



**HAL**  
open science

# Local Geometry of Nonconvex Spike Deconvolution from Low-Pass Measurements

Maxime Ferreira Da Costa, Yuejie Chi

► **To cite this version:**

Maxime Ferreira Da Costa, Yuejie Chi. Local Geometry of Nonconvex Spike Deconvolution from Low-Pass Measurements. *IEEE Journal on Selected Areas in Information Theory*, 2023, 4, pp.1-15. 10.1109/JSAIT.2023.3262689 . hal-03809877v2

**HAL Id: hal-03809877**

**<https://centralesupelec.hal.science/hal-03809877v2>**

Submitted on 16 May 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

# Local Geometry of Nonconvex Spike Deconvolution from Low-Pass Measurements

Maxime Ferreira Da Costa\*  
CentraleSupélec | Université Paris-Saclay

Yuejie Chi†  
Carnegie Mellon University

August 2022; Revised February 2023

## Abstract

Spike deconvolution is the problem of recovering the point sources from their convolution with a known point spread function, which plays a fundamental role in many sensing and imaging applications. In this paper, we investigate the local geometry of recovering the parameters of point sources—including both amplitudes and locations—by minimizing a natural nonconvex least-squares loss function measuring the observation residuals. We propose preconditioned variants of gradient descent (GD), where the search direction is scaled via some carefully designed preconditioning matrices. We begin with a simple fixed preconditioner design, which adjusts the learning rates of the locations at a different scale from those of the amplitudes, and show it achieves a linear rate of convergence—in terms of *entrywise* errors—when initialized close to the ground truth, as long as the separation between the true spikes is sufficiently large. However, the convergence rate slows down significantly when the dynamic range of the source amplitudes is large. To bridge this issue, we introduce an adaptive preconditioner design, which compensates for the learning rates of different sources in an iteration-varying manner based on the current estimate. The adaptive design provably leads to an accelerated convergence rate that is independent of the dynamic range, highlighting the benefit of adaptive preconditioning in nonconvex spike deconvolution. Numerical experiments are provided to corroborate the theoretical findings.

**Keywords:** nonconvex spike deconvolution, preconditioned gradient descent, local geometry

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>                                     | <b>2</b> |
| 1.1      | Observation model                                       | 3        |
| 1.2      | Our contributions                                       | 3        |
| 1.3      | Related work  | 4        |
| 1.4      | Notation and paper organization                         | 5        |
| <b>2</b> | <b>How does preconditioning help local convergence?</b> | <b>5</b> |
| 2.1      | Preconditioned gradient descent                         | 5        |
| 2.2      | Invariant preconditioning                               | 6        |
| 2.3      | Adaptive preconditioning                                | 8        |

---

\*M. Ferreira Da Costa is with the Laboratory of Signals and Systems (L2S) at CentraleSupélec, Université Paris-Saclay. Part of this work was done while he was with University of Southern California and Carnegie Mellon University. Email: [maxime.ferreira@centralesupelec.fr](mailto:maxime.ferreira@centralesupelec.fr).

†Y. Chi is with Department of Electrical and Computer Engineering, Carnegie Mellon University. The work of Y. Chi is supported in part by Office of Naval Research under N00014-19-1-2404, and by National Science Foundation under CAREER ECCS-1818571 and ECCS-2126634. Email: [yuejiechi@cmu.edu](mailto:yuejiechi@cmu.edu).

|          |   |           |
|----------|---|-----------|
| <b>3</b> | <b>Analysis</b>                             | <b>8</b>  |
| 3.1      | Preliminaries . . . . .                     | 8         |
| 3.2      | Proof of Theorem 1 . . . . .                | 10        |
| 3.3      | Proof of Theorem 2 . . . . .                | 12        |
| <b>4</b> | <b>Numerical experiments</b>                | <b>13</b> |
| <b>5</b> | <b>Conclusion</b>                           | <b>15</b> |
| <b>A</b> | <b>Summation bounds of the Fejér kernel</b> | <b>18</b> |
| <b>B</b> | <b>Proof of the uniform Hessian bounds</b>  | <b>22</b> |
| B.1      | Technical lemmas . . . . .                  | 22        |
| B.2      | Proof of Theorem 3 . . . . .                | 23        |
| B.3      | Proof of Theorem 4 . . . . .                | 25        |
| B.4      | Proof of Lemma 6 . . . . .                  | 28        |
| B.5      | Proof of Lemma 7 . . . . .                  | 28        |

# 1 Introduction

Spike deconvolution, also known as super-resolution [1], is the task of recovering a stream of point sources from their convolution with a point spread function (PSF). This classical problem is at the core of many sensing and imaging modalities, including but not limited to radar, sonar, optical imaging, neuroimaging, and communication systems [2, 3, 4, 5]. The PSF, which models the physical limitations of the imaging device involved in the experimental process, is commonly assumed to act as a band-limited and shift-invariant low-pass filter on the point sources [6, 7]. The sharpness of the original sources, modeled by Dirac impulses, is degraded through the convolution process, introducing undesirable ambiguity on the complex amplitudes and locations of the sources. The spike deconvolution task amounts to inverting the low-passing effects of the PSF, and to recovering the original sources as precisely as possible.

There is a rich literature on algorithmic investigations of the spike deconvolution problem, ranging from classical root-finding methods such as Prony’s method, subspace methods such as MUSIC [8, 9], ESPRIT [10] and matrix pencil [11], to more recent optimization methods such as atomic norm minimization (*a.k.a. total variation minimization*) [12, 13, 14, 15, 16] and basis pursuit [17, 18]. While classical approaches harness the algebraic properties of complex exponentials by mapping the observations onto a low-dimensional linear subspace to recover the parameters of interest, optimization methods, on the other hand, attempt to recover the parameters via minimizing some carefully-designed loss function. As such, optimization methods tend to be more versatile in adapting to different imaging modalities, as well as amenable to modern advances in large-scale optimization. Inspired by the development of compressive sensing [19, 20], initial approaches for spike deconvolution relies on a discretization of the spike locations, and then attempts to recover a sparse solution using sparsity-promoting convex relaxations such as the LASSO [21, 22]. However, the fundamental issue of basis mismatch [23] inherent to the discretization process may significantly hinder the localization performance, and increasing the grid size to reach finer precision levels leads to higher computational cost. Therefore, there has been a surge of interest in developing provably correct convex programs—such as the atomic norm minimization framework mentioned earlier—for spike deconvolution *over the continuum* in recent years, with strong performance guarantees developed under sufficient separations between the point sources [24, 25, 26, 27]. Nonetheless, the atomic norm framework requires solving a semidefinite program whose complexity scales at least cubically with respect to the signal length; and therefore is computationally expensive and memory inefficient. In addition, although it is in principle possible to examine the so-called dual polynomial to localize the sources [14], it often boils down to an additional post-processing step on the output of the convex program to recover the source parameters, which may hamper the guarantee of the overall procedure.

Motivated by the recent success of nonconvex methods, especially simple first-order methods, in various signal estimation and machine learning tasks [28, 29], we are interested in understanding the efficacy of first-order methods in nonconvex spike deconvolution. In fact, first-order methods have already been popular

empirically for spike deconvolution, but little is known about their theoretical underpinnings [30]. As a first step, this paper focuses on the local geometry and performance guarantees of recovering the parameters of point sources—including both amplitudes and locations—by minimizing a natural nonconvex least-squares loss function measuring the observation residuals.

## 1.1 Observation model

Formally, we formulate the spike deconvolution problem as follows. Consider a vector of  $2r$  parameters  $\boldsymbol{\theta}^* = [a_1^*, \dots, a_r^*, \tau_1^*, \dots, \tau_r^*]^\top$ , where  $a_\ell^* \in \mathbb{C}$  and  $\tau_\ell^* \in \mathbb{R}$  correspond to the complex amplitude and location of the  $\ell$ -th spike, respectively,  $\ell = 1, \dots, r$ . Denoting by  $\mathcal{M}$  the set of Radon measures over the reals, we assume that the point source signal  $\mu^* \in \mathcal{M}$  to resolve is of the form

$$\mu^* = \mu(\boldsymbol{\theta}^*) = \sum_{\ell=1}^r a_\ell^* \delta_{\tau_\ell^*}, \quad (1)$$

where  $\delta_\tau$  stands for the Dirac function located at  $\tau \in \mathbb{R}$ . Let us further denote the largest and the smallest amplitude of the spikes as

$$a_{\max}^* = \max_{1 \leq \ell \leq r} |a_\ell^*| = \|\mathbf{a}^*\|_\infty, \quad a_{\min}^* = \min_{1 \leq \ell \leq r} |a_\ell^*|.$$

The dynamic range of the measure  $\mu^*$ , an important quantity that will be used repetitively later, is thus defined as  $a_{\max}^*/a_{\min}^*$ . Denoting by  $g \in L_1(\mathbb{R})$  the PSF, the temporal signal  $x \in L_1(\mathbb{R})$  resulting from the convolution of the point source signal  $\mu^*$  and the PSF  $g$  reads

$$x(\tau) = (g * \mu^*)(\tau) = \sum_{\ell=1}^r a_\ell^* g(\tau - \tau_\ell^*), \quad (2)$$

where  $*$  denotes the convolution product.

A versatile observation model commonly encountered in practice considers the measurements to be taken from a uniform sampling of the Fourier transform of the temporal signal  $x$ . Denote by  $\mathcal{F}(\cdot)$  the Fourier transform of a measure  $\mu$  lying in  $\mathcal{M}$ , given by

$$\mathcal{F}(\mu)(f) = \int_{\mathbb{R}} e^{-i2\pi f\tau} d\mu(\tau), \quad \forall \mu \in \mathcal{M}, \forall f \in \mathbb{R}. \quad (3)$$

Denote by  $G = \mathcal{F}(g)$  and  $X = \mathcal{F}(x)$  the Fourier transform of  $g$  and  $x$ , respectively. We assume that the PSF is band-limited within the bandwidth  $B = 1$  so that it constrains no frequency greater than  $1/2$ , *i.e.*  $G(f) = 0$  for  $|f| > \frac{1}{2}$ .<sup>1</sup> For convenience, we assume that an odd number  $N = 2n + 1$  of measurements are taken in the Fourier domain, uniformly spaced in the bandwidth  $[-\frac{1}{2}, \frac{1}{2}]$ . The sampled signal  $\mathbf{x} \in \mathbb{C}^N$  writes

$$\begin{aligned} \mathbf{x} &= \boldsymbol{\Phi}(\mu^*) \\ &= \text{diag}(\mathbf{g}) \left[ \mathcal{F}(\mu) \left( -\frac{n}{N} \right), \dots, \mathcal{F}(\mu) \left( \frac{n}{N} \right) \right]^\top, \end{aligned} \quad (4)$$

where  $\boldsymbol{\Phi} : \mathcal{M} \rightarrow \mathbb{C}^N$  represents the observation operator and  $\mathbf{g} \in \mathbb{C}^N$  is a vector with generic term  $g_k = G\left(\frac{k}{N}\right)$  for  $k = -n, \dots, n$ . Up to a scaling, we assume  $\mathbf{g}$  to have unit Euclidean norm, *i.e.*  $\|\mathbf{g}\|_2 = 1$ . The goal of spike deconvolution is thus to recover the measure  $\mu^*$ , or equivalently, the parameter  $\boldsymbol{\theta}^*$ , from  $\mathbf{x}$ .

## 1.2 Our contributions

In the rest of this paper, we assume that the model order  $r$  is known, and consider a natural nonconvex loss function, which aims to minimize the quadratic loss of the parameters  $\boldsymbol{\theta} = [a_1, \dots, a_r, \tau_1, \dots, \tau_r]^\top$  of the Radon measures, given by

$$\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \|\boldsymbol{\Phi}(\mu(\boldsymbol{\theta})) - \mathbf{x}\|_2^2. \quad (5)$$

<sup>1</sup>The normalization  $B = 1$  is made without loss of generality, up to a rescaling of the source locations  $\tau_\ell^*$ . Assuming that  $\tau_\ell^* \in [-\frac{T}{2}, \frac{T}{2}]$ , where  $T$  is the length of the observation window, the deconvolution problem only depends on the time-bandwidth product  $T \cdot B$  [31].

Due to the nonlinear form of the parameters, the loss function  $\mathcal{L}(\boldsymbol{\theta})$  is clearly nonconvex. As a first step towards nonconvex spike deconvolution, we are interested in understanding the local geometry of the loss function (5) and its implications on the computation efficacy of first-order methods. Without loss of generality, we assume the autocorrelation function  $h$  of the PSF  $g$  to be a triangular low pass function, *i.e.*

$$h(\tau) = \int_{-\infty}^{\infty} f(u + \tau)\overline{f(u)}du = C \frac{\sin^2\left(\frac{\pi\tau}{2}\right)}{\left(\frac{\pi\tau}{2}\right)^2}, \quad (6)$$

where  $C > 0$  is a constant, and that their respective Fourier transform  $H(f)$  and  $G(f)$  are linked through the relation  $G(f) = \sqrt{H(f)}$  for all  $f \in \mathbb{R}$ . It comes that

$$g_k = G\left(\frac{k}{N}\right) = \sqrt{H\left(\frac{k}{N}\right)} = \sqrt{\frac{1}{n+1} \left(1 - \frac{|k|}{n+1}\right)}, \quad (7)$$

for all  $k = -n, \dots, n$  after rescaling with the constraint  $\|\mathbf{g}\|_2 = 1$ . Note that the triangular low-pass function and the Fejér kernel—its discrete counterpart—play an important role in the deconvolution literature and have been extensively proposed as a convolution kernel to evaluate the norm and distance between Radon measures [14]. Additionally, our main results can be re-derived with any other bandlimited PSF  $g$  by following analogous reasoning, as long as its autocorrelation  $h$  is an *absolutely integrable* function.

Concretely, we propose and analyze preconditioned variants of gradient descent (GD), where the search direction is scaled via some carefully designed preconditioning matrices. Our contributions are summarized as follows.

- We begin with a simple fixed preconditioner design, which adjusts the learning rates of the locations at a different scale from those of the amplitudes, and show it achieves a linear rate of convergence—in terms of *entrywise* errors—when initialized close to the ground truth, as long as the separation between the true spikes is sufficiently large. However, the convergence rate slows down significantly when the dynamic range of the source amplitudes is large.
- To bridge this issue, we introduce an adaptive preconditioner design, which compensates the learning rates of different sources in an iteration-varying manner based on the current estimate. The adaptive design provably leads to an accelerated convergence rate that is independent of the dynamic range, highlighting the benefit of adaptive preconditioning in nonconvex spike deconvolution.

Our result is based on understanding the geometric properties of scaled Hessian matrices via a set of novel summation bounds on the absolute sums of sampled Fejér kernels and higher-order derivatives, which might be of independent interest in other contexts.

### 1.3 Related work

The closest work to ours on recovering spike signals from low-pass observations using nonconvex optimization is [32]. The radius of the basin of attraction for gradient descent is characterized whenever the observation operator satisfies the restricted isometry property over the set of well-separated sparse measures. Although the problem setup is versatile, specializing this result in our context of low-pass measurements yields a convergence region whose size scales inversely with the number of sources, which is pessimistic when the number of sources is large. Moreover, the analysis in [32] focuses on the Euclidean error of the parameters, while we focus on the entrywise error, which is more meaningful for gauging the recovery quality of the point sources. Projected gradient methods, which merge pairs of colliding spikes at each iteration, have been proposed in [33, 34] but without theoretical convergence guarantees.

