# 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, nonconvex geometry, nonconvex optimization

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

**DOI.** 10.1137/19M1237569

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  $\mathbf{a}_0 \in \mathbb{R}^n$  ( $n \gg p_0$ ) 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) 
$$\left(\alpha \, s_{\ell}[\boldsymbol{a}_0]\right) \, * \, \left(\frac{1}{\alpha} \, s_{-\ell}[\boldsymbol{x}_0]\right) \, = \, \boldsymbol{y}.$$

Because of this symmetry, we only expect to recover  $a_0$  and  $x_0$  up to a signed shift (see

<sup>\*</sup>Received by the editors January 8, 2019; accepted for publication (in revised form) December 2, 2019; published electronically February 27, 2020.

https://doi.org/10.1137/19M1237569

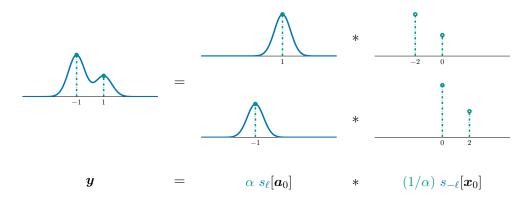
Funding: This work was funded by National Science Foundation grants NSF 1343282, NSF CCF 1527809, and NSF IIS 1546411.

<sup>&</sup>lt;sup>†</sup>Department of Electronic Engineering and Data Science Institute, Columbia University, New York, NY 10027 USA (hk2673@columbia.edu, yl3027@columbia.edu, jw2966@columbia.edu).

<sup>&</sup>lt;sup>‡</sup>Department of Computer Science, Cornell University, Ithaca, NY 14853 USA (yz2409@columbia.edu).

<sup>§</sup>Department of Applied Physics and Applied Mathematics, Columbia University, New York, NY 10027 USA

<sup>&</sup>lt;sup>1</sup>In this paper, the cyclic convolution  $a_0 * x_0$  assumes  $a_0$  to be zero-padded  $[a_0, 0^{n-p_0}]$  to length n.



**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 shift  $\ell$  samples. The ground truth signals  $\mathbf{a}_0$  and  $\mathbf{x}_0$  can only be identified up to a scaled shift.

Figure 1). Our problem of interest can be stated more formally as follows.

Problem 1.1 (short-and-sparse deconvolution). Given the cyclic convolution<sup>2</sup>  $y = a_0 * x_0 \in \mathbb{R}^n$  of  $a_0 \in \mathbb{R}^{p_0}$  short  $(p_0 \ll n)$  and  $x_0 \in \mathbb{R}^n$  sparse, recover  $a_0$  and  $x_0$ , up to a scaled shift.

Despite a long history and many applications, until recently very little algorithmic theory was available for SaS deconvolution. Much of this difficulty can be attributed to the scale-shift symmetry: natural convex relaxations fail,<sup>3</sup> and nonconvex formulations exhibit a complicated optimization landscape, with many equivalent global minimizers (scaled shifts of the ground truth), additional local minimizers (scaled shift truncations of the ground truth), and a variety of critical points [63, 64]. Currently available theory guarantees approximate recovery of a truncation<sup>4</sup> of a shift  $s_{\ell}[a_0]$ , rather than guaranteeing recovery of  $a_0$  as a whole, and requires certain (complicated) conditions on the convolution matrix associated with  $a_0$  [63].

In this paper, we describe an algorithm which, under simpler conditions, exactly recovers a scaled shift of the pair  $(a_0, x_0)$ . Our algorithm is based on a formulation first introduced in [64], which casts the deconvolution problem as (nonconvex) optimization over the sphere. We characterize the geometry of this objective function and show that near a certain union of subspaces, every local minimizer is very close to a signed shift of  $a_0$ . Based on this geometric analysis, we give provable methods for SaS deconvolution that exactly recover a scaled shift of  $(a_0, x_0)$  whenever  $a_0$  is shift-incoherent and  $x_0$  is a sufficiently sparse random vector. Our geometric analysis highlights the role of symmetry in shaping the objective landscape for SaS deconvolution.

The remainder of this paper is organized as follows. Section 2 introduces our optimization

<sup>&</sup>lt;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  $\mathbf{a}_0 \in \mathbb{R}^n$ ; then  $\mathbf{y}$  can also be expressed as circular convolution between  $\mathbf{a}_0$  and  $[\mathbf{x}_0; \mathbf{0}^{p_0-1}]$ .

<sup>3</sup>Such as matrix lifting relaxation [2, 39], in which  $\mathbf{a}_0$  or  $\mathbf{x}_0$  resides in random subspaces without shift symmetry.

<sup>&</sup>lt;sup>4</sup>That is, the portion of the shifted signal  $s_{\ell}[a_0]$  that falls in the window  $\{0,\ldots,p_0-1\}$ .

approach and modeling assumptions. Section 3 introduces our main results—both geometric and algorithmic—and compares them to the literature. Sections 4 and 5 describes the main ideas of our analysis. Section 6 demonstrates the experimental performance of the analyzed algorithm. Finally, section 7 discusses two main limitations of our analysis and describes directions for future work.

## 2. Formulation and assumptions.

2.1. Nonconvex SaS over the sphere. Our starting point is the (natural) formulation

(2.1) 
$$\min_{\boldsymbol{a}, \boldsymbol{x}} \ \frac{1}{2} \| \boldsymbol{a} * \boldsymbol{x} - \boldsymbol{y} \|_{2}^{2} + \lambda \| \boldsymbol{x} \|_{1} \quad \text{s.t.} \quad \| \boldsymbol{a} \|_{2} = 1.$$

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

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

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

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

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

The problem (2.3) is defined in terms of the optimal Lasso cost. This function is challenging to analyze, especially far away from  $a_0$ . The article [64] analyzes the local minima of a simplification of (2.3), obtained by approximating<sup>5</sup> the data fidelity term as

(2.4) 
$$\frac{1}{2} \|\boldsymbol{a} * \boldsymbol{x} - \boldsymbol{y}\|_{2}^{2} = \frac{1}{2} \|\boldsymbol{a} * \boldsymbol{x}\|_{2}^{2} - \langle \boldsymbol{a} * \boldsymbol{x}, \boldsymbol{y} \rangle + \frac{1}{2} \|\boldsymbol{y}\|_{2}^{2} \\ \approx \frac{1}{2} \|\boldsymbol{x}\|_{2}^{2} - \langle \boldsymbol{a} * \boldsymbol{x}, \boldsymbol{y} \rangle + \frac{1}{2} \|\boldsymbol{y}\|_{2}^{2}.$$

This yields a simpler objective function,

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

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

(2.6) 
$$\rho(\boldsymbol{x}) = \sum_{i} (\boldsymbol{x}_{i}^{2} + \delta^{2})^{1/2}.$$

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

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

<sup>&</sup>lt;sup>7</sup>This particular surrogate is sometimes called the pseudo-Huber function.

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

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

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

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

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

**2.2.** Analysis setting and assumptions. For convenience, we assume that  $a_0$  has unit  $\ell^2$  norm, i.e.,  $a_0 \in \mathbb{S}^{p_0-1}$ . Our analysis makes two main assumptions, on the short motif  $a_0$  and the sparse map  $x_0$ , respectively:

The first is that distinct shifts  $a_0$  have small inner product. We define the *shift coherence* of  $\mu(a_0)$  to be the largest inner product between distinct shifts:

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

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

- Spiky. When  $a_0$  is close to the Dirac delta  $\delta_0$ , the shift coherence  $\mu(a_0) \approx 0.9$  Here, the observed signal y consists of a superposition of sharp pulses. This is arguably the easiest instance of SaS deconvolution.
- Generic. If  $\mathbf{a}_0$  is chosen uniformly at random from the sphere  $\mathbb{S}^{p_0-1}$ , its coherence is bounded as  $\mu(\mathbf{a}_0) \lesssim \sqrt{1/p_0}$  with high probability.
- Tapered generic low-pass. Here,  $\mathbf{a}_0$  is generated by taking a random conjugate symmetric superposition of the first L length- $p_0$  discrete Fourier transform (DFT) basis signals, windowing (e.g., with a Hamming window) and normalizing to unit  $\ell^2$  norm. When  $L = p_0 \sqrt{1-\beta}$ , with high probability  $\mu(\mathbf{a}_0) \lesssim \beta$ . In this model,  $\mu$  does not have to diminish as  $p_0$  grows—it can be a fixed constant.<sup>10</sup>

<sup>&</sup>lt;sup>8</sup>This is purely a technical convenience. Our theory guarantees recovery of a signed shift  $(\pm s_{\ell}[\boldsymbol{a}_0], \pm s_{-\ell}[\boldsymbol{x}_0])$  of the truth. If  $\boldsymbol{a}_0$  does not have unit norm, identical reasoning implies that our method recovers a scaled shift  $(\alpha s_{\ell}[\boldsymbol{a}_0], \alpha^{-1} s_{-\ell}[\boldsymbol{x}_0])$  with  $\alpha = \pm \frac{1}{\|\boldsymbol{a}_0\|_{\infty}}$ .

<sup>&</sup>lt;sup>9</sup>The use of "≈" here suppresses constant and logarithmic factors.

<sup>&</sup>lt;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].

