

# Geometry and Symmetry in Short-and-Sparse Deconvolution\*

Han-Wen Kuo<sup>†</sup>, Yuqian Zhang<sup>‡</sup>, Yenson Lau<sup>†</sup>, and John Wright<sup>†§</sup>

**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  $\mu$ , 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.

**Key words.** Signal reconstruction, blind deconvolution, non-convex geometry, non-convex optimization.

**AMS subject classifications.** 94A12, 90C26, 65Y20.

**1. Introduction.** Datasets in a wide range of areas, including neuroscience [37], microscopy [15] and astronomy [49], 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 [14], 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  $\alpha$  and cyclic shift  $s_\ell[\cdot]$ ,

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

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<sup>2</sup>  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 \in \mathbb{R}^n$  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.

\*Submitted to the editors Jan/08/2019; revised Sep/20/2019.

**Funding:** This work was funded by NSF 1343282, NSF CCF 1527809, and NSF IIS 1546411

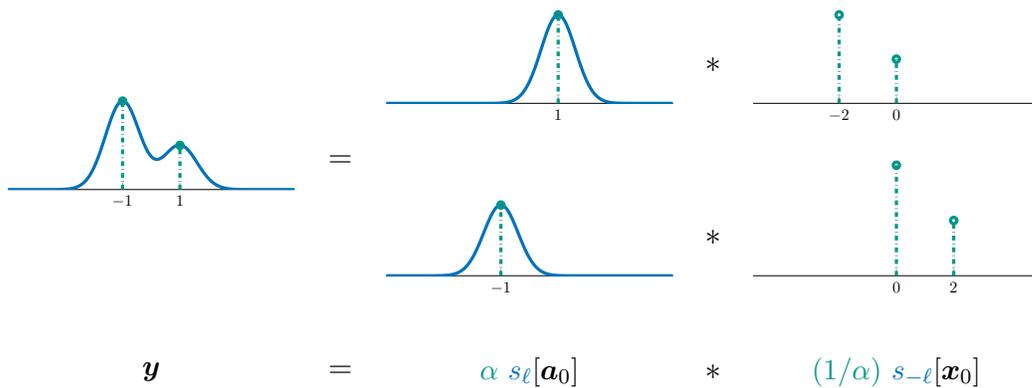
<sup>†</sup>Department of Electronic Engineering and Data Science Institute, Columbia University.

<sup>‡</sup>Department of Computer Science, Cornell University.

<sup>§</sup>Department of Applied Physics and Applied Mathematics, Columbia University.

<sup>1</sup>In this paper, the cyclic convolution  $\mathbf{a}_0 * \mathbf{x}_0$  assumes  $\mathbf{a}_0$  to be zero-padded  $[\mathbf{a}_0, \mathbf{0}^{n-p_0}]$  to length  $n$ .

<sup>2</sup>Our result can be applied to recovering direct convolutions. Let  $\mathbf{y} \in \mathbb{R}^{p_0+n-1}$  be the direct convolution between  $\mathbf{a}_0 \in \mathbb{R}^{p_0}$  and  $\mathbf{x}_0 \in \mathbb{R}^n$ , then  $\mathbf{y}$  can also be expressed as circular convolution between  $\mathbf{a}_0$  and  $[\mathbf{x}_0; \mathbf{0}^{p_0-1}]$ .



**Figure 1. Shift symmetry in Short-and-Sparse deconvolution.** An observation  $\mathbf{y}$  (left) which is a convolution of a short signal  $\mathbf{a}_0$  and a sparse signal  $\mathbf{x}_0$  (top right) can be equivalently expressed as a convolution of  $s_\ell[\mathbf{a}_0]$  and  $s_{-\ell}[\mathbf{x}_0]$ , where  $s_\ell[\cdot]$  denotes a shift  $\ell$  samples. The ground truth signals  $\mathbf{a}_0$  and  $\mathbf{x}_0$  can only be identified up to a scaled shift.

32 Despite a long history and many applications, until recently very little algorithmic theory  
 33 was available for SaS deconvolution. Much of this difficulty can be attributed to the scale-shift  
 34 symmetry: natural convex relaxations fail<sup>3</sup>, and nonconvex formulations exhibit a complicated  
 35 optimization landscape, with many equivalent global minimizers (scaled shifts of the ground  
 36 truth) and additional local minimizers (scaled shift truncations of the ground truth), and a  
 37 variety of critical points [63, 64]. Currently available theory guarantees approximate recovery  
 38 of a truncation<sup>4</sup> of a shift  $s_\ell[\mathbf{a}_0]$ , rather than guaranteeing recovery of  $\mathbf{a}_0$  as a whole, and  
 39 requires certain (complicated) conditions on the convolution matrix associated with  $\mathbf{a}_0$  [63].

40 In this paper, we describe an algorithm which, under simpler conditions, *exactly* recovers a  
 41 scaled shift of the pair  $(\mathbf{a}_0, \mathbf{x}_0)$ . Our algorithm is based on a formulation first introduced in  
 42 [64], which casts the deconvolution problem as (nonconvex) optimization over the sphere. We  
 43 characterize the geometry of this objective function, and show that near a certain union of  
 44 subspaces, every local minimizer is very close to a signed shift of  $\mathbf{a}_0$ . Based on this geometric  
 45 analysis, we give provable methods for SaS deconvolution that exactly recover a scaled shift  
 46 of  $(\mathbf{a}_0, \mathbf{x}_0)$  whenever  $\mathbf{a}_0$  is *shift-incoherent* and  $\mathbf{x}_0$  is a sufficiently sparse random vector. Our  
 47 geometric analysis highlights the role of symmetry in shaping the objective landscape for SaS  
 48 deconvolution.

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

<sup>3</sup>Such as matrix lifting relaxation [2, 39], in which  $\mathbf{a}_0$  or  $\mathbf{x}_0$  resides in random subspaces w/o shift symmetry.

<sup>4</sup>I.e., the portion of the shifted signal  $s_\ell[\mathbf{a}_0]$  that falls in the window  $\{0, \dots, p_0 - 1\}$ .

54 **2. Formulation and Assumptions.**

 55 **2.1. Nonconvex SaS over the Sphere.** Our starting point is the (natural) formulation

56 (2.1) 
$$\min_{\mathbf{a}, \mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2}_{\text{Data Fidelity}} + \lambda \underbrace{\|\mathbf{x}\|_1}_{\text{Sparsity}} \quad \text{s.t.} \quad \|\mathbf{a}\|_2 = 1.$$

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

59 (2.2) 
$$\varphi_{\text{lasso}}(\mathbf{a}) \equiv \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}$$

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

61 (2.3) 
$$\min_{\mathbf{a}} \varphi_{\text{lasso}}(\mathbf{a}) \quad \text{s.t.} \quad \|\mathbf{a}\|_2 = 1.$$

 62 In (2.1)-(2.3), we constrain  $\mathbf{a}$  to have unit  $\ell^2$  norm. This constraint breaks the scale ambi-  
 63 guity between  $\mathbf{a}$  and  $\mathbf{x}$ . Moreover, the choice of constraint manifold has surprisingly strong  
 64 implications for computation: if  $\mathbf{a}$  is instead constrained to the simplex, the problem admits  
 65 trivial global minimizers. In contrast, local minima of the sphere-constrained formulation often  
 66 correspond to shifts (or shift truncations [64]) of the ground truth  $\mathbf{a}_0$ .

 67 The problem (2.3) is defined in terms of the optimal Lasso cost. This function is challenging  
 68 to analyze, especially far away from  $\mathbf{a}_0$ . [64] analyzes the local minima of a simplification of  
 69 (2.3), obtained by approximating<sup>5</sup> the data fidelity term as

70 (2.4) 
$$\begin{aligned} \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 &= \frac{1}{2} \|\mathbf{a} * \mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2, \\ &\approx \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2. \end{aligned}$$

72 This yields a simpler objective function

73 (2.5) 
$$\varphi_{\ell^1}(\mathbf{a}) = \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}.$$

 74 We make one further simplification to this problem, replacing the nondifferentiable penalty  
 75  $\|\cdot\|_1$  with a smooth approximation  $\rho(\mathbf{x})$ .<sup>6</sup> Our analysis allows for a variety of smooth sparsity  
 76 surrogates  $\rho(\mathbf{x})$ ; for concreteness, we state our main results for the particular penalty<sup>7</sup>

77 (2.6) 
$$\rho(\mathbf{x}) = \sum_i (\mathbf{x}_i^2 + \delta^2)^{1/2}.$$

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

80 (2.7) 
$$\varphi_{\rho}(\mathbf{a}) = \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x}\|_2^2 - \langle \mathbf{a} * \mathbf{x}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|_2^2 + \lambda \rho(\mathbf{x}) \right\}.$$

---

<sup>5</sup>For a generic  $\mathbf{a}$ , we have  $\langle s_i[\mathbf{a}], s_j[\mathbf{a}] \rangle \approx 0$  and hence  $\|\mathbf{a} * \mathbf{x}\|_2^2 = \mathbf{x}^* \mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}} \mathbf{x} \approx \mathbf{x}^* \mathbf{I} \mathbf{x} = \|\mathbf{x}\|_2^2$ . The use of  $\varphi_{\rho}$  performs not as ideal comparing to bilinear Lasso when this approximation is inexact, see Section 7.

<sup>6</sup> $\varphi_{\ell^1}$  is not twice differentiable everywhere hence can't be minimized with conventional second order methods.

<sup>7</sup>This particular surrogate is sometimes being named as the pseudo-Huber function.

81 As in [64], we optimize  $\varphi_\rho(\mathbf{a})$  over the sphere  $\mathbb{S}^{p-1}$ :

82 (2.8)

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

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

89 **2.2. Analysis Setting and Assumptions.** For convenience, we assume that  $\mathbf{a}_0$  has unit  $\ell^2$   
 90 norm, i.e.,  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$ .<sup>8</sup> Our analysis makes two main assumptions, on the short motif  $\mathbf{a}_0$  and  
 91 the sparse map  $\mathbf{x}_0$ , respectively:

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

94 (2.9)

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

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

- 98 • *Spiky.* When  $\mathbf{a}_0$  is close to the Dirac delta  $\delta_0$ , the shift coherence  $\mu(\mathbf{a}_0) \approx 0$ .<sup>9</sup> Here,  
 99 the observed signal  $\mathbf{y}$  consists of a superposition of sharp pulses. This is arguably the  
 100 easiest instance of SaS deconvolution.
- 101 • *Generic.* If  $\mathbf{a}_0$  is chosen uniformly at random from the sphere  $\mathbb{S}^{p_0-1}$ , its coherence is  
 102 bounded as  $\mu(\mathbf{a}_0) \lesssim \sqrt{1/p_0}$  with high probability.
- 103 • *Tapered Generic Lowpass.* Here,  $\mathbf{a}_0$  is generated by taking a random conjugate  
 104 symmetric superposition of the first  $L$  length- $p_0$  Discrete Fourier Transform (DFT)  
 105 basis signals, windowing (e.g., with a Hamming window) and normalizing to unit  $\ell^2$   
 106 norm. When  $L = p_0\sqrt{1 - \beta}$ , with high probability  $\mu(\mathbf{a}_0) \lesssim \beta$ . In this model,  $\mu$  does  
 107 not have to diminish as  $p_0$  grows – it can be a fixed constant<sup>10</sup>.

108 Intuitively speaking, problems with smaller  $\mu$  are easier to solve, a claim which will be made  
 109 precise in our technical results.

110 We assume that  $\mathbf{x}_0$  is a sparse random vector. More precisely, we assume that  $\mathbf{x}_0$  is  
 111 Bernoulli-Gaussian, with rate  $\theta$ :

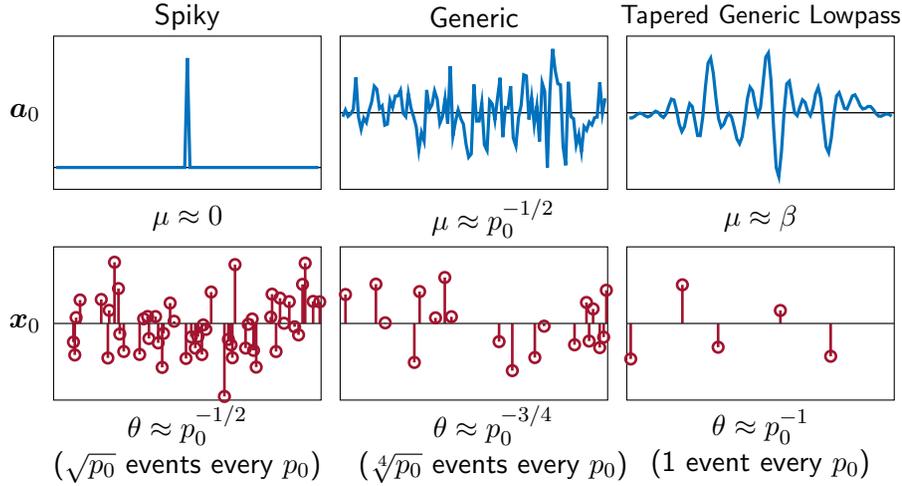
112 (2.10)

$$\mathbf{x}_{0i} = \omega_i \mathbf{g}_i,$$

<sup>8</sup>This 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  $\alpha = \pm \frac{1}{\|\mathbf{a}_0\|_2}$ .

<sup>9</sup>The use of “ $\approx$ ” here suppresses constant and logarithmic factors.

<sup>10</sup>The 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. See [34].



**Figure 2. Sparsity-coherence tradeoff:** *Top: three families of motifs  $\mathbf{a}_0$  with varying coherence  $\mu$ . Bottom: maximum allowable sparsity  $\theta$  and number of copies  $\theta p_0$  within each length- $p_0$  window. Here, we suppress constants and logarithmic factors. When the target motif has smaller shift-coherence  $\mu$ , our result allows larger  $\theta$ , 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.*

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

$$115 \quad (2.11) \quad \mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta).$$

116 Here,  $\theta$  is the probability that a given entry  $\mathbf{x}_{0i}$  is nonzero. Problems with smaller  $\theta$  are easier  
 117 to solve. In the extreme case, when  $\theta \ll 1/p_0$ , the observation  $\mathbf{y}$  contains many isolated copies  
 118 of the motif  $\mathbf{a}_0$ , and  $\mathbf{a}_0$  can be determined by direct inspection. Our analysis will focus on the  
 119 nontrivial scenario, when  $\theta \gtrsim 1/p_0$ .

120 Our technical results will articulate *sparsity-coherence* tradeoffs, in which smaller coherence  
 121  $\mu$  enables larger  $\theta$ , and vice-versa. More specifically, in our main theorem, the sparsity-coherence  
 122 relationship is captured in the form

$$123 \quad (2.12) \quad \theta \lesssim 1/(p_0\sqrt{\mu} + \sqrt{p_0}).$$

125 When the target  $\mathbf{a}_0$  is very shift-incoherent ( $\mu \approx 0$ ), our method succeeds when each length- $p_0$   
 126 window contains about  $\sqrt{p_0}$  copies of  $\mathbf{a}_0$ . When  $\mu$  is larger (as in the generic lowpass model),  
 127 our method succeeds as long as relatively few copies of  $\mathbf{a}_0$  overlap in the observed signal. In  
 128 [Figure 2](#), we illustrate these tradeoffs for the three models described above.

129 **3. Main Results: Geometry and Algorithms.** In this section, we introduce our main  
 130 results – on the geometry of  $\varphi_\rho$  ([Subsection 3.1](#)) and its algorithmic implications ([Subsection 3.2](#)).  
 131 Finally, in [Subsection 3.3](#), we compare these results with the literature on deconvolution.

132 **3.1. Geometry of the Objective  $\varphi_\rho$ .** The goal in SaS de-  
 133 convolution is to recover  $\mathbf{a}_0$  (and  $\mathbf{x}_0$ ) up to a signed shift — i.e.,  
 134 we wish to recover some  $\pm s_\ell[\mathbf{a}_0]$ . The shifts  $\pm s_\ell[\mathbf{a}_0]$  play a key  
 135 role in shaping the landscape of  $\varphi_\rho$ . In particular, we will argue  
 136 that over a certain subset of the sphere, *every local minimum*  
 137 *of  $\varphi_\rho$  is close to some  $\pm s_\ell[\mathbf{a}_0]$ .*

138 To gain intuition into the properties of  $\varphi_\rho$ , we first visualize  
 139 this function in the vicinity of a single shift  $s_\ell[\mathbf{a}_0]$  of the ground  
 140 truth  $\mathbf{a}_0$ . In Figure 3, we plot the function value of  $\varphi_\rho$  over

$$141 \quad \mathcal{B}_{\ell^2, r}(s_\ell[\mathbf{a}_0]) \cap \mathbb{S}^{p-1},$$

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

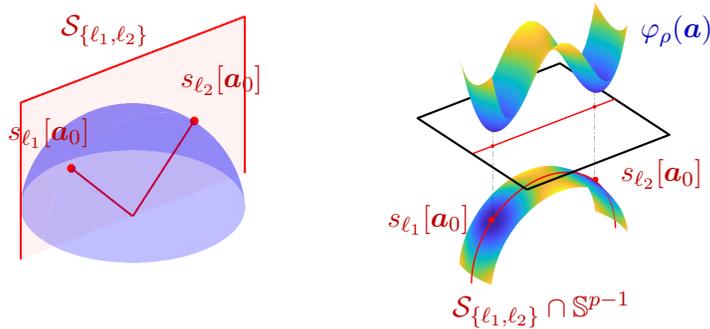
- 145 • The objective function  $\varphi_\rho$  is strongly convex in this
- 146 neighborhood of  $s_\ell[\mathbf{a}_0]$ .
- 147 • There is a local minimizer very close to  $s_\ell[\mathbf{a}_0]$ .

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

$$151 \quad \mathcal{S}_{\{\ell_1, \ell_2\}} = \{ \alpha_1 s_{\ell_1}[\mathbf{a}_0] + \alpha_2 s_{\ell_2}[\mathbf{a}_0] \mid \alpha_1, \alpha_2 \in \mathbb{R} \}.$$

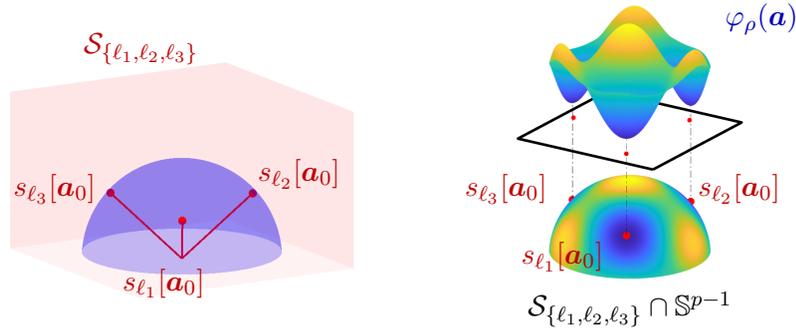
152 We make three observations:

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



**Figure 3.** Geometry of  $\varphi_\rho$  near a shift of  $\mathbf{a}_0$ . *Bottom:* a portion of the sphere  $\mathbb{S}^{p-1}$ , colored according to  $\varphi_\rho$ . *Top:*  $\varphi_\rho$  visualized as height.  $\varphi_\rho$  is strongly convex in this region, and it has a minimizer very close to  $s_\ell[\mathbf{a}_0]$ .

**Figure 4.** Geometry of  $\varphi_\rho$  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  $\varphi_\rho$  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  $\varphi_\rho$  is convex in direction away from  $\mathcal{S}_{\{\ell_1, \ell_2\}}$ .



**Figure 5.** Geometry of  $\varphi_\rho$  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,  $\varphi_\rho$  has local minimizers near each of the  $s_{\ell_i}[\mathbf{a}_0]$ , and are the only minimizers near  $\mathcal{S}_{\ell_1, \ell_2, \ell_3}$ .

159 Finally, we visualize  $\varphi_\rho$  over the intersection (Figure 5, left) of the sphere  $\mathbb{S}^{p-1}$  with the  
 160 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  $\mathbf{a}_0$ :

$$161 \quad \mathcal{S}_{\{\ell_1, \ell_2, \ell_3\}} = \{ \alpha_1 s_{\ell_1}[\mathbf{a}_0] + \alpha_2 s_{\ell_2}[\mathbf{a}_0] + \alpha_3 s_{\ell_3}[\mathbf{a}_0] \mid \alpha_1, \alpha_2, \alpha_3 \in \mathbb{R} \}$$

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

165 Our main geometric result will show that these properties are obtained from *every* subspace  
 166 spanned by a few shifts of  $\mathbf{a}_0$ . Indeed, for each subset

$$167 \quad (3.1) \quad \tau \subseteq \{-p_0 + 1, \dots, p_0 - 1\},$$

168 define a linear subspace

$$169 \quad (3.2) \quad \mathcal{S}_\tau = \left\{ \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] \mid \alpha_{-p_0+1}, \dots, \alpha_{p_0-1} \in \mathbb{R} \right\}.$$

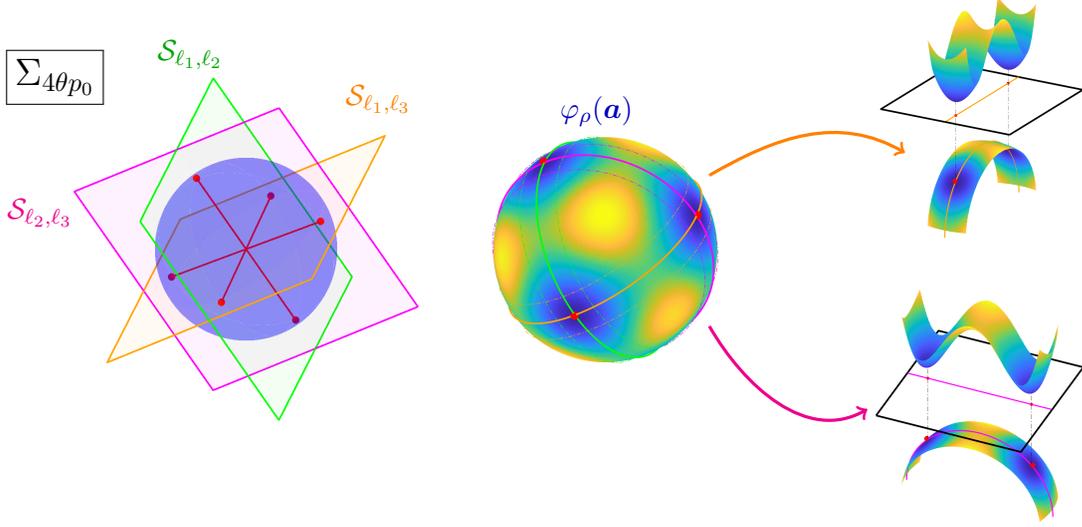
170 The subspace  $\mathcal{S}_\tau$  is the linear span of the shifts  $s_\ell[\mathbf{a}_0]$  indexed by  $\ell$  in the set  $\tau$ . Our geometric  
 171 theory will show that with high probability the function  $\varphi_\rho$  has no spurious local minimizers  
 172 near any  $\mathcal{S}_\tau$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

$$174 \quad (3.3) \quad \Sigma_{4\theta p_0} = \bigcup_{|\tau| \leq 4\theta p_0} \mathcal{S}_\tau.$$

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

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

178 Our main result formalizes the above observations, under two key assumptions: first, that  
 179 the sparsity rate  $\theta$  is sufficiently small (relative to the shift coherence  $\mu$  of  $p_0$ ), and, second,  
 180 the signal length  $n$  is sufficiently large:



**Figure 6.** Geometry of  $\varphi_\rho$  over the union of subspaces  $\Sigma_{2\theta p_0}$ . Left: schematic representation of the union of subspaces  $\Sigma_{4\theta p_0}$ . For each set  $\tau$  of at most  $4\theta p_0$  shifts, we have a subspace  $\mathcal{S}_\tau$ . Right:  $\varphi_\rho$  has good geometry near this union of subspaces.

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

$$183 \quad (3.4) \quad \theta \in \left[ \frac{c_1}{p_0}, \frac{c_2}{p_0 \sqrt{\mu} + \sqrt{p_0}} \right] \cdot \frac{1}{\log^2 p_0}.$$

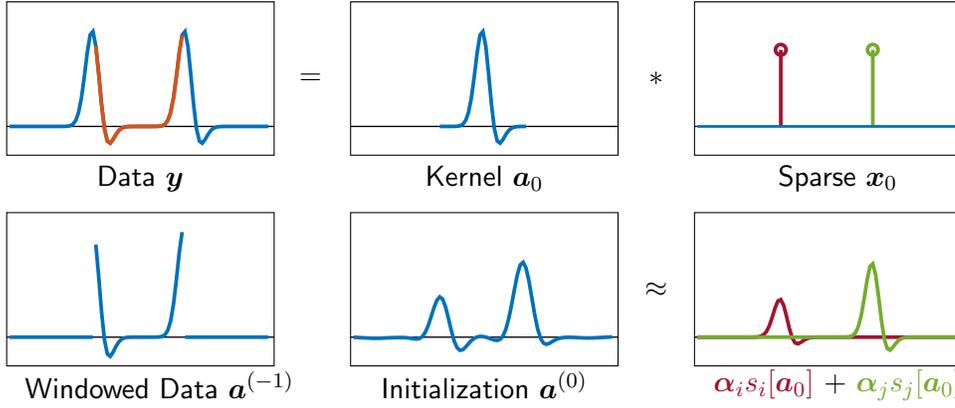
185 Chose  $\rho(x) = \sqrt{x^2 + \delta^2}$  and set  $\lambda = 0.1/\sqrt{p_0\theta}$  in  $\varphi_\rho$ . Then there exists  $\delta > 0$  and numerical  
 186 constant  $c$  such that if  $n \geq \text{poly}(p_0)$ , with high probability, every local minimizer  $\bar{\mathbf{a}}$  of  $\varphi_\rho$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

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

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

193 At a technical level, our proof of [Theorem 3.1](#) shows that (i)  $\varphi_\rho(\mathbf{a})$  is strongly convex in  
 194 the vicinity of each signed shift, and that at every other point  $\mathbf{a}$  near  $\Sigma_{4\theta p_0}$ , there is either  
 195 (ii) a nonzero gradient or (iii) a direction of strict negative curvature; furthermore (iv) the  
 196 function  $\varphi_\rho$  grows away from  $\Sigma_{4\theta p_0}$ . Points (ii)-(iii) imply that near  $\Sigma_{4\theta p_0}$  there are no “flat”  
 197 saddles: every saddle point has a direction of strict negative curvature. We will leverage these  
 198 properties to propose an efficient algorithm for finding a local minimizer near  $\Sigma_{4\theta p_0}$ . Moreover,  
 199 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  
 200  $s_\ell[\mathbf{a}_0]$ : we will give a refinement algorithm that produces  $(\pm s_\ell[\mathbf{a}_0], \pm s_{-\ell}[\mathbf{x}_0])$ .

<sup>11</sup>Typically it is possible to provide an overestimate  $p'_0 \geq p_0$ . Our theory and algorithm can be applied directly to the overestimate  $p'_0$ , with the caveat that the sparsity rate  $\theta$  now scales with  $p'_0$  rather than  $p_0$ .



**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).

201 **3.2. Provable Algorithm for SaS Deconvolution.** The objective function  $\varphi_\rho$  has good  
 202 geometric properties on (and near!) the union of subspaces  $\Sigma_{4\theta p_0}$ . In this section, we show  
 203 how to use give an efficient method that exactly recovers  $\mathbf{a}_0$  and  $\mathbf{x}_0$ , up to shift symmetry.  
 204 Although our geometric analysis only controls  $\varphi_\rho$  near  $\Sigma_{4\theta p_0}$ , we will give a descent method  
 205 which, with appropriate initialization  $\mathbf{a}^{(0)}$ , produces iterates  $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(k)}, \dots$  that remain  
 206 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  
 207 finding a local minimizer  $\bar{\mathbf{a}}$ , we refine it to produce a signed shift of  $(\mathbf{a}_0, \mathbf{x}_0)$  using alternating  
 208 minimization.

209 The next two paragraphs give the main ideas behind the main steps of the algorithm. We  
 210 then describe its components in more detail ([Algorithm 3.1](#)) and state our main algorithmic  
 211 result ([Theorem 3.2](#)), which asserts that under appropriate conditions this method produces a  
 212 signed shift of  $(\mathbf{a}_0, \mathbf{x}_0)$ .

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

$$217 \quad (3.5) \quad \mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0 = \sum_{\ell \in \text{supp}(\mathbf{x}_0)} \mathbf{x}_{0\ell} s_\ell[\mathbf{a}_0].$$

218 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  $2\theta p_0 \ll$   
 219  $4\theta p_0$  shifts  $s_\ell[\mathbf{a}_0]$ .

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

$$222 \quad (3.6) \quad [\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}].$$

224 Zero padding provides enough space to accommodate any shift  
 225  $s_\ell[\mathbf{a}_0]$  with  $\ell \in \tau$ . We then perform one step of the generalized  
 226 power method<sup>12</sup>, writing

$$227 \quad (3.7) \quad \mathbf{a}^{(0)} = -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\ell^1} (\mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{0}^{p_0-1}; \mathbf{y}_0; \dots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}]),$$

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

236 The next key observation is that the function  $\varphi_\rho$  grows as we move away from the subspace  
 237  $\mathcal{S}_\tau$  – see Figure 8. Because of this, a small-stepping descent method will not move far away  
 238 from  $\Sigma_{4\theta p_0}$ . For concreteness, we will analyze a variant of the curvilinear search method [23, 24]  
 239 , which moves in a linear combination of the negative gradient direction  $-\mathbf{g}$  and a negative  
 240 curvature direction  $-\mathbf{v}$ . At the  $k$ -th iteration, the algorithm updates  $\mathbf{a}^{(k+1)}$  as

$$241 \quad (3.8) \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)}]$$

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

247 The second step of our algorithm *rounds* the local minimizer  $\bar{\mathbf{a}} \approx \sigma s_\ell[\mathbf{a}_0]$  to produce an  
 248 exact solution  $\hat{\mathbf{a}} = \sigma s_\ell[\mathbf{a}_0]$ . As a byproduct, it also exactly recovers the corresponding signed  
 249 shift of the true sparse signal,  $\hat{\mathbf{x}} = \sigma s_{-\ell}[\mathbf{x}_0]$ .

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

$$256 \quad (3.9) \quad \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i \notin I^{(k)}} |\mathbf{x}_i|.$$

257 This can be viewed as an extreme form of *reweighting* [11]. Second, our algorithm gradually  
 258 decreases penalty variable  $\lambda$  to 0, so that eventually

$$259 \quad (3.10) \quad \hat{\mathbf{a}} * \hat{\mathbf{x}} \approx \mathbf{y}.$$

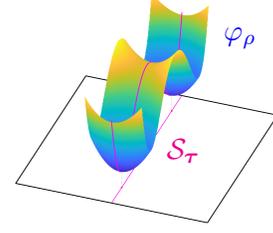
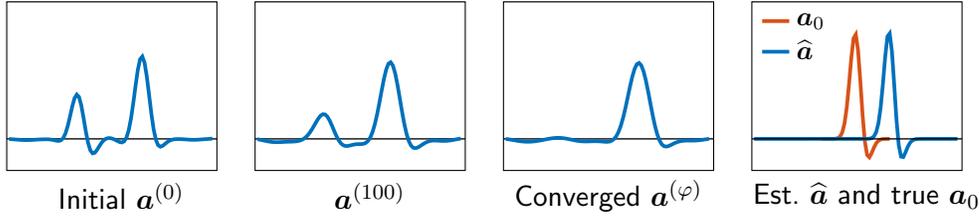


Figure 8. Growth of  $\varphi_\rho$  away from  $\mathcal{S}_\tau$ . Because  $\varphi_\rho$  grows away from  $\mathcal{S}_\tau$ , small-stepping descent methods stay near  $\mathcal{S}_\tau$ .

<sup>12</sup>The 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  $\varphi$  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  $\varphi_\rho$ .



**Figure 9. Local minimization and refinement.** *Left: data-driven initialization  $\mathbf{a}^{(0)}$  consisting of a near-superposition of two shifts. Middle: minimizing  $\varphi_\rho$  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  $\mathbf{a}_0$ .*

261 This can be viewed as a *homotopy* or *continuation* method [46, 19]. For concreteness, at  $k$ -th  
 262 iteration the algorithm reads:

$$263 \quad (3.11) \quad \text{Update } \mathbf{x}: \quad \mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|,$$

$$264 \quad (3.12) \quad \text{Update } \mathbf{a}: \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \underset{\mathbf{a}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right],$$

$$265 \quad (3.13) \quad \text{Update } \lambda \text{ and } I: \quad \lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \operatorname{supp}(\mathbf{x}^{(k+1)}).$$

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

269 Our overall algorithm is summarized as [Algorithm 3.1](#). [Figure 9](#) illustrates the main  
 270 steps of this algorithm. Our main algorithmic result states that under essentially the same  
 271 hypotheses as above, [Algorithm 3.1](#) produces a signed shift of the ground truth  $(\mathbf{a}_0, \mathbf{x}_0)$ :

272 **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  $\mu$ -*  
 273 *truncated shift coherent<sup>14</sup> such that  $\max_{i \neq j} |\langle \mathbf{t}_{p_0}^* s_i[\mathbf{a}_0], \mathbf{t}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu$  and  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in$*   
 274  *$\mathbb{R}^n$  with  $\theta, \mu$  satisfying*

$$275 \quad (3.19) \quad \theta \in \left[ \frac{c_1}{p_0}, \frac{c_2}{(p_0 \sqrt{\mu} + \sqrt{p_0}) \log^2 p_0} \right], \quad \mu \leq \frac{c_3}{\log^2 n}$$

277 *for some constant  $c_1, c_2, c_3 > 0$ . If the signal lengths  $n, p_0$  satisfy  $n > \text{poly}(p_0)$  and  $p_0 >$*   
 278  *$\text{polylog}(n)$ , then there exist  $\delta, \eta_v > 0$  such that with high probability, [Algorithm 3.1](#) produces  
 279  *$(\hat{\mathbf{a}}, \hat{\mathbf{x}})$  that are equal to the ground truth up to signed shift symmetry:**

$$280 \quad (3.20) \quad \left\| (\hat{\mathbf{a}}, \hat{\mathbf{x}}) - \sigma(s_\ell[\mathbf{a}_0], s_{-\ell}[\mathbf{x}_0]) \right\|_2 \leq \varepsilon$$

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

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

283 When solving SaS deconvolution via minimizing bilinear Lasso objective (2.2) in practice,  
 284 the algorithm is analogous to the provable method introduced in [Algorithm 3.1](#), where the

<sup>14</sup>The truncated shift coherence is a stronger condition than natural shift coherence. The statement appears mainly due to the limitation of prove strategy for algorithm.

**Algorithm 3.1** Short and Sparse Deconvolution

**Input:** Observation  $\mathbf{y}$ , motif length  $p_0$ , sparsity  $\theta$ , shift-coherence  $\mu$ , 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  $\lambda = 0.1/\sqrt{p_0\theta}$ <sup>13</sup> and  $\delta > 0$  in  $\varphi_\rho$ . For  $k = 1, 2, \dots, K_1$ , let

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

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

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

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

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

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

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

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

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

285 curvilinear descent and the refinement step can be realized as alternating gradient descent of  
 286 both variables  $\mathbf{a}, \mathbf{x}$  in (2.2). Unlike Algorithm 3.1, this alternating gradient method has yet  
 287 come with theoretical guarantees, but shown to be an effective and efficient method for SaS  
 288 deconvolution problems both in simulation and in reality [34].

289 **3.3. Relationship to the Literature.** Blind deconvolution is a classical problem in signal  
 290 processing [54, 12], and has been studied under a variety of hypotheses. In this section, we first  
 291 discuss the relationship between our results and the existing literature on the short-and-sparse  
 292 version of this problem, and then briefly discuss other deconvolution variants in the theoretical  
 293 literature.

294 The short-and-sparse model arises in a number of applications. One class of applications  
 295 involves finding basic motifs (repeated patterns) in datasets. This *motif discovery* problem  
 296 arises in extracellular spike sorting [37, 20] and calcium imaging [48], where the observed signal  
 297 exhibits repetitive *short* neuron excitation patterns occurring *sparsely* across time and/or  
 298 space. Similarly, electron microscopy images [15] arising in study of nanomaterials often exhibit  
 299 repeated motifs.

300 Another significant application of SaS deconvolution is *image deblurring*. Typically, the  
 301 blur kernel is small relative to the image size (*short*) [3, 62, 13, 35, 36]. In natural image

<sup>14</sup>In practice, we suggest setting  $\lambda = c_\lambda/\sqrt{p_0\theta}$  with  $c_\lambda \in [0.5, 0.8]$ .

302 deblurring, the target image is often assumed to have relatively few sharp edges [21, 27, 36],  
 303 and hence have *sparse* derivatives. In scientific image deblurring, e.g., in astronomy [33, 25, 9]  
 304 and geophysics [28], the target image is often sparse, either in the spatial or wavelet domains,  
 305 again leading to variants of the SaS model. The literature on blind image deconvolution is  
 306 large; see, e.g., [31, 10] for surveys.

307 Variants of the SaS deconvolution problem arise in many other areas of engineering as well.  
 308 Examples include *blind equalization* in communications [50, 51, 26], *dereverberation* in sound  
 309 engineering [44, 45] and image *super-resolution* [4, 53, 61].

310 These applications have motivated a great deal of algorithmic work on variants of the  
 311 SaS problem [32, 8, 6, 31, 43, 10, 56]. In contrast, relatively little theory is available to  
 312 explain when and why algorithms succeed. Our algorithm minimizes  $\varphi_\rho$  as an approximation  
 313 to the Lasso cost over the sphere. Our formulation and results have strong precedent in  
 314 the literature. Lasso-like objective functions have been widely used in image deblurring  
 315 [62, 14, 21, 35, 52, 60, 18, 30, 36, 59, 47, 64]. A number of insights have been obtained into the  
 316 geometry of sparse deconvolution – in particular, into the effect of various constraints on  $\mathbf{a}$  on  
 317 the presence or absence of spurious local minimizers. In image deblurring, a simplex constraint  
 318 ( $\mathbf{a} \geq \mathbf{0}$  and  $\|\mathbf{a}\|_1 = 1$ ) arises naturally from the physical structure of the problem [62, 14].  
 319 Perhaps surprisingly, simplex-constrained deconvolution admits trivial global minimizers, at  
 320 which the recovered kernel  $\mathbf{a}$  is a spike, rather than the target blur kernel [7, 36].

321 [59] imposes the  $\ell^2$  regularization on  $\mathbf{a}$  and observes that this alternative constraint gives  
 322 more reliable algorithm. [64] studies the geometry of the simplified objective  $\varphi_{\ell^1}$  over the  
 323 sphere, and proves that in the dilute limit in which  $\mathbf{x}_0$  has one nonzero entry, all strict local  
 324 minima of  $\varphi_{\ell^1}$  are close to signed shifts truncations of  $\mathbf{a}_0$ . By adopting a different objective  
 325 function (based on  $\ell^4$  maximization) over the sphere, [63] proves that on a certain region of  
 326 the sphere every local minimum is near a *truncated* signed shift of  $\mathbf{a}_0$ , i.e., the restriction of  
 327  $s_\ell[\mathbf{a}_0]$  to the window  $\{0, \dots, p_0 - 1\}$ . The analysis of [63] allows the sparse sequence  $\mathbf{x}_0$  to be  
 328 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 [64]  
 329 and [63] guarantee *approximate* recovery of a portion of  $s_\ell[\mathbf{a}_0]$ , under complicated conditions  
 330 on the kernel  $\mathbf{a}_0$ . Our core optimization problem is very similar to [64]. However, we obtain  
 331 *exact* recovery of both  $\mathbf{a}_0$  and relatively dense  $\mathbf{x}_0$ , under the much simpler assumption of shift  
 332 incoherence.

333 Other aspects of the SaS problem have been studied theoretically. One basic question is  
 334 under what circumstances the problem is identifiable, up to the scaled shift ambiguity. [17]  
 335 shows that the problem ill-posed for worst case  $(\mathbf{a}_0, \mathbf{x}_0)$  – in particular, for certain support  
 336 patterns in which  $\mathbf{x}_0$  does not have any isolated nonzero entries. This demonstrates that *some*  
 337 modeling assumptions on the support of the sparse term are needed. At the same time, this  
 338 worst case structure is unlikely to occur, either under the Bernoulli model, or in practical  
 339 deconvolution problems.

340 Motivated by a variety of applications, many low-dimensional deconvolution models have  
 341 been studied in the theoretical literature. In communication applications, the signals  $\mathbf{a}_0$  and  
 342  $\mathbf{x}_0$  either live in known low-dimensional subspaces, or are sparse in some known dictionary  
 343 [2, 16, 29, 39, 40, 41, 42]. These theoretical works assume that the subspace / dictionary are  
 344 chosen at random. This low-dimensional deconvolution model does not exhibit the signed

345 shift ambiguity; nonconvex formulations for this model exhibit a different structure from that  
 346 studied here. In fact, the variant in which both signals belong to known subspaces can be solved  
 347 by convex relaxation [2]. The SaS model does not appear to be amenable to convexification,  
 348 and exhibits a more complicated nonconvex geometry, due to the shift ambiguity. The main  
 349 motivation for tackling this model lies in the aforementioned applications in imaging and data  
 350 analysis.

351 [38, 57] study the related *multi-instance* sparse blind deconvolution problem (MISBD),  
 352 where there are  $K$  observations  $\mathbf{y}_i = \mathbf{a}_0 * \mathbf{x}_i$  consisting of multiple convolutions  $i = 1, \dots, K$  of  
 353 a kernel  $\mathbf{a}_0$  and different sparse vectors  $\mathbf{x}_i$ . Both works develop provable algorithms. There are  
 354 several key differences with our work. First, both the proposed algorithms and their analysis  
 355 require the kernel to be invertible. Second, despite the apparent similarity between the SaS  
 356 model and MISBD, these problems are not equivalent. It might seem possible to reduce SaS  
 357 to MISBD by dividing the single observation  $\mathbf{y}$  into  $K$  pieces; this apparent reduction fails  
 358 due to boundary effects.

359 **3.4. Notations.** All vectors/matrices are written in bold font  $\mathbf{a}/\mathbf{A}$ ; indexed values are writ-  
 360 ten 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  
 361 as  $\mathbf{e}_i$ . The indices for vectors/matrices all start from 0 and is taking modulo- $n$ , thus a vector  
 362 of length  $n$  should have its indices labeled as  $\{0, 1, \dots, n-1\}$ . We write  $[n] = \{0, \dots, n-1\}$ .  
 363 We often use capital italic symbols  $I, J$  for subsets of  $[n]$ . We abuse notation slightly and write  
 364  $[-p] = \{n-p+1, \dots, n-1, 0\}$  and  $[\pm p] = \{n-p+1, \dots, n-1, 0, 1, \dots, p-1\}$ . Index sets  
 365 can be labels for vectors;  $\mathbf{a}_I \in \mathbb{R}^{|I|}$  denotes the restriction of the vector  $\mathbf{a}$  to coordinates  $I$ .  
 366 Also, we use check symbol for reversal operator on index set  $\check{I} = -I$  and vectors  $\check{\mathbf{a}}_i = \mathbf{a}_{-i}$ .

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

373 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  
 374 convolution operator works on index set:  $I * J = \text{supp}(\mathbf{1}_I * \mathbf{1}_J)$ . Similarly, the shift operator  
 375  $s_\ell[\cdot] : \mathbb{R}^p \rightarrow \mathbb{R}^n$  is circular with modulo- $n$  without specification:  $(s_\ell[\mathbf{a}])_j = (\boldsymbol{\iota}_{[p]} \mathbf{a})_{j-\ell}$ . Notice  
 376 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- $n$ )  
 377 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  
 378  $\mathbf{b} \in \mathbb{R}^n$ ,

~~380~~ (3.21) 
$$\mathbf{C}_\mathbf{a} \mathbf{b} = \mathbf{a} * \mathbf{b}.$$

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

~~384~~ (3.22) 
$$\mathbf{C}_\mathbf{a}^* \mathbf{b} = \check{\mathbf{a}} * \mathbf{b}, \quad \check{\mathbf{C}}_\mathbf{a} \mathbf{b} = \mathbf{a} * \check{\mathbf{b}}.$$

386 **4. Geometry of  $\varphi_\rho$  in Shift Space.** Underlying our main geometric and algorithmic  
 387 results is a relationship between the geometry of the function  $\varphi_\rho$  and the symmetries of the

388 deconvolution problem. In this section, we describe this relationship at a more technical level,  
 389 by interpreting the gradient and hessian of the function  $\varphi_\rho$  in terms of the shifts  $s_\ell[\mathbf{a}_0]$  and  
 390 stating a key lemma which asserts that a certain neighborhood of the union of subspaces  $\Sigma_{4\theta p_0}$   
 391 can be decomposed into regions of negative curvature, strong gradient, and strong convexity  
 392 near the target solutions  $\pm s_\ell[\mathbf{a}_0]$ .

393 **4.1. Shifts and Correlations.** The set  $\Sigma_{4\theta p_0}$  is a union of subspaces. Any point  $\mathbf{a}$  in one  
 394 of these subspaces  $\mathcal{S}_\tau$  is a superposition of shifts of  $\mathbf{a}_0$ :

$$395 \quad (4.1) \quad \mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0].$$

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

$$397 \quad (4.2) \quad \mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] + \sum_{\ell \notin \tau} \alpha_\ell s_\ell[\mathbf{a}_0].$$

398 The vector  $\alpha$  can be viewed as the coefficients of a decomposition of  $\mathbf{a}$  into different shifts  
 399 of  $\mathbf{a}_0$ . This representation is not unique. For  $\mathbf{a}$  close to  $\mathcal{S}_\tau$ , we can choose a particular  $\alpha$  for  
 400 which  $\alpha_{\tau^c}$  is small, a notion that we will formalize below.

401 For convenience, we introduce a closely related vector  $\beta \in \mathbb{R}^n$ , whose entries are the inner  
 402 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  $\mathbf{C}_{\mathbf{a}_0}$  are the  
 403 shifts of  $\mathbf{a}_0$ , we can write

$$404 \quad (4.3) \quad \beta = \mathbf{C}_{\mathbf{a}_0}^* \iota \mathbf{a}$$

$$405 \quad (4.4) \quad = \mathbf{C}_{\mathbf{a}_0}^* \iota \mathbf{C}_{\mathbf{a}_0} \alpha =: \mathbf{M} \alpha.$$

407 The matrix  $\mathbf{M}$  is the Gram matrix of the truncated shifts:  $M_{ij} = \langle \iota^* s_i[\mathbf{a}_0], \iota^* s_j[\mathbf{a}_0] \rangle$ . When  $\mu$   
 408 is small, the off-diagonal elements of  $\mathbf{M}$  are small. In particular, on  $\mathcal{S}_\tau$  we may take  $\alpha_{\tau^c} = \mathbf{0}$ ,  
 409 and  $\beta \approx \alpha$ , in the sense that  $\beta_\tau \approx \alpha_\tau$  and the entries of  $\beta_{\tau^c}$  are small. For detailed elaboration,  
 410 see [Section SM2](#).

411 **4.2. Shifts and the Calculus of  $\varphi_{\ell^1}$ .** Our main geometric claims pertain to the function  
 412  $\varphi_\rho$ , which is based on a smooth sparsity surrogate  $\rho(\cdot) \approx \|\cdot\|_1$ . In this section, we sketch the  
 413 main ideas of the proof as if  $\rho(\cdot) = \|\cdot\|_1$ , by relating the geometry of the function  $\varphi_{\ell^1}$  to the  
 414 vectors  $\alpha, \beta$  introduced above. Working with  $\varphi_{\ell^1}$  simplifies the exposition; it is also faithful to  
 415 the structure of our proof, which relates the derivatives of the smooth function  $\varphi_\rho$  to similar  
 416 quantities associated with the nonsmooth function  $\varphi_{\ell^1}$ .

417 The function  $\varphi_{\ell^1}$  has a relatively simple closed form:

$$418 \quad (4.5) \quad \varphi_{\ell^1}(\mathbf{a}) = -\frac{1}{2} \|\mathcal{S}_\lambda[\check{\mathbf{y}} * \mathbf{a}]\|_2^2.$$

419 Here,  $\mathcal{S}_\lambda$  is the *soft thresholding operator*, which is defined for scalars  $t$  as

$$420 \quad (4.6) \quad \mathcal{S}_\lambda[t] = \text{sign}(t) \max\{|t| - \lambda, 0\},$$

422 and is extended to vectors by applying it elementwise. The operator  $\mathcal{S}_\lambda[\mathbf{x}]$  shrinks the elements  
 423 of  $\mathbf{x}$  towards zero. Small elements become identically zero, resulting in a sparse vector.

