



**HAL**  
open science

# The basins of attraction of the global minimizers of the non-convex sparse spike estimation problem

Yann Traonmilin, Jean-François Aujol

► **To cite this version:**

Yann Traonmilin, Jean-François Aujol. The basins of attraction of the global minimizers of the non-convex sparse spike estimation problem. *Inverse Problems*, 2020, 10.1088/1361-6420/ab5aa3 . hal-01938239v3

**HAL Id: hal-01938239**

**<https://hal.science/hal-01938239v3>**

Submitted on 5 Sep 2019

**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.

# The basins of attraction of the global minimizers of the non-convex sparse spike estimation problem

Yann Traonmilin<sup>1,2,\*</sup> and Jean-François Aujol<sup>2</sup>

<sup>1</sup>CNRS,

<sup>2</sup>Univ. Bordeaux, Bordeaux INP, CNRS, IMB, UMR 5251, F-33400 Talence, France.

## Abstract

The sparse spike estimation problem consists in estimating a number of off-the-grid impulsive sources from under-determined linear measurements. Information theoretic results ensure that the minimization of a non-convex functional is able to recover the spikes for adequately chosen measurements (deterministic or random). To solve this problem, methods inspired from the case of finite dimensional sparse estimation where a convex program is used have been proposed. Also greedy heuristics have shown nice practical results. However, little is known on the ideal non-convex minimization method. In this article, we study the shape of the global minimum of this non-convex functional: we give an explicit basin of attraction of the global minimum that shows that the non-convex problem becomes easier as the number of measurements grows. This has important consequences for methods involving descent algorithms (such as the greedy heuristic) and it gives insights for potential improvements of such descent methods.

## 1 Introduction

### 1.1 Context

Sums of sparse off-the-grid spikes can be used to model impulsive sources in signal processing (e.g. in astronomy, microscopy,...). Estimating such signals from a finite number of Fourier measurements is known as the super-resolution problem [9]. Also, the estimation of spikes from random Fourier measurements is at the core of the compressive  $K$ -means algorithm where  $k$ -means cluster centers are estimated from a compressed database [21]. In the space  $\mathcal{M}$  of finite signed measure over  $\mathbb{R}^d$ , we aim at recovering  $x_0 = \sum_{i=1,k} a_i \delta_{t_i}$  from the measurements

$$y = Ax_0 + e, \tag{1}$$

---

\*Contact author : [yann.traonmilin@math.u-bordeaux.fr](mailto:yann.traonmilin@math.u-bordeaux.fr)

where  $\delta_{t_i}$  is the Dirac measure at position  $t_i$ , the operator  $A$  is a linear observation operator,  $y \in \mathbb{C}^m$  are the  $m$  noisy measurements and  $e$  is a finite energy observation noise. Recent works have shown that it is possible to estimate spikes from a finite number of adequately chosen Fourier measurements as long as their locations are sufficiently separated, using convex minimization based variational methods in the space of measures [8, 2, 25, 12, 14]. Other general studies on inverse problems have shown that an ideal non-convex method (unfortunately computationally inefficient) can be used to recover these signals as long as the linear measurement operator has a restricted isometry property (RIP) [5]. In the case of super-resolution, adequately chosen random compressive measurements have been shown to meet the sufficient RIP conditions for separated spikes, thus guaranteeing the success of the ideal non-convex decoder [19]. These RIP results are based on an adequate kernel metric on  $\mathcal{M}$ . It must be noted that, according to the work of [5], the success of the convex decoders as described in [8] for regular Fourier sampling implies a (lower) restricted isometry property of  $A$  with respect to such a kernel metric (and not with the natural total variation metric: in this case no RIP is possible with finite regular Fourier measurements, see e.g. [6]). Greedy heuristics have also been proposed to approach the non-convex minimization problem and they have shown good practical utility [20, 21, 27].