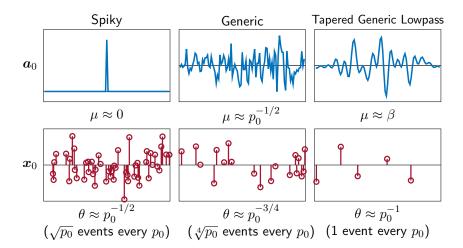


Figure 2. Sparsity-coherence tradeoff. Top: three families of motifs  $\mathbf{a}_0$  with varying coherence  $\mu$ . Bottom: maximum allowable sparsity  $\theta$  and number of copies  $\theta p_0$  within each length- $p_0$  window. Here, we suppress constants and logarithmic factors. When the target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

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

We assume that  $x_0$  is a sparse random vector. More precisely, we assume that  $x_0$  is Bernoulli–Gaussian, with rate  $\theta$ :

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

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

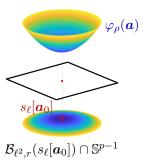
$$(2.11) x_0 \sim_{\text{i.i.d.}} BG(\theta).$$

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

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

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

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



**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  $s_{\ell}[\mathbf{a}_0]$ .

- 3. Main results: Geometry and algorithms. In this section, we introduce our main results—on the geometry of  $\varphi_{\rho}$  (subsection 3.1) and its algorithmic implications (subsection 3.2). Finally, in subsection 3.3, we compare these results with the literature on deconvolution.
- 3.1. Geometry of the objective  $\varphi_{\rho}$ . The goal in SaS deconvolution is to recover  $a_0$  (and  $x_0$ ) up to a signed shift; i.e., we wish to recover some  $\pm s_{\ell}[a_0]$ . The shifts  $\pm s_{\ell}[a_0]$  play a key role in shaping the landscape of  $\varphi_{\rho}$ . In particular, we will argue that over a certain subset of the sphere, every local minimum of  $\varphi_{\rho}$  is close to some  $\pm s_{\ell}[a_0]$ .

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

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

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

- The objective function  $\varphi_{\rho}$  is strongly convex in this neighborhood of  $s_{\ell}[a_0]$ .
- There is a local minimizer very close to  $s_{\ell}[\boldsymbol{a}_0]$ .

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

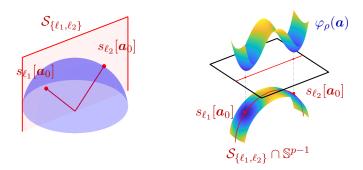
$$\mathcal{S}_{\{\ell_1,\ell_2\}} = \{ \boldsymbol{\alpha}_1 s_{\ell_1}[\boldsymbol{a}_0] + \boldsymbol{\alpha}_2 s_{\ell_2}[\boldsymbol{a}_0] \mid \boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2 \in \mathbb{R} \}.$$

We make three observations:

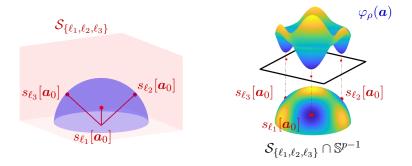
- Again, there is a local minimizer near each shift  $s_{\ell}[a_0]$ .
- These are the *only* local minimizers in the vicinity of  $S_{\{\ell_1,\ell_2\}}$ . In particular, the objective function  $\varphi$  exhibits *negative curvature* along  $S_{\{\ell_1,\ell_2\}}$  at any superposition  $\alpha_1 s_{\ell_1}[a_0] + \alpha_2 s_{\ell_2}[a_0]$  whose weights  $\alpha_1$  and  $\alpha_2$  are balanced, i.e.,  $|\alpha_1| \approx |\alpha_2|$ .
- Furthermore, the function  $\varphi_{\rho}$  exhibits *positive curvature* in directions away from the subspace  $\mathcal{S}_{\ell_1,\ell_2}$ .

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

$$\mathcal{S}_{\{\ell_1,\ell_2,\ell_3\}} = \left\{ \; \boldsymbol{\alpha}_1 s_{\ell_1}[\boldsymbol{a}_0] + \boldsymbol{\alpha}_2 s_{\ell_2}[\boldsymbol{a}_0] + \boldsymbol{\alpha}_3 s_{\ell_3}[\boldsymbol{a}_0] \; | \; \boldsymbol{\alpha}_1,\boldsymbol{\alpha}_2,\boldsymbol{\alpha}_3 \in \mathbb{R} \; \right\}.$$



**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  $s_{\ell_1}[\mathbf{a}_0]$  and  $s_{\ell_2}[\mathbf{a}_0]$ , 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  $\mathcal{S}_{\ell_1,\ell_2,\ell_3}$ .

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

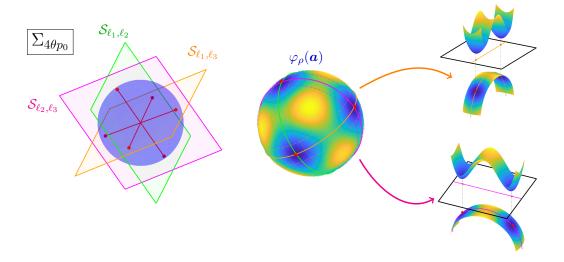
Our main geometric result will show that these properties are obtained from *every* subspace spanned by a few shifts of  $a_0$ . Indeed, for each subset

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

define a linear subspace

(3.2) 
$$S_{\tau} = \left\{ \sum_{\ell \in \tau} \alpha_{\ell} s_{\ell}[\boldsymbol{a}_{0}] \,\middle|\, \boldsymbol{\alpha}_{-p_{0}+1}, \dots, \boldsymbol{\alpha}_{p_{0}-1} \in \mathbb{R} \right\}.$$

The subspace  $S_{\tau}$  is the linear span of the shifts  $s_{\ell}[a_0]$  indexed by  $\ell$  in the set  $\tau$ . Our geometric theory will show that with high probability the function  $\varphi_{\rho}$  has no spurious local minimizers



**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.

near any  $S_{\tau}$  for which  $\tau$  is not too large—say,  $|\tau| \leq 4\theta p_0$ . Combining all of these subspaces into a single geometric object, define the union of subspaces

$$\Sigma_{4\theta p_0} = \bigcup_{|\tau| < 4\theta p_0} \mathcal{S}_{\tau}.$$

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

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

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

Theorem 3.1 (main geometric theorem). Let  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$  with  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$   $\mu$ -shift coherent<sup>11</sup> and  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \mathrm{BG}(\theta) \in \mathbb{R}^n$  with sparsity rate

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

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

<sup>&</sup>lt;sup>11</sup>Typically it is possible to provide an overestimate  $p'_0 \ge p_0$ . 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  $p_0$ .

*Proof.* This follows from Theorem 4.1.

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

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

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

The next two paragraphs give the main ideas behind the principal steps of the algorithm. We then describe its components in more detail (Algorithm 3.1) and state our main algorithmic result (Theorem 3.2), which asserts that under appropriate conditions this method produces a signed shift of  $(a_0, x_0)$ .

Our algorithm starts with an initialization scheme which generates  $a^{(0)}$  near the union of subspaces  $\Sigma_{4\theta p_0}$ , which consists of linear combinations of just a few shifts of  $a_0$ . How can we find a point near this union? Notice that the data y also consists of a linear combination of just a few shifts of  $a_0$  Indeed,

$$(3.5) y = a_0 * x_0 = \sum_{\ell \in \text{supp}(\boldsymbol{x}_0)} x_{0\ell} s_{\ell}[a_0].$$

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

Many of these copies of  $a_0$  are truncated by the restriction to  $\{0, \ldots, p_0 - 1\}$ . A relatively simple remedy is as follows: First, we zero-pad  $y_{0,\ldots,p_0-1}$  to length  $p = 3p_0 - 2$ , giving

$$[\mathbf{0}^{p_0-1}; \mathbf{y}_0; \cdots; \mathbf{y}_{p_0-1}; \mathbf{0}^{p_0-1}].$$

Zero-padding provides enough space to accommodate any shift  $s_{\ell}[a_0]$  with  $\ell \in \tau$ . We then

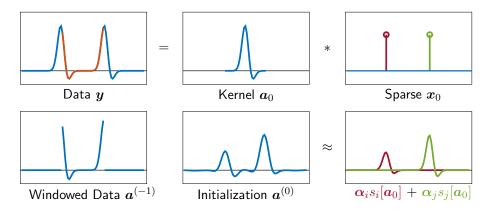


Figure 7. Data-driven initialization. Using a piece of the observed data  $\mathbf{y}$  to generate an initial point  $\mathbf{a}^{(0)}$  that is close to a superposition of shifts  $s_{\ell}[\mathbf{a}_0]$  of the ground truth. Top: data  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$  is a superposition 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).

perform one step of the generalized power method, 12 writing

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

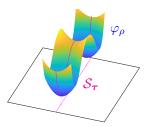
where  $P_{\mathbb{S}^{p-1}}$  projects onto the sphere. The reasoning behind this construction may seem obscure. We will explain it at a more technical level in section 5 after interpreting the gradient  $\nabla \varphi_{\rho}$  in terms of its action on the shifts  $s_{\ell}[\mathbf{a}_0]$  in section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

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

(3.8) 
$$a^{(k+1)} \leftarrow P_{\mathbb{S}^{p-1}} [a^{(k)} - tg^{(k)} - t^2 v^{(k)}]$$

with appropriately chosen step size t. The inclusion of a negative curvature direction allows the method to avoid stagnation near saddle points. Indeed, we will prove that, starting from initialization  $a^{(0)}$ , this method produces a sequence  $a^{(1)}, a^{(2)}, \ldots$  which efficiently converges to a local minimizer  $\bar{a}$  that is near some signed shift  $\pm s_{\ell}[a_0]$  of the ground truth.

<sup>&</sup>lt;sup>12</sup>The power method for minimizing a quadratic form  $\xi(\boldsymbol{a}) = \frac{1}{2}\boldsymbol{a}^*\boldsymbol{M}\boldsymbol{a}$  over the sphere consists of the iteration  $\boldsymbol{a} \mapsto -\boldsymbol{P}_{\mathbb{S}^{p-1}}\boldsymbol{M}\boldsymbol{a}$ . Notice that in this mapping,  $-\boldsymbol{M}\boldsymbol{a} = -\nabla\xi(\boldsymbol{a})$ . The generalized power method, for minimizing a function  $\varphi$  over the sphere, consists of repeatedly projecting  $-\nabla\varphi$  onto the sphere, giving the iteration  $\boldsymbol{a} \mapsto -\boldsymbol{P}_{\mathbb{S}^{p-1}}\nabla\varphi(\boldsymbol{a})$ . Equation (3.7) can be interpreted as one step of the generalized power method for the objective function  $\varphi_{\rho}$ .