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

$$\begin{aligned} 429 \quad \nabla \varphi_{\ell^1}(\mathbf{a}) &= -\iota^* \mathbf{C}_{\mathbf{a}_0} \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \iota \mathbf{a} \right] \\ 430 &= -\iota^* \mathbf{C}_{\mathbf{a}_0} \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}_0} \boldsymbol{\beta} \right] \\ 431 \quad (4.7) \quad &= -\iota^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\chi}[\boldsymbol{\beta}], \end{aligned}$$

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

$$436 \quad (4.8) \quad -\nabla \varphi_{\ell^1}(\mathbf{a}) = \sum_{\ell} \boldsymbol{\chi}[\boldsymbol{\beta}]_{\ell} s_{\ell}[\mathbf{a}_0].$$

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

$$440 \quad (4.9) \quad \frac{1}{n\theta} \cdot \boldsymbol{\chi}[\boldsymbol{\beta}]_{\ell} \approx \begin{cases} \beta_{\ell} - \lambda, & \beta_{\ell} > \lambda \\ \beta_{\ell} + \lambda, & \beta_{\ell} < -\lambda \\ 0, & \text{otherwise} \end{cases}.$$

442 We show this rigorously below, in the proof of our main theorems. Here, we support this  
 443 claim pictorially, by plotting the  $\ell$ -th entry  $\boldsymbol{\chi}[\boldsymbol{\beta}]_{\ell}$  as  $\beta_{\ell}$  varies – see [Figure 10](#) (middle left)  
 444 and compare to [Figure 10](#) (left). Because  $\boldsymbol{\chi}[\boldsymbol{\beta}]$  suppresses small entries of  $\boldsymbol{\beta}$ , the strongest  
 445 contributions to  $-\nabla \varphi_{\ell^1}$  in (4.8) will come from shifts  $s_{\ell}[\mathbf{a}_0]$  with large  $\beta_{\ell}$ . *In particular, the*  
 446 *Euclidean gradient is large whenever there is a single preferred shift  $s_{\ell}[\mathbf{a}_0]$ , i.e., the largest*  
 447 *entry of  $\boldsymbol{\beta}$  is significantly larger than the second largest entry.*

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

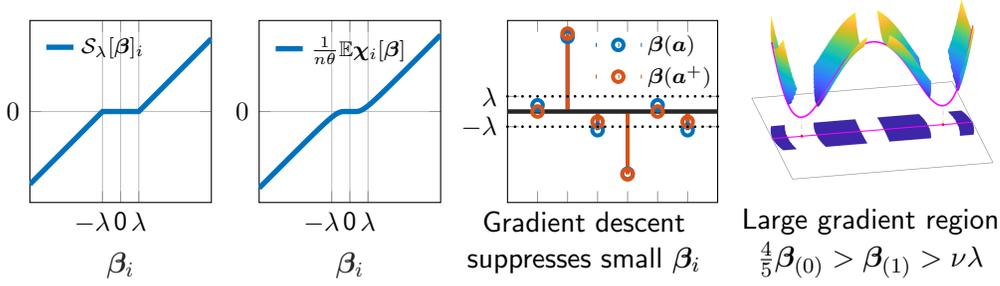
$$\begin{aligned} 450 \quad \text{grad}[\varphi_{\ell^1}](\mathbf{a}) &= \mathbf{P}_{\mathbf{a}^{\perp}} \nabla \varphi_{\ell^1}(\mathbf{a}) \\ 451 \quad (4.10) \quad &= -\mathbf{P}_{\mathbf{a}^{\perp}} \sum_{\ell} \boldsymbol{\chi}_{\ell}[\boldsymbol{\beta}] s_{\ell}[\mathbf{a}_0]. \end{aligned}$$

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

- 455 (i) **Negative gradient points to one particular shift:** there is a single preferred shift  
 456  $s_{\ell}[\mathbf{a}_0]$  so that the Euclidean gradient is large *and*
- 457 (ii)  **$\mathbf{a}$  is not too close to any shift:** it is possible to move in the tangent space in the  
 458 direction of this shift.<sup>15</sup> Since the tangent space consists of those vectors orthogonal to

---

<sup>15</sup>...so the projection of the Euclidean gradient onto the tangent space does not vanish.



**Figure 10. Gradient Sparsifies Correlations.** *Left: the soft thresholding operator  $S_\lambda[\beta]$  shrinks the entries of  $\beta$  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  $\chi_\ell[\beta] \approx S_\lambda[\beta]_\ell$ . Because of this, gradient descent sparsifies  $\beta$ . 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  $\beta$  are shrunk towards zero. Right: the gradient  $\text{grad}[\varphi_{\ell^1}](\mathbf{a})$  is large whenever it is easy to sparsify  $\beta$ ; in particular, when the largest entry  $\beta_{(0)} \gg \beta_{(1)} \gg 0$ .*

459  $\mathbf{a}$ , 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  
 460  $s_\ell[\mathbf{a}_0]$ .

461 Our technical lemma quantifies this situation in terms of the ordered entries of  $\beta$ . Write  
 462  $|\beta_{(0)}| \geq |\beta_{(1)}| \geq \dots$ , with corresponding shifts  $s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0], \dots$ . There is a strong gradient  
 463 whenever  $|\beta_{(0)}|$  is significantly larger than  $|\beta_{(1)}|$  and  $|\beta_{(1)}|$  is not too small compared to  $\lambda$ : in  
 464 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  
 465  $s_{(0)}[\mathbf{a}_0]$ , reducing  $|\beta_{(1)}|, \dots$ , and making the vector  $\beta$  sparser. We establish the technical claim  
 466 that the (Euclidean) gradient of  $\varphi_{\ell^1}$  sparsifies vectors in shift space in [Section SM3](#).

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

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

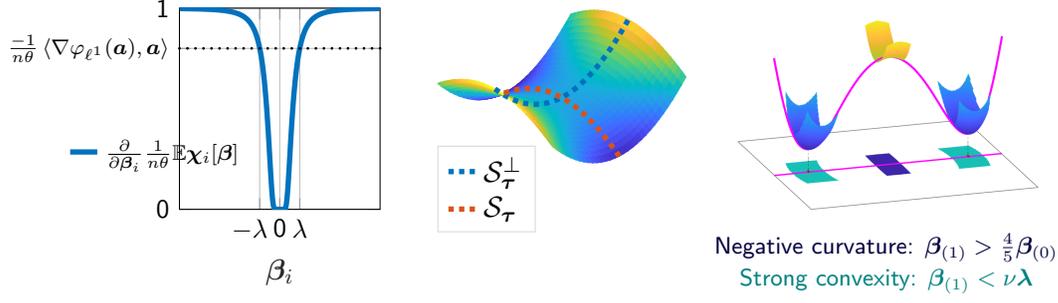
$$475 \quad (4.11) \quad \tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) = -\iota^* \mathbf{C}_{\mathbf{a}_0} \check{\mathbf{C}}_{x_0} \mathbf{P}_I \check{\mathbf{C}}_{x_0}^* \mathbf{C}_{\mathbf{a}_0}^* \iota,$$

477 with  $I = \text{supp} \left( S_\lambda \left[ \check{\mathbf{C}}_{\mathbf{y}} \iota \mathbf{a} \right] \right)$ . We (formally) extend this expression to *every*  $\mathbf{a} \in \mathbb{R}^n$ , terming  
 478  $\tilde{\nabla}^2 \varphi_{\ell^1}$  the *pseudo-Hessian* of  $\varphi_{\ell^1}$ . For appropriately chosen smooth sparsity surrogate  $\rho$ , we  
 479 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  $\tilde{\nabla}^2 \varphi_{\ell^1}$   
 480 yields useful information about the curvature of  $\varphi_\rho$ .

481 As with the gradient, the Hessian is complicated, but becomes simpler when the sample  
 482 size is large. The following approximation

$$483 \quad (4.12) \quad \tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \approx - \sum_{\ell} s_\ell[\mathbf{a}_0] s_\ell[\mathbf{a}_0]^* \left( \frac{\partial}{\partial \beta_\ell} \chi_\ell[\beta] \right)$$

484



**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  $\mathcal{S}_\tau$  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.

485 can be obtained from (4.8) 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  
 486  $j \neq \ell$ , and that

$$487 \quad (4.13) \quad \frac{1}{n\theta} \cdot \frac{\partial \chi_\ell[\beta]}{\partial \beta_\ell} \approx \begin{cases} 0 & |\beta_\ell| \ll \lambda \\ 1 & |\beta_\ell| \gg \lambda \end{cases}$$

488 Again, we corroborate this approximation pictorially – see Figure 11.

489 From this approximation, we can see that the quadratic form  $\mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1} \mathbf{v}$  takes on a large  
 490 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  
 491 linear combination of such shifts. In particular, if for some  $j$ ,  $|\beta_{(0)}|, |\beta_{(1)}|, \dots, |\beta_{(j)}| \gg \lambda$ , then  
 492  $\varphi_{\ell^1}$  will exhibit negative curvature in any direction  $\mathbf{v} \in \text{span}(s_{(0)}[\mathbf{a}_0], s_{(1)}[\mathbf{a}_0], \dots, s_{(j)}[\mathbf{a}_0])$ .

493 The (Euclidean) Hessian measures the curvature of the function  $\varphi_{\ell^1}$  over  $\mathbb{R}^n$ . The Rie-  
 494 mannian Hessian

$$495 \quad (4.14) \quad \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) = \mathbf{P}_{\mathbf{a}^\perp} \left( \begin{array}{c} \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \\ \text{Curvature of } \varphi_{\ell^1} \end{array} + \begin{array}{c} \langle -\nabla \varphi_{\ell^1}(\mathbf{a}), \mathbf{a} \rangle \cdot \mathbf{I} \\ \text{Curvature of the sphere} \end{array} \right) \mathbf{P}_{\mathbf{a}^\perp}.$$

496 measures the curvature of  $\varphi_{\ell^1}$  over the sphere. The projection  $\mathbf{P}_{\mathbf{a}^\perp}$  restricts its action to  
 497 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  
 498 for the curvature of the sphere. This term is always positive. The net effect is that directions  
 499 of strong negative curvature of  $\varphi_{\ell^1}$  over  $\mathbb{R}^n$  become directions of moderate negative curvature  
 500 over the sphere. Directions of nearly zero curvature over  $\mathbb{R}^n$  become directions of positive  
 501 curvature over the sphere. This has three implications for the geometry of  $\varphi_{\ell^1}$  over the sphere:

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

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

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

503 (ii) **Positive curvature in directions away from  $\mathcal{S}_\tau$** : The Euclidean Hessian  
 504 quadratic form  $\mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1} \mathbf{v}$  takes on relatively small values in directions orthogonal to  
 505 the subspace  $\mathcal{S}_\tau$ . The Riemannian Hessian is positive in these directions, creating  
 506 positive curvature orthogonal to the subspace  $\mathcal{S}_\tau$ .  
 507 (iii) **Strong convexity around minimizers**: Around a minimizer  $s_\ell[\mathbf{a}_0]$ , only a single  
 508 entry  $\beta_\ell$  is large. Any tangent direction  $\mathbf{v} \perp \mathbf{a}$  is nearly orthogonal to the subspace  
 509  $\text{span}(s_\ell[\mathbf{a}_0])$ , and hence is a direction of positive (Riemannian) curvature. The objective  
 510 function  $\varphi_\rho$  is strongly convex around the target solutions  $\pm s_\ell[\mathbf{a}_0]$ .  
 511 Figure 11 visualizes these regions of negative and positive curvature, and the technical claim  
 512 of positivity/negativity of curvature in shift space is presented in detail in Section SM4.

513 **4.3. Any Local Minimizer is a Near Shift.** We close this section by stating a key theorem,  
 514 which makes the above discussion precise. We will show that a certain neighborhood of any  
 515 subspace  $\mathcal{S}_\tau$  can be covered by regions of *negative curvature*, *large gradient*, and regions of  
 516 *strong convexity* containing target solutions  $\pm s_\ell[\mathbf{a}_0]$ . Furthermore, at the boundary of this  
 517 neighborhood, the negative gradient points back—*retracts*—toward the subspace  $\mathcal{S}_\tau$ , due to  
 518 the (directional) convexity of  $\varphi_\rho$  away from the subspace.

519 To formally state the result, we need a way of measuring how close  $\mathbf{a}$  is to the subspace  
 520  $\mathcal{S}_\tau$ . For technical reasons, it turns out to be convenient to do this in terms of the coefficients  
 521  $\boldsymbol{\alpha}$  in the representation

$$522 \quad (4.15) \quad \mathbf{a} = \sum_{\ell \in \tau} \alpha_\ell s_\ell[\mathbf{a}_0] + \sum_{\ell' \in \tau^c} \alpha_{\ell'} s_{\ell'}[\mathbf{a}_0].$$

523 If  $\mathbf{a} \in \mathcal{S}_\tau$ , we can take  $\boldsymbol{\alpha}$  with  $\alpha_{\tau^c} = \mathbf{0}$ . We can view the energy  $\|\boldsymbol{\alpha}_{\tau^c}\|_2$  as a measure of the  
 524 distance from  $\mathbf{a}$  to  $\mathcal{S}_\tau$ . A technical wrinkle arises, because the representation (4.15) is not  
 525 unique. We resolve this issue by choosing the  $\boldsymbol{\alpha}$  that minimizes  $\|\boldsymbol{\alpha}_{\tau^c}\|_2$ , writing:

$$526 \quad (4.16) \quad d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = \inf \{ \|\boldsymbol{\alpha}_{\tau^c}\|_2 : \sum_\ell \alpha_\ell s_\ell[\mathbf{a}_0] = \mathbf{a} \}.$$

528 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  
 529  $\varphi_\rho$  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  
 530 which contains all points that are close to some  $\mathcal{S}_\tau$ , in the above sense:

$$531 \quad (4.17) \quad \Sigma_{4\theta p_0}^\gamma := \bigcup_{|\tau| \leq 4\theta p_0} \{ \mathbf{a} : d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma \}.$$

533 The aforementioned geometric properties hold over this set:

534 **Theorem 4.1 (Geometry of  $\varphi_\rho$  over UoS).** *Suppose that  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$  where  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$  is*  
 535  *$\mu$ -shift coherent and  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta) \in \mathbb{R}^n$  satisfying*

$$536 \quad (4.18) \quad \theta \in \left[ \frac{c'}{p_0}, \frac{c}{p_0 \sqrt{\mu} + \sqrt{p_0}} \right] \cdot \frac{1}{\log^2 p_0}$$

538 *for some constants  $c', c > 0$ . Set  $\lambda = 0.1/\sqrt{p_0\theta}$  in  $\varphi_\rho$  where  $\rho(x) = \sqrt{x^2 + \delta^2}$ . There exist*  
 539 *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*  
 540 *with probability at least  $1 - c'''/n$ , for every  $\mathbf{a} \in \Sigma_{4\theta p_0}^\gamma$ , we have:*

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

$$543 \quad (4.19) \quad \lambda_{\min}(\text{Hess}[\varphi_\rho](\mathbf{a})) \leq -c_1 n \theta \lambda;$$

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

$$545 \quad (4.20) \quad \|\text{grad}[\varphi_\rho](\mathbf{a})\|_2 \geq c_2 n \theta \frac{\lambda^2}{\log^2 \theta^{-1}};$$

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

$$548 \quad (4.21) \quad \text{Hess}[\varphi_\rho](\mathbf{a}) \succ c_3 n \theta \mathbf{P}_{\mathbf{a}_\perp};$$

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

$$552 \quad (4.22) \quad \langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq c_4 \|\zeta_{\tau^c}\|_2 \|\alpha_{\tau^c}\|_2;$$

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

$$555 \quad (4.23) \quad \min_{\substack{\ell \in \{\pm p\} \\ \sigma \in \{\pm 1\}}} \|\mathbf{a} - \sigma s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max\{\mu, p_0^{-1}\},$$

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

558 *Proof.* See [Subsection SM6.5](#). ■

559 The retraction property elaborated in (4.22) implies that the negative gradient at  $\mathbf{a}$  points in  
560 a direction that decreases  $d_\alpha(\mathbf{a}, \mathcal{S}_\tau)$ . This is a consequence of positive curvature away from  $\mathcal{S}_\tau$ .  
561 It essentially implies that the gradient is monotone in  $\alpha_{\tau^c}$  space: choose any  $\underline{\mathbf{a}} \in \mathcal{S}_\tau \cap \mathbb{S}^{p-1}$ ,  
562 write  $\underline{\mathbf{a}}$  to be its coefficient, and let  $\underline{\zeta}$  be the coefficient of  $\text{grad}[\varphi_\rho](\underline{\mathbf{a}})$ . Then  $\underline{\alpha}_{\tau^c} = \mathbf{0}$ ,  $\underline{\zeta}_{\tau^c} \approx \mathbf{0}$   
563 and

$$564 \quad \langle \zeta_{\tau^c} - \underline{\zeta}_{\tau^c}, \alpha_{\tau^c} - \underline{\alpha}_{\tau^c} \rangle \approx \langle \zeta_{\tau^c} - \mathbf{0}, \alpha_{\tau^c} - \mathbf{0} \rangle = \langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle > 0.$$

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

569 **5. Provable Algorithm.** In light of [Theorem 4.1](#), in this section we introduce a two-part  
570 algorithm [Algorithm 3.1](#), which first applies the curvilinear descent method to find a local min-  
571 imum of  $\varphi_\rho$  within  $\Sigma_{4\theta p_0}^\gamma$ , followed by refinement algorithm that uses alternating minimization  
572 to exactly recover the ground truth. This algorithm exactly solves SaS deconvolution problem.

573 **5.1. Minimization.** There are three major issues in finding a local minimizer within  $\Sigma_{4\theta p_0}^\gamma$ .

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

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

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

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

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

$$\begin{aligned}
 587 \quad -\nabla \varphi_{\rho^1}(\mathbf{a}^{(-1)}) &\approx \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda [\check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_{[p_0]}(\mathbf{a}_0 * \tilde{\mathbf{x}}_0)] \\
 588 \quad &= \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathcal{X} [\mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_{[p_0]} \mathbf{C}_{\mathbf{a}_0} \tilde{\mathbf{x}}_0] \\
 589 \quad &\approx \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathcal{X} [\tilde{\mathbf{x}}_0] \\
 590 \quad &\approx n\theta \cdot \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathcal{S}_\lambda [\tilde{\mathbf{x}}_0],
 \end{aligned}$$

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

$$\begin{aligned}
 593 \quad (5.1) \quad &\max_{i \neq j} |\langle \boldsymbol{\iota}_{p_0}^* s_i[\mathbf{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu \ll 1. \\
 594
 \end{aligned}$$

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

$$\begin{aligned}
 597 \quad (5.2) \quad &\mathbf{a}^{(0)} \approx \mathbf{P}_{\mathbb{S}^{p-1}} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathcal{S}_\lambda [\tilde{\mathbf{x}}_0], \\
 598
 \end{aligned}$$

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

$$\begin{aligned}
 602 \quad (5.3) \quad &\mathbf{a}^{(0)} \in \Sigma_{4\theta p_0}^\gamma. \\
 603
 \end{aligned}$$

604 Furthermore, since  $\tilde{\mathbf{x}}_0$  is normalized, the largest magnitude for entries of  $|\tilde{\mathbf{x}}_0|$  is likely to be  
 605 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  
 606 our initializer  $\mathbf{a}^{(0)}$  will become  $\mathbf{0}$ ), the ideal  $\lambda$  should be slightly less than the largest magnitude  
 607 of  $|\tilde{\mathbf{x}}_0|$ . We suggest setting  $\lambda$  in  $\varphi_\rho$  as

$$\begin{aligned}
 608 \quad (5.4) \quad &\lambda = \frac{c}{\sqrt{p_0\theta}}. \\
 609
 \end{aligned}$$

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

611 Many methods have been proposed to optimize functions whose saddle points exhibit strict  
 612 negative curvature, including the noisy gradient method [22], trust region methods [1, 55] and  
 613 curvilinear search [58]. Any of the above methods can be adapted to minimize  $\varphi_\rho$ . In this  
 614 paper, we use *curvilinear method with restricted stepsize* to demonstrate how to analyze an  
 615 optimization problem using the geometric properties of  $\varphi_\rho$  over  $\Sigma_{4\theta p_0}^\gamma$  – in particular, negative  
 616 curvature in symmetry-breaking directions and positive curvature away from  $\mathcal{S}_\tau$ .

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

$$620 \quad (5.5) \quad \mathbf{a}^+ \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}]$$

622 For small  $t$ ,

$$623 \quad (5.6) \quad \varphi(\mathbf{a}^+) \approx \varphi(\mathbf{a}) + \langle \mathbf{g}, \boldsymbol{\xi} \rangle + \frac{1}{2} \boldsymbol{\xi}^* \mathbf{H} \boldsymbol{\xi}.$$

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

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

$$630 \quad (5.7) \quad |d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) - d_\alpha(\mathbf{a}, \mathcal{S}_\tau)| \leq \frac{\gamma}{2}.$$

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

$$633 \quad (5.8) \quad d_\alpha(\mathbf{a}^+, \mathcal{S}_\tau) \leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau).$$

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

$$636 \quad (5.9) \quad \forall k = 1, 2, \dots, \quad \mathbf{a}^{(k)} \in \Sigma_{4\theta p_0}^\gamma.$$

638 We conclude this section with the following theorem:

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

$$645 \quad (5.10) \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)}]$$

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

$$648 \quad (5.11) \quad \varphi_\rho(\mathbf{a}^{(k+1)}) < \varphi_\rho(\mathbf{a}^{(k)}) - \frac{1}{2} \left( t \|\mathbf{g}^{(k)}\|_2^2 + \frac{1}{2} t^4 \eta_v \|\mathbf{v}^{(k)}\|_2^2 \right),$$

650 *then with probability at least  $1 - 1/c$ , there exists some signed shift  $\bar{\mathbf{a}} = \pm s_i[\mathbf{a}_0]$  where  $i \in [\pm p_0]$*   
 651 *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*  
 652  *$c' < c_1$  in Theorem 4.1.*

653 *Proof.* See Subsection SM7.2. ■

654 **5.2. Local Refinement.** In this section, we describe and analyze an algorithm which  
 655 refines an estimate  $\bar{\mathbf{a}} \approx \mathbf{a}_0$  of the kernel to exactly recover  $(\mathbf{a}_0, \mathbf{x}_0)$ . Set

$$656 \quad (5.12) \quad \mathbf{a}^{(0)} \leftarrow \bar{\mathbf{a}}, \quad \lambda^{(0)} \leftarrow C(p\theta + \log n)(\mu + 1/p), \quad I^{(0)} \leftarrow \text{supp}(\mathcal{S}_\lambda[\mathbf{C}_{\bar{\mathbf{a}}}^* \mathbf{y}]).$$

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

$$659 \quad (5.13) \quad \mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|,$$

$$660 \quad (5.14) \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \underset{\mathbf{a}}{\text{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right],$$

$$661 \quad (5.15) \quad \lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}).$$

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

$$667 \quad (5.16) \quad \min_{\sigma \in \pm 1, \ell \in [\pm p_0]} \|\mathbf{a}^{(k)} - \sigma \cdot s_\ell[\mathbf{a}_0]\|_2 \leq C' 2^{-k}.$$

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

$$673 \quad (5.17) \quad \kappa_I := 6 \max \{ \theta p, \log n \},$$

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

$$676 \quad (5.18) \quad I = \text{supp}(\mathbf{x}_0), \quad I^{(k)} = \text{supp}(\mathbf{x}^{(k)}),$$

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

$$679 \quad (5.19) \quad \max_{\ell} |I \cap ([p] + \ell)| \leq \kappa_I.$$

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

$$684 \quad (5.20) \quad (\mu + 1/p) \cdot \kappa_I^2 < c$$

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

687 Recall that (4.23) in [Theorem 3.1](#) provides that

$$688 \quad (5.21) \quad \|\bar{\mathbf{a}} - \mathbf{a}_0\|_2 \leq (\mu + 1/p),$$

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

693 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  
 694 observe that the solution to the reweighted Lasso problem (5.13) can be written as

$$695 \quad (5.22) \quad \mathbf{x}^{(k+1)} = \boldsymbol{\iota}_J \left( \boldsymbol{\iota}_J^* \mathbf{C}_{\mathbf{a}^{(k)}}^* \mathbf{C}_{\mathbf{a}^{(k)}} \boldsymbol{\iota}_J \right)^{-1} \boldsymbol{\iota}_J^* \left( \mathbf{C}_{\mathbf{a}^{(k)}}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \lambda^{(k)} \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}^{(k+1)} \right),$$

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

$$698 \quad (5.23) \quad \mathbf{a}^{(k+1)} = \left( \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}^{(k+1)}} \boldsymbol{\iota} \right)^{-1} \left( \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 \right).$$

700 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  
 701 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  $I \approx J \approx T$ . Then  
 702 (5.22) gives

$$703 \quad (5.24) \quad \mathbf{x}^{(k+1)} \approx \mathbf{x}_0,$$

$$704 \quad (\mathbf{x}^{(k+1)} - \mathbf{x}_0) \approx \mathbf{P}_I \left( \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}^{(k)}} \mathbf{x}_0 \right)$$

$$705 \quad (5.25) \quad \approx \mathbf{P}_I \left[ \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}_0 - \mathbf{a}^{(k)}) \right],$$

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

$$708 \quad (\mathbf{a}^{(k+1)} - \mathbf{a}_0) \approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}^{(k+1)}}^* \mathbf{C}_{\mathbf{x}^{(k+1)}} \boldsymbol{\iota} \mathbf{a}_0$$

$$709 \quad \approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{x}_0 - \mathbf{x}^{(k+1)})$$

$$710 \quad (5.26) \quad \approx (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}^{(k)} - \mathbf{a}_0).$$

712 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  
 713 a contraction mapping with fixed point  $\mathbf{a}_0$ , as follows:

$$714 \quad (n\theta)^{-1} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \approx \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{e}_0 \mathbf{e}_0^* \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}$$

$$715 \quad (5.27) \quad \approx \mathbf{a}_0 \mathbf{a}_0^*.$$

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

722 **Theorem 5.2 (Linear rate convergence of alternating minimization).** *Suppose  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$*   
 723 *where  $\mathbf{a}_0$  is  $\mu$ -shift coherent and  $\mathbf{x}_0 \sim \text{BG}(\theta)$ , then there exists some constants  $C, c, c_\mu$  such*  
 724 *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*  
 725 *starting point  $\mathbf{a}^{(0)}$  and  $\lambda^{(0)}, I^{(0)}$  such that*

$$726 \quad (5.28) \quad \|\mathbf{a}^{(0)} - \mathbf{a}_0\|_2 \leq \mu + 1/p, \quad \lambda^{(0)} = 5\kappa_I(\mu + 1/p), \quad I^{(0)} = \text{supp}(\mathbf{C}_{\mathbf{a}^{(0)}}^* \mathbf{y}),$$

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

$$729 \quad (5.29) \quad \mathbf{x}^{(k+1)} \leftarrow \operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\mathbf{x}_i|,$$

$$730 \quad (5.30) \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \operatorname{argmin}_{\mathbf{a}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right],$$

$$731 \quad (5.31) \quad \lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \operatorname{supp}(\mathbf{x}^{(k+1)})$$

733 then

$$734 \quad (5.32) \quad \|\mathbf{a}^{(k+1)} - \mathbf{a}_0\|_2 \leq (\mu + 1/p) 2^{-k}$$

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

737 *Proof.* See [Subsection SM8.3](#). ■

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

739 **6. Experiments.** We demonstrate that the tradeoffs between the motif length  $p_0$  and  
 740 sparsity rate  $\theta$  produce a transition region for successful SaS deconvolution under generic  
 741 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  
 742 instances of synthetic data by choosing  $\mathbf{a}_0 \sim \operatorname{Unif}(\mathbb{S}^{p_0-1})$  and  $\mathbf{x}_0 \in \mathbb{R}^n$  with  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \operatorname{BG}(\theta)$   
 743 where  $n = 5 \times 10^5$ . Note that choosing  $\mathbf{a}_0$  this way implies  $\mu(\mathbf{a}_0) \approx \frac{1}{\sqrt{p_0}}$ .

744 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  
 745 ease of computation, we modify [Algorithm 3.1](#) by replacing curvilinear search with *accelerated*  
 746 *Riemannian gradient descent* method ([Algorithm 6.1](#)), which is an adaptation of accelerated  
 747 gradient descent [5] to the sphere. In particular, we apply momentum and increment by the  
 748 Riemannian gradient via the exponential and logarithmic operators

$$749 \quad (6.1) \quad \operatorname{Exp}_{\mathbf{a}}(\mathbf{u}) := \cos(\|\mathbf{u}\|_2) \cdot \mathbf{a} + \sin(\|\mathbf{u}\|_2) \cdot \frac{\mathbf{u}}{\|\mathbf{u}\|_2},$$

$$750 \quad (6.2) \quad \operatorname{Log}_{\mathbf{a}}(\mathbf{b}) := \arccos(\langle \mathbf{a}, \mathbf{b} \rangle) \cdot \frac{\mathbf{P}_{\mathbf{a}^\perp}(\mathbf{b} - \mathbf{a})}{\|\mathbf{P}_{\mathbf{a}^\perp}(\mathbf{b} - \mathbf{a})\|_2},$$

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

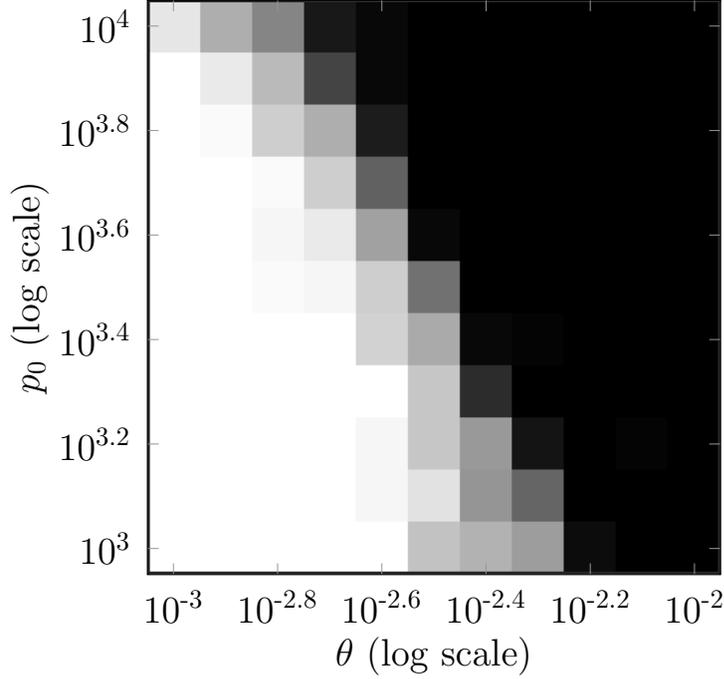
755 For each recovery instance, we say the local minimizer  $\mathbf{a}_{\min}$  generated from [Algorithm 6.1](#)  
 756 is sufficiently close to a solution of SaS deconvolution problem, if

$$757 \quad (6.3) \quad \operatorname{success}(\mathbf{a}_{\min}; \mathbf{a}_0) := \{ \max_{\ell} |\langle s_{\ell}[\mathbf{a}_0], \mathbf{a}_{\min} \rangle| > 0.95 \}.$$

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

759 [https://github.com/sbdsphere/sbd\\_experiments.git](https://github.com/sbdsphere/sbd_experiments.git)

760 **7. Discussion.** In this section, we close by discussing the most important limitations of our  
 761 results when  $\mathbf{a}_0$  is coherent, about scenarios when the signal setting breaches our assumption,  
 762 especially when  $\mathbf{x}_0$  is either highly sparse or non-symmetric, and highlighting corresponding  
 763 directions for future work.



**Figure 12.** Success probability of SaS deconvolution under generic  $\mathbf{a}_0, \mathbf{x}_0$  with varying kernel length  $p_0$ , and sparsity rate  $\theta$ . 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 6.1](#) works with high probability when  $\theta \lesssim \frac{1}{\sqrt{p_0}}$  in generic case.

---

**Algorithm 6.1** SaS deconvolution with Accelerated Riemannian gradient descent

---

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

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

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

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

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

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

**end for**

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

---

764     The main drawback of our proposed method is that it does not succeed when the target  
765 motif  $\mathbf{a}_0$  has shift coherence very close to 1. For instance, a common scenario in image blind  
766 deconvolution involves deblurring an image with a smooth, low-pass point spread function  
767 (e.g., Gaussian blur). Both our analysis and numerical experiments show that in this situation  
768 minimizing  $\varphi_\rho$  does not find the generating signal pairs  $(\mathbf{a}_0, \mathbf{x}_0)$  consistently—the minimizer of  
769  $\varphi_\rho$  is often spurious and is not close to any particular shift of  $\mathbf{a}_0$ . We do not suggest minimizing  
770  $\varphi_\rho$  in this situation. On the other hand, minimizing the bilinear lasso objective  $\varphi_{\text{lasso}}$  over the  
771 sphere often succeeds even if the true signal pair  $(\mathbf{a}_0, \mathbf{x}_0)$  is coherent and dense.

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

778 Our model assume  $\mathbf{x}_0$  to be Bernoulli-Gaussian vector, which are sparse and symmetric  
 779 iid random variables. When  $\mathbf{x}_0$  is sparse but non-symmetric, (e.g. Bernoulli), one can apply  
 780 our result with a simple symmetrization trick, by using the concatenated observation vectors  
 781  $[\mathbf{y}, -\mathbf{y}]$  as an input to our algorithm.

782 When  $\mathbf{x}_0$  is highly sparse and if  $\mathbf{y}$  is noiseless, it is possible to identify a short copy of  $\mathbf{a}_0$   
 783 via looking for a shortest consecutive non-zero entries within  $\mathbf{y}$ . When  $\theta \ll 1/p_0$ , these isolated  
 784 copies are very common. Once  $\theta$  exceeds  $1/p_0$ , or when support  $\mathbf{x}_0$  is not Bernoulli random  
 785 while being more clustered, they become very uncommon. In particular, the probability  
 786 of an isolated copy is small unless  $n \gtrsim \exp(p_0\theta)$ . Our proposed approach succeeds when  
 787  $n \geq \text{poly}(p_0)$ .

788 In applications involving noisy data, optimization approaches often outperform direct  
 789 inspection, even for samples with isolated copies of  $\mathbf{a}_0$ . An intuition for this is that optimization  
 790 methods aggregate information across the sample. One practical avenue for obtaining the best  
 791 of both worlds is to try to optimize the choice of data segment used for initialization. This can  
 792 be a potential improvement for our data-driven initialization scheme, both in theory and in  
 793 practice.

794 Finally, there are several directions in which our analysis could be improved. Our lower  
 795 bounds on the length  $n$  of the random vector  $\mathbf{x}_0$  required for success are clearly suboptimal. We  
 796 also suspect our sparsity-coherence tradeoff between  $\mu, \theta$  (roughly,  $\theta \lesssim 1/(\sqrt{\mu}p_0)$ ) is suboptimal,  
 797 even for the  $\varphi_p$  objective. Articulating optimal sparsity-coherence tradeoffs for is another  
 798 interesting direction in this line of work. Extending our current result for cases when  $\mathbf{y}$  is  
 799 affected by noise can also be a natural next step for future work.

800 **Acknowledgement.** The authors gratefully acknowledge support from NSF  
 801 1343282, NSF CCF 1527809, and NSF IIS 1546411.

802

## REFERENCES

- 803 [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton  
 804 University Press, 2009.
- 805 [2] A. AHMED, B. RECHT, AND J. ROMBERG, *Blind deconvolution using convex programming*, IEEE Transac-  
 806 tions on Information Theory, 60 (2014), pp. 1711–1732.
- 807 [3] G. AYERS AND J. C. DAINTY, *Iterative blind deconvolution method and its applications*, Optics letters, 13  
 808 (1988), pp. 547–549.
- 809 [4] S. BAKER AND T. KANADE, *Limits on super-resolution and how to break them*, IEEE Transactions on  
 810 Pattern Analysis and Machine Intelligence, 24 (2002), pp. 1167–1183.
- 811 [5] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*,  
 812 SIAM journal on imaging sciences, 2 (2009), pp. 183–202.
- 813 [6] A. J. BELL AND T. J. SEJNOWSKI, *An information-maximization approach to blind separation and blind*  
 814 *deconvolution*, Neural computation, 7 (1995), pp. 1129–1159.

- 815 [7] A. BENICHOX, E. VINCENT, AND R. GRIBONVAL, *A fundamental pitfall in blind deconvolution with*  
816 *sparse and shift-invariant priors*, in ICASSP-38th International Conference on Acoustics, Speech, and  
817 Signal Processing-2013, 2013.
- 818 [8] P. BONES, C. PARKER, B. SATHERLEY, AND R. WATSON, *Deconvolution and phase retrieval with use of*  
819 *zero sheets*, JOSA A, 12 (1995), pp. 1842–1857.
- 820 [9] D. BRIERS, D. D. DUNCAN, E. R. HIRST, S. J. KIRKPATRICK, M. LARSSON, W. STEENBERGEN,  
821 T. STROMBERG, AND O. B. THOMPSON, *Laser speckle contrast imaging: theoretical and practical*  
822 *limitations*, Journal of biomedical optics, 18 (2013), p. 066018.
- 823 [10] P. CAMPISI AND K. EGIAZARIAN, *Blind image deconvolution: theory and applications*, CRC press, 2016.
- 824 [11] E. J. CANDES, M. B. WAKIN, AND S. P. BOYD, *Enhancing sparsity by reweighted  $\ell_1$  minimization*,  
825 Journal of Fourier analysis and applications, 14 (2008), pp. 877–905.
- 826 [12] M. CANNON, *Blind deconvolution of spatially invariant image blurs with phase*, IEEE Transactions on  
827 Acoustics, Speech, and Signal Processing, 24 (1976), pp. 58–63.
- 828 [13] A. S. CARASSO, *Direct blind deconvolution*, SIAM Journal on Applied Mathematics, 61 (2001), pp. 1980–  
829 2007.
- 830 [14] T. F. CHAN AND C.-K. WONG, *Total variation blind deconvolution*, IEEE transactions on Image Processing,  
831 7 (1998), pp. 370–375.
- 832 [15] S. CHEUNG, Y. LAU, Z. CHEN, J. SUN, Y. ZHANG, J. WRIGHT, AND A. PASUPATHY, *Beyond the fourier*  
833 *transform: A nonconvex optimization approach to microscopy analysis*, Submitted, (2017).
- 834 [16] Y. CHI, *Guaranteed blind sparse spikes deconvolution via lifting and convex optimization*, IEEE Journal of  
835 Selected Topics in Signal Processing, 10 (2016), pp. 782–794.
- 836 [17] S. CHOUDHARY AND U. MITRA, *Fundamental limits of blind deconvolution part ii: Sparsity-ambiguity*  
837 *trade-offs*, arXiv preprint arXiv:1503.03184, (2015).
- 838 [18] W. DONG, L. ZHANG, G. SHI, AND X. WU, *Image deblurring and super-resolution by adaptive sparse*  
839 *domain selection and adaptive regularization*, IEEE Transactions on Image Processing, 20 (2011),  
840 pp. 1838–1857.
- 841 [19] B. EFRON, T. HASTIE, I. JOHNSTONE, R. TIBSHIRANI, ET AL., *Least angle regression*, The Annals of  
842 statistics, 32 (2004), pp. 407–499.
- 843 [20] C. EKANADHAM, D. TRANCHINA, AND E. P. SIMONCELLI, *A blind sparse deconvolution method for neural*  
844 *spike identification*, in Advances in Neural Information Processing Systems 24, 2011, pp. 1440–1448.
- 845 [21] R. FERGUS, B. SINGH, A. HERTZMANN, S. T. ROWEIS, AND W. T. FREEMAN, *Removing camera shake*  
846 *from a single photograph*, in ACM transactions on graphics (TOG), vol. 25, ACM, 2006, pp. 787–794.
- 847 [22] R. GE, F. HUANG, C. JIN, AND Y. YUAN, *Escaping from saddle points—online stochastic gradient for*  
848 *tensor decomposition*, in Conference on Learning Theory, 2015, pp. 797–842.
- 849 [23] D. GOLDFARB, *Curvilinear path steplength algorithms for minimization which use directions of negative*  
850 *curvature*, Mathematical programming, 18 (1980), pp. 31–40.
- 851 [24] D. GOLDFARB, C. MU, J. WRIGHT, AND C. ZHOU, *Using negative curvature in solving nonlinear programs*,  
852 Computational Optimization and Applications, 68 (2017), pp. 479–502.
- 853 [25] S. HARMELING, M. HIRSCH, S. SRA, AND B. SCHOLKOPF, *Online blind deconvolution for astronomical*  
854 *imaging*, in 2009 IEEE International Conference on Computational Photography (ICCP 2009), IEEE,  
855 2009, pp. 1–7.
- 856 [26] R. JOHNSON, P. SCHNITER, T. J. ENDRES, J. D. BEHM, D. R. BROWN, AND R. A. CASAS, *Blind*  
857 *equalization using the constant modulus criterion: A review*, Proceedings of the IEEE, 86 (1998),  
858 pp. 1927–1950.
- 859 [27] N. JOSHI, R. SZELISKI, AND D. J. KRIEGMAN, *Psf estimation using sharp edge prediction*, in Computer  
860 Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on, IEEE, 2008, pp. 1–8.
- 861 [28] K. F. KAARESEN AND T. TAXT, *Multichannel blind deconvolution of seismic signals*, Geophysics, 63  
862 (1998), pp. 2093–2107.
- 863 [29] M. KECH AND F. KRAHMER, *Optimal injectivity conditions for bilinear inverse problems with applications*  
864 *to identifiability of deconvolution problems*, SIAM Journal on Applied Algebra and Geometry, 1 (2017),  
865 pp. 20–37.
- 866 [30] D. KRISHNAN, T. TAY, AND R. FERGUS, *Blind deconvolution using a normalized sparsity measure*, in  
867 Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE, 2011, pp. 233–  
868 240.

- 869 [31] D. KUNDUR AND D. HATZINAKOS, *Blind image deconvolution*, IEEE signal processing magazine, 13 (1996),  
870 pp. 43–64.
- 871 [32] R. LANE AND R. BATES, *Automatic multidimensional deconvolution*, JOSA A, 4 (1987), pp. 180–188.
- 872 [33] R. G. LANE, *Blind deconvolution of speckle images*, JOSA A, 9 (1992), pp. 1508–1514.
- 873 [34] Y. LAU, Q. QU, H.-W. KUO, P. ZHOU, Y. ZHANG, AND J. WRIGHT, *Short-and-sparse deconvolution—a*  
874 *geometric approach*, arXiv preprint arXiv:1908.10959, (2019).
- 875 [35] A. LEVIN, R. FERGUS, F. DURAND, AND W. T. FREEMAN, *Deconvolution using natural image priors*,  
876 Massachusetts Institute of Technology, Computer Science and Artificial Intelligence Laboratory, 3  
877 (2007).
- 878 [36] A. LEVIN, Y. WEISS, F. DURAND, AND W. T. FREEMAN, *Understanding blind deconvolution algorithms*,  
879 IEEE transactions on pattern analysis and machine intelligence, 33 (2011), pp. 2354–2367.
- 880 [37] M. S. LEWICKI, *A review of methods for spike sorting: the detection and classification of neural action*  
881 *potentials*, Network: Computation in Neural Systems, 9 (1998), pp. R53–R78.
- 882 [38] Y. LI AND Y. BRESLER, *Global geometry of multichannel sparse blind deconvolution on the sphere*, arXiv  
883 preprint arXiv:1404.4104, (2018).
- 884 [39] Y. LI, K. LEE, AND Y. BRESLER, *Identifiability in blind deconvolution with subspace or sparsity constraints*,  
885 IEEE Transactions on Information Theory, 62 (2016), pp. 4266–4275.
- 886 [40] Y. LI, K. LEE, AND Y. BRESLER, *Identifiability and stability in blind deconvolution under minimal*  
887 *assumptions*, IEEE Transaction of Information Theory, (2017).
- 888 [41] S. LING AND T. STROHMER, *Self-calibration and biconvex compressive sensing*, Inverse Problems, 31  
889 (2015), p. 115002.
- 890 [42] S. LING AND T. STROHMER, *Blind deconvolution meets blind demixing: Algorithms and performance*  
891 *bounds*, IEEE Transactions on Information Theory, 63 (2017), pp. 4497–4520.
- 892 [43] J. MARKHAM AND J.-A. CONCHELLO, *Parametric blind deconvolution: a robust method for the simultaneous*  
893 *estimation of image and blur*, JOSA A, 16 (1999), pp. 2377–2391.
- 894 [44] M. MIYOSHI AND Y. KANEDA, *Inverse filtering of room acoustics*, IEEE Transactions on acoustics, speech,  
895 and signal processing, 36 (1988), pp. 145–152.
- 896 [45] P. A. NAYLOR AND N. D. GAUBITCH, *Speech dereverberation*, Springer Science & Business Media, 2010.
- 897 [46] M. R. OSBORNE, B. PRESNELL, AND B. A. TURLACH, *A new approach to variable selection in least*  
898 *squares problems*, IMA journal of numerical analysis, 20 (2000), pp. 389–403.
- 899 [47] D. PERRONE AND P. FAVARO, *Total variation blind deconvolution: The devil is in the details*, in Proceedings  
900 of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 2909–2916.
- 901 [48] E. A. PNEVMATIKAKIS, D. SOUDRY, Y. GAO, T. A. MACHADO, J. MEREL, D. PFAU, T. REARDON,  
902 Y. MU, C. LACEFIELD, W. YANG, ET AL., *Simultaneous denoising, deconvolution, and demixing of*  
903 *calcium imaging data*, Neuron, 89 (2016), pp. 285–299.
- 904 [49] S. K. SAHA, *Diffraction-limited imaging with large and moderate telescopes*, World Scientific, 2007.
- 905 [50] Y. SATO, *A method of self-recovering equalization for multilevel amplitude-modulation systems*, IEEE  
906 Transactions on communications, 23 (1975), pp. 679–682.
- 907 [51] O. SHALVI AND E. WEINSTEIN, *New criteria for blind deconvolution of nonminimum phase systems*  
908 *(channels)*, IEEE Transactions on information theory, 36 (1990), pp. 312–321.
- 909 [52] Q. SHAN, J. JIA, AND A. AGARWALA, *High-quality motion deblurring from a single image*, in Acm  
910 transactions on graphics (tog), vol. 27, ACM, 2008, p. 73.
- 911 [53] G. SHTENDEL, J. A. GALBRAITH, C. G. GALBRAITH, J. LIPPINCOTT-SCHWARTZ, J. M. GILLETTE,  
912 S. MANLEY, R. SOUGRAT, C. M. WATERMAN, P. KANCHANAWONG, M. W. DAVIDSON, ET AL.,  
913 *Interferometric fluorescent super-resolution microscopy resolves 3d cellular ultrastructure*, Proceedings  
914 of the National Academy of Sciences, 106 (2009), pp. 3125–3130.
- 915 [54] T. G. STOCKHAM, T. M. CANNON, AND R. B. INGEBRETSEN, *Blind deconvolution through digital signal*  
916 *processing*, Proceedings of the IEEE, 63 (1975), pp. 678–692.
- 917 [55] J. SUN, Q. QU, AND J. WRIGHT, *Complete dictionary recovery over the sphere ii: Recovery by riemannian*  
918 *trust-region method*, IEEE Transactions on Information Theory, 63 (2017), pp. 885–914.
- 919 [56] P. WALK, P. JUNG, G. E. PFANDER, AND B. HASSIBI, *Blind deconvolution with additional autocorrelations*  
920 *via convex programs*, arXiv preprint arXiv:1701.04890, (2017).
- 921 [57] L. WANG AND Y. CHI, *Blind deconvolution from multiple sparse inputs*, IEEE Signal Processing Letters,  
922 23 (2016), pp. 1384–1388.

- 923 [58] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Mathematical  
924 Programming, 142 (2013), pp. 397–434.
- 925 [59] D. WIPF AND H. ZHANG, *Revisiting bayesian blind deconvolution*, The Journal of Machine Learning  
926 Research, 15 (2014), pp. 3595–3634.
- 927 [60] L. XU AND J. JIA, *Two-phase kernel estimation for robust motion deblurring*, in European conference on  
928 computer vision, Springer, 2010, pp. 157–170.
- 929 [61] J. YANG, J. WRIGHT, T. S. HUANG, AND Y. MA, *Image super-resolution via sparse representation*, IEEE  
930 transactions on image processing, 19 (2010), pp. 2861–2873.
- 931 [62] Y.-L. YOU AND M. KAVEH, *Anisotropic blind image restoration*, in Image Processing, 1996. Proceedings.,  
932 International Conference on, vol. 2, IEEE, 1996, pp. 461–464.
- 933 [63] Y. ZHANG, H.-W. KUO, AND J. WRIGHT, *Structured local optima in sparse blind deconvolution*, arXiv  
934 preprint arXiv:1806.00338, (2018).
- 935 [64] Y. ZHANG, Y. LAU, H.-w. KUO, S. CHEUNG, A. PASUPATHY, AND J. WRIGHT, *On the global geometry  
936 of sphere-constrained sparse blind deconvolution*, in Proceedings of the IEEE Conference on Computer  
937 Vision and Pattern Recognition, 2017, pp. 4894–4902.

## SUPPLEMENTARY MATERIALS: Geometry and Symmetry in Short-and-Sparse Deconvolution\*

Han-Wen Kuo<sup>†</sup>, Yuqian Zhang<sup>‡</sup>, Yenson Lau<sup>†</sup>, and John Wright<sup>†§</sup>

**SM1. 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 SM1.1 (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*

$$(SM1.1) \quad ||I_0| - n\theta| \leq 2\sqrt{n\theta} \log \varepsilon^{-1}.$$

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

$$(SM1.2) \quad \forall t \in [2p] \setminus \{0\}, \quad \frac{1}{2}n\theta^2 \leq |I_0 \cap (I_0 + t)| \leq 2n\theta^2$$

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 SM10.4](#) 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 [\(SM1.2\)](#), let  $J_t := I_0 \cap (I_0 + t)$ . The cardinality of  $J_t$  is an inner product between shifts of  $\boldsymbol{\omega}$ :

$$(SM1.3) \quad |J_t| = \sum_{k \in [n]} \omega_k \omega_{k-t},$$

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

$$(SM1.4) \quad \begin{cases} J_{t1} = J_t \cap \mathcal{K}_1, & \mathcal{K}_1 := [n] \cap \{0, \dots, t-1, 2t, \dots, 3t-1, \dots\} \\ J_{t2} = J_t \cap \mathcal{K}_2, & \mathcal{K}_2 := [n] \cap \{t, \dots, 2t-1, 3t, \dots, 4t-1, \dots\} \end{cases}.$$

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

\*Submitted to the editors Jan/08/2019; revised Sep/20/2019.

**Funding:** This work was funded by NSF 1343282, NSF CCF 1527809, and NSF IIS 1546411

<sup>†</sup>Department of Electronic Engineering and Data Science Institute, Columbia University.

<sup>‡</sup>Department of Computer Science, Cornell University.

<sup>§</sup>Department of Applied Physics and Applied Mathematics, Columbia University.