While giving theoretical recovery guarantees, the convex-based method is non-convex in the space of parameters (amplitudes and locations) due to a polynomial root finding step. Also, it is difficult to implement in dimensions larger than one in practice [13]. Greedy heuristics based on orthogonal matching pursuit are implemented in higher dimension (they can practically be used up to  $d = 50$ ), but they still miss theoretical recovery guarantees [20]. It would be possible to overcome the limitations of such methods if it were possible to perform the ideal non-convex minimization:

$$x^* \in \underset{x \in \Sigma}{\operatorname{argmin}} \|Ax - y\|_2 \quad (2)$$

where  $\Sigma$  is a low-dimensional set modeling the separation constraints on the  $k$  Diracs. Theoretical recovery guarantees for this minimization have been given in [19]. While simple in its formulation, properties of this minimization procedure have not yet been thoroughly studied.

In this article, as a first important step towards the understanding of the non-convex sparse spike estimation problem (2), we study its formulation in the parameter space (the space of amplitudes and locations of the Diracs). We observe that a smooth non-convex optimization can be performed.

We place ourselves in a context where the number of measurements, either deterministic or random, guarantees the success of the ideal non-convex decoder with respect to a kernel metric  $\|\cdot\|_h$ , i.e. when we can ensure that:

$$\|x^* - x_0\|_h \leq C\|e\|_2, \quad (3)$$

where  $C$  is an absolute constant with respect to  $e$  and  $x_0 \in \Sigma_{k,\epsilon}$ , the set  $\Sigma_{k,\epsilon}$  is the set of sums of  $k$  spikes separated by  $\epsilon$  on a given bounded domain. Qualitatively, the kernel

metric can be viewed as a measure of the energy at a given resolution set by a kernel  $h$  (see Section 2.3).

The bound (3) is guaranteed by a restricted isometry property of  $A$  defined using such kernel metric [19]. This RIP setting is verified in the deterministic (see Section 2.3) and random weighted Fourier measurement contexts [19]. We link this RIP of measurement operators with the conditioning of the Hessian of the global minimum, and we give an explicit basin of attraction of the global minimum in the parameter space. This study has direct consequences for the theoretical study of greedy approaches. Indeed a basin of attraction permits to give recovery guarantees for the gradient descent (the initialization must fall within the basin), which is a step in the iterations of the greedy approach.

## 1.2 Parametrization of the model set $\Sigma$

Let  $\Sigma \subset \mathcal{M}$  be a model set (union of subspaces) and  $x_0 \in \Sigma$ . Let  $f(x) = \|Ax - y\|_2$ .

**Definition 1.1** (Parametrization of  $\Sigma$ ). *A parametrization of  $\Sigma$  is a function  $\phi$  such that  $\Sigma \subset \phi(\mathbb{R}^d) = \{\phi(\theta) : \theta \in \mathbb{R}^d\}$ .*

**Definition 1.2** (Local minimum). *The point  $\theta \in \mathbb{R}^d$  is a local minimum of  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  if there is  $\epsilon > 0$  such that for any  $\theta' \in \mathbb{R}^d$  such that  $\|\theta - \theta'\|_2 \leq \epsilon$ , we have  $g(\theta) \leq g(\theta')$ .*

In the following, we consider the model of  $\epsilon$ -separated Diracs with  $\epsilon > 0$ :

$$\Sigma = \Sigma_{k,\epsilon} := \left\{ \phi(\theta) = \sum_{r=1,k} a_r \delta_{t_r} : \theta = (a, t_1, \dots, t_k) \in \mathbb{R}^{k(d+1)}, a \in \mathbb{R}^k, t_r \in \mathbb{R}^d, \right. \\ \left. \forall r \neq l, \|t_r - t_l\|_2 > \epsilon, t_r \in \mathcal{B}_2(R) \right\}, \quad (4)$$

where

$$\mathcal{B}_2(R) = \{t \in \mathbb{R}^d : \|t\|_2 \leq R\}. \quad (5)$$

Note that, in this paper, the Dirac distributions could be supported on any compact set. We use  $\mathcal{B}_2(R)$  for the sake of simplicity. For  $t_r \in \mathbb{R}^d$ , we write  $t_r = (t_{r,j})_{j=1,d}$ .

