(v, \mathcal{S}_\tau)$) We will show there are numerical constants c_g, c_v such that

$$d_\alpha(g, \mathcal{S}_\tau) \leq c_g n \theta \gamma \quad \text{and} \quad d_\alpha(v, \mathcal{S}_\tau) \leq c_v n \theta p. \quad (\text{G.28})$$

Define

$$\chi_{\ell^1}[\beta] = \check{C}_{x_0} \text{prox}_{\lambda \ell^1}[\check{a} * y], \quad \chi_\rho[\beta] = \check{C}_{x_0} \text{prox}_{\lambda \rho}[\check{a} * y],$$

then the gradient can be written as (F.71)

$$\text{grad}[\varphi_{\ell^1}](a) = \iota^* C_{a_0} (\beta^* \chi_{\ell^1}[\beta] \alpha - \chi_{\ell^1}[\beta]), \quad (\text{G.29})$$

$$\text{grad}[\varphi_\rho](a) = \iota^* C_{a_0} (\beta^* \chi_\rho[\beta] \alpha - \chi_\rho[\beta]). \quad (\text{G.30})$$

Use the following inequalities:

$$\begin{aligned} \frac{1}{2} n \theta &\leq |\beta^* \chi_{\ell^1}[\beta]| \leq \frac{3}{2} n \theta, \\ \|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 &\leq \frac{1}{20} n \theta \gamma, \\ \|I - \alpha \beta^*\|_2 &\leq 4\sqrt{p}, \\ \|\chi_{\ell^1}[\beta] - \chi_\rho[\beta]\|_2 &\leq n \theta^4, \end{aligned}$$

where the first and second bounds of $\chi_{\ell^1}[\beta]$ based on event (F.59); the third by observing $\|\alpha\|_2 \leq 2$ and $\|\beta\|_2 \leq 2 + c_\mu \sqrt{p}$; the last from (E.21) of Lemma E.6 when δ is sufficiently small. Hence, by definition of $d_\alpha(\cdot, \mathcal{S}_\tau)$ (4.15) and knowing a is close to subspace $\|\alpha_{\tau^c}\|_2 \leq \gamma$, via triangle inequality, we get

$$\begin{aligned} d_\alpha(g, \mathcal{S}_\tau) &\leq d_\alpha(\text{grad}[\varphi_{\ell^1}](a), \mathcal{S}_\tau) + d_\alpha(\text{grad}[\varphi_\rho](a) - \text{grad}[\varphi_{\ell^1}](a), \mathcal{S}_\tau) \\ &\leq \|\beta^* \chi_{\ell^1}[\beta] \alpha_{\tau^c} - \chi_{\ell^1}[\beta]_{\tau^c}\|_2 + \|(I - \alpha \beta^*)(\chi_\rho[\beta] - \chi_{\ell^1}[\beta])\|_2 \\ &\leq \frac{3}{2} n \theta \gamma + \frac{1}{20} n \theta \gamma + 4\sqrt{p} n \theta^4 \\ &\leq 3 n \theta \gamma. \end{aligned} \quad (\text{G.31})$$

To bound the d_α norm of least eigenvector v , note that $\beta^* \chi_\rho[\beta] > 0$, we can conclude

$$v^* \nabla^2 \varphi_\rho(a) v \leq v^* P_{a^\perp} \nabla^2 \varphi_\rho(a) P_{a^\perp} v + \beta^* \chi_\rho[\beta] = v^* H v < -\eta_v,$$

expand $\nabla^2 \varphi_\rho(a)$ as in (E.8), and since v is the eigenvector of smallest eigenvalue $\lambda_{\min} < -\eta_v$,

$$P_{a^\perp} \nabla^2 \varphi_\rho(a) P_{a^\perp} v = (I - a a^*) \iota^* C_{a_0} \check{C}_{x_0} \text{prox}_{\lambda \rho}[\check{a} * y] \check{C}_{x_0} C_{a_0}^* \iota v = \lambda_{\min} v, \quad (\text{G.32})$$

hence there exists $\alpha(v)$ satisfies $v = \iota^* C_{a_0} \alpha(v)$ and

$$\alpha(v) = \lambda_{\min}^{-1} \left[\check{C}_{x_0} \text{prox}_{\lambda \rho}[\check{a} * y] \check{C}_{x_0} C_{a_0}^* \iota v - \left(\beta^* \check{C}_{x_0} \text{prox}_{\lambda \rho}[\check{a} * y] \check{C}_{x_0} C_{a_0}^* \iota v \right) \alpha \right].$$

Now since $\text{prox}_{\lambda \rho}[\check{a} * y]$ is a diagonal matrix with entries in $[0, 1]$,

$$d_\alpha(v, \mathcal{S}_\tau) \leq \|\alpha(v)\|_2 \leq |\lambda_{\min}|^{-1} \|\iota C_{a_0}\|_2 \|x_0\|_1^2 \|v\|_2 (1 + \|\alpha\|_2 \|\beta\|_2) < c_v n \theta p, \quad (\text{G.33})$$

where we use upper bound of $\|x_0\|_1 < c n \theta$ from Lemma A.2 and $|\lambda_{\min}| > \eta_v > c n \theta \lambda$ from Corollary F.2.

3. (Iterates stay within widened subspace). Suppose (G.22) holds. We will show that whenever

$$t \leq T' = \frac{1}{10 n \theta}, \quad (\text{G.34})$$

then setting $\mathbf{a}^+ = \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}]$, we have

$$|d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) - d_\alpha(\mathbf{a}, \mathcal{S}_\tau)| \leq \frac{\gamma}{2}, \quad (\text{G.35})$$

and whenever $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \in [\frac{\gamma}{2}, \gamma]$

$$d_\alpha^2(\mathbf{a}^+, \mathcal{S}_\tau) \leq d_\alpha^2(\mathbf{a}, \mathcal{S}_\tau) - t \cdot c'n\theta\gamma^2. \quad (\text{G.36})$$

If both (G.35) and (G.36) hold, then all iterates $\mathbf{a}^{(k)}$ will stay near the subspace: $d_\alpha(\mathbf{a}^{(k)}, \mathcal{S}_\tau) < \gamma$.

To derive (G.35), since both $\mathbf{g} \perp \mathbf{a}$ and $\mathbf{v} \perp \mathbf{a}$ we have $\|\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}\|_2^2 = \|\mathbf{a}\|_2^2 + \|t\mathbf{g} + t^2\mathbf{v}\|_2^2 > 1$, and since $d_\alpha(\cdot, \mathcal{S}_\tau)$ is a seminorm Lemma B.2:

$$\begin{aligned} d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) &= d_\alpha(\mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}], \mathcal{S}_\tau) \leq d_\alpha(\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}, \mathcal{S}_\tau) \\ &\leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau) + td_\alpha(\mathbf{g}, \mathcal{S}_\tau) + t^2d_\alpha(\mathbf{v}, \mathcal{S}_\tau) \end{aligned} \quad (\text{G.37})$$

suggests (G.35) holds via (G.28) and let $n > Cp^5\theta^{-2}$, we have

$$td_\alpha(\mathbf{g}, \mathcal{S}_\tau) + t^2d_\alpha(\mathbf{v}, \mathcal{S}_\tau) \leq \frac{c_g n\theta\gamma}{10n\theta} + \frac{c_v n\theta p}{(10n\theta)^2} < \frac{\gamma}{2} \quad (\text{G.38})$$

with sufficiently large C .

To derive (G.36), let $\alpha(\mathbf{a})$ to be a coefficient vector satisfying $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = \|\alpha(\mathbf{a})_{\tau^c}\|_2$, and based on (G.30) and (G.33), define

$$\alpha(\mathbf{g}) = \beta^* \chi_\rho[\beta] \alpha(\mathbf{a}) - \chi_\rho[\beta] \quad (\text{G.39})$$

$$\alpha(\mathbf{v}) = \lambda_{\min}^{-1} \check{\mathbf{C}}_{x_0} \nabla_{\text{prox}_{\lambda\rho}} [\check{\mathbf{a}} * \mathbf{y}] \check{\mathbf{C}}_{x_0}^* \mathbf{C}_{a_0}^* \mathbf{v}. \quad (\text{G.40})$$

By the retraction property and norm bounds,

$$\langle \alpha(\mathbf{a})_{\tau^c}, \alpha(\mathbf{g})_{\tau^c} \rangle \geq \frac{1}{6n\theta} \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 \quad (\text{G.41})$$

$$\|\alpha(\mathbf{a})_{\tau^c}\|_2 \leq \gamma \quad (\text{G.42})$$

$$\|\alpha(\mathbf{v})\|_2 \leq c_v n\theta p. \quad (\text{G.43})$$

Since $\|\alpha_{\tau^c}\|_2 > \frac{\gamma}{2}$,

$$\begin{aligned} \|\alpha(\mathbf{g})_{\tau^c}\|_2 &\geq \|\beta^* \chi_{\ell^1}[\beta] \alpha_{\tau^c} - \chi_{\ell^1}[\beta]_{\tau^c}\|_2 - \|(\mathbf{I} - \alpha\beta^*)(\chi_\rho[\beta] - \chi_{\ell^1}[\beta])\|_2 \\ &\geq |\beta^* \chi_{\ell^1}[\beta]| \|\alpha_{\tau^c}\|_2 - \|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 - \|(\mathbf{I} - \alpha\beta^*)(\chi_\rho[\beta] - \chi_{\ell^1}[\beta])\|_2 \\ &\geq \frac{1}{2}n\theta \times \frac{\gamma}{2} - \frac{1}{20}n\theta\gamma + 2n\theta^4 \\ &\geq \frac{1}{10}n\theta\gamma. \end{aligned} \quad (\text{G.44})$$

Finally, we can bound $d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau)$ as

$$\begin{aligned} d_\alpha^2(\mathbf{a}^+, \mathcal{S}_\tau) &\leq d_\alpha^2(\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}, \mathcal{S}_\tau) \\ &\leq \|[\alpha(\mathbf{a}) - t\alpha(\mathbf{g}) - t^2\alpha(\mathbf{v})]_{\tau^c}\|_2^2 \\ &= \|\alpha(\mathbf{a})_{\tau^c}\|_2^2 - 2t \langle \alpha(\mathbf{a})_{\tau^c}, [\alpha(\mathbf{g}) + t\alpha(\mathbf{v})]_{\tau^c} \rangle + t^2 \|[\alpha(\mathbf{g}) + t\alpha(\mathbf{v})]_{\tau^c}\|_2^2 \\ &\leq \|\alpha(\mathbf{a})_{\tau^c}\|_2^2 - 2t \langle \alpha(\mathbf{a})_{\tau^c}, \alpha(\mathbf{g})_{\tau^c} \rangle + 2t^2 \|\alpha(\mathbf{a})_{\tau^c}\|_2 \|\alpha(\mathbf{v})\|_2 + 2t^2 \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 + 2t^4 \|\alpha(\mathbf{v})\|_2^2 \\ &\leq d^2(\mathbf{a}, \mathcal{S}_\tau) - 2t \left[\left(\frac{1}{3n\theta} - t \right) \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 - tn\theta p\gamma - t^3(c_v n\theta p)^2 \right] \\ &\leq d^2(\mathbf{a}, \mathcal{S}_\tau) - t \cdot c'n\theta\gamma^2 \end{aligned} \quad (\text{G.45})$$

where the last inequality holds when $t < \frac{0.1}{n\theta}$ with sufficiently large n .

4. (Polynomial time convergence) The iterates $\mathbf{a}^{(k)}$ remain within a γ neighborhood of S_τ for all k . At any iteration k , $\mathbf{a}^{(k)}$ is in at least one of three regions: strong gradient, negative curvature, or strong convexity. In the gradient and curvature regions, we obtain a decrease in the function value which is at least some (nonzero) rational function of n and p . On the strongly convex region, the function value does not increase. The suboptimality at initialization is bounded by a polynomial in n and p , $\text{poly}(n, p)$, and hence the total number of steps in the gradient and curvature regions is bounded by a polynomial in n, p . After the iterates reach the strongly convex region, the number of additional steps required to achieve $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \leq \varepsilon$ is bounded by $\text{poly}(n, p) \log \varepsilon^{-1}$. In particular, the number of iterations required to achieve $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \leq \mu + 1/p$ is bounded by a polynomial in (n, p) , as claimed. ■

H Analysis of algorithm — local refinement

In this section, we describe and analyze an algorithm which refines an estimate $\mathbf{a}^{(0)} \approx \mathbf{a}_0$ of the kernel to exactly recover $(\mathbf{a}_0, \mathbf{x}_0)$. Set

$$\lambda^{(0)} \leftarrow 5\kappa_I \tilde{\mu} \quad \text{and} \quad I^{(0)} \leftarrow \text{supp}(\mathcal{S}_\lambda [\mathbf{C}_{\mathbf{a}^{(0)}}^* \mathbf{y}]), \quad (\text{H.1})$$

where as each iteration of the algorithm consists of the following key steps:

- **Sparse Estimation using Reweighted Lasso:** Set

$$\mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \sum_{i \notin I^{(k)}} \lambda^{(k)} |\mathbf{x}_i|; \quad (\text{H.2})$$

- **Kernel Estimation using Least Squares:** Set

$$\mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[\underset{\mathbf{a}}{\text{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right]; \quad (\text{H.3})$$

- **Continuation and reweighting by decreasing sparsity regularizer:** Set

$$\lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)} \quad \text{and} \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}). \quad (\text{H.4})$$

Our analysis will show that $\mathbf{a}^{(k)}$ converges to \mathbf{a}_0 at a linear rate. In the remainder of this section, we describe the assumptions of our analysis. In subsequent sections, we prove key lemmas analyzing each of the three main steps of the algorithm.