Lemma SM10.4 with  $n > C\theta^{-2} \log p$  as

$$\begin{aligned}
\mathbb{P} \left[ \max_{t \in [2p] \setminus \{0\}} |J_{t_1}| \geq n\theta^2 \right] &= \mathbb{P} \left[ \max_{t \in [2p] \setminus \{0\}} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k-t} \geq n\theta^2 \right] \\
&\leq 2p \cdot \mathbb{P} \left[ \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \geq n\theta^2 \right] \\
&\leq 2p \cdot \mathbb{P} \left[ \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} - \mathbb{E} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \geq n\theta^2 - 0.6n\theta^2 \right] \\
&\leq 4p \cdot \exp \left( \frac{-(0.4n\theta^2)^2}{2 \cdot 0.6n\theta^2 + 2 \cdot 0.4n\theta^2} \right) = \exp(\log(4p) - 0.08n\theta^2) \\
\text{(SM1.5)} \quad &\leq 1/n,
\end{aligned}$$

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

$$\begin{aligned}
\mathbb{P} \left[ \min_{t \in [2p] \setminus \{0\}} |J_{t_1}| \leq n\theta^2/4 \right] &= \mathbb{P} \left[ \min_{t \in [2p] \setminus \{0\}} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k-t} \leq n\theta^2/4 \right] \\
&\leq 2p \cdot \mathbb{P} \left[ \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \leq n\theta^2/4 \right] \\
&\leq 2p \cdot \mathbb{P} \left[ \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} - \mathbb{E} \sum_{k \in \mathcal{K}_1} \omega_k \omega_{k+1} \leq n\theta^2/4 - 0.4n\theta^2 \right] \\
&\leq 4p \cdot \exp \left( \frac{-(0.15n\theta^2)^2}{2 \cdot 0.6n\theta^2 + 2 \cdot 0.15n\theta^2} \right) \\
&= \exp(\log(4p) - 0.0015n\theta^2) \leq 1/n.
\end{aligned}$$

The bound for  $|J_2|$  can be derived similarly to (SM1.5)-(1). ■

**Lemma SM1.2 (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  $1 - \varepsilon$ ,

$$\text{(SM1.6)} \quad \left| \|\mathbf{x}_0\|_1 - \sqrt{2/\pi}n\theta \right| \leq 2\sqrt{n\theta} \log \varepsilon^{-1}, \quad \left| \|\mathbf{x}_0\|_2^2 - n\theta \right| \leq 3\sqrt{n\theta} \log \varepsilon^{-1}$$

*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 SM10.2](#), 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 SM10.4](#),  $n\theta \geq 10$  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. ■

**Lemma SM1.3 (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*

$$(SM1.7) \quad \max_{\substack{U=[2p]+j \\ j \in [n]}} \|\mathbf{P}_U \mathbf{x}_0\|_2^2 \leq 2p\theta + 6 \left( \sqrt{p\theta} + \log n \right)$$

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

$$(SM1.8) \quad \max_{\substack{U=[p]+j \\ j \in [n]}} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 \leq p\theta + \log n.$$

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

$$(SM1.9) \quad \mathbb{P} \left[ \max_{\substack{U=[2p]+j \\ j \in [n]}} \|\mathbf{P}_U \mathbf{x}_0\|_2^2 \geq 2p\theta + 6 \left( \sqrt{p\theta} + \log n \right) \right] \leq 2n \exp \left( -\frac{36 \left( \sqrt{p\theta} + \log n \right)^2}{6p\theta + 12 \left( \sqrt{p\theta} + \log n \right)} \right) \\ \leq 2 \exp \left( \log n - \frac{36t^2}{6t^2 + 12t} \right) \leq \frac{2}{n}.$$

For the second inequality, first we know calculate the expectation

$$(SM1.10) \quad \mathbb{E} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 = \mathbb{E} [\mathbf{x}_0^* \mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0] \\ = \theta \cdot \text{tr} (\mathbf{C}_{\mathbf{a}_0}^* \mathbf{P}_U \mathbf{C}_{\mathbf{a}_0}) \|\mathbf{a}_0\|_2^2 + \theta \cdot \sum_{i=1}^{p-1} \|\boldsymbol{\iota}^* s_i[\mathbf{a}_0]\|_2^2 \\ = p\theta.$$

Then apply Henson Wright inequality [Lemma SM10.6](#) 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

$$(SM1.11) \quad \mathbb{P} \left[ \max_{\substack{U=[p]+j \\ j \in [n]}} \|\mathbf{P}_U [\mathbf{a}_0 * \mathbf{x}_0]\|_2^2 \geq p\theta + \log n \right] \leq n \exp \left( -\min \left\{ \frac{\log^2 n}{64\theta^2 p(1 + \mu p)}, \frac{\log n}{8\sqrt{2}\theta(1 + \mu p)} \right\} \right) \\ \leq \exp \left( \log n - \min \left\{ \frac{\log^2 n}{128c_\mu}, \frac{\log n}{32c_\mu} \right\} \right) \leq \frac{1}{n}$$

when  $c_\mu < \frac{1}{300}$ . ■

**Lemma SM1.4 (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:*

$$(SM1.12) \quad \max_{i \neq j \in [2p]} \langle s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle \leq 6\sqrt{n\theta^2 \log n};$$

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

$$(SM1.13) \quad \max_{i \neq j \in [2p]} \langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle \leq 4n\theta^2.$$

*Proof.* We will start from proving (SM1.13). 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 (SM1.13) can be written as summation of Gaussian r.v.s. on intersection of support set between shifts:

$$(SM1.14) \quad \langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle = \sum_{k \in I_0 \cap (I_0 + t)} |\mathbf{g}_k| |\mathbf{g}_{k-t}|$$

Define  $J_t := I_0 \cap (I_0 + t) = J_{t1} \uplus J_{t2}$  same as (SM1.4). 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  $\sum_{j \in J_{ti}} |\mathbf{g}_j| |\mathbf{g}'_j|$  where  $\mathbf{g}, \mathbf{g}'$  are independent Gaussian vectors. We condition on the following event

$$(SM1.15) \quad \mathcal{E}_J := \{\forall t \in [2p] \setminus \{0\}, n\theta^2/4 \leq |J_{t1}|, |J_{t2}| \leq n\theta^2\},$$

which holds w.p. at least  $1 - 2/n$  from Lemma SM1.1. 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 SM10.3 and union bound to obtain

$$(SM1.16) \quad \begin{aligned} \mathbb{P} \left[ \max_{t \in [2p] \setminus \{0\}} \sum_{j \in J_{t1}} |\mathbf{g}_j \mathbf{g}'_j| > 2|J_{t1}| \right] &\leq 2p \cdot \mathbb{P} \left[ \|\mathbf{g}_{J_{t1}}\|_2 \|\mathbf{g}'_{J_{t1}}\|_2 - \mathbb{E} \|\mathbf{g}_{J_{t1}}\|_2 \|\mathbf{g}'_{J_{t1}}\|_2 > |J_{t1}| \right] \\ &\leq 4p \cdot \mathbb{P} \left[ \|\mathbf{g}_{J_{t1}}\|_2 - \mathbb{E} \|\mathbf{g}_{J_{t1}}\|_2 > \sqrt{|J_{t1}|}/3 \right] \\ &\leq 4p \exp(-(|J_{t1}|/9)/2) \leq 4p \exp(-n\theta^2/72) \leq 1/n \end{aligned}$$

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  $i \neq j \in [2p]$ ,

$$(SM1.17) \quad \langle s_i[\mathbf{x}], s_j[\mathbf{x}] \rangle = \sum_{k \in J_{t1}} |\mathbf{g}_k| |\mathbf{g}'_{k-t}| + \sum_{k \in J_{t2}} |\mathbf{g}_k| |\mathbf{g}'_{k-t}| \leq 2(|J_{t1}| + |J_{t2}|) \leq 4n\theta^2.$$

For (SM1.12) similarly condition on event  $\mathcal{E}_J$ , using Bernstein inequality Lemma SM10.4 with  $(\sigma^2, R) = (1, 1)$ :

$$(SM1.18) \quad \mathbb{P} \left[ \max_{t \in [2p] \setminus \{0\}} \left| \sum_{j \in J_{t1}} \mathbf{g}_j \mathbf{g}'_j \right| > 3\sqrt{n\theta^2 \log n} \right] \leq p \cdot \exp \left( \frac{-9n\theta^2 \log n}{2|J_{t1}| + 6\sqrt{n\theta^2 \log n}} \right) \\ \leq p \cdot \exp \left( \frac{-9n\theta^2 \log n}{3n\theta^2} \right) \leq \frac{1}{n}$$

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

$$(SM1.19) \quad |\langle s_i[\mathbf{x}_0], s_j[\mathbf{s}_0] \rangle| \leq \left| \sum_{k \in J_{t1}} \mathbf{g}_k \mathbf{g}'_{k-t} \right| + \left| \sum_{k \in J_{t2}} \mathbf{g}_k \mathbf{g}'_{k-t} \right| \leq 6\sqrt{n\theta^2 \log n}.$$

Finally, both (SM1.17), (SM1.19) holds simultaneously with probability at least

$$(SM1.20) \quad 1 - 2/n - 1/n - 1/n = 1 - 4/n \quad \blacksquare$$

**Lemma SM1.5 (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  $\mu$ -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:*

$$(SM1.21) \quad \|\mathbf{C}_y \boldsymbol{\iota}\|_2^2 \leq 3(1 + \mu p)n\theta$$

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

$$(SM1.22) \quad \|\mathbf{P}_J \mathbf{C}_y \boldsymbol{\iota}\|_2^2 \leq 14|J|(1 + \mu p)(p\theta + \log n)$$

*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 SM1.2 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 SM1.4 we get:

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

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  $\|\mathbf{P}_J \mathbf{C}_y \boldsymbol{\iota} \mathbf{a}\|_2^2$ . Simply apply Lemma SM1.3 and utilize norm bound of  $\|\boldsymbol{\beta}\|_2^2$ , with probability at least  $1 - 2/n$  we have:

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

Finally apply Lemma SM2.4 and Gershgorin disc theorem obtain

$$(SM1.23) \quad \|\boldsymbol{\beta}\|_2^2 = \|\mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{a}\|_2^2 \leq \|\mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}\|_2^2 = \sigma_{\max}(\mathbf{M}) \leq 1 + \mu p.$$

**Remark SM1.6.** When  $\mathbf{a}_0$  is a basis vector  $\mathbf{e}_0$ , the result of Lemma SM1.5 gives upper bound of  $\|\mathbf{C}_{\mathbf{x}_0}\|_2 < 3n\theta$ , whose lower bound can be derived similarly with  $\|\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}\|_2 \geq \frac{2}{3}n\theta$

**SM2. Vectors in shift space.** In this section, we will establish a number of properties of the coefficient vectors  $\alpha$  and correlation vector  $\beta$ . Generally speaking, when  $\mathbf{a}$  is close to the subspace  $\mathcal{S}_\tau$ , then both vectors  $\alpha, \beta$  have most of their energy concentrated on the entries  $\tau$ . In this section, we derive upper bounds on  $\alpha_{\tau^c}$  and  $\beta_{\tau^c}$  under various assumptions.

In particular, we will introduce a relationship between the sparsity rate  $\theta$ , coherence  $\mu$  and size  $|\tau|$ , which we term the sparsity-coherence condition. In [Lemma SM2.2](#) we prove that measuring the distance from  $\mathbf{a}$  to subspace  $\mathcal{S}_\tau$  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  $\mathcal{S}_\tau$ . Later, in [Lemma SM2.4](#) we illustrate the relationship between  $\alpha$  and  $\beta$ , where  $\beta = C_{\mathbf{a}_0}^* \iota^* C_{\mathbf{a}_0} \alpha$ . Finally in [Lemma SM2.5](#) and [Corollary SM2.6](#), controls the magnitude of  $\alpha_{\tau^c}$  and  $\beta_{\tau^c}$  near  $\mathcal{S}_\tau$ .

**Definition SM2.1 (Sparsity-coherence condition).** Let  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$  with shift coherence  $\mu$ . We say that  $(\mathbf{a}_0, \theta, |\tau|)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$  with constant  $c_\mu$ , if

$$(SM2.1) \quad \theta \in \left[ \frac{1}{p}, \frac{c_\mu}{4 \max\{|\tau|, \sqrt{p}\}} \right] \cdot \frac{1}{\log^2 \theta^{-1}}, \quad \mu \cdot \max\{|\tau|^2, p^2 \theta^2\} \cdot \log^2 \theta^{-1} \leq \frac{c_\mu}{4},$$

where  $p = 3p_0 - 2$ .

**Lemma SM2.2 ( $d_\alpha$  is a seminorm).** For every solution subspace  $\mathcal{S}_\tau$ , the function  $d_\alpha(\cdot, \mathcal{S}_\tau) : \mathbb{R}^p \rightarrow \mathbb{R}_+$  defined as

$$(SM2.2) \quad d_\alpha(\mathbf{a}, \mathcal{S}_\tau) = \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha \}.$$

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  $d(\cdot, \mathcal{S}_\tau)$  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  $d_\alpha$ , for all  $\mathbf{a}_1, \mathbf{a}_2$  we have

$$\begin{aligned} d_\alpha(\mathbf{a}_1 + \mathbf{a}_2, \mathcal{S}_\tau) &= \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a}_1 + \mathbf{a}_2 = \iota^* C_{\mathbf{a}_0} \alpha \} \\ &= \inf \{ \|\alpha_{1\tau^c} + \alpha_{2\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{\mathbf{a}_0} \alpha_1, \quad \mathbf{a}_2 = \iota^* C_{\mathbf{a}_0} \alpha_2 \} \\ &\leq \inf \{ \|\alpha_{1\tau^c}\|_2 + \|\alpha_{2\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{\mathbf{a}_0} \alpha_1, \quad \mathbf{a}_2 = \iota^* C_{\mathbf{a}_0} \alpha_2 \} \\ &= \inf \{ \|\alpha_{1\tau^c}\|_2 \mid \mathbf{a}_1 = \iota^* C_{\mathbf{a}_0} \alpha_1 \} + \inf \{ \|\alpha_{2\tau^c}\|_2 \mid \mathbf{a}_2 = \iota^* C_{\mathbf{a}_0} \alpha_2 \} \\ &= d_\alpha(\mathbf{a}_1, \mathcal{S}_\tau) + d_\alpha(\mathbf{a}_2, \mathcal{S}_\tau). \end{aligned}$$

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

$$\begin{aligned} d_\alpha(c \cdot \mathbf{a}, \mathcal{S}_\tau) &= \inf \{ \|\alpha'_{\tau^c}\|_2 \mid c \cdot \mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha' \} = \inf \{ \|c \cdot \alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha \} \\ &= |c| \cdot \inf \{ \|\alpha_{\tau^c}\|_2 \mid \mathbf{a} = \iota^* C_{\mathbf{a}_0} \alpha \} = |c| \cdot d_\alpha(\mathbf{a}, \mathcal{S}_\tau), \end{aligned}$$

which completes the proof that  $d_\alpha$  is a seminorm. ■

**Definition SM2.3 (Widened subspace).** For subspace  $\mathcal{S}_\tau$  let

$$(SM2.3) \quad \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu)) := \{ \mathbf{a} \in \mathbb{S}^{p-1} \mid d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma \}$$

denote its widening by  $\gamma$ , in the seminorm  $d_\alpha$ .

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  $c_\mu$ , via

$$(SM2.4) \quad \gamma(c_\mu) = \frac{c_\mu}{4 \log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\tau|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\tau|} \right\}$$

**Lemma SM2.4 (Properties of  $\mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0}$ ).** *Let  $\mathbf{M} = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0}$ , with  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$   $\mu$ -shift coherent. The diagonal entries of  $\mathbf{M}$  satisfy*

$$(SM2.5) \quad \begin{cases} \mathbf{M}_{ii} = 1 & i \in [-p_0 + 1, p_0 - 1] = [\pm p_0], \\ 0 \leq \mathbf{M}_{ii} \leq 1 & i \in [-2p_0 + 2, -p_0] \cup [p_0, 2p_0 - 2], \\ \mathbf{M}_{ii} = 0 & \text{otherwise,} \end{cases}$$

and the off-diagonal entries satisfy

$$(SM2.6) \quad \begin{cases} |\mathbf{M}_{ij}| \leq \mu & 0 < |i - j| < p_0, \quad \{i \in [-p_0 + 1, p_0 - 1]\} \cup \{j \in [-p_0 + 1, p_0 - 1]\} \\ |\mathbf{M}_{ij}| < 1 & \{i, j \in [-2p_0 + 2, -p_0]\} \cup \{i, j \in [p_0, 2p_0 - 2]\} \\ 0 & \text{otherwise} \end{cases}.$$

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

$$(SM2.7) \quad \begin{cases} 1 - \mu |\tau| \leq \sigma_{\min}(\boldsymbol{\nu}_\tau^* \mathbf{M} \boldsymbol{\nu}_\tau) \leq \sigma_{\max}(\boldsymbol{\nu}_\tau^* \mathbf{M} \boldsymbol{\nu}_\tau) \leq 1 + \mu |\tau| \\ \sigma_{\max}(\boldsymbol{\nu}_{\tau^c}^* \mathbf{M} \boldsymbol{\nu}_{\tau^c}) \leq \mu \sqrt{p |\tau|} \\ \sigma_{\max}(\boldsymbol{\nu}_{\tau^c}^* \mathbf{M} \boldsymbol{\nu}_{\tau^c}) \leq 1 + \mu p \end{cases}$$

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

$$(SM2.8) \quad \left| \boldsymbol{\nu}_{[\pm 2p_0 - 1]}^* \mathbf{M} \boldsymbol{\nu}_{[\pm 2p_0 - 1]} \right| \leq \begin{bmatrix} \mathbf{J} & \mu \cdot \mathbf{1} & \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} & \mathbf{0} & \mathbf{0} \\ \mu \cdot \mathbf{1} & \mathbf{I} + \mu \cdot \mathbf{1}_o & \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix} & \mu \cdot \mathbf{1} & \mathbf{0} \\ [0 \cdots 0] & [\mu \cdots \mu] & 1 & [\mu \cdots \mu] & [0 \cdots 0] \\ \mathbf{0} & \mu \cdot \mathbf{1} & \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix} & \mathbf{I} + \mu \cdot \mathbf{1}_o & \mu \cdot \mathbf{1} \\ \mathbf{0} & \mathbf{0} & \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} & \mu \cdot \mathbf{1} & \mathbf{J} \end{bmatrix}.$$

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  $|\mathbf{J}_{ij}| < 1$  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:

$$(SM2.9) \quad \sigma_{\max}(\iota_{\tau^c}^* \mathbf{M} \iota_{\tau}) \leq \mu \sqrt{(2p_0 - 1)|\tau|} < \mu \sqrt{p|\tau|}.$$

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

$$(SM2.10) \quad \sigma_{\max}(\iota_{\tau^c}^* \mathbf{M} \iota_{\tau^c}) \leq \|\mathbf{C}_{a_0}^* \iota^* \mathbf{C}_{a_0}\|_2 = \|\iota^* \mathbf{C}_{a_0} \mathbf{C}_{a_0}^* \iota\|_2 = \|\iota^* \mathbf{C}_{a_0}^* \mathbf{C}_{a_0} \iota\|_2 \leq 1 + \mu p$$

where the third equality is derived via commutativity of convolution.  $\blacksquare$

**Lemma SM2.5 (Shift space vectors in widened subspace).** *Let  $(\mathbf{a}_0, \theta, |\tau|)$  satisfy the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Then for every  $\mathbf{a} \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ , every  $\alpha$  satisfying  $\mathbf{a} = \iota^* \mathbf{C}_{a_0} \alpha$  and  $\|\alpha_{\tau^c}\|_2 \leq \gamma(c_\mu)$  has*

$$(SM2.11) \quad \left| \|\alpha_\tau\|_2 - 1 \right| \leq c_\mu;$$

moreover,  $\beta = \mathbf{C}_{a_0}^* \iota \alpha$  satisfies

$$(SM2.12) \quad 1 - 3c_\mu \leq \|\beta_\tau\|_2^2 \leq 1 + \frac{c_\mu}{|\tau| \log^2 \theta^{-1}}$$

$$(SM2.13) \quad \|\beta_{\tau^c}\|_\infty \leq \frac{c_\mu}{\sqrt{|\tau|} \log^2 \theta^{-1}}$$

$$(SM2.14) \quad \|\beta_{\tau^c}\|_2 \leq \frac{c_\mu}{|\tau| \theta \log \theta^{-1}} \min \left\{ \sqrt{\theta}, \gamma(c_\mu) \right\}.$$

*Proof.* Write  $-1/\log \theta = \theta_{\log}$  and  $\gamma = \gamma(c_\mu)$  for convenience. First, by using bounds on  $\gamma$  in (SM2.4) and  $\mu|\tau| < 1$  we obtain:

$$(SM2.15) \quad \begin{cases} \gamma \cdot \sqrt{1 + \mu p} \leq \gamma(1 + \sqrt{\mu p}) \leq c_\mu \theta_{\log}^2 / 2 \\ \gamma \cdot \sqrt{1 + \mu^2 p} \leq \gamma(1 + \sqrt{\mu^2 p}) \leq \frac{c_\mu \theta_{\log}^2}{4} \left( \frac{1}{\sqrt{|\tau|}} + \sqrt{\mu} \right) \leq \frac{c_\mu \theta_{\log}^2}{2\sqrt{|\tau|}} \\ \gamma \cdot \mu \sqrt{p|\tau|} \leq \gamma \cdot \sqrt{\mu p} \cdot \sqrt{\mu|\tau|} \leq c_\mu \theta_{\log}^2 / 4 \end{cases}$$

Let  $\mathbf{a} = \iota^* \mathbf{C}_{a_0} \alpha$  with  $\|\alpha_{\tau^c}\|_2 < \gamma$ . Utilize properties of  $\iota^* \mathbf{C}_{a_0}$  from Lemma SM2.4 and  $\mu|\tau| < c_\mu/4$  and (SM2.15), we have:

$$(SM2.16) \quad \begin{aligned} \|\alpha_\tau\|_2 &\geq \|\iota^* \mathbf{C}_{a_0} \iota_\tau\|_2^{-1} (\|\mathbf{a}\|_2 - \|\iota^* \mathbf{C}_{a_0} \alpha_{\tau^c}\|_2) \\ &\geq \|\iota^* \mathbf{C}_{a_0} \iota_\tau\|_2^{-1} (1 - \|\iota^* \mathbf{C}_{a_0}\|_2 \|\alpha_{\tau^c}\|_2) \\ &\geq \frac{1}{\sqrt{1 + \mu|\tau|}} \left( 1 - \gamma \cdot \sqrt{1 + \mu p} \right) \geq \frac{1 - c_\mu/2}{\sqrt{1 + c_\mu/4}} \geq 1 - c_\mu, \end{aligned}$$

and similarly, the upper bound can be derived as:

$$\begin{aligned}
 \|\alpha_\tau\|_2 &\leq \sigma_{\min}^{-1}(\iota^* \mathbf{C}_{a_0} \iota_\tau) (\|\mathbf{a}\|_2 + \|\iota^* \mathbf{C}_{a_0} \alpha_{\tau^c}\|_2) \\
 &\leq \sigma_{\min}^{-1}(\iota^* \mathbf{C}_{a_0} \iota_\tau) (1 + \|\iota^* \mathbf{C}_{a_0}\|_2 \|\alpha_{\tau^c}\|_2) \\
 \text{(SM2.17)} \quad &\leq \frac{1}{\sqrt{1 - \mu|\tau|}} \left(1 + \gamma \cdot \sqrt{1 + \mu p}\right) \leq \frac{1 + c_\mu/2}{\sqrt{1 - c_\mu/4}} \leq 1 + c_\mu.
 \end{aligned}$$

The bound of  $\|\beta_\tau\|_2^2$  can be simply obtained using  $\mu|\tau| < c_\mu/4$  and  $\gamma$  bound from (SM2.15) as:

$$\begin{aligned}
 \text{(SM2.18)} \quad \|\beta_\tau\|_2^2 &\leq \sigma_{\max}^2(\iota_\tau^* \mathbf{C}_{a_0} \iota) \leq 1 + \mu|\tau| \leq 1 + \frac{c_\mu \theta_{\log}^2}{|\tau|} \\
 \|\beta_\tau\|_2^2 &\geq (\sigma_{\min}(\iota_\tau^* \mathbf{M} \iota_\tau) \|\alpha_\tau\|_2 - \sigma_{\max}(\iota_\tau^* \mathbf{M} \iota_{\tau^c}) \|\alpha_{\tau^c}\|_2)^2 \\
 \text{(SM2.19)} \quad &\geq \left( (1 - \mu|\tau|)(1 - c_\mu) - \mu\sqrt{p|\tau|} \cdot \gamma \right)^2 \geq 1 - 3c_\mu.
 \end{aligned}$$

As for the upper bound of and  $\|\beta_{\tau^c}\|_\infty$ , follow from (SM2.15), we have:

$$\begin{aligned}
 \|\beta_{\tau^c}\|_\infty &\leq \|\iota_{\tau^c}^* \mathbf{M} \alpha_\tau\|_\infty + \|\iota_{\tau^c}^* \mathbf{M} \alpha_{\tau^c}\|_\infty \leq \mu\sqrt{|\tau|} \|\alpha_\tau\|_2 + \sqrt{1 + \mu^2 p} \|\alpha_{\tau^c}\|_2 \\
 \text{(SM2.20)} \quad &\leq \frac{c_\mu \theta_{\log}^2 (1 + c_\mu)}{4|\tau|} + \gamma \cdot \sqrt{1 + \mu^2 p} \leq \frac{c_\mu \theta_{\log}^2}{\sqrt{|\tau|}};
 \end{aligned}$$

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

$$\text{(SM2.21)} \quad \|\beta_{\tau^c}\|_2 \leq \|\iota_{\tau^c}^* \mathbf{M} \alpha_\tau\|_2 + \|\iota_{\tau^c}^* \mathbf{M} \alpha_{\tau^c}\|_2 \leq \mu\sqrt{p|\tau|} \|\alpha_\tau\|_2 + (1 + \mu p) \|\alpha_{\tau^c}\|_2,$$

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  $\gamma$  and  $\theta|\tau| \leq c_\mu \theta_{\log}^2 / 4$  we have:

$$\begin{aligned}
 \text{(SM2.21)} &\leq \frac{\mu\sqrt{p\theta} |\tau|^{3/2}}{\sqrt{\theta} |\tau|} (1 + c_\mu) + \frac{\sqrt{\theta} |\tau| \cdot \sqrt{|\tau|} \gamma}{\sqrt{\theta} |\tau|} + \frac{\mu p \sqrt{\theta} |\tau| \gamma}{\sqrt{\theta} |\tau|} \\
 \text{(SM2.22)} &\leq \frac{2c_\mu \theta_{\log}^2 + c_\mu \theta_{\log}^3 + c_\mu \theta_{\log}^2}{4\sqrt{\theta} |\tau|} \leq \frac{c_\mu \theta_{\log}^2}{\sqrt{\theta} |\tau|},
 \end{aligned}$$

and similarly for the second inequality, use both conditions of  $\mu$ , we have:

$$\begin{aligned}
& \text{(SM2.21)} \\
& \leq \frac{\gamma}{\theta |\boldsymbol{\tau}|} \cdot \frac{\mu \sqrt{p} \theta |\boldsymbol{\tau}|^{3/2}}{\gamma} (1 + c_\mu) + \gamma + \mu p \gamma \\
& \leq \frac{\gamma}{\theta |\boldsymbol{\tau}|} \cdot \frac{4\mu \sqrt{p} \theta |\boldsymbol{\tau}|^{3/2}}{c_\mu \theta_{\log}^2} \cdot \max \left\{ \sqrt{|\boldsymbol{\tau}|}, \sqrt{\mu p}, \mu p \sqrt{\theta} |\boldsymbol{\tau}| \right\} \\
& \quad + \frac{\gamma}{\theta |\boldsymbol{\tau}|} \cdot \theta |\boldsymbol{\tau}| + \frac{\gamma}{\theta |\boldsymbol{\tau}|} \cdot \mu p \theta |\boldsymbol{\tau}| \\
& \leq \frac{\gamma}{\theta |\boldsymbol{\tau}|} \cdot \left( \frac{4}{c_\mu \theta_{\log}^2} \cdot \max \left\{ \mu |\boldsymbol{\tau}|^2 \cdot \sqrt{p} \theta, \mu (p\theta) |\boldsymbol{\tau}| \cdot \sqrt{\mu} |\boldsymbol{\tau}|, \right. \right. \\
& \quad \left. \left. \mu \sqrt{p} \theta |\boldsymbol{\tau}|^{3/2} \cdot \mu p \theta |\boldsymbol{\tau}| \right\} + \frac{c_\mu \theta_{\log}^2}{4} + \frac{c_\mu \theta_{\log}^2}{4} \right) \\
\text{(SM2.23)} \quad & \leq \frac{\gamma}{\theta |\boldsymbol{\tau}|} \left( \frac{c_\mu \theta_{\log}}{4} + \frac{c_\mu \theta_{\log}^2}{4} + \frac{c_\mu \theta_{\log}^2}{4} \right) \leq \frac{c_\mu \theta_{\log} \gamma}{\theta |\boldsymbol{\tau}|},
\end{aligned}$$

which completes the proof. ■

**Corollary SM2.6** ( $|\langle \boldsymbol{\beta}_{\boldsymbol{\tau}^c}, \mathbf{x}_{0, \boldsymbol{\tau}^c} \rangle|$  is small). *Given  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$  and  $|\boldsymbol{\tau}|, c_\mu$  such that  $(\mathbf{a}_0, \theta, |\boldsymbol{\tau}|)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Write  $\lambda = c_\lambda / \sqrt{|\boldsymbol{\tau}|}$  with some  $c_\lambda \geq 1/5$ , then if  $c_\mu \leq \frac{c_\lambda}{25}$ ,*

$$\text{(SM2.24)} \quad \mathbb{P} \left[ \left| \sum_{i \in \boldsymbol{\tau}^c} \boldsymbol{\beta}_i \mathbf{x}_{0i} \right| > \frac{\lambda}{10} \right] \leq 2\theta, \quad \mathbb{P} \left[ \left| \sum_i \boldsymbol{\beta}_i \mathbf{x}_{0i} \right| > \frac{\lambda}{10} \right] \leq \theta |\boldsymbol{\tau}| + 2\theta.$$

*Proof.* We bound tail probability of the first result with Gaussian moments [Lemma SM10.2](#) and Bernstein inequality [Lemma SM10.4](#). Via Hölder's inequality,  $\sum_{i \in \boldsymbol{\tau}^c} \mathbb{E}(\beta_i x_i)^q = \mathbb{E} x_0^q \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_q^q \leq \theta (q-1)!! \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_2^2 \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_\infty^{q-2}$ , thus

$$\text{(SM2.25)} \quad \mathbb{P} \left[ \left| \sum_{i \in \boldsymbol{\tau}^c} \boldsymbol{\beta}_i \mathbf{x}_{0i} \right| > \lambda/10 \right] \leq 2 \exp \left( \frac{-(\lambda/10)^2}{2\theta \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_2^2 + 2(\lambda/10) \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_\infty} \right)$$

Write  $\theta_{\log} = -\frac{1}{\log \theta}$ , [Lemma SM2.5](#) implies when  $c_\mu \leq \frac{c_\lambda}{25}$ , we have  $\theta \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_2^2 \leq \frac{c_\mu^2 \theta_{\log}^2}{|\boldsymbol{\tau}|^2} \leq \frac{\theta_{\log} \lambda^2}{625}$  and  $\|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_\infty \leq \frac{c_\mu \theta_{\log}}{\sqrt{|\boldsymbol{\tau}|}} \leq \frac{\theta_{\log} \lambda}{25}$ , therefore,

$$\begin{aligned}
\text{(SM2.25)} & \leq 2 \exp \left( \frac{-\lambda^2/100}{2\theta_{\log} \lambda^2/625 + 2(\theta_{\log} \lambda/25) \cdot (\lambda/10)} \right) \\
\text{(SM2.26)} & \leq 2 \exp(\log \theta) \leq 2\theta
\end{aligned}$$

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

$$\begin{aligned}
 \mathbb{P} \left[ \left| \sum_i \beta_i \mathbf{x}_{0i} \right| > \frac{\lambda}{10} \right] &\leq \mathbb{P} [|\beta_{\tau}^* \mathbf{x}_{\tau}| + |\beta_{\tau^c}^* \mathbf{x}_{\tau^c}| > \lambda/10] \\
 &\leq \mathbb{P} [\mathbf{x}_{\tau} \neq \mathbf{0}] + \mathbb{P} [\mathbf{x}_{\tau} = \mathbf{0}] \cdot \mathbb{P} [|\beta_{\tau^c}^* \mathbf{x}_{\tau^c}| > \lambda/10] \\
 \text{(SM2.27)} \quad &\leq \theta |\tau| + 2\theta. \quad \blacksquare
 \end{aligned}$$

**Corollary SM2.7** (*|\langle \beta\_{\tau \setminus \{0\}}, \mathbf{x}\_{0, \tau \setminus \{0\}} \rangle| is small near shifts*). Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $|\tau|, c_{\mu}$  such that  $(\mathbf{a}_0, \theta, |\tau|)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_{\mu})$ , then if  $c_{\mu} \leq \frac{1}{10}$ , for any  $\mathbf{a}$  such that  $|\beta_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}$ , we have

$$\text{(SM2.28)} \quad \mathbb{P} \left[ \left| \sum_{i \in \tau \setminus \{0\}} \beta_i \mathbf{x}_{0i} \right| > \frac{2\lambda}{5} \right] \leq 2\theta$$

*Proof.* For the last tail bound, write  $\mathbf{x} = \boldsymbol{\omega} \circ \mathbf{g}$ . Wlog define  $\beta_0$  be the largest correlation  $\beta_{(0)}$ , define random variables  $s' = \langle \beta_{\tau \setminus \{0\}}, \mathbf{x}_{\tau \setminus \{0\}} \rangle$ . Firstly most of the entries of  $\mathbf{x}_{\tau}$  would be zero since via Bernstein inequality with  $\theta |\tau| < 0.1$ :

$$\begin{aligned}
 \mathbb{P} \left[ \sum_{i \in \tau} \omega_i > \log \theta^{-1} \right] &\leq \mathbb{P} \left[ \sum_{i \in \tau} \omega_i > \theta |\tau| + 0.9 \log \theta^{-1} \right] \\
 \text{(SM2.29)} \quad &\leq \exp \left( \frac{-0.9^2 \log^2 \theta^{-1}}{2(\theta |\tau| + 0.9 \log \theta^{-1}/3)} \right) \leq \theta
 \end{aligned}$$

thus with probability at least  $1 - \theta$ , we can write  $s'$  as a Gaussian r.v. with variation bounded as  $\mathbb{E} s'^2 \leq \mathbb{E} \left[ \sum_{i=1}^{\log \theta^{-1}} \beta_i \mathbf{g}_i \right]^2 = \log \theta^{-1} \beta_{(1)}^2$ , then via Gaussian tail bound [Lemma SM10.1](#):

$$\begin{aligned}
 \mathbb{P} [ |s'| > 0.4\lambda ] &\leq \mathbb{P} \left[ |g| > \frac{0.4\lambda}{\sqrt{\log \theta^{-1}} |\beta_{(1)}|} \right] + \mathbb{P} \left[ \sum_{i \in \tau} \omega_i > \log \theta^{-1} \right] \\
 \text{(SM2.30)} \quad &\leq \frac{2}{\sqrt{2\pi}} \exp(-1.2 \log \theta^{-1}) + \theta \leq 2\theta, \quad \blacksquare
 \end{aligned}$$

**SM3. Euclidean gradient as soft-thresholding in shift space.** In this section, we will study the Euclidean gradient (4.7), by deriving bounds showing that the  $\chi$  operator approximates a soft-thresholding function in shift space ([Lemma SM3.2](#) and [Corollary SM3.4](#)). Furthermore, we will show the operator  $\chi[\beta_i]$  is monotone in  $|\beta_i|$  from [Lemma SM3.3](#). A figure of visualized  $\chi$  operator is shown in [Figure SM1](#).

To understand the  $\chi$  operator, we shall first consider a simple case—when  $\mathbf{x}_0$  is highly sparse. By definition of  $\beta$  from (4.3) we can see that  $\beta$  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:

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

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 SM3.1 (Gaussian smoothed soft-thresholding).** *Let  $g \sim \mathcal{N}(0, 1)$ . Then for every  $b, s \in \mathbb{R}$  and  $\lambda > 0$ ,*

$$(SM3.1) \quad \mathbb{E}_g \left[ g \mathcal{S}_\lambda [b \cdot g + s] \right] = b(1 - \text{erf}_b(\lambda, s)),$$

where

$$(SM3.2) \quad \text{erf}_b(\lambda, s) = \frac{1}{2} \text{erf} \left( \frac{\lambda + s}{\sqrt{2}|b|} \right) + \frac{1}{2} \text{erf} \left( \frac{\lambda - s}{\sqrt{2}|b|} \right).$$

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

$$(SM3.3) \quad \sigma \mathcal{S}_{\nu'_2 \lambda} [b] \leq \sigma \mathbb{E}_g \left[ g \mathcal{S}_\lambda [b \cdot g] \right] \leq \sigma \mathcal{S}_{\nu'_1(\varepsilon) \lambda} [b] + \varepsilon$$

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

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

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

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

To show (SM3.3), via definition of error function, for  $x > 0$ , we know:

$$(SM3.4) \quad \min \left\{ 1 - \varepsilon, \frac{1 - \varepsilon}{\sqrt{\log(1/\varepsilon)}} x \right\} \leq \operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \leq \frac{2x}{\sqrt{\pi}}$$

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

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

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

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

Lastly plug (SM3.4) into (SM3.1) and apply condition  $|b| \leq 1$  and  $\varepsilon < 1/4$  we have

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

which completes the proof. ■

This lemma establishes when  $\mathbf{x}_0$  is separated, then  $\chi$  is soft thresholding operator on  $\beta$  with threshold about  $\lambda/2$ . This phenomenon extends beyond the separated case, as long as when  $\mathbf{x}_0$  is sufficiently sparse (when Definition SM2.1 holds). Recall that  $\chi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is defined as

$$(SM3.5) \quad \chi[\beta] = \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}_0} \beta \right].$$

The following lemma bounds its expectation:

**Lemma SM3.2 (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  $\chi$  as in (SM3.5), then*

$$(SM3.6) \quad n^{-1} \mathbb{E} \chi[\beta]_i = \theta \beta_i (1 - \mathbb{E}_{\mathbf{s}_i} \operatorname{erf}_{\beta_i}(\lambda, \mathbf{s}_i))$$

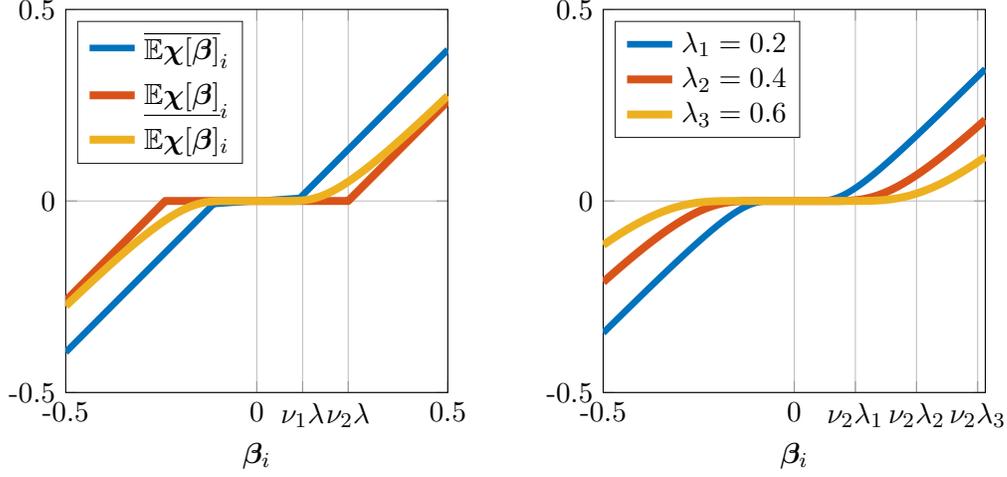
where  $\mathbf{s}_i = \sum_{\ell \neq i} \beta_\ell \mathbf{x}_{0\ell}$ . Suppose  $(\mathbf{a}_0, \theta, |\tau|)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$  and  $\lambda = c_\lambda / \sqrt{|\tau|}$  for some  $c_\lambda > 1/5$  and  $\sigma_i = \operatorname{sign}(\beta_i)$ , then there exists some numerical constant  $\bar{c}$  such that if  $c_\mu \leq \bar{c}$  then for every  $\mathbf{a} \in \mathfrak{A}(\mathcal{S}_\tau, \gamma(c_\mu))$  and every  $i \in [n]$ , (SM3.6) has upper bound

$$(SM3.7) \quad \sigma_i n^{-1} \mathbb{E} \chi[\beta]_i \leq \sigma_i n^{-1} \overline{\mathbb{E} \chi[\beta]_i} := \begin{cases} 4\theta^2 |\tau| |\beta_i| & |\beta_i| < \nu_1 \lambda \\ \theta (|\beta_i| - \nu_1 \lambda / 2) & |\beta_i| \geq \nu_1 \lambda \end{cases},$$

and lower bound

$$(SM3.8) \quad \sigma_i n^{-1} \mathbb{E} \chi[\beta]_i \geq \sigma_i n^{-1} \underline{\mathbb{E} \chi[\beta]_i} =: \theta \mathcal{S}_{\nu_2 \lambda} [|\beta_i|],$$

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



**Figure SM1.** A numerical example of  $\mathbb{E}\chi[\beta]_i$ . We provide figures for the expectation of  $\chi$  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 (SM3.1), and the blue/red lines are its upper/lower bound (SM3.3) utilized in the analysis respectively. Right: functions of  $\beta_i \rightarrow \beta_i (1 - \text{erf}_{\beta_i}(\lambda, 0))$  with different  $\lambda$ , 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  $\chi$  operator as follows:

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

$$\begin{aligned} \chi[\beta]_i &= e_i^* \check{\mathcal{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda [\check{\mathcal{C}}_{\mathbf{x}_0} \beta] = \langle s_{-i}[\mathbf{x}_0], \mathcal{S}_\lambda [\mathbf{x}_0 * \check{\beta}] \rangle \\ &\equiv_d \langle s_{-j}[\mathbf{x}_0], \mathcal{S}_\lambda [s_{i-j}[\mathbf{x}_0] * \check{\beta}] \rangle = \chi[s_{j-i}[\beta]]_j \end{aligned}$$

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

$$(SM3.9) \quad Z_j(\beta) = \mathbf{x}_j \mathcal{S}_\lambda \left[ \beta_0 \mathbf{x}_j + \sum_{\ell \in [\pm p] \setminus 0} \beta_\ell \mathbf{x}_{j+\ell} \right],$$

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 SM3.1, we can conclude

$$(SM3.10) \quad n^{-1} \mathbb{E}\chi[\beta]_0 = \mathbb{E}_{\mathbf{x}_0, \mathbf{s}_0} \mathbf{x}_0 \mathcal{S}_\lambda [\beta_0 \mathbf{x}_0 + \mathbf{s}_0] = \theta \beta_0 (1 - \mathbb{E}_{\mathbf{s}_0} \text{erf}_{\beta_0}(\lambda, \mathbf{s}_0))$$

so that (SM3.6) 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  $\mathbb{E}\mathbf{x}_0\mathcal{S}_\lambda[0 \cdot \mathbf{x}_0 + \mathbf{s}_0] = 0$ , we have

$$\begin{aligned}
 \mathbb{E}Z(\boldsymbol{\beta}) &\leq \beta_0 \sup_{\beta \in [0, \beta_0]} \frac{d}{d\beta} \mathbb{E}_{\mathbf{x}_0, \mathbf{s}_0} \mathbf{x}_0 \mathcal{S}_\lambda[\beta \mathbf{x}_0 + \mathbf{s}_0] \\
 &= \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \frac{d}{d\beta} \int_{|\beta g + \mathbf{s}_0| > \lambda} g(\beta g + \mathbf{s}_0 - \text{sign}(\beta g + \mathbf{s}_0) \cdot \lambda) d\mu(g) d\mu(\mathbf{s}_0) \\
 &= \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \mathbb{E}_{g, \mathbf{s}_0} [g^2 \mathbf{1}_{\{|\beta g + \mathbf{s}_0| > \lambda\}}] \\
 &\leq \theta \beta_0 \sup_{\beta \in [0, \beta_0]} \mathbb{E}_{g, \mathbf{s}_0} \left[ g^2 \left( \mathbf{1}_{\{|\beta g| > \frac{9\lambda}{10}\}} + \mathbf{1}_{\{|\mathbf{s}_0| > \frac{\lambda}{10}\}} \right) \right] \\
 \text{(SM3.11)} \quad &\leq \theta \beta_0 \left( (\mathbb{E}g^6)^{1/3} \mathbb{P}[|\beta_0 g| > (9\lambda/10)]^{2/3} + \mathbb{P}[|\mathbf{s}_0| > \lambda/10] \right)
 \end{aligned}$$

We bound the tail probability of  $\mathbf{s}_0$  using Corollary SM2.6 where

$$\text{(SM3.12)} \quad \mathbb{P}[|\mathbf{s}_0| > \lambda/10] \leq \mathbb{P}[\|\sum_i \beta_i \mathbf{x}_i\| > \lambda/10] \leq \theta |\boldsymbol{\tau}| + 2\theta \leq 3\theta |\boldsymbol{\tau}|.$$

On the other hand, the first term in (SM3.11) can be derived by pdf of Gaussian r.v. Lemma SM10.1 as:

$$\begin{aligned}
 (\mathbb{E}g^6)^{1/3} \mathbb{P}[|\beta_0 g| > (9\lambda/10)]^{2/3} &\leq \sqrt[3]{15} \left( \frac{10\beta_0}{9\lambda\sqrt{2\pi}} \right)^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right) \\
 \text{(SM3.13)} \quad &\leq \frac{3}{2} \left( \frac{\beta_0}{\lambda} \right)^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right).
 \end{aligned}$$

Combine (SM2.26), (SM3.13), when  $\beta_0 < \nu_1 \lambda$ , we know  $e^{-\frac{\lambda^2}{4\beta_0^2}} \leq e^{\log \theta} \leq \theta |\boldsymbol{\tau}|$ . The first type of upper bound  $\mathbb{E}Z$  is derived as

$$\text{(SM3.14)} \quad \forall \beta_0 \in [0, \nu_1 \lambda], \quad \mathbb{E}Z(\boldsymbol{\beta}) \leq \theta \beta_0 \left( \frac{3}{2} \nu_1^{2/3} \exp\left(-\frac{\lambda^2}{4\beta_0^2}\right) + 3\theta |\boldsymbol{\tau}| \right) \leq 4\theta^2 |\boldsymbol{\tau}| \beta_0.$$

(2). The second type of upper bound can be derived directly from Lemma SM3.1:

$$\begin{aligned}
 \mathbb{E}Z(\boldsymbol{\beta}) &\leq \mathbb{E}_{\mathbf{x}_0} \mathbb{E}_{\mathbf{s}_0} \mathbf{x}_0 \mathcal{S}_\lambda[\beta_0 \mathbf{x}_0 + \mathbf{s}_0] \leq \mathbb{E}_{\mathbf{x}_0} \mathbf{x}_0 \mathcal{S}_\lambda[\beta_0 \mathbf{x}_0] + \mathbb{E}_{\mathbf{x}_0} |\mathbf{x}_0| \mathbb{E}_{\mathbf{s}_0} |\mathbf{s}_0| \\
 \text{(SM3.15)} \quad &\leq \theta \cdot \left( \mathcal{S}_{\nu_1 \lambda}[\beta_0] + \varepsilon + \sqrt{2/\pi} \cdot \mathbb{E}|\mathbf{s}_0| \right),
 \end{aligned}$$

where  $\mathbb{E}|\mathbf{s}_0|$  can be bounded with  $\|\boldsymbol{\beta}\|_2$  and  $\theta |\boldsymbol{\tau}| < c_\mu \theta_{\log}$  from Lemma SM2.5. When  $c_\mu < \frac{1}{10}$ , observe that

$$\text{(SM3.16)} \quad \mathbb{E}|\mathbf{s}_0| \leq \sqrt{\sum_{\ell} \mathbb{E}\mathbf{x}_\ell^2 \beta_\ell^2} \leq \sqrt{\theta} (\|\boldsymbol{\beta}_\tau\|_2 + \|\boldsymbol{\beta}_{\tau^c}\|_2) \leq \sqrt{\theta} (1 + c_\mu) + \frac{c_\mu \theta_{\log}}{|\boldsymbol{\tau}|} \leq \frac{2c_\mu \theta_{\log}}{\sqrt{|\boldsymbol{\tau}|}}.$$

Now choose  $\varepsilon = \theta \leq \frac{c_\mu \theta_{\log}}{|\boldsymbol{\tau}|}$ , so that  $\nu'_1 = \nu_1 = \frac{\sqrt{\theta_{\log}}}{2}$  in (SM3.15). Since  $c_\mu < \frac{c_\lambda}{25}$  we gain

$$\begin{aligned} \mathbb{E}Z(\boldsymbol{\beta}) &\leq \theta \left( \mathcal{S}_{\nu_1 \lambda}[\boldsymbol{\beta}_0] + \frac{c_\mu \theta_{\log}}{|\boldsymbol{\tau}|} + \sqrt{\frac{2}{\pi}} \cdot \frac{2c_\mu \theta_{\log}}{\sqrt{|\boldsymbol{\tau}|}} \right) \leq \theta \left( \mathcal{S}_{\nu_1 \lambda}[\boldsymbol{\beta}_0] + \frac{3c_\mu \theta_{\log}}{\sqrt{|\boldsymbol{\tau}|}} \right) \\ \text{(SM3.17)} \quad &\leq \theta \left( \mathcal{S}_{\nu_1 \lambda}[\boldsymbol{\beta}_0] + \frac{\sqrt{\theta_{\log}} \lambda}{5} \right) \leq \theta \left( \mathcal{S}_{\nu_1 \lambda}[\boldsymbol{\beta}_0] + \frac{1}{2} \nu_1 \lambda \right) \end{aligned}$$

(3). Combine both (SM3.14) and (SM3.17), we can thus conclude that

$$\text{(SM3.18)} \quad \mathbb{E}Z(\boldsymbol{\beta}) := \overline{\mathbb{E}Z(\boldsymbol{\beta})} \leq \begin{cases} 4\theta^2 |\boldsymbol{\tau}| \boldsymbol{\beta}_0 & \boldsymbol{\beta}_0 \leq \nu_1 \lambda \\ \theta (\boldsymbol{\beta}_0 - \frac{\nu_1}{2} \lambda) & \boldsymbol{\beta}_0 > \nu_1 \lambda \end{cases}.$$

2. (Lower bound of  $\mathbb{E}Z$ ) On the other hand, for the lower bound for  $\mathbb{E}Z$ , use the fact that  $\text{erf}_{\boldsymbol{\beta}}(\lambda, \mathbf{s})$  is concave in  $\mathbf{s}_0$ , we have

$$\begin{aligned} \mathbb{E}Z(\boldsymbol{\beta}) &= \mathbb{E}_{\mathbf{s}_0} \mathbb{E}_{\mathbf{x}_0} \mathcal{S}_\lambda[\boldsymbol{\beta}_0 \mathbf{x}_0 + \mathbf{s}_0] \\ &= \theta \cdot \mathbb{E}_{\mathbf{s}_0} \left[ \boldsymbol{\beta}_0 - \frac{\boldsymbol{\beta}_0}{2} \cdot \text{erf} \left( \frac{\lambda - \mathbf{s}_0}{\sqrt{2} |\boldsymbol{\beta}_0|} \right) - \frac{\boldsymbol{\beta}_0}{2} \cdot \text{erf} \left( \frac{\lambda + \mathbf{s}_0}{\sqrt{2} |\boldsymbol{\beta}_0|} \right) \right] \\ \text{(SM3.19)} \quad &\geq \theta \left( \boldsymbol{\beta}_0 - \boldsymbol{\beta}_0 \cdot \text{erf} \left( \frac{\lambda}{\sqrt{2} |\boldsymbol{\beta}_0|} \right) \right) \geq \theta \cdot \mathcal{S}_{\nu_2 \lambda}[\boldsymbol{\beta}_0] =: \underline{\mathbb{E}Z(\boldsymbol{\beta})}. \end{aligned}$$

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

Another convenient fact of  $\mathbb{E}\boldsymbol{\chi}[\boldsymbol{\beta}]_i$  is that it is monotone increasing w.r.t.  $|\boldsymbol{\beta}_i|$ . The monotonicity is clear in Figure SM1; it is demonstrated rigorously with the following lemma:

**Lemma SM3.3 (Monotonicity of  $\mathbb{E}\boldsymbol{\chi}(\boldsymbol{\beta})$ ).** *Suppose  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $|\boldsymbol{\tau}|, c_\mu$  such that  $(\mathbf{a}_0, \theta, |\boldsymbol{\tau}|)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{|\boldsymbol{\tau}|}$  in  $\varphi_{\ell^1}$  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  $|\mathbb{E}[\boldsymbol{\chi}[\boldsymbol{\beta}]_i]|$  is monotone increasing in  $|\boldsymbol{\beta}_i|$ . In other words, if  $|\boldsymbol{\beta}_i| > |\boldsymbol{\beta}_j|$  then*

$$\text{(SM3.20)} \quad \sigma_i \mathbb{E}\boldsymbol{\chi}[\boldsymbol{\beta}]_i \geq \sigma_j \mathbb{E}\boldsymbol{\chi}[\boldsymbol{\beta}]_j$$

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

The proof first operate simple calculus and then followed by studying cases of  $|\boldsymbol{\beta}_i| - |\boldsymbol{\beta}_j|$  when either it is smaller or larger than  $\lambda$ .