We consider the following parametrization of  $\Sigma_{k,\epsilon}$ :  $\sum_{i=1,k} a_i \delta_{t_i} = \phi(\theta)$  with  $\theta = (a_1, \dots, a_k, t_1, \dots, t_k)$ . We define

$$\Theta_{k,\epsilon} := \phi^{-1}(\Sigma_{k,\epsilon}). \quad (6)$$

We consider the problem

$$\theta^* \in \arg \min_{\theta \in E} g(\theta) = \arg \min_{\theta \in E} \|A\phi(\theta) - y\|_2. \quad (7)$$

where  $E = \mathbb{R}^{k(d+1)}$  or  $E = \Theta_{k,\epsilon}$  and  $g(\theta) = f(\phi(\theta))$ .

Note that when  $E = \Theta_{k,\epsilon}$ , performing minimization (7) allows to recover the minima of the ideal minimization (2), yielding stable recovery guarantees under a RIP assumption. Hence we are particularly interested in this case. When  $E = \mathbb{R}^{k(d+1)}$ , we speak about unconstrained minimization for minimization (7).

The objective of this paper is to study the shape of the basin of attraction of the global minimum of (7) when  $E = \Theta_{k,\epsilon}$ .

### 1.3 Basin of attraction and descent algorithms

In this work, we are interested in minimizing  $g$  defined in (7). Since  $g$  is a smooth function, a classical method to minimize  $g$  is to consider a fixed step gradient descent. The algorithm is the following. Consider an initial point  $\theta_0 \in \mathbb{R}^d$  and a step size  $\tau > 0$ . We define by recursion the sequence  $\theta_n$  by

$$\theta_{n+1} = \theta_n - \tau \nabla g(\theta_n) \quad (8)$$

Such algorithm is used as a refinement step in the greedy heuristic based on orthogonal matching pursuit [21, Algorithm 1, Step 5] in the practical setting of compressive statistical learning.

We now give the definition of basin of attraction that we will use in this paper.

**Definition 1.3** (Basin of attraction). *We say that a set  $\Lambda \subset \mathbb{R}^d$  is a basin of attraction of  $g$  if there exists  $\theta^* \in \Lambda$  and  $\tau > 0$ , such that if  $\theta_0 \in \Lambda$  then the sequence  $\theta_n$  defined by (8) converges to  $\theta^*$ .*

This definition of basin of attraction is related to the following classical optimization result (see e.g. [11]):

**Proposition 1.1.** *Assume  $g$  to be a smooth coercive convex function, whose gradient is  $L$  Lipschitz. Let  $\theta_0 \in \mathbb{R}^d$ . Then, if  $\tau < \frac{1}{L}$ , there exists  $\theta^* \in \mathbb{R}^d$  such that the sequence  $\theta_n$  defined by (8) converges to  $\theta^*$ .*

An immediate consequence of the previous proposition is the following corollary.

**Corollary 1.1.** *Assume  $g$  to be a smooth function. Assume that  $g$  has a minimizer  $\theta^* \in \mathbb{R}^d$ . Assume that there exists an open set  $\Lambda \subset \mathbb{R}^d$  such that  $\theta^* \in \Lambda$ ,  $g$  is convex on  $\Lambda$  with  $L$  Lipschitz gradient. Assume also that the sequence  $\theta_n$  generated by the descent algorithm remains in  $\Lambda$ . Then, if  $\theta_0 \in \Lambda$  and  $\tau < \frac{1}{L}$ , the sequence  $\theta_n$  defined by (8) converges to  $\theta^*$ .*