Modified coherence and rate assumptions Below, we will write

$$\tilde{\mu} = \max \{ \mu, p^{-1} \}. \quad (\text{H.5})$$

Our refinement algorithm will demand an initialization satisfying

$$\|\mathbf{a}^{(0)} - \mathbf{a}_0\|_2 \leq \tilde{\mu}. \quad (\text{H.6})$$

Support density of \mathbf{x}_0 Our goal is to show that the proposed annealing algorithm exactly solves the SaS deconvolution problem, i.e., exactly recovers $(\mathbf{a}_0, \mathbf{x}_0)$ up to a signed shift. We will denote the support sets of true sparse vector \mathbf{x}_0 and recovered $\mathbf{x}^{(k)}$ in the intermediate k -th steps as

$$I = \text{supp}(\mathbf{x}_0), \quad I^{(k)} = \text{supp}(\mathbf{x}^{(k)}). \quad (\text{H.7})$$

It should be clear that exact recovery is unlikely if \mathbf{x}_0 contains many consecutive nonzero entries: in this situation, even *non-blind* deconvolution fails. We introduce the notation κ_I as an upper bound for number of nonzero entries of \mathbf{x}_0 in a length- p window:

$$\kappa_I = 6 \max \{ \theta p, \log n \}, \quad (\text{H.8})$$

then in the Bernoulli-Gaussian model, with high probability,

$$\max_{\ell} |I \cap ([p] + \ell)| \leq \kappa_I. \quad (\text{H.9})$$

Here, indexing and addition should be interpreted modulo n . The $\log n$ term reflects the fact that as n becomes enormous (exponential in p) eventually it becomes likely that some length- p window of \mathbf{x}_0 is densely

occupied. In our main theorem statement, we preclude this possibility by putting an upper bound on n (w.r.t $\tilde{\mu}$). We find it useful to also track the maximum ℓ^2 norm of \mathbf{x}_0 over any length- p window:

$$\|\mathbf{x}_0\|_{\square} := \max_{\ell} \|\mathbf{P}_{([p]+\ell)} \mathbf{x}_0\|_2. \quad (\text{H.10})$$

Below, we will sometimes work with the \square -induced operator norm:

$$\|\mathbf{M}\|_{\square \rightarrow \square} = \sup_{\|\mathbf{x}\|_{\square} \leq 1} \|\mathbf{M}\mathbf{x}\|_{\square} \quad (\text{H.11})$$

For now, we note that in the Bernoulli-Gaussian model, $\|\mathbf{x}_0\|_{\square}$ is typically not large

$$\|\mathbf{x}_0\|_{\square} \leq \sqrt{\kappa_I}. \quad (\text{H.12})$$

H.1 Reweighted Lasso finds the large entries of \mathbf{x}_0

The following lemma asserts that when \mathbf{a} is close to \mathbf{a}_0 , the reweighted Lasso finds all of the large entries of \mathbf{x}_0 . Our reweighted Lasso is modified version from [CWB08], we only penalize \mathbf{x} on entries outside of its known support subset. We write T to be the subset of true support I , and define the sparsity surrogate as

$$\sum_{i \in T^c} |\mathbf{x}_i| \quad (\text{H.13})$$

The reweighted Lasso recovers more accurate \mathbf{x} on set T compares to the vanilla Lasso problem, it turns out to be very helpful in our analysis which proves convergence of the proposed alternating minimization.

Lemma H.1 (Accuracy of reweighted Lasso estimate). *Suppose that $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ with \mathbf{a}_0 is $\tilde{\mu}$ -shift coherent and $\|\mathbf{x}_0\|_{\square} \leq \sqrt{\kappa_I}$ with $\kappa_I \geq 1$. If $\tilde{\mu}\kappa_I^2 \leq c_{\mu}$, then for every $T \subseteq I$ and \mathbf{a} satisfying $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$, the solution \mathbf{x}^+ to the optimization problem*

$$\min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i \in T^c} |\mathbf{x}_i| \right\}, \quad (\text{H.14})$$

with

$$\lambda > 5\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2, \quad (\text{H.15})$$

is unique with the form

$$\mathbf{x}^+ = \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a}J}^* \mathbf{C}_{\mathbf{a}J})^{-1} \boldsymbol{\iota}_J^* (\mathbf{C}_{\mathbf{a}}^* \mathbf{y} - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}) \quad (\text{H.16})$$

where $\boldsymbol{\sigma} = \text{sign}(\mathbf{x}^+)$. Its support set J satisfies

$$(T \cup I_{\geq 3\lambda}) \subseteq J \subseteq I \quad (\text{H.17})$$

and its entrywise error is bounded as

$$\|\mathbf{x}^+ - \mathbf{x}_0\|_{\infty} \leq 3\lambda. \quad (\text{H.18})$$

Above, $c_{\mu} > 0$ is a positive numerical constant.

We prove [Lemma H.1](#) below. The proof relies heavily on the fact that when \mathbf{a}_0 is shift-incoherent and $\mathbf{a} \approx \mathbf{a}_0$, \mathbf{a} is also shift-incoherent, an observation which is formalized in a sequence of calculations in [Appendix H.4](#).

Proof. 1. (Restricted support Lasso problem). We first consider the restricted problem

$$\min_{\mathbf{w} \in \mathbb{R}^{|I|}} \left\{ \frac{1}{2} \|\mathbf{a} * \boldsymbol{\iota}_I \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{i \in T^c} |(\boldsymbol{\iota}_I \mathbf{w})_i| \right\}. \quad (\text{H.19})$$

Under our assumptions, provided $c < \frac{1}{9}$, [Lemma H.6](#) implies that

$$\iota_I^* C_a^* C_a \iota_I = [C_a^* C_a]_{I,I} \succ \mathbf{0}, \quad (\text{H.20})$$

and the restricted problem is strongly convex and its solution is unique. The KKT conditions imply that a vector \mathbf{w}_* is the unique optimal solution to this problem if and only if

$$\iota_I^* C_a^* C_a \iota_I \mathbf{w}_* \in \iota_I^* C_a^* \mathbf{y} - \lambda \partial \| \mathbf{P}_{T^c} [\cdot] \|_1 (\mathbf{w}_*). \quad (\text{H.21})$$

Writing $J = \text{supp}(\iota_I \mathbf{w}_*) \subseteq I$, $C_{a,J} = C_a \iota_J$, $\mathbf{w}_J = \iota_J^* \mathbf{w}_*$ the corresponding sub-vector containing the nonzero entries of \mathbf{w}_* and $\sigma_{J \setminus T} = \iota_J^* \mathbf{P}_{T^c} [\text{sign}(\iota_I \mathbf{w}_*)]$, the condition [\(H.21\)](#) is satisfied if and only if

$$C_{a,J}^* C_{a,J} \mathbf{w}_J = C_{a,J}^* \mathbf{y} - \lambda \sigma_{J \setminus T}, \quad (\text{H.22})$$

$$\| C_{a, I \setminus J}^* (C_{a,J} \mathbf{w}_J - \mathbf{y}) \|_\infty \leq \lambda. \quad (\text{H.23})$$

We will argue that under our assumptions, J necessarily contains all of the large entries of \mathbf{x}_0 :

$$I_{>3\lambda} = \{\ell \in I \mid |\mathbf{x}_{0\ell}| > 3\lambda\} \subseteq J. \quad (\text{H.24})$$

We show this by contradiction – namely, if some large entry of \mathbf{x}_0 is not in J , then the dual condition [\(H.23\)](#) is violated, contradicting the optimality of \mathbf{w}_* . To this end, note that by [Corollary H.7](#), $C_{a,J}^* C_{a,J}$ has full rank. From [\(H.22\)](#),

$$\mathbf{w}_J = [C_{a,J}^* C_{a,J}]^{-1} [C_{a,J}^* \mathbf{y} - \lambda \sigma_{J \setminus T}]. \quad (\text{H.25})$$

Write $\mathbf{x}_{0,J} = \iota_J^* \mathbf{x}_0$ and $(\mathbf{x}_0)_{I \setminus J} = \mathbf{P}_{I \setminus J} \mathbf{x}_0$. We can further notice that

$$\begin{aligned} C_{a,J} \mathbf{w}_J - \mathbf{y} &= \left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) \mathbf{y} - \lambda C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} \sigma_{J \setminus T} \\ &= \left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) C_{a_0 J} \mathbf{x}_{0,J} + \left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) C_{a_0 I \setminus J} (\mathbf{x}_0)_{I \setminus J} \\ &\quad - \lambda C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} \sigma_{J \setminus T} \\ &= \left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) C_{a_0 - a, J} \mathbf{x}_{0,J} + \left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) C_{a_0 I \setminus J} (\mathbf{x}_0)_{I \setminus J} \\ &\quad - \lambda C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} \sigma_{J \setminus T}, \end{aligned} \quad (\text{H.26})$$

where in the final line we have used that

$$\left(C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* - \mathbf{I} \right) C_{a,J} = \mathbf{0}. \quad (\text{H.27})$$

Suppose that J is a strict subset of I (otherwise, if $J = I$, we are done). Take any $i \in I \setminus J$ such that $|\mathbf{x}_{0i}| = \|(\mathbf{x}_0)_{I \setminus J}\|_\infty$, and let $\xi = \text{sign}(\mathbf{x}_{0i})$. Using [\(H.26\)](#), [Corollary H.7](#) and [Lemma H.8](#), we have

$$\begin{aligned} -\xi s_i[\mathbf{a}]^* (C_{a,J} \mathbf{w}_J - \mathbf{y}) &= \xi s_i[\mathbf{a}]^* \left(\mathbf{I} - C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* \right) s_i[\mathbf{a}_0] \mathbf{x}_{0i} \\ &\quad + \xi s_i[\mathbf{a}]^* \left(\mathbf{I} - C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* \right) C_{a_0} (\mathbf{x}_0)_{I \setminus (J \cup \{i\})} \\ &\quad + \xi s_i[\mathbf{a}]^* \left(\mathbf{I} - C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} C_{a,J}^* \right) C_{a_0 - a, J} \mathbf{x}_{0,J} \\ &\quad + \xi \lambda s_i[\mathbf{a}]^* C_{a,J} [C_{a,J}^* C_{a,J}]^{-1} \sigma_{J \setminus T} \\ &\geq \left(\langle s_i[\mathbf{a}], s_i[\mathbf{a}_0] \rangle - \|s_i[\mathbf{a}]^* C_{a,J}\|_1 \left\| [C_{a,J}^* C_{a,J}]^{-1} \right\|_{\infty \rightarrow \infty} \|C_{a,J}^* s_i[\mathbf{a}_0]\|_\infty \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\ &\quad - \left(\left\| s_i[\mathbf{a}]^* C_{a_0 I \setminus \{i\}} \right\|_1 + \|s_i[\mathbf{a}]^* C_{a,J}\|_1 \left\| [C_{a,J}^* C_{a,J}]^{-1} \right\|_{\infty \rightarrow \infty} \|C_{a,J}^* C_{a_0 I \setminus J}\|_{\infty \rightarrow \infty} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\ &\quad - \left(\|s_i[\mathbf{a}]^* C_{a_0 - a, J}\|_2 + \|s_i[\mathbf{a}]^* C_{a,J}\|_2 \left\| [C_{a,J}^* C_{a,J}]^{-1} \right\|_{\square \rightarrow \square} \|C_{a,J}^* C_{a_0 - a, J}\|_{\square \rightarrow \square} \right) \sqrt{2} \|\mathbf{x}_0\|_\square \end{aligned} \quad (\text{H.28})$$

$$\begin{aligned}
& -\lambda \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty \rightarrow \infty} \|\boldsymbol{\sigma}_{J \setminus T}\|_\infty \\
\geq & \left((1 - \|\mathbf{a} - \mathbf{a}_0\|_2) - C_1 \kappa_I \tilde{\mu} \times 1 \times \tilde{\mu} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
& - C_2 \left(\kappa_I \tilde{\mu} + \kappa_I \tilde{\mu} \times 1 \times \kappa_I \tilde{\mu} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
& - \left(2\sqrt{\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 + C_3 \sqrt{\kappa_I} \tilde{\mu} \times 1 \times \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 \right) \|\mathbf{x}_0\|_\square \\
& - \lambda C_4 \kappa_I \tilde{\mu}
\end{aligned} \tag{H.29}$$

$$\begin{aligned}
& \geq \left(1 - C_1' \kappa_I \tilde{\mu} - C_2 (\kappa_I \tilde{\mu})^2 \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
& - 2\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 - \left(C_3 \kappa_I^{3/2} \tilde{\mu} \right) \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 - (C_4 \kappa_I \tilde{\mu}) \lambda
\end{aligned} \tag{H.30}$$

$$\begin{aligned}
& \geq \left(1 - C_1' \kappa_I \tilde{\mu} - C_2 (\kappa_I \tilde{\mu})^2 \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
& - 2\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 - \left(C_3 \kappa_I^{3/2} \tilde{\mu} \right) \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 - (C_4 \kappa_I \tilde{\mu}) \lambda
\end{aligned} \tag{H.31}$$

$$\geq \frac{1}{2} \|(\mathbf{x}_0)_{I \setminus J}\|_\infty - \lambda/2, \tag{H.32}$$

where the last line holds provided $\tilde{\mu} \kappa_I^2 \leq c_\mu$ to be a sufficiently small numerical constants. If $\|(\mathbf{x}_0)_{I \setminus J}\|_\infty > 3\lambda$, this is strictly larger than λ , implying that $|\mathbf{a}_i^* (\mathbf{C}_{aJ} \mathbf{w}_J - \mathbf{y})| > \lambda$, and contradicting the KKT conditions for the restricted problem. Hence, under our assumptions

$$\|(\mathbf{x}_0)_{I \setminus J}\|_\infty \leq 3\lambda. \tag{H.33}$$

2. (Solution of Full Lasso problem) We next argue that the solution of the restricted support Lasso problem, \mathbf{w}_J , when extended to \mathbb{R}^n as $\mathbf{x}^+ = \iota_J \mathbf{w}_J$, is the unique optimal solution to the *full* Lasso problem

$$\min_{\mathbf{x}} \varphi_{\text{lasso}}(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i \in T^c} |\mathbf{x}_i|. \tag{H.34}$$

To prove that \mathbf{x}^+ is the unique optimal solution, it suffices to show that for every $i \in I^c$,

$$|s_i[\mathbf{a}]^* (\mathbf{a} * \mathbf{x}^+ - \mathbf{y})| < \lambda. \tag{H.35}$$

Indeed, suppose that this inequality is in force. Write $\varepsilon = \lambda - \max_{i \in I^c} |s_i[\mathbf{a}]^* (\mathbf{a} * \mathbf{x}^+ - \mathbf{y})|$, and notice that from the KKT conditions for the restricted problem,

$$\mathbf{0} \in \mathbf{P}_I \partial_{\mathbf{x}} \varphi_{\text{lasso}}(\mathbf{x}) \tag{H.36}$$

Combining with (H.35), we have that for every vector $\boldsymbol{\zeta}$ with $\text{supp}(\boldsymbol{\zeta}) \subseteq I^c$ and $\|\boldsymbol{\zeta}\|_\infty \leq 1$, then $\varepsilon \boldsymbol{\zeta} \in \partial \varphi_{\text{lasso}}(\mathbf{x}^+)$. Let \mathbf{x}' be any vector with $\mathbf{x}'_{I^c} \neq \mathbf{0}$ and set $\boldsymbol{\zeta} = \mathcal{P}_{I^c} \text{sign}(\mathbf{x}')$, then from the subgradient inequality,

$$\begin{aligned}
\varphi_{\text{lasso}}(\mathbf{x}') & \geq \varphi_{\text{lasso}}(\mathbf{x}^+) + \langle \varepsilon \boldsymbol{\zeta}, \mathbf{x}' - \mathbf{x}^+ \rangle \\
& \geq \varphi_{\text{lasso}}(\mathbf{x}^+) + \varepsilon \|\mathbf{x}'_{I^c}\|_1,
\end{aligned} \tag{H.37}$$

which is strictly larger than $\varphi_{\text{lasso}}(\mathbf{x}^+)$. Hence, when (H.35) holds, any optimal solution $\bar{\mathbf{x}}$ to the full Lasso problem must satisfy $\text{supp}(\bar{\mathbf{x}}) \subseteq I$. By strong convexity of the restricted problem, the solution to (H.34) is unique and equal to \mathbf{x}^+ .