*Proof.* 1. (Monotonicity by gradient negativity) Wlog assume  $\boldsymbol{\beta}_i > \boldsymbol{\beta}_j > 0$ , and from Lemma SM3.2 we can write  $(n\theta)^{-1} \mathbb{E}\boldsymbol{\chi}[\boldsymbol{\beta}]_i = \boldsymbol{\beta}_i (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\boldsymbol{\beta}_i}(\lambda, \mathbf{s}_i))$ . Consider  $t \in [0, 1]$  and define  $\ell(t) = t\boldsymbol{\beta}_i - t\boldsymbol{\beta}_j$ . Write the random variable  $\mathbf{s}_{ij} = \sum_{\ell \neq i, j} \boldsymbol{\beta}_\ell \mathbf{x}_\ell$ . Define  $h$  as a function of  $t$  such that

$$\begin{aligned} h(t) &= \mathbb{E}_{\mathbf{x}, \mathbf{s}_{ij}} \left[ \left( (1-t)\boldsymbol{\beta}_i + t\boldsymbol{\beta}_j \right) \left( 1 - \text{erf}_{(1-t)\boldsymbol{\beta}_i + t\boldsymbol{\beta}_j}(\lambda, ((1-t)\boldsymbol{\beta}_j + t\boldsymbol{\beta}_i)x + \mathbf{s}_{ij}) \right) \right] \\ \text{(SM3.21)} \quad &= \mathbb{E}_{\mathbf{x}, \mathbf{s}_{ij}} \left[ \left( \boldsymbol{\beta}_i - \ell(t) \right) \left( 1 - \text{erf}_{\boldsymbol{\beta}_i - \ell(t)}(\lambda, x \cdot (\boldsymbol{\beta}_j + \ell(t)) + \mathbf{s}_{ij}) \right) \right]. \end{aligned}$$

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:

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

(SM3.22)

$$= (\beta_j - \beta_i) \mathbb{E}_{x, \mathbf{s}_{ij}} \left[ 1 - \int_0^{z_{\lambda_+}} f(z) dz - \int_0^{z_{\lambda_-}} f(z) dz + (z_{\lambda_+} + x)f(z_{\lambda_+}) + (z_{\lambda_-} - x)f(z_{\lambda_-}) \right].$$

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

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

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 (SM3.22) using tower property and from assumption  $\beta_j - \beta_i < 0$  we obtain

$$\begin{aligned}
(\text{SM3.22}) &\leq (\beta_j - \beta_i) \left( 1 - \mathbb{P}[z_{\lambda_+} > 0] \cdot \mu_{++} \right. \\
&\quad \left. - \mathbb{P}[z_{\lambda_-} > 0] \cdot \mu_{--} + \mathbb{E}_{x, s_{ij}} [x(f(z_{\lambda_+}) - f(z_{\lambda_-}))] \right) \\
&\leq (\beta_j - \beta_i) \left( 1 - \min \left\{ \frac{\mathbb{P}[z_{\lambda_+} > 0]}{2}, \frac{\mathbb{E}|z_{\lambda_+}|}{\sqrt{2\pi}} \right\} \right. \\
(\text{SM3.23}) \quad &\quad \left. - \min \left\{ \frac{\mathbb{P}[z_{\lambda_-} > 0]}{2}, \frac{\mathbb{E}|z_{\lambda_-}|}{\sqrt{2\pi}} \right\} - \frac{\theta}{\sqrt{2\pi}} \cdot \mathbb{E}|g| \right),
\end{aligned}$$

where  $g$  is standard Gaussian r.v..

2. (Cases of varying  $\beta_i, \beta_j$ ) 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  $\beta_*^2 = \|\beta_\tau\|_2^2 > \frac{1-3c_\mu}{|\tau|}$ . On the other hand if at least one of  $i, j$  is in  $\tau$  and all other  $\beta_\tau$  entries are small where  $\|\beta_{\tau \setminus \{i, j\}}\|_\infty^2 < \frac{1-3c_\mu}{|\tau|}$ , then we know via norm inequalities,

$$(\text{SM3.24}) \quad (\beta_i + \beta_j)^2 > \beta_i^2 + \beta_j^2 > \|\beta_\tau\|_2^2 - (|\tau| - 1) \|\beta_{\tau \setminus \{i, j\}}\|_\infty^2 \geq \frac{1 - 3c_\mu}{|\tau|},$$

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

$$(\text{SM3.25}) \quad \beta_* = \beta_j + \ell(t) = (\beta_i + \beta_j) - (\beta_i - \ell(t)) \geq \frac{\sqrt{1-3c_\mu}}{\sqrt{|\tau|}} - \frac{1}{4\sqrt{|\tau|}} \geq \frac{0.72}{\sqrt{|\tau|}}.$$

In this case, adopt result from Corollary SM2.6 such that  $\mathbb{P}[\sum \beta_\ell x_\ell > \lambda/10] \leq 3\theta |\tau| \leq .01$ , we have

$$\begin{aligned}
(\text{SM3.26}) \quad \mathbb{P}[z_{\lambda_-} > 0] &= \mathbb{P}[z_{\lambda_+} > 0] = 1 - \mathbb{P}[x(\beta_j + \ell(t)) + s_{ij} < -\lambda] \\
&\leq 1 - \mathbb{P}[\mathbf{x}_* \beta_* < -11\lambda/10] \cdot \mathbb{P}[x(\beta_j + \ell(t)) + s_{ij} - \mathbf{x}_* \beta_* < \lambda/10] \\
&\leq 1 - \theta \cdot \mathbb{P}\left[\mathbf{g}_* \cdot \frac{0.72}{\sqrt{|\tau|}} < \frac{-11c_\lambda}{10\sqrt{|\tau|}}\right] \cdot \left(1 - \mathbb{P}\left[\sum \beta_\ell x_\ell > \frac{\lambda}{10}\right]\right) \\
&\leq 1 - \theta \cdot \mathbb{P}[0.72 \cdot \mathbf{g}_* \leq -1.1 \cdot 0.25] \cdot (1 - 3c_\mu) \\
&\leq 1 - 0.35\theta.
\end{aligned}$$

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:

$$\begin{aligned}
(\text{SM3.27}) \quad \mathbb{E}_{x, s_{ij}} |z_{\lambda_-}| &= \mathbb{E}_{x, s_{ij}} |z_{\lambda_+}| \leq \mathbb{E}_{x, s_{ij}} \frac{\lambda + |x(\beta_j + \ell(t)) - s_{ij}|}{\beta_i - \ell(t)} \\
&\leq 1 + 4\sqrt{|\tau|} \cdot \left(\mathbb{E}_{x, s_{ij}} |x(\beta_j + \ell(t)) - s_{ij}|^2\right)^{1/2} \\
&\leq 1 + 4\sqrt{|\tau|\theta} \|\beta\|_2 \leq 1 + 4\sqrt{|\tau|\theta} \left(1 + c_\mu + \frac{c_\mu}{\sqrt{\theta|\tau|}}\right) \leq 1.2.
\end{aligned}$$

Combine (SM3.23), (SM3.26) we have

$$(SM3.28) \quad h'(t) \leq (\beta_j - \beta_i) \left( 1 - 2 \cdot \frac{(1 - 0.35\theta)}{2} - \frac{\theta}{\sqrt{2\pi}} \cdot \sqrt{\frac{2}{\pi}} \right) \leq 0.03\theta(\beta_j - \beta_i) < 0,$$

and combine (SM3.23), (SM3.27) and  $\theta < c_\mu$  we have

$$(SM3.29) \quad h'(t) \leq (\beta_j - \beta_i) \left( 1 - 2 \cdot \frac{1.2}{\sqrt{2\pi}} - \frac{\theta}{\sqrt{2\pi}} \cdot \sqrt{\frac{2}{\pi}} \right) \leq 0.03(\beta_j - \beta_i) < 0,$$

which proves the monotonicity. ■

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

**Corollary SM3.4 (Finite sample deviation of  $\chi(\beta)$ ).** *Suppose  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda/\sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  and every  $i \in [n]$ , we have:*

$$(SM3.30) \quad |n^{-1}\chi[\beta]_i - n^{-1}\mathbb{E}\chi[\beta]_i| \leq c\theta/p^{3/2},$$

*Proof.* See Subsection SM9.1 ■

**SM4. Euclidean Hessian as logic function in shift space.** We can express the (pseudo) curvature (4.11) in direction  $\mathbf{v} \in \mathbb{S}^{p-1}$  in terms of the correlation  $\gamma = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\nu} \mathbf{v}$  between  $\mathbf{v}$  and  $\mathbf{a}_0$ , giving

$$\mathbf{v}^* \tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} = -\gamma^* \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{P}_I \check{\mathbf{C}}_{\mathbf{x}_0} \gamma,$$

where

$$(SM4.1) \quad I(\mathbf{a}) = \text{supp} \left( \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\nu} \mathbf{a} \right] \right) = \left\{ i \in [n] \mid \left| \mathbf{x}_0 * \check{\boldsymbol{\beta}} \right|_i > \lambda \right\}.$$

The  $i$ -th diagonal entry of  $\check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}_0}$  is

$$(SM4.2) \quad -\mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{e}_i = - \left\| \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{e}_i \right\|_2^2 = - \left\| \mathbf{P}_{I(\mathbf{a})} \mathbf{s}_{-i}[\mathbf{x}_0] \right\|_2^2,$$

which is the core component for us to study the curvature of objective  $\varphi_{\ell^1}$ . We illustrate the expectation of diagonal term of Hessian in Lemma SM4.2 and Corollary SM4.3, whose figure of visualized  $\left\| \mathbf{P}_{I(\mathbf{a})} \mathbf{s}_{-i}[\mathbf{x}_0] \right\|_2$  is shown in Figure SM1. Lastly, we also prove the off-diagonal terms  $\mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{e}_j$  of Hessian is likely inconsequential in calculation of curvature in Lemma SM4.4.

We expect the Hessian to have stronger negative component in the  $s_i[\mathbf{a}_0]$  direction as  $\left\| \mathbf{P}_{I(\mathbf{a})} \mathbf{s}_{-i}[\mathbf{x}_0] \right\|_2^2$  becomes larger. This term can be tremendously simplified when  $\mathbf{x}_0$  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 (SM4.1), we can derive

$$(SM4.3) \quad -\|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 = -\sum_{j \in I_0} \mathbf{x}_{0j}^2 \mathbf{1}_{\{|\sum_{\ell} \beta_{\ell} \mathbf{x}_{0(\ell+j-i)}| > \lambda\}} \stackrel{\text{sep.}}{=} -\sum_{j \in I_0} \mathbf{g}_j^2 \mathbf{1}_{\{|\beta_i \mathbf{g}_j| > \lambda\}},$$

where  $\mathbf{1}$  is the indicator function and  $\mathbf{g}_j$  are independent standard Gaussian r.v.s.. In expectation, the summands in (SM4.3) acts like a smoothed logic function on entry  $\beta_i$ :

**Lemma SM4.1 (Gaussian smoothed indicator).** *Let  $g \sim \mathcal{N}(0, 1)$ , then for any  $b, s \in \mathbb{R}$  and  $\lambda > 0$ .*

$$(SM4.4) \quad \mathbb{E}_g [g^2 \mathbf{1}_{\{|b \cdot g + s| > \lambda\}}] = 1 - \text{erf}_b(\lambda, s) + f_b(\lambda, s),$$

where

$$(SM4.5) \quad f_b(\lambda, s) = \frac{1}{\sqrt{2\pi}} \left[ \left( \frac{\lambda + s}{|b|} \right) e^{-\frac{(\lambda+s)^2}{2b^2}} + \left( \frac{\lambda - s}{|b|} \right) e^{-\frac{(\lambda-s)^2}{2b^2}} \right].$$

*Proof.* The proof can be derived via same calculation of integrals in Lemma SM3.1.  $\blacksquare$

Although the definition (SM4.4) 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 be an increasing function of  $|b|$ . Moreover by assigning different values for  $|b|$  we obtain:

$$(SM4.6) \quad \mathbb{E} g^2 \mathbf{1}_{\{|b \cdot g| > \lambda\}} \approx \begin{cases} 1, & |b| \approx 1 \\ 1/2, & |b| \approx \sqrt{2/\pi}\lambda \\ 0, & |b| \approx 0 \end{cases}$$

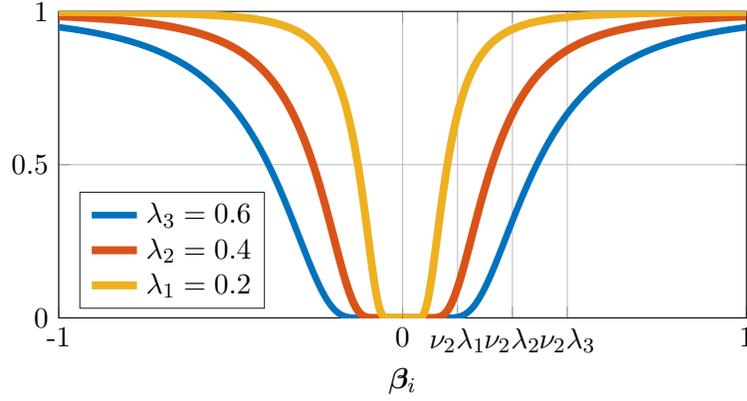
Relate (SM4.6) to (SM4.3), when  $|\beta_i|$  is close to 1 then we expect  $-\frac{1}{n\theta} \|\mathbf{P}_{I} s_{-i}[\mathbf{x}_0]\|_2^2$  to be close to  $-1$ , and it increases to 0 as  $|\beta_i|$  decreases, suggests that the Euclidean Hessian at point  $\mathbf{a}$  has stronger negative component at  $s_i[\mathbf{a}_0]$  direction if  $|\langle \mathbf{a}, s_i[\mathbf{a}_0] \rangle|$  is larger. See Figure SM2 for a numerical example. This phenomenon can be extend beyond the idealistic separating case as follows:

**Lemma SM4.2 (Expected Hessian diagonals).** *Let  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  and  $\lambda > 0$ , define the set  $I(\mathbf{a})$  in (SM4.1), write  $\mathbf{s}_i = \sum_{\ell \neq i} \beta_{\ell} \mathbf{x}_{0\ell}$ , then for every  $\mathbf{a} \in \mathbb{S}^{p-1}$  and  $i \in [n]$ :*

$$(SM4.7) \quad n^{-1} \mathbb{E} \|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 = \theta [1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i) + \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i)]$$

*Proof.* Write  $\mathbf{x}_0$  as  $\mathbf{x}$ . Observe that  $\mathbf{y} * \check{\mathbf{a}} = \mathbf{x}_0 * \check{\beta} = \sum_{\ell} \beta_{\ell} s_{-\ell}[\mathbf{x}_0]$ . Thus for any  $j \in [n]$  and  $i \in [\pm p]$ :

$$(SM4.8) \quad (\mathbf{y} * \check{\mathbf{a}})_{j-i} = \left( \beta_i s_{-i}[\mathbf{x}] + \sum_{\ell \neq i} \beta_{\ell} s_{-\ell}[\mathbf{x}] \right)_{j-i} = \beta_i \mathbf{x}_j + \sum_{\ell \neq i} \beta_{\ell} \mathbf{x}_{j+\ell-i} =: \beta_i \mathbf{x}_j + \mathbf{s}_j,$$



**Figure SM2.** A numerical example for  $\mathbb{E} \|\mathbf{P}_{I(\mathbf{a})} s_i[\mathbf{x}_0]\|_2^2$ . We provide a figure to illustrate the expectation of  $-\frac{1}{n\theta} \|\mathbf{P}_{I(\mathbf{a})} s_i[\mathbf{x}_0]\|_2^2$  when entries of  $\mathbf{x}_0$  are  $2p$ -separated, as a function plot of  $\beta_i \rightarrow 1 - \text{erf}_{\beta_i}(\lambda, 0) + f_{\beta_i}(\lambda, 0)$  from (SM4.4) with different  $\lambda$ . 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 than  $\lambda$  its value grow to 1, implies there is a negative curvature at  $s_i[\mathbf{a}_0]$  direction. Similarly if  $|\beta_i|$  is much smaller than  $\lambda$  the function value is 0 thus the curvature is positive in  $s_i[\mathbf{a}_0]$  direction.

where  $\mathbf{x}_j$  is independent of  $\mathbf{s}_j$ , and both  $\mathbf{x}_j, \mathbf{s}_j$  are symmetric and identically distributed for all  $j \in [n]$ . Rewrite the random variable using (SM4.1) as

$$\begin{aligned} \|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 &= \left\| \mathbf{P}_{I(\mathbf{a})} \sum_{j \in [n]} (\mathbf{x}_{0j} \mathbf{e}_{j-i}) \right\|_2^2 = \sum_{j \in [n]} \mathbf{x}_{0j}^2 \mathbf{1}_{\{|\mathbf{y}^* \tilde{\mathbf{a}}|_{j-i} > \lambda\}} \\ &= \sum_{j \in [n]} \mathbf{x}_{0j}^2 \mathbf{1}_{\{|\beta_i \mathbf{x}_{0j} + \mathbf{s}_j| > \lambda\}} \end{aligned}$$

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

$$\mathbb{E} \|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 = n\theta \cdot \mathbb{E} \mathbf{g}_0^2 \mathbf{1}_{\{|\beta_i \mathbf{g}_0 + \mathbf{s}_0| > \lambda\}} = n\theta \cdot \mathbb{E} (1 - \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i) + f_{\beta_i}(\lambda, \mathbf{s}_i))$$

where  $\mathbf{s}_i = \sum_{\ell \neq i} \mathbf{x}_{0\ell} \beta_i$  with  $\mathbf{x}_{0i} \sim_{\text{i.i.d.}} \text{BG}(\theta)$ , yielding the claimed expression.  $\blacksquare$

When the signal length of  $\mathbf{y}$  is sufficiently large, then  $i$ -th diagonal term for Hessian  $\|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2$  will be close enough to its expected value.

**Corollary SM4.3 (Large sample deviation of curvature).** Suppose  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  and every  $i \in [n]$ , we have:

$$(SM4.9) \quad \left| n^{-1} \|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 - n^{-1} \mathbb{E} \|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2^2 \right| \leq c\theta/p$$

*Proof.* See Subsection SM9.2.  $\blacksquare$

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 SM4.4 (Hessian off-diagonal term near solution).** *Suppose  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ , where  $|\beta_{(1)}| \leq \frac{1}{4 \log \theta^{-1}} \lambda$  and every  $i \neq j \in [\pm p] \setminus \{(0)\}$ , we have*

$$(SM4.10) \quad |s_i[\mathbf{x}_0]^* | \mathbf{P}_{I(\mathbf{a})} | s_j[\mathbf{x}_0]| < 8n\theta^3$$

*Proof.* Write  $\theta_{\log} = -1/\log \theta$  and  $\mathbf{x}_0$  as  $\mathbf{x} = \boldsymbol{\omega} \circ \mathbf{g}$ . Wlog let  $\beta_0$  be the largest correlation  $\beta_{(0)}$ . Define random variables  $s' = \langle \boldsymbol{\beta}_{\tau \setminus \{0, i, j\}}, \mathbf{x}_{\tau \setminus \{0, i, j\}} \rangle$ . Firstly via [Corollary SM2.7](#) we have  $\mathbb{P}[|s'| > 0.4\lambda] \leq 2\theta$ ; also define  $s = \langle \boldsymbol{\beta}_{\tau^c \setminus \{0, i, j\}}, \mathbf{x}_{\tau^c \setminus \{0, i, j\}} \rangle$ , and base on [Corollary SM2.6](#) we have  $\mathbb{P}[|s| > \lambda/10] \leq 2\theta$ . Expand the  $(-i, -j)$ -th cross term with  $\theta < 0.1$  we have:

$$(SM4.11) \quad \begin{aligned} \mathbb{E} |s_{-i}[\mathbf{x}]^* | \mathbf{P}_{I(\mathbf{a})} | s_{-j}[\mathbf{x}]| &= \mathbb{E} \sum_{k \in [n]} |\mathbf{x}_{k+i} \mathbf{x}_{k+j}| \mathbf{1}_{\{|\beta_0 \mathbf{x}_k + \beta_i \mathbf{x}_{k+i} + \beta_j \mathbf{x}_{k+j} + s + s'| > \lambda\}} \\ &= n\theta^2 \cdot \mathbb{E} |\mathbf{g}_i \mathbf{g}_j| \mathbf{1}_{\{|\beta_0 \mathbf{x}_0 + \beta_i \mathbf{g}_i + \beta_j \mathbf{g}_j + s + s'| > \lambda\}} \\ &\leq n\theta^2 \cdot \mathbb{E} [|\mathbf{g}_i \mathbf{g}_j| (2\mathbf{1}_{\{|\beta_i \mathbf{g}_i| > \lambda/4\}} \\ &\quad + \mathbb{P}[\mathbf{x}_0 \neq 0] + \mathbb{P}[|s| > 0.1\lambda] + \mathbb{P}[|s'| > 0.4\lambda])] \\ &\leq n\theta^2 \cdot (\exp(-\log^2 \theta^{-1}) + \theta + 2\theta + 2\theta) \\ &\leq 6n\theta^3. \end{aligned}$$

Write (SM4.10) as two summation of independent random variables with  $t = j - i$  by separating sum into two sets  $J_{t1}, J_{t2}$  defined in (SM1.4) with both  $|J_{t1}|, |J_{t2}| < n\theta^2$  with probability at least  $1 - 2/n$  from [Lemma SM1.1](#)

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

whose first summands can be upper bounded with high probability via Bernstein inequality [Lemma SM10.4](#) with  $(\sigma^2, R) = (1, 1)$  and writes  $\mathcal{C} := \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu)) \cap \left\{ \mathbf{a} \mid |\beta_{(1)}| \leq \frac{1}{4 \log \theta^{-1}} \lambda \right\}$ ,

then we have

$$\begin{aligned}
 & \mathbb{P} \left[ \max_{\substack{i \neq j \in [\pm p] \setminus \{0\} \\ \mathbf{a} \in \mathcal{C}}} \left( \sum_{(k-i) \in I(\mathbf{a}) \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| - \mathbb{E} \sum_{(k-i) \in I(\mathbf{a}) \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| \right) \geq n\theta^3 \right] \\
 & \mathbb{P} \left[ \max_{i \neq j \in [\pm p] \setminus \{0\}} \left( \sum_{(k-i) \in \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| - \mathbb{E} \sum_{(k-i) \in \cap J_{t1}} |\mathbf{g}_k| |\mathbf{g}_{k+t}| \right) \geq n\theta^3 \right] \\
 \text{(SM4.12)} \quad & \leq 4p^2 \cdot \exp \left( \frac{-n^2\theta^6}{2|J_{t1}| + 2n\theta^3} \right) \leq \exp \left( 8 \log p - \frac{-n^2\theta^6}{3n\theta^2} \right) \leq \exp \left( -\frac{n\theta^4}{10} \right) \leq \frac{1}{n}
 \end{aligned}$$

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 (SM4.11) and (SM4.12) 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 SM1.1 and (SM4.12). ■

**SM5. Geometric relation between  $\rho$  and  $\ell^1$ -norm.** In this section, we discuss how to ensure that the smooth sparsity surrogate  $\rho$  approximates  $\|\cdot\|_1$  accurately enough that guarantees  $\varphi_\rho$  inherits the good properties of  $\varphi_{\ell^1}$ . We prove several lemmas which allow us to transfer properties of  $\varphi_{\ell^1}$  to  $\varphi_\rho$ . 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 SM5.2 that is smooth and well approximates  $\ell^1$  when the proper smoothing parameter  $\delta$  is chosen from the result of Lemma SM5.6. In particular we ask the regularizer  $\rho_\delta(x)$  to be uniformly bounded to  $|x|$  by  $\delta/2$ :

$$\text{(SM5.1)} \quad \forall x \in \mathbb{R}, \quad |\rho_\delta(x) - |x|| \leq \delta/2$$

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

$$\text{(SM5.2)} \quad \|\text{prox}_{\lambda\ell^1}[\check{\mathbf{a}} * \mathbf{y}] - \text{prox}_{\lambda\rho_\delta}[\check{\mathbf{a}} * \mathbf{y}]\|_2 \rightarrow 0,$$

$$\text{(SM5.3)} \quad \|\nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_{\rho_\delta}(\mathbf{a})\|_2 \rightarrow 0,$$

$$\text{(SM5.4)} \quad \|\tilde{\nabla}^2\varphi_{\ell^1}(\mathbf{a}) - \nabla^2\varphi_{\rho_\delta}(\mathbf{a})\|_2 \rightarrow 0.$$

An example choices of eligible smooth sparse surrogate is demonstrated in Table SM1.

The marginal minimizer over  $\mathbf{x}$  in (2.7) can be expressed in terms of the proximal operator [SM2] of  $\rho$  at point  $\check{\mathbf{a}} * \mathbf{y}$ :

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

Plugging in, we obtain

$$(SM5.5) \quad \varphi_\rho(\mathbf{a}) = \lambda\rho(\text{prox}_{\lambda\rho}[\check{\mathbf{a}} * \mathbf{y}]) + \frac{1}{2} \|\check{\mathbf{a}} * \mathbf{y} - \text{prox}_{\lambda\rho}[\check{\mathbf{a}} * \mathbf{y}]\|_2^2 - \frac{1}{2} \|\check{\mathbf{a}} * \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2$$

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

$$(SM5.6) \quad \varphi_\rho(\mathbf{a}) = \epsilon(\lambda\rho)(\check{\mathbf{a}} * \mathbf{y}) - \frac{1}{2} \|\check{\mathbf{a}} * \mathbf{y}\|_2^2 + \frac{1}{2} \|\mathbf{y}\|_2^2,$$

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 SM5.1 (Derivative of Moreau envelope, [SM2], 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,  $\varphi_\rho$  is twice differentiable whenever  $\text{prox}_{\lambda\rho}$  is differentiable. In this case, the (Euclidean) gradient and hessian of  $\varphi_\rho$  are given by

$$(SM5.7) \quad \nabla\varphi_\rho(\mathbf{a}) = -\boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \text{prox}_{\lambda\rho} [\check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}],$$

$$(SM5.8) \quad \nabla^2\varphi_\rho(\mathbf{a}) = -\boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \nabla \text{prox}_{\lambda\rho} [\check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}] \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota}.$$

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

$$(SM5.9) \quad \text{grad}[\varphi_\rho](\mathbf{a}) = -\mathbf{P}_{\mathbf{a}^\perp} \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \text{prox}_{\lambda\rho} [\check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}],$$

$$(SM5.10) \quad \text{Hess}[\varphi_\rho](\mathbf{a}) = -\mathbf{P}_{\mathbf{a}^\perp} \left( \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \nabla \text{prox}_{\lambda\rho} [\check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a}] \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} - \langle \nabla\varphi_\rho(\mathbf{a}), \mathbf{a} \rangle \mathbf{I} \right) \mathbf{P}_{\mathbf{a}^\perp}.$$

Our analysis accommodates any sufficiently accurate smooth approximation  $\rho$  to the  $\ell^1$  function. The requisite sense of approximation is captured in the following definition:

**Definition SM5.2 ( $\delta$ -smoothed  $\ell^1$  function).** We call an additively separable function  $\rho(\mathbf{x}) = \sum_{i=1}^n \rho_i(x_i) : \mathbb{R}^n \rightarrow \mathbb{R}$ , a  $\delta$ -smoothed  $\ell^1$  function with  $\delta > 0$  if for each  $i \in [n]$ ,  $\rho_i$  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  $x \in \mathbb{R}$ :

$$(SM5.11) \quad |\rho_i(x) - |x|| + c \leq \delta/2$$

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

**Definition SM5.3 ( $\sqrt{\delta}$ -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

$$(SM5.12) \quad 0 \leq \text{sign}(z) \left( \mathcal{S}_\lambda^\delta[z] - \mathcal{S}_\lambda[z] \right) \leq \sqrt{\lambda\delta}, \quad \forall z \in \mathbb{R}$$

| Surrogate class       | $\rho_i(x)$   | $\nabla\rho_i(x)$                             | $\nabla^2\rho_i(x)$                                  |
|-----------------------|---|---|--|
| Log hyperbolic cosine | $\frac{\delta}{2} \log \left( e^{2x/\delta} + e^{-2x/\delta} \right)$ | $\frac{e^{4x/\delta} - 1}{e^{4x/\delta} + 1}$ | $\frac{4e^{4x/\delta}}{\delta(e^{4x/\delta} + 1)^2}$ |
| 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 SM1**

**Classes of smooth sparse surrogate  $\rho$  and how to set its parameter.** *Three common classes are listed with parameter  $\delta$  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 (SM5.11). Also its second order derivatives  $\nabla^2\rho_i(x)$  are monotone decreasing w.r.t.  $|x|$ , hence are certified to be eligible  $\delta$ -smoothed  $\ell^1$  surrogates.*

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

$$(SM5.13) \quad \left| \nabla \mathcal{S}_\lambda^\delta[z] - \nabla \mathcal{S}_\lambda[z] \right| \leq \sqrt{\lambda\delta}/B, \quad ||z| - \lambda| \geq B.$$

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

**Lemma SM5.4 (Proximal operator for smoothed  $\ell^1$ ).** *Suppose  $\rho$  is a  $\delta$ -smoothed  $\ell^1$  function, then  $z_i \mapsto \text{prox}_{\lambda\rho}[z]_i$  is a  $\sqrt{\delta}$ -smoothed soft threshold function.*

*Proof.* We know that

$$(SM5.14) \quad \mathbf{x}_z := \text{prox}_{\lambda\rho}[z] = \underset{\mathbf{x} \in \mathbb{R}^n}{\text{argmin}} \lambda\rho(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - z\|_2^2.$$

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

$$(SM5.15) \quad \mathbf{x}_{zi} = (\text{Id} + \lambda\nabla\rho_i)^{-1}(z_i).$$

Since  $\rho_i$  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

$$(SM5.16) \quad \nabla\mathbf{x}_{zi} = \nabla_i(\text{Id} + \lambda\nabla\rho_i)^{-1}(z_i) = \frac{1}{\lambda\nabla^2\rho_i(\mathbf{x}_{zi}) + 1} < 1.$$

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

Now we are left to show that (SM5.12) and (SM5.13) 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

$$(SM5.17) \quad |\nabla\rho_i(x)| \leq 1 \quad \forall x \in \mathbb{R}, \quad \nabla\rho_i(x) \geq 1 - \sqrt{\delta/\lambda} \quad \forall x \geq \sqrt{\lambda\delta}.$$

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 (SM5.12), when  $0 \leq \mathbf{z}_i \leq \lambda$ , then  $\mathcal{S}_\lambda[\mathbf{z}_i] = 0$  and  $\mathbf{x}_{z_i} \leq \sqrt{\lambda\delta}$  since if  $\mathbf{x}_{z_i} > \sqrt{\lambda\delta}$ , by (SM5.17):

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

then  $\mathbf{x}_{z_i}$  violate the stationary condition in (SM5.15), resulting  $0 \leq \mathbf{x}_{z_i} - \mathcal{S}_\lambda[\mathbf{z}_i] \leq \sqrt{\lambda\delta}$  whenever  $0 \leq \mathbf{z}_i \leq \lambda$ . Likewise in the case of  $\mathbf{z}_i \geq \lambda$  where  $\mathcal{S}_\lambda[\mathbf{z}_i] = \mathbf{z}_i - \lambda$ , (SM5.17) provides:

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

again violates (SM5.15) and therefore (SM5.12) holds for all  $\mathbf{z}_i \in \mathbb{R}$ .

Lastly (SM5.13) is a direct result of (SM5.12). For all  $\mathbf{z}_i \leq \lambda - B$ , recall that  $\nabla\mathbf{x}_{z_i}$  is monotone increasing in  $\mathbf{z}_i$ :

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

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

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

implies (SM5.13) holds. ■

Based on (SM5.9)-(SM5.10) and denote  $\check{\mathbf{C}}_{\mathbf{y}}\boldsymbol{\iota}\mathbf{a} = \check{\mathbf{a}} * \mathbf{y}$ , the only differences of Riemannian gradient and Hessian between  $\varphi_\rho$  and  $\varphi_{\ell^1}$  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  $\varphi_\rho$  with that of objective  $\varphi_{\ell^1}$ , we may apply both Definition SM5.3 and Lemma SM5.4, together suggest if  $\rho$  is a  $\delta$ -smoothed  $\ell^1$  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 SM5.5 (Gradients discontinuity entries).** For each  $\mathbf{a} \in \mathbb{S}^{p-1}$ , let

$$(SM5.18) \quad J_B(\mathbf{a}) := \left\{ i \mid \left( \widetilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota}_{\mathbf{a}} \right)_i \in [-\lambda - B, -\lambda + B] \cup [\lambda - B, \lambda + B] \right\}.$$

Suppose the subspace dimension is at most  $k$  and signal  $\mathbf{y}$  satisfies [Definition SM2.1](#). 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

$$(SM5.19) \quad |J_B(\mathbf{a})| \leq \frac{24c'n\theta^2}{p \log n}$$

*Proof.* See [Subsection SM9.3](#). ■

The geometric approximation between  $\varphi_{\ell^1}$  and  $\varphi_\rho$  necessarily consists of three parts: the gradient, the Hessian, and the coefficients. Here we conclude the approximation result with the following lemma:

**Lemma SM5.6 ( $\varphi_{\ell^1}$  approximates  $\varphi_\rho$ ).** Suppose  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  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  $\delta$ -smoothed  $\ell^1$  function with

$$(SM5.20) \quad \lambda = \frac{c_\lambda}{\sqrt{k}}, \quad \delta \leq \frac{c'^4 \theta^8}{p^2 \log^2 n} \lambda$$

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

$$(SM5.21) \quad \left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \text{prox}_{\lambda \ell^1}[\check{\mathbf{a}} * \mathbf{y}] - \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}] \right\|_2 \leq c'n\theta^4.$$

(2). The gradient has norm difference

$$(SM5.22) \quad \|\nabla \varphi_{\ell^1}(\mathbf{a}) - \nabla \varphi_\rho(\mathbf{a})\|_2 \leq c'n\theta^4.$$

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

$$(SM5.23) \quad \forall \mathbf{v} \in \mathbb{S}^{p-1}, \quad \left| \mathbf{v}^* \left( \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) - \text{Hess}[\varphi_\rho](\mathbf{a}) \right) \mathbf{v} \right| \leq 200c'n\theta^2.$$

*Proof.* 1. (Coefficients) From [Lemma SM5.4](#), the proximal  $\delta$ -smoothed  $\ell^1$  function satisfies

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

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

$$(SM5.24) \quad \left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda[\check{\mathbf{a}} * \mathbf{y}] - \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda^\delta[\check{\mathbf{a}} * \mathbf{y}] \right\|_2 \leq \sqrt{\lambda \delta n} \cdot \left\| \boldsymbol{\iota}_{[\pm p]}^* \widetilde{\mathbf{C}}_{\mathbf{x}_0} \right\|_2.$$

Apply [Lemma SM1.5](#) 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

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

and we gain

$$\text{(SM5.24)} \leq \sqrt{\lambda\delta n} \cdot \sqrt{3n\theta} \leq n\sqrt{3\lambda\theta\delta} \leq c'n\theta^4.$$

2. ([Gradient](#)) From definition of Riemannian gradient ([SM5.9](#)) and apply similar norm bound of ([SM5.24](#)), and condition on the following events of [Lemma SM1.5](#) holds, obtain

$$\text{(SM5.25)} \quad \left\| \nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a}) \right\|_2 \leq \sqrt{\lambda\delta n} \cdot \left\| \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \right\|_2 \leq n\sqrt{3\lambda\theta(1+\mu p)\delta} \leq c'n\theta^4.$$

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 SM5.5](#), condition on the event such that

$$\text{(SM5.26)} \quad B \leq \frac{c'\lambda\theta^2}{p \log n}, \quad |J| \leq \frac{24c'n\theta^2}{p \log n};$$

and rewrite  $J_B(\mathbf{a})$  as  $J$ . Also condition on the event using [Lemma SM1.5](#) and  $(1+\mu p)\theta \log \theta^{-1} < 1$

$$\text{(SM5.27)} \quad \left\| \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \right\|_2 \leq \sqrt{3n}, \quad \left\| \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \mathbf{P}_J \right\|_2 \leq \sqrt{8|J|p \log n},$$

then the difference of Hessian ([SM5.10](#)), in direction  $\mathbf{v} \in \mathbb{S}^{p-1}$  can be bounded as

$$\begin{aligned} & \left| \mathbf{v}^* \left( \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) - \text{Hess}[\varphi_\rho](\mathbf{a}) \right) \mathbf{v} \right| \\ \text{(SM5.28)} \quad & \leq \left| \mathbf{v}^* \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \left( \mathbf{P}_{I(\mathbf{a})} - \text{diag} \left[ \nabla \mathcal{S}_\lambda^\delta \left[ \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a} \right] \right] \right) \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{v} \right| + \left\| \nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a}) \right\|_2 \end{aligned}$$

where  $I(\mathbf{a})$  is defined in ([SM4.1](#)). Let  $\mathbf{D} = \mathbf{P}_{I(\mathbf{a})} - \text{diag} \left[ \nabla \mathcal{S}_\lambda^\delta \left[ \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a} \right] \right]$  and notice that  $\mathbf{D}$  is a diagonal matrix, which suggests ([SM5.28](#)) 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  $\ell^1$  function [Lemma SM5.4](#):

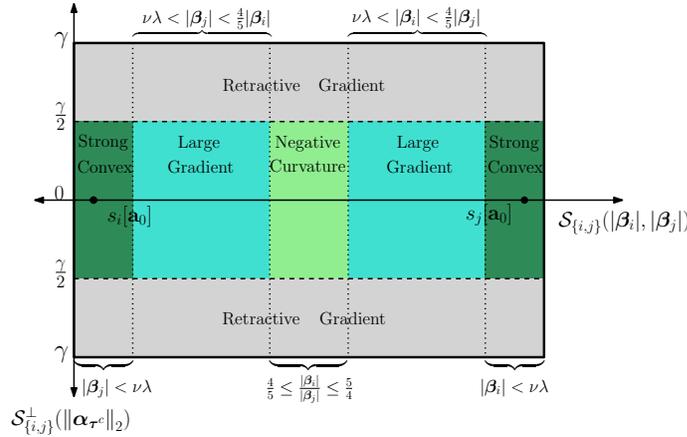
$$\max_j |\mathbf{P}_J \mathbf{D} \mathbf{P}_J|_{jj} \leq 1, \quad \max_j |\mathbf{P}_{J^c} \mathbf{D} \mathbf{P}_{J^c}|_{jj} \leq \sqrt{\lambda\delta}/B.$$

Finally, once again apply  $\delta$  bound from ([SM5.20](#)) and bounds for  $B, |J|, \mathbf{y}$  from ([SM5.26](#))-([SM5.27](#)), we gain

$$\begin{aligned} \text{(SM5.28)} & \leq \left\| \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \mathbf{P}_J \right\|_2^2 + \frac{\sqrt{\lambda\delta}}{B} \left\| \boldsymbol{\iota}^* \check{\mathbf{C}}_{\mathbf{y}} \right\|_2^2 + \left\| \nabla\varphi_{\ell^1}(\mathbf{a}) - \nabla\varphi_\rho(\mathbf{a}) \right\|_2^2 \\ & \leq 8|J|p \log n + \frac{3n\sqrt{\lambda\delta}}{B} + c'n\theta^2 \\ & \leq 8 \cdot \frac{24c'n\theta^2}{p \log n} \cdot p \log n + \frac{3n(c'^4 \lambda^2 \theta^8 / p^2 \log^2 n)^{1/2}}{c' \lambda \theta^2 / p \log p} + c'n\theta^2 \\ & \leq 200c'n\theta^2, \end{aligned}$$

where all above result holds with probability at least  $1 - 10/n$  from [Lemma SM5.5](#) and [Lemma SM1.5](#). ■

**SM6. Analysis of geometry.** In this section we prove major geometrical result in [Theorem 4.1](#). This lemma consists of three parts of geometry of  $\varphi_\rho$ ; including the negative curvature region [Corollary SM6.2](#), large gradient region [Corollary SM6.4](#), strong convexity region near shift [Corollary SM6.6](#), and retraction to subspace [Corollary SM6.8](#), which are respectively base on geometric properties of  $\varphi_{\ell^1}$  in [Lemma SM6.1](#), [Lemma SM6.3](#), [Lemma SM6.5](#) and [Lemma SM6.7](#). 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:

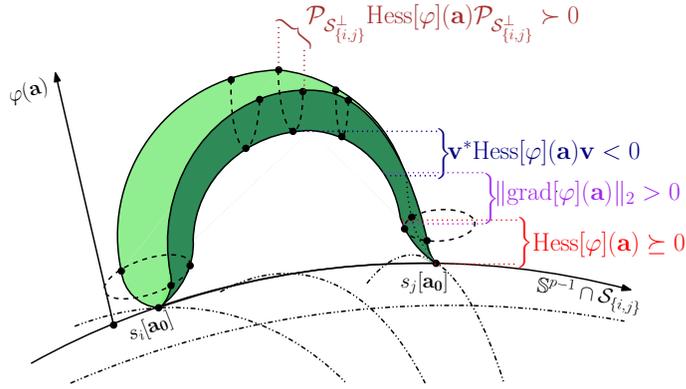


**Figure SM3.** The top view of geometry over subspace  $\mathcal{S}_{\{i,j\}}$ . We display the geometric properties in the neighborhood of subspace  $\mathcal{S}_{\{i,j\}}$  (horizontal axis) which contains the solutions  $s_i[\mathbf{a}_0]$  and  $s_j[\mathbf{a}_0]$ . When  $\mathbf{a}$  lies near 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).

**SM6.1. Negative curvature .** For any  $\mathbf{a} \in \mathbb{S}^{p-1}$  near the subspace  $\mathcal{S}_\tau$  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  $\varphi_{\ell^1}$  in [Lemma SM6.1](#), then show  $\varphi_\rho$  approximates  $\varphi_{\ell^1}$  in terms of Hessian in [Corollary SM6.2](#) when  $\rho$  is properly defined as in [Section SM5](#).