**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  $S_{\tau}$ .

The second step of our algorithm rounds the local minimizer  $\bar{a} \approx \sigma s_{\ell}[a_0]$  to produce an exact solution  $\hat{a} = \sigma s_{\ell}[a_0]$ . As a by-product, it also exactly recovers the corresponding signed shift of the true sparse signal,  $\hat{x} = \sigma s_{-\ell}[x_0]$ .

Our rounding algorithm is an alternating minimization scheme, which alternates between minimizing the Lasso cost over a with x fixed, and minimizing the Lasso cost over x with afixed. We make two modifications to this basic idea, both of which are important for obtaining exact recovery. First, unlike the standard Lasso cost, which penalizes all of the entries of x, we maintain a running estimate  $I^{(k)}$  of the support of  $x_0$  and only penalize those entries that are not in  $I^{(k)}$ :

(3.9) 
$$\frac{1}{2} \|\boldsymbol{a} * \boldsymbol{x} - \boldsymbol{y}\|_{2}^{2} + \lambda \sum_{i \neq I^{(k)}} |\boldsymbol{x}_{i}|.$$

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

$$\widehat{\boldsymbol{a}} * \widehat{\boldsymbol{x}} \approx \boldsymbol{y}.$$

This can be viewed as a homotopy or continuation method [46, 19]. For concreteness, at the kth iteration the algorithm reads

(3.11) Update 
$$\boldsymbol{x}$$
:  $\boldsymbol{x}^{(k+1)} \leftarrow \underset{\boldsymbol{x}}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{a}^{(k)} * \boldsymbol{x} - \boldsymbol{y}\|_{2}^{2} + \lambda^{(k)} \sum_{i \notin I^{(k)}} |\boldsymbol{x}_{i}|;$ 
(3.12) Update  $\boldsymbol{a}$ :  $\boldsymbol{a}^{(k+1)} \leftarrow \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[\underset{\boldsymbol{a}}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{a} * \boldsymbol{x}^{(k+1)} - \boldsymbol{y}\|_{2}^{2}\right];$ 
(3.13) Update  $\lambda$  and  $I$ :  $\lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \operatorname{supp} \left(\boldsymbol{x}^{(k+1)}\right).$ 

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

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

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

Our overall algorithm is summarized as Algorithm 3.1. Figure 9 illustrates the main steps of this algorithm. Our main algorithmic result states that under essentially the same hypotheses as above, Algorithm 3.1 produces a signed shift of the ground truth  $(a_0, x_0)$ .

### **Algorithm 3.1.** Short-and-sparse deconvolution.

**Input:** Observation  $\boldsymbol{y}$ , motif length  $p_0$ , sparsity  $\theta$ , shift-coherence  $\mu$ , and curvature threshold  $-\eta_v$ .

## Minimization:

Set 
$$\boldsymbol{a}^{(0)} \leftarrow -\boldsymbol{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\rho} \left( \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \boldsymbol{0}^{p_0-1}; \boldsymbol{y}_0; \cdots; \boldsymbol{y}_{p_0-1}; \boldsymbol{0}^{p_0-1} \right] \right)$$
.  
Set  $\lambda = 0.1/\sqrt{p_0 \theta}^{13}$  and  $\delta > 0$  in  $\varphi_{\rho}$ . For  $k = 1, 2, \dots, K_1$ , let
$$\boldsymbol{a}^{(k+1)} \leftarrow \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \boldsymbol{a}^{(k)} - t \boldsymbol{a}^{(k)} - t^2 \boldsymbol{v}^{(k)} \right],$$

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

(3.15) 
$$\varphi_{\rho}(\boldsymbol{a}^{(k+1)}) < \varphi_{\rho}(\boldsymbol{a}^{(k)}) - \frac{1}{2}t\|\boldsymbol{g}^{(k)}\|_{2}^{2} - \frac{1}{4}t^{4}\eta_{v}\|\boldsymbol{v}^{(k)}\|_{2}^{2}$$

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

#### Refinement:

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

(3.16) 
$$x^{(k+1)} \leftarrow \operatorname{argmin}_{x} \frac{1}{2} \|a^{(k)} * x - y\|_{2}^{2} + \lambda^{(k)} \sum_{i \notin I^{(k)}} |x_{i}|,$$

(3.17) 
$$\boldsymbol{a}^{(k+1)} \leftarrow \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \operatorname{argmin}_{\boldsymbol{a}} \frac{1}{2} \| \boldsymbol{a} * \boldsymbol{x}^{(k+1)} - \boldsymbol{y} \|_{2}^{2} \right],$$

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

to obtain  $(\widehat{\boldsymbol{a}}, \widehat{\boldsymbol{x}}) \leftarrow (\boldsymbol{a}^{(K_2)}, \boldsymbol{x}^{(K_2)}).$ 

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

Theorem 3.2 (main algorithmic theorem). Suppose  $\mathbf{y} = \mathbf{a}_0 * \mathbf{x}_0$ , where  $\mathbf{a}_0 \in \mathbb{S}^{p_0-1}$  is  $\mu$ -truncated shift coherent<sup>14</sup> such that  $\max_{i \neq j} \left| \left\langle \boldsymbol{\iota}_{p_0}^* s_i[\mathbf{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\mathbf{a}_0] \right\rangle \right| \leq \mu$  and  $\mathbf{x}_0 \sim_{\text{i.i.d.}} \mathrm{BG}(\theta) \in \mathbb{R}^n$  with  $\theta, \mu$  satisfying

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

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

(3.20) 
$$\|(\widehat{\boldsymbol{a}}, \widehat{\boldsymbol{x}}) - \sigma(s_{\ell}[\boldsymbol{a}_0], s_{-\ell}[\boldsymbol{x}_0])\|_2 \le \varepsilon$$

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

*Proof.* See Theorems 5.1 and 5.2.

<sup>&</sup>lt;sup>13</sup>In practice, we suggest setting  $\lambda = c_{\lambda}/\sqrt{p_0\theta}$  with  $c_{\lambda} \in [0.5, 0.8]$ .

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

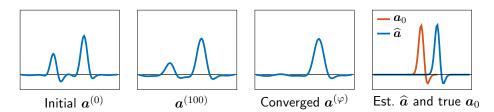


Figure 9. Local minimization and refinement. Left: data-driven initialization  $\mathbf{a}^{(0)}$  consisting of a near superposition of two shifts. Middle: minimizing  $\varphi_{\rho}$  produces a near shift of  $\mathbf{a}_0$ . Right: rounded solution  $\hat{\mathbf{a}}$  using the Lasso.  $\hat{\mathbf{a}}$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

When solving SaS deconvolution via minimizing Bilinear Lasso objective (2.2) in practice, the algorithm is analogous to the provable method introduced in Algorithm 3.1, where the curvilinear descent and the refinement step can be realized as alternating gradient descent of both variables a, x in (2.2). Unlike Algorithm 3.1, this alternating gradient method has yet to come with theoretical guarantees, but has shown to be an effective and efficient method for SaS deconvolution problems both in simulation and in reality [34].

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

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

Another significant application of SaS deconvolution is *image deblurring*. Typically, the blur kernel is small relative to the image size (*short*) [3, 62, 13, 35, 36]. In natural image deblurring, the target image is often assumed to have relatively few sharp edges [21, 27, 36], and hence have *sparse* derivatives. In scientific image deblurring, e.g., in astronomy [33, 25, 9] and geophysics [28], the target image is often sparse, either in the spatial or wavelet domains, again leading to variants of the SaS model. The literature on blind image deconvolution is large; see, e.g., [31, 10] for surveys.

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

These applications have motivated a great deal of algorithmic work on variants of the SaS problem [32, 8, 6, 31, 43, 10, 56]. In contrast, relatively little theory is available to explain when and why algorithms succeed. Our algorithm minimizes  $\varphi_{\rho}$  as an approximation to the Lasso cost over the sphere. Our formulation and results have strong precedent in the literature. Lasso-like objective functions have been widely used in image deblurring [62, 14, 21, 35, 52, 60, 18, 30, 36, 59, 47, 64]. A number of insights have been obtained into the

geometry of sparse deconvolution—in particular, into the effect of various constraints on  $\boldsymbol{a}$  on the presence or absence of spurious local minimizers. In image deblurring, a simplex constraint  $(\boldsymbol{a} \geq \boldsymbol{0} \text{ and } \|\boldsymbol{a}\|_1 = 1)$  arises naturally from the physical structure of the problem [62, 14]. Perhaps surprisingly, simplex-constrained deconvolution admits trivial global minimizers, at which the recovered kernel  $\boldsymbol{a}$  is a spike, rather than the target blur kernel [7, 36].

The work [59] imposes the  $\ell^2$  regularization on  $\boldsymbol{a}$  and observes that this alternative constraint gives a more reliable algorithm. In [64], the geometry of the simplified objective  $\varphi_{\ell^1}$  over the sphere is studied, and it is proved that in the dilute limit in which  $\boldsymbol{x}_0$  has one nonzero entry, all strict local minima of  $\varphi_{\ell^1}$  are close to signed shifts truncations of  $\boldsymbol{a}_0$ . By adopting a different objective function (based on  $\ell^4$  maximization) over the sphere, [63] proves that on a certain region of the sphere every local minimum is near a truncated signed shift of  $\boldsymbol{a}_0$ , i.e., the restriction of  $s_{\ell}[\boldsymbol{a}_0]$  to the window  $\{0,\ldots,p_0-1\}$ . The analysis of [63] allows the sparse sequence  $\boldsymbol{x}_0$  to be denser ( $\theta \sim p_0^{-2/3}$  for a generic kernel  $\boldsymbol{a}_0$ , as opposed to  $\theta \lesssim p_0^{-3/4}$  in our result). Both [64] and [63] guarantee approximate recovery of a portion of  $s_{\ell}[\boldsymbol{a}_0]$ , under complicated conditions on the kernel  $\boldsymbol{a}_0$ . Our core optimization problem is very similar to that of [64]. However, we obtain exact recovery of both  $\boldsymbol{a}_0$  and relatively dense  $\boldsymbol{x}_0$  under the much simpler assumption of shift incoherence.