We finish by showing (H.35). Using the same expansion as above, we obtain

$$\begin{aligned}
|s_i[\mathbf{a}]^* (\mathbf{C}_{aJ} \mathbf{w}_J - \mathbf{y})| & \leq \left| s_i[\mathbf{a}]^* \left(\mathbf{I} - \mathbf{C}_{aJ} [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \mathbf{C}_{aJ}^* \right) \mathbf{C}_{a_0 I \setminus J} (\mathbf{x}_0)_{I \setminus J} \right| \\
& + \left| s_i[\mathbf{a}]^* \left(\mathbf{I} - \mathbf{C}_{aJ} [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \mathbf{C}_{aJ}^* \right) \mathbf{C}_{a_0 - a_J} \mathbf{x}_{0J} \right| \\
& + \lambda \left| s_i[\mathbf{a}]^* \mathbf{C}_{aJ} [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \right| \\
& \leq \left(\left\| s_i[\mathbf{a}]^* \mathbf{C}_{a_0 I \setminus J} \right\|_1 + \left\| s_i[\mathbf{a}]^* \mathbf{C}_{aJ} \right\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty \rightarrow \infty} \left\| \mathbf{C}_{aJ}^* \mathbf{C}_{a_0 I \setminus J} \right\|_{\infty \rightarrow \infty} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_\infty
\end{aligned} \tag{H.38}$$

$$\begin{aligned}
& + \left(\|s_i[\mathbf{a}]^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}_J}\|_2 + \|s_i[\mathbf{a}]^* \mathbf{C}_{\mathbf{a}_J}\|_2 \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \right\|_{\square \rightarrow \square} \|\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}_J}\|_{\square \rightarrow \square} \right) \sqrt{2} \|\mathbf{x}_0\|_{\square} \\
& + \lambda \|s_i[\mathbf{a}]^* \mathbf{C}_{\mathbf{a}_J}\|_1 \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \right\|_{\infty \rightarrow \infty} \|\boldsymbol{\sigma}_{J \setminus T}\|_{\infty} \tag{H.39} \\
\leq & C_1 (\tilde{\mu} \kappa_I + \tilde{\mu} \kappa_I \times 1 \times \tilde{\mu} \kappa_I) \times 2\lambda \\
& + (2\sqrt{\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 + C_2 \sqrt{\kappa_I} \tilde{\mu} \times 1 \times \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2) \times \sqrt{\kappa_I} \\
& + \lambda C_3 \times \tilde{\mu} \kappa_I \tag{H.40} \\
\leq & ((C_1 + C_3) \tilde{\mu} \kappa_I + C_1 (\tilde{\mu} \kappa_I)^2) \lambda + (2 + C_2 \tilde{\mu} \kappa_I) \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 \tag{H.41} \\
< & \lambda, \tag{H.42}
\end{aligned}$$

where the last line holds as long as c_μ is a sufficiently small numerical constant. This establishes that \mathbf{x}^+ is the unique optimal solution to the full Lasso problem.

3. (Entrywise difference to \mathbf{x}_0) Finally we will be controlling $\|\mathbf{x}_J^+ - (\mathbf{x}_0)_J\|_{\infty}$. Indeed, from [Lemma H.8](#),

$$\begin{aligned}
\|\mathbf{x}_J^+ - (\mathbf{x}_0)_J\|_{\infty} &= \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \lambda [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} - (\mathbf{x}_0)_J \right\|_{\infty} \\
&\leq \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}_J} (\mathbf{x}_0)_J \right\|_{\infty} + \lambda \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \right\|_{\infty} \\
&\quad + \left\| [\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}]^{-1} \mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_{I \setminus J}} (\mathbf{x}_0)_{I \setminus J} \right\|_{\infty} \\
&\leq 2 \|\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}_J}\|_{\square \rightarrow \infty} \|(\mathbf{x}_0)_J\|_{\square} + 2\lambda + 2 \|\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_{I \setminus J}}\|_{\infty \rightarrow \infty} \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \\
&\leq 2\sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 \|\mathbf{x}_0\|_{\square} + 2\lambda + 2 \times 3\tilde{\mu} \times 2\kappa_{I \setminus J} \times 3\lambda \\
&\leq 3\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 + 2\lambda + 36\lambda \tilde{\mu} \kappa_I \\
&\leq 3\lambda, \tag{H.43}
\end{aligned}$$

establishing the claim. \blacksquare

H.2 Least squares solution $\mathbf{a}^{(k)}$ contracts

Approximation of least squares solution. In this section, given \mathbf{x} to be the solution to the reweighted Lasso from \mathbf{a} , we will show the solution of the least squares problem

$$\mathbf{a}^+ \leftarrow \underset{\mathbf{a}' \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{a}' * \mathbf{x} - \mathbf{y}\|_2^2 \tag{H.44}$$

is closer to \mathbf{a}_0 compared to \mathbf{a} . Observe that in [Lemma H.1](#), the solution of (H.16)

$$\mathbf{x} = \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} \boldsymbol{\iota}_J^* (\mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}), \tag{H.45}$$

by assuming $\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J} \approx \mathbf{I}$, $\mathbf{a} \approx \mathbf{a}_0$ and $J \setminus T \approx \emptyset$, is a good approximation to the true sparse map \mathbf{x}_0

$$\mathbf{x} \approx \mathbf{I}(\mathbf{x}_0 - \mathbf{0}) = \mathbf{x}_0; \tag{H.46}$$

furthermore, its difference to the true sparse map $\|\mathbf{x}_0 - \mathbf{x}\|_2$ is proportional to $\|\mathbf{a}_0 - \mathbf{a}\|_2$ as

$$\mathbf{x} - \mathbf{x}_0 \approx \mathbf{P}_I (\mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}} \mathbf{x}_0) \approx \mathbf{P}_I [\mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}(\mathbf{a}_0 - \mathbf{a})]. \tag{H.47}$$

To this end, since we know the solution of least square problem \mathbf{a}^+ is simply

$$\mathbf{a}^+ = (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0), \tag{H.48}$$

this implies the difference between the new \mathbf{a}^+ and \mathbf{a}_0 , has the relationship with $\mathbf{a} - \mathbf{a}_0$ roughly

$$\begin{aligned}
\mathbf{a}^+ - \mathbf{a}_0 &= (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota} \mathbf{a}_0) \approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{x}_0 - \mathbf{x}) \\
&\approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a} - \mathbf{a}_0). \tag{H.49}
\end{aligned}$$

To make this point precise, we introduce the following lemma:

Lemma H.2 (Approximation of least square estimate). *Given $\mathbf{a}_0 \in \mathbb{R}^{p_0}$ to be $\tilde{\mu}$ -shift coherent and $\mathbf{x}_0 \sim \text{BG}(\theta) \in \mathbb{R}^n$. There exists some constants C, C', c, c', c_μ such that if $\lambda < c' \tilde{\mu} \kappa_I$, $\tilde{\mu} \kappa_I^2 \leq c_\mu$ and $n > Cp^2 \log p$, then with probability at least $1 - c/n$, for every \mathbf{a} satisfying $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$ and \mathbf{x} of the form*

$$\mathbf{x} = \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a}, J}^* \mathbf{C}_{\mathbf{a}, J})^{-1} \boldsymbol{\iota}_J^* (\mathbf{C}_{\mathbf{a}}^* \mathbf{y} - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}) \quad (\text{H.50})$$

where the set J, T satisfies $I_{>6\lambda} \subseteq T \subseteq J \subseteq I$, we have

$$\frac{1}{n\theta} \left\| \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}_0 - \mathbf{a}) \right\|_2 \leq C' \lambda \left(\tilde{\lambda} + \tilde{\mu} \kappa_I \right) + \frac{1}{32} \|\mathbf{a} - \mathbf{a}_0\|_2 \quad (\text{H.51})$$

with $\tilde{\lambda} = \lambda + \frac{\log n}{\sqrt{n\theta^2}}$.

Proof. We will begin with listing the conditions we use for both \mathbf{x} and \mathbf{x}_0 . First, we know from [Lemma H.1](#) and our assumptions on the set T , then \mathbf{x} approximates \mathbf{x}_0 in the sense that

$$\|\mathbf{x} - \mathbf{x}_0\|_\infty \leq 3\lambda \quad (\text{H.52})$$

$$\|(\mathbf{x}_0)_{I \setminus J}\|_\infty \leq 3\lambda \quad (\text{H.53})$$

$$\|(\mathbf{x}_0)_{I \setminus T}\|_\infty \leq 6\lambda. \quad (\text{H.54})$$

Write $\mathbf{x}_0 = \mathbf{g} \circ \boldsymbol{\omega}$ with \mathbf{g} iid standard normal, $\boldsymbol{\omega}$ iid Bernoulli and \mathbf{g} and $\boldsymbol{\omega}$ independent. From (H.53) we know $|I \setminus J| = |\{i \mid |\mathbf{g}_i| \leq 3\lambda, \boldsymbol{\omega}_i \neq 0\}|$. Since $\mathbb{P}[\boldsymbol{\omega}_i \neq 0] = \theta$ and $\mathbb{P}[|\mathbf{g}_i| \leq 3\lambda] \leq 3\lambda$, [Lemma A.1](#) implies that with probability at least $1 - 2/n$:

$$|I \setminus J| \leq 3\lambda n\theta + 6\sqrt{\lambda n\theta} \log n \leq 3\tilde{\lambda} n\theta \quad (\text{H.55})$$

$$|I \setminus T| \leq 6\lambda n\theta + 12\sqrt{\lambda n\theta} \log n \leq 6\tilde{\lambda} n\theta, \quad (\text{H.56})$$

and

$$|(I \setminus J) \cap s_\ell[I]| \leq 3\lambda n\theta^2 + 6\sqrt{\lambda n\theta^2} \log n \leq 3\tilde{\lambda} n\theta^2; \quad (\text{H.57})$$

together with base on properties of Bernoulli-Gaussian vector \mathbf{x}_0 from [Appendix A](#) and we conclude with probability at least $1 - c/n$, all the following events hold:

$$\frac{1}{2}n\theta \leq |I| \leq 2n\theta, \quad (\text{H.58})$$

$$\max_{\ell \neq 0} |I \cap s_\ell[I]| \leq 2n\theta^2 \quad (\text{H.59})$$

$$\max_{\ell \neq 0} |(I \setminus J) \cap s_\ell[I]| \leq 6\tilde{\lambda} n\theta^2, \quad (\text{H.60})$$

$$\|\mathbf{x}_0\|_\square^2 \leq \kappa_I, \quad (\text{H.61})$$

$$\|\tilde{\mathbf{a}}_0 * \mathbf{x}_0\|_\square^2 \leq \kappa_I, \quad (\text{H.62})$$

$$\|\mathbf{x}_0\|_2^2 \leq 2n\theta, \quad (\text{H.63})$$

$$\|\mathbf{x}_0\|_1 \leq 2n\theta, \quad (\text{H.64})$$

$$\max_{\ell \neq 0} \|\mathbf{P}_{I \cap s_\ell[I]} \mathbf{x}_0\|_2^2 \leq 2n\theta^2, \quad (\text{H.65})$$

$$\max_{\ell \neq 0} \|\mathbf{P}_{I \cap s_\ell[I \setminus J]} \mathbf{x}_0\|_1 \leq 12\tilde{\lambda} n\theta^2, \quad (\text{H.66})$$

$$\|\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}\|_2^2 \leq 3n\theta, \quad (\text{H.67})$$

provided by $n \geq C\theta^{-2} \log p$ for sufficiently large constant C .

1. (Approximate $\mathbf{C}_{\mathbf{x}}$ with $\mathbf{C}_{\mathbf{x}_0}$) Since

$$\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{x}-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 + \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}-\mathbf{x}_0}^* \mathbf{C}_{\mathbf{x}-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 \quad (\text{H.68})$$

where

$$\begin{aligned}
\|\iota^* C_{x-x_0}^* C_{x-x_0} \iota a_0\|_2 &\leq \|a_0\|_2 \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \|C_{a_0} \iota\|_2 \sqrt{2p} \max_{\ell \neq 0} |\langle s_\ell[\mathbf{x} - \mathbf{x}_0], \mathbf{x} - \mathbf{x}_0 \rangle| \\
&\leq \|\mathbf{x} - \mathbf{x}_0\|_\infty^2 \times |I| + \sqrt{2\tilde{\mu}p^2} \left(\|\mathbf{x} - \mathbf{x}_0\|_\infty^2 \times \max_{\ell \neq 0} |I \cap s_\ell[I]| \right) \\
&\leq C_1 \left(\lambda^2 n \theta + \sqrt{2\tilde{\mu}p^2} (\lambda^2 n \theta^2) \right) \\
&\leq 2C_1 \lambda^2 n \theta,
\end{aligned} \tag{H.69}$$

we have that

$$\|\iota^* C_{\mathbf{x}}^* C_{\mathbf{x}-\mathbf{x}_0} \iota a_0 - \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{x}-\mathbf{x}_0} \iota a_0\|_2 \leq 2C_1 \lambda^2 n \theta. \tag{H.70}$$