Our work can be viewed as falling into a growing line of research on developing provably efficient nonconvex methods—especially first-order methods—for high-dimensional signal estimation, examples including phase retrieval [35, 36], low-rank matrix estimation [37, 28], blind (sparse) deconvolution [38, 39, 40], dictionary learning [41, 42], multi-channel sparse deconvolution [43, 44], and so on. In particular, the preconditioned gradient methods considered in this paper are motivated by [45, 46, 47], which demonstrated that preconditioning can efficiently accelerate the convergence of gradient descent in ill-conditioned low-rank estimation.

## 1.4 Notation and paper organization

Vectors and matrices are denoted by boldface and capital boldface letters, respectively. Vectors  $\mathbf{x} \in \mathbb{C}^N$  with odd dimension  $N = 2n + 1$  are indexed between  $-n$  and  $n$ , so that  $\mathbf{x} = [x_{-n}, \dots, x_n]^\top$  for convenience. Transpose and Hermitian transpose of a vector or a matrix  $\mathbf{A}$  are denoted by  $\mathbf{A}^\top$  and  $\mathbf{A}^H$ , respectively. Furthermore, the adjoint of the operator  $\Phi$  is denoted by  $\Phi^*$ . We write  $\mathbf{1}_d$  and  $\mathbf{0}_d$  the all-one and null vector (or matrix) in dimension  $d$ , respectively. With a slight abuse of notation, we denote by  $|\mathbf{a}|$ ,  $|\mathbf{a}|^2$ ,  $\mathbf{a}^{-1}$  the vector with entries equal to the modulus, the squared modulus, and the inverse of the entries of  $\mathbf{a}$ , respectively. The element-wise product between two vectors  $\mathbf{a}$  and  $\mathbf{a}'$  is written as  $\mathbf{a} \odot \mathbf{a}'$ . We denote by  $\mathcal{D}_\ell$  the space of  $\ell$ -times differentiable functions of the real variable. For any function  $h \in \mathcal{D}_\ell$ , we write its  $\ell$ th derivative  $h^{(\ell)}$ . We denote by  $\langle \cdot, \cdot \rangle$  and  $\langle \cdot, \cdot \rangle_{\mathbb{R}} = \Re(\langle \cdot, \cdot \rangle)$  the usual inner product and real inner product between Radon measure, respectively. Additionally, we let  $\delta^{(\ell)} \in \mathcal{M}$  be the functional which satisfies

$$g^{(\ell)}(\tau) = \langle \delta_\tau^{(\ell)}, g \rangle, \quad \forall g \in \mathcal{D}_\ell, \forall \tau \in \mathbb{R}. \quad (8)$$

**Fejér kernel** We denote by  $F_N(\cdot)$  the normalized Fejér kernel of order  $N = 2n + 1$  defined by

$$\begin{aligned} F_N(t) &= \frac{1}{n+1} \sum_{k=-n}^n \left(1 - \frac{|k|}{n+1}\right) e^{i2\pi kt} \\ &= \begin{cases} \frac{\sin^2(\pi(n+1)t)}{(n+1)^2 \sin^2(\pi t)} & \text{if } t \notin \mathbb{Z} \\ 1 & \text{otherwise.} \end{cases} \end{aligned} \quad (9)$$

The Fejér kernel is a trigonometric polynomial, hence it is infinitely differentiable. We point out that the second derivative of  $F_N(\cdot)$  at the origin satisfies

$$\begin{aligned} F_N''(0) &= \frac{1}{n+1} \sum_{k=-n}^n -4\pi^2 k^2 \left(1 - \frac{|k|}{n+1}\right) \\ &= -\frac{2}{3}\pi^2 n(n+2) < 0. \end{aligned} \quad (10)$$

Some of its properties, key to this paper, are derived and discussed in Appendix A. The Fejér kernel plays an important role in the sequel as the Gramian  $\Phi^* \Phi : \mathcal{M} \rightarrow \mathcal{M}$  of the observation operator  $\Phi$  is a convolution product with  $F_N$ , *i.e.*

$$\Phi^* \Phi(\mu) = F_N * \mu, \quad \forall \mu \in \mathcal{M}. \quad (11)$$

**Wrap-around distance** For any set of  $r$  points  $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_r\} \subset \mathbb{T}$ , we denote by  $\Delta(\boldsymbol{\tau})$  is minimal *wrap-around distance*, defined by

$$\Delta(\boldsymbol{\tau}) \triangleq \min_{\ell \neq \ell'} \inf_{p \in \mathbb{Z}} |\tau_\ell - \tau_{\ell'} + p|. \quad (12)$$

The rest of this paper is organized as follows. Section 2 starts by defining the preconditioned gradient methods and two of its designs with provable local convergence guarantees, using a fixed preconditioner and an adaptive preconditioner in Section 2.2 and Section 2.3, respectively. Section 3 provides the analysis of the main theorems by controlling the conditioning of the scaled Hessian matrix of the loss function in a neighborhood of the ground truth. Numerical experiments are provided in Section 4 to corroborate our findings. Finally, a brief conclusion is drawn in Section 5.

## 2 How does preconditioning help local convergence?

### 2.1 Preconditioned gradient descent

Recognizing that the parameters corresponding to the amplitudes and locations may require different treatments, we consider iterates of preconditioned gradient descent (GD) to recover the ground truth parameters, where the preconditioner can possibly be iteration-varying. Given an initialization point  $\boldsymbol{\theta}_0 \in \mathbb{C}^{2r}$ ,

the update sequence of preconditioned GD is obtained by successively moving oppositely along the direction of a *linear transform* of the gradient. More specifically, the update rule reads

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{P}_k \nabla \mathcal{L}(\boldsymbol{\theta}_k), \quad (13)$$

where  $\boldsymbol{\theta}_k = [\mathbf{a}_k^\top, \boldsymbol{\tau}_k^\top]^\top$  is the  $k$ -th iterate,  $\mathbf{P}_k \in \mathbb{C}^{2r \times 2r}$  is a preconditioning matrix (also called preconditioner) that can vary at each iteration; the choice of  $\mathbf{P}_k$  will be detailed momentarily. Here, it is worth noticing that there are no additional learning rates in (13), which can be thought of as already absorbed and set within the preconditioner  $\mathbf{P}_k$ . By analogy with the celebrated Newton-Raphson method, which selects  $\mathbf{P}_k = \nabla^2 \mathcal{L}(\boldsymbol{\theta}_k)^{-1}$  (which might however be computationally expensive), the role of the preconditioning matrix  $\mathbf{P}_k$  is to balance the local optimization landscape towards a quadratic function to improve the convergence rate towards a local minimum over the vanilla gradient method. By basic calculation, the gradient  $\nabla \mathcal{L}(\boldsymbol{\theta})$  at point  $\boldsymbol{\theta} = [a_1, \dots, a_r, \tau_1, \dots, \tau_r]^\top$  is given by

$$\begin{aligned} \frac{d\mathcal{L}(\boldsymbol{\theta})}{da_j} &= \langle \boldsymbol{\Phi}(\delta_{\tau_j}), \boldsymbol{\Phi}(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \rangle \\ &= \langle \delta_{\tau_j}, \boldsymbol{\Phi}^* \boldsymbol{\Phi}(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \rangle \\ &= \langle \delta_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \rangle \\ &= \sum_{\ell=1}^r a_\ell F_N(\tau_j - \tau_\ell) - \sum_{\ell=1}^r a_\ell^* F_N(\tau_j - \tau_\ell^*) \end{aligned} \quad (14a)$$

for  $j = 1, \dots, r$ , and similarly,

$$\begin{aligned} \frac{d\mathcal{L}(\boldsymbol{\theta})}{d\tau_j} &= \left\langle \boldsymbol{\Phi}(a_j \delta'_{\tau_j}), \boldsymbol{\Phi}(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle_{\mathbb{R}} \\ &= \left\langle a_j \delta'_{\tau_j}, \boldsymbol{\Phi}^* \boldsymbol{\Phi}(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle_{\mathbb{R}} \\ &= \left\langle a_j \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle_{\mathbb{R}} \\ &= \Re \left( \bar{a}_j \left( \sum_{\ell=1}^r a_\ell F'_N(\tau_j - \tau_\ell) - \sum_{\ell=1}^r a_\ell^* F'_N(\tau_j - \tau_\ell^*) \right) \right) \end{aligned} \quad (14b)$$

for  $j = 1, \dots, r$ .

In this paper, we are particularly interested in preconditioning matrices  $\mathbf{P}_k$  that are diagonally structured, so they do not add computation overhead compared with vanilla gradient methods. In the sequel, we study the basin of attraction and the convergence rate of preconditioned GD for two different preconditioning strategies. The first consists of selecting a time-invariant, diagonal preconditioning matrix  $\mathbf{P} = \mathbf{P}_k$  whose role is to judiciously renormalize the learning rates between the amplitudes  $\mathbf{a} \in \mathbb{C}^r$  and the locations  $\boldsymbol{\tau} \in \mathbb{R}^r$  which are of different units. The second strategy seeks to dynamically update the preconditioning matrix  $\mathbf{P}_k$  based on the current iterate to better approximate the inverse of the Hessian matrix around the point  $\boldsymbol{\theta}_k$  and accelerate convergence.

## 2.2 Invariant preconditioning

In this section, we seek to recover the ground truth parameter  $\boldsymbol{\theta}^*$  from an instance of the preconditioned GD algorithm (13) where the sequence of preconditioning matrices is constant, *i.e.*  $\mathbf{P} = \mathbf{P}_k$  for all  $k \in \mathbb{N}$ . We fix

$$\mathbf{P} = \text{diag} \left( \begin{bmatrix} \mathbf{1}_r \\ -F''_N(0)^{-1} A^{-2} \mathbf{1}_r \end{bmatrix} \right), \quad (15)$$

where  $A > 0$  is an input parameter that controls the ratio between the learning rate applied to the amplitudes  $\mathbf{a}_k$  of the sources and that applied to the locations  $\boldsymbol{\tau}_k$  of the sources throughout the iterative process.

**Performance metric** To gauge the performance, we define by  $\mathbf{S} \in \mathbb{C}^{2r \times 2r}$  the weighting matrix

$$\mathbf{S} = \text{diag} \left( \left[ \frac{\mathbf{a}^{\star -1}}{\sqrt{-F_N''(0)} \mathbf{1}_r} \right] \right), \quad (16)$$

where  $F_N''(0)$  is given in (10), and study the convergence properties of preconditioned GD in terms of the infinity norm weighted by the matrix  $\mathbf{S}$ , *i.e.*

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty = \max_j \left\{ \frac{|a_{k,j} - a_j^*|}{|a_j^*|}, \sqrt{-F_N''(0)} |\tau_{k,j} - \tau_j^*| \right\}. \quad (17)$$

Intuitively, the role of this scaling is to analyze a unitless metric that decorrelates the error with the dynamic range of the sources and with the problem dimension, as we have  $\sqrt{-F_N''(0)} = \mathcal{O}(n)$  and that the error on the source locations is expected to be inversely proportional to the number of observation:  $\|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty = \mathcal{O}(n^{-1})$ .

The following theorem establishes the linear convergence of preconditioned GD with a fixed preconditioning matrix  $\mathbf{P}$  whenever the input parameter  $A$  is properly set, and the initial point  $\boldsymbol{\theta}_0$  is close enough to the ground truth  $\boldsymbol{\theta}^*$ , as long as the true spikes are sufficiently separated.

**Theorem 1** (Linear convergence with invariant preconditioner). *Suppose that  $n \geq 2$  and that the input parameter  $A$  satisfies  $\|\mathbf{a}^*\|_\infty \leq \frac{3}{2}A$ . Moreover, assume that*

$$\eta := 276.21 \frac{A^2 \|\mathbf{a}^*\|_\infty}{(a_{\min}^*)^3} ((n+1) \Delta(\boldsymbol{\tau}^*))^{-2} < 1, \quad (18)$$

then if the initial point  $\boldsymbol{\theta}_0 = [\mathbf{a}_0^\top, \boldsymbol{\tau}_0^\top]^\top$  satisfies

$$\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \leq \frac{1}{2}, \quad (19)$$

the iterates  $\{\boldsymbol{\theta}_k\}$  of preconditioned GD (13) with a fixed preconditioner (15) converge towards  $\boldsymbol{\theta}^*$  according to

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \left( 1 - \frac{1}{4} \frac{(a_{\min}^*)^2}{A^2} (1 - \eta) \right)^k \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \quad (20)$$

for all  $k \in \mathbb{N}$ .

Theorem 1 indicates that preconditioned GD admits a linear rate of convergence as long as the separation condition  $\Delta(\boldsymbol{\tau}^*)$  is sufficiently large with respect to the dynamic range, *i.e.*

$$(n+1) \Delta(\boldsymbol{\tau}^*) \gtrsim \left( \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \right)^{3/2}. \quad (21)$$

Additionally, our finding is independent of the number of sources  $r$  of the input measure, both in terms of the size of the basin of attraction (cf. (19)) and the convergence rate. Faster convergence rate are achieved for smaller values of the parameter  $0 < \eta < 1$ , when the separation  $\Delta(\boldsymbol{\tau}^*)$  of the true spikes is larger or the dynamic range  $\frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*}$  of the amplitudes is smaller. However, even for small values of  $\eta$ , the convergence rate  $\rho$  predicted by Theorem 1 is lower bounded by

$$\rho \geq 1 - \frac{1}{4} \left( \frac{a_{\min}^*}{A} \right)^2 \geq 1 - \frac{9}{16} \left( \frac{a_{\min}^*}{\|\mathbf{a}^*\|_\infty} \right)^2. \quad (22)$$

This suggests that a high dynamic range will lead to a slow convergence rate, independently of the separation  $\Delta(\boldsymbol{\tau}^*)$ . Additionally, the convergence guarantees established in Theorem 1 demand to adjust the input parameter  $A$  as a function of  $\|\mathbf{a}^*\|_\infty$ , which can be impractical in scenarios with no postulate on the norm of the source amplitudes.



## 2.3 Adaptive preconditioning

In order to mitigate the limitations of the fixed preconditioning strategy presented in Section 2.2, we propose to study an instance of preconditioned GD where the preconditioner  $\mathbf{P}_k$  varies at each iteration and is selected as a function of the current iterate  $\boldsymbol{\theta}_k$ . In particular, we fix

$$\mathbf{P}_k = \text{diag} \left( \begin{bmatrix} \mathbf{1}_r \\ -F_N''(0)^{-1} |\mathbf{a}_k|^{-2} \end{bmatrix} \right). \quad (23)$$

Similar to Theorem 1, the next theorem guarantees a linear convergence rate of the iterates towards the ground truth  $\boldsymbol{\theta}^*$ , provided a good enough initialization point  $\boldsymbol{\theta}_0$ , as long as the true spikes are sufficiently separated.

**Theorem 2** (Linear convergence with adaptive preconditioner). *Suppose that  $n \geq 2$ , and assume that*

$$\gamma := 11.60 \frac{\|\mathbf{a}\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} < \frac{1}{2}, \quad (24)$$

then if the initial point  $\boldsymbol{\theta}_0 = [\mathbf{a}_0^\top, \boldsymbol{\tau}_0^\top]^\top$  satisfies

$$\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \leq 1 - \sqrt{\frac{2}{3}}, \quad (25)$$

the iterates  $\{\boldsymbol{\theta}_k\}$  of preconditioned GD (13) with an adaptive preconditioner (23) converge towards  $\boldsymbol{\theta}^*$  according to

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \left(\frac{1}{2} + \gamma\right)^k \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \quad (26)$$

for all  $k \in \mathbb{N}$ .