**Lemma SM6.1 (Negative curvature for  $\varphi_{\ell^1}$ ).** Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Set  $\lambda = c_\lambda/\sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{v} \in \mathcal{S}_{\{(0),(1)\}} \cap \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$ ,

$$(SM6.1) \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} \leq -cn\theta\lambda.$$



**Figure SM4.** 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.

*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  $\beta$ 's

$$(SM6.2) \quad 0.79 \geq \frac{|\beta_{(0)}|}{\sqrt{\beta_{(0)}^2 + \beta_{(1)}^2}} \|\beta_\tau\|_2 \geq |\beta_{(0)}| \geq |\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}| \geq \frac{4}{5} \cdot \frac{\|\beta_\tau\|_2}{\sqrt{|\tau|}} \geq \frac{0.79}{\sqrt{|\tau|}}$$

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

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

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

$$(SM6.3) \quad \gamma_{(0)}^2 \geq \frac{(1 - \mu) \cdot 16}{25 + 16} \geq 0.385$$

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

$$(SM6.4) \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} = \mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} - \langle \nabla \varphi_{\ell^1}(\mathbf{a}), \mathbf{a} \rangle = -\gamma^* \iota_J^* \widetilde{\mathbf{M}} \widetilde{\mathbf{C}}_x \mathbf{P}_{I(\mathbf{a})} \widetilde{\mathbf{C}}_x \mathbf{M} \iota_J \boldsymbol{\gamma} + \beta^* \boldsymbol{\chi}[\boldsymbol{\beta}]$$

expand the first term of (SM6.4) we obtain

$$\begin{aligned}
 & -\gamma^* \boldsymbol{\iota}_J^* \mathbf{M} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{M} \boldsymbol{\iota}_J \gamma \\
 &= -\gamma^* \boldsymbol{\iota}_J^* \mathbf{M} (\mathbf{P}_{(0)} + \mathbf{P}_{(1)} + \mathbf{P}_{J^c}) \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} (\mathbf{P}_{(0)} + \mathbf{P}_{(1)} + \mathbf{P}_{J^c}) \mathbf{M} \boldsymbol{\iota}_J \gamma \\
 &\leq -\sum_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_i \right\|_2^2 (e_i^* \mathbf{M} \boldsymbol{\iota}_J \gamma)^2 + 2 \sum_{\substack{(i,j) \in \{J, J^c\} \\ (i,j) = (0),(1)}} \left| e_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| |(e_i^* \mathbf{M} \boldsymbol{\iota}_J \gamma) (e_j^* \mathbf{M} \boldsymbol{\iota}_J \gamma)| \\
 &\leq -\sum_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_i \right\|_2^2 (|\gamma_i| - \mu)^2 \\
 \text{(SM6.5)} \quad & + 2 \max_{i \neq j \in [\pm p]} \left| e_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| (\|\boldsymbol{\iota}_J^* \mathbf{M} \boldsymbol{\iota}_J \gamma\|_1 \|\boldsymbol{\iota}_{J^c}^* \mathbf{M} \boldsymbol{\iota}_J \gamma\|_1 + (|\gamma_{(0)}| + \mu) (|\gamma_{(1)}| + \mu))
 \end{aligned}$$

Consider the following events

$$\begin{aligned}
 \text{(SM6.6)} \quad & \left\{ \mathcal{E}_{\text{cross}} := \left\{ \forall \mathbf{a} \in \mathbb{S}^{p-1}, \max_{i \neq j \in [\pm p]} \left| e_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| < 4n\theta^2 \right\} \right. \\
 & \left. \left\{ \mathcal{E}_{\text{ncurv}} := \left\{ \forall \mathbf{a} \in \mathfrak{R}(\mathcal{S}_{\boldsymbol{\tau}}, \gamma(c_{\mu})), \min_{i \in J} \left\| \mathbf{P}_{I(\mathbf{a})} \mathbf{s}_{-i}[\mathbf{x}] \right\|_2^2 \geq n\theta (1 - \mathbb{E}_{\mathbf{s}_i}(\lambda, \mathbf{s}_i) + \mathbb{E}_{\mathbf{s}_i}(\lambda, \mathbf{s}_i)) - \frac{c_{\mu}n\theta}{p} \right\} \right. \right\} ,
 \end{aligned}$$

and from Lemma SM2.4 we know

$$\|\boldsymbol{\iota}_J^* \mathbf{M} \boldsymbol{\iota}_J \gamma\|_1 \leq \|\gamma\|_1 + 2\mu \leq 1.5, \quad \|\boldsymbol{\iota}_{J^c}^* \mathbf{M} \boldsymbol{\iota}_J \gamma\|_1 \leq \mu p \|\gamma\|_1 \leq 1.5\mu p,$$

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

$$\begin{aligned}
 & -\gamma^* \boldsymbol{\iota}_J^* \mathbf{M} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{M} \boldsymbol{\iota}_J \gamma \\
 \text{(SM6.7)} \quad & \leq \underbrace{-n\theta \cdot \sum_{i \in J} (|\gamma_i| - \mu)^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i) + \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i))}_{g_1(\boldsymbol{\beta})} + (18\mu p + 8) n\theta^2 + \frac{2c_{\mu}n\theta}{\sqrt{|\boldsymbol{\tau}|}}
 \end{aligned}$$

Meanwhile, for the latter term of (SM6.4), consider the following event  $\mathcal{E}_{\bar{\chi}}$  where we write  $\sigma_i = \text{sign}(\boldsymbol{\beta}_i)$  as:

$$\text{(SM6.8)} \quad \mathcal{E}_{\bar{\chi}} := \left\{ \sigma_i \boldsymbol{\chi}[\boldsymbol{\beta}]_i \leq \begin{cases} n\theta \cdot |\boldsymbol{\beta}_i| (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) + \frac{c_{\mu}n\theta}{p}, & \forall i \in \boldsymbol{\tau} \\ n\theta \cdot |\boldsymbol{\beta}_i| 4\theta |\boldsymbol{\tau}| + \frac{c_{\mu}n\theta}{p}, & \forall i \in \boldsymbol{\tau}^c \end{cases} \right\},$$

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

$$\begin{aligned}
 & \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \leq n\theta \cdot \sum_{i \in \boldsymbol{\tau}} \boldsymbol{\beta}_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) + 4n\theta^2 |\boldsymbol{\tau}| \|\boldsymbol{\beta}_{\boldsymbol{\tau}^c}\|_2^2 + \frac{c_{\mu}n\theta}{p} \|\boldsymbol{\beta}\|_1 \\
 \text{(SM6.9)} \quad & \leq \underbrace{n\theta \cdot \sum_{i \in \boldsymbol{\tau}} \boldsymbol{\beta}_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i))}_{g_2(\boldsymbol{\beta})} + \frac{5c_{\mu}n\theta}{\sqrt{|\boldsymbol{\tau}|}}.
 \end{aligned}$$

2. (Lower bound  $\mathbb{E}f_{\beta_i}$ ) Combine the first term from each of the (SM6.7) and (SM6.9). Use  $\mu \leq c_\mu \leq \frac{1}{300}$  and (SM6.3) to obtain  $(|\gamma_{(0)}| - \mu)^2 > 0.38$ , we have

$$(SM6.10) \quad \begin{aligned} \frac{1}{n\theta} (g_1(\boldsymbol{\beta}) + g_2(\boldsymbol{\beta})) &\leq - \sum_{i \in J} \left[ (|\gamma_i| - \mu)^2 - \beta_i^2 \right] (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) \\ &+ \sum_{i \in \tau \setminus J} \beta_i^2 (1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)) - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i), \end{aligned}$$

now use Taylor expansion <sup>1</sup> for  $f_{\beta_i}$ , and apply the upper bound where  $\mathbb{E} \mathbf{s}_i^2 \leq \theta \|\boldsymbol{\beta}\|_2^2 \leq \theta \left( 1 + \frac{c_\mu}{\sqrt{|\tau|}} + \frac{c_\mu}{\theta|\tau|^2} \right) \leq \frac{3c_\mu}{|\tau|}$ ,

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

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 (SM6.2), and also apply assumption  $c_\lambda \leq \frac{1}{3}$  and  $c_\mu \leq \frac{1}{300}$ , we gain

$$(SM6.11) \quad \begin{aligned} 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) &\geq 0.3 \min_{\beta = \frac{0.79}{\sqrt{|\tau|}}, 0.79} f(\beta) \\ &\geq 0.3 \min \left\{ \frac{2c_\lambda}{0.79} - \frac{c_\lambda^3 + 9c_\mu c_\lambda}{0.79^3}, \lambda \left( \frac{2}{0.79} - \frac{c_\lambda^2 + 9c_\mu}{0.79^3} \right) \right\} \\ &\geq 0.3 \min \{2c_\lambda, 2\lambda\} \geq 0.6\lambda. \end{aligned}$$

3. (Upper bound  $\mathbb{E}\chi[\beta]_i$ ) When  $\beta_{(0)}^2 = (|\gamma_{(0)}| - \mu)^2 - \eta$  for some  $\eta > 0$ . With monotonicity Lemma SM3.3, which implies:

$$(SM6.12) \quad \begin{aligned} \left( 1 - \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)}) \right) &\geq \left( 1 - \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) \right) \\ &\geq \left( 1 - \mathbb{E}_{\mathbf{s}_i} \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i) \right), \end{aligned}$$

---

<sup>1</sup> Apply  $\exp[-x^2/2] > 1 - x^2/2$

then combine (22)-(SM6.12) and use  $\mu \leq \frac{c_\mu}{4\sqrt{|\tau|}}$  from Lemma SM2.5

$$\begin{aligned}
 \text{(SM6.10)} &\leq - \underbrace{\left( (|\gamma_{(0)}|^2 - \mu)^2 - \beta_{(0)}^2 - \eta \right)}_{=0} \left( 1 - \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)}) \right) \\
 &\quad + \left( \sum_{i \in \tau \setminus (0)} \beta_i^2 - (|\gamma_{(1)}| - \mu)^2 - \eta \right) \underbrace{\left( 1 - \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) \right)}_{<1} \\
 &\quad - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \\
 &\leq \left( \|\beta_\tau\|_2^2 - \|\gamma\|_2^2 + 2\mu \|\gamma\|_1 \right) - 0.6\lambda \\
 \text{(SM6.13)} &\leq \frac{2c_\mu}{\sqrt{|\tau|}} - 0.6\lambda.
 \end{aligned}$$

On the other hand, when  $\beta_{(0)}^2 \geq (|\gamma_{(0)}| - \mu)^2 > 0.38$ , combining (22)-(SM6.12) gives:

$$\begin{aligned}
 \text{(SM6.10)} &\leq \left( \|\beta_\tau\|_2^2 - \|\gamma\|_2^2 + 2\mu \|\gamma\|_1 \right) + \left( (|\gamma_{(0)}| - \mu)^2 - \beta_{(0)}^2 \right) \mathbb{E}_{\mathbf{s}_{(0)}} \text{erf}_{\beta_{(0)}}(\lambda, \mathbf{s}_{(0)}) \\
 &\quad + \left( (|\gamma_{(1)}| - \mu)^2 - \sum_{i \in \tau \setminus (0)} \beta_i^2 \right) \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) - 0.38 \sum_{i \in J} \mathbb{E}_{\mathbf{s}_i} f_{\beta_i}(\lambda, \mathbf{s}_i) \\
 \text{(SM6.14)} &\leq \left( \frac{c_\mu}{\sqrt{|\tau|}} + 4\mu \right) + \left( \gamma_{(1)}^2 - \|\beta_\tau\|_2^2 + \beta_{(0)}^2 \right) \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) - 0.6\lambda,
 \end{aligned}$$

where Lemma SM3.2 provides the upper bound for  $\mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)})$  as

$$\begin{aligned}
 \mathbb{E}_{\mathbf{s}_{(1)}} \text{erf}_{\beta_{(1)}}(\lambda, \mathbf{s}_{(1)}) &= 1 - \frac{1}{n\theta\beta_{(1)}} \mathbb{E}\chi[\beta]_{(1)} \leq 1 - \frac{\sigma(1)}{n\theta|\beta_{(1)}|} \mathbb{E}\chi[\beta]_{(1)} \\
 \text{(SM6.15)} &= 1 - \frac{1}{|\beta_{(1)}|} \left( |\beta_{(1)}| - \sqrt{\frac{2}{\pi}}\lambda \right) \leq \sqrt{\frac{2}{\pi}} \cdot \frac{\lambda}{|\beta_{(1)}|},
 \end{aligned}$$

then calculate the constant for the second term in (SM6.14) 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

$$\begin{aligned}
 \frac{(\gamma_{(1)}^2 - 1) + c_\mu + \beta_{(0)}^2}{|\beta_{(1)}|} &\leq -\frac{\kappa}{(\kappa^2 + 1)|\beta_{(0)}|} + \kappa |\beta_{(0)}| + \frac{\mu + c_\mu}{0.3} \\
 \text{(SM6.16)} &\leq \frac{\kappa^2 - 1}{\sqrt{\kappa^2 + 1}} + \kappa \left( \|\beta_\tau\|_2^2 - 1 \right) + 4.2c_\mu \leq 0.36 + 6c_\mu,
 \end{aligned}$$

and finally combine (SM6.15)-(SM6.16), follow from (SM6.14) and use  $c_\lambda \leq \frac{1}{3}$ :

$$\begin{aligned}
(\text{SM6.10}) &\leq \frac{2c_\mu}{\sqrt{|\boldsymbol{\tau}|}} + \sqrt{\frac{2}{\pi}} \left( \gamma_{(1)}^2 - 1 + c_\mu + \beta_{(0)}^2 \right) \frac{\lambda}{|\beta_{(1)}|} - 0.6\lambda \\
&\leq \frac{2c_\mu}{\sqrt{|\boldsymbol{\tau}|}} + \sqrt{\frac{2}{\pi}} \left( 0.36\lambda + \frac{6c_\mu c_\lambda}{0.3} \right) - 0.6\lambda \\
(\text{SM6.17}) &\leq \frac{4c_\mu}{\sqrt{|\boldsymbol{\tau}|}} - 0.3\lambda
\end{aligned}$$

3. (Collect all results) Combine the components of pseudo Hessian (SM6.7), (SM6.9) with bounds for  $g_1 + g_2$  from (SM6.13) and (SM6.17), and use Lemma SM2.5 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:

$$\begin{aligned}
(\text{SM6.18}) \quad \mathbf{v}^* \widetilde{\text{Hess}}_{\varphi_{\ell^1}}[\mathbf{a}]\mathbf{v} &\leq g_1(\boldsymbol{\beta}) + g_2(\boldsymbol{\beta}) + \frac{7c_\mu n \theta}{\sqrt{|\boldsymbol{\tau}|}} + (18\mu p + 8) n \theta^2 \\
&\leq n \theta \cdot \left( \frac{4c_\mu}{\sqrt{|\boldsymbol{\tau}|}} - 0.3\lambda \right) + n \theta \cdot \frac{7c_\mu}{\sqrt{|\boldsymbol{\tau}|}} + n \theta \cdot \frac{6.5c_\mu}{|\boldsymbol{\tau}|} \\
&\leq \frac{n \theta}{\sqrt{|\boldsymbol{\tau}|}} (0.059 - 0.06) \leq -0.001 n \theta \lambda
\end{aligned}$$

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

$$(\text{SM6.19}) \quad 1 - \underbrace{\mathbb{P}[\mathcal{E}_{\text{cross}}^c]}_{\text{Lemma SM1.4}} - \underbrace{\mathbb{P}[\mathcal{E}_{\text{ncurv}}^c]}_{\text{Corollary SM4.3}} - \underbrace{\mathbb{P}[\mathcal{E}_{\frac{c}{\lambda}}^c]}_{\text{Corollary SM3.4}}. \quad \blacksquare$$

Similarly for objective  $\varphi_\rho$ , we have that

**Corollary SM6.2 (Negative curvature for  $\varphi_\rho$ ).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{k}$  in  $\varphi_\rho$  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  $\rho$  is  $\delta$ -smoothed  $\ell^1$  function where  $\delta \leq c'' \lambda \theta^8 / p^2 \log^2 n$ ,  $n > C p^5 \theta^{-2} \log p$  and  $c_\mu \leq \bar{c}$ , then with probability at least  $1 - c'/n$ , for every  $\mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_\boldsymbol{\tau}, \gamma(c_\mu))$  satisfying  $|\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}|$ : for  $\mathbf{v} \in \mathcal{S}_{\{(0), (1)\}} \cap \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$ ,*

$$(\text{SM6.20}) \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_\rho](\mathbf{a})\mathbf{v} \leq -c n \theta \lambda$$

*Proof.* Choose  $\mathbf{v} \in \mathbb{S}^{p-1}$  according to Lemma SM6.1 and (SM5.23) from Lemma SM5.6 with constant multiplier  $\delta$  satisfies  $c''^{1/4} < 10^{-3} c$ , we gain

$$(\text{SM6.21}) \quad \mathbf{v}^* \text{Hess}[\varphi_\rho](\mathbf{a})\mathbf{v} \leq -c n \theta \lambda + 200 c' n \theta^2 \leq -c n \theta \lambda / 2 \quad \blacksquare$$

**SM6.2. Large gradient.** For any  $\mathbf{a} \in \mathbb{S}^{p-1}$  near subspace and the second largest correlation  $\beta_{(1)}$  much smaller than the first correlation  $\beta_{(0)}$  while not being near 0, the negative gradient of  $\varphi_\rho(\mathbf{a})$  will point at the largest shift. We show this in [Lemma SM6.3](#), and the  $\varphi_\rho$  version in [Corollary SM6.4](#) when  $\rho$  is properly defined as in [Section SM5](#).

**Lemma SM6.3 (Large gradient for  $\varphi_{\ell^1}$ ).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  satisfying  $\frac{4}{5} |\beta_{(0)}| > |\beta_{(1)}| > \frac{1}{4 \log \theta^{-1}} \lambda$ ,*

$$(SM6.22) \quad \langle \boldsymbol{\sigma}_{(0)} \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_{\ell^1}](\mathbf{a}) \rangle \geq cn\theta (\log^{-2} \theta^{-1}) \lambda^2$$

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

*Proof.* 1. (Properties for  $\boldsymbol{\alpha}, \boldsymbol{\beta}$ ) 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  $\boldsymbol{\iota}^* s_i[\mathbf{a}_0]$  as a function of  $\beta_i$  as  $d_{\mathbb{S}}(\mathbf{a}, \pm \boldsymbol{\iota}^* s_i[\mathbf{a}_0]) = \cos^{-1}(\beta_i)$ , then by triangle inequality we have:

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

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 SM2.1](#), implies

$$(SM6.23) \quad |\beta_{(0)}| \leq 1 - \frac{\beta_{(1)}^2}{2} \left( 1 - \frac{4\mu\sqrt{|\tau|}}{\theta_{\log} c_\lambda} \right) \leq 1 - 0.49 \beta_{(1)}^2 =: \beta_{\text{ub}}.$$

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

$$(SM6.24) \quad |\beta_{(0)}| \geq \frac{5}{4} |\beta_{(1)}| =: \beta_{\text{lb}}.$$

Also since  $\boldsymbol{\beta} = \mathbf{M}\boldsymbol{\alpha}$ , based on properties of  $\mathbf{M}$  from [Lemma SM2.4](#). When  $\|\boldsymbol{\alpha}_\tau\|_2 \leq 1 + c_\mu$  and  $\|\boldsymbol{\alpha}_{\tau^c}\|_2 \leq \gamma \leq \frac{c_\mu \theta_{\log}^2}{4\mu\sqrt{p}|\tau|}$ , we gain:

$$\begin{aligned} \boldsymbol{\beta}_{(0)} &= \boldsymbol{\alpha}_{(0)} + \mathbf{e}_{(0)}^* \mathbf{M} \boldsymbol{\alpha}_{\setminus(0)} \\ \implies |\boldsymbol{\alpha}_{(0)} - \boldsymbol{\beta}_{(0)}| &\leq \mu \sqrt{|\tau|} \|\boldsymbol{\alpha}_\tau\|_2 + \mu \sqrt{p} \|\boldsymbol{\alpha}_{\tau^c}\|_2 \\ (SM6.25) \quad &\leq \frac{c_\mu \theta_{\log}^2 (1 + c_\mu)}{4|\tau|} + \mu \sqrt{p} \gamma \leq \frac{c_\mu \theta_{\log}^2}{|\tau|}. \end{aligned}$$

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

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 SM10.7](#) tells us since  $\|\beta_{\tau \setminus (0)}\|_\infty \leq \beta_{(1)}$ :

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

then condition on the following event using [Corollary SM3.4](#),

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

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  $\theta$  from [Definition SM2.1](#),  $\|\beta_{\tau^c}\|_2 \leq \frac{c_\mu \theta_{\log}}{\sqrt{\theta} |\tau|}$  from [Lemma SM2.5](#),  $|\tau| \leq \sqrt{p}$  from lemma assumption and let  $c_\mu < \frac{1}{100}$ :

$$\begin{aligned} \beta^* \chi[\beta] &\leq \chi[\beta]_{(0)} \beta_{(0)} + \sum_{i \in \tau \setminus (0)} \beta_i \chi[\beta]_i + \langle \beta_{\tau^c}, \chi[\beta]_{\tau^c} \rangle \\ &\leq \chi[\beta]_{(0)} \beta_{(0)} + n \left( \theta \sum_{i \in \tau \setminus (0)} h(\beta_i) + 4\theta^2 |\tau| \|\beta_{\tau^c}\|_2^2 \right. \\ &\quad \left. + \frac{c_\mu \theta}{p^{3/2}} \left( \sqrt{|\tau|} \|\beta_\tau\|_2 + \sqrt{p} \|\beta_{\tau^c}\|_2 \right) \right) \\ &\leq \chi[\beta]_{(0)} \beta_{(0)} + n \left( \theta \cdot \eta (1 - \beta_{(0)}^2) + \theta \cdot \frac{c_\mu \theta_{\log}^2}{|\tau|} + \frac{4\theta^2 |\tau| c_\mu^2 \theta_{\log}^2}{\theta |\tau|^2} \right. \\ &\quad \left. + c_\mu \theta \left( \frac{1 + c_\mu}{p^{3/4} |\tau|} + \frac{c_\mu \theta_{\log}}{p \sqrt{\theta} |\tau|} \right) \right) \\ \text{(SM6.26)} \quad &\leq \chi[\beta]_{(0)} \beta_{(0)} + n\theta \left( \eta (1 - \beta_{(0)}^2) + \frac{6c_\mu \theta_{\log}^2}{|\tau|} \right), \end{aligned}$$

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

3. (Align the gradient with  $\iota^* s_{(0)}[\mathbf{a}_0]$ ) Base on the definition  $\beta$ , since  $\beta_{(0)} = \langle \mathbf{a}, \iota^* s_{(0)}[\mathbf{a}_0] \rangle$ , we can expect that the negative gradient is likely aligned with direction toward one of the candidate solution  $\pm \iota^* s_{(0)}[\mathbf{a}_0]$ . Wlog assume that both  $\beta_{(0)}, \beta_{(1)}$  are positive, then expand the

gradient and use incoherent property for  $\mathbf{a}_0$  [Lemma SM2.4](#) we have:

$$\begin{aligned}
 \langle \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}_{\varphi_{\ell_1}}[\mathbf{a}] \rangle &= \langle \iota^* s_{(0)}[\mathbf{a}_0], \iota^* \mathbf{C}_{\mathbf{a}_0} (\boldsymbol{\chi}[\boldsymbol{\beta}] - \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \boldsymbol{\alpha}) \rangle \\
 \text{(SM6.27)} \qquad \qquad \qquad &\geq (\boldsymbol{\chi}[\boldsymbol{\beta}]_{(0)} - \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \boldsymbol{\alpha}_{(0)}) - \mu \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\setminus(0)} - \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \boldsymbol{\alpha}_{\setminus(0)}\|_1,
 \end{aligned}$$

where  $\setminus(0)$  is an abbreviation of the complement set  $[\pm 2p_0] \setminus (0)$ . The latter part of [\(SM6.27\)](#) has an upper bound using bounds of  $\boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] < \frac{3n\theta}{2}$ ,  $\|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau^c}\|_2 < \frac{n\theta\gamma_2}{20}$  from [\(SM6.62\)](#), and  $\|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau \setminus (0)}\|_2 \leq n\theta \|\boldsymbol{\beta}_{\tau \setminus (0)}\|_2$  in event  $\mathcal{E}_{\bar{\mathcal{X}}}$ , we obtain:

$$\begin{aligned}
 &\mu \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\setminus(0)} - \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \boldsymbol{\alpha}_{\setminus(0)}\|_1 \\
 &\leq \mu \left( \sqrt{|\boldsymbol{\tau}|} \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau \setminus (0)}\|_2 + \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \sqrt{|\boldsymbol{\tau}|} \|\boldsymbol{\alpha}_{\tau \setminus (0)}\|_2 \right. \\
 &\quad \left. + \sqrt{p} \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau^c}\|_2 + \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \sqrt{p} \|\boldsymbol{\alpha}_{\tau^c}\|_2 \text{ Big} \right) \\
 &\leq n\theta \cdot \left[ \mu \sqrt{|\boldsymbol{\tau}|} (\|\boldsymbol{\beta}_{\tau}\|_2 - |\boldsymbol{\beta}_{(0)}|) + \mu \sqrt{|\boldsymbol{\tau}|} (\|\boldsymbol{\alpha}_{\tau}\|_2 - |\boldsymbol{\alpha}_{(0)}|) \right. \\
 &\quad \left. + \frac{1}{20} \mu \sqrt{p} \gamma_2 + \frac{3}{2} \mu \sqrt{p} \gamma_2 \right] \\
 &\leq n\theta \cdot \frac{c_\mu \theta_{\log}^2}{4|\boldsymbol{\tau}|} \left[ 2(1 + c_\mu) - |\boldsymbol{\beta}_{(0)}| - |\boldsymbol{\alpha}_{(0)}| + \left( \frac{1}{20} + \frac{3}{2} \right) c_\mu \right] \\
 \text{(SM6.28)} \qquad \qquad \qquad &\leq n\theta \cdot \frac{c_\mu \theta_{\log}^2}{|\boldsymbol{\tau}|} (0.5 + c_\mu - 0.5\boldsymbol{\beta}_{(0)}).
 \end{aligned}$$

On the other hand, the former term of [\(SM6.27\)](#) possesses a lower bound using [\(SM6.25\)](#)-

(SM6.26),  $\chi[\beta]_{(0)} > n\theta \left( \beta_{(0)} - \frac{\nu_1}{2}\lambda - \frac{c_\mu}{p} \right) \geq n\theta \left( \beta_{(0)} - 0.51\nu_1\lambda \right)$  and  $\alpha_{(0)} \leq 1$ :

$$\begin{aligned}
& \chi[\beta]_{(0)} - \beta^* \chi[\beta] \alpha_{(0)} \\
& \geq (1 - \alpha_{(0)}\beta_{(0)}) \chi[\beta]_{(0)} - n\theta \cdot \left[ \eta \left( 1 - \beta_{(0)}^2 \right) + \frac{6c_\mu\theta_{\log}^2}{|\tau|} \right] \alpha_{(0)} \\
& \geq n\theta \underbrace{\left( 1 - \left( \beta_{(0)} + \frac{c_\mu\theta_{\log}^2}{|\tau|} \right) \beta_{(0)} \right)}_{(a)} \left( \beta_{(0)} - 0.51\nu_1\lambda \right) \\
& \quad - n\theta \underbrace{\left[ \eta \left( 1 - \beta_{(0)}^2 \right) \left( \beta_{(0)} + \frac{c_\mu\theta_{\log}^2}{|\tau|} \right) + \frac{6c_\mu\theta_{\log}^2}{|\tau|} \alpha_{(0)} \right]}_{(b)} \\
& \geq n\theta \left[ \underbrace{\left( 1 - \beta_{(0)}^2 \right) \left( \beta_{(0)} - 0.51\nu_1\lambda \right) - \frac{c_\mu\theta_{\log}^2\beta_{(0)}^2}{|\tau|}}_{(a)} \right. \\
& \quad \left. - \underbrace{\left( 1 - \beta_{(0)}^2 \right) \eta\beta_{(0)} - \eta \frac{c_\mu\theta_{\log} \left( 1 - \beta_{(0)}^2 \right)}{|\tau|} - \frac{6c_\mu\theta_{\log}^2}{|\tau|}}_{(b)} \right] \\
\text{(SM6.29)} \quad & \geq n\theta \left[ \left( 1 - \beta_{(0)}^2 \right) \left( (1 - \eta)\beta_{(0)} - 0.51\nu_1\lambda \right) - \frac{c_\mu\theta_{\log}^2}{|\tau|} \left( (1 - \eta)\beta_{(0)}^2 + 7 \right) \right],
\end{aligned}$$

combine (SM6.27) with (SM6.28)-(SM6.29) and  $\eta > 0$ , we have

$$\begin{aligned}
\text{(SM6.27)} \quad & \geq n\theta \left[ \left( 1 - \beta_{(0)}^2 \right) \left( (1 - \eta)\beta_{(0)} - 0.51\nu_1\lambda \right) - \frac{c_\mu\theta_{\log}^2}{|\tau|} \left( (1 - \eta)\beta_{(0)}^2 + 7 \right) \right] \\
& \quad - n\theta \cdot \frac{c_\mu\theta_{\log}^2}{|\tau|} (0.5 + c_\mu - 0.5\beta_{(0)}) \\
\text{(SM6.30)} \quad & \geq n\theta \left[ \underbrace{\left( 1 - \beta_{(0)}^2 \right) \left( \frac{\nu_1\lambda}{2\beta_{(1)}}\beta_{(0)} - 0.51\nu_1\lambda \right)}_{f(\beta)} - \frac{8c_\mu\theta_{\log}^2}{|\tau|} \text{Big} \right].
\end{aligned}$$

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 (SM6.23)-(SM6.24), 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:

$$\begin{aligned} f(\beta_{\text{lb}}) &\geq (1 - \beta_{\text{lb}}^2) \left( \frac{\nu_1\lambda}{2\beta_{(1)}}\beta_{\text{lb}} - 0.51\nu_1\lambda \right) \geq (1 - 0.616) \cdot \left( \frac{5}{8} - 0.51 \right) \nu_1\lambda \\ \text{(SM6.31)} \quad &\geq \frac{1}{16\sqrt{2}}\nu_1\lambda \geq \frac{\theta_{\log}^2}{32}\lambda^2. \end{aligned}$$

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

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

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:

$$\begin{aligned} f(\beta_{\text{ub}}) &\geq 0.49 \left( \frac{\lambda}{4\log\theta^{-1}} \right)^2 \cdot \left( \frac{1}{2} \left( 1 - 0.49 \left( \frac{\lambda}{4\log\theta^{-1}} \right)^2 \right) - 0.51\nu_1\lambda \right) \\ \text{(SM6.32)} \quad &\geq \frac{1}{6} \left( \frac{\lambda}{4\log\theta^{-1}} \right)^2 \geq \frac{\theta_{\log}^2}{96}\lambda^2. \end{aligned}$$

Finally combine (SM6.30) with the lower bound of cubic function (SM6.31)-(SM6.32) together with condition  $c_\mu < \frac{c_\lambda^2}{800}$  and  $\nu_1 = \frac{\sqrt{\theta_{\log}}}{2}$ , obtain

$$\begin{aligned} \langle \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], -\text{grad}_{\varphi_{\ell_1}}[\mathbf{a}] \rangle &\geq n\theta \cdot \left( \min\{f(\beta_{\text{ub}}), f(\beta_{\text{lb}})\} - \frac{8c_\mu\theta_{\log}^2}{|\tau|} \right) \\ \text{(SM6.33)} \quad &\geq n\theta \left( \frac{\theta_{\log}^2 c_\lambda^2}{96|\tau|} - \frac{8\theta_{\log}^2 c_\lambda^2}{800|\tau|} \right) \geq 6 \times 10^{-3} n\theta \theta_{\log}^2 c_\lambda^2. \end{aligned}$$

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

As a consequence, we have that

**Corollary SM6.4 (Large gradient for  $\varphi_\rho$ ).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda/\sqrt{k}$  in  $\varphi_\rho$  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  $\rho$  is  $\delta$ -smoothed  $\ell^1$  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  $\frac{4}{5}|\beta_{(0)}| > |\beta_{(1)}| > \frac{1}{4\log\theta^{-1}\lambda}$ ,*

$$\text{(SM6.34)} \quad \langle \boldsymbol{\sigma}_{(0)} \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq cn\theta (\log^{-2}\theta^{-1}) \lambda^2$$

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

*Proof.* Choose  $\iota^* s_{(0)}[\mathbf{a}_0]$  as in Lemma SM6.3, and apply (SM5.22) from Lemma SM5.6 with the constant multiplier of  $\delta$  satisfies  $c''^4 < c/4$ , then utilize  $\theta |\boldsymbol{\tau}| \log^2 \theta^{-1} < c_\mu$  from Definition SM2.1 we have

$$(SM6.35) \quad \langle \sigma_{(0)} \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq cn\theta(\log^{-2} \theta^{-1})\lambda - c''n\theta^2 \geq cn\theta(\log^{-2} \theta^{-1})\lambda/2 \quad \blacksquare$$

**SM6.3. 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  $\varphi_\rho$  will be strongly convex at  $\mathbf{a}$ . We show this in Lemma SM6.5, and the  $\varphi_\rho$  version in Corollary SM6.6 when  $\rho$  is properly defined as in Section SM5.

**Lemma SM6.5 (Strong convexity of  $\varphi_{\ell^1}$  near shift).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda/\sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  satisfying  $|\beta_{(1)}| < \frac{1}{4 \log \theta^{-1}} \lambda$ : for all  $\mathbf{v} \in \mathbb{S}^{p-1} \cap \mathbf{v}^\perp$ ,*

$$(SM6.36) \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} > cn\theta;$$

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

$$(SM6.37) \quad \min_{\ell} \|\bar{\mathbf{a}} - s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max\{\mu, p^{-1}\}.$$

*Proof.* 1. (Expectation of  $\boldsymbol{\chi}$  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  $\boldsymbol{\chi}$  operator shrinks all other smaller entries of correlation vector  $\boldsymbol{\beta}_{\setminus(0)}$  in an even larger shrinking ratio. Firstly we can show  $|\langle \boldsymbol{\beta}_{\setminus(0)}, \mathbf{x}_{\setminus(0)} \rangle|$  is no larger than  $\lambda/2$  with probability at least  $1 - 4\theta$ , since

$$(SM6.38) \quad \begin{aligned} & \mathbb{P} \left[ |\langle \boldsymbol{\beta}_{\setminus(0)}, \mathbf{x}_{\setminus(0)} \rangle| > \frac{\lambda}{2} \right] \\ & \leq \mathbb{P} \left[ |\langle \boldsymbol{\beta}_{\tau \setminus(0)}, \mathbf{x}_{\tau \setminus(0)} \rangle| > \frac{2\lambda}{5} \right] + \mathbb{P} \left[ |\langle \boldsymbol{\beta}_{\tau^c}, \mathbf{x}_{\tau^c} \rangle| > \frac{\lambda}{10} \right] \leq 4\theta \end{aligned}$$

via Corollary SM2.6 and Corollary SM2.7. Now recall from Lemma SM3.2 and the derivation of (SM3.10)-(SM3.11), we know for every  $i \neq (0)$ ,

$$(SM6.39) \quad \begin{aligned} \sigma_i \mathbb{E} \boldsymbol{\chi}[\boldsymbol{\beta}]_i &= n\theta |\beta_i| \mathbb{E}_{\mathbf{s}_i} [1 - \text{erf}_{\beta_i}(\lambda, \mathbf{s}_i)] \\ &\leq n\theta |\beta_i| \mathbb{E}_{g, \mathbf{x}_i} \left[ g^2 \mathbf{1}_{\{|\beta_i g + \beta_{(0)} \mathbf{x}_{(0)} + \beta_{\setminus\{(0), i}\}^* \mathbf{x}_{\setminus\{(0), i}\}}| > \lambda\}} \right] \\ &\leq n\theta |\beta_i| \left( \mathbb{E} g^2 \mathbf{1}_{\{|\beta_i g| > \frac{\lambda}{2}\}} + \mathbb{P}[\mathbf{x}_{(0)} \neq 0] \right. \\ &\quad \left. + \mathbb{P} \left[ |\langle \boldsymbol{\beta}_{\setminus\{(0), i}\}, \mathbf{x}_{\setminus\{(0), i}\} \rangle| > \frac{\lambda}{2} \right] \right) \\ &\leq n\theta |\beta_i| \left( (\mathbb{E} g^2)^{1/2} \mathbb{P} \left[ |\beta_{(1)} g| > \frac{\lambda}{2} \right]^{1/2} + \theta + 4\theta \right) \\ &\leq n\theta |\beta_i| (\exp(-\log^2 \theta^{-1}) + 5\theta) \\ &\leq 6n\theta^2 |\beta_i| \end{aligned}$$

where the third inequality is derived using union bound; the the fourth inequality is the result of (SM6.38), and the fifth inequality is derived from Gaussian tail bound Lemma SM10.1.

2. (Local strong convexity) Let  $\boldsymbol{\gamma} = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{v}$ , for any  $\|\mathbf{v}\|_2 = 1$  we have  $\|\boldsymbol{\gamma}\|_2^2 \leq 1 + \mu p$ . Furthermore:

$$\begin{aligned} |\gamma_{(0)}| &= |\langle \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], \mathbf{v} \rangle| = |\langle \mathbf{P}_{\mathbf{a}_0} \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0], \mathbf{v} \rangle| = |\langle \boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0] - \boldsymbol{\beta}_{(0)} \mathbf{a}, \mathbf{v} \rangle| \\ (SM6.40) \quad &\leq \|\boldsymbol{\iota}^* s_{(0)}[\mathbf{a}_0] - \boldsymbol{\beta}_{(0)} \mathbf{a}\|_2 \leq \sqrt{1 - \beta_{(0)}^2}. \end{aligned}$$

Consider any such  $\mathbf{v}$ , the pseudo Hessian can be lower bounded as

$$\begin{aligned} \mathbf{v}^* \tilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} &= -\boldsymbol{\gamma}^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\gamma} \\ &\geq -\gamma_{(0)}^2 \left\| \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_{(0)} \right\|_2^2 - \sum_{i \neq (0)} \left\| \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_i \right\|_2^2 \gamma_i^2 \\ &\quad - 2 \sum_{i \neq j} \left| \mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| |\gamma_i| |\gamma_j| \\ &\geq -\left(1 - \beta_{(0)}^2\right) \|\mathbf{x}\|_2^2 - \max_{i \neq (0)} \left\| \mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}] \right\|_2^2 \|\boldsymbol{\gamma}\|_2^2 \\ (SM6.41) \quad &\quad - 2 \max_{i \neq j} \left| \mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| \|\boldsymbol{\gamma}\|_1^2, \end{aligned}$$

where the second term is bounded by using its expectation derived in Lemma SM4.2, and utilize  $\mathbb{P}[|s_i| > \lambda/2] < 4\theta$  from (SM6.38),  $\mathbb{E}\boldsymbol{\chi}$  from (SM6.39) and regional condition  $|\boldsymbol{\beta}_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}$  to acquire

$$\begin{aligned} \mathbb{E} \left\| \mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}] \right\|_2^2 &= n\theta \left[ 1 - \mathbb{E}_{s_i} \operatorname{erf}_{\beta_i}(\lambda, s_i) + \mathbb{E}_{s_i} f_{\beta_i}(\lambda, s_i) \right] \\ &\leq \frac{|\mathbb{E}\boldsymbol{\chi}[\boldsymbol{\beta}]_i|}{|\beta_i|} + n\theta \cdot \left( \max_{|s_i| \leq \frac{\lambda}{2}} f_{\beta_i}(\lambda, s_i) + \mathbb{P}\left[|s_i| > \frac{\lambda}{2}\right] \right) \\ &\leq 6n\theta^2 + \frac{2n\theta}{\sqrt{2\pi}} \max_{|s_i| \leq \frac{\lambda}{2}} \left( \frac{\lambda + |s_i|}{|\beta_i|} \cdot \exp\left[-\frac{(\lambda - |s_i|)^2}{2\beta_i^2}\right] \right) + 4n\theta^2 \\ &\leq 10n\theta^2 + n\theta \cdot \log \theta^{-1} \exp(-2 \log^2 \theta^{-1}) \\ (SM6.42) \quad &\leq 11n\theta^2, \end{aligned}$$

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

$$(SM6.43) \quad \begin{cases} \mathcal{E}_{\text{pcurv}} := \left\{ \forall \mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_{\boldsymbol{\tau}}, \gamma(c_\mu)), \left\| \mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}] \right\|_2^2 \leq 11n\theta^2 + \frac{c_\mu n \theta}{p} \right\} \\ \mathcal{E}_{\text{cross}} := \left\{ \forall \mathbf{a} \in \cup_{|\boldsymbol{\tau}| \leq k} \mathfrak{R}(\mathcal{S}_{\boldsymbol{\tau}}, \gamma(c_\mu)), |\boldsymbol{\beta}_{(1)}| \leq \frac{\lambda}{4 \log \theta^{-1}}, \max_{i \neq j \in [\pm p]} \left| \mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathbf{P}_{I(\mathbf{a})} \check{\mathbf{C}}_{\mathbf{x}} \mathbf{e}_j \right| \leq 8n\theta^3 \right\} \\ \mathcal{E}_{\|\mathbf{x}\|_2} := \left\{ \|\mathbf{x}\|_2^2 \leq n\theta + 3\sqrt{n\theta} \log n \right\} \end{cases} .$$

For the Hessian term, on the event  $\mathcal{E}_{\text{pcurv}} \cap \mathcal{E}_{\text{cross}} \cap \mathcal{E}_{\|\mathbf{x}\|_2}$ , and use all  $\mu p^2 \theta^2$ ,  $\mu p \theta |\boldsymbol{\tau}|$  and  $\theta \sqrt{p}$  are all less than  $\frac{c_\mu}{4 \log^2 \theta^{-1}}$ , from Lemma SM2.5, and from lemma assumption with sufficiently

large  $C$  we have  $n > \theta^{-1} 36 \log^2 n$ , thus  $\mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v}$  can be lower bounded from (SM6.41) as

$$\begin{aligned}
\mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} &\geq - \left(1 - \beta_{(0)}^2\right) \left(n\theta + 3\sqrt{n\theta} \log n\right) \\
&\quad - (1 + \mu p) \left(11n\theta^2 + \frac{c_\mu n\theta}{p}\right) - 8p(1 + \mu p) \cdot 8n\theta^3 \\
&\geq -\frac{1}{2}n\theta \cdot (1 - \beta_{(0)}^2) - n\theta \cdot \left(\frac{11c_\mu}{4} + c_\mu^2 + \frac{64c_\mu}{4} + \frac{64c_\mu}{4}\right) \\
\text{(SM6.44)} \quad &\geq -\frac{1}{2}n\theta \cdot \left(1 - \beta_{(0)}^2 + 20c_\mu\right).
\end{aligned}$$

The bounds of  $\beta^* \chi[\beta]$  can be derive on the event whose expectation is drawn from Lemma SM3.2 and (SM6.39) as

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

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

$$\begin{aligned}
\beta^* \chi[\beta] &\geq n\theta |\beta_{(0)}| (|\beta_{(0)}| - \nu_2 \lambda) - c_\mu \|\beta\|_1 \frac{n\theta}{p} \\
&\geq n\theta \left( \beta_{(0)}^2 - \sqrt{\frac{2}{\pi}} \lambda - \frac{c_\mu}{2} \lambda \right) \\
\text{(SM6.45)} \quad &\geq n\theta \left( \beta_{(0)}^2 - \lambda \right).
\end{aligned}$$

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

$$\text{(SM6.46)} \quad \beta_{(0)}^2 \geq \|\beta_\tau\|_2^2 - |\tau| \beta_{(1)}^2 \geq 1 - 2c_\mu - 0.1^2 > 0.9,$$

then we collect all above results and obtain:

$$\text{(SM6.47)} \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_{\ell^1}](\mathbf{a}) \mathbf{v} = \mathbf{v}^* \widetilde{\nabla}^2 \varphi_{\ell^1}(\mathbf{a}) \mathbf{v} - \beta^* \chi[\beta] \geq \left(1.5\beta_{(0)}^2 - 0.5 - \lambda - 20c_\mu\right) n\theta \geq 0.3n\theta,$$

with probability at least

$$\text{(SM6.48)} \quad 1 - \underbrace{\mathbb{P}[\mathcal{E}_{\text{cross}}^c]}_{\text{Lemma SM4.4}} - \underbrace{\mathbb{P}[\mathcal{E}_{\text{pcurv}}^c]}_{\text{Corollary SM4.3}} - \underbrace{\mathbb{P}[\mathcal{E}_{\|\mathbf{x}\|_2}^c]}_{\text{Lemma SM1.2}} - \underbrace{\mathbb{P}[\mathcal{E}_\chi^c]}_{\text{Corollary SM3.4}} \geq 1 - c'/n.$$

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 (SM6.47), provides the upper bound on  $\|\mathbf{a}_* - \mathbf{a}_0\|_2^2$  via

$$\begin{aligned}
\varphi_{\ell^1}(\mathbf{a}_*) &\geq \varphi_{\ell^1}(\mathbf{a}_0) + \langle \mathbf{a}_* - \mathbf{a}_0, \text{grad}[\varphi_{\ell^1}](\mathbf{a}_0) \rangle + \frac{0.3}{2} n\theta \|\mathbf{a}_* - \mathbf{a}_0\|_2^2 \\
\text{(SM6.49)} \quad &\implies \|\text{grad}[\varphi_{\ell^1}](\mathbf{a}_0)\|_2 \geq 0.15n\theta \|\mathbf{a}_* - \mathbf{a}_0\|_2
\end{aligned}$$

Thus we only require to bound the gradient at  $\mathbf{a}_0$ , whose coefficients  $\boldsymbol{\alpha} = \mathbf{e}_0$  and correlation  $\boldsymbol{\beta}$  has properties  $\beta_0 = 1$  and  $\|\boldsymbol{\beta}_{\setminus 0}\|_\infty \leq \mu$  hence  $\|\boldsymbol{\beta}_{\setminus 0}\|_2 \leq \sqrt{2p}\mu$ . Expand the gradient term and condition on  $\mathcal{E}_\chi$ , 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

$$\begin{aligned}
 \|\text{grad}[\varphi_{\ell^1}](\mathbf{a}_0)\|_2 &= \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} (\boldsymbol{\chi}[\boldsymbol{\beta}] - \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \mathbf{e}_0)\|_2 \leq \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}\|_2 \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\setminus 0}\|_2 \\
 &\leq \sqrt{1 + \mu p} \left( 6n\theta^2 \|\boldsymbol{\beta}_{\setminus 0}\|_2 + n\theta \cdot \frac{c_\mu}{p^{3/2}} \cdot \sqrt{2p} \right) \\
 &\leq n\theta \sqrt{1 + \mu p} \left( 6\mu \sqrt{2p} \cdot \theta + \frac{2c_\mu}{p} \right) \\
 &\leq n\theta \left( 3c_\mu \mu + 6\mu \cdot \sqrt{2\mu} \cdot p\theta + \frac{2c_\mu}{p} + \frac{2c_\mu \sqrt{\mu}}{\sqrt{p}} \right) \\
 \text{(SM6.50)} \quad &\leq 7\sqrt{c_\mu} n\theta \cdot \max \left\{ \mu, \frac{1}{p} \right\}.
 \end{aligned}$$

Thus we conclude that with sufficiently small  $c_\mu$ :

$$\text{(SM6.51)} \quad \|\mathbf{a}_* - \mathbf{a}_0\|_2 \leq 50\sqrt{c_\mu} \max \left\{ \mu, p^{-1} \right\} \leq \frac{1}{2} \max \left\{ \mu, p^{-1} \right\}. \quad \blacksquare$$

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  $\varphi_\rho$  we have

**Corollary SM6.6 (Strong convexity of  $\varphi_\rho$  of near shift).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{k}$  in  $\varphi_\rho$  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  $\rho$  is  $\delta$ -smoothed  $\ell^1$  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  $|\boldsymbol{\beta}_{(1)}| < \nu_1 \lambda$ : for all  $\mathbf{v} \in \mathbb{S}^{p-1} \cap \mathbf{a}^\perp$ ,*

$$\text{(SM6.52)} \quad \mathbf{v}^* \widetilde{\text{Hess}}[\varphi_\rho](\mathbf{a}) \mathbf{v} > cn\theta;$$

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

$$\text{(SM6.53)} \quad \min_{\ell} \|\bar{\mathbf{a}} - s_\ell[\mathbf{a}_0]\|_2 \leq \frac{1}{2} \max \left\{ \mu, p^{-1} \right\}$$

*Proof.* The strong convexity (SM6.52) is derived by combining (SM6.36) and (SM5.23) by letting constant multiplier of  $\delta$  satisfies  $c^{1/4} < 10^{-3}c$ . On the other hand the local minimizer near solution (SM6.53) is derived via combining (SM6.49), (SM5.21) and utilize both  $\theta\sqrt{p} < c_\mu$  and  $\mu p^2 \theta^2 < c_\mu$  such that:

$$\begin{aligned}
 \|\text{grad}[\varphi_\rho](\mathbf{a})\|_2 &\leq \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}\|_2 \left\| \boldsymbol{\chi}[\boldsymbol{\beta}] - \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda^\delta \left[ \check{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota} \mathbf{a} \right] \right\|_2 + \|\boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}\|_2 \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\setminus 0}\|_2 \\
 &\leq \sqrt{1 + \mu p} \cdot n\theta^3 + 7\sqrt{c_\mu} n\theta \cdot \max \left\{ \mu, p^{-1} \right\} \\
 \text{(SM6.54)} \quad &\leq 8n\theta \sqrt{c_\mu} \cdot \max \left\{ \mu, p^{-1} \right\} \quad \blacksquare
 \end{aligned}$$

**SM6.4. Retraction toward subspace.** As in [Figure SM4](#), the function value grows in direction away from subspace  $\mathcal{S}_\tau$ , we will illustrate this phenomenon by proving the negative gradient direction  $-\mathbf{g}$  will point toward the subspace  $\mathcal{S}_\tau$ . To show this, we prove for every coefficients of  $\mathbf{a}$  as  $\boldsymbol{\alpha}$ , there exists coefficients of  $\mathbf{g}$  as  $\boldsymbol{\zeta}$  satisfies

$$(SM6.55) \quad \langle \boldsymbol{\alpha}_{\tau^c}(\mathbf{g}), \boldsymbol{\alpha}_{\tau^c}(\mathbf{a}) \rangle > c \|\boldsymbol{\alpha}_{\tau^c}\|_2 \|\boldsymbol{\zeta}_{\tau^c}\|_2$$

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*  $\mathcal{S}_\tau$ . This retractive phenomenon is true for gradient of both  $\varphi_{\ell^1}$  and  $\varphi_\rho$ .