2. (Extract the $a_0 - a$ term) Observe that

$$\begin{aligned}
&\iota^* C_{\mathbf{x}_0}^* C_{\mathbf{x}-\mathbf{x}_0} \iota a_0 \\
&= \iota^* C_{\mathbf{x}_0}^* C_{a_0} (\mathbf{x} - \mathbf{x}_0) \\
&= \iota^* C_{\mathbf{x}_0}^* C_{a_0} \left(\iota_J (C_{a,J}^* C_{a,J})^{-1} \iota_J^* (C_{a,J}^* C_{a_0} \mathbf{x}_0 - \lambda P_{J \setminus T} \sigma) - \iota_J (C_{a,J}^* C_{a,J})^{-1} (C_{a,J}^* C_{a,J}) (\mathbf{x}_0)_J - P_{I \setminus J} \mathbf{x}_0 \right) \\
&= \iota^* C_{\mathbf{x}_0}^* C_{a_0 J} (C_{a,J}^* C_{a,J})^{-1} C_{a,J}^* (C_{a_0} \mathbf{x}_0 - C_{a,J} (\mathbf{x}_0)_J) \\
&\quad + \iota^* C_{\mathbf{x}_0}^* C_{a_0 J} (C_{a,J}^* C_{a,J})^{-1} C_{a,J}^* (C_a \mathbf{x}_0 - C_{a,J} (\mathbf{x}_0)_J) \\
&\quad - \iota^* C_{\mathbf{x}_0}^* C_{a_0} P_{I \setminus J} \mathbf{x}_0 \\
&\quad - \lambda \iota^* C_{\mathbf{x}_0}^* C_{a_0 J} (C_{a,J}^* C_{a,J})^{-1} \iota_J^* P_{J \setminus T} \sigma,
\end{aligned} \tag{H.71}$$

where, the second term in (H.71) is bounded as

$$\begin{aligned}
&\|\iota^* C_{\mathbf{x}_0}^* C_{a_0 J} (C_{a,J}^* C_{a,J})^{-1} C_{a,J}^* (C_a \mathbf{x}_0 - C_{a,J} (\mathbf{x}_0)_J)\|_2 \\
&\leq \|C_{\mathbf{x}_0} \iota\|_2 \times \|C_{a_0 J}\|_2 \|(C_{a,J}^* C_{a,J})^{-1}\|_2 \times \|C_{a,J}^* C_{a I \setminus J}\|_2 \times \|(\mathbf{x}_0)_{I \setminus J}\|_2 \\
&\leq C_2 \left(\sqrt{n\theta} \times 3 \times \tilde{\mu} \kappa_I \times \lambda \sqrt{\tilde{\lambda} n \theta} \right) \\
&\leq 3C_2 \tilde{\mu} \kappa_I \lambda n \theta;
\end{aligned} \tag{H.72}$$

the third term in (H.71) is bounded as

$$\begin{aligned}
\|\iota^* C_{\mathbf{x}_0}^* C_{a_0} P_{I \setminus J} \mathbf{x}_0\|_2 &= \|\iota^* C_{a_0} (P_{[\pm p] \setminus 0} + e_0 e_0^*) C_{\mathbf{x}_0}^* P_{I \setminus J} \mathbf{x}_0\|_2 \\
&\leq \|a_0\|_2 \|(\mathbf{x}_0)_{I \setminus J}\|_2^2 + \|C_{a_0} \iota\|_2 \times \sqrt{2p} \times \max_{\ell \neq 0} \|P_{I \cap s_\ell[I \setminus J]} \mathbf{x}_0\|_1 \times \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
&\leq C_3 \left(\lambda^2 \times \tilde{\lambda} n \theta + \sqrt{\tilde{\mu} p^2} \times \tilde{\lambda} n \theta^2 \times \lambda \right) \\
&\leq 2C_3 \tilde{\lambda} \lambda n \theta;
\end{aligned} \tag{H.73}$$

and finally, write $\Delta = (C_{a,J}^* C_{a,J})^{-1} - I$, then the forth term in (H.71) is bounded as

$$\begin{aligned}
&\lambda \|\iota^* C_{\mathbf{x}_0}^* C_{a_0} \iota_J (C_{a,J}^* C_{a,J})^{-1} \iota_J^* P_{J \setminus T} \sigma\|_2 \\
&= \lambda \|\iota^* C_{a_0} (P_{[\pm p] \setminus 0} + e_0 e_0^*) C_{\mathbf{x}_0}^* \iota_J (I + \Delta) \iota_J^* P_{J \setminus T} \sigma\|_2 \\
&\leq \lambda \|C_{a_0}^* \iota\|_2 \sqrt{2p} \max_{\ell \neq 0} \|P_{I \cap s_\ell[I \setminus T]} \mathbf{x}_0\|_1 + \lambda \|a_0\|_2 \|P_{I \setminus T} \mathbf{x}_0\|_1 \\
&\quad + \lambda \|C_{a_0}^* \iota\|_2 \sqrt{2p} \|P_{I \cap s_\ell[I]} \mathbf{x}_0\|_1 \|\Delta\|_{\infty \rightarrow \infty} + \lambda \|a_0\|_2 \|\mathbf{x}_0\|_2 \|\Delta\|_2 \sqrt{|J \setminus T|} \\
&\leq C_4 \lambda \left(\sqrt{\tilde{\mu} p^2} \times \tilde{\lambda} n \theta^2 + \lambda \tilde{\lambda} n \theta + \sqrt{\tilde{\mu} p^2} \times n \theta^2 \times \tilde{\mu} \kappa_I + \sqrt{n\theta} \times \tilde{\mu} \kappa_I \sqrt{\tilde{\lambda} n \theta} \right)
\end{aligned}$$

$$\leq 2C_4 \left(\tilde{\lambda} + \tilde{\mu}\kappa_I \right) \lambda n \theta. \quad (\text{H.74})$$

Therefore, combining (H.72)-(H.74) we obtain

$$\left\| \iota^* C_{x_0}^* C_{x-x_0} \iota a_0 - \iota^* C_{x_0}^* C_{a_0 J} (C_{aJ}^* C_{aJ})^{-1} C_{aJ}^* C_{a_0-a} x_0 \right\|_2 \leq C_5 \left(\tilde{\lambda} + \tilde{\mu}\kappa_I \right) \lambda n \theta. \quad (\text{H.75})$$

3. (Extract the set J) Lastly, we will further simplify the term with $a - a_0$ in (H.75) by extracting the set J :

$$\begin{aligned} & \iota^* C_{x_0}^* C_{a_0 J} (C_{aJ}^* C_{aJ})^{-1} C_{aJ}^* C_{a_0-a} x_0 \\ &= \iota^* C_{x_0}^* C_{a_0 J} (I + \Delta) C_{a_0+(a-a_0)J}^* C_{x_0} \iota (a_0 - a) \\ &= \iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota (a_0 - a) \\ & \quad + \iota^* C_{x_0}^* C_{a_0 J} \Delta C_{a_0 J}^* C_{x_0} \iota (a_0 - a) + \iota^* C_{x_0}^* C_{a_0 J} (C_{aJ}^* C_{aJ})^{-1} C_{a-a_0 J}^* C_{x_0} \iota (a_0 - a) \\ & \quad - \iota^* C_{x_0}^* C_{a_0} P_{I \setminus J} C_{a_0}^* C_{x_0} \iota (a_0 - a), \end{aligned} \quad (\text{H.76})$$

where, the latter terms in (H.76) are bounded as

$$\begin{aligned} \left\| \iota^* C_{x_0}^* C_{a_0 J} \Delta C_{a_0 J}^* C_{x_0} \iota \right\|_2 &\leq \|C_{x_0} \iota\|_2^2 \|C_{a_0 J}\|_2^2 \|\Delta\|_2 \leq C_6 \tilde{\mu} \kappa_I n \theta \\ \left\| \iota^* C_{x_0}^* C_{a_0 J} (C_{aJ}^* C_{aJ})^{-1} C_{a-a_0 J}^* C_{x_0} \iota \right\|_2 &\leq \|C_{x_0} \iota\|_2^2 \|C_{a_0 J}\|_2 \|(C_{aJ}^* C_{aJ})^{-1}\|_2 \|C_{a_0-a} \iota\|_2 \leq C_7 \tilde{\mu} \sqrt{\kappa_I} n \theta \\ \left\| P_{I \setminus J} C_{a_0}^* C_{x_0} \iota \right\|_2^2 &\leq |I \setminus J| \|\check{a}_0 * x_0\|_{\square}^2 \leq C_8 \tilde{\lambda} n \theta \times \kappa_I \leq C_8 \left(\lambda \kappa_I + \frac{\kappa_I \log n}{\sqrt{n \theta^2}} \right) n \theta, \end{aligned} \quad (\text{H.77})$$

whence we conclude, that since $c_\mu \kappa_I^2 \leq c_\mu$ and $\lambda \kappa_I \leq 5c_\mu$, as long as $c_\mu < \frac{1}{100} \left(\frac{1}{C_6} + \frac{1}{C_7} + \frac{1}{5C_8} \right)$ and $n > 10^6 C_8^2 \theta^{-2} \kappa_I^2 \log^2 n$, we gain:

$$\begin{aligned} & \left\| \iota^* C_{x_0}^* C_{a_0 J} (C_{aJ}^* C_{aJ})^{-1} C_{aJ}^* C_{a_0-a} x_0 - \iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota (a_0 - a) \right\|_2 \\ & \leq \left(\frac{3}{100} + \frac{1}{1000} \right) n \theta \|a_0 - a\|_2 \\ & \leq \frac{1}{32} n \theta \|a_0 - a\|_2. \end{aligned} \quad (\text{H.78})$$

The claimed result therefore is followed by combining (H.70), (H.75) and (H.78). \blacksquare

Contraction of least square estimate of a toward a_0 . The next thing is to show the operator

$$(n\theta)^{-1} \left(\iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota \right) \quad (\text{H.79})$$

contracts a toward a_0 . We first will show that

$$(n\theta)^{-1} \left(\iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota \right) \approx a_0 a_0^* \quad (\text{H.80})$$

by seeing $\iota^* C_{x_0}^* P_I C_{x_0} \iota \approx (n\theta) e_0 e_0^*$ via sparsity of x_0 . Finally since the local perturbation on sphere is close to a quadratic function in ℓ^2 -norm of difference, we have

$$|\langle a_0, a - a_0 \rangle| \leq \frac{1}{2} \|a - a_0\|_2^2. \quad (\text{H.81})$$

Again, we introduce the following lemma to solidify our claim:

Lemma H.3 (Contraction of a to a_0). *Given $a_0 \in \mathbb{R}^{p_0}$ to be $\tilde{\mu}$ -shift coherent and $x_0 \sim \text{BG}(\theta) \in \mathbb{R}^n$. There exists some constants C, C', c, c', c_μ such that if $\lambda < c' \tilde{\mu} \kappa_I$, $\tilde{\mu} \kappa_I^2 \leq c_\mu$ and $n > C \theta^{-2} p^2 \log p$, then with probability at least $1 - c/n$, for every $\|a - a_0\|_2 \leq \tilde{\mu}$,*

$$\left\| \iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota (a_0 - a) \right\|_2 \leq \frac{1}{32} \|a - a_0\|_2 n \theta. \quad (\text{H.82})$$

Proof. Since $\mathbb{E} \langle P_I s_i[x_0], s_j[x_0] \rangle = 0$ for all $i \neq j$ and set I , we calculate

$$\begin{aligned} \mathbb{E} [\iota_{[\pm p]}^* C_{x_0}^* P_I C_{x_0} \iota_{[\pm p]}] &= \sum_{i \in [\pm p]} \mathbb{E} [e_i^* C_{x_0}^* P_I C_{x_0} e_i] e_i e_i^* = \mathbb{E} \|x_0\|_2^2 e_0 e_0^* + \sum_{i \in [\pm p] \setminus 0} \mathbb{E} \|P_I s_i[x_0]\|_2^2 e_i e_i^* \\ &= n\theta e_0 e_0^* + n\theta^2 P_{[\pm p] \setminus 0} = n\theta^2 I + n\theta(1-\theta) e_0 e_0^*. \end{aligned} \quad (\text{H.83})$$

whence

$$\mathbb{E} [\iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota] = \iota^* C_{a_0}^* \mathbb{E} [C_{x_0}^* P_I C_{x_0}] C_{a_0} \iota = n\theta^2 \iota^* C_{a_0}^* C_{a_0} \iota + n\theta(1-\theta) a_0 a_0^*, \quad (\text{H.84})$$

implying the expectation is a contraction mapping for $a_0 - a$ when $c_\mu < \frac{1}{200}$:

$$\begin{aligned} \|\mathbb{E} [\iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota] (a_0 - a)\|_2 &\leq n\theta^2 \|\iota^* C_{a_0}^* C_{a_0} \iota\|_2 \|a_0 - a\|_2 + n\theta \|a_0\|_2 |\langle a_0, a_0 - a \rangle| \\ &\leq n\theta^2 \times 2\tilde{\mu}p \times \|a_0 - a\|_2 + \frac{1}{2}n\theta \|a_0 - a\|_2^2 \\ &\leq (2c_\mu + \frac{1}{2}c_\mu) \|a_0 - a\|_2 n\theta \\ &\leq \frac{1}{64} \|a_0 - a\|_2 n\theta. \end{aligned} \quad (\text{H.85})$$

For each entry of $C_{x_0}^* P_I C_{x_0}$, again from [Appendix A](#) we know with probability at least $1 - c/n$:

$$|e_i^* C_{x_0}^* P_I C_{x_0} e_j - \mathbb{E} [e_i^* C_{x_0}^* P_I C_{x_0} e_j]| \leq \begin{cases} C' \sqrt{n\theta \log n} & i = j = 0 \\ C' \sqrt{n\theta^2 \log n} & \text{otherwise} \end{cases}.$$

Thus via Gershgorin disc theorem, when $n > 10^3 C'^2 \theta^{-2} p^2 \log n$:

$$\lambda_{\max} \left(\iota_{[\pm p]}^* C_{x_0}^* P_I C_{x_0} \iota_{[\pm p]} - \mathbb{E} [\iota_{[\pm p]}^* C_{x_0}^* P_I C_{x_0} \iota_{[\pm p]}] \right) \leq C' p \sqrt{n\theta^2 \log n} \leq \frac{1}{64} n\theta^2. \quad (\text{H.86})$$

Finally we combine (H.85), (H.86) and get

$$\|\iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota (a_0 - a)\|_2 \leq \left(\frac{1}{64} n\theta + \frac{1}{64} n\theta^2 \|C_{a_0} \iota_{\pm p}\|_2^2 \right) \|a_0 - a\|_2 \leq \frac{1}{32} \|a_0 - a\|_2 n\theta. \quad (\text{H.87})$$

■

[Lemma H.1-H.3](#) together implies the single iterate of alternating minimization contracts a toward a_0 . We show it with the following lemma:

Lemma H.4 (Contraction of least square estimate). *Given $a_0 \in \mathbb{R}^{p_0}$ to be $\tilde{\mu}$ -shift coherent and $x_0 \sim \text{BG}(\theta) \in \mathbb{R}^n$. There exists some constants C, C', c, c_μ such that if $\tilde{\mu}\kappa_I^2 \leq c_\mu$ and $n > C\theta^{-2}p^2 \log n$, then with probability at least $1 - c/n$, for every λ and a satisfying*

$$5\tilde{\mu}\kappa_I \geq \lambda \geq 5\kappa_I \|a - a_0\|_2, \quad (\text{H.88})$$

and suppose x^+ has the form of (H.16), then the solution a^+ to

$$\min_{a' \in \mathbb{R}^p} \left\{ \|a' * x^+ - y\|_2^2 \right\} \quad (\text{H.89})$$

is unique and satisfies

$$\|P_{\mathbb{S}^{p-1}} [a^+] - a_0\|_2 \leq \frac{1}{2} \|a - a_0\|_2. \quad (\text{H.90})$$

Proof. Write x as x^+ , then

$$\begin{aligned} \lambda_p (\iota^* C_x^* C_x \iota) &= \sigma_{\min}^2 (C_{x_0} \iota + C_{x-x_0} \iota) \\ &\geq \left[\sigma_{\min}(C_{x_0} \iota) - \|C_{x-x_0} \iota\| \right]_+^2 \end{aligned}$$

$$\begin{aligned}
&\geq \left[\sigma_{\min}(\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}) - 2\sqrt{\kappa_I} \|\mathbf{x} - \mathbf{x}_0\|_2 \right]_+^2 \\
&\geq \left[\frac{2}{3}\sqrt{\theta n} - 8\lambda\sqrt{\kappa_I}\sqrt{\theta n} \right]_+^2 \\
&\geq \frac{1}{2}\theta n,
\end{aligned} \tag{H.91}$$

where the fourth inequality is derived from using the upper bound of sparse convolution matrix from [Remark A.6](#), and the last line holds by knowing $\lambda < 5c_\mu\kappa_I^{-1}$. From (H.91) we know the least square problem of (H.89) has unique solution \mathbf{a}^+ , written as

$$\mathbf{a}^+ = (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{y}, \tag{H.92}$$

whence

$$\mathbf{a}^+ - \mathbf{a}_0 = (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}) \mathbf{a}_0 - \mathbf{a}_0 = (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}) \mathbf{a}_0. \tag{H.93}$$