Theorem 2 guarantees that preconditioned GD with an adaptive preconditioner achieves a *constant* linear rate of convergence in a similar basin of attraction, provided that the separation condition  $\Delta(\boldsymbol{\tau}^*)$  is sufficiently large with respect to the dynamic range, *i.e.*

$$(n+1)\Delta(\boldsymbol{\tau}^*) \gtrsim \left(\frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*}\right)^{1/2}, \quad (27)$$

which is much weaker than the requirement for the case using a fixed preconditioner, as indicated in Theorem 1. Consequently, this highlights the benefit of adaptive preconditioning in accelerating the convergence in the presence of high dynamic ranges for nonconvex spike deconvolution.

**Remark.** *We have not attempted to fully optimize the constants in the above theorems. Therefore, their values are set in a quite pessimistic fashion; see Section 4 for numerical experiments.*

## 3 Analysis

This section is devoted to proving the two main results of this paper comprised in Theorem 1 and Theorem 2. Before entering the core of the proofs, we first provide some warm-up analysis that will be required in the latter proofs.

### 3.1 Preliminaries

#### 3.1.1 Contraction of entrywise errors

The convergence analysis of the preconditioned gradient method presented in Theorem 1 and Theorem 2 calls for understanding of the contraction properties of the sequence  $\{\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty\}$ . Starting from the update rule (13), leveraging  $\nabla \mathcal{L}(\boldsymbol{\theta}_*) = \mathbf{0}$ , and applying the fundamental theorem of calculus, we have

$$\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_*) = \mathbf{S}(\boldsymbol{\theta}_k - \mathbf{P}_k \nabla \mathcal{L}(\boldsymbol{\theta}_k) - \boldsymbol{\theta}_*)$$

$$\begin{aligned}
&= \mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*) - \mathbf{S}\mathbf{P}_k (\nabla\mathcal{L}(\boldsymbol{\theta}_k) - \nabla\mathcal{L}(\boldsymbol{\theta}_*)) \\
&= \mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*) - \mathbf{S}\mathbf{P}_k \left( \int_0^1 \nabla^2\mathcal{L}(\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)) du \right) (\boldsymbol{\theta}_k - \boldsymbol{\theta}_*) \\
&= \left[ \mathbf{I} - \mathbf{S}\mathbf{P}_k \left( \int_0^1 \nabla^2\mathcal{L}(\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)) du \right) \mathbf{S}^{-1} \right] \mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*). \tag{28}
\end{aligned}$$

Let  $\mathcal{S}_k = \{\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*) \mid u \in [0, 1]\}$  be the line segment that connects  $\boldsymbol{\theta}^*$  and  $\boldsymbol{\theta}_k$  in  $\mathbb{C}^{2r}$ , and denote by  $\rho_k$  the quantity

$$\rho_k \triangleq \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \|\mathbf{I} - \mathbf{S}\mathbf{P}_k \mathbf{H}(\boldsymbol{\theta}) \mathbf{S}^{-1}\|_\infty, \tag{29}$$

where  $\mathbf{H}(\boldsymbol{\theta}) = \nabla^2\mathcal{L}(\boldsymbol{\theta})$  is the Hessian of the loss function  $\mathcal{L}$  at point  $\boldsymbol{\theta}$ . Continuing to bound (28) yields

$$\begin{aligned}
\|\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_*)\|_\infty &\leq \left\| \mathbf{I} - \mathbf{S}\mathbf{P}_k \left( \int_0^1 \mathbf{H}(\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)) du \right) \mathbf{S}^{-1} \right\|_\infty \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*)\|_\infty \\
&= \left\| \int_0^1 (\mathbf{I} - \mathbf{S}\mathbf{P}_k \mathbf{H}(\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)) \mathbf{S}^{-1}) du \right\|_\infty \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*)\|_\infty \\
&\leq \left( \int_0^1 \|\mathbf{I} - \mathbf{S}\mathbf{P}_k \mathbf{H}(\boldsymbol{\theta}^* + u(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)) \mathbf{S}^{-1}\|_\infty du \right) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*)\|_\infty \\
&\leq \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \{\|\mathbf{I} - \mathbf{S}\mathbf{P}_k \mathbf{H}(\boldsymbol{\theta}) \mathbf{S}^{-1}\|_\infty\} \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*)\|_\infty \\
&= \rho_k \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_*)\|_\infty. \tag{30}
\end{aligned}$$

Hence, the crux of the convergence analysis is to show that  $\rho_k < 1$  (and control the size of  $\rho_k$ ) uniformly over the segment  $\mathcal{S}_k$  whenever the point  $\boldsymbol{\theta}_k$  lies in an appropriate region centered around the ground truth  $\boldsymbol{\theta}^*$ . Further analysis towards that goal requires an explicit derivation of the Hessian matrix  $\mathbf{H}(\boldsymbol{\theta})$ , which is done next.

### 3.1.2 Hessian decomposition

Recall that  $\mathbf{H}(\boldsymbol{\theta}) = \nabla^2\mathcal{L}(\boldsymbol{\theta}) \in \mathbb{C}^{2r \times 2r}$  is the Hessian matrix of the loss function  $\mathcal{L}$  in (5) at the point  $\boldsymbol{\theta} = [\mathbf{a}^\top, \boldsymbol{\tau}^\top]^\top$ . We decompose  $\mathbf{H}(\boldsymbol{\theta})$  as

$$\mathbf{H}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{H}_{a,a}(\boldsymbol{\theta}) & \mathbf{H}_{a,\tau}(\boldsymbol{\theta}) \\ \mathbf{H}_{a,\tau}^H(\boldsymbol{\theta}) & \mathbf{H}_{\tau,\tau}(\boldsymbol{\theta}) \end{bmatrix}, \tag{31}$$

where each block is of size  $r \times r$ , with generic terms

$$[\mathbf{H}_{a,a}(\boldsymbol{\theta})]_{(i,j)} = \frac{d^2f(\boldsymbol{\theta})}{da_i da_j}, \quad [\mathbf{H}_{a,\tau}(\boldsymbol{\theta})]_{(i,j)} = \frac{d^2f(\boldsymbol{\theta})}{da_i d\tau_j}, \quad [\mathbf{H}_{\tau,\tau}(\boldsymbol{\theta})]_{(i,j)} = \frac{d^2f(\boldsymbol{\theta})}{d\tau_i d\tau_j}$$

for  $i, j = 1 \dots, r$ . A direct calculation of the Hessian matrix  $\mathbf{H}(\boldsymbol{\theta})$  (see, *e.g.* [32]) yields a decomposition of the form

$$\mathbf{H}(\boldsymbol{\theta}) = \mathbf{G}(\boldsymbol{\theta}) + \mathbf{E}(\boldsymbol{\theta}), \tag{32}$$

where the terms  $\mathbf{G}(\boldsymbol{\theta}) \in \mathbb{C}^{2r \times 2r}$  and  $\mathbf{E}(\boldsymbol{\theta}) \in \mathbb{C}^{2r \times 2r}$  are described in the sequel.

**Structure of  $\mathbf{G}(\boldsymbol{\theta})$**  The matrix  $\mathbf{G}(\boldsymbol{\theta})$  can be written as

$$\mathbf{G}(\boldsymbol{\theta}) = \text{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)} \mathbf{a}} \right]^H \mathbf{D}(\boldsymbol{\tau}) \text{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)} \mathbf{a}} \right] \right) \right), \tag{33}$$

with  $\mathbf{D}(\boldsymbol{\tau}) \in \mathbb{C}^{2r \times 2r}$  given with a block structure

$$\mathbf{D}(\boldsymbol{\tau}) = \begin{bmatrix} \mathbf{D}_0(\boldsymbol{\tau}) & \mathbf{D}_1(\boldsymbol{\tau}) \\ \mathbf{D}_1(\boldsymbol{\tau})^H & \mathbf{D}_2(\boldsymbol{\tau}) \end{bmatrix}. \tag{34}$$

The entries of the blocks  $\mathbf{D}_0(\boldsymbol{\tau})$ ,  $\mathbf{D}_1(\boldsymbol{\tau})$ ,  $\mathbf{D}_2(\boldsymbol{\tau}) \in \mathbb{C}^{r \times r}$  are composed of

$$[\mathbf{D}_0(\boldsymbol{\tau})]_{(i,j)} = F_N(\tau_i - \tau_j), \quad (35a)$$

$$[\mathbf{D}_1(\boldsymbol{\tau})]_{(i,j)} = -F'_N(\tau_i - \tau_j) / \sqrt{-F''_N(0)}, \quad (35b)$$

$$[\mathbf{D}_2(\boldsymbol{\tau})]_{(i,j)} = F''_N(\tau_i - \tau_j) / F''_N(0) \quad (35c)$$

for all  $i, j = 1, \dots, r$ . As shall be seen, the matrix  $\mathbf{G}(\boldsymbol{\theta})$  is a relatively well-conditioned matrix whose spectrum can be controlled as a function of the separation parameter between the spikes  $(n+1)\Delta(\boldsymbol{\tau}^*)$ , the dynamic range of the amplitudes  $\frac{\|\mathbf{a}\|_\infty}{a_{\min}^*}$ , and the distance of  $\boldsymbol{\theta}$  to the ground truth parameter  $\boldsymbol{\theta}^*$ .

**Structure of  $\mathbf{E}(\boldsymbol{\theta})$**  The matrix  $\mathbf{E}(\boldsymbol{\theta})$  is given by the block structure decomposition

$$\mathbf{E}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{0}_{r \times r} & \mathbf{E}_1(\boldsymbol{\theta}) \\ \mathbf{E}_1(\boldsymbol{\theta})^H & \mathbf{E}_2(\boldsymbol{\theta}) \end{bmatrix}, \quad (36)$$

where the entries of  $\mathbf{E}_1(\boldsymbol{\theta})$ ,  $\mathbf{E}_2(\boldsymbol{\theta}) \in \mathbb{C}^{r \times r}$  are given as

$$[\mathbf{E}_1(\boldsymbol{\theta})]_{(i,j)} = \begin{cases} \langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \rangle & i = j; \\ 0 & i \neq j; \end{cases} \quad (37a)$$

$$[\mathbf{E}_2(\boldsymbol{\theta})]_{(i,j)} = \begin{cases} a_j \langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \rangle & i = j; \\ 0 & i \neq j. \end{cases} \quad (37b)$$

It is worth noting that  $\mathbf{E}(\boldsymbol{\theta}_*) = \mathbf{0}$ , hence  $\mathbf{E}$  can be interpreted as a perturbation term that grows as  $\boldsymbol{\theta}$  deviates from  $\boldsymbol{\theta}_*$ .

With the above preliminaries, we are now ready to prove the main results of this paper. The following two sections present our convergence analyses for preconditioned GD using a fixed preconditioner and an adaptive preconditioner, respectively.

### 3.2 Proof of Theorem 1

We recall that  $\mathbf{P}_k = \mathbf{P} = \text{diag} \left( \begin{bmatrix} \mathbf{1}_r \\ -F''_N(0)^{-1} A^{-2} \mathbf{1}_r \end{bmatrix} \right)$  for all  $k \in \mathbb{N}$  when the preconditioner is fixed. The contraction analysis (30) presented in Section 3.1.1 suggests that the contraction rate in Theorem 1 is controlled by the quantity  $\|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$  in neighborhood around the ground truth  $\boldsymbol{\theta}^*$ . The next theorem, whose proof is deferred to Appendix B, provides a uniform bound on this quantity on a neighborhood of the ground truth.

**Theorem 3** (Uniform bound of the Hessian). *Suppose that  $n \geq 2$ ,  $(n+1)\Delta(\boldsymbol{\tau}^*) \geq 16.5$  and  $\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty \leq \frac{1}{4}\Delta(\boldsymbol{\tau}^*)$ . Let  $\boldsymbol{\theta}_k = [\mathbf{a}_k^\top, \boldsymbol{\tau}_k^\top]^\top$  and assume that  $A \geq \max\{\frac{3}{2}\|\mathbf{a}^*\|_\infty, \|\mathbf{a}_k\|_\infty\}$  then there exist two positive constants*

$$K_\Delta = 2.13, \quad (38a)$$

$$K_\theta = 44.42, \quad (38b)$$

such that for all  $\boldsymbol{\theta} = [\mathbf{a}^\top, \boldsymbol{\tau}^\top]^\top \in \mathcal{S}_k$  with  $\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty < 1$ , we have that

$$\begin{aligned} \|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty &\leq 1 - \left(\frac{a_{\min}^*}{A}\right)^2 (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\ &\quad + (4K_\Delta + K_\theta \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2. \end{aligned} \quad (39)$$

Theorem 3 provides a bound on  $\|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$  depending on the quantity  $\frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}$ , which can be made small enough under the hypothesis of Theorem 1. We proceed with the rest of the proof by induction.

For the base case, it is trivial that (20) holds for  $k = 0$ . We start the induction by assuming that

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \left(1 - \frac{1}{4} \left(\frac{a_{\min}^*}{A}\right)^2 (1 - \eta)\right)^k \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \quad (40)$$

holds for some  $k \in \mathbb{N}$ . We begin by verifying the assumptions of Theorem 3.

- First, as the dynamic range  $\frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \geq 1$ , it is easy to see that the hypothesis (18) immediately implies that  $(n+1)\Delta(\boldsymbol{\tau}^*) \geq 16.6$ .
- By the definition (17) and the induction hypothesis (40), it follows that

$$\begin{aligned} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty &\leq \frac{1}{\sqrt{-F_N''(0)}} \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \frac{1}{\sqrt{-F_N''(0)}} \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \frac{\sqrt{3}}{2\pi} \frac{1}{n+1} \leq \frac{1}{n+1} \leq \frac{\Delta(\boldsymbol{\tau}^*)}{4}. \end{aligned} \quad (41)$$

- Furthermore, we have that

$$\begin{aligned} \|\mathbf{a}_k\|_\infty &\leq \|\mathbf{a}^*\|_\infty + \|\mathbf{a}_k - \mathbf{a}^*\|_\infty \\ &\leq \|\mathbf{a}^*\|_\infty (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \\ &< \|\mathbf{a}^*\|_\infty (1 + \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty) \\ &\leq \frac{3}{2} \|\mathbf{a}^*\|_\infty \leq A. \end{aligned} \quad (42)$$

Hence, the assumptions of Theorem 3 hold, which yields

$$\begin{aligned} \rho_k &= \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq 1 - \frac{(a_{\min}^*)^2}{4A^2} + \frac{9}{4} \left(4K_\Delta + \frac{1}{2}K_\theta\right) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} \\ &\leq 1 - \frac{(a_{\min}^*)^2}{4A^2} \left(1 - 9 \left(4K_\Delta + \frac{1}{2}K_\theta\right) \frac{A^2 \|\mathbf{a}^*\|_\infty}{a_{\min}^{*3}} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}\right) \\ &\leq 1 - \frac{(a_{\min}^*)^2}{4A^2} \left(1 - 276.21 \frac{A^2 \|\mathbf{a}^*\|_\infty}{a_{\min}^{*3}} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}\right) \\ &\leq 1 - \frac{(a_{\min}^*)^2}{4A^2} (1 - \eta), \end{aligned} \quad (43)$$

where we substituted the definition (18) of  $\eta$  in the last line. It results from the iterative analysis (30) that the next update  $\boldsymbol{\theta}_{k+1}$  obeys

$$\begin{aligned} \|\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}^*)\|_\infty &\leq \rho_k \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \left(1 - \frac{(a_{\min}^*)^2}{4A^2} (1 - \eta)\right) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \left(1 - \frac{(a_{\min}^*)^2}{4A^2} (1 - \eta)\right)^{k+1} \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty, \end{aligned} \quad (44)$$

which concludes the proof of Theorem 1.