**Remark 1.1.** *Assume that  $g$  is in  $\mathcal{C}^2$ . Let  $\lambda_{\max}(t)$  the largest eigenvalue of the Hessian matrix of  $g(t)$ . Let  $\Theta \subset \mathbb{R}^d$  an open set. If there exists  $L > 0$  such that for all  $t$  in  $\Theta$ ,  $\lambda_{\max}(t) \leq L$ , then  $g$  has a  $L$  Lipschitz gradient in  $\Theta$ .*

It is not obvious that the unconstrained gradient descent defined in iterations (8) and the corresponding notion of basin of attraction is suitable to perform constrained minimization (7). In fact, we show in this paper (essentially through Lemma 3.1) that the global minimum of constrained minimization (7) has a basin of attraction.

### 1.4 Related work

While original for the sparse spike estimation problem, it must be noted that the study of non-convex optimization schemes for linear inverse problems has gained attraction recently for different kinds of low-dimensional models. For low-rank matrix estimation, a smooth parametrization of the problem is possible and it has been shown that a

RIP guarantees the absence of spurious minima [29, 3]. In [28], a model for phase recovery with alternated projections and smart initialization is considered. Conditions on the number of measurements guarantee the success of the technique. In the area of blind deconvolution and bi-convex programming, recent works have exploited similar ideas [22, 7].

In the case of super-resolution, the idea of gradient descent has been studied in an asymptotic regime ( $k \rightarrow \infty$ ) in [10] with theoretical conditions based on Wasserstein gradient flow for the initialization. In our case, we study the particular super-resolution problem with a fixed number of impulsions and we place ourselves in conditions where stable recovery is guaranteed, leading to explicit conditions on the initialization.

The objective of this article is to investigate to what extent these ideas can be applied to the theoretical study of the case of spike super-resolution estimation.

The question of projected gradient descent raised in the last Section has been explored for general low-dimensional models [4]. It has been shown that the RIP guarantees the convergence of such algorithms with an ideal (often non practical) projection. Approached projected gradient descents have also been studied and shown to be successful for some particular applications [17]. The spikes super-resolution problem adds the parametrization step to these problems.

## 1.5 Contributions and organization of the paper

After a precise description of the setting, the definition of the kernel metric of interest and the associated restricted isometry for the spike estimation problem at the beginning of Section 2, this article gives the following original results:

1. A bound on the conditioning of the Hessian at a global minimum of the minimization in the parameter space is given in Section 2. This bound shows that the better RIP constants are (RIP constants improve with respect to the number of measurements), the better the non-convex minimization problem behaves. It also shows that there is a basin of attraction of the global optimum where no separation constraints are needed (for descent algorithms with an initialization close to the minimum, separation constraints can be discarded).
2. An explicit shape of the basin of attraction of global minima is given in Section 3. The size of the basin of attraction increases when the RIP constant gets better.

To conclude, we discuss the role of the separation constraint in descent algorithms in Section 4, and we explain why enforcing a separation might improve them.

## 2 Conditioning of the Hessian

This section is devoted to the study of the Hessian matrix of  $g$ . In particular, we provide a bound on the conditioning of the Hessian at a global minimum of the minimization in the parameter space.

## 2.1 Notations

The operator  $A$  is a linear operator modeling  $m$  measurements in  $\mathbb{C}^m$  ( $\text{Im}A \subset \mathbb{C}^m$ ) on the space of measures on  $\mathbb{R}^d$  defined by: for  $l = 1, m$ ,

$$(Au)_l = \int_{\mathbb{R}^d} \alpha_l(t) du(t) \quad (9)$$

where  $(\alpha_l)_l$  is a collection of functions in  $\mathcal{C}^2(\mathcal{B}_2(R))$  (twice continuously differentiable functions on  $\mathcal{B}_2(R)$  defined in (5)).