Combine [Lemma H.2](#) and [Lemma H.3](#), we know

$$\|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}\|_2 \leq \left(C_1 \lambda (\tilde{\lambda} + \tilde{\mu} \kappa_I) + \frac{1}{16} \|\mathbf{a} - \mathbf{a}_0\|_2 \right) n\theta \tag{H.94}$$

for some constant C_1 . Combine (H.91), (H.93), (H.94) and since $\lambda < \tilde{\mu} \kappa_I$, by letting $c_\mu < \frac{1}{4C_1}$, we gain

$$\|\mathbf{a}^+ - \mathbf{a}_0\|_2 \leq \frac{\|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}\|_2}{\lambda_p(\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x}} \boldsymbol{\iota})} \leq 2C_1 \lambda (\tilde{\lambda} + \tilde{\mu} \kappa_I) + \frac{1}{8} \|\mathbf{a} - \mathbf{a}_0\|_2 \leq \frac{1}{4}. \tag{H.95}$$

For the final bound,

$$\begin{aligned}
\left\| \frac{\mathbf{a}^+}{\|\mathbf{a}^+\|_2} - \mathbf{a}_0 \right\|_2 &\leq \frac{\|\mathbf{a}^+ - \mathbf{a}_0\|_2 + \|\mathbf{a}^+\|_2 - 1}{\|\mathbf{a}^+\|_2} \leq \frac{2\|\mathbf{a}^+ - \mathbf{a}_0\|_2}{1 - \|\mathbf{a}^+ - \mathbf{a}_0\|_2} \leq \frac{8}{3} \|\mathbf{a}^+ - \mathbf{a}_0\|_2, \\
&\leq C_2 \lambda (\tilde{\lambda} + \tilde{\mu} \kappa_I) + \frac{1}{3} \|\mathbf{a} - \mathbf{a}_0\|_2,
\end{aligned} \tag{H.96}$$

and since $\lambda > \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2$, finally we gain

$$\begin{aligned}
\text{(H.96)} &\leq C_2 \left(\lambda \kappa_I + \frac{p \kappa_I \log n}{n\theta} + \tilde{\mu} \kappa_I^2 \right) \|\mathbf{a} - \mathbf{a}_0\|_2 + \frac{1}{3} \|\mathbf{a} - \mathbf{a}_0\|_2 \\
&\leq \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2
\end{aligned} \tag{H.97}$$

as long as $n > 20C_2\theta^{-1}p\kappa_I \log n$ and $c_\mu < \frac{1}{20C_2}$. ■

H.3 Linear convergence of alternating minimization (Proof of [Theorem 5.2](#))

In the first two sections we have shown the iterate contract \mathbf{a} toward \mathbf{a}_0 , under our signal assumption. We tie up these result by showing the following theorem which proves that the iterates produced by alternating minimization converge linearly to \mathbf{a}_0 :

Proof. We will prove our claim by induction on k . Clearly, when $k = 0$, we have $5\kappa_I \|\mathbf{a}^{(0)} - \mathbf{a}_0\|_2 \leq \lambda^{(0)} = 5\tilde{\mu}\kappa_I$ and $I^{(0)} = \{i : |s_i[\mathbf{a}^{(0)}]^* \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0| > \lambda^{(0)}\}$. Then for all $|\mathbf{x}_j| > 6\lambda^{(0)}$, we have

$$\begin{aligned}
|s_j[\mathbf{a}^{(0)}]^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0| &\geq \left(1 - |\langle \mathbf{a}^{(0)} | \mathbf{a}_0 \rangle|\right) |\mathbf{x}_j| - \left\| \mathbf{P}_{[\pm p] \setminus \{j\}} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} s_j[\mathbf{a}^{(0)}] \right\|_2 \times \sqrt{2} \|\mathbf{x}_0\|_\square \\
&\geq (1 - 2\tilde{\mu}) 6\lambda^{(0)} - 2\tilde{\mu}\sqrt{\kappa_I} \times \sqrt{2\kappa_I}
\end{aligned}$$

$$\begin{aligned}
&\geq 5\lambda^{(0)} - 4\lambda^{(0)} \\
&= \lambda^{(0)}.
\end{aligned} \tag{H.98}$$

hence $I_{>6\lambda^{(0)}} \subseteq I^{(0)}$, therefore the condition of [Lemma H.4](#) is satisfied, implies (5.32) holds for $k = 0$.

Suppose it is true for $1, 2, \dots, k-1$, such that

$$\kappa_I \|\mathbf{a}^{(k)} - \mathbf{a}_0\|_2 \leq \frac{1}{2} \lambda^{(k-1)} = \lambda^{(k)}, \quad \text{and} \quad I_{>3\lambda^{(k-1)}} \subseteq I^{(k)} \tag{H.99}$$

and since $I_{>6\lambda^{(k)}} = I_{>3\lambda^{(k-1)}} \subseteq I^{(k)}$, we can again apply [Lemma H.4](#), resulting

$$\kappa_I \|\mathbf{a}^{(k+1)} - \mathbf{a}\|_2 \leq \frac{1}{2} \kappa_I \|\mathbf{a}^{(k)} - \mathbf{a}_0\|_2 \leq \frac{1}{2} \lambda^{(k)} \tag{H.100}$$

as claimed. \blacksquare

H.4 Supporting lemmas for refinement

The following lemma controls the shift coherence of \mathbf{a} :

Lemma H.5 (Coherence of \mathbf{a} near \mathbf{a}_0). *Suppose that \mathbf{a}_0 is $\tilde{\mu}$ -shift coherent, and $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$. Then*

$$\|\text{off}[C_{\mathbf{a}}^* C_{\mathbf{a}_0}]\|_{\infty} \leq 2\tilde{\mu} \tag{H.101}$$

$$\|\text{off}[C_{\mathbf{a}}^* C_{\mathbf{a}}]\|_{\infty} \leq 3\tilde{\mu} \tag{H.102}$$

Proof. Notice that for any $\ell \neq 0$, $|\langle \mathbf{a}, s_{\ell}[\mathbf{a}_0] \rangle| \leq |\langle \mathbf{a}_0, s_{\ell}[\mathbf{a}_0] \rangle| + |\langle \mathbf{a} - \mathbf{a}_0, s_{\ell}[\mathbf{a}_0] \rangle| \leq \tilde{\mu} + \|\mathbf{a}_0 - \mathbf{a}\|_2 \leq 2\tilde{\mu}$. Similarly, $|\langle \mathbf{a}, s_{\ell}[\mathbf{a}] \rangle| \leq |\langle \mathbf{a} - \mathbf{a}_0, s_{\ell}[\mathbf{a}_0] \rangle| + |\langle \mathbf{a}, s_{\ell}[\mathbf{a}_0] \rangle| \leq \|\mathbf{a} - \mathbf{a}_0\|_2 + 2\tilde{\mu} \leq 3\tilde{\mu}$, as claimed. \blacksquare

From this we obtain the following spectral control on $C_{\mathbf{a}}^* C_{\mathbf{a}}$, to simply the notations, we will write

$$C_{\mathbf{a}I}^* C_{\mathbf{a}I} = \iota_I^* C_{\mathbf{a}}^* C_{\mathbf{a}} \iota_I = [C_{\mathbf{a}}^* C_{\mathbf{a}}]_{I,I} \tag{H.103}$$

in the latter part of this section.

Lemma H.6 (Off-diagonals of $[C_{\mathbf{a}}^* C_{\mathbf{a}}]_{I,I}$). *Suppose that \mathbf{a}_0 is $\tilde{\mu}$ -shift coherent and $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$. Then*

$$\left\| [C_{\mathbf{a}}^* C_{\mathbf{a}} - I]_{I,I} \right\|_2 \leq 9\kappa_I \tilde{\mu}. \tag{H.104}$$

We prove this lemma by noting that $C_{\mathbf{a}}^* C_{\mathbf{a}} = C_{\mathbf{r}_{\mathbf{a},\mathbf{a}}}$ is the convolution matrix associated with the autocorrelation $\mathbf{r}_{\mathbf{a},\mathbf{a}}$ of \mathbf{a} . Since $\text{supp}(\mathbf{r}_{\mathbf{a},\mathbf{a}}) \subseteq \{-p+1, \dots, p-1\}$ is confined to a (cyclic) stripe of width $2p-1$, we can tightly control the norm of this matrix by dividing it into three block-diagonal submatrices with blocks of size $p \times p$. Formally:

Proof. Divide I into $r = \lceil n/p \rceil$ subsets I_0, \dots, I_{r-1} such that for all $\ell = 0, \dots, r-1$:

$$I_{\ell} = I \cap \{p\ell, p\ell+1, \dots, p\ell+(p-1)\} = I \cap ([p] + p\ell).$$

Notice that for each ℓ :

$$\text{supp}([C_{\mathbf{a}}^* C_{\mathbf{a}}]_{I_{\ell}, I}) \subseteq I_{\ell} \times (I_{\ell-1} \uplus I_{\ell} \uplus I_{\ell+1}),$$

where $\ell+1$ and $\ell-1$ are interpreted cyclically modulo r .

For an arbitrary $\mathbf{v} \in \mathbb{R}^{|I|}$, we calculate

$$\left\| [C_{\mathbf{a}}^* C_{\mathbf{a}} - I]_{I,I} \mathbf{v} \right\|_2^2 = \sum_{\ell=0}^{r-1} \left\| [C_{\mathbf{a}}^* C_{\mathbf{a}} - I]_{I_{\ell}, I} \mathbf{v} \right\|_2^2 \tag{H.105}$$

$$= \sum_{\ell=0}^{r-1} \left\| [C_a^* C_a - I]_{I_\ell, I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} v_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_2^2 \quad (\text{H.106})$$

$$\leq \sum_{\ell=0}^{r-1} \left\| [C_a^* C_a - I]_{I_\ell, I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_F^2 \|v_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}}\|_2^2 \quad (\text{H.107})$$

$$\leq 3\kappa_I^2 \times (3\tilde{\mu})^2 \times \sum_{\ell=0}^{r-1} \|v_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}}\|_2^2 \quad (\text{H.108})$$

$$\leq 3\kappa_I^2 \times 9\tilde{\mu}^2 \times 3 \|v\|_2^2, \quad (\text{H.109})$$

giving the claimed result. \blacksquare

As a consequence, we have that

Corollary H.7 (Inverse of $[C_a^* C_a]_{J,J}$). *Suppose that \mathbf{a}_0 is μ -shift coherent, that $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$ and that $\kappa_I \tilde{\mu} < \frac{1}{18}$. Then for every $J \subseteq I$ and any norm $\|\cdot\|_\diamond \in \{\|\cdot\|_{\square \rightarrow \square}, \|\cdot\|_{\infty \rightarrow \infty}, \|\cdot\|_2\}$, we have*

$$\left\| [C_a^* C_a - I]_{J,J} \right\|_\diamond \leq 9\kappa_I \tilde{\mu} \quad (\text{H.110})$$

$$\left\| [C_a^* C_a]_{J,J}^{-1} - I \right\|_\diamond \leq 18\kappa_I \tilde{\mu} \quad (\text{H.111})$$

$$\left\| [C_a^* C_a]_{J,J}^{-1} \right\|_\diamond \leq 2. \quad (\text{H.112})$$

Proof. First we prove

$$\left\| [C_a^* C_a - I]_{J,J} \right\|_2 \leq 9\kappa_I \tilde{\mu}, \quad \left\| [C_a^* C_a - I]_{J,J} \right\|_{\infty \rightarrow \infty} \leq 6\kappa_I \tilde{\mu}, \quad \left\| [C_a^* C_a - I]_{J,J} \right\|_{\square \rightarrow \square} \leq 6\kappa_I \tilde{\mu} \quad (\text{H.113})$$

Where the first claim follows from [Lemma H.6](#). The second follows by noting that the ℓ^∞ operator norm is the maximum row ℓ^1 norm, and that each row has at most $2\kappa_I$ entries, of size at most $3\tilde{\mu}$. The last follows by noting that

$$\begin{aligned} \left\| [C_a^* C_a - I]_{J,J} \right\|_{\square \rightarrow \square} &\leq \max_{\ell, \ell'} \left\| [C_a^* C_a - I]_{J \cap ([p] + \ell), J \cap ([2p] + \ell')} \right\|_F \\ &\leq 6\kappa_I \tilde{\mu}. \end{aligned} \quad (\text{H.114})$$

Then we prove

$$\left\| [C_a^* C_a]_{J,J}^{-1} - I \right\|_2 \leq 18\kappa_I \tilde{\mu}, \quad \left\| [C_a^* C_a]_{J,J}^{-1} - I \right\|_{\infty \rightarrow \infty} \leq 12\kappa_I \tilde{\mu}, \quad \left\| [C_a^* C_a]_{J,J}^{-1} - I \right\|_{\square \rightarrow \square} \leq 12\kappa_I \tilde{\mu}, \quad (\text{H.115})$$

which are followed from the fact that if $\|\cdot\|_\diamond$ is a matrix norm and $\|\Delta\|_\diamond < 1$, then

$$\|(I + \Delta)^{-1} - I\|_\diamond \leq \frac{\|\Delta\|_\diamond}{1 - \|\Delta\|_\diamond}.$$

Finally, [\(H.112\)](#) follows from the triangle inequality. \blacksquare

Also, we need to bound the convolution of $\mathbf{a}_0 - \mathbf{a}$ with $\|\mathbf{a}_0 - \mathbf{a}\|_2$ requiring for bounds of the lasso solution:

Lemma H.8 (Convolution of $\mathbf{a}_0 - \mathbf{a}$). *Suppose that \mathbf{a}_0 is μ -shift coherent and $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$, then for every $J \subseteq I$,*

$$\|[C_a^* C_{\mathbf{a}_0 - \mathbf{a}}]_{J,J}\|_{\square \rightarrow \infty} \leq \sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 \quad (\text{H.116})$$

$$\|[C_a^* C_{\mathbf{a}_0 - \mathbf{a}}]_{J,J}\|_{\square \rightarrow \square} \leq \sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 \quad (\text{H.117})$$