Other aspects of the SaS problem have been studied theoretically. One basic question is under what circumstances the problem is identifiable up to the scaled shift ambiguity. The paper [17] shows that the problem is ill-posed for worst case  $(a_0, x_0)$ , particularly for certain support patterns in which  $x_0$  does not have any isolated nonzero entries. This demonstrates that *some* modeling assumptions on the support of the sparse term are needed. At the same time, this worst-case structure is unlikely to occur, either under the Bernoulli model or in practical deconvolution problems.

Motivated by a variety of applications, much research has focused on low-dimensional deconvolution models in the theoretical literature. In communication applications, the signals  $a_0$  and  $x_0$  either live in known low-dimensional subspaces or are sparse in some known dictionary [2, 16, 29, 39, 40, 41, 42]. These theoretical works assume that the subspace/dictionary are chosen at random. This low-dimensional deconvolution model does not exhibit the signed shift ambiguity; nonconvex formulations for this model exhibit a different structure from that studied here. In fact, the variant in which both signals belong to known subspaces can be solved by convex relaxation [2]. The SaS model does not appear to be amenable to convexification and exhibits a more complicated nonconvex geometry due to the shift ambiguity. The main motivation for tackling this model lies in the aforementioned applications in imaging and data analysis.

In [38, 57] the related multi-instance sparse blind deconvolution problem (MISBD) is studied, where there are K observations  $y_i = a_0 * x_i$  consisting of multiple convolutions i = 1, ..., K of a kernel  $a_0$  and different sparse vectors  $x_i$ . Both works develop provable algorithms. There are several key differences with our work. First, both the proposed algorithms and their analysis require the kernel to be invertible. Second, despite the apparent similarity between the SaS model and MISBD, these problems are not equivalent. It might seem possible to reduce SaS to MISBD by dividing the single observation y into K pieces; this apparent reduction fails due to boundary effects.

**3.4. Notation.** All vectors/matrices are written in bold font, a/A; indexed values are written as  $a_i$ ,  $A_{ij}$ . Zero or one vectors are defined as  $\mathbf{0}$  or  $\mathbf{1}$ , and ith canonical basis vector defined as  $e_i$ . The indices for vectors/matrices all start from 0 and are taken modulo-n, and thus a vector of length n should have its indices labeled as  $\{0,1,\ldots,n-1\}$ . We write  $[n] = \{0,\ldots,n-1\}$ . We often use the capital italic symbols I,J for subsets of [n]. We abuse notation slightly and write  $[-p] = \{n-p+1,\ldots,n-1,0\}$  and  $[\pm p] = \{n-p+1,\ldots,n-1,0,1,\ldots,p-1\}$ . Index sets can be labels for vectors;  $a_I \in \mathbb{R}^{|I|}$  denotes the restriction of the vector a to coordinates I. Also, we use a check symbol to denote the reversal operator on index set I = -I and vectors  $a = a_i$ .

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

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

$$(3.21) C_{a}b = a * b.$$

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

(3.22) 
$$C_a^*b = \widecheck{a} * b, \quad \widecheck{C}_ab = a * \widecheck{b}.$$

- 4. Geometry of  $\varphi_{\rho}$  in shift space. Underlying our main geometric and algorithmic results is a relationship between the geometry of the function  $\varphi_{\rho}$  and the symmetries of the deconvolution problem. In this section, we describe this relationship at a more technical level by interpreting the gradient and Hessian of the function  $\varphi_{\rho}$  in terms of the shifts  $s_{\ell}[a_0]$  and stating a key lemma which asserts that a certain neighborhood of the union of subspaces  $\Sigma_{4\theta p_0}$  can be decomposed into regions of negative curvature, strong gradient, and strong convexity near the target solutions  $\pm s_{\ell}[a_0]$ .
- **4.1. Shifts and correlations.** The set  $\Sigma_{4\theta p_0}$  is a union of subspaces. Any point  $\boldsymbol{a}$  in one of these subspaces  $\mathcal{S}_{\tau}$  is a superposition of shifts of  $\boldsymbol{a}_0$ :

(4.1) 
$$a = \sum_{\ell \in \tau} \alpha_{\ell} s_{\ell}[a_0].$$

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

(4.2) 
$$a = \sum_{\ell \in \tau} \alpha_{\ell} s_{\ell}[a_0] + \sum_{\ell \notin \tau} \alpha_{\ell} s_{\ell}[a_0].$$

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

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

$$\beta = C_{a_0}^* \iota a$$

$$= C_{a_0}^* \iota \iota^* C_{a_0} \alpha =: M\alpha.$$

The matrix M is the Gram matrix of the truncated shifts:  $M_{ij} = \langle \iota^* s_i[a_0], \iota^* s_j[a_0] \rangle$ . When  $\mu$  is small, the off-diagonal elements of M are small. In particular, on  $S_{\tau}$  we may take  $\alpha_{\tau^c} = 0$ , 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, see section SM2 in the supplementary material.

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

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

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

Here,  $S_{\lambda}$  is the soft thresholding operator, which is defined for scalars t as

(4.6) 
$$S_{\lambda}[t] = \operatorname{sign}(t) \max\{|t| - \lambda, 0\}$$

and is extended to vectors by applying it elementwise. The operator  $S_{\lambda}[x]$  shrinks the elements of x toward zero. Small elements become identically zero, resulting in a sparse vector.

**Gradient:** Sparsifying the correlations  $\beta$ . Our goal is to understand the local minimizers of the function  $\varphi_{\ell^1}$  over the sphere. The function  $\varphi_{\ell^1}$  is differentiable. Clearly, any point a at which its gradient (over the sphere) is nonzero cannot be a local minimizer. We first give an expression for the gradient of  $\varphi_{\ell^1}$  over Euclidean space  $\mathbb{R}^p$ , and then extend it to the sphere  $\mathbb{S}^{p-1}$ . Using  $y = a_0 * x_0$  and calculus gives

$$\nabla \varphi_{\ell^{1}}(\boldsymbol{a}) = -\iota^{*} C_{\boldsymbol{a}_{0}} \widecheck{C}_{\boldsymbol{x}_{0}} \mathcal{S}_{\lambda} \left[ \widecheck{C}_{\boldsymbol{x}_{0}} C_{\boldsymbol{a}_{0}}^{*} \iota \boldsymbol{a} \right]$$

$$= -\iota^{*} C_{\boldsymbol{a}_{0}} \widecheck{C}_{\boldsymbol{x}_{0}} \mathcal{S}_{\lambda} \left[ \widecheck{C}_{\boldsymbol{x}_{0}} \boldsymbol{\beta} \right]$$

$$= -\iota^{*} C_{\boldsymbol{a}_{0}} \chi[\boldsymbol{\beta}],$$

$$(4.7)$$

where we have simplified the notation by introducing an operator  $\chi : \mathbb{R}^n \to \mathbb{R}^n$  as  $\chi[\beta] = \widetilde{C}_{x_0} \mathcal{S}_{\lambda}[\widecheck{C}_{x_0}\beta]$ . This representation exhibits the (negative) gradient as a superposition of shifts of  $a_0$  with coefficients given by the entries of  $\chi[\beta]$ :

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

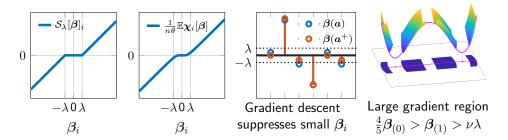


Figure 10. Gradient sparsifies correlations. Left: the soft thresholding operator  $S_{\lambda}[\beta]$  shrinks the entries of  $\beta$  toward zero, making it sparser. Middle left: the negative gradient  $-\nabla \varphi_{\ell^1}$  is a superposition of shifts  $s_{\ell}[\mathbf{a}_0]$ , with coefficients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

The operator  $\chi$  appears complicated. However, its effect is relatively simple: when  $x_0$  is a long random vector,  $\chi[\beta]$  acts like a soft thresholding operator on the vector  $\beta$ . That is,

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

We show this rigorously below in the proof of our main theorems. Here, we support this claim pictorially by plotting the  $\ell$ th entry  $\chi[\beta]_{\ell}$  as  $\beta_{\ell}$  varies; see Figure 10 (middle left) and compare to Figure 10 (left). Because  $\chi[\beta]$  suppresses small entries of  $\beta$ , the strongest contributions to  $-\nabla \varphi_{\ell^1}$  in (4.8) will come from shifts  $s_{\ell}[a_0]$  with large  $\beta_{\ell}$ . In particular, the Euclidean gradient is large whenever there is a single preferred shift  $s_{\ell}[a_0]$ , i.e., the largest entry of  $\beta$  is significantly larger than the second largest entry.

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

(4.10) 
$$\begin{aligned} \operatorname{grad}[\varphi_{\ell^{1}}](\boldsymbol{a}) &= \boldsymbol{P}_{\boldsymbol{a}^{\perp}} \nabla \varphi_{\ell^{1}}(\boldsymbol{a}) \\ &= -\boldsymbol{P}_{\boldsymbol{a}^{\perp}} \sum_{\ell} \boldsymbol{\chi}_{\ell}[\boldsymbol{\beta}] \, s_{\ell}[\boldsymbol{a}_{0}]. \end{aligned}$$

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

- (i) the negative gradient points to one particular shift: there is a single preferred shift  $s_{\ell}[\mathbf{a}_0]$  so that the Euclidean gradient is large; and
- (ii)  $\boldsymbol{a}$  is not too close to any shift: it is possible to move in the tangent space in the direction of this shift. Since the tangent space consists of those vectors orthogonal to  $\boldsymbol{a}$ , this is possible whenever  $s_{\ell}[\boldsymbol{a}_0]$  is not too aligned with  $\boldsymbol{a}$ , i.e.,  $\boldsymbol{a}$  is not too close to  $s_{\ell}[\boldsymbol{a}_0]$ .