Notice that the integral used in (9) is in fact a duality product  $\langle u, \alpha_l \rangle$  between a function in  $\mathcal{C}^2(\mathcal{B}_2(R))$  and a finite signed measure over  $\mathbb{R}^d$ . As the  $\alpha_l$  are in  $\mathcal{C}^2(\mathcal{B}_2(R))$ , we can similarly apply  $A$  to distributions of order 1 and 2 with support included in the relative interior of  $\mathcal{B}_2(R)$  which we note  $\text{rint}\mathcal{B}_2(R)$ .

While a lot of results for spike super-resolution are expressed on the  $d$ -dimensional Torus  $\mathbb{T}^d$ , we prefer the setting of Diracs with bounded support on  $\mathbb{R}^d$  which is often closer to the physics of the considered phenomenon. However, our work is directly extended to the Torus setting by replacing  $\mathbb{R}^d$  by  $\mathbb{T}^d$  and  $\mathcal{B}^2(R)$  by  $\mathbb{T}^d$ .

In  $\mathbb{C}^m$ , we consider the Hermitian product  $\langle x, y \rangle = \sum x_i \bar{y}_i$ . An example of such measurement operator is the (weighted) Fourier sampling:  $(Au)_l = \frac{1}{\sqrt{m}} \int_{\mathbb{R}^d} c_l e^{-j(\omega_l, t)} du(t)$  for some chosen frequencies  $\omega_l \in \mathbb{R}^d$  and frequency dependent weights  $c_l \in \mathbb{R}$ .

Let  $x = \sum_{i=1, k} a_i \delta_{t_i}$ . By linearity of  $A$ , we have

$$(Ax)_l = \sum_{i=1}^k (A\delta_{t_i})_l = \sum_{i=1}^k a_i \alpha_l(t_i). \quad (10)$$

With  $g(\theta) = f(\phi(\theta)) = \|A\phi(\theta) - y\|_2^2$ , we get:

$$g(\theta) = \sum_{l=1}^m \left| \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right|^2. \quad (11)$$

In the following, the notion of directional derivative will be important.

**Definition 2.1** (Directional derivatives). *Let  $f$  be a  $\mathcal{C}^1$  function, and  $v \in \mathbb{R}^d$  such that  $\|v\|_2 = 1$ . Then we can define the directional derivative of  $f$  in direction  $v$  by:*

$$f'_v(t) := \langle v, \nabla f(t) \rangle = \lim_{h \rightarrow 0^+} \frac{f(t + hv) - f(t)}{h} \quad (12)$$

*Let  $f$  be a  $\mathcal{C}^2$  function, and  $(v_1, v_2) \in \mathbb{R}^{2d}$  such that  $\|v_1\|_2 = \|v_2\|_2 = 1$ . Then we can define the second order directional derivative of  $f$  in directions  $v_1$  and  $v_2$  by:*

$$f''_{v_1, v_2}(t) := \langle v_1, \nabla^2 f(t) v_2 \rangle \quad (13)$$

*Notice that of course  $f''_{v_1, v_2}(t) = f''_{v_2, v_1}(t)$ . If  $v_1 = v_2$ , we write  $f''_{v_1}(t) := f''_{v_1, v_1}(t)$*

In particular, they permit to introduce derivatives of Dirac measures supported on  $\mathbb{R}^d$ .

**Definition 2.2** (Directional derivatives of Dirac). *Let  $v \in \mathbb{R}^d$  such that  $\|v\|_2 = 1$ . The distribution  $\delta'_{t_0,v}$  is defined by  $\langle \delta'_{t_0,v}, f \rangle = -f'_v(t_0)$ . It is the limit of  $\nu_\eta = -\frac{\delta_{t_0+\eta v} - \delta_{t_0}}{\eta}$  for  $\eta \rightarrow 0^+$  in the distributional sense : for all  $h \in \mathcal{C}^1(\mathbb{R}^d)$ ,  $\int_{\mathbb{R}} h(t) d\nu_\eta(t) \rightarrow_{\eta \rightarrow 0^+} \langle \delta'_{t_0,v}, h \rangle$ .*