Proof. For the first inequality, we have

$$\begin{aligned}
\| [C_{\mathbf{a}}^* C_{\mathbf{a}_0 - \mathbf{a}}]_{J, J} \mathbf{v} \|_{\square \rightarrow \infty} &= \max_{j \in J, \|\mathbf{v}\|_{\square} = 1} |\langle s_j[\mathbf{a}], (\mathbf{a}_0 - \mathbf{a}) * \mathbf{v} \rangle| \\
&\leq \max_{j \in [n], \|\mathbf{v}\|_{\square} = 1} \| \mathbf{P}_{[p]+j} [(\mathbf{a}_0 - \mathbf{a}) * \mathbf{v}] \|_2 \\
&\leq \| \mathbf{a} - \mathbf{a}_0 \|_2 \times \max_{j \in [n], \|\mathbf{v}\|_{\square} = 1} \| \mathbf{P}_{[\pm p]+j} \mathbf{v} \|_1 \\
&\leq \sqrt{2\kappa_I} \| \mathbf{a}_0 - \mathbf{a} \|_2
\end{aligned} \tag{H.118}$$

The second inequality is derived by

$$\begin{aligned}
\| [C_{\mathbf{a}}^* C_{\mathbf{a}_0 - \mathbf{a}}]_{J, J} \|_{\square \rightarrow \square} &\leq \max_{\ell, \ell'} \| [C_{\mathbf{a}}^* C_{\mathbf{a}_0 - \mathbf{a}}]_{J \cap ([p] + \ell), J \cap ([2p] + \ell')} \|_F \\
&\leq \sqrt{2\kappa_I^2 \max_{i,j} |\langle s_i[\mathbf{a}], s_j[\mathbf{a}_0 - \mathbf{a}] \rangle|^2} \\
&\leq \sqrt{2\kappa_I} \| \mathbf{a} - \mathbf{a}_0 \|_2,
\end{aligned} \tag{H.119}$$

finishing the proof. \blacksquare

Again, using a variant of the argument for [Lemma H.6](#), we have the following:

Lemma H.9 (Off-diagonal of submatrix of $C_{\mathbf{a}}^* C_{\mathbf{a}_0}$). *Suppose that \mathbf{a}_0 is μ -shift coherent and $\| \mathbf{a} - \mathbf{a}_0 \|_2 \leq \tilde{\mu}$. For any $J \subset I$, if*

$$\kappa_J = \max_{\ell} |J \cap \{\ell, \ell + 1, \dots, \ell + p - 1\}| \tag{H.120}$$

$$\kappa_{I \setminus J} = \max_{\ell} |(I \setminus J) \cap \{\ell, \ell + 1, \dots, \ell + p - 1\}| \tag{H.121}$$

Then

$$\| [C_{\mathbf{a}}^* C_{\mathbf{a}_0}]_{J, I \setminus J} \|_2 \leq 6\sqrt{\kappa_J \kappa_{I \setminus J}} \tilde{\mu}. \tag{H.122}$$

Proof. Take $r = \lceil n/p \rceil$ and for $\ell = 0, \dots, r - 1$, write

$$J_{\ell} = J \cap ([p] + p\ell), \quad L_{\ell} = (I \setminus J) \cap ([p] + p\ell),$$

Take $\mathbf{v} \in \mathbb{R}^{|I \setminus J|}$ arbitrary and notice that

$$\begin{aligned}
\| [C_{\mathbf{a}}^* C_{\mathbf{a}_0}]_{J, I \setminus J} \mathbf{v} \|_2^2 &= \sum_{\ell=0}^{r-1} \| [C_{\mathbf{a}}^* C_{\mathbf{a}_0}]_{J_{\ell}, I \setminus J} \mathbf{v} \|_2^2 \\
&= \sum_{\ell=0}^{r-1} \| [C_{\mathbf{a}}^* C_{\mathbf{a}_0}]_{J_{\ell}, L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \mathbf{v}_{L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \|_2^2 \\
&\leq 4\tilde{\mu}^2 \times \kappa_J \times 3\kappa_{I \setminus J} \times \sum_{\ell=0}^{r-1} \| \mathbf{v}_{L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \|_2^2 \\
&\leq 4\tilde{\mu}^2 \times \kappa_J \times 3\kappa_{I \setminus J} \times 3\| \mathbf{v} \|_2^2,
\end{aligned} \tag{H.123}$$

giving the result. \blacksquare

Lemma H.10 (Perturbation of vector over sphere). *If both \mathbf{a}, \mathbf{a}_0 are unit vectors in inner product space, then*

$$|\langle \mathbf{a}, \mathbf{a} - \mathbf{a}_0 \rangle| \leq \frac{1}{2} \| \mathbf{a} - \mathbf{a}_0 \|_2^2. \tag{H.124}$$

Proof. Via simple norm inequalities:

$$\frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2^2 = 1 - \langle \mathbf{a}, \mathbf{a}_0 \rangle = 1 - \langle \mathbf{a}, \mathbf{a}_0 - \mathbf{a} + \mathbf{a} \rangle = \langle \mathbf{a}, \mathbf{a} - \mathbf{a}_0 \rangle > 0 \quad (\text{H.125})$$

■

Lemma H.11 (Convolution of short and sparse). *Suppose $\delta \in \mathbb{R}^p$, and $\mathbf{v} \in \mathbb{R}^n$ where $\text{supp}(\mathbf{v}) = I$ satisfies*

$$\max_{\ell \in [n]} |I \cap ([p] + \ell)| \leq \kappa \quad (\text{H.126})$$

then

$$\|\delta * \mathbf{v}\|_2 \leq \sqrt{2\kappa} \|\delta\|_2 \|\mathbf{v}\|_2 \quad (\text{H.127})$$

Proof. Since every p -contiguous segment of I has at most κ elements, by splitting $I = I_1 \uplus I_2 \uplus \dots \uplus I_\kappa \uplus R$ such that each sets I_i are p -separated:

$$\begin{aligned} I_1 &= \{i_1, i_{\kappa+1}, i_{2\kappa+1}, \dots\} \cap \{0, \dots, n-p-1\}, \\ I_2 &= \{i_2, i_{\kappa+2}, i_{2\kappa+2}, \dots\} \cap \{0, \dots, n-p-1\}, \\ &\vdots \\ I_\kappa &= \{i_\kappa, i_{2\kappa}, i_{3\kappa}, \dots\} \cap \{0, \dots, n-p-1\}, \end{aligned} \quad (\text{H.128})$$

$$R = I \cap \{n-p, \dots, n-1\}. \quad (\text{H.129})$$

Then the p -separating property gives $\|\delta * P_{I_i} \mathbf{v}\|_2 = \|\delta\|_2 \|P_{I_i} \mathbf{v}\|_2$. Hence:

$$\begin{aligned} \|\delta * P_I \mathbf{v}\|_2 &= \left\| \sum_{i \in \kappa} \delta * P_{I_i} \mathbf{v} + \delta * P_R \mathbf{v} \right\|_2 \leq \sum_{i \in \kappa} \|\delta * P_{I_i} \mathbf{v}\|_2 + \|\delta * P_R \mathbf{v}\|_2 \\ &= \|\delta\|_2 \sum_{i \in \kappa} \|P_{I_i} \mathbf{v}\|_2 + \|\delta\|_2 \|P_R \mathbf{v}\|_1 \\ &\leq \sqrt{\kappa} \|\mathbf{v}_{I_1 \uplus \dots \uplus I_\kappa}\|_2 \|\delta\|_2 + \sqrt{\kappa} \|\mathbf{v}_R\|_2 \|\delta\|_2 \\ &\leq \sqrt{2\kappa} \|\mathbf{v}\|_2 \|\delta\|_2, \end{aligned} \quad (\text{H.130})$$

where the last two inequalities were coming from Cauchy-Schwartz. ■

I Finite sample approximation

In this section we collect several major components of proof about large sample deviation. In particular, the concentration for shift space gradient $\chi(\beta)_i$, shift space Hessian diagonals $\|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2$, and the set of gradients discontinuity entries $|J_B(\mathbf{a})|$.

I.1 Proof of Corollary C.4

Proof. 1. (ε -net) Write \mathbf{x} as \mathbf{x}_0 and $\|\beta\|_2 = \eta$ through out this proof, firstly from Definition B.1 for every $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, we know $\eta \leq 1 + c_\mu + \frac{c_\mu}{\sqrt{\theta k \log \theta^{-1}}} \leq \sqrt{p}$. Define $\varepsilon = \frac{c_2}{2n^{3/2}p^{3/2}}$ and consider the ε -net \mathcal{N}_ε for sphere of radius η . From Lemma J.5 we know for any $c_2 < 1$:

$$|\mathcal{N}_\varepsilon| \leq \left(\frac{3\eta}{\varepsilon}\right)^{2p} \leq \left(\frac{3n^{3/2}p^2}{c_2}\right)^{2p} \leq \left(\frac{3np^2}{c_2}\right)^{3p} \quad (\text{I.1})$$

for each $i \in [n]$ define such net as $\mathcal{N}_{\varepsilon,i}$, and define an event such that all center of subsets in $\mathcal{N}_{\varepsilon,i}$ are being well-behaved:

$$\mathcal{E}_{\text{Net}} := \left\{ \forall i \in [n], \quad \sigma_i n^{-1} \chi[\beta_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta_\varepsilon]}_i < \frac{c_1 \theta}{p^{3/2}} \quad \forall \beta_\varepsilon \in \mathcal{N}_{\varepsilon,i}, \right\} \quad (\text{I.2})$$

2. (Lipschitz constant) The Lipschitz constant L of $\chi[\cdot]_i$ w.r.t β is bounded in terms of \mathbf{x} regardless of entry i :

$$\begin{aligned} |\chi[\beta]_i - \chi[\beta']_i| &\leq \left| \mathbf{e}_i^* \tilde{\mathbf{C}}_{\mathbf{x}} \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta] - \mathbf{e}_i^* \tilde{\mathbf{C}}_{\mathbf{x}} \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta'] \right| \leq \|\mathbf{x}\|_2 \left\| \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta] - \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta'] \right\|_2 \\ &\leq \|\mathbf{x}\|_2 \sqrt{\sum_{j \in [n]} \left| \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta]_j - \mathcal{S}_\lambda [\tilde{\mathbf{C}}_{\mathbf{x}} \beta']_j \right|^2} \leq \|\mathbf{x}\|_2 \left\| \tilde{\mathbf{C}}_{\mathbf{x}} \beta - \tilde{\mathbf{C}}_{\mathbf{x}} \beta' \right\|_2 \\ &\leq \|\mathbf{x}\|_2 \cdot \|\mathbf{x}\|_1 \cdot \|\beta - \beta'\|_2 =: L \|\beta - \beta'\|_2 \end{aligned} \quad (\text{I.3})$$

Define the event that $\chi[\beta]_i$ that has small Lipschitz constant as

$$\mathcal{E}_{\text{Lip}} := \left\{ L < 2n^{3/2}\theta \right\} \quad (\text{I.4})$$

on the event \mathcal{E}_{Lip} , for every points in $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and $i \in [n]$, there exists some $\beta_\varepsilon \in \mathcal{N}_{\varepsilon,i}$ such that

$$\left| \left(\sigma_i n^{-1} \chi[\beta]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta]}_i \right) - \left(\sigma_i n^{-1} \chi[\beta_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta_\varepsilon]}_i \right) \right| \leq 2L\varepsilon \leq \frac{c_2 \theta}{p^{3/2}} \quad (\text{I.5})$$

On event $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$, (I.2), (I.5) implies $\chi[\beta]$ is well concentrated entrywise and anywhere in $\cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$:

$$\left| \sigma_i n^{-1} \chi[\beta]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta]}_i \right| \leq \frac{(c_1 + c_2)\theta}{p^{3/2}}, \quad \forall \mathbf{a} \in \cup_{k \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu)), \quad \forall i \in [n] \quad (\text{I.6})$$

as desired, where, using Lemma A.2,

$$\mathbb{P} [\mathcal{E}_{\text{Lip}}^c] \leq \mathbb{P} \left[\|\mathbf{x}\|_2^2 > 2n\theta \right] \leq 1/n; \quad (\text{I.7})$$

and using union bound,

$$\mathbb{P} [\mathcal{E}_{\text{Net}}^c] \leq \mathbb{P} \left[\max_{\substack{\mathbf{a}_\varepsilon \in \mathcal{N}_{\varepsilon,i} \\ i \in [n]}} \sigma_i n^{-1} \chi[\beta_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta_\varepsilon]}_i > \frac{c_1 \theta}{p^{3/2}} \right]$$

$$\leq n |\mathcal{N}_\varepsilon| \mathbb{P} \left[\sigma_0 n^{-1} \chi[\beta_\varepsilon]_0 - \sigma_0 n^{-1} \mathbb{E} \chi[\beta_\varepsilon]_0 > \frac{c_1 \theta}{p^{3/2}} \right]. \quad (\text{I.8})$$

3. (Bound $\mathbb{P}[\mathcal{E}_{\text{Net}}^c]$) Wlog write $n = t \cdot (2p)$ for some integer t and $2p \geq 4p_0 - 3$ and replace \mathbf{x}_0 with \mathbf{x} . Observe that $\mathbf{Z}_j(\beta)$ from (C.9) is independent of $\mathbf{Z}_{j+2p}(\beta)$ for all $j \in [n]$ while all \mathbf{Z}_j are identical distributed. We write $\chi[\beta]_0$ as sum of iid r.v.s. as

$$\chi[\beta]_0 = \sum_{j \in [n]} \mathbf{Z}_j(\beta) = \sum_{k \in [2p]} \left(\sum_{t=0}^{n/2p-1} \mathbf{Z}_{k+2tp}(\beta) \right)$$

wlog let $\sigma_0 = 1$ and split the independent r.v.s, write $\mathbb{E} \mathbf{Z}_0 = \mathbb{E} \mathbf{Z}$, bound the tail probability of $\chi[\beta]_0$ as

$$\mathbb{P} \left[n^{-1} \chi[\beta]_0 > n^{-1} \overline{\mathbb{E} \chi(\beta)}_0 + \frac{c_1 \theta}{p^{3/2}} \right] \leq 2p \cdot \mathbb{P} \left[\sum_{t=0}^{n/2p-1} \mathbf{Z}_{2tp}(\beta) > \frac{n}{2p} \mathbb{E} \mathbf{Z}(\beta) + \frac{c_1 n \theta}{2p^{5/2}} \right] \quad (\text{I.9})$$

The moments of \mathbf{Z}_0 can be bounded by using $|\mathbf{Z}_0(\beta)| \leq |\mathbf{x}_0| |\beta_0 \mathbf{x}_0 + \mathbf{s}_0| \leq \beta_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|$ where $\mathbf{s}_0 = \sum_{\ell \neq 0} \mathbf{x}_\ell \beta_\ell$, write $\mathbf{x} = \boldsymbol{\omega} \circ \mathbf{g} \sim_{\text{i.i.d.}} \text{BG}(\theta)$. For the 2-norm we know

$$\mathbb{E} |\mathbf{s}_0|^2 = \mathbb{E} \left| \sum_{\ell} \mathbf{x}_\ell \beta_\ell \right|^2 \leq \theta \|\beta\|_2^2 \leq \theta \left(1 + c_\mu + \frac{c_\mu}{\theta k^2} \right) \leq \frac{1}{2} \quad (\text{I.10})$$