### 3.3 Proof of Theorem 2

We proceed with the proof of Theorem 2 analogously to the proof of Theorem 1 presented in Section 3.2. First, we establish the following intermediate theorem that controls the conditioning of the scaled Hessian matrix  $\mathbf{S}\mathbf{P}_k(\boldsymbol{\theta})\mathbf{H}\mathbf{S}^{-1}$  uniformly over the segment  $\mathcal{S}_k$  as a function of the weighted infinity-norm distance  $\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty$ . The proof of Theorem 4 is deferred to Appendix B.3.

**Theorem 4** (Uniform bound of the Hessian). *Suppose that  $n \geq 2$ ,  $(n+1)\Delta(\boldsymbol{\tau}^*) \geq 4.7$  and  $\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty \leq \frac{1}{4}\Delta(\boldsymbol{\tau}^*)$  then there exists two positive constants*

$$K_\Delta \leq 2.32, \quad (45a)$$

$$K_\theta \leq 75.80, \quad (45b)$$

such that for all  $\boldsymbol{\theta} = [\mathbf{a}^\top, \boldsymbol{\tau}^\top]^\top \in \mathcal{S}_k$  satisfying  $\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty < 1$  we have that

$$\begin{aligned} \|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty &\leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \\ &\quad + (4K_\Delta + K_\theta \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}. \end{aligned} \quad (46)$$

The rest of the proof follows similarly by induction. For the base case, it is trivial that the initial point  $\boldsymbol{\theta}_0$  verifies (26). We now assume that  $\boldsymbol{\theta}_k$  satisfies

$$\|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \left(\frac{1}{2} + \gamma\right)^k \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \quad (47)$$

for some  $k \in \mathbb{N}$ . Let's verify the assumptions of Theorem 4.

- First, as the dynamic range  $\frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \geq 1$ , the assumption (24) easily implies that  $(n+1)\Delta(\boldsymbol{\tau}^*) \geq 4.7$ .
- By the definition (17) and the induction hypothesis (47), it follows that

$$\begin{aligned} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty &\leq \frac{1}{\sqrt{-F_N''(0)}} \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \frac{1}{\sqrt{-F_N''(0)}} \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \frac{\sqrt{3}}{\pi} \left(1 - \sqrt{\frac{2}{3}}\right) \frac{1}{n+1} \\ &\leq \frac{1}{n+1} \leq \frac{\Delta(\boldsymbol{\tau}^*)}{4}. \end{aligned} \quad (48)$$

Hence, the assumptions of Theorem 4 hold. Noticing that the function  $f(u) := \frac{1}{(1-u)^2}$  is increasing over  $[0, 1)$ , we have that  $\frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \leq \frac{1}{2}$ . Together with Theorem 4, this yields the bound

$$\begin{aligned} \rho_k &= \max_{\boldsymbol{\theta} \in \mathcal{S}_k} \|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \frac{1}{2} + \left(4K_\Delta + K_\theta \left(1 - \sqrt{\frac{2}{3}}\right)\right) \frac{\|\mathbf{a}^*\|_\infty}{2a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} \\ &\leq \frac{1}{2} + 11.60 \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} \\ &\leq \frac{1}{2} + \gamma < 1, \end{aligned} \quad (49)$$

where we substituted the definition (24) of  $\gamma$  in the third inequality. It results from the iterative analysis (30) that the next update  $\boldsymbol{\theta}_{k+1}$  satisfies

$$\begin{aligned} \|\mathbf{S}(\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}^*)\|_\infty &\leq \rho_k \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \left(\frac{1}{2} + \gamma\right) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \\ &\leq \left(\frac{1}{2} + \gamma\right)^{k+1} \|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_\infty, \end{aligned} \quad (50)$$

which concludes the proof of Theorem 2.

## 4 Numerical experiments

This section provides a numerical validation of Theorem 1 and Theorem 2. In the following experiments, the signal length is set to  $N = 65$  (*i.e.*  $n = 32$ ). The ground truth signal is composed of  $r = 6$  sources placed in the interval  $[-\frac{1}{2}, \frac{1}{2})$  while ensuring that  $(n + 1)\Delta(\boldsymbol{\tau}^*) \geq 2$ , which is a more optimistic separation condition than what the theorems' statements suggest. Additionally, the dynamic range is denoted by  $\kappa = \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*}$ . The complex amplitudes  $\mathbf{a}^* \in \mathbb{C}^r$  are selected independently and uniformly at random in a complex annulus with bounds  $1 \leq |a_\ell^*| \leq \kappa$ . The input parameter of the invariant preconditioning scheme is set at  $A = \frac{3}{2} \|\mathbf{a}^*\|_\infty$ .

**Size of the basin of attraction** Of critical importance in the analysis of Theorem 1 and Theorem 2 is the distance between the initial parameter  $\boldsymbol{\theta}_0$  and the ground truth  $\boldsymbol{\theta}_*$ . We start by comparing the success rates of both preconditioning schemes on reconstructing the ground truth as a function of the initialization distance  $\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*)\|_\infty$ . In each experiment, the starting point  $\boldsymbol{\theta}_0$  is drawn uniformly over the set of points equidistant to  $\boldsymbol{\theta}_*$ . An experiment is labeled as a success if  $\|\mathbf{S}(\boldsymbol{\theta}_{200} - \boldsymbol{\theta}_*)\|_\infty \leq 10^{-2}$  after 200 iterations. Figure 1 suggests that, for both schemes, the size of the basin of attraction is independent of the dynamic range  $\kappa$ , and is around the order of magnitude  $\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*)\|_\infty \simeq 1$ . This suggests that the numerical constants ( $1/2$  and  $1 - \sqrt{2/3} \sim 0.184$ , respectively) set forth in Theorem 1 and Theorem 2, respectively, are pessimistic and nonconvex spike deconvolution performs in a much more benign manner than predicted by our theory, indicating room for further refinements.

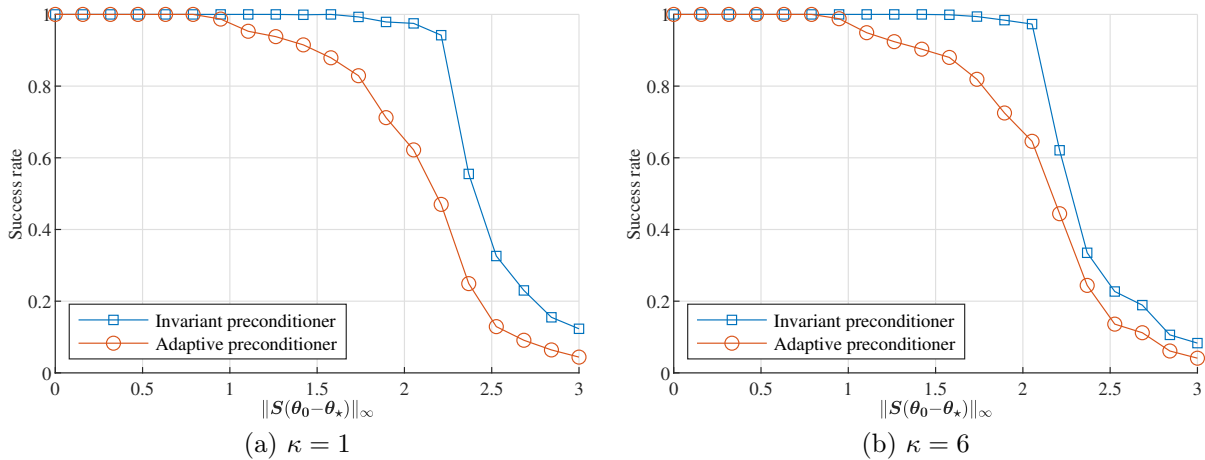


Figure 1: Success rate of the invariant and adaptive preconditioning schemes on reconstructing the ground truth  $\boldsymbol{\theta}_*$  as a function of the initialization distance  $\|\mathbf{S}(\boldsymbol{\theta}_0 - \boldsymbol{\theta}_*)\|_\infty$  for different dynamic ranges  $\kappa$ . The results are averaged over 1000 randomized trials.

**Linear convergence using a spectral initialization** In practice, several ad hoc initialization methods could be envisaged to produce an initial point  $\boldsymbol{\theta}_0$  that falls in the basin of attraction of the preconditioned

gradient descent methods. Herein, we proceed by uniformly discretizing the spectral domain over  $N$  elements. Given the knowledge of the ground truth model order  $r$ , the initial locations  $\boldsymbol{\tau}_0$  are selected as the  $r$  elements of the discrete grid whose weighted Fourier transform best describes the observation  $\boldsymbol{x}$ . Mathematically, consider the following optimization problem

$$\mathbf{u}_0 = \arg \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{x} - \text{diag}(\mathbf{g})\mathbf{F}_N\mathbf{u}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{u}\|_0 \leq r, \quad (51)$$

where  $\mathbf{F}_N \in \mathbb{C}^{N \times N}$  is a discrete Fourier transform matrix, and the  $\ell_0$ -norm denotes the cardinality of the support. Writing  $\mathcal{I}_0 = \text{supp}(\mathbf{u}_0) \subset [-n, \dots, n]$  the support of the solution  $\mathbf{u}_0$  of (51), the parameter  $\boldsymbol{\theta}_0$  is constructed in a second stage by selecting  $\boldsymbol{\tau}_0 = [\frac{k_1}{N}, \dots, \frac{k_r}{N}]^\top$  where  $k_\ell \in \mathcal{I}_0$ ,  $\ell = 1, \dots, r$  and  $\mathbf{a}_0 = \mathbf{u}|_{\mathcal{I}_0}$  as the restriction of  $\mathbf{u}$  to the elements in  $\mathcal{I}_0$ . The program (51) is itself a non-convex sparse reconstruction problem, which we approximate the solution using the orthogonal matching pursuit algorithm [48]. The proposed initialization procedure offers several benefits over more classical methods: It is highly scalable, robust to high dynamic range, and does not involve any polynomial root finding subroutine.

Figure 2 pictures the convergence rate of preconditioned GD under the invariant and adaptive preconditioning schemes, respectively. For both schemes,  $\boldsymbol{\theta}_0$  is selected according to the previously described initialization procedure. It can be seen that, although both preconditioning schemes ensure a linear converge of the iterate sequence, the convergence rate with a fixed preconditioner degrades as the dynamic range of the sources increases. In contrast, the one with an adaptive preconditioner remains unchanged. Additionally, the adaptive preconditioning scheme benefits from faster convergence rates for a given dynamic range. These experimental results corroborate the theoretical findings presented in Section 2.

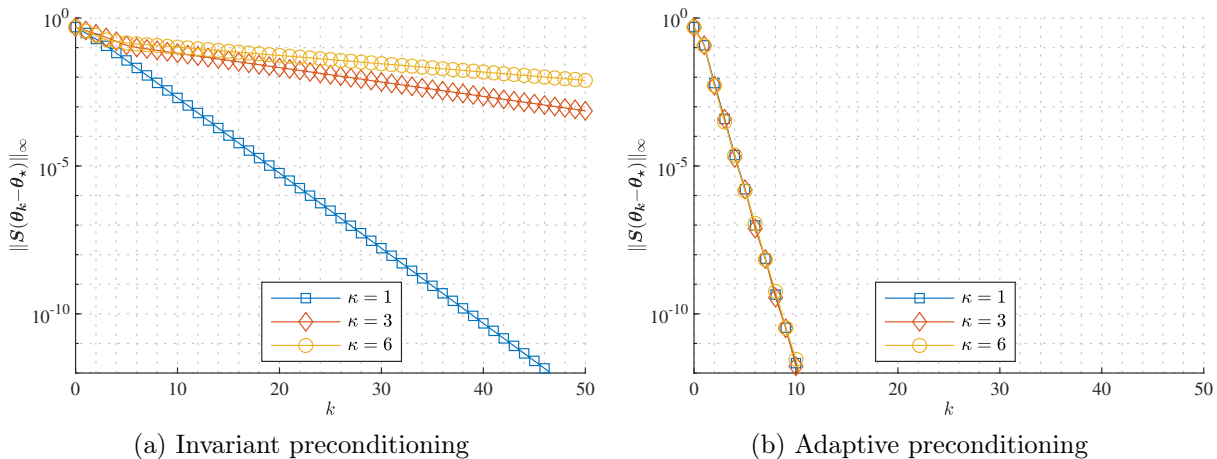


Figure 2: Convergence rates of the iterate sequence of preconditioned GD towards the ground truth as the dynamic range  $\kappa$  varies for: (a) the invariant preconditioning scheme; (b) the adaptive preconditioning scheme.

**Noisy recovery** We next examine the performance of preconditioned GD in the presence of noise. We assume observations of the form  $\boldsymbol{x} = \boldsymbol{\Phi}(\boldsymbol{\mu}^*) + \boldsymbol{w}$ , where  $\boldsymbol{w}$  is white Gaussian noise, and estimate  $\boldsymbol{\mu}^*$  by minimizing (5) starting from an initial point  $\boldsymbol{\theta}_0$  obtained by the spectral initialization procedure described above. Figure 3 draws the statistical error  $\|\mathcal{S}(\boldsymbol{\theta}_{200} - \boldsymbol{\theta}_*)\|_\infty$  of both preconditioning schemes after 200 iterations — when convergence is reached — as a function of the signal-to-noise ratio (SNR), defined as  $\text{SNR} = \|\boldsymbol{\Phi}(\boldsymbol{\mu}^*)\|_2^2 / \|\boldsymbol{w}\|_2^2$ . The results are benchmarked against the Cramér-Rao bound (CRB) [49]. Both statistical errors remain close to the CRB under a sufficiently large SNR, providing an empirical validation of the robustness of the proposed algorithms.

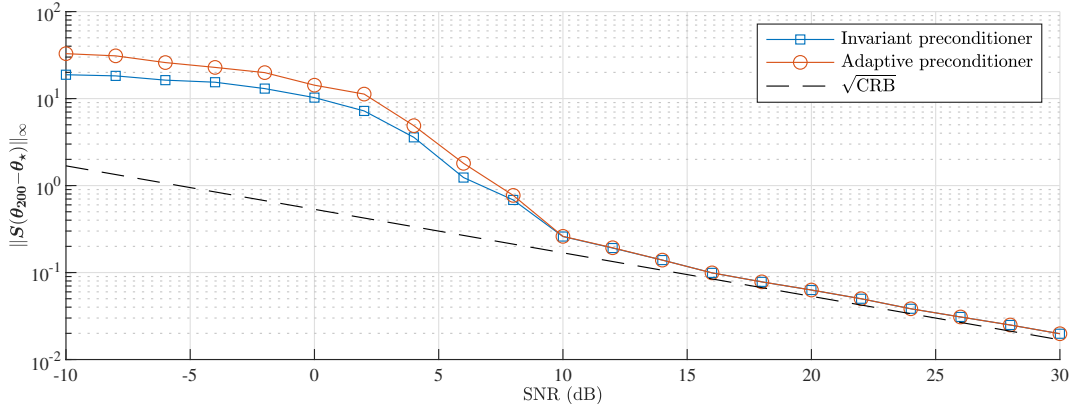


Figure 3: Statistical error  $\|\mathcal{S}(\theta_{200} - \theta_*)\|_{\infty}$  of both preconditioned GD schemes as a function of the SNR. The dynamic range is set to  $\kappa = 3$ , and the results are averaged over 1000 randomized trials.