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

Our technical lemma quantifies this situation in terms of the ordered entries of  $\boldsymbol{\beta}$ . Write  $|\boldsymbol{\beta}_{(0)}| \geq |\boldsymbol{\beta}_{(1)}| \geq \cdots$ , with corresponding shifts  $s_{(0)}[\boldsymbol{a}_0], s_{(1)}[\boldsymbol{a}_0], \ldots$ . There is a strong gradient whenever  $|\boldsymbol{\beta}_{(0)}|$  is significantly larger than  $|\boldsymbol{\beta}_{(1)}|$  and  $|\boldsymbol{\beta}_{(1)}|$  is not too small compared to  $\lambda$ : in particular, when  $\frac{4}{5}|\boldsymbol{\beta}_{(0)}| > |\boldsymbol{\beta}_{(1)}| > \frac{\lambda}{4\log^2\theta^{-1}}$ . In this situation, gradient descent drives  $\boldsymbol{a}$  toward  $s_{(0)}[\boldsymbol{a}_0]$ , reducing  $|\boldsymbol{\beta}_{(1)}|, \ldots$ , and making the vector  $\boldsymbol{\beta}$  sparser. We establish the technical claim that the (Euclidean) gradient of  $\varphi_{\ell^1}$  sparsifies vectors in shift space in section SM3.

Hessian: Negative curvature breaks symmetry. When there is no single preferred shift, i.e., when  $|\beta_{(1)}|$  is close to  $|\beta_{(0)}|$ , the gradient can be small. Similarly, when a is very close to  $\pm s_{(0)}[a_0]$ , the gradient can be small. In either of these situations, we need to study the curvature of the function  $\varphi$  to determine whether there are local minimizers.

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

(4.11) 
$$\widetilde{\nabla}^2 \varphi_{\ell^1}(\boldsymbol{a}) = -\iota^* \boldsymbol{C}_{\boldsymbol{a}_0} \widecheck{\boldsymbol{C}}_{\boldsymbol{x}_0} \boldsymbol{P}_I \widecheck{\boldsymbol{C}}_{\boldsymbol{x}_0} \boldsymbol{C}_{\boldsymbol{a}_0}^* \iota,$$

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

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

(4.12) 
$$\widetilde{\nabla}^2 \varphi_{\ell^1}(\boldsymbol{a}) \approx -\sum_{\ell} s_{\ell}[\boldsymbol{a}_0] s_{\ell}[\boldsymbol{a}_0]^* \left( \frac{\partial}{\partial \boldsymbol{\beta}_{\ell}} \boldsymbol{\chi}_{\ell}[\boldsymbol{\beta}] \right)$$

can be obtained from (4.8) by noting that  $\frac{\partial}{\partial a} \chi_{\ell}[\beta] = \sum_{j} s_{j}[a_{0}] \frac{\partial}{\partial \beta_{j}} \chi_{\ell}[\beta]$ , that  $\frac{\partial}{\partial \beta_{j}} \chi_{\ell}[\beta] \approx 0$  for  $j \neq \ell$ , and that

(4.13) 
$$\frac{1}{n\theta} \cdot \frac{\partial \mathbf{\chi}_{\ell}[\boldsymbol{\beta}]}{\partial \boldsymbol{\beta}_{\ell}} \approx \begin{cases} 0, & |\boldsymbol{\beta}_{\ell}| \ll \lambda, \\ 1, & |\boldsymbol{\beta}_{\ell}| \gg \lambda. \end{cases}$$

Again, we corroborate this approximation pictorially; see Figure 11.

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

The (Euclidean) Hessian measures the curvature of the function  $\varphi_{\ell^1}$  over  $\mathbb{R}^n$ . The Riemannian Hessian

$$(4.14) \qquad \widetilde{\operatorname{Hess}}[\varphi_{\ell^{1}}](\boldsymbol{a}) = \boldsymbol{P}_{\boldsymbol{a}^{\perp}} \begin{pmatrix} \widetilde{\nabla}^{2} \varphi_{\ell^{1}}(\boldsymbol{a}) & + & \langle -\nabla \varphi_{\ell^{1}}(\boldsymbol{a}), \boldsymbol{a} \rangle \cdot \boldsymbol{I} \\ \operatorname{Curvature of } \varphi_{\ell^{1}} & \operatorname{Curvature of the sphere} \end{pmatrix} \boldsymbol{P}_{\boldsymbol{a}^{\perp}}$$

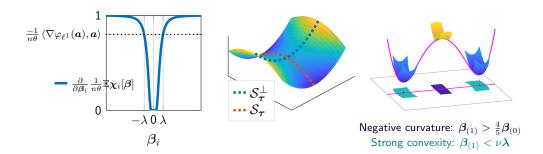


Figure 11. Hessian breaks symmetry. Left: contribution of  $-s_i[\mathbf{a}_0]s_i[\mathbf{a}_0]^*$  to the Euclidean Hessian. If  $|\beta_i| \gg \lambda$ , the Euclidean Hessian exhibits a strong negative component in the  $s_i[\mathbf{a}_0]$  direction. The Riemannian Hessian exhibits negative curvature 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  $S_{\tau}$  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.

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

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

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

of the corresponding shifts of  $a_0$ .

- (ii) Positive curvature in directions away from  $S_{\tau}$ : The Euclidean Hessian quadratic form  $v^*\widetilde{\nabla}^2\varphi_{\ell^1}v$  takes on relatively small values in directions orthogonal to the subspace  $S_{\tau}$ . The Riemannian Hessian is positive in these directions, creating positive curvature orthogonal to the subspace  $S_{\tau}$ .
- (iii) Strong convexity around minimizers: Around a minimizer  $s_{\ell}[a_0]$ , only a single entry  $\beta_{\ell}$  is large. Any tangent direction  $\mathbf{v} \perp \mathbf{a}$  is nearly orthogonal to the subspace span $(s_{\ell}[a_0])$ , and hence is a direction of positive (Riemannian) curvature. The objective function  $\varphi_{\rho}$  is strongly convex around the target solutions  $\pm s_{\ell}[a_0]$ .

Figure 11 visualizes these regions of negative and positive curvatures, and the technical claim of positivity/negativity of curvature in shift space is presented in detail in section SM4.

**4.3.** Any local minimizer is a near shift. We close this section by stating a key theorem, which makes the above discussion precise. We will show that a certain neighborhood of any subspace  $S_{\tau}$  can be covered by regions of negative curvature, of large gradient, and of strong convexity containing target solutions  $\pm s_{\ell}[a_0]$ . Furthermore, at the boundary of this neighborhood, the negative gradient points back—retracts—toward the subspace  $S_{\tau}$ , due to

the (directional) convexity of  $\varphi_{\rho}$  away from the subspace.

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

(4.15) 
$$\mathbf{a} = \sum_{\ell \in \tau} \alpha_{\ell} s_{\ell}[\mathbf{a}_{0}] + \sum_{\ell' \in \tau^{c}} \alpha_{\ell'} s_{\ell'}[\mathbf{a}_{0}].$$

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

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

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

(4.17) 
$$\Sigma_{4\theta p_0}^{\gamma} := \bigcup_{|\boldsymbol{\tau}| \le 4\theta p_0} \left\{ \boldsymbol{a} : d_{\alpha}(\boldsymbol{a}, \mathcal{S}_{\boldsymbol{\tau}}) \le \gamma \right\}.$$

The aforementioned geometric properties hold over this set.

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

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

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 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  $n > Cp_0^5 \theta^{-2} \log p_0$ , then with probability at least 1 - c'''/n, for every  $\mathbf{a} \in \Sigma_{4\theta p_0}^{\gamma}$ , we have the following:

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

$$\lambda_{\min} \left( \text{Hess}[\varphi_a](\boldsymbol{a}) \right) < -c_1 n \theta \lambda.$$

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

(4.20) 
$$\|\operatorname{grad}[\varphi_{\rho}](\boldsymbol{a})\|_{2} \geq c_{2}n\theta \frac{\lambda^{2}}{\log^{2}\theta^{-1}}.$$

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

(4.21) 
$$\operatorname{Hess}[\varphi_o](\boldsymbol{a}) \succ c_3 n \theta \boldsymbol{P}_{\boldsymbol{a}^{\perp}}.$$

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

$$\langle \boldsymbol{\zeta}_{\boldsymbol{\tau}^c}, \boldsymbol{\alpha}_{\boldsymbol{\tau}^c} \rangle \geq c_4 \| \boldsymbol{\zeta}_{\boldsymbol{\tau}^c} \|_2 \| \boldsymbol{\alpha}_{\boldsymbol{\tau}^c} \|_2.$$

(Local minimizers.) If a is a local minimizer,

(4.23) 
$$\min_{\substack{\ell \in [\pm p] \\ \sigma \in \{\pm 1\}}} \|\boldsymbol{a} - \sigma s_{\ell}[\boldsymbol{a}_{0}]\|_{2} \leq \frac{1}{2} \max \left\{\mu, p_{0}^{-1}\right\},$$

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

*Proof.* See subsection SM6.5.

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

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

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

- **5. Provable algorithm.** In light of Theorem 4.1, in this section we introduce a two-part algorithm, Algorithm 3.1, which first applies the curvilinear descent method to find a local minimum of  $\varphi_{\rho}$  within  $\Sigma_{4\theta p_0}^{\gamma}$ , followed by a refinement algorithm that uses alternating minimization to exactly recover the ground truth. This algorithm exactly solves the SaS deconvolution problem.
  - **5.1. Minimization.** There are three major issues in finding a local minimizer within  $\Sigma_{4\theta p_0}^{\gamma}$ :
    - (i) Initialization. The initializer  $a^{(0)}$  to reside within  $\Sigma_{4\theta p_0}^{\gamma}$
    - (ii) Negative curvature. The method to avoid stagnating near saddle points of  $\varphi_{\rho}$ .
    - (iii) No exit. The descent method to remain inside  $\Sigma_{4\theta p_0}^{\gamma}$ .