As for the q -norm, use the moment generating function bound, such that for all $t \geq 0$:

$$\begin{aligned} \mathbb{E} |\mathbf{s}_0|^q &\leq q! t^{-q} \mathbb{E} \exp[t |\mathbf{s}_0|] \leq q! t^{-q} \prod_{\ell} \mathbb{E}_{\boldsymbol{\omega}_\ell, \mathbf{g}_\ell} \exp[t \boldsymbol{\omega}_\ell |\mathbf{g}_\ell| |\beta_\ell|] \leq 2q! t^{-q} \prod_{\ell} \mathbb{E}_{\boldsymbol{\omega}_\ell} \exp[\boldsymbol{\omega}_\ell t^2 \beta_\ell^2 / 2] \\ &\leq 2q! t^{-q} \prod_{\ell} (1 - \theta + \theta \exp[t^2 \beta_\ell^2 / 2]) \end{aligned} \quad (\text{I.11})$$

notice that the entrywise twice derivative of (I.11) w.r.t. β_ℓ^2 's are always positive, this function is convex for all β_ℓ^2 . Constrain on the polytope $\sum_{\ell} \beta_\ell^2 \leq \|\beta\|_2^2$, the maximizer of (I.11) w.r.t. β_ℓ^2 's occurs at a vertex point where $\beta_0^2 = \|\beta\|_2^2$. Thus

$$(\text{I.11}) \leq 2q! t^{-q} \left(1 - \theta + \theta \exp \left[t^2 \|\beta\|_2^2 / 2 \right] \right) \prod_{\ell \neq 0} (1 - \theta + \theta e^0) \leq 2q! t^{-q} (1 + \theta \exp[\|\beta\|_2^2 t^2 / 2]).$$

Choose $t = \sqrt{q} / \|\beta\|_2$, use $q!! > (q!/2) \cdot (e/q)^{q/2}$, we have

$$\mathbb{E} |\mathbf{s}_0|^q \leq 2q! q^{-q/2} \|\beta\|_2^q (1 + \theta \exp[q/2]) \leq 8 \|\beta\|_2^q \max \left\{ e^{-q/2}, \theta \right\} q!! \quad (\text{I.12})$$

Apply Jensen's inequality $\left(\sum_{i=1}^N z_i \right)^q \leq N^{q-1} \sum_{i=1}^N z_i^q$, use Gaussian moment Lemma J.2, (I.10) and (I.12), obtain for $q \geq 3$,

$$\begin{aligned} \mathbb{E} Z(\beta)^2 &\leq \mathbb{E} (\beta_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|)^2 \leq 2\mathbb{E} [\beta_0^2 \mathbf{x}_0^4 + \mathbf{x}_0^2 \mathbf{s}_0^2] \leq 6\theta + 2\theta^2 \|\beta\|_2^2 \leq 7\theta, \\ \mathbb{E} Z(\beta)^q &\leq \mathbb{E} (\beta_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|)^q \leq 2^{q-1} \left(\mathbb{E} \mathbf{x}_0^{2q} + \mathbb{E} |\mathbf{x}_0|^q \mathbb{E} |\mathbf{s}_0|^q \right) \\ &\leq \theta 2^{q-1} (2q-1)!! + \theta 2^{q-1} (q-1)!! \left(8 \|\beta\|_2^q \max \left\{ e^{-q/2}, \theta \right\} q!! \right) \\ &\leq \theta 4^q q! + \theta 2^q \|\beta\|_2^q q!. \end{aligned}$$

Thus, recall that $\|\beta\|_2 = \eta$, use $(\sigma^2, R) = (8\theta\eta^2, 4\eta)$, from (I.8)-(I.9), apply Bernstein inequality [Lemma J.4](#) with $n \geq Cp^5\theta^{-2} \log p$, and $c_1, c_2 \in [0, 1]$ we have

$$\begin{aligned} \mathbb{P}[\mathcal{E}_{\text{Net}}^c] &\leq 2np |\mathcal{N}_\varepsilon| \cdot \mathbb{P} \left[\sum_{t=0}^{n/2p-1} \mathbf{Z}_{2tp}(\beta) > \frac{n}{2p} \mathbb{E} \mathbf{Z}(\beta) + \frac{c_1 n \theta}{2p^{5/2}} \right] \leq 2np \left(\frac{3np^2}{c_2} \right)^{3p} \exp \left(\frac{-(c_1 n \theta / 2p^{5/2})^2}{16n\theta\eta^2/2p + 8\eta c_1 n \theta / 2p^{5/2}} \right) \\ &\leq \exp \left(4p \log \left(\frac{3np^2}{c_2} \right) - \frac{(c_1 n \theta / 2p^{5/2})^2}{16n\theta\eta^2/p} \right) \leq \exp \left(4p \log \left(\frac{3np^2}{c_2} \right) - \frac{c_1^2 n \theta^2}{64p^4} \right) \\ &\leq \exp \left(\frac{-c_1^2 n \theta^2}{100p^4} \right) \leq \frac{1}{n} \end{aligned} \quad (\text{I.13})$$

when $\frac{C}{\log C} > \frac{10^5}{c_1^2 c_2}$. The proof of lower bound and negative β_0 is derived in the same manner. \blacksquare

I.2 Proof of [Corollary D.3](#)

Proof. Write \mathbf{x} as \mathbf{x}_0 though our this proof. Write $\beta_i \mathbf{x}_j + \mathbf{s}_j = \sum_{\ell \in [\pm p]} \beta_\ell \mathbf{x}_{\ell-i+j} = \langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle$, and the support w.r.t. some \mathbf{a} as $I(\beta)$. Define the random variable $\mathbf{Z}_{ij}(\beta)$ as

$$\|P_{I(\beta)} \mathbf{s}_{-i}[\mathbf{x}]\|_2^2 = \sum_{j \in [n]} \mathbf{x}_j^2 \mathbf{1}_{\{|\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| > \lambda\}} =: \sum_{j \in [n]} \mathbf{Z}_{ij}(\beta) \quad (\text{I.14})$$

and define $\{\bar{\mathbf{Z}}_{ij}(\beta)\}_{j \in [n]}$ that are independent r.v.s. and as a upper bounding function of $\mathbf{Z}_{ij}(\beta)$ as

$$\bar{\mathbf{Z}}_{ij}(\beta) := \begin{cases} \mathbf{x}_j^2, & |\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| > \lambda \\ 0, & |\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| < \lambda/2, \\ \frac{\mathbf{x}_j^2}{\lambda/2} (|\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| - \lambda/2), & \text{otherwise} \end{cases} \quad (\text{I.15})$$

Similar to proof of [Corollary C.4](#). Let $\|\beta\|_2 \leq \eta \leq \sqrt{p}$. Define $\varepsilon = \frac{c'_2 \lambda}{24np\sqrt{p\theta \log n \log \theta^{-1}}}$ for some $c'_2 > 0$ and consider the ε -net \mathcal{N}_ε for sphere of radius η . From [Lemma J.5](#) we know

$$|\mathcal{N}_\varepsilon| \leq \left(\frac{3\eta}{\varepsilon} \right)^{2p} \leq \left(\frac{72}{c'_2 c_\lambda} np^2 \sqrt{\theta |\tau| \log n \log \theta^{-1}} \right)^{2p} \leq \left(\frac{72}{c'_2 c_\lambda} np^2 \log n \right)^{2p}, \quad (\text{I.16})$$

for each $i \in [n]$ define such net as $\mathcal{N}_{\varepsilon,i}$, and define an event such that all center of subsets in $\mathcal{N}_{\varepsilon,i}$ are being well-behaved:

$$\mathcal{E}_{\text{Net}} := \left\{ \forall i \in [n], \left| n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta_\varepsilon) - \mathbb{E} \bar{\mathbf{Z}}_i(\beta_\varepsilon) \right| \leq \frac{c'_1 \theta}{p} \quad \forall \beta_\varepsilon \in \mathcal{N}_{\varepsilon,i} \right\}, \quad (\text{I.17})$$

Also, $\sum_j \bar{\mathbf{Z}}_{ij}(\beta)$ is a Lipchitz function over β for every $i \in [n]$ as

$$\begin{aligned} \left| \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta) - \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta') \right| &\leq \sum_{j \in [n]} \frac{\mathbf{x}_j^2}{\lambda/2} |\langle \beta - \beta', \mathbf{x}_{[\pm p]-i+j} \rangle| \leq \sum_{j \in [n]} \frac{\mathbf{x}_j^2 \|\mathbf{x}_{[\pm p]-i+j}\|_2}{\lambda/2} \|\beta - \beta'\|_2, \\ &\leq \frac{1}{\lambda/2} \|\mathbf{x}\|_2^2 \cdot \max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 \cdot \|\beta - \beta'\|_2 := L \|\beta - \beta'\|_2, \end{aligned} \quad (\text{I.18})$$

and define event \mathcal{E}_{Lip} such that the Lipchitz constant is bounded as

$$\mathcal{E}_{\text{Lip}} := \left\{ L \leq 12n\theta \sqrt{p\theta \log n \log \theta^{-1}} \lambda^{-1} \right\}, \quad (\text{I.19})$$

then on event \mathcal{E}_{Lip} , for any points β in $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and $i \in [n]$, there exists some β_ε in $\mathcal{N}_{\varepsilon, i}$ with $\|\beta - \beta_\varepsilon\|_2 \leq \varepsilon$, and thus

$$\left| \left(n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta) - \mathbb{E} \bar{\mathbf{Z}}_i(\beta) \right) - \left(n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta_\varepsilon) - \mathbb{E} \bar{\mathbf{Z}}_i(\beta_\varepsilon) \right) \right| \leq 2L\varepsilon \leq \frac{c'_2 \theta}{p}. \quad (\text{I.20})$$

On event $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$, from (I.17), (I.20), we can conclude that for all $\beta \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and $i \in [n]$ that:

$$n^{-1} \|\mathbf{P}_{I(\beta)} s_{-i}[\mathbf{x}_0]\|_2^2 - n^{-1} \mathbb{E} \|\mathbf{P}_{I(\beta)} s_{-i}[\mathbf{x}_0]\|_2^2 \leq n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\beta) - \mathbb{E} \bar{\mathbf{Z}}_i(\beta) \leq \frac{(c'_1 + c'_2)\theta}{p} \quad (\text{I.21})$$

as desired, where the error probability of $\mathcal{E}_{\text{Lip}}^c$ is bounded using Lemma A.2 and Lemma A.3, which give

$$\mathbb{P}[\mathcal{E}_{\text{Lip}}^c] \leq \mathbb{P}[\|\mathbf{x}\|_2^2 > 2n\theta] + \mathbb{P}\left[\max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 > 3\sqrt{p\theta \log n \log \theta^{-1}}\right] \leq 3/n, \quad (\text{I.22})$$

when $n > 10^3 \theta^{-1}$. As for $\mathcal{E}_{\text{Net}}^c$ use union bound and split the r.v.s since $\mathbf{Z}_j, \mathbf{Z}_{j+2p}$ are independent for all j :

$$\mathbb{P}[\mathcal{E}_{\text{Net}}^c] \leq 2np \cdot |\mathcal{N}_\varepsilon| \cdot \mathbb{P}\left[\left|\sum_k^{n/2p} \bar{\mathbf{Z}}_{i,2kj}(\beta) - \frac{n}{2p} \mathbb{E} \bar{\mathbf{Z}}_i(\beta)\right| \geq \frac{c'_1 n \theta}{2p^2}\right].$$

Now we calculate the variance and L^q -norm of $\sum_k \bar{\mathbf{Z}}_{i,2kj}$ for $q \geq 3$:

$$\begin{cases} \mathbb{E} \bar{\mathbf{Z}}_{i,j}^2 \leq \mathbb{E} \mathbf{x}_j^4 \leq 3\theta \\ \mathbb{E} \bar{\mathbf{Z}}_{i,j}^q \leq \mathbb{E} \mathbf{x}_j^{2q} \leq \theta(2q-1)!! \leq \frac{1}{2} \cdot (3\theta) \cdot 2^{q-2} q! \end{cases} \quad (\text{I.23})$$

and apply Bernstein inequality with $(\sigma^2, R) = (3\theta, 2)$, then use $n \geq Cp^4 \theta^{-1} \log p$ and $c'_1, c'_2 < 1$ to obtain

$$\begin{aligned} 2np |\mathcal{N}_\varepsilon| \mathbb{P}\left[\left|\sum_k^{n/2p} \bar{\mathbf{Z}}_{i,2kj}(\beta) - \frac{n}{2p^2} \mathbb{E} \bar{\mathbf{Z}}_i\right| \geq \frac{c'_1 n \theta}{2p^2}\right] &\leq \exp\left[\log(2np) + 2p \log\left(\frac{72}{c'_2 c_\lambda} np^2 \log n\right) - \frac{(c'_1 n \theta / 2p^2)^2}{6n\theta/2p + 4c'_1 n \theta / 2p^2}\right] \\ &\leq \exp\left[3p \log\left(\frac{72}{c'_2 c_\lambda} np^2 \log n\right) - \frac{c'_1{}^2 n \theta}{24p^3}\right] \\ &\leq \exp[-c'_1{}^2 n \theta / (50p^3)] \leq 1/n, \end{aligned} \quad (\text{I.24})$$

where the last two inequalities holds when $\frac{C}{\log C} \geq \frac{10^5}{c'_1{}^2 c'_2 c_\lambda}$. The other side of inequality of (D.9) can be derived by defining $\underline{\mathbf{Z}}_{ij}$ as

$$\underline{\mathbf{Z}}_{ij}(\beta) := \begin{cases} \mathbf{x}_j^2, & |\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| > 3\lambda/2 \\ 0, & |\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| < \lambda \\ \frac{\mathbf{x}_j^2}{\lambda/2} (|\langle \beta, \mathbf{x}_{[\pm p]-i+j} \rangle| - \lambda), & \text{otherwise} \end{cases}, \quad (\text{I.25})$$

and define $\mathcal{E}_{\text{Net}}, \mathcal{E}_{\text{Lip}}$ similarly, such that on intersection of these events,

$$n^{-1} \|\mathbf{P}_{I(\beta)} s_{-i}[\mathbf{x}]\|_2^2 - n^{-1} \mathbb{E} \|\mathbf{P}_{I(\beta)} s_{-i}[\mathbf{x}]\|_2^2 \geq n^{-1} \sum_{j \in [n]} \underline{\mathbf{Z}}_{ij}(\beta) - \mathbb{E} \underline{\mathbf{Z}}_i(\beta) \geq \frac{(c'_1 + c'_2)\theta}{p} \quad (\text{I.26})$$

as desired. ■

I.3 Proof of Lemma E.5

Proof. 1. (Expectation upper bound) We will write \mathbf{x} as \mathbf{x}_0 . Similar to proof of Corollary C.4 let $\|\beta\|_2 \leq \eta \leq \sqrt{p}$. For each $i \in [n]$, define the random variable

$$\mathbf{X}_i(\beta) = \mathbf{1}_{\{|\langle s_i[\mathbf{x}], \beta \rangle - \lambda| \leq B\}} + \mathbf{1}_{\{|\langle s_i[\mathbf{x}], \beta \rangle + \lambda| \leq B\}}, \quad (\text{I.27})$$

then number of indices for vector $\mathbf{x} * \check{\beta}$ that are within B of $\pm \lambda$ is a random variable $\sum_{i \in [n]} \mathbf{X}_i(\beta)$. For each of the $\mathbf{X}_i(\beta)$'s consider an upper bound $\overline{\mathbf{X}}_i(\beta)$ defined as

$$\overline{\mathbf{X}}_i(\beta) = \begin{cases} \frac{1}{M} (|\langle s_i[\mathbf{x}], \beta \rangle| - (\lambda - B - M)) & |\langle s_i[\mathbf{x}], \beta \rangle| \in [\lambda - B - M, \lambda - B] \\ 1 & |\langle s_i[\mathbf{x}], \beta \rangle| \in [\lambda - B, \lambda + B] \\ \frac{1}{M} ((\lambda + B + M) - |\langle s_i[\mathbf{x}], \beta \rangle|) & |\langle s_i[\mathbf{x}], \beta \rangle| \in [\lambda + B, \lambda + B + M] \\ 0 & \text{else} \end{cases} \quad (\text{I.28})$$

where $B < M = c\lambda\theta^2 / (p \log n) \leq \lambda/4$ for some constant $0 < c < 1$.