## 5 Conclusion

This work proposed and analyzed preconditioned gradient methods for nonconvex spike deconvolution using both fixed and adaptive preconditioners, and demonstrated that for ground truth with sufficiently separated spikes, the proposed methods achieve a linear rate of convergence that is independent of the number of spikes, as long as a close enough initialization is provided near the ground truth. In particular, by designing the preconditioner to compensate adaptively for the amplitude profile of the spikes, it is possible to accelerate the convergence rate to be dimension-free and independent of the dynamic range, while the convergence using a fixed preconditioner slows down when the dynamic range is large. Our work thus highlights the importance of preconditioning in accelerating convergence in nonconvex spike deconvolution.

As a first step towards understanding the efficacy of first-order methods for spike deconvolution, this work opens up several interesting directions for further investigation.

- *Initialization schemes.* One immediate direction is to analyze initialization schemes that produce initial estimates that fall into the basin of attraction, which we suspect the procedure described in Section 4 is a good candidate.
- *Model order.* For simplicity, it is assumed that the model order  $r$  is known perfectly, which might not hold in practice. It is of great interest to develop modified algorithms when the model order is overspecified, which has recently been examined comprehensively in [50] for low-rank estimation from small random initializations.
- *General observations.* Another direction is to extend the analysis to more general observation operators, possibly including random sampling, missing data, as well as corruptions. This may necessarily require a reformulation of the loss function, such as a nonsmooth and nonconvex formulation using the least absolute deviation [47] to improve robustness.
- *Separation condition.* Last but not least, it is of great importance to study to what extent it is possible to relax the success condition in terms of the separation condition, possibly with additional positive constraints of the source amplitudes.

## References

- [1] D. L. Donoho, “Superresolution via sparsity constraints,” *SIAM Journal on Mathematical Analysis*, vol. 23, no. 5, pp. 1309–1331, 1992.
- [2] L. C. Potter, E. Ertin, J. T. Parker, and M. Cetin, “Sparsity and compressed sensing in radar imaging,” *Proceedings of the IEEE*, vol. 98, no. 6, pp. 1006–1020, 2010.



- [3] Z. Zhu, G. Tang, P. Setlur, S. Gogineni, M. B. Wakin, and M. Rangaswamy, “Super-resolution in SAR imaging: Analysis with the atomic norm,” in *2016 IEEE Sensor Array and Multichannel Signal Processing Workshop (SAM)*. IEEE, 2016, pp. 1–5.
- [4] L. Zhu, W. Zhang, D. Elnatan, and B. Huang, “Faster STORM using compressed sensing,” *Nature methods*, vol. 9, no. 7, pp. 721–723, 2012.
- [5] C. Berger, S. Zhou, J. Preisig, and P. Willett, “Sparse channel estimation for multicarrier underwater acoustic communication: From subspace methods to compressed sensing,” *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1708–1721, Mar. 2010.
- [6] J. Lindberg, “Mathematical concepts of optical superresolution,” *Journal of Optics - IOP Publishing*, vol. 14, no. 8, p. 83001, 2012.
- [7] M. Born and E. Wolf, *Principles of optics: electromagnetic theory of propagation, interference and diffraction of light*. Elsevier, 2013.
- [8] R. Schmidt, “Multiple emitter location and signal parameter estimation,” *IEEE Transactions on Antennas and Propagation*, vol. 34, no. 3, pp. 276–280, 1986.
- [9] W. Liao and A. Fannjiang, “MUSIC for single-snapshot spectral estimation: Stability and super-resolution,” *Applied and Computational Harmonic Analysis*, vol. 40, no. 1, pp. 33–67, 2016.
- [10] R. Roy and T. Kailath, “ESPRIT-estimation of signal parameters via rotational invariance techniques,” *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 37, no. 7, pp. 984–995, Jul. 1989.
- [11] A. Moitra, “Super-resolution, extremal functions and the condition number of Vandermonde matrices,” in *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, 2015, pp. 821–830.
- [12] Y. De Castro and F. Gamboa, “Exact reconstruction using Beurling minimal extrapolation,” *Journal of Mathematical Analysis and Applications*, vol. 395, no. 1, pp. 336–354, 2012.
- [13] G. Tang, B. N. Bhaskar, P. Shah, and B. Recht, “Compressed sensing off the grid,” *IEEE Transactions On Information Theory*, vol. 59, no. 11, pp. 7465–7490, 2013.
- [14] E. J. Candès and C. Fernandez-Granda, “Towards a mathematical theory of super-resolution,” *Communications on Pure and Applied Mathematics*, vol. 67, no. 6, pp. 906–956, 2014.
- [15] Y. Chi and M. Ferreira Da Costa, “Harnessing sparsity over the continuum: Atomic norm minimization for superresolution,” *IEEE Signal Processing Magazine*, vol. 37, no. 2, pp. 39–57, March 2020.
- [16] R. Heckel, V. I. Morgenshtern, and M. Soltanolkotabi, “Super-resolution radar,” *Information and Inference: A Journal of the IMA*, vol. 5, no. 1, pp. 22–75, 2016.
- [17] D. Malioutov, M. Cetin, and A. S. Willsky, “A sparse signal reconstruction perspective for source localization with sensor arrays,” *IEEE transactions on signal processing*, vol. 53, no. 8, pp. 3010–3022, 2005.
- [18] M. Herman and T. Strohmer, “High-resolution radar via compressed sensing,” *IEEE Transactions on Signal Processing*, vol. 57, no. 6, pp. 2275–2284, Jun. 2009.
- [19] E. Candès and T. Tao, “Near-Optimal Signal Recovery From Random Projections: Universal Encoding Strategies?” *IEEE Transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, Dec. 2006.
- [20] D. Donoho, “Compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 52, no. 2, pp. 1289–1306, Feb. 2006.
- [21] L. Zhu, W. Zhang, D. Elnatan, and B. Huang, “Faster storm using compressed sensing,” *Nature methods*, vol. 9, no. 7, pp. 721–723, 2012.
- [22] G. Tang, B. N. Bhaskar, and B. Recht, “Sparse recovery over continuous dictionaries-just discretize,” in *2013 Asilomar Conference on Signals, Systems and Computers*, Nov. 2013, pp. 1043–1047.

- [23] Y. Chi, L. L. Scharf, A. Pezeshki, and R. Calderbank, “Sensitivity to basis mismatch of compressed sensing for spectrum analysis and beamforming,” in *Proc. 6th U.S./Australia Joint Workshop on Defence Applications of Signal Processing (DASP)*, Sep. 2009, event-place: Lihue, HI.
- [24] C. Fernandez-Granda, “Super-resolution of point sources via convex programming,” *Information and Inference*, vol. 5, no. 3, pp. 251–303, Sep. 2016.
- [25] M. Ferreira Da Costa and W. Dai, “A tight converse to the spectral resolution limit via convex programming,” in *2018 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2018, pp. 901–905.
- [26] Q. Li and G. Tang, “Approximate support recovery of atomic line spectral estimation: A tale of resolution and precision,” *Applied and Computational Harmonic Analysis*, 2018.
- [27] M. Ferreira Da Costa and Y. Chi, “On the stable resolution limit of total variation regularization for spike deconvolution,” *IEEE Transactions on Information Theory*, vol. 66, no. 11, pp. 7237–7252, 2020.
- [28] Y. Chi, Y. M. Lu, and Y. Chen, “Nonconvex optimization meets low-rank matrix factorization: An overview,” *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5239–5269, 2019.
- [29] L. Chizat and F. Bach, “On the global convergence of gradient descent for over-parameterized models using optimal transport,” *Advances in neural information processing systems*, vol. 31, 2018.
- [30] J. Huang, M. Sun, J. Ma, and Y. Chi, “Super-resolution image reconstruction for high-density three-dimensional single-molecule microscopy,” *IEEE Transactions on Computational Imaging*, vol. 3, no. 4, pp. 763–773, 2017.
- [31] Y. Chi, “Guaranteed Blind Sparse Spikes Deconvolution via Lifting and Convex Optimization,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 782–794, 2016.
- [32] Y. Traonmilin and J.-F. Aujol, “The basins of attraction of the global minimizers of the non-convex sparse spike estimation problem,” *Inverse Problems*, vol. 36, no. 4, p. 045003, 2020.
- [33] Y. Traonmilin, J.-F. Aujol, and A. Leclaire, “Projected gradient descent for non-convex sparse spike estimation,” *IEEE Signal Processing Letters*, vol. 27, pp. 1110–1114, 2020.
- [34] P.-J. Bénard, Y. Traonmilin, and J.-F. Aujol, “Fast off-the-grid sparse recovery with over-parametrized projected gradient descent,” *arXiv preprint arXiv:2202.13757*, 2022.
- [35] E. Candès, X. Li, and M. Soltanolkotabi, “Phase retrieval via Wirtinger flow: Theory and algorithms,” *Information Theory, IEEE Transactions on*, vol. 61, no. 4, pp. 1985–2007, 2015.
- [36] Y. Chen and E. Candès, “Solving random quadratic systems of equations is nearly as easy as solving linear systems,” *Communications on Pure and Applied Mathematics*, vol. 70, no. 5, pp. 822–883, 2017.
- [37] R. Sun and Z.-Q. Luo, “Guaranteed matrix completion via non-convex factorization,” *IEEE Transactions on Information Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.
- [38] X. Li, S. Ling, T. Strohmer, and K. Wei, “Rapid, robust, and reliable blind deconvolution via nonconvex optimization,” *Applied and computational harmonic analysis*, vol. 47, no. 3, pp. 893–934, 2019.
- [39] Y. Chen, J. Fan, B. Wang, and Y. Yan, “Convex and nonconvex optimization are both minimax-optimal for noisy blind deconvolution under random designs,” *Journal of the American Statistical Association*, pp. 1–11, 2021.
- [40] K. Lee, Y. Li, M. Junge, and Y. Bresler, “Blind recovery of sparse signals from subsampled convolution,” *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 802–821, 2016.
- [41] J. Sun, Q. Qu, and J. Wright, “Complete dictionary recovery over the sphere I: Overview and the geometric picture,” *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 853–884, 2016.

- [42] D. Gilboa, S. Buchanan, and J. Wright, “Efficient dictionary learning with gradient descent,” *ICML Workshop on Modern Trends in Nonconvex Optimization for Machine Learning*, 2018.
- [43] L. Shi and Y. Chi, “Manifold gradient descent solves multi-channel sparse blind deconvolution provably and efficiently,” *IEEE Transactions on Information Theory*, vol. 67, no. 7, pp. 4784–4811, 2021.
- [44] Q. Qu, X. Li, and Z. Zhu, “Exact recovery of multichannel sparse blind deconvolution via gradient descent,” *SIAM Journal on Imaging Sciences*, vol. 13, no. 3, pp. 1630–1652, 2020.
- [45] T. Tong, C. Ma, and Y. Chi, “Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent,” *Journal of Machine Learning Research*, vol. 22, pp. 1–63, 2021.
- [46] T. Tong, C. Ma, A. Prater-Bennette, E. Tripp, and Y. Chi, “Scaling and scalability: Provable nonconvex low-rank tensor estimation from incomplete measurements,” *Journal of Machine Learning Research*, vol. 23, no. 163, pp. 1–77, 2022.
- [47] T. Tong, C. Ma, and Y. Chi, “Low-rank matrix recovery with scaled subgradient methods: Fast and robust convergence without the condition number,” *IEEE Transactions on Signal Processing*, vol. 69, pp. 2396–2409, 2021.
- [48] J. A. Tropp and A. C. Gilbert, “Signal recovery from random measurements via orthogonal matching pursuit,” *IEEE Transactions on information theory*, vol. 53, no. 12, pp. 4655–4666, 2007.
- [49] L. L. Scharf and L. T. McWhorter, “Geometry of the Cramér-Rao bound,” *Signal Processing*, vol. 31, no. 3, pp. 301–311, 1993.
- [50] X. Xu, Y. Shen, Y. Chi, and C. Ma, “The power of preconditioning in overparameterized low-rank matrix sensing,” *arXiv preprint arXiv:2302.01186*, 2023.

## A Summation bounds of the Fejér kernel

The purpose of this section is to present Lemma 5, which delivers fundamental bounds on the absolute sum of the Fejér kernel and its derivatives at sampled points of interest. Although specific to the Fejér kernel, Lemma 5 could be adapted to any other absolutely integrable point spread function without a significant change in the proof structure.

**Lemma 5** (Uniform bounds on the Fejér kernel). *Suppose that  $n \geq 2$ . Let  $\boldsymbol{\tau} = \{\tau_1, \dots, \tau_r\} \subset \mathbb{T}$ , and let  $\alpha > 0$  be such that  $(n+1)\Delta(\boldsymbol{\tau}) \geq \alpha$ . Let  $\mathbf{u} = \{u_{i,j}\}_{i \neq j} \subset \mathbb{R}$  be a set of  $\frac{r(r-1)}{2}$  real numbers that are absolutely bounded by  $\beta$  such that*

$$(n+1) \max_{i \neq j} \{|u_{i,j}|\} \triangleq \beta < \frac{\alpha}{2}. \quad (52)$$

Then the inequalities

$$\max_i \sum_{j \neq i} \left| F_N^{(\ell)}(\tau_j - \tau_i + u_{i,j}) \right| \leq C_\ell (n+1)^\ell ((n+1)\Delta(\boldsymbol{\tau}))^{-2} \quad (53)$$

hold for  $\ell = 0, 1, 2, 3$ , where the constants  $C_\ell$  only depend on  $\alpha$  and  $\beta$  and are given by

$$C_0 = \frac{4}{\pi^2} (\alpha - 2\beta)^{-1} \alpha, \quad (54a)$$

$$C_1 = \left( \frac{4}{\pi} (\alpha - 2\beta)^{-1} + \frac{8}{\pi^2} (\alpha - 2\beta)^{-2} \right) \alpha, \quad (54b)$$

$$C_2 = \left( \frac{80}{9} (\alpha - 2\beta)^{-1} + \frac{16}{\pi} (\alpha - 2\beta)^{-2} + \frac{64}{3\pi^2} (\alpha - 2\beta)^{-3} \right) \alpha, \quad (54c)$$

$$C_3 = \left( \frac{16}{3} (\alpha - 2\beta)^{-1} \pi + \frac{1488}{27} (\alpha - 2\beta)^{-2} + \frac{192}{\pi} (\alpha - 2\beta)^{-3} + \frac{192}{\pi^2} (\alpha - 2\beta)^{-4} \right) \alpha. \quad (54d)$$

*Proof.* First, the function  $F_N$  is a trigonometric polynomial, hence infinitely differentiable. Let us begin by examining the general expression of  $F_N(t)$  and its derivatives up to the third order. Assuming  $t \notin \mathbb{Z}$ , by basic calculation, we have that

$$F_N(t) = \frac{\sin^2(\pi(n+1)t)}{(n+1)^2 \sin^2(\pi t)}, \quad (55a)$$

$$F'_N(t) = \frac{2\pi}{n+1} \cos((n+1)\pi t) \sin((n+1)\pi t) \csc^2(\pi t) - \frac{2\pi}{(n+1)^2} \sin((n+1)\pi t) \cot(\pi t) \csc^2(\pi t), \quad (55b)$$

$$F''_N(t) = 2\pi^2 (\cos^2((n+1)\pi t) - \sin^2((n+1)\pi t)) \csc^2(\pi t) - \frac{8\pi^2}{n+1} \cos((n+1)\pi t) \sin((n+1)\pi t) \cot(\pi t) \csc^2(\pi t) + \frac{2\pi^2}{(n+1)^2} \sin^2(\pi(n+1)t) (2 \cot^2(\pi t) + 1) \csc^2(\pi t) \quad (55c)$$

$$F'''_N(t) = -(n+1)8\pi^3 \cos((n+1)\pi t) \sin((n+1)\pi t) \csc^2(\pi t) - 12\pi^3 (\cos^2((n+1)\pi t) - \sin^2((n+1)\pi t)) \cot(\pi t) \csc^2(\pi t) + \frac{12\pi^3}{n+1} \cos((n+1)\pi t) \sin((n+1)\pi t) (3 \cot^2(\pi t) + 1) \csc^2(\pi t) - \frac{8\pi^3}{(n+1)^2} \sin^2((n+1)\pi t) (3 \cot^3(\pi t) + 2 \cot(\pi t)) \csc^2(\pi t). \quad (55d)$$