*Similarly, the distribution  $\delta''_{t_0,v}$  is defined by  $\langle \delta''_{t_0,v}, f \rangle = f''_v(t_0)$  for  $f \in \mathcal{C}^2(\mathbb{R}^d)$  and the distribution  $\delta''_{t_0,v_1,v_2}$  is defined by  $\langle \delta''_{t_0,v_1,v_2}, f \rangle = f''_{v_1,v_2}(t_0)$  for  $f \in \mathcal{C}^2(\mathbb{R}^d)$  where  $f''_{v_1,v_2}$  is the derivative of  $f$  in direction  $v_1$  chained with the derivative of  $f$  in direction  $v_2$ .*

*When  $v = e_i$  is a vector of the canonical basis of  $\mathbb{R}^d$ , we write  $\delta'_{t_0,i} = \delta'_{t_0,e_i}$  and  $\delta''_{t_0,i} = \delta''_{t_0,e_i,e_i}$ .*

We now have the necessary tools to start the study of the Hessian of  $g$ .

## 2.2 Gradient and Hessian of the objective function $g$

We calculate the gradient and Hessian of  $g$  in the two following propositions. We start with the gradient of  $g$ .

**Proposition 2.1.** *For any  $\theta \in \mathbb{R}^{2k}$ , we have:*

$$\frac{\partial g(\theta)}{\partial a_r} = 2\mathcal{R}e\langle A\delta_{t_r}, A\phi(\theta) - y \rangle, \quad (14)$$

$$\frac{\partial g(\theta)}{\partial t_{r,j}} = -2a_r \mathcal{R}e\langle A\delta'_{t_r,j}, A\phi(\theta) - y \rangle. \quad (15)$$

*Proof.* See Appendix A.1. □

The next proposition gives the values of the Hessian matrix of  $g$  which has a simple expression with the use of derivatives of Diracs.

**Proposition 2.2.** *For any  $\theta \in \mathbb{R}^{k(d+1)}$*

$$H_{1,r,s} = \frac{\partial^2 g(\theta)}{\partial a_r \partial a_s} = 2\mathcal{R}e\langle A\delta_{t_r}, A\delta_{t_s} \rangle. \quad (16)$$

$$H_{2,r,j_1,s,j_2} = \frac{\partial^2 g(\theta)}{\partial t_{r,j_1} \partial t_{s,j_2}} = 2a_r a_s \mathcal{R}e\langle A\delta'_{t_r,j_1}, A\delta'_{t_s,j_2} \rangle + \mathbf{1}(r=s) 2a_r \mathcal{R}e\langle A\delta''_{t_r,j_1,j_2}, A\phi(\theta) - y \rangle. \quad (17)$$

$$H_{12,r,s,j} = \frac{\partial^2 g(\theta)}{\partial a_r \partial t_{s,j}} = -2a_s \mathcal{R}e\langle A\delta_{t_r}, A\delta'_{t_s,j} \rangle - \mathbf{1}(r=s) 2\mathcal{R}e\langle A\delta'_{t_s,j}, A\phi(\theta) - y \rangle. \quad (18)$$



Hence the Hessian can be decomposed as the sum of two matrices  $H = G + F$  with

$$\begin{aligned} G_{1,r,s} &= 2\mathcal{R}e\langle A\delta_{t_r}, A\delta_{t_s}\rangle, \\ G_{2,r,j_1,s,j_2} &= 2a_r a_s \mathcal{R}e\langle A\delta'_{t_r,j_1}, A\delta'_{t_s,j_2}\rangle, \\ G_{12,r,s,j} &= -2a_s \mathcal{R}e\langle A\delta_{t_r}, A\delta'_{t_s,j}\rangle. \end{aligned} \quad (19)$$

and

$$\begin{aligned} F_{1,r,s} &= 0, \\ F_{2,r,j_1,s,j_2} &= \mathbf{1}(r = s)2a_r \mathcal{R}e\langle A\delta''_{t_r,j_1,j_2}, A\phi(\theta) - y\rangle, \\ F_{12,r,s,j} &= -\mathbf{1}(r = s)2\mathcal{R}e\langle A\delta'_{t_s,j}, A\phi(\theta) - y\rangle. \end{aligned} \quad (20)$$