**Lemma SM6.7 (Retraction of  $\varphi_{\ell^1}$  toward subspace).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda/\sqrt{k}$  in  $\varphi_{\ell^1}$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathcal{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  such that if*

$$(SM6.56) \quad d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \geq \gamma(c_\mu)/2$$

then for every  $\boldsymbol{\alpha}$  satisfying  $\mathbf{a} = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}$ , there exists some  $\boldsymbol{\zeta}$  satisfying  $\text{grad}[\varphi_{\ell^1}](\mathbf{a}) = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\zeta}$  that

$$(SM6.57) \quad \langle \boldsymbol{\zeta}_{\tau^c}, \boldsymbol{\alpha}_{\tau^c} \rangle \geq \frac{1}{4n\theta} \|\boldsymbol{\zeta}_{\tau^c}\|_2^2.$$

*Proof.* Write  $\gamma = \gamma(c_\mu)$  Recall the gradient can be derived as

$$(SM6.58) \quad \text{grad}[\varphi_{\ell^1}](\mathbf{a}) = -\mathbf{P}_{\mathbf{a}_0^\perp} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\chi}[\boldsymbol{\beta}] = (\mathbf{a}\mathbf{a}^* - \mathbf{I}) \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\chi}[\boldsymbol{\beta}] = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} (\boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \boldsymbol{\alpha} - \boldsymbol{\chi}[\boldsymbol{\beta}]),$$

for every  $\boldsymbol{\alpha}$  satisfies  $\mathbf{a} = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}$ . Now via [Corollary SM3.4](#), condition on the event:

$$(SM6.59) \quad \mathcal{E}_\chi := \left\{ \sigma_i \boldsymbol{\chi}[\boldsymbol{\beta}]_i \leq \begin{cases} n\theta \cdot |\boldsymbol{\beta}_i| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau \\ n\theta \cdot |\boldsymbol{\beta}_i| 4\theta |\tau| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau^c \end{cases}, \quad \sigma_i \boldsymbol{\chi}[\boldsymbol{\beta}]_i \geq n\theta \cdot \mathcal{S}_{\sqrt{2/\pi}\lambda} [|\boldsymbol{\beta}_i|] \right\},$$

and on this event, utilize [Lemma SM2.5](#), bounds of  $\boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}]$  and  $\|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau^c}\|_2$  can be derived with  $c_\mu < \frac{1}{100}$  as:

$$(SM6.60) \quad \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \leq n\theta \left( \|\boldsymbol{\beta}_\tau\|_2^2 + 4\theta |\tau| \|\boldsymbol{\beta}_{\tau^c}\|_2^2 + c_\mu \right) \geq n\theta (1 + c_\mu + 4c_\mu^2 + c_\mu) \leq \frac{3}{2}n\theta$$

$$(SM6.61) \quad \boldsymbol{\beta}^* \boldsymbol{\chi}[\boldsymbol{\beta}] \geq n\theta \left( \|\boldsymbol{\beta}_\tau\|_2^2 - \sqrt{2/\pi}\lambda \|\boldsymbol{\beta}_\tau\|_1 - c_\mu \right) \geq n\theta \left( 1 - 4c_\mu - \sqrt{2/\pi}c_\lambda - c_\mu \right) \geq \frac{1}{2}n\theta$$

$$(SM6.62) \quad \|\boldsymbol{\chi}[\boldsymbol{\beta}]_{\tau^c}\|_2 \leq 4n\theta^2 |\tau| \|\boldsymbol{\beta}_{\tau^c}\|_2 + \frac{c_\mu n\theta}{p} \sqrt{p} \leq n\theta (4c_\mu \gamma + c_\mu \gamma) \leq \frac{1}{20}n\theta \gamma.$$

Let  $\alpha(\mathbf{g}) = \beta^* \chi[\beta] \alpha - \chi[\beta]$ , derive

$$\begin{aligned}
 & \langle \alpha(\mathbf{g})_{\tau^c}, \alpha_{\tau^c} \rangle - \frac{1}{4n\theta} \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 \\
 &= \beta^* \chi[\beta] \|\alpha_{\tau^c}\|_2^2 - \langle \alpha_{\tau^c}, \chi[\beta]_{\tau^c} \rangle \\
 &\quad - \frac{1}{4n\theta} \|\beta^* \chi[\beta] \alpha_{\tau^c} - \chi[\beta]_{\tau^c}\|_2^2 \\
 &\geq \beta^* \chi[\beta] \|\alpha_{\tau^c}\|_2^2 - \|\alpha_{\tau^c}\|_2 \|\chi[\beta]_{\tau^c}\|_2 \\
 &\quad - \frac{1}{2n\theta} |\beta^* \chi[\beta]|^2 \|\alpha_{\tau^c}\|_2^2 - \frac{1}{2n\theta} \|\chi[\beta]_{\tau^c}\|_2^2 \\
 \text{(SM6.63)} \quad &\geq (\beta^* \chi[\beta] - \frac{1}{2n\theta} (\beta^* \chi[\beta])^2) \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20} n\theta\gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000} n\theta\gamma^2,
 \end{aligned}$$

notice that this is a quadratic function of  $\beta^* \chi[\beta]$  with negative leading coefficient and zeros at  $\{0, 2n\theta\}$ , hence (SM6.63) is minimized when  $\beta^* \chi[\beta] = \frac{1}{2}n\theta$ . Plugging in,

$$\text{(SM6.64)} \quad \text{(SM6.63)} \geq \frac{3}{8}n\theta \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20}n\theta\gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000}n\theta\gamma^2$$

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 (SM6.64) is minimized at  $\|\alpha_{\tau^c}\|_2 = \frac{\gamma}{2}$ . Plugging in again,

$$\text{(SM6.65)} \quad \text{(SM6.64)} \geq \frac{3}{8}n\theta \|\alpha_{\tau^c}\|_2^2 - \frac{1}{20}n\theta\gamma \|\alpha_{\tau^c}\|_2 - \frac{1}{1000}n\theta\gamma^2 \geq (\frac{3}{32} - \frac{1}{80} - \frac{1}{1000}) n\theta\gamma^2 > 0$$

which concludes our proof. ■

As a consequence, we have that

**Corollary SM6.8 (Retraction of  $\varphi_\rho$  toward the subspace).** *Suppose that  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \text{BG}(\theta)$  in  $\mathbb{R}^n$ , and  $k, c_\mu$  such that  $(\mathbf{a}_0, \theta, k)$  satisfies the sparsity-coherence condition  $\text{SCC}(c_\mu)$ . Define  $\lambda = c_\lambda / \sqrt{|k|}$  in  $\varphi_\rho$  with  $c_\lambda \in (0, \frac{1}{3}]$ , then there exists some numerical constants  $C, c, c', c'', \bar{c} > 0$  such that if  $\rho$  is  $\delta$ -smoothed  $\ell^1$  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  $\mathbf{a} \in \cup_{|\tau| \leq k} \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  such that if*

$$\text{(SM6.66)} \quad d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \geq \gamma(c_\mu)/2$$

then for every  $\alpha$  satisfying  $\mathbf{a} = \iota^* \mathbf{C}_{\mathbf{a}_0} \alpha$ , there exists some  $\zeta$  satisfying  $\text{grad}[\varphi_\rho](\mathbf{a}) = \iota^* \mathbf{C}_{\mathbf{a}_0} \zeta$  that

$$\text{(SM6.67)} \quad \langle \zeta_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{1}{6n\theta} \|\zeta_{\tau^c}\|_2^2.$$

*Proof.* Write  $\gamma = \gamma(c_\mu)$ . Define

$$\chi_{\ell^1}[\beta] = \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda [\check{\mathbf{a}} * \mathbf{y}], \quad \chi_\rho[\beta] = \check{\mathbf{C}}_{\mathbf{x}_0} \mathcal{S}_\lambda^\delta [\check{\mathbf{a}} * \mathbf{y}],$$

which, and on event (SM6.59) and Lemma SM5.6, we know

$$\text{(SM6.68)} \quad \beta^* \chi_{\ell^1}[\beta] \leq \frac{3}{2}n\theta,$$

$$\text{(SM6.69)} \quad \|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 \leq \frac{1}{20}n\theta\gamma,$$

$$\text{(SM6.70)} \quad \|\chi_{\ell^1}[\beta] - \chi_\rho[\beta]\|_2 \leq c_1 n\theta^4,$$

for some constant  $c_1 > 0$ . Now given any  $\alpha$  satisfies  $\mathbf{a} = \iota^* \mathbf{C}_{a_0} \alpha$ , the gradient of both objective can be derived as

$$\begin{aligned} \text{grad}[\varphi_{\ell^1}](\mathbf{a}) &= -\mathbf{P}_{\mathbf{a}^\perp} \iota^* \mathbf{C}_{a_0} \text{prox}_{\lambda \|\cdot\|_1}[\check{\mathbf{a}} * \mathbf{y}] = (\mathbf{a}\mathbf{a}^* - \mathbf{I}) \iota^* \mathbf{C}_{a_0} \chi_{\ell^1}[\beta] \\ \text{(SM6.71)} \quad &= \iota^* \mathbf{C}_{a_0} (\beta^* \chi_{\ell^1}[\beta] \alpha - \chi_{\ell^1}[\beta]), \end{aligned}$$

$$\begin{aligned} \text{grad}[\varphi_\rho](\mathbf{a}) &= -\mathbf{P}_{\mathbf{a}^\perp} \iota^* \mathbf{C}_{a_0} \text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}] = (\mathbf{a}\mathbf{a}^* - \mathbf{I}) \iota^* \mathbf{C}_{a_0} \chi_\rho[\beta] \\ \text{(SM6.72)} \quad &= \iota^* \mathbf{C}_{a_0} (\beta^* \chi_\rho[\beta] \alpha - \chi_\rho[\beta]). \end{aligned}$$

In the same spirit, define the coefficient of each gradient vector

$$\text{(SM6.73)} \quad \zeta_{\ell^1} = \beta^* \chi_{\ell^1}[\beta] \alpha - \chi_{\ell^1}[\beta],$$

$$\text{(SM6.74)} \quad \zeta_\rho = \beta^* \chi_\rho[\beta] \alpha - \chi_\rho[\beta],$$

which, by norm inequality from (SM6.68)-(SM6.70) and Lemma SM6.7, we can derive

$$\text{(SM6.75)} \quad \|\zeta_{\ell^1} - \zeta_\rho\|_2 \leq \|(\mathbf{I} - \alpha\beta^*)(\chi_\rho[\beta] - \chi_{\ell^1}[\beta])\|_2 \leq c_1 n \theta^4,$$

$$\text{(SM6.76)} \quad \|(\zeta_{\ell^1})_{\tau^c}\|_2 \geq |\beta^* \chi_{\ell^1}[\beta]| \|\alpha_{\tau^c}\|_2 - \|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 \geq \frac{1}{5} n \theta \gamma,$$

$$\text{(SM6.77)} \quad \langle (\zeta_{\ell^1})_{\tau^c}, \alpha_{\tau^c} \rangle \geq \frac{1}{4n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2,$$

where the first inequality is derived by observing  $(\mathbf{I} - \alpha\beta^*)$  is a projection operator, as such:

$$\begin{aligned} \beta^* \alpha &= \mathbf{a}^* \iota^* \mathbf{C}_{a_0} \alpha = \mathbf{a}^* \mathbf{a} = 1, \\ (\mathbf{I} - \alpha\beta^*)^2 &= \mathbf{I} - 2\alpha\beta^* + \alpha(\beta^* \alpha)\beta^* = \mathbf{I} - \alpha\beta^*. \end{aligned}$$

Now we are ready to derive (SM6.57):

$$\begin{aligned} \langle (\zeta_\rho)_{\tau^c}, \alpha_{\tau^c} \rangle &\geq \langle (\zeta_{\ell^1})_{\tau^c}, \alpha_{\tau^c} \rangle - \|\alpha_{\tau^c}\|_2 \|\zeta_\rho - \zeta_{\ell^1}\|_2 \\ &\geq \frac{1}{4n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2 - c_1 n \theta^4 \gamma \\ &\geq \frac{1}{12n\theta} \|(\zeta_{\ell^1})_{\tau^c}\|_2^2 \\ &\quad + \frac{1}{6n\theta} \left( \|(\zeta_\rho)_{\tau^c}\|_2^2 - 2 \|(\zeta_{\ell^1})_{\tau^c}\|_2 \|\zeta_{\ell^1} - \zeta_\rho\|_2 - \|\zeta_{\ell^1} - \zeta_\rho\|_2^2 \right) - c_1 n \theta^4 \gamma \\ &\geq \frac{1}{6n\theta} \|(\zeta_\rho)_{\tau^c}\|_2^2 + \frac{1}{12n\theta} \left( \frac{1}{5} n \theta \gamma \right)^2 - \frac{1}{3n\theta} \left( \frac{1}{5} n \theta \gamma \right) (c_1 n \theta^4) \\ &\quad - \frac{1}{6n\theta} (c_1 n \theta^4)^2 - c_1 n \theta^4 \gamma \\ \text{(SM6.78)} \quad &\geq \frac{1}{6n\theta} \|(\zeta_\rho)_{\tau^c}\|_2^2. \end{aligned}$$

where the last inequality is true since  $\theta^3 \ll \gamma$ . ■

**SM6.5. 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  $\mathcal{S}_\tau$ , either one of the following in true

$$\text{(SM6.79)} \quad \lambda_{\min}(\text{Hess}[\varphi_\rho](\mathbf{a})) \leq -c_1 n \theta \lambda,$$

$$\text{(SM6.80)} \quad \langle \sigma_{(0)} \iota^* s_{(0)}[\mathbf{a}_0], -\text{grad}[\varphi_\rho](\mathbf{a}) \rangle \geq c_2 n \theta (\log^{-2} \theta^{-1}) \lambda^2,$$

$$\text{(SM6.81)} \quad \text{Hess}[\varphi_\rho](\mathbf{a}) \succ c_3 n \theta \cdot \mathbf{P}_{\mathbf{a}^\perp},$$

all local minimizer  $\bar{\mathbf{a}}$  satisfies for some  $\mathbf{a}_* \in \{\pm \boldsymbol{\nu}^* s_\ell[\mathbf{a}] \mid \ell \in [\pm p_0]\}$ ,

$$(SM6.82) \quad \|\bar{\mathbf{a}} - \mathbf{a}_*\|_2 \leq c_4 \sqrt{c_\mu} \max\{\mu, p_0^{-1}\},$$

and whenever  $\frac{\gamma}{2} \leq d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma$ , coefficient of  $\mathbf{a}$  and its gradient  $\mathbf{g}$ ,  $\boldsymbol{\alpha}$ , written as  $\boldsymbol{\zeta}$ , satisfies

$$(SM6.83) \quad \langle \boldsymbol{\zeta}_{\tau^c}, \boldsymbol{\alpha}_{\tau^c} \rangle \geq \frac{c_5}{n\theta} \|\boldsymbol{\zeta}_{\tau^c}\|_2^2.$$

To connect the geometric results introduced in [Lemma SM6.1](#), [Lemma SM6.3](#), [Lemma SM6.5](#) and [Lemma SM6.7](#), we are only required to prove the required signal condition claimed in [Theorem 4.1](#) is necessary from [Definition SM2.1](#). 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  $\delta$  in the pseudo-Huber penalty  $\rho(x) = \sqrt{x^2 + \delta^2}$  approximate  $|x|$  sufficiently well, hence results of [Corollary SM6.2](#), [Corollary SM6.4](#), [Corollary SM6.6](#) and [Corollary SM6.8](#) also holds.

*Proof.* Firstly we will show when largest solution subspace dimension  $k = 4p_0\theta$ , the signal condition of [Definition SM2.1](#) will be satisfied. Recall that the signal condition of [Theorem 4.1](#) requests

$$(SM6.84) \quad \frac{2}{p_0 \log^2 p_0} \leq \theta \leq \frac{c}{(p_0 \sqrt{\mu} + \sqrt{p_0}) \log^2 p_0},$$

since  $p = 3p_0 - 2$ , this implies the lower bounds for sparsity  $\theta$  as:

$$(SM6.85) \quad \theta \geq \frac{1}{2p_0 (\frac{1}{2} \log p_0)^2} \geq \frac{1}{p \log^2 \theta^{-1}};$$

the upper bound of  $\theta$  via  $\theta \sqrt{p_0} \log^2 p_0 \leq c$ :

$$(SM6.86) \quad \theta \leq \frac{9c}{\sqrt{p_0} (3 \log p_0)^2} \leq \frac{16c}{\sqrt{p} \log^2 \theta^{-1}}, \quad \theta \leq \frac{4c^2}{k \log^4 p_0} \leq \frac{36c^2}{k (3 \log p_0)^2} \leq \frac{36c^2}{k \log^2 \theta^{-1}};$$

and the upper bound for coherence  $\mu$  as:

$$(SM6.87) \quad \begin{aligned} \mu \max\{k^2, (p\theta)^2\} \log^2 \theta^{-1} &\leq \mu \max\{16(p_0\theta)^2, 9(p_0\theta)^2\} \log^2 \theta^{-1} \\ &\leq 16 (\sqrt{\mu} p_0 \theta)^2 \log^2 p_0 \leq 16c. \end{aligned}$$

Therefore [Definition SM2.1](#) holds if  $\max\{16c, 36c^2\} \leq c_\mu/4$  via [\(SM6.85\)](#)-[\(SM6.87\)](#).

Furthermore, we know from lemma assumption all interested  $\mathbf{a}$  are near subspace  $\mathcal{S}_\tau$  by

$$(SM6.88) \quad \begin{aligned} d_\alpha(\mathbf{a}, \mathcal{S}_\tau) &\leq \frac{c}{\sqrt{p_0} \log^2 \theta^{-1}} \cdot \min \left\{ \frac{1}{\sqrt{\theta}}, \frac{1}{\sqrt{\mu}} \cdot \frac{1}{\mu (p_0\theta)^{3/2}} \right\} \\ &\leq \frac{c}{\log^2 \theta^{-1}} \min \left\{ \frac{2}{\sqrt{k}}, \frac{1}{\sqrt{p_0\mu}}, \frac{4}{\mu p_0 \sqrt{\theta} k} \right\} \leq \gamma \end{aligned}$$

where  $\gamma$  is defined in [Definition SM2.3](#) 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  $\ell^1$  smoothed sparse surrogate defined in [Definition SM5.2](#), by observing that it is convex, smooth, even, whose second order derivative (according to [Table SM1](#))  $\nabla^2 \rho(x) = \frac{\delta^2}{(x^2 + \delta^2)^{3/2}}$  is monotone decreasing in  $|x|$ . More importantly

$$(SM6.89) \quad \sup_{x \in \mathbb{R}} |\rho(x) - |x|| = |\rho(0) - |0|| = \delta.$$

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  $\varphi_\rho$ . We obtain the geometrical results in [Corollary SM6.2](#) when  $|\beta_{(1)}| \geq \frac{4}{5} |\beta_{(0)}|$ , [Corollary SM6.4](#) when  $\frac{4}{5} |\beta_{(0)}| \geq |\beta_{(1)}| \geq \frac{\lambda}{4 \log^2 \theta^{-1}}$  and [Corollary SM6.6](#) when  $\frac{\lambda}{4 \log^2 \theta^{-1}} \geq |\beta_{(1)}|$ , and the retraction result in [Corollary SM6.8](#).  $\blacksquare$

**SM7. Analysis of algorithm — minimization within widened subspace.** In this section, we prove convergence of the first part of our algorithm—minimization of  $\varphi_\rho$  near  $\mathcal{S}_\tau$ . We begin by proving the initialization method guarantees that  $\mathbf{a}^{(0)}$  is near  $\mathcal{S}_\tau$ , in the sense that

$$(SM7.1) \quad d_\alpha(\mathbf{a}^{(0)}, \mathcal{S}_\tau) \leq \gamma,$$

where the distance  $d_\alpha$  is defined in [\(4.16\)](#). We then demonstrate that small-stepping curvilinear search converges to a desired local minimum of  $\varphi_\rho$  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  $\mathcal{S}_\tau$  and (ii) the geometric properties of  $\varphi_\rho$  near  $\mathcal{S}_\tau$ .

**SM7.1. 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

$$(SM7.2) \quad \mathbf{a}^{(-1)} = \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \sum_{\ell \in \tau} \mathbf{x}_{0\ell} \boldsymbol{\iota}_{p_0}^* s_\ell[\mathbf{a}_0] \right],$$

and is very close to the subspace  $\mathcal{S}_\tau$ :

**Lemma SM7.1 (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*

$$(SM7.3) \quad \mathbf{a}^{(0)} = -\mathbf{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\ell^1} \left( \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \mathbf{0}^{p_0-1}; [\bar{\mathbf{y}}_0; \dots; \bar{\mathbf{y}}_{p_0-1}]; \mathbf{0}^{p_0-1} \right] \right),$$

with  $\lambda = 0.2/\sqrt{p\theta}$  in  $\varphi_1$ . 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  $\max_{i \neq j} |\langle \boldsymbol{\iota}_{p_0}^* s_i[\mathbf{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\mathbf{a}_0] \rangle| \leq \mu$ . 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

$$(SM7.4) \quad d_\alpha(\mathbf{a}^{(0)}, \mathcal{S}_\tau) \leq \frac{c_\mu}{4 \log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\tau|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\tau|} \right\}.$$

*Proof.* 1. (Distance to  $\mathcal{S}_\tau$  from  $\mathbf{a}^{(0)}$ ) Let  $\eta = \|\boldsymbol{\iota}_{p_0}^*(\mathbf{a}_0 * \mathbf{x})\|_2 = \|\boldsymbol{\iota}_{p_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}\|_2$  and  $\gamma = \gamma(c_\mu)$ , as in [\(SM7.4\)](#). Expand the expression of  $\mathbf{a}^{(0)}$  from [\(SM7.3\)](#) we have

$$(SM7.5) \quad \begin{aligned} \mathbf{a}^{(0)} &= \mathbf{P}_{\mathbb{S}^{p-1}} \boldsymbol{\iota}^* \tilde{\mathbf{C}}_{\mathbf{y}} \mathcal{S}_\lambda \left[ \tilde{\mathbf{C}}_{\mathbf{y}} \boldsymbol{\iota}_{p_0} \mathbf{P}_{\mathbb{S}^{p_0-1}} \boldsymbol{\iota}_{p_0}^* (\mathbf{a}_0 * \mathbf{x}) \right] \\ &= \mathbf{P}_{\mathbb{S}^{p-1}} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\chi} \left[ \frac{1}{\eta} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}_{p_0} \boldsymbol{\iota}_{p_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x} \right] \end{aligned}$$

To relate  $\mathbf{a}^{(0)}$  to its coefficient, introduce the truncated autocorrelation matrix  $\widetilde{\mathbf{M}} = \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\nu}_{p_0} \boldsymbol{\nu}_{p_0}^* \mathbf{C}_{\mathbf{a}_0}$ , define  $\widetilde{\boldsymbol{\alpha}}, \widetilde{\boldsymbol{\beta}}$  as

$$(SM7.6) \quad \widetilde{\boldsymbol{\beta}} = \frac{1}{\eta} \widetilde{\mathbf{M}} \mathbf{x}, \quad \widetilde{\boldsymbol{\alpha}} = \boldsymbol{\chi} \left[ \frac{1}{\eta} \widetilde{\mathbf{M}} \mathbf{x} \right] = \boldsymbol{\chi}[\widetilde{\boldsymbol{\beta}}]$$

and note that  $\widetilde{\mathbf{M}}$  is bounded entrywise as

$$(SM7.7) \quad \left| \widetilde{\mathbf{M}}_{ij} \right| \leq \begin{cases} 1 & i = j \in [-p_0 + 1, p_0 - 1] \\ \mu & i \neq j \in [-p_0 + 1, p_0 - 1], |i - j| < p_0 \\ 0 & \text{otherwise} \end{cases}$$

From (SM7.5), we can write  $\mathbf{a}^{(0)} = \mathbf{P}_{\mathbb{S}^{p-1}} \boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0} \widetilde{\boldsymbol{\alpha}}$ , meaning that the normalized version of  $\widetilde{\boldsymbol{\alpha}}$  is a valid coefficient vector for  $\mathbf{a}^{(0)}$ . Let  $\boldsymbol{\tau}^c = [\pm 2p_0] \setminus \boldsymbol{\tau}$ . The distance  $d_{\alpha}$  to subspace  $\mathcal{S}_{\boldsymbol{\tau}}$  (4.16) is upper bounded as

$$\begin{aligned} d_{\alpha}(\mathbf{a}^{(0)}, \mathcal{S}_{\boldsymbol{\tau}}) &\leq \frac{\|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2}{\|\boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0} \widetilde{\boldsymbol{\alpha}}\|_2} \leq \frac{\|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2}{\|\boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0} \widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}}\|_2 - \|\boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0} \widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2} \\ &\leq \frac{\|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2}{\sqrt{1 - \mu |\boldsymbol{\tau}|} \|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}}\|_2 - \sqrt{1 + \mu p} \|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2} \end{aligned}$$

where the last inequality is derived with Lemma SM2.4. Therefore, it is sufficient to show

$$(SM7.8) \quad \left(1 + \gamma \sqrt{1 + \mu p}\right) \|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2 \leq \gamma \sqrt{1 - \mu |\boldsymbol{\tau}|} \|\widetilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}}\|_2$$

to complete the proof that  $d_{\alpha}(\mathbf{a}^{(0)}, \mathcal{S}_{\boldsymbol{\tau}}) \leq \gamma$ .

2. (Bound  $\eta$ ) Condition on the following two events

$$(SM7.9) \quad \mathcal{E}_{\boldsymbol{\tau}} := \{|\boldsymbol{\tau}| < 4p_0\theta\}, \quad \mathcal{E}_{\|\mathbf{x}\|_2} := \left\{ \sqrt{p_0\theta} \leq \|\mathbf{x}\|_2 \leq \sqrt{3p_0\theta} \right\}$$

and utilize  $\mu$  bound from Lemma SM2.5 such that  $\mu |\boldsymbol{\tau}| < 0.1$ . An upper bound on  $\eta$  can be obtained using properties of  $\widetilde{\mathbf{M}}$  of (SM7.7):

$$(SM7.10) \quad \eta = \|\boldsymbol{\nu}_{p_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}\|_2 \leq \|\boldsymbol{\nu}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}\|_2 \leq \sqrt{1 + \mu |\boldsymbol{\tau}|} \|\mathbf{x}\|_2 \leq 2\sqrt{p_0\theta}$$

To lower bound  $\eta$ , use  $\eta^2 = \mathbf{g}^* \mathbf{P}_{\boldsymbol{\tau}} \widetilde{\mathbf{M}} \mathbf{P}_{\boldsymbol{\tau}} \mathbf{g}$  where  $\mathbf{g}$  is the standard Gaussian vector. Observe the submatrix of  $\widetilde{\mathbf{M}}$  is diagonal dominant:

$$(SM7.11) \quad \begin{cases} \widetilde{\mathbf{M}}_{ii} = \|\boldsymbol{\nu}_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 \in [0, 1] \\ \text{tr}(\widetilde{\mathbf{M}}) = \sum_{i \in [\pm p_0]} \|\boldsymbol{\nu}_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 = \|\mathbf{a}_0\|_2^2 + \sum_{i=1}^{p_0-1} \left( \|\boldsymbol{\nu}_{p_0}^* s_i[\mathbf{a}_0]\|_2^2 + \|\boldsymbol{\nu}_{p_0}^* s_{i-p_0}[\mathbf{a}_0]\|_2^2 \right) = p_0 \end{cases}$$

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  $\mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau$  can be written as sum of independent r.v.s as:

$$\text{tr} \left( \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \right) = \sum_{i \in [\pm p_0]} w_i \left\| \boldsymbol{\nu}_{p_0}^* s_i [\mathbf{a}_0] \right\|_2^2,$$

Bernstein inequality [Lemma SM10.4](#) and [\(SM7.11\)](#) gives

$$\begin{aligned} \mathbb{P} \left[ \text{tr} \left( \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \right) < \frac{3p_0\theta}{4} \right] &\leq \mathbb{P} \left[ \text{tr} \left( \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \right) - p_0\theta \leq -\frac{p_0\theta}{4} \right] \\ \text{(SM7.12)} \quad &\leq 2 \exp \left( \frac{-(p_0\theta/4)^2}{2p_0\theta + p_0\theta/2} \right) \leq 2 \exp \left( \frac{-p_0\theta}{40} \right), \end{aligned}$$

thus condition on  $\boldsymbol{\omega}$  satisfies  $\text{tr} \left( \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \right) \geq 3p_0\theta/4$  and  $\mathcal{E}_\tau$ , expectation  $\eta^2$  has lower bound

$$\mathbb{E}_{\mathbf{g}|\mathbf{w}} \eta^2 = \mathbb{E}_{\mathbf{g}|\mathbf{w}} \left[ \mathbf{g}^* \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \mathbf{g} \right] = \text{tr} \left( \mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau \right) \geq \frac{3p_0\theta}{4}$$

then apply Bernstein inequality again by first writing svd of  $\mathbf{P}_\tau \widetilde{\mathbf{M}} \mathbf{P}_\tau = \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^*$  with  $\boldsymbol{\Sigma}$  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  $\eta^2$  to be small as

$$\begin{aligned} \mathbb{P}_{\mathbf{g}|\mathbf{w}} \left[ \eta^2 < \frac{p_0\theta}{2} \right] &\leq \mathbb{P}_{\mathbf{g}|\mathbf{w}} \left[ \eta^2 - \mathbb{E}_{\mathbf{g}|\mathbf{w}} \eta^2 < -\frac{p_0\theta}{4} \right] \\ \text{(SM7.13)} \quad &\leq 2 \exp \left( \frac{-(p_0\theta/4)^2}{2(1 + \mu |\tau|)(12p_0\theta + p_0\theta/2)} \right) \leq 2 \exp \left( \frac{-p_0\theta}{440} \right) \end{aligned}$$

by applying moment bounds  $(\sigma^2, R) = (12p_0\theta(1 + \mu |\tau|), 2(1 + \mu |\tau|))$ . We thereby define event

$$\text{(SM7.14)} \quad \mathcal{E}_\eta = \left\{ \sqrt{p_0\theta/2} \leq \eta \leq 2\sqrt{p_0\theta} \right\},$$

which holds w.h.p. based on [\(SM7.9\)](#), [\(SM7.12\)](#) and [\(SM7.13\)](#).

3. ([Bound  \$\tilde{\boldsymbol{\alpha}}\$](#) ) Condition on  $\mathcal{E}_\eta \cap \mathcal{E}_{\|\mathbf{x}\|_2} \cap \mathcal{E}_\tau$ . Use definition  $\tilde{\boldsymbol{\beta}} = \frac{1}{\eta} \widetilde{\mathbf{M}} \mathbf{x}$  from [\(SM7.6\)](#), and properties of  $\widetilde{\mathbf{M}}$  from [\(SM7.7\)](#) we can obtain:

$$\text{(SM7.15)} \quad \begin{cases} \|\tilde{\boldsymbol{\beta}}_{\tau^c}\|_2 \leq \frac{1}{\eta} \left\| \boldsymbol{\nu}_{\tau^c}^* \widetilde{\mathbf{M}} \boldsymbol{\nu}_{\tau} \right\|_2 \|\mathbf{x}\|_2 \leq \frac{\mu \sqrt{p_0 |\tau|}}{\sqrt{p_0\theta/2}} \cdot \sqrt{3p_0\theta} \leq 3\mu \sqrt{p_0 |\tau|} \\ \|\tilde{\boldsymbol{\beta}}_\tau\|_2 \geq \frac{1}{\eta} \left\| \boldsymbol{\nu}_\tau^* \widetilde{\mathbf{M}} \boldsymbol{\nu}_\tau \right\|_2 \|\mathbf{x}\|_2 \geq \frac{\sqrt{1-\mu |\tau|}}{2\sqrt{p_0\theta}} \cdot \sqrt{p_0\theta} \geq 0.45 \end{cases}.$$

Use definition  $\|\tilde{\boldsymbol{\alpha}}\|_2 = \|\boldsymbol{\chi}[\tilde{\boldsymbol{\beta}}]\|_2$ , condition on event

$$\mathcal{E}_\chi := \left\{ \begin{cases} \sigma_i \boldsymbol{\chi}[\boldsymbol{\beta}]_i \geq n\theta \mathcal{S}_{\nu_2\lambda} [|\boldsymbol{\beta}_i|] - \frac{c_\mu^2 n\theta}{p}, & \forall i \in \tau \\ \sigma_i \boldsymbol{\chi}[\boldsymbol{\beta}]_i \leq 4n\theta^2 |\tau| |\boldsymbol{\beta}_i| + \frac{c_\mu n\theta}{p}, & \forall i \in \tau^c \end{cases} \right\},$$

also from [Definition SM2.1](#) we have  $\mu (p\theta)^{1/2} |\boldsymbol{\tau}|^{3/2} < \frac{c_\mu}{4 \log^2 \theta^{-1}}$  and from lemma assumption  $\lambda = \frac{1}{5\sqrt{p\theta}}$ , provides bounds of  $\|\tilde{\boldsymbol{\alpha}}\|_2$  via triangle inequality as:

$$(SM7.16) \quad \begin{cases} \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2 \leq 4n\theta^2 |\boldsymbol{\tau}| \cdot \|\tilde{\boldsymbol{\beta}}_{\boldsymbol{\tau}^c}\|_2 + \frac{c_\mu n\theta}{p} \cdot \sqrt{2p_0} \leq 3c_\mu n\theta \left( \frac{\sqrt{\theta}}{\log^2 \theta^{-1}} + \frac{c_\mu}{p} \right) \\ \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}}\|_2 \geq n\theta \left( \|\tilde{\boldsymbol{\beta}}_{\boldsymbol{\tau}}\|_2 - \nu_2 \lambda \sqrt{|\boldsymbol{\tau}|} - \frac{c_\mu}{p} \sqrt{|\boldsymbol{\tau}|} \right) \geq n\theta \left( 0.45 - \sqrt{\frac{2}{\pi}} \cdot \frac{1}{5} - c_\mu \right) \geq 0.2n\theta \end{cases},$$

since both  $\theta |\boldsymbol{\tau}|$ ,  $\mu p\theta |\boldsymbol{\tau}| < c_\mu$ , we have

$$\begin{cases} \sqrt{1 + \mu p} \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2 \leq 3c_\mu n\theta \sqrt{1 + \mu p} (\sqrt{\theta} + p^{-1}) \leq 6c_\mu n\theta \\ \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2 \leq \frac{6c_\mu^{3/2} n\theta}{\log^2 \theta^{-1}} \min \left\{ \frac{1}{\sqrt{|\boldsymbol{\tau}|}}, \frac{1}{\sqrt{\mu p}}, \frac{1}{\mu p \sqrt{\theta} |\boldsymbol{\tau}|} \right\} \leq 24\sqrt{c_\mu} n\theta \gamma \end{cases},$$

which satisfies [\(SM7.8\)](#), since  $\mu |\boldsymbol{\tau}| < c_\mu < \frac{1}{1000}$ ,

$$(SM7.17) \quad \begin{aligned} (1 + \gamma \sqrt{1 + \mu p}) \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}^c}\|_2 &\leq (24\sqrt{c_\mu} + 6c_\mu) n\theta \gamma \leq 0.1n\theta \gamma \\ &\leq \gamma \sqrt{1 - \mu |\boldsymbol{\tau}|} \|\tilde{\boldsymbol{\alpha}}_{\boldsymbol{\tau}}\|_2. \end{aligned}$$

Finally, given  $p_0\theta > 1000c$ , this result holds with probability at least

$$(SM7.18) \quad \begin{aligned} 1 - \underbrace{\mathbb{P}[\mathcal{E}_\tau^c]}_{\text{Lemma SM1.1}} - \underbrace{\mathbb{P}[\mathcal{E}_{\|\mathbf{x}\|_2}^c]}_{\text{Lemma SM1.2}} - \underbrace{\mathbb{P}[\mathcal{E}_\eta^c]}_{\text{(SM7.14)}} - \underbrace{\mathbb{P}[\mathcal{E}_\chi^c]}_{\text{Corollary SM3.4}} \\ \geq 1 - \frac{2}{p_0\theta} - \frac{1}{n} - 4 \exp\left(\frac{-p_0\theta}{440}\right) \geq 1 - \frac{1}{c} \end{aligned} \quad \blacksquare$$

**SM7.2. Minimization near subspace (Proof of Theorem 5.1)** . Before we start the proof of theorem, writing  $\mathbf{g} = \text{grad}[\varphi_\rho](\mathbf{a})$  and  $\mathbf{H} = \text{Hess}[\varphi_\rho](\mathbf{a})$ , we will first restate the results of [Theorem 4.1](#) in simplified terms. The theorem shows that for any  $\mathbf{a} \in \mathbb{S}^{p-1}$  whose distance to subspace  $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \leq \gamma$ , then at least one of the the following statement hold:

$$(SM7.19) \quad \|\mathbf{g}\|_2 \geq \eta_g$$

$$(SM7.20) \quad \lambda_{\min}(\mathbf{H}) \leq -\eta_v$$

$$(SM7.21) \quad \mathbf{H} \succ \eta_c \cdot \mathbf{P}_{\mathbf{a}^\perp}.$$

Furthermore,  $\varphi_\rho$  is retractive near  $\mathcal{S}_\tau$ : wherever  $d_\alpha(\mathbf{a}, \mathcal{S}_\tau) \geq \frac{\gamma}{2}$ , writing  $\boldsymbol{\alpha}(\mathbf{a})$ ,  $\boldsymbol{\alpha}(\mathbf{g})$  to be the coefficient of  $\mathbf{a}$ ,  $\mathbf{g}$ , we have

$$(SM7.22) \quad \langle \boldsymbol{\alpha}(\mathbf{a})_{\boldsymbol{\tau}^c}, \boldsymbol{\alpha}(\mathbf{g})_{\boldsymbol{\tau}^c} \rangle \geq \eta_r \|\boldsymbol{\alpha}(\mathbf{g})_{\boldsymbol{\tau}^c}\|_2.$$

Also, the the gradient, Hessian and the third order derivative are all bounded as follows:

*Remark SM7.2.* 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 SM7.2](#) 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

$$(SM7.23) \quad \mathbf{a}^{(k+1)} = \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \mathbf{a}^{(k)} - t\mathbf{g}^{(k)} - t^2\mathbf{v}^{(k)} \right]$$

converges toward the minimizer  $\bar{\mathbf{a}}$  in  $\ell^2$ -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  $\mu$ -truncated shift coherent implies that  $\mathbf{a}_0$  is at most  $2\mu$ -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  $\varphi_\rho$  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  $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \leq (1 - t\eta_c)^{-(k-\bar{k})} \|\mathbf{a}^{(\bar{k})} - \bar{\mathbf{a}}\|_2$  [[SM1](#)] (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  $\mathcal{S}_\tau$  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  $\varphi_\rho$  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 [[SM1](#)] (Section 7.1.3)): for any  $\boldsymbol{\xi} \perp \mathbf{a}$ , writing  $\mathbf{a}^+ = \mathbf{P}_{\mathbb{S}^{p-1}}[\mathbf{a} + \boldsymbol{\xi}]$ ,

$$(SM7.24) \quad |\varphi_\rho(\mathbf{a}^+) - (\varphi_\rho(\mathbf{a}) + \langle \text{grad}[\varphi_\rho](\mathbf{a}), \boldsymbol{\xi} \rangle + \frac{1}{2}\boldsymbol{\xi}^*\text{Hess}[\varphi_\rho](\mathbf{a})\boldsymbol{\xi})| \leq \bar{\eta} \|\boldsymbol{\xi}\|_2^3,$$

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

$$(SM7.25) \quad \varphi_\rho(\mathbf{a}^+) > \varphi_\rho(\mathbf{a}) - \frac{1}{2} \left( t \|\mathbf{g}\|_2^2 + \frac{1}{2}t^4\eta_v \|\mathbf{v}\|_2^2 \right).$$

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

$$\begin{aligned} & \langle \mathbf{g}, -t\mathbf{g} - t^2\mathbf{v} \rangle + \frac{1}{2}(t\mathbf{g} + t^2\mathbf{v})^*\mathbf{H}(t\mathbf{g} + t^2\mathbf{v}) \\ & \quad + c \|t\mathbf{g} + t^2\mathbf{v}\|_2^3 > -\frac{1}{2} \left( t \|\mathbf{g}\|_2^2 + \frac{1}{2}t^4\eta_v \|\mathbf{v}\|_2^2 \right) \\ \implies & -\frac{1}{2}t \|\mathbf{g}\|_2^2 + \frac{1}{2}t^2\mathbf{g}^*\mathbf{H}\mathbf{g} + t^3\mathbf{v}^*\mathbf{H}\mathbf{g} \\ & \quad - \frac{1}{4}t^4\eta_v \|\mathbf{v}\|_2^2 + 4\bar{\eta}t^3 \|\mathbf{g}\|_2^3 + 4\bar{\eta}t^6 \|\mathbf{v}\|_2^3 > 0 \\ \implies & -\frac{1}{2}t \|\mathbf{g}\|_2^2 + t^2 \left( \frac{1}{2}\bar{\eta} \|\mathbf{g}\|_2^2 + t\bar{\eta} \|\mathbf{v}\|_2 \|\mathbf{g}\|_2 + 4\bar{\eta}t \|\mathbf{g}\|_2^3 \right) \\ (SM7.26) \quad & -\frac{1}{4}t^4\eta_v \|\mathbf{v}\|_2^2 + 4\bar{\eta}t^6 \|\mathbf{v}\|_2^3 > 0. \end{aligned}$$

If

$$(SM7.27) \quad t < T = \min \left\{ \frac{\|\mathbf{g}\|_2}{\bar{\eta} \|\mathbf{g}\|_2 + 2\bar{\eta}t \|\mathbf{v}\|_2 + 8\bar{\eta}t \|\mathbf{g}\|_2^2}, \sqrt{\frac{\eta_v}{16\bar{\eta} \|\mathbf{v}\|_2}} \right\},$$

then (SM7.26)  $< 0$  contradicting (SM7.25). Using our bounds on  $\|\mathbf{g}\|_2$ ,  $\bar{\eta}$ ,  $\eta_v$  and  $\|\mathbf{v}\|$ , it follows that  $T$  is lower bounded by a polynomial  $\text{poly}(n^{-1}, p^{-1})$ .

2. (Bounds on  $d_\alpha(\mathbf{g}, \mathcal{S}_\tau)$ ,  $d_\alpha(\mathbf{v}, \mathcal{S}_\tau)$ ) We will show there are numerical constants  $c_g, c_v$  such that

$$(SM7.28) \quad d_\alpha(\mathbf{g}, \mathcal{S}_\tau) \leq c_g n \theta \gamma \quad \text{and} \quad d_\alpha(\mathbf{v}, \mathcal{S}_\tau) \leq c_v n \theta p.$$

Define

$$\chi_{\ell^1}[\boldsymbol{\beta}] = \check{\mathbf{C}}_{x_0} \text{prox}_{\lambda \ell^1}[\check{\mathbf{a}} * \mathbf{y}], \quad \chi_\rho[\boldsymbol{\beta}] = \check{\mathbf{C}}_{x_0} \text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}],$$

then the gradient can be written as (SM6.58)

$$(SM7.29) \quad \text{grad}[\varphi_{\ell^1}](\mathbf{a}) = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} (\boldsymbol{\beta}^* \chi_{\ell^1}[\boldsymbol{\beta}] \boldsymbol{\alpha} - \chi_{\ell^1}[\boldsymbol{\beta}]),$$

$$(SM7.30) \quad \text{grad}[\varphi_\rho](\mathbf{a}) = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} (\boldsymbol{\beta}^* \chi_\rho[\boldsymbol{\beta}] \boldsymbol{\alpha} - \chi_\rho[\boldsymbol{\beta}]).$$

Use the following inequalities:

$$\begin{aligned} \frac{1}{2} n \theta &\leq |\boldsymbol{\beta}^* \chi_{\ell^1}[\boldsymbol{\beta}]| \leq \frac{3}{2} n \theta, \\ \|\chi_{\ell^1}[\boldsymbol{\beta}]_{\tau^c}\|_2 &\leq \frac{1}{20} n \theta \gamma, \\ \|\mathbf{I} - \boldsymbol{\alpha} \boldsymbol{\beta}^*\|_2 &\leq 4\sqrt{p}, \\ \|\chi_{\ell^1}[\boldsymbol{\beta}] - \chi_\rho[\boldsymbol{\beta}]\|_2 &\leq n \theta^4, \end{aligned}$$

where the first and second bounds of  $\chi_{\ell^1}[\boldsymbol{\beta}]$  based on event (SM6.59); the third by observing  $\|\boldsymbol{\alpha}\|_2 \leq 2$  and  $\|\boldsymbol{\beta}\|_2 \leq 2 + c_\mu \sqrt{p}$ ; the last from (SM5.21) of Lemma SM5.6 when  $\delta$  is sufficiently small. Hence, by definition of  $d_\alpha(\cdot, \mathcal{S}_\tau)$  (4.16) and knowing  $\mathbf{a}$  is close to subspace  $\|\boldsymbol{\alpha}_{\tau^c}\|_2 \leq \gamma$ , via triangle inequality, we get

$$\begin{aligned} d_\alpha(\mathbf{g}, \mathcal{S}_\tau) &\leq d_\alpha(\text{grad}[\varphi_{\ell^1}](\mathbf{a}), \mathcal{S}_\tau) + d_\alpha(\text{grad}[\varphi_\rho](\mathbf{a}) - \text{grad}[\varphi_{\ell^1}](\mathbf{a}), \mathcal{S}_\tau) \\ &\leq \|\boldsymbol{\beta}^* \chi_{\ell^1}[\boldsymbol{\beta}] \boldsymbol{\alpha}_{\tau^c} - \chi_{\ell^1}[\boldsymbol{\beta}]_{\tau^c}\|_2 + \|(\mathbf{I} - \boldsymbol{\alpha} \boldsymbol{\beta}^*)(\chi_\rho[\boldsymbol{\beta}] - \chi_{\ell^1}[\boldsymbol{\beta}])\|_2 \\ &\leq \frac{3}{2} n \theta \gamma + \frac{1}{20} n \theta \gamma + 4\sqrt{p} n \theta^4 \\ (SM7.31) \quad &\leq 3n \theta \gamma. \end{aligned}$$

To bound the  $d_\alpha$  norm of least eigenvector  $\mathbf{v}$ , note that  $\boldsymbol{\beta}^* \chi_\rho[\boldsymbol{\beta}] > 0$ , we can conclude

$$\mathbf{v}^* \nabla^2 \varphi_\rho(\mathbf{a}) \mathbf{v} \leq \mathbf{v}^* \mathbf{P}_{\mathbf{a}^\perp} \nabla^2 \varphi_\rho(\mathbf{a}) \mathbf{P}_{\mathbf{a}^\perp} \mathbf{v} + \boldsymbol{\beta}^* \chi_\rho[\boldsymbol{\beta}] = \mathbf{v}^* \mathbf{H} \mathbf{v} < -\eta_v,$$

expand  $\nabla^2 \varphi_\rho(\mathbf{a})$  as in (SM5.8), and since  $\mathbf{v}$  is the eigenvector of smallest eigenvalue  $\lambda_{\min} < -\eta_v$ ,

$$(SM7.32) \quad \mathbf{P}_{\mathbf{a}^\perp} \nabla^2 \varphi_\rho(\mathbf{a}) \mathbf{P}_{\mathbf{a}^\perp} \mathbf{v} = (\mathbf{I} - \mathbf{a} \mathbf{a}^*) \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \check{\mathbf{C}}_{x_0} \nabla \text{prox}_{\lambda \rho}[\check{\mathbf{a}} * \mathbf{y}] \check{\mathbf{C}}_{x_0} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{v} = \lambda_{\min} \mathbf{v},$$

hence there exists  $\boldsymbol{\alpha}(\mathbf{v})$  satisfies  $\mathbf{v} = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\alpha}(\mathbf{v})$  and

$$\boldsymbol{\alpha}(\mathbf{v}) = \lambda_{\min}^{-1} \left[ \check{\mathbf{C}}_{\mathbf{x}_0} \nabla_{\text{prox}_{\lambda\rho}} [\check{\mathbf{a}} * \mathbf{y}] \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{v} - \left( \boldsymbol{\beta}^* \check{\mathbf{C}}_{\mathbf{x}_0} \nabla_{\text{prox}_{\lambda\rho}} [\check{\mathbf{a}} * \mathbf{y}] \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{v} \right) \boldsymbol{\alpha} \right].$$

Now since  $\nabla_{\text{prox}_{\lambda\rho}} [\check{\mathbf{a}} * \mathbf{y}]$  is a diagonal matrix with entries in  $[0, 1]$ ,

$$(SM7.33) \quad d_{\alpha}(\mathbf{v}, \mathcal{S}_{\tau}) \leq \|\boldsymbol{\alpha}(\mathbf{v})\|_2 \leq |\lambda_{\min}|^{-1} \|\boldsymbol{\iota} \mathbf{C}_{\mathbf{a}_0}\|_2 \|\mathbf{x}_0\|_1^2 \|\mathbf{v}\|_2 (1 + \|\boldsymbol{\alpha}\|_2 \|\boldsymbol{\beta}\|_2) < c_v n \theta p,$$

where we use upper bound of  $\|\mathbf{x}_0\|_1 < cn\theta$  from [Lemma SM1.2](#) and  $|\lambda_{\min}| > \eta_v > cn\theta\lambda$  from [Corollary SM6.2](#).

3. (Iterates stay within widened subspace). Suppose [\(SM7.22\)](#) holds. We will show that whenever

$$(SM7.34) \quad t \leq T' = \frac{1}{10n\theta},$$

then setting  $\mathbf{a}^+ = \mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}]$ , we have

$$(SM7.35) \quad |d_{\alpha}(\mathbf{a}^+, \mathcal{S}_{\tau}) - d_{\alpha}(\mathbf{a}, \mathcal{S}_{\tau})| \leq \frac{\gamma}{2},$$

and whenever  $d_{\alpha}(\mathbf{a}, \mathcal{S}_{\tau}) \in [\frac{\gamma}{2}, \gamma]$

$$(SM7.36) \quad d_{\alpha}^2(\mathbf{a}^+, \mathcal{S}_{\tau}) \leq d_{\alpha}^2(\mathbf{a}, \mathcal{S}_{\tau}) - t \cdot c' n \theta \gamma^2.$$

If both [\(SM7.35\)](#) and [\(SM7.36\)](#) hold, then all iterates  $\mathbf{a}^{(k)}$  will stay near the subspace:  $d_{\alpha}(\mathbf{a}^{(k)}, \mathcal{S}_{\tau}) < \gamma$ .

To derive [\(SM7.35\)](#), 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 SM2.2](#):

$$(SM7.37) \quad \begin{aligned} d_{\alpha}(\mathbf{a}^+, \mathcal{S}_{\tau}) &= d_{\alpha}(\mathbf{P}_{\mathbb{S}^{p-1}} [\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}], \mathcal{S}_{\tau}) \leq d_{\alpha}(\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}, \mathcal{S}_{\tau}) \\ &\leq d_{\alpha}(\mathbf{a}, \mathcal{S}_{\tau}) + td_{\alpha}(\mathbf{g}, \mathcal{S}_{\tau}) + t^2 d_{\alpha}(\mathbf{v}, \mathcal{S}_{\tau}) \end{aligned}$$

suggests [\(SM7.35\)](#) holds via [\(SM7.28\)](#) and let  $n > Cp^5\theta^{-2}$ , we have

$$(SM7.38) \quad td_{\alpha}(\mathbf{g}, \mathcal{S}_{\tau}) + t^2 d_{\alpha}(\mathbf{v}, \mathcal{S}_{\tau}) \leq \frac{c_g n \theta \gamma}{10n\theta} + \frac{c_v n \theta p}{(10n\theta)^2} < \frac{\gamma}{2}$$

with sufficiently large  $C$ .

To derive [\(SM7.36\)](#), let  $\boldsymbol{\alpha}(\mathbf{a})$  to be a coefficient vector satisfying  $d_{\alpha}(\mathbf{a}, \mathcal{S}_{\tau}) = \|\boldsymbol{\alpha}(\mathbf{a})_{\tau^c}\|_2$ , and based on [\(SM7.30\)](#) and [\(32\)](#), define

$$(SM7.39) \quad \boldsymbol{\alpha}(\mathbf{g}) = \boldsymbol{\beta}^* \boldsymbol{\chi}_{\rho}[\boldsymbol{\beta}] \boldsymbol{\alpha}(\mathbf{a}) - \boldsymbol{\chi}_{\rho}[\boldsymbol{\beta}]$$

$$(SM7.40) \quad \boldsymbol{\alpha}(\mathbf{v}) = \lambda_{\min}^{-1} \check{\mathbf{C}}_{\mathbf{x}_0} \nabla_{\text{prox}_{\lambda\rho}} [\check{\mathbf{a}} * \mathbf{y}] \check{\mathbf{C}}_{\mathbf{x}_0} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota} \mathbf{v}.$$

By the retraction property and norm bounds,

$$(SM7.41) \quad \langle \boldsymbol{\alpha}(\mathbf{a})_{\tau^c}, \boldsymbol{\alpha}(\mathbf{g})_{\tau^c} \rangle \geq \frac{1}{6n\theta} \|\boldsymbol{\alpha}(\mathbf{g})_{\tau^c}\|_2^2$$

$$(SM7.42) \quad \|\boldsymbol{\alpha}(\mathbf{a})_{\tau^c}\|_2 \leq \gamma$$

$$(SM7.43) \quad \|\boldsymbol{\alpha}(\mathbf{v})\|_2 \leq c_v n \theta p.$$

Since  $\|\alpha_{\tau^c}\|_2 > \frac{\gamma}{2}$ ,

$$\begin{aligned}
 \|\mathbf{a}(\mathbf{g})_{\tau^c}\|_2 &\geq \|\beta^* \chi_{\ell^1}[\beta] \alpha_{\tau^c} - \chi_{\ell^1}[\beta]_{\tau^c}\|_2 - \|(\mathbf{I} - \alpha\beta^*)(\chi_{\rho}[\beta] - \chi_{\ell^1}[\beta])\|_2 \\
 &\geq |\beta^* \chi_{\ell^1}[\beta]| \|\alpha_{\tau^c}\|_2 - \|\chi_{\ell^1}[\beta]_{\tau^c}\|_2 - \|(\mathbf{I} - \alpha\beta^*)\|_2 \|(\chi_{\rho}[\beta] - \chi_{\ell^1}[\beta])\|_2 \\
 &\geq \frac{1}{2}n\theta \times \frac{\gamma}{2} - \frac{1}{20}n\theta\gamma + 2n\theta^4 \\
 \text{(SM7.44)} \quad &\geq \frac{1}{10}n\theta\gamma.
 \end{aligned}$$

Finally, we can bound  $d_{\alpha}(\mathbf{a}^+, \mathcal{S}_{\tau})$  as

$$\begin{aligned}
 d_{\alpha}^2(\mathbf{a}^+, \mathcal{S}_{\tau}) &\leq d_{\alpha}^2(\mathbf{a} - t\mathbf{g} - t^2\mathbf{v}, \mathcal{S}_{\tau}) \\
 &\leq \|[\alpha(\mathbf{a}) - t\alpha(\mathbf{g}) - t^2\alpha(\mathbf{v})]_{\tau^c}\|_2^2 \\
 &= \|\alpha(\mathbf{a})_{\tau^c}\|_2^2 - 2t \langle \alpha(\mathbf{a})_{\tau^c}, [\alpha(\mathbf{g}) + t\alpha(\mathbf{v})]_{\tau^c} \rangle + t^2 \|[\alpha(\mathbf{g}) + t\alpha(\mathbf{v})]_{\tau^c}\|_2^2 \\
 &\leq \|\alpha(\mathbf{a})_{\tau^c}\|_2^2 - 2t \langle \alpha(\mathbf{a})_{\tau^c}, \alpha(\mathbf{g})_{\tau^c} \rangle + 2t^2 \|\alpha(\mathbf{a})_{\tau^c}\|_2 \|\alpha(\mathbf{v})\|_2 \\
 &\quad + 2t^2 \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 + 2t^4 \|\alpha(\mathbf{v})\|_2^2 \\
 &\leq d^2(\mathbf{a}, \mathcal{S}_{\tau}) - 2t \left[ \left( \frac{1}{3n\theta} - t \right) \|\alpha(\mathbf{g})_{\tau^c}\|_2^2 - tn\theta p\gamma - t^3(c_v n\theta p)^2 \right] \\
 \text{(SM7.45)} \quad &\leq d^2(\mathbf{a}, \mathcal{S}_{\tau}) - t \cdot c' n\theta\gamma^2
 \end{aligned}$$

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  $\gamma$  neighborhood of  $\mathcal{S}_{\tau}$  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.  $\blacksquare$

**SM8. 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

$$\text{(SM8.1)} \quad \lambda^{(0)} \leftarrow 5\kappa_I \tilde{\mu} \quad \text{and} \quad I^{(0)} \leftarrow \text{supp}(\mathcal{S}_{\lambda}[\mathbf{C}_{\mathbf{a}^{(0)}}^* \mathbf{y}]),$$

where as each iteration of the algorithm consists of the following key steps:

- **Sparse Estimation using Reweighted Lasso:** Set

$$\text{(SM8.2)} \quad \mathbf{x}^{(k+1)} \leftarrow \underset{\mathbf{x}}{\text{argmin}} \frac{1}{2} \|\mathbf{a}^{(k)} * \mathbf{x} - \mathbf{y}\|_2^2 + \sum_{i \notin I^{(k)}} \lambda^{(k)} |\mathbf{x}_i|;$$

- **Kernel Estimation using Least Squares:** Set

$$\text{(SM8.3)} \quad \mathbf{a}^{(k+1)} \leftarrow \mathbf{P}_{\mathbb{S}^{p-1}} \left[ \underset{\mathbf{a}}{\text{argmin}} \frac{1}{2} \|\mathbf{a} * \mathbf{x}^{(k+1)} - \mathbf{y}\|_2^2 \right];$$

- **Continuation and reweighting by decreasing sparsity regularizer:** Set

$$(SM8.4) \quad \lambda^{(k+1)} \leftarrow \frac{1}{2}\lambda^{(k)} \quad \text{and} \quad I^{(k+1)} \leftarrow \text{supp}(\mathbf{x}^{(k+1)}).$$

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.