Using the four trigonometric bounds  $|\sin(a)| \leq 1$ ,  $|\cos(a)| \leq 1$ ,  $|\cos(a) \sin(a)| \leq \frac{1}{2}$ , and  $|\cos^2(a) - \sin^2(a)| \leq 1$  for  $a \in \mathbb{R}$ , and the triangle inequality, (55) can be further absolutely bounded as

$$|F_N(t)| \leq \frac{1}{(n+1)^2} \csc^2(\pi t), \quad (56a)$$

$$|F'_N(t)| \leq \frac{\pi}{n+1} \csc^2(\pi t) + \frac{2\pi}{(n+1)^2} |\cot(\pi t)| \csc^2(\pi t), \quad (56b)$$

$$|F''_N(t)| \leq 2\pi^2 \csc^2(\pi t) + \frac{4\pi^2}{n+1} |\cot(\pi t)| \csc^2(\pi t) + \frac{2\pi^2}{(n+1)^2} (2 \cot^2(\pi t) + 1) \csc^2(\pi t) \quad (56c)$$

$$|F'''_N(t)| \leq (n+1)4\pi^3 \csc^2(\pi t) + 12\pi^3 |\cot(\pi t)| \csc^2(\pi t) + \frac{12\pi^3}{n+1} (3 \cot^2(\pi t) + 1) \csc^2(\pi t) + \frac{8\pi^3}{(n+1)^2} (3 |\cot^3(\pi t)| + 2 |\cot(\pi t)|) \csc^2(\pi t). \quad (56d)$$

To continue, observe that a basic building block in the bound (56) takes the following function form, which is denoted by

$$h_\ell(t) = \cot^\ell(\pi t) \csc^2(\pi t), \quad \ell = 0, 1, 2. \quad (57)$$

The functions  $h_\ell$ ,  $\ell = 0, 1, 2$  are even, 1-periodic, and continuous, non-negative, decreasing and convex over the interval  $(0, \frac{1}{2}]$ . Define the three sums  $S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i)$ ,  $\ell = 0, 1, 2$  by

$$S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i) = \sum_{j \neq i} h_\ell(\tau_j - \tau_i + u_{i,j}). \quad (58)$$

It is straightforward to see that the terms of interest can be bounded in terms of  $S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i)$  in view of (56). We shall claim the following bound of  $S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i)$  holds, which will be proven at the end of proof:

$$S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i) \leq \left( \frac{2}{\pi \Delta(\boldsymbol{\tau})} \right)^{\ell+2} \frac{1}{(\ell+1) (1 - 2\alpha^{-1}\beta)^{\ell+1}}, \quad \ell = 0, 1, 2. \quad (59)$$

Interestingly, the above bound on  $S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i)$  does not depend on  $n$  or  $i$ . We are now ready to establish the bound (53) of interest. First, we have

$$\begin{aligned}
\max_i \sum_{j \neq i} |F_N(\tau_j - \tau_i + u_{i,j})| &\leq \sum_{j \neq i} \frac{1}{(n+1)^2} h_0(\tau_j - \tau_i + u_{i,j}) \\
&= \frac{1}{(n+1)^2} S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\leq \frac{4}{\pi^2} \frac{\Delta(\boldsymbol{\tau})^{-2}}{(n+1)^2 (1-2\alpha^{-1}\beta)} \\
&= \frac{4}{\pi^2} (\alpha - 2\beta)^{-1} \alpha ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \tag{60}
\end{aligned}$$

Similarly, for the case of the first derivative, we have

$$\begin{aligned}
\max_i \sum_{j \neq i} |F'_N(\tau_j - \tau_i + u_{i,j})| &\leq \sum_{j \neq i} \left( \frac{\pi}{n+1} h_0(\tau_j - \tau_i + u_{i,j}) + \frac{2\pi}{(n+1)^2} h_1(\tau_j - \tau_i + u_{i,j}) \right) \\
&= \frac{\pi}{n+1} S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{2\pi}{(n+1)^2} S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\leq \frac{4\Delta(\boldsymbol{\tau})^{-2}}{\pi(n+1)(1-2\alpha^{-1}\beta)} + \frac{8\Delta(\boldsymbol{\tau})^{-3}}{\pi^2(n+1)^2(1-2\alpha^{-1}\beta)^2} \\
&\leq \left( \frac{4}{\pi} (\alpha - 2\beta)^{-1} + \frac{8}{\pi^2} (\alpha - 2\beta)^{-2} \right) \alpha(n+1) ((n+1)\Delta(\boldsymbol{\tau}))^{-2}, \tag{61}
\end{aligned}$$

where the last line uses  $(n+1)\Delta(\boldsymbol{\tau}) \geq \alpha$ . Moving onto the second derivative, it follows

$$\begin{aligned}
\max_i \sum_{j \neq i} |F''_N(\tau_j - \tau_i + u_{i,j})| &\leq \sum_{j \neq i} \left( 2\pi^2 h_0(\tau_j - \tau_i + u_{i,j}) + \frac{4\pi^2}{n+1} h_1(\tau_j - \tau_i + u_{i,j}) + \frac{2\pi^2}{(n+1)^2} (2h_2(\tau_j - \tau_i + u_{i,j}) + h_0(\tau_j - \tau_i + u_{i,j})) \right) \\
&= 2\pi^2 S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{4\pi^2}{n+1} S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{2\pi^2}{(n+1)^2} (2S_2(n, \boldsymbol{\tau}, \mathbf{u}, i) + S_0(n, \boldsymbol{\tau}, \mathbf{u}, i)) \\
&\leq 2\pi^2 \left( 1 + \frac{1}{(n+1)^2} \right) S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{4\pi^2}{n+1} S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{4\pi^2}{(n+1)^2} S_2(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\leq \frac{20\pi^2}{9} S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{4\pi^2}{n+1} S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{4\pi^2}{(n+1)^2} S_2(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\leq \frac{80\Delta(\boldsymbol{\tau})^{-2}}{9(1-2\alpha^{-1}\beta)} + \frac{16\Delta(\boldsymbol{\tau})^{-3}}{\pi(n+1)(1-2\alpha^{-1}\beta)^2} + \frac{64\Delta(\boldsymbol{\tau})^{-4}}{3\pi^2(n+1)^2(1-2\alpha^{-1}\beta)^3} \\
&\leq \left( \frac{80}{9} (\alpha - 2\beta)^{-1} + \frac{16}{\pi} (\alpha - 2\beta)^{-2} + \frac{64}{3\pi^2} (\alpha - 2\beta)^{-3} \right) \alpha(n+1)^2 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \tag{62}
\end{aligned}$$

Finally, for the third derivative, it holds

$$\begin{aligned}
\max_i \sum_{j \neq i} |F'''_N(\tau_j - \tau_i + u_{i,j})| &\leq \sum_{j \neq i} \left( (n+1)4\pi^3 h_0(\tau_j - \tau_i + u_{i,j}) + 12\pi^3 h_1(\tau_j - \tau_i + u_{i,j}) \right. \\
&\quad \left. + \frac{12\pi^3}{n+1} (3h_2(\tau_j - \tau_i + u_{i,j}) + h_0(\tau_j - \tau_i + u_{i,j})) + \frac{8\pi^3}{(n+1)^2} (3h_3(\tau_j - \tau_i + u_{i,j}) + 2h_1(\tau_j - \tau_i + u_{i,j})) \right)
\end{aligned}$$

$$\begin{aligned}
&= (n+1)4\pi^3 S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + 12\pi^3 S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{12\pi^3}{n+1} (3S_2(n, \boldsymbol{\tau}, \mathbf{u}, i) + S_0(n, \boldsymbol{\tau}, \mathbf{u}, i)) \\
&\quad + \frac{8\pi^3}{(n+1)^2} (3S_3(n, \boldsymbol{\tau}, \mathbf{u}, i) + 2S_1(n, \boldsymbol{\tau}, \mathbf{u}, i)) \\
&\leq (n+1)4\pi^3 \left(1 + \frac{3}{(n+1)^2}\right) S_0(n, \boldsymbol{\tau}, \mathbf{u}, i) + 12\pi^3 \left(1 + \frac{4}{3(n+1)^2}\right) S_1(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\quad + \frac{36\pi^3}{n+1} S_2(n, \boldsymbol{\tau}, \mathbf{u}, i) + \frac{24\pi^3}{(n+1)^2} S_3(n, \boldsymbol{\tau}, \mathbf{u}, i) \\
&\leq (n+1) \frac{64}{3} \pi \frac{\Delta(\boldsymbol{\tau})^{-2}}{1-2\alpha^{-1}\beta} + \frac{2488}{27} \frac{\Delta(\boldsymbol{\tau})^{-3}}{(1-2\alpha^{-1}\beta)^2} + \frac{192}{\pi(n+1)} \frac{\Delta(\boldsymbol{\tau})^{-4}}{(1-2\alpha^{-1}\beta)^3} + \frac{192}{\pi^2(n+1)^2} \frac{\Delta(\boldsymbol{\tau})^{-5}}{(1-2\alpha^{-1}\beta)^4} \\
&\leq \left(\frac{16}{3} (\alpha-2\beta)^{-1} \pi + \frac{1488}{27} (\alpha-2\beta)^{-2} + \frac{192}{\pi} (\alpha-2\beta)^{-3} + \frac{192}{\pi^2} (\alpha-2\beta)^{-4}\right) \alpha(n+1)^3 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}.
\end{aligned} \tag{63}$$

The proof is thus completed if we can prove (59), which is the focus of the rest of the proof.

**Proof of (59)** We fix the index  $i$  and take the convention  $u_{i,i} = 0$ . As the functions  $\{h_\ell\}$  are 1-periodic, the quantity  $S_\ell(n, \boldsymbol{\tau}, \beta, i)$  is invariant by integer translations of  $\{\tau_j\}$ 's. Therefore, one can make the assumption, up to a modulo considerations and a reordering of the indices that the sequence  $\{\tau_j - \tau_i + u_{i,j}\}_j$  is within the range  $[-\frac{1}{2}, \frac{1}{2})$  and in an ascending order, so that

$$-\frac{1}{2} \leq \tau_1 - \tau_i + u_{i,1} < \tau_2 - \tau_i + u_{i,2} < \cdots < \tau_r - \tau_i + u_{i,r} < \frac{1}{2}.$$

Denote by  $r_+$  and  $r_-$  the number of positive and negative elements in the set  $\{\tau_j - \tau_i + u_{i,j}\}_j$ , respectively. As  $\tau_j - \tau_i + u_{i,j} = 0$  if and only if  $i = j$ , we have that  $r_+ + r_- = r - 1$ . Using the separation condition, and as  $h_\ell$  is decreasing over  $(0, \frac{1}{2}]$  and even, we have that

$$\begin{aligned}
0 &\leq h_\ell(\tau_{i+j} - \tau_i + u_{i,j}) \leq h_\ell\left(j\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}\right), \quad j = 1, \dots, r_+; \\
0 &\leq h_\ell(\tau_{i-j} - \tau_i + u_{i,j}) \leq h_\ell\left(-j\Delta(\boldsymbol{\tau}) + \frac{\beta}{n+1}\right) = h_\ell\left(j\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}\right), \quad j = 1, \dots, r_-.
\end{aligned}$$

We can subsequently bound the sum (58) as

$$\begin{aligned}
S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i) &= \sum_{j=1}^{r_+} h_\ell(\tau_{i+j} - \tau_i + u_{i,(i+j)}) + \sum_{j=1}^{r_-} h_\ell(\tau_{i-j} - \tau_i + u_{i,(i-j)}) \\
&\leq \sum_{j=1}^{r_+} h_\ell\left(j\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}\right) + \sum_{j=1}^{r_-} h_\ell\left(j\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}\right) \\
&\leq 2 \sum_{j=1}^{\lceil \frac{r-1}{2} \rceil} h_\ell\left(j\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}\right).
\end{aligned} \tag{64}$$

Identifying and associating the right-hand side of (64) to a Riemann sum with a mid-point rule and recalling that  $h_\ell$  is decreasing and convex over  $(0, \frac{1}{2}]$  leads to the majorant

$$\begin{aligned}
S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i) &\leq 2 \int_{\frac{\Delta(\boldsymbol{\tau})}{2} - \frac{\beta}{n+1}}^{(\lceil \frac{r-1}{2} \rceil + \frac{1}{2})\Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1}} h_\ell(u) du \\
&\leq 2 \int_{\frac{\Delta(\boldsymbol{\tau})}{2} - \frac{\beta}{n+1}}^{\frac{1}{2}} h_\ell(u) du
\end{aligned}$$

$$= 2 \left( H_\ell \left( \frac{1}{2} \right) - H_\ell \left( \frac{\Delta(\boldsymbol{\tau})}{2} - \frac{\beta}{n+1} \right) \right), \quad (65)$$

where we used the non-negativity of  $h_\ell$  and the inequality  $(\lceil \frac{r-1}{2} \rceil + \frac{1}{2}) \Delta(\boldsymbol{\tau}) - \frac{\beta}{n+1} \leq \frac{1}{2}$  in the second inequality. Here,  $\{H_\ell\}$ 's are primitives of the functions  $\{h_\ell\}$  over the interval  $(0, \frac{1}{2}]$ . Moreover, we have

$$H_\ell(t) = -\frac{1}{\pi(\ell+1)} \cot^{\ell+1}(\pi t) + c \quad (66)$$

for  $\ell = 0, 1, 2$ , where  $c$  is an arbitrary constant. This yields, with the inequality  $0 \leq \cot(\pi u) \leq \frac{1}{\pi u}$  for all  $0 < u < \frac{1}{2}$ , a further simplification of (65):

$$\begin{aligned} \Delta(\boldsymbol{\tau}) S_\ell(n, \boldsymbol{\tau}, \mathbf{u}, i) &\leq \frac{2}{\pi(\ell+1)} \cot^{\ell+1} \left( \pi \left( \frac{\Delta(\boldsymbol{\tau})}{2} - \frac{\beta}{n+1} \right) \right) \\ &\leq \left( \frac{2}{\pi} \right)^{\ell+2} \Delta(\boldsymbol{\tau})^{-\ell-1} \frac{1}{\left( 1 - \frac{2\beta\Delta(\boldsymbol{\tau})^{-1}}{n+1} \right)^{\ell+1}} \\ &\leq \left( \frac{2}{\pi} \right)^{\ell+2} \Delta(\boldsymbol{\tau})^{-\ell-1} \frac{1}{(1 - 2\alpha^{-1}\beta)^{\ell+1}}, \end{aligned} \quad (67)$$

which immediately leads to the claimed bound (59).  $\square$

## B Proof of the uniform Hessian bounds

This section is dedicated to establish Theorem 3 and Theorem 4.

### B.1 Technical lemmas

Lemma 6 and Lemma 7 provide bounds on core quantities that are involved in the decomposition of the Hessian matrix  $\mathbf{H}(\boldsymbol{\theta})$  characterized in Section 3.1.2. Their proofs are presented in Appendix B.4 and Appendix B.5, respectively.