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

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

$$\begin{split} \boldsymbol{a}^{(0)} &= -\boldsymbol{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\rho} \left( \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \boldsymbol{0}^{p_0-1}; \boldsymbol{y}_0; \cdots; \boldsymbol{y}_{p_0-1}; \boldsymbol{0}^{p_0-1} \right] \right) \\ &= -\boldsymbol{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\rho} \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \boldsymbol{P}_{[p_0]} (\boldsymbol{a}_0 * \boldsymbol{x}_0) \right] \\ &\approx -\boldsymbol{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\rho} \left[ \boldsymbol{P}_{[p_0]} (\boldsymbol{a}_0 * \boldsymbol{\tilde{x}}_0) \right] \end{split}$$

is the normalized gradient vector from a chunk of data  $\mathbf{a}^{(-1)} := \mathbf{P}_{[p_0]}(\mathbf{a}_0 * \widetilde{\mathbf{x}}_0)$  with  $\widetilde{\mathbf{x}}_0$  a normalized Bernoulli–Gaussian random vector of length  $2p_0 - 1$ . Since  $\nabla \varphi_\rho \approx \nabla \varphi_{\ell^1}$ , expanding the gradient  $\nabla \varphi_{\ell^1}$  and rewriting the gradient  $\nabla_{\ell^1}(\mathbf{a}^{(-1)})$  in shift space gives us

$$\begin{split} -\nabla \varphi_{\rho^1}(\boldsymbol{a}^{(-1)}) &\approx \iota^* \boldsymbol{C}_{\boldsymbol{a}_0} \widecheck{\boldsymbol{C}}_{\boldsymbol{x}_0} \mathcal{S}_{\lambda} \left[ \widecheck{\boldsymbol{C}}_{\boldsymbol{x}_0} \boldsymbol{C}_{\boldsymbol{a}_0}^* \boldsymbol{P}_{[p_0]}(\boldsymbol{a}_0 * \widetilde{\boldsymbol{x}}_0) \right] \\ &= \iota^* \boldsymbol{C}_{\boldsymbol{a}_0} \boldsymbol{\chi} \left[ \boldsymbol{C}_{\boldsymbol{a}_0}^* \boldsymbol{P}_{[p_0]} \boldsymbol{C}_{\boldsymbol{a}_0} \widetilde{\boldsymbol{x}}_0 \right] \\ &\approx \iota^* \boldsymbol{C}_{\boldsymbol{a}_0} \boldsymbol{\chi} \left[ \widetilde{\boldsymbol{x}}_0 \right] \\ &\approx n \theta \cdot \iota^* \boldsymbol{C}_{\boldsymbol{a}_0} \mathcal{S}_{\lambda} \left[ \widetilde{\boldsymbol{x}}_0 \right], \end{split}$$

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

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

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

(5.2) 
$$a^{(0)} \approx P_{\mathbb{S}^{p-1}} \iota^* C_{a_0} S_{\lambda} [\widetilde{x}_0],$$

approximate  $S_{\lambda}[\tilde{\boldsymbol{x}}_0]$ , which resides near the subspace  $S_{\tau}$ , in which  $\tau$  contains the nonzero entries of  $\tilde{\boldsymbol{x}}_0$  on  $\{-p_0+1,\ldots,p_0-1\}$ . With high probability, the number of nonzero entries is  $|\tau| \lesssim 4\theta p_0$ , and we therefore conclude that our initializer  $\boldsymbol{a}^{(0)}$  satisfies

(5.3) 
$$\boldsymbol{a}^{(0)} \in \Sigma^{\gamma}_{4\theta p_0}.$$

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

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

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

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

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

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

For small t,

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

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

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

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

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

(5.8) 
$$d_{\alpha}\left(\boldsymbol{a}^{+}, \mathcal{S}_{\tau}\right) \leq d_{\alpha}\left(\boldsymbol{a}, \mathcal{S}_{\tau}\right).$$

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

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

We conclude this section with the following theorem.

Theorem 5.1 (convergence of retractive curvilinear search). Suppose signals  $a_0, x_0$  satisfy the conditions of Theorem 4.1,  $\theta > 10^3 c/p_0$  (c > 1), and  $\mathbf{a}_0$  is  $\mu$ -truncated shift coherent  $\max_{i\neq j} \left| \left\langle \boldsymbol{\iota}_{p_0}^* s_i[\boldsymbol{a}_0], \boldsymbol{\iota}_{p_0}^* s_j[\boldsymbol{a}_0] \right\rangle \right| \leq \mu. \text{ Write } \boldsymbol{g} = \operatorname{grad}[\varphi_{\rho}](\boldsymbol{a}) \text{ and } \boldsymbol{H} = \operatorname{Hess}[\varphi_{\rho}](\boldsymbol{a}). \text{ When the }$ smallest eigenvalue of H is strictly smaller than  $-\eta_v$ , let v be the unit eigenvector of smallest eigenvalue, scaled so  $v^*g \geq 0$ ; otherwise let v = 0. Define a sequence  $\{a^{(k)}\}_{k \in \mathbb{N}}$  where  $a^{(0)}$ equals (3.7) and for  $k = 1, 2, ..., K_1$ 

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

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

(5.11) 
$$\varphi_{\rho}(\boldsymbol{a}^{(k+1)}) < \varphi_{\rho}(\boldsymbol{a}^{(k)}) - \frac{1}{2} \left( t \| \boldsymbol{g}^{(k)} \|_{2}^{2} + \frac{1}{2} t^{4} \eta_{v} \| \boldsymbol{v}^{(k)} \|_{2}^{2} \right).$$

Then with probability at least 1-1/c, there exists some signed shift  $\bar{a} = \pm s_i[a_0]$ , where  $i \in [\pm p_0]$ , such that  $\|\mathbf{a}^{(k)} - \bar{\mathbf{a}}\|_2 \le \mu + 1/p$  for all  $k \ge K_1 = \text{poly}(n, p)$ . Here,  $\eta_v = c' n\theta \lambda$  for some  $c' < c_1$  in Theorem 4.1.

*Proof.* See subsection SM7.2.

**5.2.** Local refinement. In this section, we describe and analyze an algorithm which refines an estimate  $\bar{a} \approx a_0$  of the kernel to exactly recover  $(a_0, x_0)$ . Set

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

We alternatively minimize the Lasso objective with respect to a and x:

(5.13) 
$$x^{(k+1)} \leftarrow \underset{x}{\operatorname{argmin}} \frac{1}{2} ||a^{(k)} * x - y||_{2}^{2} + \lambda^{(k)} \sum_{i \notin I^{(k)}} |x_{i}|,$$

(5.14) 
$$a^{(k+1)} \leftarrow P_{\mathbb{S}^{p-1}} \left[ \underset{1}{\operatorname{argmin}} \frac{1}{2} \| a * x^{(k+1)} - y \|_2^2 \right]$$

(5.14) 
$$a^{(k+1)} \leftarrow P_{\mathbb{S}^{p-1}} \left[ \underset{a}{\operatorname{argmin}} \frac{1}{2} \| a * x^{(k+1)} - y \|_{2}^{2} \right],$$
  
(5.15)  $\lambda^{(k+1)} \leftarrow \frac{1}{2} \lambda^{(k)}, \quad I^{(k+1)} \leftarrow \operatorname{supp} (x^{(k+1)}).$ 

One departure from standard alternating minimization procedures is our use of a continuation method, which (i) decreases  $\lambda$ , and (ii) maintains a running estimate  $I^{(k)}$  of the support set. Our analysis will show that  $a^{(k)}$  converges to one of the signed shifts of  $a_0$  at a linear rate, in the sense that

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

It should be clear that exact recovery is unlikely if  $x_0$  contains many consecutive nonzero entries: in fact in this situation, even *nonblind* deconvolution fails. Therefore to obtain exact recovery it is necessary to put an upper bound on signal dimension n. Here, we introduce the notation  $\kappa_I$  as an upper bound for the number of nonzero entries of  $x_0$  in a length-p window:

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

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

(5.18) 
$$I = \operatorname{supp}(\boldsymbol{x}_0), \qquad I^{(k)} = \operatorname{supp}(\boldsymbol{x}^{(k)}).$$

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

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

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  $x_0$  is densely occupied. In our main theorem statement, we preclude this possibility by putting an upper bound on signal length n with respect to window length p and shift coherence p. We will assume

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

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

Recall that (4.23) in Theorem 3.1 provides that

(5.21) 
$$\|\bar{a} - a_0\|_2 \le (\mu + 1/p),$$

which is sufficiently close to  $a_0$  as long as (5.19) holds true. Here, we will elaborate upon this by showing that a single iteration of alternating minimization algorithm (5.13)–(5.15) is a contraction mapping for a toward  $a_0$ .

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

$$(5.22) x^{(k+1)} = \iota_J \left( \iota_J^* C_{\boldsymbol{a}^{(k)}}^* C_{\boldsymbol{a}^{(k)}} \iota_J \right)^{-1} \iota_J^* \left( C_{\boldsymbol{a}^{(k)}}^* C_{\boldsymbol{a}_0} x_0 - \lambda^{(k)} P_{J \setminus T} \sigma^{(k+1)} \right),$$

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

(5.23) 
$$a^{(k+1)} = \left(\iota^* C_{x^{(k+1)}}^* C_{x^{(k+1)}} \iota\right)^{-1} \left(\iota^* C_{x^{(k+1)}}^* C_{x_0} \iota a_0\right).$$

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

(5.24) 
$$x^{(k+1)} \approx x_0,$$
  $(x^{(k+1)} - x_0) \approx P_I \left( C_{a_0}^* C_{a_0} x_0 - C_{a_0}^* C_{a^{(k)}} x_0 \right)$   $\approx P_I \left[ C_{a_0}^* C_{x_0} \iota(a_0 - a^{(k)}) \right],$ 

which implies, while assuming  $C_{x_0}^* C_{x_0} \approx n\theta I$ , that from (5.23),

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

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

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

Now since  $C_{x_0}^* P_I C_{x_0} \approx n\theta \, e_0 e_0^*$ , this suggests that  $(n\theta)^{-1} \iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota$  approximates a contraction mapping with fixed point  $a_0$ , as follows:

$$(n\theta)^{-1} \iota^* C_{x_0}^* C_{a_0} P_I C_{a_0}^* C_{x_0} \iota \approx \iota^* C_{a_0} e_0 e_0^* C_{a_0}^* \iota$$

$$\approx a_0 a_0^*.$$
(5.27)

Hence, if we can ensure all of the above approximation is sufficiently and increasingly accurate as the iterate proceeds, the alternating minimization essentially is a power method which finds the leading eigenvector of matrix  $a_0 a_0^*$ —and the solution to this algorithm is apparently  $a_0$ . Indeed, we prove that the iterates produced by this sequence of operations converge to the ground truth at a linear rate, as long as it is initialized sufficiently nearby.