Below, we will write

$$(SM8.5) \quad \tilde{\mu} = \max \{ \mu, p^{-1} \}.$$

Our refinement algorithm will demand an initialization satisfying

$$(SM8.6) \quad \|\mathbf{a}^{(0)} - \mathbf{a}_0\|_2 \leq \tilde{\mu}.$$

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

$$(SM8.7) \quad I = \text{supp}(\mathbf{x}_0), \quad I^{(k)} = \text{supp}(\mathbf{x}^{(k)}).$$

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  $\kappa_I$  as an upper bound for number of nonzero entries of  $\mathbf{x}_0$  in a length- $p$  window:

$$(SM8.8) \quad \kappa_I = 6 \max \{ \theta p, \log n \},$$

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

$$(SM8.9) \quad \max_{\ell} |I \cap ([p] + \ell)| \leq \kappa_I.$$

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  $\ell^2$  norm of  $\mathbf{x}_0$  over any length- $p$  window:

$$(SM8.10) \quad \|\mathbf{x}_0\|_{\square} := \max_{\ell} \|\mathbf{P}_{([p] + \ell)} \mathbf{x}_0\|_2.$$

Below, we will sometimes work with the  $\square$ -induced operator norm:

$$(SM8.11) \quad \|\mathbf{M}\|_{\square \rightarrow \square} = \sup_{\|\mathbf{x}\|_{\square} \leq 1} \|\mathbf{M}\mathbf{x}\|_{\square}$$

For now, we note that in the Bernoulli-Gaussian model,  $\|\mathbf{x}_0\|_{\square}$  is typically not large

$$(SM8.12) \quad \|\mathbf{x}_0\|_{\square} \leq \sqrt{\kappa_I}.$$

**SM8.1. 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 [SM3], 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

$$(SM8.13) \quad \sum_{i \in T^c} |\mathbf{x}_i|$$

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 SM8.1 (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*

$$(SM8.14) \quad \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{a} * \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i \in T^c} |\mathbf{x}_i| \right\},$$

with

$$(SM8.15) \quad \lambda > 5\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2,$$

is unique with the form

$$(SM8.16) \quad \mathbf{x}^+ = \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a},J}^* \mathbf{C}_{\mathbf{a},J})^{-1} \boldsymbol{\iota}_J^* (\mathbf{C}_{\mathbf{a}}^* \mathbf{y} - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma})$$

where  $\boldsymbol{\sigma} = \text{sign}(\mathbf{x}^+)$ . Its support set  $J$  satisfies

$$(SM8.17) \quad (T \cup I_{\geq 3\lambda}) \subseteq J \subseteq I$$

and its entrywise error is bounded as

$$(SM8.18) \quad \|\mathbf{x}^+ - \mathbf{x}_0\|_{\infty} \leq 3\lambda.$$

Above,  $c_{\mu} > 0$  is a positive numerical constant.

We prove [Lemma SM8.1](#) 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 [Subsection SM8.4](#).

*Proof.* 1. (Restricted support Lasso problem). We first consider the restricted problem

$$(SM8.19) \quad \min_{\mathbf{w} \in \mathbb{R}^{|I|}} \left\{ \frac{1}{2} \|\mathbf{a} * \boldsymbol{\iota}_I \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{i \in T^c} |(\boldsymbol{\iota}_I \mathbf{w})_i| \right\}.$$

Under our assumptions, provided  $c < \frac{1}{9}$ , [Lemma SM8.6](#) implies that

$$(SM8.20) \quad \iota_I^* \mathbf{C}_a^* \mathbf{C}_a \iota_I \mathbf{w}_* = [\mathbf{C}_a^* \mathbf{C}_a]_{I,I} \succ \mathbf{0},$$

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

$$(SM8.21) \quad \iota_I^* \mathbf{C}_a^* \mathbf{C}_a \iota_I \mathbf{w}_* \in \iota_I^* \mathbf{C}_a^* \mathbf{y} - \lambda \partial \| \mathbf{P}_{T^c} [\cdot] \|_1 (\mathbf{w}_*).$$

Writing  $J = \text{supp}(\iota_I \mathbf{w}_*) \subseteq I$ ,  $\mathbf{C}_{a,J} = \mathbf{C}_a \iota_J$ ,  $\mathbf{w}_J = \iota_J^* \iota_I \mathbf{w}_*$  the corresponding sub-vector containing the nonzero entries of  $\mathbf{w}_*$  and  $\boldsymbol{\sigma}_{J \setminus T} = \iota_J^* \mathbf{P}_{T^c} [\text{sign}(\iota_I \mathbf{w}_*)]$ , the condition [\(SM8.21\)](#) is satisfied if and only if

$$(SM8.22) \quad \mathbf{C}_{a,J}^* \mathbf{C}_{a,J} \mathbf{w}_J = \mathbf{C}_{a,J}^* \mathbf{y} - \lambda \boldsymbol{\sigma}_{J \setminus T},$$

$$(SM8.23) \quad \| \mathbf{C}_{a, I \setminus J}^* (\mathbf{C}_{a,J} \mathbf{w}_J - \mathbf{y}) \|_\infty \leq \lambda.$$

We will argue that under our assumptions,  $J$  necessarily contains all of the large entries of  $\mathbf{x}_0$ :

$$(SM8.24) \quad I_{>3\lambda} = \{\ell \in I \mid |\mathbf{x}_{0\ell}| > 3\lambda\} \subseteq J.$$

We show this by contradiction – namely, if some large entry of  $\mathbf{x}_0$  is not in  $J$ , then the dual condition [\(SM8.23\)](#) is violated, contradicting the optimality of  $\mathbf{w}_*$ . To this end, note that by [Corollary SM8.7](#),  $\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}$  has full rank. From [\(SM8.22\)](#),

$$(SM8.25) \quad \mathbf{w}_J = [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} [\mathbf{C}_{a,J}^* \mathbf{y} - \lambda \boldsymbol{\sigma}_{J \setminus T}].$$

Write  $\mathbf{x}_{0J} = \iota_J^* \mathbf{x}_0$  and  $(\mathbf{x}_0)_{I \setminus J} = \mathbf{P}_{I \setminus J} \mathbf{x}_0$ . We can further notice that

$$\begin{aligned} \mathbf{C}_{a,J} \mathbf{w}_J - \mathbf{y} &= \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{y} - \lambda \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \\ &= \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{C}_{a_0 J} \mathbf{x}_{0J} \\ &\quad + \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{C}_{a_0 I \setminus J} (\mathbf{x}_0)_{I \setminus J} \\ &\quad - \lambda \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \\ &= \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{C}_{a_0 - a J} \mathbf{x}_{0J} \\ &\quad + \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{C}_{a_0 I \setminus J} (\mathbf{x}_0)_{I \setminus J} \\ (SM8.26) \quad &\quad - \lambda \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \boldsymbol{\sigma}_{J \setminus T}, \end{aligned}$$

where in the final line we have used that

$$(SM8.27) \quad \left( \mathbf{C}_{a,J} [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* - \mathbf{I} \right) \mathbf{C}_{a,J} = \mathbf{0}.$$

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 [\(SM8.26\)](#), [Corollary SM8.7](#) and

Lemma SM8.8, and simplify the induced norms  $\|\cdot\|_{\infty \rightarrow \infty}$  and  $\|\cdot\|_{\square \rightarrow \square}$  as  $\|\cdot\|_{\infty}$  and  $\|\cdot\|_{\square}$ , we have

$$\begin{aligned}
 (SM8.28) \quad & -\xi s_i[\mathbf{a}]^* (\mathbf{C}_{aJ} \mathbf{w}_J - \mathbf{y}) = \xi s_i[\mathbf{a}]^* \left( \mathbf{I} - \mathbf{C}_{aJ} [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \mathbf{C}_{aJ}^* \right) s_i[\mathbf{a}_0] \mathbf{x}_{0i} \\
 & \quad + \xi 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}(\mathbf{x}_0)_{I \setminus (J \cup \{i\})} \\
 & \quad + \xi 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 - aJ} \mathbf{x}_{0J} \\
 & \quad + \xi \lambda s_i[\mathbf{a}]^* \mathbf{C}_{aJ} [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \\
 & \geq \left( \langle s_i[\mathbf{a}], s_i[\mathbf{a}_0] \rangle \right. \\
 & \quad - \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty} \| \mathbf{C}_{aJ}^* s_i[\mathbf{a}_0] \|_{\infty} \left. \right) \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \\
 & \quad - \left( \left\| s_i[\mathbf{a}]^* \mathbf{C}_{a_0 I \setminus \{i\}} \right\|_1 \right. \\
 & \quad \left. + \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty} \left\| \mathbf{C}_{aJ}^* \mathbf{C}_{a_0 I \setminus J} \right\|_{\infty} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \\
 & \quad - \left( \|s_i[\mathbf{a}]^* \mathbf{C}_{a_0 - aJ}\|_2 \right. \\
 & \quad \left. + \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_2 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\square} \| \mathbf{C}_{aJ}^* \mathbf{C}_{a_0 - aJ} \|_{\square} \right) \sqrt{2} \| \mathbf{x}_0 \|_{\square} \\
 (SM8.29) \quad & - \lambda \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\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} \\
 & \quad - 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} \\
 & \quad - \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} \\
 (SM8.30) \quad & - \lambda C_4 \kappa_I \tilde{\mu} \\
 & \geq \left( 1 - C'_1 \kappa_I \tilde{\mu} - C_2 (\kappa_I \tilde{\mu})^2 \right) \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \\
 (SM8.31) \quad & - 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 \\
 (SM8.32) \quad & \geq \frac{1}{2} \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} - \lambda/2,
 \end{aligned}$$

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  $\lambda$ , 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

$$(SM8.33) \quad \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \leq 3\lambda.$$

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  $\mathbf{x}^+ = \boldsymbol{\iota}_J \mathbf{w}_J$ , is the unique optimal solution to the *full* Lasso problem

$$(SM8.34) \quad \min_{\mathbf{x}} \varphi_{\text{lasso}}(\mathbf{x}) \equiv \frac{1}{2} \| \mathbf{a} * \mathbf{x} - \mathbf{y} \|_2^2 + \lambda \sum_{i \in T^c} |x_i|.$$

To prove that  $\mathbf{x}^+$  is the unique optimal solution, it suffices to show that for every  $i \in I^c$ ,

$$(SM8.35) \quad |s_i[\mathbf{a}]^*(\mathbf{a} * \mathbf{x}^+ - \mathbf{y})| < \lambda.$$

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,

$$(SM8.36) \quad \mathbf{0} \in \mathbf{P}_I \partial_{\mathbf{x}} \varphi_{\text{lasso}}(\mathbf{x})$$

Combining with (SM8.35), 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,

$$(SM8.37) \quad \begin{aligned} \varphi_{\text{lasso}}(\mathbf{x}') &\geq \varphi_{\text{lasso}}(\mathbf{x}^+) + \langle \varepsilon \boldsymbol{\zeta}, \mathbf{x}' - \mathbf{x}^+ \rangle \\ &\geq \varphi_{\text{lasso}}(\mathbf{x}^+) + \varepsilon \|\mathbf{x}'_{I^c}\|_1, \end{aligned}$$

which is strictly larger than  $\varphi_{\text{lasso}}(\mathbf{x}^+)$ . Hence, when (SM8.35) 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 (SM8.34) is unique and equal to  $\mathbf{x}^+$ .

We finish by showing (SM8.35). Using the same expansion as above, we obtain

$$(SM8.38) \quad \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| \\ &\quad + \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 - aJ} \mathbf{x}_{0J} \right| \\ &\quad + \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 \right. \\ &\quad \left. + \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty} \left\| \mathbf{C}_{aJ}^* \mathbf{C}_{a_0 I \setminus J} \right\|_{\infty} \right) \|(\mathbf{x}_0)_{I \setminus J}\|_{\infty} \\ &\quad + \left( \|s_i[\mathbf{a}]^* \mathbf{C}_{a_0 - aJ}\|_2 \right. \\ &\quad \left. + \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_2 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\square} \left\| \mathbf{C}_{aJ}^* \mathbf{C}_{a_0 - aJ} \right\|_{\square} \right) \sqrt{2} \|\mathbf{x}_0\|_{\square} \\ &\quad + \lambda \|s_i[\mathbf{a}]^* \mathbf{C}_{aJ}\|_1 \left\| [\mathbf{C}_{aJ}^* \mathbf{C}_{aJ}]^{-1} \right\|_{\infty} \|\boldsymbol{\sigma}_{J \setminus T}\|_{\infty} \\ (SM8.39) \quad &\leq C_1 (\tilde{\mu} \kappa_I + \tilde{\mu} \kappa_I \times 1 \times \tilde{\mu} \kappa_I) \times 2\lambda \\ &\quad + (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} \\ (SM8.40) \quad &\quad + \lambda C_3 \times \tilde{\mu} \kappa_I \\ (SM8.41) \quad &\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 \\ (SM8.42) \quad &< \lambda, \end{aligned}$$

where the last line holds as long as  $c_{\mu}$  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 [Corollary SM8.7](#), [Lemma SM8.8](#),

$$\begin{aligned}
 \|\mathbf{x}_J^+ - (\mathbf{x}_0)_J\|_\infty &= \left\| [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* \mathbf{C}_{a_0} \mathbf{x}_0 - \lambda [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} - (\mathbf{x}_0)_J \right\|_\infty \\
 &\leq \left\| [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* \mathbf{C}_{a_0 - a_J} (\mathbf{x}_0)_J \right\|_\infty + \lambda \left\| [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \boldsymbol{\sigma}_{J \setminus T} \right\|_\infty \\
 &\quad + \left\| [\mathbf{C}_{a,J}^* \mathbf{C}_{a,J}]^{-1} \mathbf{C}_{a,J}^* \mathbf{C}_{a_{I \setminus J}} (\mathbf{x}_0)_{I \setminus J} \right\|_\infty \\
 &\leq 2 \|\mathbf{C}_{a,J}^* \mathbf{C}_{a_0 - a_J}\|_{\square \rightarrow \infty} \|(\mathbf{x}_0)_J\|_\square + \\
 &\quad 2\lambda + 2 \|\mathbf{C}_{a,J}^* \mathbf{C}_{a_{I \setminus J}}\|_\infty \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \\
 &\leq 2\sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2 \|\mathbf{x}_0\|_\square + 2\lambda + 2 \times 3\tilde{\mu} \times 2\kappa_{I \setminus J} \times 3\lambda \\
 &\leq 3\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2 + 2\lambda + 36\lambda\tilde{\mu}\kappa_I \\
 \text{(SM8.43)} \quad &\leq 3\lambda,
 \end{aligned}$$

establishing the claim. ■

**SM8.2. Least squares solution  $\mathbf{a}^{(k)}$  contracts.** 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

$$\text{(SM8.44)} \quad \mathbf{a}^+ \leftarrow \operatorname{argmin}_{\mathbf{a}' \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{a}' * \mathbf{x} - \mathbf{y}\|_2^2$$

is closer to  $\mathbf{a}_0$  compared to  $\mathbf{a}$ . Observe that in [Lemma SM8.1](#), the solution of [\(SM8.16\)](#)

$$\text{(SM8.45)} \quad \mathbf{x} = \iota_J (\mathbf{C}_{a,J}^* \mathbf{C}_{a,J})^{-1} \iota_J^* (\mathbf{C}_a^* \mathbf{C}_{a_0} \mathbf{x}_0 - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}),$$

by assuming  $\mathbf{C}_{a,J}^* \mathbf{C}_{a,J} \approx \mathbf{I}$ ,  $\mathbf{a} \approx \mathbf{a}_0$  and  $J \setminus T \approx \emptyset$ , is a good approximation to the true sparse map  $\mathbf{x}_0$

$$\text{(SM8.46)} \quad \mathbf{x} \approx \mathbf{I} (\mathbf{x}_0 - \mathbf{0}) = \mathbf{x}_0;$$

furthermore, its difference to the true sparse map  $\|\mathbf{x}_0 - \mathbf{x}\|_2$  is proportional to  $\|\mathbf{a}_0 - \mathbf{a}\|_2$  as

$$\text{(SM8.47)} \quad \mathbf{x} - \mathbf{x}_0 \approx \mathbf{P}_I (\mathbf{C}_a^* \mathbf{C}_{a_0} \mathbf{x}_0 - \mathbf{C}_a^* \mathbf{C}_a \mathbf{x}_0) \approx \mathbf{P}_I [\mathbf{C}_{a_0}^* \mathbf{C}_{x_0} \iota (\mathbf{a}_0 - \mathbf{a})].$$

To this end, since we know the solution of least square problem  $\mathbf{a}^+$  is simply

$$\text{(SM8.48)} \quad \mathbf{a}^+ = (\iota^* \mathbf{C}_x^* \mathbf{C}_x \iota)^{-1} (\iota^* \mathbf{C}_x^* \mathbf{C}_{x_0} \iota \mathbf{a}_0),$$

this implies the difference between the new  $\mathbf{a}^+$  and  $\mathbf{a}_0$ , has the relationship with  $\mathbf{a} - \mathbf{a}_0$  roughly

$$\begin{aligned}
 \mathbf{a}^+ - \mathbf{a}_0 &= (\iota^* \mathbf{C}_x^* \mathbf{C}_x \iota)^{-1} (\iota^* \mathbf{C}_x^* \mathbf{C}_{x_0} \iota \mathbf{a}_0 - \iota^* \mathbf{C}_x^* \mathbf{C}_x \iota \mathbf{a}_0) \\
 &\approx (n\theta)^{-1} \iota^* \mathbf{C}_{x_0}^* \mathbf{C}_{a_0} (\mathbf{x}_0 - \mathbf{x}) \\
 \text{(SM8.49)} \quad &\approx (n\theta)^{-1} \iota^* \mathbf{C}_{x_0}^* \mathbf{C}_{a_0} \mathbf{P}_I \mathbf{C}_{a_0}^* \mathbf{C}_{x_0} \iota (\mathbf{a} - \mathbf{a}_0).
 \end{aligned}$$

To make this point precise, we introduce the following lemma:

**Lemma SM8.2 (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  $C, C', c, c', c_\mu$  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*

$$(SM8.50) \quad \mathbf{x} = \iota_J (\mathbf{C}_{\mathbf{a}J}^* \mathbf{C}_{\mathbf{a}J})^{-1} \iota_J^* (\mathbf{C}_{\mathbf{a}J}^* \mathbf{y} - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma})$$

where the set  $J, T$  satisfies  $I_{>6\lambda} \subseteq T \subseteq J \subseteq I$ , we have

$$(SM8.51) \quad \begin{aligned} & \frac{1}{n\theta} \left\| \iota^* \mathbf{C}_{\mathbf{x}}^* \mathbf{C}_{\mathbf{x} - \mathbf{x}_0} \iota \mathbf{a}_0 - \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \iota (\mathbf{a}_0 - \mathbf{a}) \right\|_2 \\ & \leq C' \lambda \left( \tilde{\lambda} + \tilde{\mu} \kappa_I \right) + \frac{1}{32} \|\mathbf{a} - \mathbf{a}_0\|_2 \end{aligned}$$

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 SM8.1](#) and our assumptions on the set  $T$ , then  $\mathbf{x}$  approximates  $\mathbf{x}_0$  in the sense that

$$(SM8.52) \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq 3\lambda$$

$$(SM8.53) \quad \|(\mathbf{x}_0)_{I \setminus J}\|_\infty \leq 3\lambda$$

$$(SM8.54) \quad \|(\mathbf{x}_0)_{I \setminus T}\|_\infty \leq 6\lambda.$$

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 [\(SM8.53\)](#) 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 SM1.1](#) implies that with probability at least  $1 - 2/n$ :

$$(SM8.55) \quad |I \setminus J| \leq 3\lambda n \theta + 6\sqrt{\lambda n \theta} \log n \leq 3\tilde{\lambda} n \theta$$

$$(SM8.56) \quad |I \setminus T| \leq 6\lambda n \theta + 12\sqrt{\lambda n \theta} \log n \leq 6\tilde{\lambda} n \theta,$$

and

$$(SM8.57) \quad |(I \setminus J) \cap s_\ell[I]| \leq 3\lambda n \theta^2 + 6\sqrt{\lambda n \theta^2} \log n \leq 3\tilde{\lambda} n \theta^2;$$

together with base on properties of Bernoulli-Gaussian vector  $\mathbf{x}_0$  from [Section SM1](#) and we

conclude with probability at least  $1 - c/n$ , all the following events hold:

$$(SM8.58) \quad \frac{1}{2}n\theta \leq |I| \leq 2n\theta,$$

$$(SM8.59) \quad \max_{\ell \neq 0} |I \cap s_\ell[I]| \leq 2n\theta^2$$

$$(SM8.60) \quad \max_{\ell \neq 0} |(I \setminus J) \cap s_\ell[I]| \leq 6\tilde{\lambda}n\theta^2,$$

$$(SM8.61) \quad \|\mathbf{x}_0\|_{\square}^2 \leq \kappa_I,$$

$$(SM8.62) \quad \|\check{\mathbf{a}}_0 * \mathbf{x}_0\|_{\square}^2 \leq \kappa_I,$$

$$(SM8.63) \quad \|\mathbf{x}_0\|_2^2 \leq 2n\theta,$$

$$(SM8.64) \quad \|\mathbf{x}_0\|_1 \leq 2n\theta,$$

$$(SM8.65) \quad \max_{\ell \neq 0} \|\mathbf{P}_{I \cap s_\ell[I]} \mathbf{x}_0\|_2^2 \leq 2n\theta^2,$$

$$(SM8.66) \quad \max_{\ell \neq 0} \|\mathbf{P}_{I \cap s_\ell[I \setminus J]} \mathbf{x}_0\|_1 \leq 12\tilde{\lambda}n\theta^2,$$

$$(SM8.67) \quad \|\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}\|_2^2 \leq 3n\theta,$$

provided by  $n \geq C\theta^{-2} \log p$  for sufficiently large constant  $C$ .

1. (Approximate  $\mathbf{C}_x$  with  $\mathbf{C}_{\mathbf{x}_0}$ ) Since

$$(SM8.68) \quad \boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 = \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 + \boldsymbol{\iota}^* \mathbf{C}_{x-\mathbf{x}_0}^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0$$

where

$$(SM8.69) \quad \begin{aligned} \|\boldsymbol{\iota}^* \mathbf{C}_{x-\mathbf{x}_0}^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0\|_2 &\leq \|\mathbf{a}_0\|_2 \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \|\mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}\|_2 \sqrt{2p} \max_{\ell \neq 0} |\langle s_\ell[\mathbf{x} - \mathbf{x}_0], \mathbf{x} - \mathbf{x}_0 \rangle| \\ &\leq \|\mathbf{x} - \mathbf{x}_0\|_\infty^2 \times |I| + \sqrt{2\tilde{\mu}p^2} \left( \|\mathbf{x} - \mathbf{x}_0\|_\infty^2 \times \max_{\ell \neq 0} |I \cap s_\ell[I]| \right) \\ &\leq C_1 \left( \lambda^2 n\theta + \sqrt{2\tilde{\mu}p^2} (\lambda^2 n\theta^2) \right) \\ &\leq 2C_1 \lambda^2 n\theta, \end{aligned}$$

we have that

$$(SM8.70) \quad \|\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0\|_2 \leq 2C_1 \lambda^2 n\theta.$$

2. (Extract the  $\mathbf{a}_0 - \mathbf{a}$  term) Observe that

$$(SM8.71) \quad \begin{aligned} \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{x-\mathbf{x}_0} \boldsymbol{\iota} \mathbf{a}_0 &= \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{x} - \mathbf{x}_0) \\ &= \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \left( \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} \boldsymbol{\iota}_J^* (\mathbf{C}_{\mathbf{a}}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0 - \lambda \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}) \right. \\ &\quad \left. - \boldsymbol{\iota}_J (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J}) (\mathbf{x}_0)_J - \mathbf{P}_{I \setminus J} \mathbf{x}_0 \right) \\ &= \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0 J} (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} \mathbf{C}_{\mathbf{a}_J}^* (\mathbf{C}_{\mathbf{a}_0 - \mathbf{a}} \mathbf{x}_0) \\ &\quad + \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0 J} (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} \mathbf{C}_{\mathbf{a}_J}^* (\mathbf{C}_{\mathbf{a}} \mathbf{x}_0 - \mathbf{C}_{\mathbf{a}_J} (\mathbf{x}_0)_J) \\ &\quad - \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_{I \setminus J} \mathbf{x}_0 \\ &\quad - \lambda \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0 J} (\mathbf{C}_{\mathbf{a}_J}^* \mathbf{C}_{\mathbf{a}_J})^{-1} \boldsymbol{\iota}_J^* \mathbf{P}_{J \setminus T} \boldsymbol{\sigma}, \end{aligned}$$

where, the second term in (SM8.71) is bounded as

$$\begin{aligned}
& \left\| \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0 J} (\mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} J})^{-1} \mathbf{C}_{\mathbf{a} J}^* (\mathbf{C}_{\mathbf{a}} \mathbf{x}_0 - \mathbf{C}_{\mathbf{a} J} (\mathbf{x}_0)_J) \right\|_2 \\
& \leq \left\| \mathbf{C}_{\mathbf{x}_0} \iota \right\|_2 \times \left\| \mathbf{C}_{\mathbf{a}_0 J} \right\|_2 \left\| (\mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} J})^{-1} \right\|_2 \\
& \quad \times \left\| \mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} I \setminus J} \right\|_2 \times \left\| (\mathbf{x}_0)_{I \setminus J} \right\|_2 \\
& \leq C_2 \left( \sqrt{n\theta} \times 3 \times \tilde{\mu} \kappa_I \times \lambda \sqrt{\tilde{\lambda} n \theta} \right) \\
\text{(SM8.72)} \quad & \leq 3C_2 \tilde{\mu} \kappa_I \lambda n \theta;
\end{aligned}$$

the third term in (SM8.71) is bounded as

$$\begin{aligned}
& \left\| \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_{I \setminus J} \mathbf{x}_0 \right\|_2 = \left\| \iota^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{P}_{[\pm p] \setminus 0} + \mathbf{e}_0 \mathbf{e}_0^*) \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_{I \setminus J} \mathbf{x}_0 \right\|_2 \\
& \leq \left\| \mathbf{a}_0 \right\|_2 \left\| (\mathbf{x}_0)_{I \setminus J} \right\|_2^2 \\
& \quad + \left\| \mathbf{C}_{\mathbf{a}_0} \iota \right\|_2 \times \sqrt{2p} \times \max_{\ell \neq 0} \left\| \mathbf{P}_{I \cap s_\ell [I \setminus J]} \mathbf{x}_0 \right\|_1 \times \left\| (\mathbf{x}_0)_{I \setminus J} \right\|_\infty \\
& \leq C_3 \left( \lambda^2 \times \tilde{\lambda} n \theta + \sqrt{\tilde{\mu} p^2} \times \tilde{\lambda} n \theta^2 \times \lambda \right) \\
\text{(SM8.73)} \quad & \leq 2C_3 \tilde{\lambda} \lambda n \theta;
\end{aligned}$$

and finally, write  $\mathbf{\Delta} = (\mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} J})^{-1} - \mathbf{I}$ , then the fourth term in (SM8.71) is bounded as

$$\begin{aligned}
& \lambda \left\| \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \iota_J (\mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} J})^{-1} \iota_J^* \mathbf{P}_{J \setminus T} \boldsymbol{\sigma} \right\|_2 \\
& = \lambda \left\| \iota^* \mathbf{C}_{\mathbf{a}_0} (\mathbf{P}_{[\pm p] \setminus 0} + \mathbf{e}_0 \mathbf{e}_0^*) \mathbf{C}_{\mathbf{x}_0}^* \iota_J (\mathbf{I} + \mathbf{\Delta}) \iota_J^* \mathbf{P}_{J \setminus T} \boldsymbol{\sigma} \right\|_2 \\
& \leq \lambda \left\| \mathbf{C}_{\mathbf{a}_0}^* \iota \right\|_2 \sqrt{2p} \max_{\ell \neq 0} \left\| \mathbf{P}_{I \cap s_\ell [I \setminus T]} \mathbf{x}_0 \right\|_1 + \lambda \left\| \mathbf{a}_0 \right\|_2 \left\| \mathbf{P}_{I \setminus T} \mathbf{x}_0 \right\|_1 \\
& \quad + \lambda \left\| \mathbf{C}_{\mathbf{a}_0}^* \iota \right\|_2 \sqrt{2p} \left\| \mathbf{P}_{I \cap s_\ell [I]} \mathbf{x}_0 \right\|_1 \left\| \mathbf{\Delta} \right\|_{\infty \rightarrow \infty} \\
& \quad + \lambda \left\| \mathbf{a}_0 \right\|_2 \left\| \mathbf{x}_0 \right\|_2 \left\| \mathbf{\Delta} \right\|_2 \sqrt{|J \setminus T|} \\
& \leq C_4 \lambda \left( \sqrt{\tilde{\mu} p^2} \times \tilde{\lambda} n \theta^2 + \lambda \tilde{\lambda} n \theta \right. \\
& \quad \left. + \sqrt{\tilde{\mu} p^2} \times n \theta^2 \times \tilde{\mu} \kappa_I + \sqrt{n\theta} \times \tilde{\mu} \kappa_I \sqrt{\tilde{\lambda} n \theta} \right) \\
\text{(SM8.74)} \quad & \leq 2C_4 \left( \tilde{\lambda} + \tilde{\mu} \kappa_I \right) \lambda n \theta.
\end{aligned}$$

Therefore, combining (SM8.72)-(SM8.74) we obtain

$$\begin{aligned}
& \left\| \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{x} - \mathbf{x}_0} \iota \mathbf{a}_0 - \iota^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0 J} (\mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a} J})^{-1} \mathbf{C}_{\mathbf{a} J}^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}} \mathbf{x}_0 \right\|_2 \\
\text{(SM8.75)} \quad & \leq C_5 \left( \tilde{\lambda} + \tilde{\mu} \kappa_I \right) \lambda n \theta.
\end{aligned}$$

3. (Extract the set  $J$ ) Lastly, we will further simplify the term with  $\mathbf{a} - \mathbf{a}_0$  in (SM8.75) by

extracting the set  $J$ :

$$\begin{aligned}
 & \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} (C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 J})^{-1} C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 - \mathbf{a}} \mathbf{x}_0 \\
 &= \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} (\mathbf{I} + \Delta) C_{\mathbf{a}_0 + (\mathbf{a} - \mathbf{a}_0) J}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}) \\
 &= \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0} P_I C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}) \\
 &\quad + \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} \Delta C_{\mathbf{a}_0 J}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}) \\
 &\quad + \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} (C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 J})^{-1} C_{\mathbf{a}_0 - \mathbf{a} J}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}) \\
 &\quad - \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0} P_{I \setminus J} C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}),
 \end{aligned} \tag{SM8.76}$$

where, the latter terms in (SM8.76) are bounded as

$$\begin{aligned}
 & \left\| \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} \Delta C_{\mathbf{a}_0 J}^* C_{\mathbf{x}_0} \iota \right\|_2 \leq \|C_{\mathbf{x}_0} \iota\|_2^2 \|C_{\mathbf{a}_0 J}\|_2^2 \|\Delta\|_2 \leq C_6 \tilde{\mu} \kappa_I n \theta \\
 & \left\| \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} (C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 J})^{-1} C_{\mathbf{a}_0 - \mathbf{a} J}^* C_{\mathbf{x}_0} \iota \right\|_2 \\
 & \leq \|C_{\mathbf{x}_0} \iota\|_2^2 \|C_{\mathbf{a}_0 J}\|_2 \left\| (C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 J})^{-1} \right\|_2 \|C_{\mathbf{a}_0 - \mathbf{a} J}\|_2 \leq C_7 \tilde{\mu} \sqrt{\kappa_I} n \theta \\
 & \left\| P_{I \setminus J} C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota \right\|_2^2 \leq |I \setminus J| \|\check{\mathbf{a}}_0 * \mathbf{x}_0\|_{\square}^2 \\
 & \leq C_8 \tilde{\lambda} n \theta \times \kappa_I \leq C_8 \left( \lambda \kappa_I + \frac{\kappa_I \log n}{\sqrt{n \theta^2}} \right) n \theta,
 \end{aligned} \tag{SM8.77}$$

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:

$$\begin{aligned}
 & \left\| \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0 J} (C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 J})^{-1} C_{\mathbf{a}_0 J}^* C_{\mathbf{a}_0 - \mathbf{a}} \mathbf{x}_0 \right. \\
 & \quad \left. - \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0} P_I C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota(\mathbf{a}_0 - \mathbf{a}) \right\|_2 \\
 & \leq \left( \frac{3}{100} + \frac{1}{1000} \right) n \theta \|\mathbf{a}_0 - \mathbf{a}\|_2 \\
 & \leq \frac{1}{32} n \theta \|\mathbf{a}_0 - \mathbf{a}\|_2.
 \end{aligned} \tag{SM8.78}$$

The claimed result therefore is followed by combining (SM8.70), (SM8.75) and (SM8.78).  $\blacksquare$

The next thing is to show the operator

$$(n\theta)^{-1} \left( \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0} P_I C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota \right) \tag{SM8.79}$$

contracts  $\mathbf{a}$  toward  $\mathbf{a}_0$ . We first will show that

$$(n\theta)^{-1} \left( \iota^* C_{\mathbf{x}_0}^* C_{\mathbf{a}_0} P_I C_{\mathbf{a}_0}^* C_{\mathbf{x}_0} \iota \right) \approx \mathbf{a}_0 \mathbf{a}_0^* \tag{SM8.80}$$

by seeing  $\iota^* C_{\mathbf{x}_0}^* P_I C_{\mathbf{x}_0} \iota \approx (n\theta) \mathbf{e}_0 \mathbf{e}_0^*$  via sparsity of  $\mathbf{x}_0$ . Finally since the local perturbation on sphere is close to a quadratic function in  $\ell^2$ -norm of difference, we have

$$|\langle \mathbf{a}_0, \mathbf{a} - \mathbf{a}_0 \rangle| \leq \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2^2. \tag{SM8.81}$$

Again, we introduce the following lemma to solidify our claim:

**Lemma SM8.3 (Contraction of  $\mathbf{a}$  to  $\mathbf{a}_0$ ).** *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  $C, C', c, c', c_\mu$  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  $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$ ,*

$$(SM8.82) \quad \left\| \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}_0 - \mathbf{a}) \right\|_2 \leq \frac{1}{32} \|\mathbf{a} - \mathbf{a}_0\|_2 n\theta.$$

*Proof.* Since  $\mathbb{E} \langle \mathbf{P}_I s_i[\mathbf{x}_0], s_j[\mathbf{x}_0] \rangle = 0$  for all  $i \neq j$  and set  $I$ , we calculate

$$(SM8.83) \quad \begin{aligned} \mathbb{E} \left[ \boldsymbol{\iota}_{[\pm p]}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}_{[\pm p]} \right] &= \sum_{i \in [\pm p]} \mathbb{E} \left[ \mathbf{e}_i^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \mathbf{e}_i \right] \mathbf{e}_i \mathbf{e}_i^* \\ &= \mathbb{E} \|\mathbf{x}_0\|_2^2 \mathbf{e}_0 \mathbf{e}_0^* + \sum_{i \in [\pm p] \setminus 0} \mathbb{E} \|\mathbf{P}_I s_i[\mathbf{x}_0]\|_2^2 \mathbf{e}_i \mathbf{e}_i^* \\ &= n\theta \mathbf{e}_0 \mathbf{e}_0^* + n\theta^2 \mathbf{P}_{[\pm p] \setminus 0} \\ &= n\theta^2 \mathbf{I} + n\theta(1 - \theta) \mathbf{e}_0 \mathbf{e}_0^*. \end{aligned}$$

whence

$$(SM8.84) \quad \begin{aligned} \mathbb{E} \left[ \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \right] &= \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}^* \mathbb{E} \left[ \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \right] \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota} \\ &= n\theta^2 \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota} + n\theta(1 - \theta) \mathbf{a}_0 \mathbf{a}_0^*, \end{aligned}$$

implying the expectation is a contraction mapping for  $\mathbf{a}_0 - \mathbf{a}$  when  $c_\mu < \frac{1}{200}$ :

$$(SM8.85) \quad \begin{aligned} &\left\| \mathbb{E} \left[ \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} \right] (\mathbf{a}_0 - \mathbf{a}) \right\|_2 \\ &\leq n\theta^2 \left\| \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota} \right\|_2 \|\mathbf{a}_0 - \mathbf{a}\|_2 + n\theta \|\mathbf{a}_0\|_2 |\langle \mathbf{a}_0, \mathbf{a}_0 - \mathbf{a} \rangle| \\ &\leq n\theta^2 \times 2\tilde{\mu}p \times \|\mathbf{a}_0 - \mathbf{a}\|_2 + \frac{1}{2}n\theta \|\mathbf{a}_0 - \mathbf{a}\|_2^2 \\ &\leq (2c_\mu + \frac{1}{2}c_\mu) \|\mathbf{a}_0 - \mathbf{a}\|_2 n\theta \\ &\leq \frac{1}{64} \|\mathbf{a}_0 - \mathbf{a}\|_2 n\theta. \end{aligned}$$

For each entry of  $\mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0}$ , again from [Section SM1](#) we know with probability at least  $1 - c/n$ :

$$\left| \mathbf{e}_i^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \mathbf{e}_j - \mathbb{E} \left[ \mathbf{e}_i^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \mathbf{e}_j \right] \right| \leq \begin{cases} C' \sqrt{n\theta \log n} & i = j = 0 \\ C' \sqrt{n\theta^2 \log n} & \text{otherwise} \end{cases}.$$

Thus via Gershgorin disc theorem, when  $n > 10^3 C'^2 \theta^{-2} p^2 \log n$ :

$$(SM8.86) \quad \begin{aligned} \lambda_{\max} \left( \boldsymbol{\iota}_{[\pm p]}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}_{[\pm p]} - \mathbb{E} \left[ \boldsymbol{\iota}_{[\pm p]}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{P}_I \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}_{[\pm p]} \right] \right) &\leq C' p \sqrt{n\theta^2 \log n} \\ &\leq \frac{1}{64} n\theta^2. \end{aligned}$$

Finally we combine [\(SM8.85\)](#), [\(SM8.86\)](#) and get

$$(SM8.87) \quad \begin{aligned} \left\| \boldsymbol{\iota}^* \mathbf{C}_{\mathbf{x}_0}^* \mathbf{C}_{\mathbf{a}_0} \mathbf{P}_I \mathbf{C}_{\mathbf{a}_0}^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} (\mathbf{a}_0 - \mathbf{a}) \right\|_2 &\leq \left( \frac{1}{64} n\theta + \frac{1}{64} n\theta^2 \|\mathbf{C}_{\mathbf{a}_0} \boldsymbol{\iota}_{\pm p}\|_2^2 \right) \|\mathbf{a}_0 - \mathbf{a}\|_2 \\ &\leq \frac{1}{32} \|\mathbf{a}_0 - \mathbf{a}\|_2 n\theta. \end{aligned} \quad \blacksquare$$

Lemma SM8.1-SM8.3 together implies the single iterate contract of alternating minimization contracts  $\mathbf{a}$  toward  $\mathbf{a}_0$ . We show it with the following lemma:

**Lemma SM8.4 (Contraction 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  $C, C', c, c_\mu$  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  $\lambda$  and  $\mathbf{a}$  satisfying*

$$(SM8.88) \quad 5\tilde{\mu}\kappa_I \geq \lambda \geq 5\kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2,$$

and suppose  $\mathbf{x}^+$  has the form of (SM8.16), then the solution  $\mathbf{a}^+$  to

$$(SM8.89) \quad \min_{\mathbf{a}' \in \mathbb{R}^p} \left\{ \|\mathbf{a}' * \mathbf{x}^+ - \mathbf{y}\|_2^2 \right\}$$

is unique and satisfies

$$(SM8.90) \quad \|\mathbf{P}_{\mathbb{S}^{p-1}}[\mathbf{a}^+] - \mathbf{a}_0\|_2 \leq \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2.$$

*Proof.* Write  $\mathbf{x}$  as  $\mathbf{x}^+$ , then

$$(SM8.91) \quad \begin{aligned} \lambda_p(\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_x \boldsymbol{\iota}) &= \sigma_{\min}^2(\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota} + \mathbf{C}_{\mathbf{x} - \mathbf{x}_0} \boldsymbol{\iota}) \\ &\geq \left[ \sigma_{\min}(\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}) - \|\mathbf{C}_{\mathbf{x} - \mathbf{x}_0} \boldsymbol{\iota}\| \right]_+^2 \\ &\geq \left[ \sigma_{\min}(\mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}) - 2\sqrt{\kappa_I} \|\mathbf{x} - \mathbf{x}_0\|_2 \right]_+^2 \\ &\geq \left[ \frac{2}{3}\sqrt{\theta n} - 8\lambda\sqrt{\kappa_I}\sqrt{\theta n} \right]_+^2 \\ &\geq \frac{1}{2}\theta n, \end{aligned}$$

where the fourth inequality is derived from using the upper bound of sparse convolution matrix from Remark SM1.6, and the last line holds by knowing  $\lambda < 5c_\mu\kappa_I^{-1}$ . From (SM8.91) we know the least square problem of (SM8.89) has unique solution  $\mathbf{a}^+$ , written as

$$(SM8.92) \quad \mathbf{a}^+ = (\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_x \boldsymbol{\iota})^{-1} \boldsymbol{\iota} \mathbf{C}_x^* \mathbf{y},$$

whence

$$(SM8.93) \quad \begin{aligned} \mathbf{a}^+ - \mathbf{a}_0 &= (\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_x \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{\mathbf{x}_0} \boldsymbol{\iota}) \mathbf{a}_0 - \mathbf{a}_0 \\ &= (\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_x \boldsymbol{\iota})^{-1} (\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}) \mathbf{a}_0. \end{aligned}$$

Combine Lemma SM8.2 and Lemma SM8.3, we know

$$(SM8.94) \quad \|\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}\|_2 \leq \left( C_1 \lambda (\tilde{\lambda} + \tilde{\mu}\kappa_I) + \frac{1}{16} \|\mathbf{a} - \mathbf{a}_0\|_2 \right) n\theta$$

for some constant  $C_1$ . Combine (SM8.91), (SM8.93), (SM8.94) and since  $\lambda < \tilde{\mu}\kappa_I$ , by letting  $c_\mu < \frac{1}{4C_1}$ , we gain

$$(SM8.95) \quad \begin{aligned} \|\mathbf{a}^+ - \mathbf{a}_0\|_2 &\leq \frac{\|\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_{\mathbf{x}_0 - \mathbf{x}} \boldsymbol{\iota}\|_2}{\lambda_p(\boldsymbol{\iota}^* \mathbf{C}_x^* \mathbf{C}_x \boldsymbol{\iota})} \\ &\leq 2C_1 \lambda (\tilde{\lambda} + \tilde{\mu}\kappa_I) + \frac{1}{8} \|\mathbf{a} - \mathbf{a}_0\|_2 \leq \frac{1}{4}. \end{aligned}$$

For the final bound,

$$\begin{aligned}
\left\| \frac{\mathbf{a}^+}{\|\mathbf{a}^+\|_2} - \mathbf{a}_0 \right\|_2 &\leq \frac{\|\mathbf{a}^+ - \mathbf{a}_0\|_2 + \|\mathbf{a}^+\|_2 - 1}{\|\mathbf{a}^+\|_2} \\
&\leq \frac{2\|\mathbf{a}^+ - \mathbf{a}_0\|_2}{1 - \|\mathbf{a}^+ - \mathbf{a}_0\|_2} \leq \frac{8}{3} \|\mathbf{a}^+ - \mathbf{a}_0\|_2, \\
(\text{SM8.96}) \quad &\leq C_2 \lambda \left( \tilde{\lambda} + \tilde{\mu} \kappa_I \right) + \frac{1}{3} \|\mathbf{a} - \mathbf{a}_0\|_2,
\end{aligned}$$

and since  $\lambda > \kappa_I \|\mathbf{a} - \mathbf{a}_0\|_2$ , finally we gain

$$\begin{aligned}
(\text{SM8.96}) &\leq C_2 \left( \lambda \kappa_I + \frac{p \kappa_I \log n}{n \theta} + \tilde{\mu} \kappa_I^2 \right) \|\mathbf{a} - \mathbf{a}_0\|_2 + \frac{1}{3} \|\mathbf{a} - \mathbf{a}_0\|_2 \\
(\text{SM8.97}) \quad &\leq \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2 \quad \blacksquare
\end{aligned}$$

as long as  $n > 20C_2\theta^{-1}p\kappa_I \log n$  and  $c_\mu < \frac{1}{20C_2}$ .

**SM8.3. 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

$$\begin{aligned}
|s_j[\mathbf{a}^{(0)}]^* \mathbf{C}_{\mathbf{a}_0} \mathbf{x}_0| &\geq \left(1 - |\langle \mathbf{a}^{(0)}, \mathbf{a}_0 \rangle|\right) |\mathbf{x}_j| - \left\| \mathbf{P}_{[\pm p] \setminus \{j\}} \mathbf{C}_{\mathbf{a}_0}^* \boldsymbol{\iota}_{s_j}[\mathbf{a}^{(0)}] \right\|_2 \times \sqrt{2} \|\mathbf{x}_0\|_\square \\
&\geq (1 - 2\tilde{\mu}) 6\lambda^{(0)} - 2\tilde{\mu}\sqrt{\kappa_I} \times \sqrt{2\kappa_I} \\
&\geq 5\lambda^{(0)} - 4\lambda^{(0)} \\
(\text{SM8.98}) \quad &= \lambda^{(0)}.
\end{aligned}$$

hence  $I_{>6\lambda^{(0)}} \subseteq I^{(0)}$ , therefore the condition of Lemma SM8.4 is satisfied, implies (5.32) holds for  $k = 0$ .

Suppose it is true for  $1, 2, \dots, k-1$ , such that

$$(\text{SM8.99}) \quad \kappa_I \|\mathbf{a}^{(k)} - \mathbf{a}_0\|_2 \leq \frac{1}{2} \lambda^{(k-1)} = \lambda^{(k)}, \quad \text{and} \quad I_{>3\lambda^{(k-1)}} \subseteq I^{(k)}$$

and since  $I_{>6\lambda^{(k)}} = I_{>3\lambda^{(k-1)}} \subseteq I^{(k)}$ , we can again apply Lemma SM8.4, resulting

$$(\text{SM8.100}) \quad \kappa_I \|\mathbf{a}^{(k+1)} - \mathbf{a}\|_2 \leq \frac{1}{2} \kappa_I \|\mathbf{a}^{(k)} - \mathbf{a}_0\|_2 \leq \frac{1}{2} \lambda^{(k)}$$

as claimed. \blacksquare

**SM8.4. Supporting lemmas for refinement.** The following lemma controls the shift coherence of  $\mathbf{a}$ :

**Lemma SM8.5 (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*

$$(SM8.101) \quad \|\text{off} [C_{\mathbf{a}}^* C_{\mathbf{a}_0}]\|_{\infty} \leq 2\tilde{\mu}$$

$$(SM8.102) \quad \|\text{off} [C_{\mathbf{a}}^* C_{\mathbf{a}}]\|_{\infty} \leq 3\tilde{\mu}$$

*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. ■

From this we obtain the following spectral control on  $C_{\mathbf{a}}^* C_{\mathbf{a}}$ , to simply the notations, we will write

$$(SM8.103) \quad C_{\mathbf{a}I}^* C_{\mathbf{a}I} = \iota_I^* C_{\mathbf{a}}^* C_{\mathbf{a}} \iota_I = [C_{\mathbf{a}}^* C_{\mathbf{a}}]_{I,I}$$

in the latter part of this section.