**Lemma 6.** *Suppose that  $n \geq 2$  and let  $\boldsymbol{\tau} \subset \mathbb{T}$  be such that  $(n+1)\Delta(\boldsymbol{\tau}) \geq \alpha$  for some  $\alpha > 0$ , then there exists a constant  $K_\Delta$  with*

$$K_\Delta := \max \left\{ C_0 + \frac{3\sqrt{3}}{4\pi} C_1, \frac{3\sqrt{3}}{4\pi} C_1 + \frac{27}{16\pi^2} C_2 \right\} \quad (68)$$

where the constants  $C_0, C_1, C_2 > 0$  are defined in in (54) with parameters  $\alpha$  and  $\beta = 0$  such that

$$\|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \leq K_\Delta ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \quad (69)$$

**Lemma 7.** *Suppose that  $n \geq 2$  and let  $\boldsymbol{\tau}_* = [\tau_1^*, \dots, \tau_r^*]^\top$  and  $\boldsymbol{\tau} = [\tau_1, \dots, \tau_r]^\top$  be two vectors of points around the torus. Assume that  $(n+1)\Delta(\boldsymbol{\tau}_*) \geq \alpha$ . As long as  $(n+1)\|\boldsymbol{\tau} - \boldsymbol{\tau}_*\|_\infty \leq \beta < \frac{\alpha}{2}$ , we have*

$$\left| \left\langle \Phi(\delta'_{\tau_j}), \Phi(\mu(\boldsymbol{\theta})) - \mu(\boldsymbol{\theta}_*) \right\rangle \right| \leq (C_1 \|\mathbf{a} - \mathbf{a}^*\|_\infty + C_2 \|\mathbf{a}^*\|_\infty (n+1) \|\boldsymbol{\tau} - \boldsymbol{\tau}_*\|_\infty) (n+1) ((n+1)\Delta(\boldsymbol{\tau}))^{-2}, \quad (70a)$$

$$\left| \left\langle \Phi(\delta''_{\tau_j}), \Phi(\mu(\boldsymbol{\theta})) - \mu(\boldsymbol{\theta}_*) \right\rangle \right| \leq (C_2 \|\mathbf{a} - \mathbf{a}^*\|_\infty + C_3 \|\mathbf{a}^*\|_\infty (n+1) \|\boldsymbol{\tau} - \boldsymbol{\tau}_*\|_\infty) (n+1)^2 ((n+1)\Delta(\boldsymbol{\tau}))^{-2} \quad (70b)$$

for all  $j = 1, \dots, r$ , where the constants  $C_1, C_2, C_3 > 0$  are defined in (54) with parameters  $(\alpha, \beta)$ .

## B.2 Proof of Theorem 3

Recalling the expression of the Hessian in (32), it follows that

$$\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} = \mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} + \mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}. \quad (71)$$

We proceed to bound  $\|\mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$  and  $\|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$  separately, and then combine them via the triangle inequality.

**Step 1: bound  $\|\mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$**  From (16) and (33), we have that

$$\begin{aligned} \mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} &= \mathbf{SP} \operatorname{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)\mathbf{a}}} \right] \right)^H \mathbf{D}(\boldsymbol{\tau}) \operatorname{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)\mathbf{a}}} \right] \right) \mathbf{S}^{-1} - \mathbf{I} \\ &= \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star-1}}{A^{-2}\mathbf{a}} \right] \right)^H \mathbf{D}(\boldsymbol{\tau}) \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star}}{\mathbf{a}} \right] \right) - \mathbf{I} \\ &= \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star-1}}{A^{-2}\mathbf{a}} \right] \right)^H (\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}) \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star}}{\mathbf{a}} \right] \right) + \operatorname{diag} \left( \left[ \frac{\mathbf{1}_r}{A^{-2}|\mathbf{a}|^2} \right] \right) - \mathbf{I}. \end{aligned} \quad (72)$$

This immediately yields from the triangle inequality that

$$\begin{aligned} &\|\mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \left\| \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star-1}}{A^{-2}\odot\mathbf{a}} \right] \right)^H (\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}) \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star}}{\mathbf{a}} \right] \right) \right\|_\infty + \left\| \operatorname{diag} \left( \left[ \frac{\mathbf{1}_r}{A^{-2}\odot|\mathbf{a}|^2} \right] \right) - \mathbf{I} \right\|_\infty \\ &\leq \left\| \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star-1}}{A^{-2}\odot\mathbf{a}} \right] \right)^H \right\|_\infty \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \left\| \operatorname{diag} \left( \left[ \frac{\mathbf{a}^{\star}}{\mathbf{a}} \right] \right) \right\|_\infty + \left\| \operatorname{diag} \left( \left[ \frac{\mathbf{1}_r}{A^{-2}\odot|\mathbf{a}|^2} \right] \right) - \mathbf{I} \right\|_\infty \\ &\leq \max_j \left\{ \frac{1}{|a_j^*|}, \frac{|a_j|}{A^2} \right\} \max_j \{|a_j^*|, |a_j|\} \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + \max_j \left\{ \left| 1 - \frac{|a_j|^2}{A^2} \right| \right\}. \end{aligned} \quad (73)$$

The three maxima in (73) can be controlled using the basic relation

$$1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty \leq \frac{|a_j^*| - |a_{k,j} - a_j^*|}{|a_j^*|} \leq \frac{|a_j|}{|a_j^*|} \leq \frac{|a_j^*| + |a_{k,j} - a_j^*|}{|a_j^*|} \leq 1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty, \quad (74)$$

for  $j = 1, \dots, r$  and by exploiting the assumption  $\|\mathbf{a}^*\|_\infty \leq A$  as follows

$$\max_j \left\{ \frac{1}{|a_j^*|}, \frac{|a_j|}{A^2} \right\} \leq \max_j \left\{ \frac{1}{|a_j^*|}, \frac{|a_j|}{A^2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \right\} \leq \frac{1}{a_{\min}^*} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty), \quad (75a)$$

$$\max_j \{|a_j^*|, |a_j|\} \leq \max_j \{|a_j^*|, |a_j^*| (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)\} \leq \|\mathbf{a}^*\|_\infty (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty), \quad (75b)$$

$$\max_j \left\{ \left| 1 - \frac{|a_j|^2}{A^2} \right| \right\} \leq 1 - \min_j \left\{ \frac{|a_j|^2}{A^2} \right\} \leq 1 - \frac{(a_{\min}^*)^2}{A^2} (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty). \quad (75c)$$

The bounds (75) and Lemma 6 imply with (73) that

$$\begin{aligned} &\|\mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + 1 - \frac{(a_{\min}^*)^2 (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}{A^2} \\ &\leq \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 K_\Delta ((n+1)\Delta(\boldsymbol{\tau}))^{-2} + 1 - \frac{(a_{\min}^*)^2 (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}{A^2}, \end{aligned} \quad (76)$$

where the constant  $K_\Delta$  defined in Lemma 6 is numerically evaluated to  $K_\Delta \leq 2.32$  by setting the parameter  $\alpha = 16.5$ .



**Step 2: bound  $\|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$**  Using the block diagonal structure of the matrix  $\mathbf{E}(\boldsymbol{\theta})$  defined in (36), we have that

$$\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1} = \begin{bmatrix} \mathbf{0} & \frac{1}{\sqrt{-F_N''(0)}} \text{diag}(\mathbf{a}^{*-1}) \mathbf{E}_1(\boldsymbol{\theta}) \\ \frac{1}{\sqrt{-F_N''(0)}} \text{diag}(A^{-2}\mathbf{a}^*)^H \mathbf{E}_1(\boldsymbol{\theta}) & -\frac{1}{F_N''(0)} A^{-2} \mathbf{E}_2(\boldsymbol{\theta}) \end{bmatrix}. \quad (77)$$

Therefore, it follows

$$\begin{aligned} \|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty &\leq \frac{1}{\sqrt{-F_N''(0)}} \max_j \left\{ \max \left\{ \frac{1}{|a_j^*|}, \frac{|a_j^*|}{A^2} \right\} \left| \left\langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right| \right\} \\ &\quad - \frac{1}{F_N''(0)} \max_j \left\{ \frac{|a_j|}{A^2} \left| \left\langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right| \right\} \\ &\leq (a_{\min}^*)^{-1} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \left( \frac{1}{\sqrt{-F_N''(0)}} \max_j \left\{ \left| \left\langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right| \right\} \right. \\ &\quad \left. - \frac{1}{F_N''(0)} \max_j \left\{ \left| \left\langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right| \right\} \right), \end{aligned} \quad (78)$$

where we used the bounds (75) in the second line. From (78), it can be seen that bounding the quantity of interest amounts to controlling  $\left| \left\langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right|$  and  $\left| \left\langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^*)) \right\rangle \right|$ , which can be achieved by applying Lemma 7. Substituting the expression provided by Lemma 7 into (78) leads to

$$\begin{aligned} \|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty &\leq \left( \left( C_1 \frac{(n+1)}{\sqrt{-F_N''(0)}} + C_2 \frac{(n+1)^2}{-F_N''(0)} \right) \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} \right. \\ &\quad \left. + \left( C_2 \frac{(n+1)^2}{-F_N''(0)} + C_3 \frac{(n+1)^3}{(-F_N''(0))^{3/2}} \right) \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \\ &\quad \cdot \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty). \end{aligned} \quad (79)$$

Evaluating the constants  $C_1 \leq 2.75$ ,  $C_2 \leq 19.08$  and  $C_3 \leq 48.74$  defined in (54) with parameters  $\alpha = 16.5$  and  $\beta = \frac{\alpha}{4} = 4.125$ , altogether with the inequality  $\frac{(n+1)^2}{-F_N''(0)} \leq \frac{27}{16\pi^2}$  under the assumption  $n \geq 2$  yields

$$\begin{aligned} \|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty &\leq \left( K_a \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} + K_\tau \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \\ &\quad \cdot \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty), \end{aligned} \quad (80)$$

where the constants  $K_a$  and  $K_\tau$  are given by

$$K_a = C_1 \sqrt{\frac{27}{16\pi^2}} + C_2 \frac{27}{16\pi^2} \leq 4.40, \quad (81a)$$

$$K_\tau = C_2 \frac{27}{16\pi^2} + C_3 \left( \frac{27}{16\pi^2} \right)^{3/2} \leq 6.71. \quad (81b)$$

**Step 3: combine the bounds** The scaled Hessian matrix  $\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1}$  can be controlled over  $\mathcal{S}_k$  by the triangle inequality as follows

$$\begin{aligned} \|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty &\leq \|\mathbf{SP}(\mathbf{G}(\boldsymbol{\theta}) + \mathbf{E}(\boldsymbol{\theta}))\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \|\mathbf{SPG}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty + \|\mathbf{SPE}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty. \end{aligned} \quad (82)$$

Furthermore, we note that the inequality  $|\Delta(\boldsymbol{\tau}) - \Delta(\boldsymbol{\tau}^*)| \leq 2\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty$  holds for every  $\boldsymbol{\tau}, \boldsymbol{\tau}^* \in \mathbb{R}$ . This implies, with the assumption  $\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty \leq \frac{1}{4}\Delta(\boldsymbol{\tau}^*)$ , that  $\Delta(\boldsymbol{\tau}) \geq \frac{1}{2}\Delta(\boldsymbol{\tau}^*)$ . Substituting the bounds given in (76) and (80) yields

$$\begin{aligned}
& \|\mathbf{SPH}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\
& \leq 1 - \frac{a_{\min}^*{}^2}{A^2} (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \quad + \left( K_\Delta + K_a \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} + K_\tau \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \leq 1 - \frac{a_{\min}^*{}^2}{A^2} (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \quad + (K_\Delta + (K_a + K_\tau) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \leq 1 - \frac{a_{\min}^*{}^2}{A^2} (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \quad + (K_\Delta + (K_a + K_\tau) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} 4((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \leq 1 - \frac{a_{\min}^*{}^2}{A^2} (1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2 \\
& \quad + (4K_\Delta + K_\theta \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} ((n+1)\Delta(\boldsymbol{\tau}^*))^{-2} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2, \tag{83}
\end{aligned}$$

where we defined in the last line the constant  $K_\theta$  as

$$K_\theta = 4(K_a + K_\tau) \leq 44.42. \tag{84}$$

This concludes the proof of the theorem.  $\square$

### B.3 Proof of Theorem 4

We proceed analogously to the proof of Theorem 3 presented in Appendix B.2. We start from the expansion (32) of the Hessian matrix to get

$$\mathbf{SP}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} = \mathbf{SP}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} + \mathbf{SP}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}, \tag{85}$$

and proceed to bound  $\|\mathbf{SP}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$  and  $\|\mathbf{SP}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty$  individually before recombining them via the triangle inequality.

**Step 1: bound  $\|\mathbf{SP}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty$**  From (16) and (33), we have that

$$\begin{aligned}
\mathbf{SP}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I} &= \mathbf{SP}_k \text{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)}\mathbf{a}} \right] \right)^H \mathbf{D}(\boldsymbol{\tau}) \text{diag} \left( \left[ \frac{\mathbf{1}_r}{\sqrt{-F_N''(0)}\mathbf{a}^H} \right] \right) \mathbf{S}^{-1} - \mathbf{I} \\
&= \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^{*-1} \\ |\mathbf{a}_k|^{-2} \odot \mathbf{a} \end{array} \right] \right)^H \mathbf{D}(\boldsymbol{\tau}) \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^* \\ \mathbf{a} \end{array} \right] \right) - \mathbf{I} \\
&= \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^{*-1} \\ |\mathbf{a}_k|^{-2} \odot \mathbf{a} \end{array} \right] \right)^H (\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}) \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^* \\ \mathbf{a} \end{array} \right] \right) + \text{diag} \left( \left[ \begin{array}{c} \mathbf{1}_r \\ |\mathbf{a}_k|^{-2} \odot |\mathbf{a}|^2 \end{array} \right] \right) - \mathbf{I}. \tag{86}
\end{aligned}$$

This immediately yields from the triangle inequality

$$\begin{aligned}
& \|\mathbf{SP}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\
& \leq \left\| \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^{*-1} \\ |\mathbf{a}_k|^{-2} \odot \mathbf{a} \end{array} \right] \right)^H (\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}) \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^* \\ \mathbf{a} \end{array} \right] \right) \right\|_\infty + \left\| \text{diag} \left( \left[ \begin{array}{c} \mathbf{1}_r \\ |\mathbf{a}_k|^{-2} \odot |\mathbf{a}|^2 \end{array} \right] \right) - \mathbf{I} \right\|_\infty
\end{aligned}$$

$$\begin{aligned}
&\leq \left\| \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^{\star-1} \\ |\mathbf{a}_k|^{-2} \odot \mathbf{a} \end{array} \right] \right)^{\text{H}} \right\|_{\infty} \left\| \mathbf{D}(\boldsymbol{\tau}) - \mathbf{I} \right\|_{\infty} \left\| \text{diag} \left( \left[ \begin{array}{c} \mathbf{a}^{\star} \\ \mathbf{a} \end{array} \right] \right) \right\|_{\infty} + \left\| \text{diag} \left( \left[ \begin{array}{c} \mathbf{1}_r \\ |\mathbf{a}_k|^{-2} \odot |\mathbf{a}|^2 \end{array} \right] \right) - \mathbf{I} \right\|_{\infty} \\
&\leq \max_j \left\{ \frac{1}{|a_j^{\star}|}, \frac{|a_j^{\star}|}{|a_{k,j}|^2} \right\} \max_j \{ |a_j^{\star}|, |a_j| \} \left\| \mathbf{D}(\boldsymbol{\tau}) - \mathbf{I} \right\|_{\infty} + \max_j \left\{ \left| \frac{|a_j^{\star}|^2}{|a_{k,j}|^2} - 1 \right| \right\}
\end{aligned} \tag{87}$$