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

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

and for k = 1, 2, ...,

(5.29) 
$$x^{(k+1)} \leftarrow \underset{x}{\operatorname{argmin}} \frac{1}{2} || a^{(k)} * x - y ||_{2}^{2} + \lambda^{(k)} \sum_{i \neq I(k)} |x_{i}|,$$

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

(5.31) 
$$\lambda^{(k+1)} \leftarrow \frac{1}{2}\lambda^{(k)}, \qquad I^{(k+1)} \leftarrow \text{supp}\left(\boldsymbol{x}^{(k+1)}\right),$$

then

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

for every k = 0, 1, 2, ...

*Proof.* See subsection SM8.3.

Remark 5.3. The estimates  $x^{(k)}$  also converges to the ground truth  $x_0$  at a linear rate.

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

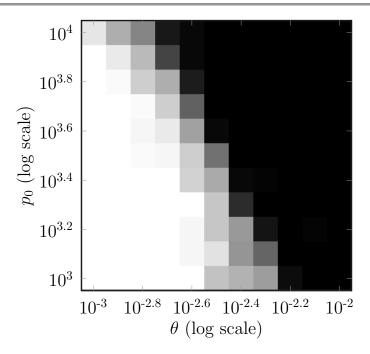


Figure 12. Success probability of SaS deconvolution under generic  $a_0$ ,  $x_0$  with varying kernel length  $p_0$  and sparsity rate  $\theta$ . When sparsity rate decreases sufficiently with respect to kernel length, successful recovery becomes very likely (brighter), and vice versa (darker). A transition 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 the generic case.

For each instance, we recover  $a_0$  and  $x_0$  from  $y = a_0 * x_0$  by minimizing problem (2.5). For ease of computation, we modify Algorithm 3.1 by replacing curvilinear search with the accelerated Riemannian gradient descent method (Algorithm 6.1), which is an adaptation of accelerated gradient descent [5] to the sphere. In particular, we apply momentum and increment by the Riemannian gradient via the exponential and logarithmic operators

(6.1) 
$$\operatorname{Exp}_{a}(u) := \cos(\|u\|_{2}) \cdot a + \sin(\|u\|_{2}) \cdot \frac{u}{\|u\|_{2}}$$

(6.1) 
$$\operatorname{Exp}_{\boldsymbol{a}}(\boldsymbol{u}) := \cos(\|\boldsymbol{u}\|_{2}) \cdot \boldsymbol{a} + \sin(\|\boldsymbol{u}\|_{2}) \cdot \frac{\boldsymbol{u}}{\|\boldsymbol{u}\|_{2}},$$
(6.2) 
$$\operatorname{Log}_{\boldsymbol{a}}(\boldsymbol{b}) := \arccos(\langle \boldsymbol{a}, \boldsymbol{b} \rangle) \cdot \frac{P_{\boldsymbol{a}^{\perp}}(\boldsymbol{b} - \boldsymbol{a})}{\|P_{\boldsymbol{a}^{\perp}}(\boldsymbol{b} - \boldsymbol{a})\|_{2}},$$

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

For each recovery instance, we say the local minimizer  $a_{\min}$  generated from Algorithm 6.1 is sufficiently close to a solution of the SaS deconvolution problem if

$$(6.3) \qquad \operatorname{success}(\boldsymbol{a}_{\min}, ; \boldsymbol{a}_0) := \{ \max_{\ell} |\langle s_{\ell} | \boldsymbol{a}_0 |, \boldsymbol{a}_{\min} \rangle| > 0.95 \}.$$

The result is shown in Figure 12. Our source code can be accessed via the following address:

https://github.com/sbdsphere/sbd\_experiments.git

Algorithm 6.1. SaS deconvolution with accelerated Riemannian gradient descent.

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

Initialize \boldsymbol{a}^{(0)} \leftarrow -\boldsymbol{P}_{\mathbb{S}^{p-1}} \nabla \varphi_{\rho} \left( \boldsymbol{P}_{\mathbb{S}^{p-1}} \left[ \boldsymbol{0}^{p_0-1}; [\boldsymbol{y}_0, \ldots, \boldsymbol{y}_{p_0-1}]; \boldsymbol{0}^{p_0-1} \right] \right), for k = 1, 2, \ldots, K do

Get momentum: \boldsymbol{w} \leftarrow \operatorname{Exp}_{\boldsymbol{a}^{(k)}} \left( \eta \cdot \operatorname{Log}_{\boldsymbol{a}^{(k-1)}}(\boldsymbol{a}^{(k)}) \right).