Notice that $\mathbf{x} \sim_{\text{i.i.d.}} \text{BG}(\theta)$ is equal in distribution to $\mathbf{P}_{I(\mathbf{a})}\mathbf{g}$, where $\mathbf{g} \sim_{\text{i.i.d.}} \mathcal{N}(0, 1)$, and $I(\mathbf{a}) \subseteq [n]$ is an independent Bernoulli subset. Conditioned on $I(\mathbf{a})$, $\langle \mathbf{x}, \beta \rangle = \langle \mathbf{g}, \mathbf{P}_{I(\mathbf{a})}\beta \rangle \sim \mathcal{N}(0, \|\mathbf{P}_{I(\mathbf{a})}\beta\|_2^2)$. For all realizations of $I(\mathbf{a})$, the variance $\|\mathbf{P}_{I(\mathbf{a})}\beta\|_2^2$ is bounded by $\|\mathbf{P}_{I(\mathbf{a})}\beta\|_2^2 \leq \|\beta\|_2^2 \leq p$. Using these observations, and letting $f_\sigma(t) = (\sqrt{2\pi}\sigma)^{-1} \exp(-t^2/2\sigma^2)$ denote the pdf of an $\mathcal{N}(0, \sigma^2)$ random variable, the expectation of $\sum_i \overline{\mathbf{X}}_i(\beta)$ can be upper bounded as

$$\begin{aligned} \sum_{i \in [n]} \mathbb{E} [\overline{\mathbf{X}}_i(\beta)] &\leq (2n) \cdot \mathbb{P}[\langle \mathbf{x}, \beta \rangle \in [\lambda - B - M, \lambda + B + M]] \\ &\leq (2n) \cdot 2(B + M) \sup_{\sigma^2 \in (0, p]} \max_{t \in [\lambda - B - M, \lambda + B + M]} f_\sigma(t) \\ &\leq 4n(B + M) \sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda - B - M) \\ &\leq 4n(B + M) \sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda/2). \end{aligned} \quad (\text{I.29})$$

Notice that

$$\frac{d}{d\sigma} f_\sigma\left(\frac{\lambda}{2}\right) = \frac{d}{d\sigma} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\lambda^2}{8\sigma^2}\right) = \frac{\lambda^2 - 4\sigma^2}{4\sqrt{2\pi}\sigma^4} \exp\left(-\frac{\lambda^2}{8\sigma^2}\right),$$

and hence $f_\sigma(\lambda/2)$ is maximized at either $\sigma^2 = 0$, $\sigma^2 = p$ or $\sigma^2 = \lambda^2/4$. Comparing values at these points, we obtain that

$$\sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda/2) \leq f_{\lambda/2}(\lambda/2) \leq \frac{1}{\sqrt{2\pi}(\lambda/2)} \exp\left(-\frac{1}{2}\right) \leq \frac{1}{2\lambda}, \quad (\text{I.30})$$

whence, by letting $B \leq c\lambda\theta^2 / (p \log n)$, the upper bound of expectation become:

$$\sum_{i \in [n]} \mathbb{E} [\overline{\mathbf{X}}_i(\beta)] \leq \frac{4n}{2\lambda} (B + M) \leq \frac{4cn\theta^2}{p \log n} =: n\overline{\mathbb{E}\mathbf{X}}(\beta). \quad (\text{I.31})$$

2. (ε -net) Define $\varepsilon = \frac{c^2\lambda\theta^{3.5}}{3p^{2.5}\log^{2.5}n\log^{0.5}\theta^{-1}}$. Write $\lambda = c_\lambda/\sqrt{|\tau|}$ and consider the ε -net \mathcal{N}_ε for sphere of radius $\eta \leq \sqrt{p}$. From Lemma J.5 we know

$$|\mathcal{N}_\varepsilon| \leq \left(\frac{3\eta}{\varepsilon}\right)^{2p} \leq \left(\frac{81|\tau|p^6\log^5n\log\theta^{-1}}{c^4c_\lambda^2\theta^7}\right)^p \leq \left(\frac{2p\log n}{c \cdot c_\lambda}\right)^{13p} \quad (\text{I.32})$$

and define an event such that all center of subsets in \mathcal{N}_ε are being well-behaved:

$$\mathcal{E}_{\text{Net}} := \left\{ \sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta_\varepsilon) - n\mathbb{E}\bar{\mathbf{X}}(\beta_\varepsilon) < \frac{18cn\theta^2}{p \log n} \quad \forall \beta_\varepsilon \in \mathcal{N}_\varepsilon, \right\} \quad (\text{I.33})$$

3. (Lipschitz constant) Furthermore, the function $\sum_i^n \bar{\mathbf{X}}_i(\beta)$ is Lipchitz over β such that

$$\left| \sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta) - \sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta') \right| \leq \sum_{i \in [n]} \frac{1}{M} |\langle s_i[\mathbf{x}], \beta - \beta' \rangle| \leq \frac{n}{M} \max_{i \in [n]} \|\mathbf{P}_{[\pm p]+i} \mathbf{x}\|_2 \|\beta - \beta'\|_2 =: L \|\beta - \beta'\|_2$$

define the set \mathcal{N}_ε where Lipschitz constant is well bounded:

$$\mathcal{E}_{\text{Lip}} := \left\{ L \leq \frac{3n\sqrt{p\theta \log n \log \theta^{-1}}}{M} \right\},$$

then on event \mathcal{E}_{Lip} , for every β in $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$, there exists some β_ε in $\mathcal{N}_{\varepsilon,i}$ with $\|\beta - \beta_\varepsilon\|_2 \leq \varepsilon$, thus

$$\left| \left(\sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta) - n\mathbb{E}\bar{\mathbf{X}}(\beta) \right) - \left(\sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta_\varepsilon) - n\mathbb{E}\bar{\mathbf{X}}(\beta_\varepsilon) \right) \right| \leq 2L\varepsilon \leq \frac{2cn\theta^2}{p \log n}. \quad (\text{I.34})$$

On event $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$, from (I.31), (I.33) and (I.34), we can conclude that for every $\beta \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ and $i \in [n]$,

$$\sum_{i \in [n]} \bar{\mathbf{X}}_i(\beta) \leq \frac{24cn\theta^2}{p \log n} \quad (\text{I.35})$$

as desired, where the error probability of $\mathcal{E}_{\text{Lip}}^c$ is bounded using Lemma A.3, which gives

$$\mathbb{P}[\mathcal{E}_{\text{Lip}}^c] \leq \mathbb{P} \left[\max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 > 3\sqrt{p\theta \log n \log \theta^{-1}} \right] \leq 2/n, \quad (\text{I.36})$$

4. (Bound $\mathbb{P}[\mathcal{E}_{\text{Net}}^c]$) Wlog let us assume that $2p$ divides n . By applying union bound and observing that $\bar{\mathbf{X}}_i(\beta)$ is independent of $\bar{\mathbf{X}}_{i+2p}(\beta)$ for any $i \in [n]$, we split $\sum_i \bar{\mathbf{X}}_i(\beta)$ into $n/2p$ independent sums of r.v.s, we have

$$\mathbb{P}[\mathcal{E}_{\text{Net}}^c] \leq 2p |\mathcal{N}_\varepsilon| \cdot \mathbb{P} \left[\sum_{j=0}^{n/2p-1} (\bar{\mathbf{X}}_{2pj}(\beta) - \mathbb{E}[\bar{\mathbf{X}}(\beta)]) > \frac{9cn\theta^2}{p^2 \log n} \right],$$

where each summand has bounded variance and L^q -norm derived similarly as its expectation such that

$$\mathbb{E} \bar{\mathbf{X}}_i(\beta)^q \leq 2 \cdot \mathbb{P}[\langle s_i[\mathbf{x}], \beta \rangle \in [\lambda - B - M, \lambda + B + M]] \leq 2 \cdot \frac{1}{2\lambda} \cdot 2(B + M) \leq \frac{4c\theta^2}{p^2 \log n},$$

and apply Bernstein inequality Lemma J.4 with $(\sigma^2, R) = (4c\theta^2 / (p \log n), 1)$, obtains

$$\mathbb{P} \left[\sum_{j=0}^{n/2p-1} (\bar{\mathbf{X}}_{2pj}(\beta) - \mathbb{E}[\bar{\mathbf{X}}(\beta)]) > \frac{9cn\theta^2}{p^2 \log n} \right] \leq \exp \left[\frac{-(9cn\theta^2/p^2 \log n)^2}{2cn\theta^2/p^2 \log n + 2(9cn\theta^2/p^2 \log n)} \right] \leq \exp \left[\frac{-4cn\theta^2}{p^2 \log n} \right],$$

thus when $n = Cp^5\theta^{-2} \log p$:

$$\mathbb{P}[\mathcal{E}_{\text{Net}}^c] \leq \exp \left[\log(2p) + 13p \log \left(\frac{2p \log n}{c \cdot c_\lambda} \right) - \frac{4cn\theta^2}{p^2 \log n} \right] \leq 1/n \quad (\text{I.37})$$

as long as $\frac{C}{\log C} > 10^5 / (c^2 \cdot c_\lambda)$. ■

J Tools

Lemma J.1 (Tail bound for Gaussian r.v.). *If $X \sim \mathcal{N}(0, \sigma^2)$, then its tail bound for $t > 0$ can be*

$$\mathbb{P}[X > t] \leq \frac{\sigma}{t\sqrt{2\pi}} \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad (\text{J.1})$$

Lemma J.2 (Moments of the Gaussian random variables). *If $X \sim \mathcal{N}(0, \sigma^2)$, then for all integer $p \geq 1$,*

$$\mathbb{E}[|X|^p] \leq \sigma^p (p-1)!! \quad (\text{J.2})$$

Lemma J.3 (Gaussian concentration inequality). *Let $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ be a vector of n independent standard normal variables. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an L -Lipschitz function. Then for all $t > 0$,*

$$\mathbb{P}[|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| \geq t] \leq 2 \exp\left(-\frac{t^2}{2L^2}\right) \quad (\text{J.3})$$

Lemma J.4 (Moment control Bernstein inequality for scalar r.v.s). ([FR13], Theorem 7.30) *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be independent real-valued random variables. Suppose that there exist some positive number R and σ^2 such that $\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{X}_i^2] \leq \sigma^2$ and*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[|\mathbf{x}_k|^p] \leq \frac{1}{2} \sigma^2 R^{p-2} p!, \quad \text{for all integers } p \geq 3.$$

Let $S \doteq \sum_{i=1}^n \mathbf{x}_i$, then for all $t > 0$, it holds that

$$\mathbb{P}[|S - \mathbb{E}[S]| \geq t] \leq 2 \exp\left(-\frac{t^2}{2n\sigma^2 + 2Rt}\right) \quad (\text{J.4})$$

Lemma J.5 (ε -net on sphere). [Ver10] *Let (X, d) be a metric space and let $\varepsilon > 0$. A subset \mathcal{N}_ε of X is called an ε -net of X if for every point $x \in X$ there exists some point $y \in \mathcal{N}_\varepsilon$ so that $d(x, y) \leq \varepsilon$. There exists an ε -net \mathcal{N}_ε for the sphere \mathbb{S}^{n-1} of size $|\mathcal{N}_\varepsilon| \leq (3/\varepsilon)^n$.*

Lemma J.6 (Hanson-Wright). [RV+13] *Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be independent, subgaussian random variables with subgaussian norm $\sup_{p \geq 1} p^{-1/2} (\mathbb{E}|\mathbf{x}_i^p|)^{1/p} \leq \sigma$. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$, then for every $t > 0$,*

$$\mathbb{P}[|\mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbb{E} \mathbf{x}^* \mathbf{A} \mathbf{x}| \geq t] \leq 2 \exp\left(-c \min\left(\frac{t^2}{64 \sigma^4 \|\mathbf{A}\|_F^2}, \frac{t}{8\sqrt{2} \sigma^2 \|\mathbf{A}\|_2}\right)\right) \quad (\text{J.5})$$

Lemma J.7 (Maximum of separable convex function). *Let $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a convex function of the form $f(x) = x - s(x)$ with $s : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ satisfying*

$$\frac{s(x)}{x} \leq \frac{s(y)}{y}, \quad \text{for all } x \geq y > 0.$$

Then for $n \in \mathbb{N}$ and $0 < N \leq nL$,

$$\max_{0 \leq \mathbf{x} \leq L, \|\mathbf{x}\|_1 \leq N} \sum_{i=1}^n f(\mathbf{x}_i) \leq N \left(1 - \frac{s(L)}{L}\right) \quad (\text{J.6})$$

Proof. Since the feasible set is a convex polytope; the convex function $\sum_{i=1}^n f(\mathbf{x}_i)$ is maximized at a vertex, and that its vertices consist of 0 and permutations of the vector $[\underbrace{L, \dots, L}_{\lfloor N/L \rfloor}, r, 0, \dots, 0]$, where $r =$

$N - \lfloor N/L \rfloor L \leq L$. Then the function value at the maximizing vector \mathbf{x}_* can be derived as:

$$\begin{aligned} \sum_{i=1}^n f(\mathbf{x}_{*i}) &= \lfloor \frac{N}{L} \rfloor f(L) + f(r) = \frac{N-r}{L} (L - s(L)) + (r - s(r)) \\ &= N \left(1 - \frac{s(L)}{L}\right) + r \left(\frac{s(L)}{L} - \frac{s(r)}{r}\right) \leq N \left(1 - \frac{s(L)}{L}\right) \end{aligned}$$

■