The three maxima in (87) can be controlled using the basic relation

$$\frac{|a_j^{\star}|}{|a_{k,j}|} \leq \frac{|a_j^{\star}|}{|a_j^{\star}| - |a_{k,j} - a_j^{\star}|} \leq \frac{1}{1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty}}, \tag{88}$$

as follows

$$\max_j \left\{ \frac{1}{|a_j^{\star}|}, \frac{|a_j^{\star}|}{|a_{k,j}|^2} \right\} \leq \max_j \left\{ \frac{1}{|a_j^{\star}|}, \frac{1}{|a_j^{\star}|} \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} \right\} \leq \frac{1}{a_{\min}^{\star}} \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} \tag{89a}$$

$$\max_j \{ |a_j^{\star}|, |a_j| \} \leq \max_j \{ |a_j^{\star}|, |a_j^{\star}| (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty}) \} \leq \|\mathbf{a}^{\star}\|_{\infty} (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty}) \tag{89b}$$

$$\begin{aligned}
\max_j \left\{ \left| \frac{|a_j^{\star}|^2}{|a_{k,j}|^2} - 1 \right| \right\} &\leq \max_j \left\{ 1 - \frac{|a_j^{\star}|^2}{|a_{k,j}|^2}, \frac{|a_j^{\star}|^2}{|a_{k,j}|^2} - 1 \right\} \\
&\leq \max_j \left\{ 1, \frac{|a_j^{\star}|^2}{|a_{k,j}|^2} - 1 \right\} \leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} - 1.
\end{aligned} \tag{89c}$$

The bounds (89) and Lemma 6 conclude with (87) on

$$\begin{aligned}
&\left\| \mathbf{S} \mathbf{P}_k \mathbf{G}(\boldsymbol{\theta}) \mathbf{S}^{-1} - \mathbf{I} \right\|_{\infty} \\
&\leq \frac{\|\mathbf{a}^{\star}\|_{\infty}}{a_{\min}^{\star}} \frac{1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} \left\| \mathbf{D}(\boldsymbol{\tau}) - \mathbf{I} \right\|_{\infty} + \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} - 1 \\
&\leq \frac{\|\mathbf{a}^{\star}\|_{\infty}}{a_{\min}^{\star}} \frac{1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} K_{\Delta} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} + \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} - 1,
\end{aligned} \tag{90}$$

where the constant  $K_{\Delta}$  defined in Lemma 6 is numerically evaluated to  $K_{\Delta} \leq 2.32$  by setting the parameter  $\alpha = 4.7$ .

**Step 2: bound  $\left\| \mathbf{S} \mathbf{P}_k \mathbf{E}(\boldsymbol{\theta}) \mathbf{S}^{-1} \right\|_{\infty}$**  Again, using the block diagonal structure of the matrix  $\mathbf{E}(\boldsymbol{\theta})$  defined in (36), we similarly have that

$$\mathbf{S} \mathbf{P}_k \mathbf{E}(\boldsymbol{\theta}) \mathbf{S}^{-1} = \begin{bmatrix} \mathbf{0} & \frac{1}{\sqrt{-F_N''(0)}} \text{diag}(\mathbf{a}^{\star-1}) \mathbf{E}_1(\boldsymbol{\theta}) \\ \frac{1}{\sqrt{-F_N''(0)}} \text{diag}(|\mathbf{a}_k|^{-2} \odot \mathbf{a}^{\star})^{\text{H}} \mathbf{E}_1(\boldsymbol{\theta}) & -\frac{1}{F_N''(0)} |\mathbf{a}_k|^{-2} \mathbf{E}_2(\boldsymbol{\theta}) \end{bmatrix}. \tag{91}$$

Therefore, the quantity of interest can be bounded as follows

$$\begin{aligned}
\left\| \mathbf{S} \mathbf{P}_k \mathbf{E}(\boldsymbol{\theta}) \mathbf{S}^{-1} \right\|_{\infty} &\leq \frac{1}{\sqrt{-F_N''(0)}} \max_j \left\{ \max \left\{ \frac{1}{|a_j^{\star}|}, \frac{|a_j^{\star}|}{|a_{j,k}|^2} \right\} \left| \left\langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^{\star})) \right\rangle \right| \right\} \\
&\quad - \frac{1}{F_N''(0)} \max_j \left\{ \frac{|a_j|}{|a_{j,k}|^2} \left| \left\langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^{\star})) \right\rangle \right| \right\} \\
&\leq \frac{(a_{\min}^{\star})^{-1}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^{\star})\|_{\infty})^2} \left( \frac{1}{\sqrt{-F_N''(0)}} \max_j \left\{ \left| \left\langle \delta'_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^{\star})) \right\rangle \right| \right\} \right. \\
&\quad \left. - \frac{1}{F_N''(0)} \max_j \left\{ \left| \left\langle \delta''_{\tau_j}, F_N * (\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}^{\star})) \right\rangle \right| \right\} \right),
\end{aligned} \tag{92}$$

where we used the inequalities (89) on the second line. Substituting the expression provided by Lemma 7 into (92) leads to

$$\begin{aligned} \|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty &\leq \left( \left( C_1 \frac{n+1}{\sqrt{-F_N''(0)}} + C_2 \frac{(n+1)^2}{-F_N''(0)} \right) \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} \right. \\ &\quad \left. + \left( C_2 \frac{(n+1)^2}{-F_N''(0)} + C_3 \frac{(n+1)^3}{(-F_N''(0))^{3/2}} \right) \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \\ &\quad \cdot \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}. \end{aligned} \quad (93)$$

Evaluating the constants  $C_1 \leq 3.24$ ,  $C_2 \leq 22.90$  and  $C_3 \leq 105.55$  defined in (54) with parameters  $\alpha = 4.7$  and  $\beta = \frac{\alpha}{4} = 1.2$ , altogether with the inequality  $\frac{(n+1)^2}{-F_N''(0)} \leq \frac{27}{16\pi^2}$  under the assumption  $n \geq 2$  yields

$$\|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty \leq \left( K_a \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} + K_\tau \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}, \quad (94)$$

where the constants  $K_a$  and  $K_\tau$  are given by

$$K_a = C_1 \sqrt{\frac{27}{16\pi^2}} + C_2 \frac{27}{16\pi^2} \leq 5.26, \quad (95a)$$

$$K_\tau = C_2 \frac{27}{16\pi^2} + C_3 \left( \frac{27}{16\pi^2} \right)^{3/2} \leq 11.38. \quad (95b)$$

**Step 3: combine the bounds** The scaled Hessian matrix  $\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1}$  can be controlled over  $\mathcal{S}_k$  by the triangle inequality as follows

$$\begin{aligned} \|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty &\leq \|\mathbf{S}\mathbf{P}_k(\mathbf{G}(\boldsymbol{\theta}) + \mathbf{E}(\boldsymbol{\theta}))\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \|\mathbf{S}\mathbf{P}_k\mathbf{G}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty + \|\mathbf{S}\mathbf{P}_k\mathbf{E}(\boldsymbol{\theta})\mathbf{S}^{-1}\|_\infty. \end{aligned} \quad (96)$$

Finally, the inequality  $|\Delta(\boldsymbol{\tau}) - \Delta(\boldsymbol{\tau}^*)| \leq 2\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty$  holds for every  $\boldsymbol{\tau}, \boldsymbol{\tau}^* \in \mathbb{R}$ . This implies, with the assumption  $\|\boldsymbol{\tau} - \boldsymbol{\tau}^*\|_\infty \leq \frac{1}{4}\Delta(\boldsymbol{\tau}^*)$ , that  $\Delta(\boldsymbol{\tau}) \geq \frac{1}{2}\Delta(\boldsymbol{\tau}^*)$ . Substituting the bounds given in (90) and (94) yields

$$\begin{aligned} &\|\mathbf{S}\mathbf{P}_k\mathbf{H}(\boldsymbol{\theta})\mathbf{S}^{-1} - \mathbf{I}\|_\infty \\ &\leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \\ &+ \left( K_\Delta (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) + K_a \frac{\|\mathbf{a} - \mathbf{a}^*\|_\infty}{\|\mathbf{a}^*\|_\infty} + K_\tau \sqrt{-F_N''(0)} \|\boldsymbol{\tau}_k - \boldsymbol{\tau}^*\|_\infty \right) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \\ &\leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \\ &\quad + (K_\Delta (1 + \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) + (K_a + K_\tau) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \\ &\leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \\ &\quad + (K_\Delta + (K_\Delta + K_a + K_\tau) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \\ &\leq \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 \end{aligned}$$

$$\begin{aligned}
& + (K_\Delta + (K_\Delta + K_a + K_\tau) \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{4((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} \\
\leq & \frac{1}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2} - 1 + (4K_\Delta + K_\theta \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty) \frac{\|\mathbf{a}^*\|_\infty}{a_{\min}^*} \frac{((n+1)\Delta(\boldsymbol{\tau}^*))^{-2}}{(1 - \|\mathbf{S}(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)\|_\infty)^2}, \tag{97}
\end{aligned}$$

where defined the constant  $K_\theta$  as

$$K_\theta = 4(K_\Delta + K_a + K_\tau) \leq 75.80. \tag{98}$$

This concludes the proof of the theorem.  $\square$

## B.4 Proof of Lemma 6

Leveraging the block structure (34) of the matrix  $\mathbf{D}(\boldsymbol{\tau})$ , it boils down to controlling

$$\|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \leq \max\{\|\mathbf{D}_0(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + \|\mathbf{D}_1(\boldsymbol{\tau})\|_\infty, \|\mathbf{D}_2(\boldsymbol{\tau}) - \mathbf{I}\|_\infty + \|\mathbf{D}_1(\boldsymbol{\tau})\|_\infty\}, \tag{99}$$

which can be accomplished with the aid of Lemma 5. Specifically, recalling the expressions in (35), applying Lemma 5 with parameters  $\alpha$  and  $\beta = 0$ , and noticing that the inequality  $\frac{(n+1)^2}{-F_N''(0)} \leq \frac{27}{16\pi^2}$  holds whenever  $n \geq 2$ , the quantities of interest in (99) can be controlled as follows

$$\|\mathbf{D}_0(\boldsymbol{\tau}) - \mathbf{I}\|_\infty = \max_i \sum_{j \neq i} |F_N(\tau_j - \tau_i)| \leq C_0 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}, \tag{100a}$$

$$\begin{aligned}
\|\mathbf{D}_1(\boldsymbol{\tau})\|_\infty &= \frac{1}{\sqrt{-F_N''(0)}} \max_i \sum_j |F_N'(\tau_j - \tau_i)| \leq C_1 \frac{n+1}{\sqrt{-F_N''(0)}} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} \\
&\leq \frac{3\sqrt{3}}{4\pi} C_1 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}, \tag{100b}
\end{aligned}$$

$$\begin{aligned}
\|\mathbf{D}_2(\boldsymbol{\tau}) - \mathbf{I}\|_\infty &= \frac{1}{-F_N''(0)} \max_i \sum_{j \neq i} |F_N''(\tau_j - \tau_i)| \leq C_2 \frac{(n+1)^2}{-F_N''(0)} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} \\
&\leq \frac{27}{16\pi^2} C_2 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \tag{100c}
\end{aligned}$$

Further substituting (100) into (99) leads to

$$\|\mathbf{D}(\boldsymbol{\tau}) - \mathbf{I}\|_\infty \leq \max\left\{C_0 + \frac{3\sqrt{3}}{4\pi} C_1, \frac{3\sqrt{3}}{4\pi} C_1 + \frac{27}{16\pi^2} C_2\right\} ((n+1)\Delta(\boldsymbol{\tau}))^{-2} \tag{101}$$

which leads to the desired statement.  $\square$

## B.5 Proof of Lemma 7

For the first inequality, following the definition, we have that

$$\begin{aligned}
\left| \left\langle \Phi(\delta'_{\tau_j}), \Phi(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}_*)) \right\rangle \right| &= \left| \left\langle \delta'_{\tau_j}, \Phi^* (\Phi(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}_*))) \right\rangle \right| \\
&= \left| \sum_{\ell=1}^r a_\ell F_N'(\tau_j - \tau_\ell) - \sum_{\ell=1}^{r^*} a_\ell^* F_N'(\tau_j - \tau_\ell^*) \right| \\
&= \left| \sum_{\ell=1}^r (a_\ell - a_\ell^*) F_N'(\tau_j - \tau_\ell) + \sum_{\ell=1}^r a_\ell^* (F_N'(\tau_j - \tau_\ell) - F_N'(\tau_j - \tau_\ell^*)) \right| \\
&\leq \left| \sum_{\ell=1}^r (a_\ell - a_\ell^*) F_N'(\tau_j - \tau_\ell) \right| + \left| \sum_{\ell=1}^r a_\ell^* (F_N'(\tau_j - \tau_\ell) - F_N'(\tau_j - \tau_\ell^*)) \right|, \tag{102}
\end{aligned}$$

where the last line used the triangle inequality. To proceed, we control the two terms separately. Using Hölder's inequality and Lemma 5, we obtain

$$\begin{aligned} \left| \sum_{\ell=1}^r (a_\ell - a_\ell^*) F'_N(\tau_j - \tau_\ell) \right| &\leq \| \mathbf{a} - \mathbf{a}^* \|_\infty \sum_{\ell=1}^r |F'_N(\tau_j - \tau_\ell)| \\ &\leq C_1 (n+1) \| \mathbf{a} - \mathbf{a}^* \|_\infty ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \end{aligned} \quad (103)$$

Next, the second term can be bounded by applying Hölder's inequality, the mean-value theorem and Lemma 5 with parameter  $\beta \geq (n+1) \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty$  as follows

$$\begin{aligned} \left| \sum_{\ell=1}^r a_\ell^* (F'_N(\tau_j - \tau_\ell) - F'_N(\tau_j - \tau_\ell^*)) \right| &\leq \| \mathbf{a}^* \|_\infty \sum_{\ell=1}^r |F'_N(\tau_j - \tau_\ell) - F'_N(\tau_j - \tau_\ell^*)| \\ &\leq \| \mathbf{a}^* \|_\infty \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty \sum_{\ell=1}^r \sup_{|u_\ell| \leq \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty} |F''_N(\tau_j - \tau_\ell + u_\ell)| \\ &\leq \| \mathbf{a}^* \|_\infty \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty C_2 (n+1)^2 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \end{aligned} \quad (104)$$

Plugging the previous two bounds into the inequality (102) reduces to

$$\begin{aligned} &\left| \left\langle \Phi(\delta'_{\tau_j}), \Phi(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}_*)) \right\rangle \right| \\ &\leq (C_1 \| \mathbf{a} - \mathbf{a}^* \|_\infty + C_2 \| \mathbf{a}^* \|_\infty (n+1) \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty) (n+1) ((n+1)\Delta(\boldsymbol{\tau}))^{-2}. \end{aligned} \quad (105)$$

Moreover, we may show that through analogous reasoning that

$$\begin{aligned} &\left| \left\langle \Phi(\delta''_{\tau_j}), \Phi(\mu(\boldsymbol{\theta}) - \mu(\boldsymbol{\theta}_*)) \right\rangle \right| \\ &\leq (C_2 \| \mathbf{a} - \mathbf{a}^* \|_\infty + C_3 \| \mathbf{a}^* \|_\infty (n+1) \| \boldsymbol{\tau} - \boldsymbol{\tau}^* \|_\infty) (n+1)^2 ((n+1)\Delta(\boldsymbol{\tau}))^{-2}, \end{aligned} \quad (106)$$

which concludes the proof.  $\square$