*Proof.* See Appendix A.1. □

### 2.3 Kernel, dipoles and the RIP

In order to be able to build an operator  $A$  with a RIP, we define a reproducible kernel Hilbert space (RKHS) structure on the space of measures as in [19], see also [24]. The natural metric on the space of finite signed measures, the total variation of measures, is not well suited for a RIP analysis of the spikes super-resolution problems, as it does not measure the spacing between Diracs. When using the RIP, fundamental objects appear in the calculations: dipoles of Diracs. In this section we show that the typical RIP implies a RIP on dipoles and their generalization.

**Definition 2.3** (Kernel, scalar product and norm). *For finite signed measures over  $\mathbb{R}^d$ , the Hilbert structure induced by a kernel  $h$  (a smooth function from  $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ ) is defined by the following scalar product between 2 measures  $\pi_1, \pi_2$*

$$\langle \pi_1, \pi_2 \rangle_h = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} h(t_1, t_2) d\pi_1(t_1) d\pi_2(t_2). \quad (21)$$

We can consequently define

$$\|\pi_1\|_h^2 = \langle \pi_1, \pi_1 \rangle_h. \quad (22)$$

We have the relation

$$\|\pi_1 + \pi_2\|_h^2 = \|\pi_1\|_h^2 + 2\langle \pi_1, \pi_2 \rangle_h + \|\pi_2\|_h^2. \quad (23)$$

Measuring distances with the help of  $\|\cdot\|_h$  can be viewed as measuring distances at a given resolution set by  $h$ . Typically we use Gaussian kernels where the sharper the kernel is, the more accurate it is.

The next definition is taken from [19].

**Definition 2.4** (( $\epsilon$ -)Dipole, separation). *An  $\epsilon$ -dipole (noted dipole for simplicity) is a measure  $\pi = a_1\delta_{t_1} - a_2\delta_{t_2}$  where  $\|t_1 - t_2\|_2 \leq \epsilon$ . Two dipoles  $\pi_1 = a_1\delta_{t_1} - a_2\delta_{t_2}$  and  $\pi_2 = a_3\delta_{t_3} - a_4\delta_{t_4}$  are  $\epsilon$ -separated if their support are strictly  $\epsilon$ -separated (with respect to the  $\ell^2$ -norm on  $\mathbb{R}^d$ ), i.e. if  $\|t_1 - t_3\|_2 > \epsilon$ ,  $\|t_2 - t_3\|_2 > \epsilon$  and  $\|t_1 - t_4\|_2 > \epsilon$  and  $\|t_2 - t_4\|_2 > \epsilon$ .*

Compared to [19], we need to introduce a new definition.

**Definition 2.5** (Generalized dipole). *A generalized dipole  $\nu$  is either a dipole or a distribution of order 1 of the form  $a_1\delta_t + a_2\delta'_{t,v}$ . Two generalized dipoles are  $\epsilon$ -separated if their support are strictly  $\epsilon$ -separated (with respect to the  $\ell^2$ -norm on  $\mathbb{R}^d$ ).*

In this article we use regular, symmetrical, translation invariant kernels. Most recent developments to non translation invariant kernels [23] could be considered to generalize this work, but they are out of the scope of this article for the sake of simplicity.