**Lemma SM8.6 (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*

$$(SM8.104) \quad \left\| [C_{\mathbf{a}}^* C_{\mathbf{a}} - I]_{I,I} \right\|_2 \leq 9\kappa_I \tilde{\mu}.$$

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  $\ell$ :

$$\text{supp} ([C_{\mathbf{a}}^* C_{\mathbf{a}}]_{I_{\ell}, I_{\ell}}) \subseteq I_{\ell} \times (I_{\ell-1} \uplus I_{\ell} \uplus I_{\ell+1}),$$

where  $\ell+1$  and  $\ell-1$  are interpreted cyclically modulo  $r$ .

For an arbitrary  $\mathbf{v} \in \mathbb{R}^{|I|}$ , we calculate

$$(SM8.105) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{I,I} \mathbf{v} \right\|_2^2 = \sum_{\ell=0}^{r-1} \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{I_\ell, I} \mathbf{v} \right\|_2^2$$

$$(SM8.106) \quad = \sum_{\ell=0}^{r-1} \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{I_\ell, I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \mathbf{v}_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_2^2$$

$$(SM8.107) \quad \leq \sum_{\ell=0}^{r-1} \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{I_\ell, I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_F^2 \left\| \mathbf{v}_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_2^2$$

$$(SM8.108) \quad \leq 3\kappa_I^2 \times (3\tilde{\mu})^2 \times \sum_{\ell=0}^{r-1} \left\| \mathbf{v}_{I_{\ell-1} \uplus I_\ell \uplus I_{\ell+1}} \right\|_2^2$$

$$(SM8.109) \quad \leq 3\kappa_I^2 \times 9\tilde{\mu}^2 \times 3 \left\| \mathbf{v} \right\|_2^2,$$

giving the claimed result. ■

As a consequence, we have that

**Corollary SM8.7 (Inverse of  $[\mathbf{C}_a^* \mathbf{C}_a]_{J,J}$ ).** *Suppose that  $\mathbf{a}_0$  is  $\mu$ -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*

$$(SM8.110) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J,J} \right\|_\diamond \leq 9\kappa_I \tilde{\mu}$$

$$(SM8.111) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a]_{J,J}^{-1} - \mathbf{I} \right\|_\diamond \leq 18\kappa_I \tilde{\mu}$$

$$(SM8.112) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a]_{J,J}^{-1} \right\|_\diamond \leq 2.$$

*Proof.* First we prove

$$(SM8.113) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J,J} \right\|_2 \leq 9\kappa_I \tilde{\mu},$$

$$(SM8.114) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J,J} \right\|_{\infty \rightarrow \infty} \leq 6\kappa_I \tilde{\mu},$$

$$(SM8.115) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J,J} \right\|_{\square \rightarrow \square} \leq 6\kappa_I \tilde{\mu}$$

Where the first claim follows from [Lemma SM8.6](#). The second follows by noting that the  $\ell^\infty$  operator norm is the maximum row  $\ell^1$  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

$$(SM8.116) \quad \begin{aligned} \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J,J} \right\|_{\square \rightarrow \square} &\leq \max_{\ell, \ell'} \left\| [\mathbf{C}_a^* \mathbf{C}_a - \mathbf{I}]_{J \cap ([p] + \ell), J \cap ([2p] + \ell')} \right\|_F \\ &\leq 6\kappa_I \tilde{\mu}. \end{aligned}$$

Then we prove

$$\begin{aligned}
 & \left\| [\mathbf{C}_a^* \mathbf{C}_a]_{J,J}^{-1} - \mathbf{I} \right\|_2 \leq 18\kappa_I \tilde{\mu}, \\
 & \left\| [\mathbf{C}_a^* \mathbf{C}_a]_{J,J}^{-1} - \mathbf{I} \right\|_{\infty \rightarrow \infty} \leq 12\kappa_I \tilde{\mu}, \\
 \text{(SM8.117)} \quad & \left\| [\mathbf{C}_a^* \mathbf{C}_a]_{J,J}^{-1} - \mathbf{I} \right\|_{\square \rightarrow \square} \leq 12\kappa_I \tilde{\mu},
 \end{aligned}$$

which are followed from the fact that if  $\|\cdot\|_\diamond$  is a matrix norm and  $\|\Delta\|_\diamond < 1$ , then

$$\left\| (\mathbf{I} + \Delta)^{-1} - \mathbf{I} \right\|_\diamond \leq \frac{\|\Delta\|_\diamond}{1 - \|\Delta\|_\diamond}.$$

Finally, (SM8.112) follows from the triangle inequality. ■

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 SM8.8 (Convolution of  $\mathbf{a}_0 - \mathbf{a}$ ).** *Suppose that  $\mathbf{a}_0$  is  $\mu$ -shift coherent and  $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$ , then for every  $J \subseteq I$ ,*

$$\text{(SM8.118)} \quad \left\| [\mathbf{C}_a^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}}]_{J,J} \right\|_{\square \rightarrow \infty} \leq \sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2$$

$$\text{(SM8.119)} \quad \left\| [\mathbf{C}_a^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}}]_{J,J} \right\|_{\square \rightarrow \square} \leq \sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2$$

*Proof.* For the first inequality, we have

$$\begin{aligned}
 \left\| [\mathbf{C}_a^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}}]_{J,J} \mathbf{v} \right\|_{\square \rightarrow \infty} &= \max_{j \in J, \|\mathbf{v}\|_\square = 1} |\langle s_j[\mathbf{a}], (\mathbf{a}_0 - \mathbf{a}) * \mathbf{v} \rangle| \\
 &\leq \max_{j \in [n], \|\mathbf{v}\|_\square = 1} \left\| \mathbf{P}_{[p]+j} [(\mathbf{a}_0 - \mathbf{a}) * \mathbf{v}] \right\|_2 \\
 &\leq \|\mathbf{a} - \mathbf{a}_0\|_2 \times \max_{j \in [n], \|\mathbf{v}\|_\square = 1} \left\| \mathbf{P}_{[\pm p]+j} \mathbf{v} \right\|_1 \\
 \text{(SM8.120)} \quad &\leq \sqrt{2\kappa_I} \|\mathbf{a}_0 - \mathbf{a}\|_2
 \end{aligned}$$

The second inequality is derived by

$$\begin{aligned}
 \left\| [\mathbf{C}_a^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}}]_{J,J} \right\|_{\square \rightarrow \square} &\leq \max_{\ell, \ell'} \left\| [\mathbf{C}_a^* \mathbf{C}_{\mathbf{a}_0 - \mathbf{a}}]_{J \cap ([p]+\ell), J \cap ([2p]+\ell')} \right\|_F \\
 &\leq \sqrt{2\kappa_I^2 \max_{i,j} |\langle s_i[\mathbf{a}], s_j[\mathbf{a}_0 - \mathbf{a}] \rangle|^2} \\
 \text{(SM8.121)} \quad &\leq \sqrt{2\kappa_I} \|\mathbf{a} - \mathbf{a}_0\|_2,
 \end{aligned}$$

finishing the proof. ■

Again, using a variant of the argument for Lemma SM8.6, we have the following:

**Lemma SM8.9 (Off-diagonal of submatrix of  $\mathbf{C}_a^* \mathbf{C}_{a_0}$ ).** Suppose that  $\mathbf{a}_0$  is  $\mu$ -shift coherent and  $\|\mathbf{a} - \mathbf{a}_0\|_2 \leq \tilde{\mu}$ . For any  $J \subset I$ , if

$$(SM8.122) \quad \kappa_J = \max_{\ell} |J \cap \{\ell, \ell + 1, \dots, \ell + p - 1\}|$$

$$(SM8.123) \quad \kappa_{I \setminus J} = \max_{\ell} |(I \setminus J) \cap \{\ell, \ell + 1, \dots, \ell + p - 1\}|$$

Then

$$(SM8.124) \quad \left\| [\mathbf{C}_a^* \mathbf{C}_{a_0}]_{J, I \setminus J} \right\|_2 \leq 6\sqrt{\kappa_J \kappa_{I \setminus J} \tilde{\mu}}.$$

*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

$$\begin{aligned} \left\| [\mathbf{C}_a^* \mathbf{C}_{a_0}]_{J, I \setminus J} \mathbf{v} \right\|_2^2 &= \sum_{\ell=0}^{r-1} \left\| [\mathbf{C}_a^* \mathbf{C}_{a_0}]_{J_{\ell}, I \setminus J} \mathbf{v} \right\|_2^2 \\ &= \sum_{\ell=0}^{r-1} \left\| [\mathbf{C}_a^* \mathbf{C}_{a_0}]_{J_{\ell}, L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \mathbf{v}_{L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \right\|_2^2 \\ &\leq 4\tilde{\mu}^2 \times \kappa_J \times 3\kappa_{I \setminus J} \times \sum_{\ell=0}^{r-1} \left\| \mathbf{v}_{L_{\ell-1} \cup L_{\ell} \cup L_{\ell+1}} \right\|_2^2 \\ (SM8.125) \quad &\leq 4\tilde{\mu}^2 \times \kappa_J \times 3\kappa_{I \setminus J} \times 3\|\mathbf{v}\|_2^2, \end{aligned}$$

giving the result. ■

**Lemma SM8.10 (Perturbation of vector over sphere).** If both  $\mathbf{a}, \mathbf{a}_0$  are unit vectors in inner product space, then

$$(SM8.126) \quad |\langle \mathbf{a}, \mathbf{a} - \mathbf{a}_0 \rangle| \leq \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2^2.$$

*Proof.* Via simple norm inequalities:

$$(SM8.127) \quad \frac{1}{2} \|\mathbf{a} - \mathbf{a}_0\|_2^2 = 1 - \langle \mathbf{a}, \mathbf{a}_0 \rangle = 1 - \langle \mathbf{a}, \mathbf{a}_0 - \mathbf{a} + \mathbf{a} \rangle = \langle \mathbf{a}, \mathbf{a} - \mathbf{a}_0 \rangle > 0 \quad \blacksquare$$

**Lemma SM8.11 (Convolution of short and sparse).** Suppose  $\boldsymbol{\delta} \in \mathbb{R}^p$ , and  $\mathbf{v} \in \mathbb{R}^n$  where  $\text{supp}(\mathbf{v}) = I$  satisfies

$$(SM8.128) \quad \max_{\ell \in [n]} |I \cap ([p] + \ell)| \leq \kappa$$

then

$$(SM8.129) \quad \|\boldsymbol{\delta} * \mathbf{v}\|_2 \leq \sqrt{2\kappa} \|\boldsymbol{\delta}\|_2 \|\mathbf{v}\|_2$$

*Proof.* Since every  $p$ -contiguous segment of  $I$  has at most  $\kappa$  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:

$$\begin{aligned}
 I_1 &= \{i_1, i_{\kappa+1}, i_{2\kappa+1}, \dots\} \cap \{0, \dots, n-p-1\}, \\
 I_2 &= \{i_2, i_{\kappa+2}, i_{2\kappa+2}, \dots\} \cap \{0, \dots, n-p-1\}, \\
 &\vdots \\
 \text{(SM8.130)} \quad I_\kappa &= \{i_\kappa, i_{2\kappa}, i_{3\kappa}, \dots\} \cap \{0, \dots, n-p-1\}, \\
 \text{(SM8.131)} \quad R &= I \cap \{n-p, \dots, n-1\}.
 \end{aligned}$$

Then the  $p$ -separating property gives  $\|\delta * \mathbf{P}_{I_i} \mathbf{v}\|_2 = \|\delta\|_2 \|\mathbf{P}_{I_i} \mathbf{v}\|_2$ . Hence:

$$\begin{aligned}
 \|\delta * \mathbf{P}_I \mathbf{v}\|_2 &= \left\| \sum_{i \in \kappa} \delta * \mathbf{P}_{I_i} \mathbf{v} + \delta * \mathbf{P}_R \mathbf{v} \right\|_2 \leq \sum_{i \in \kappa} \|\delta * \mathbf{P}_{I_i} \mathbf{v}\|_2 + \|\delta * \mathbf{P}_R \mathbf{v}\|_2 \\
 &= \|\delta\|_2 \sum_{i \in \kappa} \|\mathbf{v}_{I_i}\|_2 + \|\delta\|_2 \|\mathbf{P}_R \mathbf{v}\|_1 \\
 &\leq \sqrt{\kappa} \|\mathbf{v}_{I_1 \uplus \dots \uplus I_\kappa}\|_2 \|\delta\|_2 + \sqrt{\kappa} \|\mathbf{v}_R\|_2 \|\delta\|_2 \\
 \text{(SM8.132)} \quad &\leq \sqrt{2\kappa} \|\mathbf{v}\|_2 \|\delta\|_2,
 \end{aligned}$$

where the last two inequalities were coming from Cauchy-Schwartz. ■

**SM9. 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  $\chi(\boldsymbol{\beta})_i$ , shift space Hessian diagonals  $\|\mathbf{P}_{I(\mathbf{a})} s_{-i}[\mathbf{x}_0]\|_2$ , and the set of gradients discontinuity entries  $|J_B(\mathbf{a})|$ .

### SM9.1. Proof of Corollary SM3.4.

*Proof.* 1. ( $\varepsilon$ -net) Write  $\mathbf{x}$  as  $\mathbf{x}_0$  and  $\|\boldsymbol{\beta}\|_2 = \eta$  through out this proof, firstly from [Definition SM2.1](#) 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  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  for sphere of radius  $\eta$ . From [Lemma SM10.5](#) we know for any  $c_2 < 1$ :

$$\text{(SM9.1)} \quad |\mathcal{N}_\varepsilon| \leq \left(\frac{3\eta}{\varepsilon}\right)^{2p} \leq \left(\frac{3n^{3/2} p^2}{c_2}\right)^{2p} \leq \left(\frac{3np^2}{c_2}\right)^{3p}$$

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:

$$\text{(SM9.2)} \quad \mathcal{E}_{\text{Net}} := \left\{ \forall i \in [n], \quad \sigma_i n^{-1} \chi[\boldsymbol{\beta}_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}_\varepsilon]_i} < \frac{c_1 \theta}{p^{3/2}} \quad \forall \boldsymbol{\beta}_\varepsilon \in \mathcal{N}_{\varepsilon, i} \right\}$$

2. (Lipschitz constant) The Lipschitz constant  $L$  of  $\chi[\cdot]_i$  w.r.t  $\boldsymbol{\beta}$  is bounded in terms of  $\mathbf{x}$

regardless of entry  $i$ :

$$\begin{aligned}
|\chi[\boldsymbol{\beta}]_i - \chi[\boldsymbol{\beta}']_i| &\leq \left| \mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta} \right] - \mathbf{e}_i^* \check{\mathbf{C}}_{\mathbf{x}} \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta}' \right] \right| \\
&\leq \|\mathbf{x}\|_2 \left\| \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta} \right] - \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta}' \right] \right\|_2 \\
&\leq \|\mathbf{x}\|_2 \sqrt{\sum_{j \in [n]} \left| \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta} \right]_j - \mathcal{S}_\lambda \left[ \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta}' \right]_j \right|^2} \\
&\leq \|\mathbf{x}\|_2 \left\| \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta} - \check{\mathbf{C}}_{\mathbf{x}} \boldsymbol{\beta}' \right\|_2 \\
(\text{SM9.3}) \quad &\leq \|\mathbf{x}\|_2 \cdot \|\mathbf{x}\|_1 \cdot \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2 =: L \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2
\end{aligned}$$

Define the event that  $\chi[\boldsymbol{\beta}]_i$  that has small Lipschitz constant as

$$(\text{SM9.4}) \quad \mathcal{E}_{\text{Lip}} := \left\{ L < 2n^{3/2}\theta \right\}$$

on the event  $\mathcal{E}_{\text{Lip}}$ , for every points in  $\mathfrak{A}(\mathcal{S}_\tau, \gamma(c_\mu))$  and  $i \in [n]$ , there exists some  $\boldsymbol{\beta}_\varepsilon \in \mathcal{N}_{\varepsilon, i}$  such that

$$(\text{SM9.5}) \quad \left| \left( \sigma_i n^{-1} \chi[\boldsymbol{\beta}]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}]_i} \right) - \left( \sigma_i n^{-1} \chi[\boldsymbol{\beta}_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}_\varepsilon]_i} \right) \right| \leq 2L\varepsilon \leq \frac{c_2\theta}{p^{3/2}}$$

On event  $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$ , (SM9.2), (SM9.5) implies  $\chi[\boldsymbol{\beta}]$  is well concentrated entrywise and anywhere in  $\cup_{|\tau| \leq k} \mathfrak{A}(\mathcal{S}_\tau, \gamma(c_\mu))$ :

$$(\text{SM9.6}) \quad \left| \sigma_i n^{-1} \chi[\boldsymbol{\beta}]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}]_i} \right| \leq \frac{(c_1 + c_2)\theta}{p^{3/2}}, \quad \forall \mathbf{a} \in \cup_{k \leq k} \mathfrak{A}(\mathcal{S}_\tau, \gamma(c_\mu)), \quad \forall i \in [n]$$

as desired, where, using Lemma SM1.2,

$$(\text{SM9.7}) \quad \mathbb{P}[\mathcal{E}_{\text{Lip}}^c] \leq \mathbb{P}[\|\mathbf{x}\|_2^2 > 2n\theta] \leq 1/n;$$

and using union bound,

$$\begin{aligned}
(\text{SM9.8}) \quad \mathbb{P}[\mathcal{E}_{\text{Net}}^c] &\leq \mathbb{P} \left[ \max_{\substack{\mathbf{a}_\varepsilon \in \mathcal{N}_{\varepsilon, i} \\ i \in [n]}} \sigma_i n^{-1} \chi[\boldsymbol{\beta}_\varepsilon]_i - \sigma_i n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}_\varepsilon]_i} > \frac{c_1\theta}{p^{3/2}} \right] \\
&\leq n |\mathcal{N}_\varepsilon| \mathbb{P} \left[ \sigma_0 n^{-1} \chi[\boldsymbol{\beta}_\varepsilon]_0 - \sigma_0 n^{-1} \overline{\mathbb{E} \chi[\boldsymbol{\beta}_\varepsilon]_0} > \frac{c_1\theta}{p^{3/2}} \right].
\end{aligned}$$

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  $\mathbf{Z}_j(\boldsymbol{\beta})$  from (SM3.9) is independent of  $\mathbf{Z}_{j+2p}(\boldsymbol{\beta})$  for all  $j \in [n]$  while all  $\mathbf{Z}_j$  are identical distributed. We write  $\chi[\boldsymbol{\beta}]_0$  as sum of iid r.v.s. as

$$\chi[\boldsymbol{\beta}]_0 = \sum_{j \in [n]} \mathbf{Z}_j(\boldsymbol{\beta}) = \sum_{k \in [2p]} \left( \sum_{t=0}^{n/2p-1} \mathbf{Z}_{k+2tp}(\boldsymbol{\beta}) \right)$$

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  $\chi[\boldsymbol{\beta}]_0$  as

$$(SM9.9) \quad \mathbb{P} \left[ n^{-1} \chi[\boldsymbol{\beta}]_0 > n^{-1} \overline{\mathbb{E}\chi(\boldsymbol{\beta})}_0 + \frac{c_1 \theta}{p^{3/2}} \right] \leq 2p \cdot \mathbb{P} \left[ \sum_{t=0}^{n/2p-1} \mathbf{Z}_{2tp}(\boldsymbol{\beta}) > \frac{n}{2p} \mathbb{E}\mathbf{Z}(\boldsymbol{\beta}) + \frac{c_1 n \theta}{2p^{5/2}} \right]$$

The moments of  $\mathbf{Z}_0$  can be bounded by using  $|\mathbf{Z}_0(\boldsymbol{\beta})| \leq |\mathbf{x}_0| |\boldsymbol{\beta}_0 \mathbf{x}_0 + \mathbf{s}_0| \leq \boldsymbol{\beta}_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|$  where  $\mathbf{s}_0 = \sum_{\ell \neq 0} \mathbf{x}_\ell \boldsymbol{\beta}_\ell$ , write  $\mathbf{x} = \boldsymbol{\omega} \circ \mathbf{g} \sim_{\text{i.i.d.}} \text{BG}(\theta)$ . For the 2-norm we know

$$(SM9.10) \quad \mathbb{E} |\mathbf{s}_0|^2 = \mathbb{E} \left| \sum_{\ell} \mathbf{x}_\ell \boldsymbol{\beta}_\ell \right|^2 \leq \theta \|\boldsymbol{\beta}\|_2^2 \leq \theta \left( 1 + c_\mu + \frac{c_\mu}{\theta k^2} \right) \leq \frac{1}{2}$$

As for the  $q$ -norm, use the moment generating function bound, such that for all  $t \geq 0$ :

$$(SM9.11) \quad \begin{aligned} \mathbb{E} |\mathbf{s}_0|^q &\leq q! t^{-q} \mathbb{E} \exp [t |\mathbf{s}_0|] \leq q! t^{-q} \prod_{\ell} \mathbb{E}_{\boldsymbol{\omega}_\ell, \mathbf{g}_\ell} \exp [t \boldsymbol{\omega}_\ell |\mathbf{g}_\ell| |\boldsymbol{\beta}_\ell|] \\ &\leq 2q! t^{-q} \prod_{\ell} \mathbb{E}_{\boldsymbol{\omega}_\ell} \exp [\boldsymbol{\omega}_\ell t^2 \boldsymbol{\beta}_\ell^2 / 2] \\ &\leq 2q! t^{-q} \prod_{\ell} (1 - \theta + \theta \exp [t^2 \boldsymbol{\beta}_\ell^2 / 2]) \end{aligned}$$

notice that the entrywise twice derivative of (SM9.11) w.r.t.  $\boldsymbol{\beta}_\ell^2$ 's are always positive, this function is convex for all  $\boldsymbol{\beta}_\ell^2$ . Constrain on the polytope  $\sum_{\ell} \boldsymbol{\beta}_\ell^2 \leq \|\boldsymbol{\beta}\|_2^2$ , the maximizer of (SM9.11) w.r.t.  $\boldsymbol{\beta}_\ell^2$ 's occurs and a vertex point where  $\boldsymbol{\beta}_0^2 = \|\boldsymbol{\beta}\|_2^2$ . Thus

$$(SM9.11) \leq 2q! t^{-q} \left( 1 - \theta + \theta \exp [t^2 \|\boldsymbol{\beta}\|_2^2 / 2] \right) \prod_{\ell \neq 0} (1 - \theta + \theta e^0) \leq 2q! t^{-q} (1 + \theta \exp [\|\boldsymbol{\beta}\|_2^2 t^2 / 2]).$$

Choose  $t = \sqrt{q} / \|\boldsymbol{\beta}\|_2$ , use  $q!! > (q!/2) \cdot (e/q)^{q/2}$ , we have

$$(SM9.12) \quad \mathbb{E} |\mathbf{s}_0|^q \leq 2q! q^{-q/2} \|\boldsymbol{\beta}\|_2^q (1 + \theta \exp [q/2]) \leq 8 \|\boldsymbol{\beta}\|_2^q \max \left\{ e^{-q/2}, \theta \right\} q!!.$$

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 SM10.2, (SM9.10) and (SM9.12), obtain for  $q \geq 3$ ,

$$\begin{aligned} \mathbb{E}\mathbf{Z}(\boldsymbol{\beta})^2 &\leq \mathbb{E} (\boldsymbol{\beta}_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|)^2 \leq 2\mathbb{E} [\boldsymbol{\beta}_0^2 \mathbf{x}_0^4 + \mathbf{x}_0^2 \mathbf{s}_0^2] \leq 6\theta + 2\theta^2 \|\boldsymbol{\beta}\|_2^2 \leq 7\theta, \\ \mathbb{E}\mathbf{Z}(\boldsymbol{\beta})^q &\leq \mathbb{E} (\boldsymbol{\beta}_0 \mathbf{x}_0^2 + |\mathbf{x}_0| |\mathbf{s}_0|)^q \leq 2^{q-1} \left( \mathbb{E}\mathbf{x}_0^{2q} + \mathbb{E} |\mathbf{x}_0|^q \mathbb{E} |\mathbf{s}_0|^q \right) \\ &\leq \theta 2^{q-1} (2q-1)!! + \theta 2^{q-1} (q-1)!! \left( 8 \|\boldsymbol{\beta}\|_2^q \max \left\{ e^{-q/2}, \theta \right\} q!! \right) \\ &\leq \theta 4^q q! + \theta 2^q \|\boldsymbol{\beta}\|_2^q q!. \end{aligned}$$

Thus, recall that  $\|\boldsymbol{\beta}\|_2 = \eta$ , use  $(\sigma^2, R) = (8\theta\eta^2, 4\eta)$ , from (SM9.8)-(SM9.9), apply Bernstein inequality Lemma SM10.4 with  $n \geq Cp^5\theta^{-2} \log p$ , and  $c_1, c_2 \in [0, 1]$  we have

$$\begin{aligned}
\mathbb{P}[\mathcal{E}_{\text{Net}}^c] &\leq 2np |\mathcal{N}_\varepsilon| \cdot \mathbb{P} \left[ \sum_{t=0}^{n/2p-1} \mathbf{Z}_{2tp}(\boldsymbol{\beta}) > \frac{n}{2p} \mathbb{E} \mathbf{Z}(\boldsymbol{\beta}) + \frac{c_1 n \theta}{2p^{5/2}} \right] \\
&\leq 2np \left( \frac{3np^2}{c_2} \right)^{3p} \exp \left( \frac{-(c_1 n \theta / 2p^{5/2})^2}{16n\theta\eta^2/2p + 8\eta c_1 n \theta / 2p^{5/2}} \right) \\
&\leq \exp \left( 4p \log \left( \frac{3np^2}{c_2} \right) - \frac{(c_1 n \theta / 2p^{5/2})^2}{16n\theta\eta^2/p} \right) \\
&\leq \exp \left( 4p \log \left( \frac{3np^2}{c_2} \right) - \frac{c_1^2 n \theta^2}{64p^4} \right) \\
\text{(SM9.13)} \quad &\leq \exp \left( \frac{-c_1^2 n \theta^2}{100p^4} \right) \leq \frac{1}{n}
\end{aligned}$$

when  $\frac{C}{\log C} > \frac{10^5}{c_1^2 c_2}$ . The proof of lower bound and negative  $\boldsymbol{\beta}_0$  is derived in the same manner. ■

### SM9.2. Proof of Corollary SM4.3.

*Proof.* Write  $\mathbf{x}$  as  $\mathbf{x}_0$  though our this proof. Write  $\boldsymbol{\beta}_i \mathbf{x}_j + \mathbf{s}_j = \sum_{\ell \in [\pm p]} \boldsymbol{\beta}_\ell \mathbf{x}_{\ell-i+j} = \langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle$ , and the support w.r.t. some  $\mathbf{a}$  as  $I(\boldsymbol{\beta})$ . Define the random variable  $\mathbf{Z}_{ij}(\boldsymbol{\beta})$  as

$$\text{(SM9.14)} \quad \|\mathbf{P}_{I(\boldsymbol{\beta})} \mathbf{s}_{-i}[\mathbf{x}]\|_2^2 = \sum_{j \in [n]} \mathbf{x}_j^2 \mathbf{1}_{\{|\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| > \lambda\}} =: \sum_{j \in [n]} \mathbf{Z}_{ij}(\boldsymbol{\beta})$$

and define  $\{\bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta})\}_{j \in [n]}$  that are independent r.v.s. and as a upper bounding function of  $\mathbf{Z}_{ij}(\boldsymbol{\beta})$  as

$$\text{(SM9.15)} \quad \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) := \begin{cases} \mathbf{x}_j^2, & |\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| > \lambda \\ 0, & |\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| < \lambda/2, \\ \frac{\mathbf{x}_j^2}{\lambda/2} (|\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| - \lambda/2), & \text{otherwise} \end{cases}$$

Similar to proof of Corollary SM3.4. Let  $\|\boldsymbol{\beta}\|_2 \leq \eta \leq \sqrt{p}$ . Define  $\varepsilon = \frac{c'_2 \lambda}{24np\sqrt{p\theta \log n \log \theta^{-1}}}$  for some  $c'_2 > 0$  and consider the  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  for sphere of radius  $\eta$ . From Lemma SM10.5 we know

$$\text{(SM9.16)} \quad |\mathcal{N}_\varepsilon| \leq \left( \frac{3\eta}{\varepsilon} \right)^{2p} \leq \left( \frac{72}{c'_2 c_\lambda} np^2 \sqrt{\theta |\tau| \log n \log \theta^{-1}} \right)^{2p} \leq \left( \frac{72}{c'_2 c_\lambda} np^2 \log n \right)^{2p},$$

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:

$$\text{(SM9.17)} \quad \mathcal{E}_{\text{Net}} := \left\{ \forall i \in [n], \left| n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}_\varepsilon) - \mathbb{E} \bar{\mathbf{Z}}_i(\boldsymbol{\beta}_\varepsilon) \right| \leq \frac{c'_1 \theta}{p} \quad \forall \boldsymbol{\beta}_\varepsilon \in \mathcal{N}_{\varepsilon, i} \right\},$$

Also,  $\sum_j \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta})$  is a Lipschitz function over  $\boldsymbol{\beta}$  for every  $i \in [n]$  as

$$\begin{aligned}
 \left| \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) - \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}') \right| &\leq \sum_{j \in [n]} \frac{\mathbf{x}_j^2}{\lambda/2} |\langle \boldsymbol{\beta} - \boldsymbol{\beta}', \mathbf{x}_{[\pm p]-i+j} \rangle| \\
 &\leq \sum_{j \in [n]} \frac{\mathbf{x}_j^2 \|\mathbf{x}_{[\pm p]-i+j}\|_2}{\lambda/2} \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2, \\
 &\leq \frac{1}{\lambda/2} \|\mathbf{x}\|_2^2 \cdot \max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 \cdot \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2 \\
 \text{(SM9.18)} \quad &:= L \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2,
 \end{aligned}$$

and define event  $\mathcal{E}_{\text{Lip}}$  such that the Lipschitz constant is bounded as

$$\text{(SM9.19)} \quad \mathcal{E}_{\text{Lip}} := \left\{ L \leq 12n\theta \sqrt{p\theta \log n \log \theta^{-1} \lambda^{-1}} \right\},$$

then on event  $\mathcal{E}_{\text{Lip}}$ , for any points  $\boldsymbol{\beta}$  in  $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  and  $i \in [n]$ , there exists some  $\boldsymbol{\beta}_\varepsilon$  in  $\mathcal{N}_{\varepsilon,i}$  with  $\|\boldsymbol{\beta} - \boldsymbol{\beta}_\varepsilon\|_2 \leq \varepsilon$ , and thus

$$\text{(SM9.20)} \quad \left| \left( n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) - \mathbb{E} \bar{\mathbf{Z}}_i(\boldsymbol{\beta}) \right) - \left( n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}_\varepsilon) - \mathbb{E} \bar{\mathbf{Z}}_i(\boldsymbol{\beta}_\varepsilon) \right) \right| \leq 2L\varepsilon \leq \frac{c'_2 \theta}{p}.$$

On event  $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$ , from (SM9.17), (SM9.20), we can conclude that for all  $\boldsymbol{\beta} \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  and  $i \in [n]$  that:

$$\begin{aligned}
 n^{-1} \|\mathbf{P}_{I(\boldsymbol{\beta})s-i}[\mathbf{x}_0]\|_2^2 - n^{-1} \mathbb{E} \|\mathbf{P}_{I(\boldsymbol{\beta})s-i}[\mathbf{x}_0]\|_2^2 &\leq n^{-1} \sum_{j \in [n]} \bar{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) - \mathbb{E} \bar{\mathbf{Z}}_i(\boldsymbol{\beta}) \\
 \text{(SM9.21)} \quad &\leq \frac{(c'_1 + c'_2)\theta}{p}
 \end{aligned}$$

as desired, where the error probability of  $\mathcal{E}_{\text{Lip}}^c$  is bounded using Lemma SM1.2 and Lemma SM1.3, which give

$$\begin{aligned}
 \mathbb{P}[\mathcal{E}_{\text{Lip}}^c] &\leq \mathbb{P}[\|\mathbf{x}\|_2^2 > 2n\theta] + \mathbb{P}\left[\max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 > 3\sqrt{p\theta \log n \log \theta^{-1}}\right] \\
 \text{(SM9.22)} \quad &\leq 3/n,
 \end{aligned}$$

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  $j$ :

$$\mathbb{P}[\mathcal{E}_{\text{Net}}^c] \leq 2np \cdot |\mathcal{N}_\varepsilon| \cdot \mathbb{P}\left[\left|\sum_k^{n/2p} \bar{\mathbf{Z}}_{i,2kj}(\boldsymbol{\beta}) - \frac{n}{2p} \mathbb{E} \bar{\mathbf{Z}}_i(\boldsymbol{\beta})\right| \geq \frac{c'_1 n \theta}{2p^2}\right].$$

Now we calculate the variance and  $L^q$ -norm of  $\sum_k \bar{\mathbf{Z}}_{i,2kj}$  for  $q \geq 3$ :

$$(SM9.23) \quad \begin{cases} \mathbb{E} \bar{\mathbf{Z}}_{i,j}^2 \leq \mathbb{E} \mathbf{x}_j^4 \leq 3\theta \\ \mathbb{E} \bar{\mathbf{Z}}_{i,j}^q \leq \mathbb{E} \mathbf{x}_j^{2q} \leq \theta(2q-1)!! \leq \frac{1}{2} \cdot (3\theta) \cdot 2^{q-2} q! \end{cases}$$

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

$$(SM9.24) \quad \begin{aligned} & 2np |\mathcal{N}_\varepsilon| \mathbb{P} \left[ \left| \sum_k^{n/2p} \bar{\mathbf{Z}}_{i,2kj}(\boldsymbol{\beta}) - \frac{n}{2p^2} \mathbb{E} \bar{\mathbf{Z}}_i \right| \geq \frac{c'_1 n \theta}{2p^2} \right] \\ & \leq \exp \left[ \log(2np) + 2p \log \left( \frac{72}{c'_2 c_\lambda} np^2 \log n \right) - \frac{(c'_1 n \theta / 2p^2)^2}{6n\theta/2p + 4c'_1 n \theta / 2p^2} \right] \\ & \leq \exp \left[ 3p \log \left( \frac{72}{c'_2 c_\lambda} np^2 \log n \right) - \frac{c'_1{}^2 n \theta}{24p^3} \right] \\ & \leq \exp[-c'_1{}^2 n \theta / (50p^3)] \leq 1/n, \end{aligned}$$

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 (SM4.9) can be derived by defining  $\underline{\mathbf{Z}}_{ij}$  as

$$(SM9.25) \quad \underline{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) := \begin{cases} \mathbf{x}_j^2, & |\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| > 3\lambda/2 \\ 0, & |\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| < \lambda \\ \frac{\mathbf{x}_j^2}{\lambda/2} (|\langle \boldsymbol{\beta}, \mathbf{x}_{[\pm p]-i+j} \rangle| - \lambda), & \text{otherwise} \end{cases},$$

and define  $\mathcal{E}_{\text{Net}}, \mathcal{E}_{\text{Lip}}$  similarly, such that on intersection of these events,

$$(SM9.26) \quad \begin{aligned} n^{-1} \|\mathbf{P}_{I(\boldsymbol{\beta})} s_{-i}[\mathbf{x}]\|_2^2 - n^{-1} \mathbb{E} \|\mathbf{P}_{I(\boldsymbol{\beta})} s_{-i}[\mathbf{x}]\|_2^2 & \geq n^{-1} \sum_{j \in [n]} \underline{\mathbf{Z}}_{ij}(\boldsymbol{\beta}) - \mathbb{E} \underline{\mathbf{Z}}_i(\boldsymbol{\beta}) \\ & \geq \frac{(c'_1 + c'_2)\theta}{p} \end{aligned}$$

as desired. ■

### SM9.3. Proof of Lemma SM5.5 .

*Proof.* 1. (Expectation upper bound) We will write  $\mathbf{x}$  as  $\mathbf{x}_0$ . Similar to proof of Corollary SM3.4 let  $\|\boldsymbol{\beta}\|_2 \leq \eta \leq \sqrt{p}$ . For each  $i \in [n]$ , define the random variable

$$(SM9.27) \quad \mathbf{X}_i(\boldsymbol{\beta}) = \mathbf{1}_{\{|\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle - \lambda| \leq B\}} + \mathbf{1}_{\{|\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle + \lambda| \leq B\}},$$

then number of indices for vector  $\mathbf{x} * \check{\boldsymbol{\beta}}$  that are within  $B$  of  $\pm\lambda$  is a random variable  $\sum_{i \in [n]} \mathbf{X}_i(\boldsymbol{\beta})$ . For each of the  $\mathbf{X}_i(\boldsymbol{\beta})$ 's consider an upper bound  $\bar{\mathbf{X}}_i(\boldsymbol{\beta})$  defined as

$$(SM9.28) \quad \bar{\mathbf{X}}_i(\boldsymbol{\beta}) = \begin{cases} \frac{1}{M} (|\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle| - (\lambda - B - M)) & |\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle| \in [\lambda - B - M, \lambda - B] \\ 1 & |\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle| \in [\lambda - B, \lambda + B] \\ \frac{1}{M} ((\lambda + B + M) - |\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle|) & |\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle| \in [\lambda + B, \lambda + B + M] \\ 0 & \text{else} \end{cases}$$

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  $I(\mathbf{a})$ ,  $\langle \mathbf{x}, \boldsymbol{\beta} \rangle = \langle \mathbf{g}, \mathbf{P}_{I(\mathbf{a})}\boldsymbol{\beta} \rangle \sim \mathcal{N}(0, \|\mathbf{P}_{I(\mathbf{a})}\boldsymbol{\beta}\|_2^2)$ . For all realizations of  $I(\mathbf{a})$ , the variance  $\|\mathbf{P}_{I(\mathbf{a})}\boldsymbol{\beta}\|_2^2$  is bounded by  $\|\mathbf{P}_{I(\mathbf{a})}\boldsymbol{\beta}\|_2^2 \leq \|\boldsymbol{\beta}\|_2^2 \leq p$ . 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  $\sum_i \bar{\mathbf{X}}_i(\boldsymbol{\beta})$  can be upper bounded as

$$\begin{aligned}
 \sum_{i \in [n]} \mathbb{E} [\bar{\mathbf{X}}_i(\boldsymbol{\beta})] &\leq (2n) \cdot \mathbb{P}[\langle \mathbf{x}, \boldsymbol{\beta} \rangle \in [\lambda - B - M, \lambda + B + M]] \\
 &\leq (2n) \cdot 2(B + M) \sup_{\sigma^2 \in (0, p]} \max_{t \in [\lambda - B - M, \lambda + B + M]} f_\sigma(t) \\
 &\leq 4n(B + M) \sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda - B - M) \\
 \text{(SM9.29)} \quad &\leq 4n(B + M) \sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda/2).
 \end{aligned}$$

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

$$\text{(SM9.30)} \quad \sup_{\sigma^2 \in (0, p]} f_\sigma(\lambda/2) \leq f_{\lambda/2}(\lambda/2) \leq \frac{1}{\sqrt{2\pi}(\lambda/2)} \exp\left(-\frac{1}{2}\right) \leq \frac{1}{2\lambda},$$

whence, by letting  $B \leq c\lambda\theta^2 / (p \log n)$ , the upper bound of expectation become:

$$\text{(SM9.31)} \quad \sum_{i \in [n]} \mathbb{E} [\bar{\mathbf{X}}_i(\boldsymbol{\beta})] \leq \frac{4n}{2\lambda}(B + M) \leq \frac{4cn\theta^2}{p \log n} =: n\mathbb{E}\bar{\mathbf{X}}(\boldsymbol{\beta}).$$

2. ( $\varepsilon$ -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  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  for sphere of radius  $\eta \leq \sqrt{p}$ . From [Lemma SM10.5](#) we know

$$\text{(SM9.32)} \quad |\mathcal{N}_\varepsilon| \leq \left(\frac{3\eta}{\varepsilon}\right)^{2p} \leq \left(\frac{81|\tau|p^6 \log^5 n \log \theta^{-1}}{c^4 c_\lambda^2 \theta^7}\right)^p \leq \left(\frac{2p \log n}{c \cdot c_\lambda}\right)^{13p}$$

and define an event such that all center of subsets in  $\mathcal{N}_\varepsilon$  are being well-behaved:

$$\text{(SM9.33)} \quad \mathcal{E}_{\text{Net}} := \left\{ \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}_\varepsilon) - n\mathbb{E}\bar{\mathbf{X}}(\boldsymbol{\beta}_\varepsilon) < \frac{18cn\theta^2}{p \log n} \quad \forall \boldsymbol{\beta}_\varepsilon \in \mathcal{N}_\varepsilon, \right\}$$

3. (Lipschitz constant) Furthermore, the function  $\sum_i^n \bar{\mathbf{X}}_i(\boldsymbol{\beta})$  is Lipschitz over  $\boldsymbol{\beta}$  such that

$$\begin{aligned} \left| \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}) - \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}') \right| &\leq \sum_{i \in [n]} \frac{1}{M} |\langle s_i[\mathbf{x}], \boldsymbol{\beta} - \boldsymbol{\beta}' \rangle| \\ &\leq \frac{n}{M} \max_{i \in [n]} \|\mathbf{P}_{[\pm p]+i} \mathbf{x}\|_2 \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2 \\ &=: L \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2 \end{aligned}$$

define the set  $\mathcal{N}_\varepsilon$  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  $\mathcal{E}_{\text{Lip}}$ , for every  $\boldsymbol{\beta}$  in  $\mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$ , there exists some  $\boldsymbol{\beta}_\varepsilon$  in  $\mathcal{N}_{\varepsilon,i}$  with  $\|\boldsymbol{\beta} - \boldsymbol{\beta}_\varepsilon\|_2 \leq \varepsilon$ , thus

$$(SM9.34) \quad \left| \left( \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}) - n\mathbb{E}\bar{\mathbf{X}}(\boldsymbol{\beta}) \right) - \left( \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}_\varepsilon) - n\mathbb{E}\bar{\mathbf{X}}(\boldsymbol{\beta}_\varepsilon) \right) \right| \leq 2L\varepsilon \leq \frac{2cn\theta^2}{p \log n}.$$

On event  $\mathcal{E}_{\text{Lip}} \cap \mathcal{E}_{\text{Net}}$ , from (SM9.31), (SM9.33) and (SM9.34), we can conclude that for every  $\boldsymbol{\beta} \in \mathfrak{R}(\mathcal{S}_\tau, \gamma(c_\mu))$  and  $i \in [n]$ ,

$$(SM9.35) \quad \sum_{i \in [n]} \bar{\mathbf{X}}_i(\boldsymbol{\beta}) \leq \frac{24cn\theta^2}{p \log n}$$

as desired, where the error probability of  $\mathcal{E}_{\text{Lip}}^c$  is bounded using Lemma SM1.3, which gives

$$(SM9.36) \quad \mathbb{P}[\mathcal{E}_{\text{Lip}}^c] \leq \mathbb{P} \left[ \max_{j \in [n]} \|\mathbf{x}_{[\pm p]+j}\|_2 > 3\sqrt{p\theta \log n \log \theta^{-1}} \right] \leq 2/n,$$

4. (Bound  $\mathbb{P}[\mathcal{E}_{\text{Net}}^c]$ ) Wlog let us assume that  $2p$  divides  $n$ . By applying union bound and observing that  $\bar{\mathbf{X}}_i(\boldsymbol{\beta})$  is independent of  $\bar{\mathbf{X}}_{i+2p}(\boldsymbol{\beta})$  for any  $i \in [n]$ , we split  $\sum_i \bar{\mathbf{X}}_i(\boldsymbol{\beta})$  into  $n/2p$  independent sums of r.v.s, we have

$$\mathbb{P}[\mathcal{E}_{\text{Net}}^c] \leq 2p |\mathcal{N}_\varepsilon| \cdot \mathbb{P} \left[ \sum_{j=0}^{n/2p-1} (\bar{\mathbf{X}}_{2pj}(\boldsymbol{\beta}) - \mathbb{E}[\bar{\mathbf{X}}(\boldsymbol{\beta})]) > \frac{9cn\theta^2}{p^2 \log n} \right],$$

where each summand has bounded variance and  $L^q$ -norm derived similarly as its expectation such that

$$\mathbb{E}\bar{\mathbf{X}}_i(\boldsymbol{\beta})^q \leq 2 \cdot \mathbb{P}[\langle s_i[\mathbf{x}], \boldsymbol{\beta} \rangle \in [\lambda - B - M, \lambda + B + M]] \leq 2 \cdot \frac{1}{2\lambda} \cdot 2(B + M) \leq \frac{4c\theta^2}{p \log n},$$

and apply Bernstein inequality [Lemma SM10.4](#) with  $(\sigma^2, R) = (4c\theta^2 / (p \log n), 1)$ , obtains

$$\begin{aligned} \mathbb{P} \left[ \sum_{j=0}^{n/2p-1} (\bar{\mathbf{X}}_{2pj}(\boldsymbol{\beta}) - \mathbb{E} [\bar{\mathbf{X}}(\boldsymbol{\beta})]) > \frac{9cn\theta^2}{p^2 \log n} \right] &\leq \exp \left[ \frac{-(9cn\theta^2/p^2 \log n)^2}{2cn\theta^2/p^2 \log n + 2(9cn\theta^2/p^2 \log n)} \right] \\ &\leq \exp \left[ \frac{-4cn\theta^2}{p^2 \log n} \right], \end{aligned}$$

thus when  $n = Cp^5\theta^{-2} \log p$ :

$$(SM9.37) \quad \mathbb{P} [\mathcal{E}_{\text{Net}}^c] \leq \exp \left[ \log(2p) + 13p \log \left( \frac{2p \log n}{c \cdot c_\lambda} \right) - \frac{4cn\theta^2}{p^2 \log n} \right] \leq 1/n \quad \blacksquare$$

as long as  $\frac{C}{\log C} > 10^5 / (c^2 \cdot c_\lambda)$ .

### SM10. Tools.

**Lemma SM10.1 (Tail bound for Gaussian r.v.).** *If  $X \sim \mathcal{N}(0, \sigma^2)$ , then its tail bound for  $t > 0$  can be*

$$(SM10.1) \quad \mathbb{P} [X > t] \leq \frac{\sigma}{t\sqrt{2\pi}} \exp \left( -\frac{t^2}{2\sigma^2} \right)$$

**Lemma SM10.2 (Moments of the Gaussian random variables).** *If  $X \sim \mathcal{N}(0, \sigma^2)$ , then for all integer  $p \geq 1$ ,*

$$(SM10.2) \quad \mathbb{E} [|X|^p] \leq \sigma^p (p-1)!!.$$

**Lemma SM10.3 (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  $t > 0$ ,*

$$(SM10.3) \quad \mathbb{P} [|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| \geq t] \leq 2 \exp \left( -\frac{t^2}{2L^2} \right).$$

**Lemma SM10.4 (Moment control Bernstein inequality for scalar r.v.s).** ([\[SM4\]](#), 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  $\sigma^2$  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*

$$(SM10.4) \quad \mathbb{P} [|S - \mathbb{E} [S]| \geq t] \leq 2 \exp \left( -\frac{t^2}{2n\sigma^2 + 2Rt} \right).$$

**Lemma SM10.5 ( $\varepsilon$ -net on sphere).** [\[SM6\]](#) *Let  $(X, d)$  be a metric space and let  $\varepsilon > 0$ . A subset  $\mathcal{N}_\varepsilon$  of  $X$  is called an  $\varepsilon$ -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  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  for the sphere  $\mathbb{S}^{n-1}$  of size  $|\mathcal{N}_\varepsilon| \leq (3/\varepsilon)^n$ .*

**Lemma SM10.6 (Hanson-Wright).** [SM5] Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be independent, subgaussian random variables with subgaussian norm  $\sup_{p \geq 1} p^{-1/2} (\mathbb{E} |x_i^p|)^{1/p} \leq \sigma$ . Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , then for every  $t > 0$ ,

$$(SM10.5) \quad \mathbb{P} [|\mathbf{x}^* \mathbf{A} \mathbf{x} - \mathbb{E} \mathbf{x}^* \mathbf{A} \mathbf{x}| \geq t] \leq 2 \exp \left( -c \min \left( \frac{t^2}{64 \sigma^4 \|\mathbf{A}\|_F^2}, \frac{t}{8\sqrt{2} \sigma^2 \|\mathbf{A}\|_2} \right) \right).$$

**Lemma SM10.7 (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  $0 < N \leq nL$ ,

$$(SM10.6) \quad \max_{0 \leq \mathbf{x} \leq L, \|\mathbf{x}\|_1 \leq N} \sum_{i=1}^n f(\mathbf{x}_i) \leq N \left( 1 - \frac{s(L)}{L} \right)$$

*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:

$$\begin{aligned} \sum_{i=1}^n f(\mathbf{x}_{*i}) &= \lfloor \frac{N}{L} \rfloor f(L) + f(r) = \frac{N-r}{L} (L - s(L)) + (r - s(r)) \\ &= N \left( 1 - \frac{s(L)}{L} \right) + r \left( \frac{s(L)}{L} - \frac{s(r)}{r} \right) \leq N \left( 1 - \frac{s(L)}{L} \right) \end{aligned} \quad \blacksquare$$

## REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2009.
- [2] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer Publishing Company, Incorporated, 1st ed., 2011.
- [3] E. J. CANDÈS, M. B. WAKIN, AND S. P. BOYD, *Enhancing sparsity by reweighted  $\ell_1$  minimization*, Journal of Fourier analysis and applications, 14 (2008), pp. 877–905.
- [4] S. FOUCART AND H. RAUHUT, *A Mathematical Introduction to Compressive Sensing*, Springer, 2013.
- [5] M. RUDELSON, R. VERSHYNIN, ET AL., *Hanson-wright inequality and sub-gaussian concentration*, Electronic Communications in Probability, 18 (2013).
- [6] R. VERSHYNIN, *Introduction to the non-asymptotic analysis of random matrices*, arXiv preprint arXiv:1011.3027, (2010).