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

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

end for

Output: Return \boldsymbol{a}^{(K)}.
```

**7.** Discussion. In this section, we close by discussing the most important limitations of our results when  $a_0$  is coherent, regarding scenarios when the signal setting breaches our assumption, especially when  $x_0$  is either highly sparse or nonsymmetric, and highlighting corresponding directions for future work.

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

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

Our model assumes  $x_0$  to be Bernoulli–Gaussian vectors, which are sparse and symmetric i.i.d. random variables. When  $x_0$  is sparse but nonsymmetric, (e.g., Bernoulli), one can apply our result with a simple symmetrization trick, using the concatenated observation vectors [y, -y] as an input to our algorithm.

When  $x_0$  is highly sparse and if y is noiseless, it is possible to identify a short copy of  $a_0$  via looking for the shortest consecutive nonzero entries within y. When  $\theta \ll 1/p_0$ , these isolated copies are very common. Once  $\theta$  exceeds  $1/p_0$ , or when support  $x_0$  is not Bernoulli random while being more clustered, they become very uncommon. In particular, the probability of an isolated copy is small unless  $n \gtrsim \exp(p_0\theta)$ . Our proposed approach succeeds when  $n \ge \text{poly}(p_0)$ .

In applications involving noisy data, optimization approaches often outperform direct inspection, even for samples with isolated copies of  $a_0$ . An intuition for this is that optimization methods aggregate information across the sample. One practical avenue for obtaining the best of both worlds is to try to optimize the choice of data segment used for initialization. This can be a potential improvement for our data-driven initialization scheme, both in theory and in practice.

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

#### REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, Optimization Algorithms on Matrix Manifolds, Princeton University Press, 2009.
- [2] A. AHMED, B. RECHT, AND J. ROMBERG, Blind deconvolution using convex programming, IEEE Trans. Inform. Theory, 60 (2014), pp. 1711–1732.
- [3] G. AYERS AND J. C. DAINTY, Iterative blind deconvolution method and its applications, Opt. Lett., 13 (1988), pp. 547–549.
- [4] S. Baker and T. Kanade, Limits on super-resolution and how to break them, IEEE Trans. Pattern Anal. Mach. Intell., 24 (2002), pp. 1167–1183.
- [5] A. Beck and M. Teboulle, A fast iterative shrinkage-thresholding algorithm for linear inverse problems, SIAM J. Imaging Sci., 2 (2009), pp. 183–202, https://doi.org/10.1137/080716542.
- [6] A. J. Bell and T. J. Sejnowski, An information-maximization approach to blind separation and blind deconvolution, Neural Comput., 7 (1995), pp. 1129–1159.
- [7] A. BENICHOUX, E. VINCENT, AND R. GRIBONVAL, A fundamental pitfall in blind deconvolution with sparse and shift-invariant priors, in the 38th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP-2013), IEEE, 2013, pp. 6108-6112.
- [8] P. Bones, C. Parker, B. Satherley, and R. Watson, Deconvolution and phase retrieval with use of zero sheets, J. Opt. Soc. A, 12 (1995), pp. 1842–1857.
- [9] D. Briers, D. D. Duncan, E. R. Hirst, S. J. Kirkpatrick, M. Larsson, W. Steenbergen, T. Stromberg, and O. B. Thompson, Laser speckle contrast imaging: Theoretical and practical limitations, J. Biomed. Opt., 18 (2013), 066018.
- [10] P. Campisi and K. Egiazarian, Blind Image Deconvolution: Theory and Applications, CRC Press, 2016.
- [11] E. J. CANDES, M. B. WAKIN, AND S. P. BOYD, Enhancing sparsity by reweighted ℓ₁ minimization, J. Fourier Anal. Appl., 14 (2008), pp. 877–905.
- [12] M. CANNON, Blind deconvolution of spatially invariant image blurs with phase, IEEE Trans. Acoustics Speech Signal Process., 24 (1976), pp. 58–63.
- [13] A. S. CARASSO, Direct blind deconvolution, SIAM J. Appl. Math., 61 (2001), pp. 1980–2007, https://doi.org/10.1137/S0036139999362592.
- [14] T. F. CHAN AND C.-K. WONG, Total variation blind deconvolution, IEEE Trans. Image Process., 7 (1998), pp. 370–375.
- [15] S. Cheung, Y. Lau, Z. Chen, J. Sun, Y. Zhang, J. Wright, and A. Pasupathy, Beyond the Fourier transform: A nonconvex optimization approach to microscopy analysis, submitted.
- [16] Y. Chi, Guaranteed blind sparse spikes deconvolution via lifting and convex optimization, IEEE J. Selected Topics Signal Process., 10 (2016), pp. 782–794.
- [17] S. CHOUDHARY AND U. MITRA, Fundamental Limits of Blind Deconvolution Part II: Sparsity-Ambiguity Trade-Offs, preprint, https://arxiv.org/abs/1503.03184, 2015.
- [18] W. Dong, L. Zhang, G. Shi, and X. Wu, Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization, IEEE Trans. Image Process., 20 (2011), pp. 1838–1857.
- [19] B. EFRON, T. HASTIE, I. JOHNSTONE, AND R. TIBSHIRANI, Least angle regression, Ann. Statist., 32 (2004), pp. 407–499.
- [20] C. EKANADHAM, D. TRANCHINA, AND E. P. SIMONCELLI, A blind sparse deconvolution method for neural spike identification, in Advances in Neural Information Processing Systems 24, 2011, pp. 1440–1448.
- [21] R. FERGUS, B. SINGH, A. HERTZMANN, S. T. ROWEIS, AND W. T. FREEMAN, Removing camera shake from a single photograph, ACM Trans. Graphics, 25 (2006), pp. 787–794.

- [22] R. GE, F. HUANG, C. JIN, AND Y. YUAN, Escaping from saddle points—online stochastic gradient for tensor decomposition, in Conference on Learning Theory, 2015, pp. 797–842.
- [23] D. Goldfarb, Curvilinear path steplength algorithms for minimization which use directions of negative curvature, Math. Programming, 18 (1980), pp. 31–40.
- [24] D. GOLDFARB, C. Mu, J. WRIGHT, AND C. ZHOU, Using negative curvature in solving nonlinear programs, Comput. Optim. Appl., 68 (2017), pp. 479–502.
- [25] S. HARMELING, M. HIRSCH, S. SRA, AND B. SCHOLKOPF, Online blind deconvolution for astronomical imaging, in the 2009 IEEE International Conference on Computational Photography (ICCP 2009), IEEE, 2009, pp. 1–7.
- [26] R. JOHNSON, P. SCHNITER, T. J. ENDRES, J. D. BEHM, D. R. BROWN, AND R. A. CASAS, Blind equalization using the constant modulus criterion: A review, Proc. IEEE, 86 (1998), pp. 1927–1950.
- [27] N. JOSHI, R. SZELISKI, AND D. J. KRIEGMAN, PSF estimation using sharp edge prediction, in the 2008 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2008, pp. 1–8.
- [28] K. F. KAARESEN AND T. TAXT, Multichannel blind deconvolution of seismic signals, Geophys., 63 (1998), pp. 2093–2107.
- [29] M. KECH AND F. KRAHMER, Optimal injectivity conditions for bilinear inverse problems with applications to identifiability of deconvolution problems, SIAM J. Appl. Algebra Geom., 1 (2017), pp. 20–37, https://doi.org/10.1137/16M1067469.
- [30] D. KRISHNAN, T. TAY, AND R. FERGUS, Blind deconvolution using a normalized sparsity measure, in the 2011 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2011, pp. 233–240.
- [31] D. KUNDUR AND D. HATZINAKOS, Blind image deconvolution, IEEE Signal Process. Mag., 13 (1996), pp. 43–64.
- [32] R. LANE AND R. BATES, Automatic multidimensional deconvolution, J. Opt. Soc. A, 4 (1987), pp. 180-188.
- [33] R. G. Lane, Blind deconvolution of speckle images, J. Opt. Soc. A, 9 (1992), pp. 1508–1514.
- [34] Y. LAU, Q. QU, H.-W. KUO, P. ZHOU, Y. ZHANG, AND J. WRIGHT, Short-and-sparse deconvolution: A geometric approach, in the 8th International Conference on Learning Representations, 2020; preprint, https://arxiv.org/abs/1908.10959, 2019.
- [35] A. LEVIN, R. FERGUS, F. DURAND, AND W. T. FREEMAN, Deconvolution Using Natural Image Priors, Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, 2007.
- [36] A. LEVIN, Y. WEISS, F. DURAND, AND W. T. FREEMAN, Understanding blind deconvolution algorithms, IEEE Trans. Pattern Anal. Mach. Intell., 33 (2011), pp. 2354–2367.
- [37] M. S. LEWICKI, A review of methods for spike sorting: The detection and classification of neural action potentials, Network, 9 (1998), pp. R53–R78.
- [38] Y. LI AND Y. BRESLER, Global geometry of multichannel sparse blind deconvolution on the sphere, in Advances in Neural Information Processing Systems 31, Curran Associates, 2018, pp. 1132–1143.
- [39] Y. LI, K. LEE, AND Y. BRESLER, Identifiability in blind deconvolution with subspace or sparsity constraints, IEEE Trans. Inform. Theory, 62 (2016), pp. 4266–4275.
- [40] Y. LI, K. LEE, AND Y. BRESLER, Identifiability and stability in blind deconvolution under minimal assumptions, IEEE Trans. Inform. Theory, 63 (2017), pp. 4619–4633.
- [41] S. Ling and T. Strohmer, Self-calibration and biconvex compressive sensing, Inverse Problems, 31 (2015), 115002.
- [42] S. LING AND T. STROHMER, Blind deconvolution meets blind demixing: Algorithms and performance bounds, IEEE Trans. Inform. Theory, 63 (2017), pp. 4497–4520.
- [43] J. MARKHAM AND J.-A. CONCHELLO, Parametric blind deconvolution: A robust method for the simultaneous estimation of image and blur, J. Opt. Soc. A, 16 (1999), pp. 2377–2391.
- [44] M. MIYOSHI AND Y. KANEDA, Inverse filtering of room acoustics, IEEE Trans. Acoustics Speech Signal Process., 36 (1988), pp. 145–152.
- [45] P. A. NAYLOR AND N. D. GAUBITCH, Speech Dereverberation, Springer Science & Business Media, 2010.
- [46] M. R. OSBORNE, B. PRESNELL, AND B. A. TURLACH, A new approach to variable selection in least squares problems, IMA J. Numer. Anal., 20 (2000), pp. 389–403.
- [47] D. PERRONE AND P. FAVARO, Total variation blind deconvolution: The devil is in the details, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2014, pp. 2909–2916.
- [48] E. A. Pnevmatikakis, D. Soudry, Y. Gao, T. A. Machado, J. Merel, D. Pfau, T. Reardon,

- Y. Mu, C. Lacefield, W. Yang et al., Simultaneous denoising, deconvolution, and demixing of calcium imaging data, Neuron, 89 (2016), pp. 285–299.
- [49] S. K. Saha, Diffraction-Limited Imaging with Large and Moderate Telescopes, World Scientific, 2007.
- [50] Y. Sato, A method of self-recovering equalization for multilevel amplitude-modulation systems, IEEE Trans. Commun., 23 (1975), pp. 679–682.
- [51] O. Shalvi and E. Weinstein, New criteria for blind deconvolution of nonminimum phase systems (channels), IEEE Trans. Inform. Theory, 36 (1990), pp. 312–321.
- [52] Q. Shan, J. Jia, and A. Agarwala, High-quality motion deblurring from a single image, ACM Trans. Graphics, 27 (2008), pp. 1–10.
- [53] G. SHTENGEL, J. A. GALBRAITH, C. G. GALBRAITH, J. LIPPINCOTT-SCHWARTZ, J. M. GILLETTE, S. MANLEY, R. SOUGRAT, C. M. WATERMAN, P. KANCHANAWONG, M. W. DAVIDSON, R. D. FETTER, AND H. F. HESS, Interferometric fluorescent super-resolution microscopy resolves 3D cellular ultrastructure, Proc. Natl. Acad. Sci. USA, 106 (2009), pp. 3125–3130.
- [54] T. G. STOCKHAM, T. M. CANNON, AND R. B. INGEBRETSEN, Blind deconvolution through digital signal processing, Proc. IEEE, 63 (1975), pp. 678–692.
- [55] J. Sun, Q. Qu, and J. Wright, Complete dictionary recovery over the sphere II: Recovery by Riemannian trust-region method, IEEE Trans. Inform. Theory, 63 (2017), pp. 885–914.
- [56] P. WALK, P. Jung, G. E. PFANDER, AND B. HASSIBI, Blind Deconvolution with Additional Autocorrelations via Convex Programs, preprint, https://arxiv.org/abs/1701.04890, 2017.
- [57] L. WANG AND Y. CHI, Blind deconvolution from multiple sparse inputs, IEEE Signal Process. Lett., 23 (2016), pp. 1384–1388.
- [58] Z. WEN AND W. YIN, A feasible method for optimization with orthogonality constraints, Math. Program., 142 (2013), pp. 397–434.
- [59] D. WIPF AND H. ZHANG, Revisiting Bayesian blind deconvolution, J. Mach. Learn. Res., 15 (2014), pp. 3595–3634.
- [60] L. Xu and J. Jia, Two-phase kernel estimation for robust motion deblurring, in the 11th European Conference on Computer Vision, Springer, 2010, pp. 157–170.
- [61] J. Yang, J. Wright, T. S. Huang, and Y. Ma, Image super-resolution via sparse representation, IEEE Trans. Image Process., 19 (2010), pp. 2861–2873.
- [62] Y.-L. YOU AND M. KAVEH, Anisotropic blind image restoration, in Proceedings of the 3rd International Conference on Image Processing, Vol. 2, IEEE, 1996, pp. 461–464.
- [63] Y. ZHANG, H.-W. KUO, AND J. WRIGHT, Structured local optima in sparse blind deconvolution, IEEE Trans. Inform. Theory, 66 (2020), pp. 419–452.
- [64] Y. Zhang, Y. Lau, H.-W. Kuo, S. Cheung, A. Pasupathy, and J. Wright, On the global geometry of sphere-constrained sparse blind deconvolution, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2017, pp. 4894–4902.