**Assumption 2.1.** *A kernel  $h$  follows this assumption if*

- $h \in \mathcal{C}^2(\mathbb{R}^d, \mathbb{R}^d)$ .
- $h$  is symmetrical with respect to 0, translation invariant, i.e. we can write  $h(t_1, t_2) = \rho(\|t_1 - t_2\|_2)$  where  $\rho \in \mathcal{C}^2(\mathbb{R})$ .
- $h(t, t) = \rho(0) = 1 = \max_{t \in \mathbb{R}^d, s \in \mathbb{R}^d} |h(t, s)|$ ,  $\rho'(0) = 0$ , and  $\rho''(0) < 0$ .
- there is a constant  $c_h$  such that  $0 < c_h \leq \frac{\epsilon}{2}$  and  $\rho(t) \leq 1 - \frac{|\rho''(0)|}{2}t^2$  for  $t \in [0, c_h]$  (the existence of  $c_h$  is a consequence of previous assumptions).
- there is a constant  $\mu_h$  such that, for all two  $\epsilon$ -separated dipoles,  $\langle \nu_1, \nu_2 \rangle_h \leq \mu_h \|\nu_1\|_h \|\nu_2\|_h$  (mutual coherence).

Note that the assumption that  $h \in \mathcal{C}^2$  guarantees the existence of integrals with respect to finite signed measures and duality product with distribution of order 1 with bounded supports.

**Example** The now almost canonical well behaved kernel is the Gaussian kernel. From [19], for  $\epsilon = 1$ , using  $h_0(t, s) = e^{-(t-s)^2/(2\sigma_k^2)}$  with  $\sigma_k^2 = \frac{1}{2.4 \log(2k-1)+24}$ , we have that  $h_0$  follows Assumption 2.1 with  $\mu_{h_0} = \frac{3}{4(k-1)}$ .

The following Lemma and definition extend the scalar product induced by  $h$  to generalized dipoles.

**Lemma 2.1.** *Let  $\nu_1 = a_1\delta_{t_1} + b_1\delta'_{v_1, t_1}$ ,  $\nu_2 = a_2\delta_{t_2} + b_2\delta'_{v_2, t_2}$  be two generalized dipoles. Then  $\nu_1$  and  $\nu_2$  are limits (in the distributional sense) of two sequences of dipoles  $\nu_1^{\eta_1}$  and  $\nu_2^{\eta_2}$  for  $\eta_1, \eta_2 \rightarrow 0$ , the quantity  $\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h$  converges, the limit is unique (does not depend on the choice of  $\nu_1^{\eta_1}$  and  $\nu_2^{\eta_2}$ ) and*

$$\begin{aligned} \lim_{\eta_1, \eta_2 \rightarrow 0} \langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h &= a_1 a_2 f(t_1 - t_2) - a_2 b_1 f'_{v_1}(t_1 - t_2) - a_1 b_2 f'_{v_2}(t_2 - t_1) \\ &\quad - b_1 b_2 f''_{v_1, v_2}(t_1 - t_2) \end{aligned} \quad (24)$$

where  $f(t) = \rho(\|t\|_2)$ .

*Proof.* See Appendix A.2. □

**Definition 2.6.** Let  $\nu_1 = a_1\delta_{t_1} + b_1\delta'_{v_1,t_1}$ ,  $\nu_2 = a_2\delta_{t_2} + b_2\delta'_{v_2,t_2}$  be two generalized dipoles. With the previous Lemma, we define

$$\langle \nu_1, \nu_2 \rangle_h := \lim_{\eta_1, \eta_2 \rightarrow 0} \langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h \quad (25)$$

where  $\nu_1^{\eta_1}$  and  $\nu_2^{\eta_2}$  are two sequences of dipoles that converge to  $\nu_1$  and  $\nu_2$  (in the distributional sense) for  $\eta_1, \eta_2 \rightarrow 0$ .

We have the following properties that are immediate consequences of Lemma 2.1.

**Lemma 2.2.** Let  $h$  be a kernel meeting Assumption 2.1. We have the following properties for any  $t \in \mathbb{R}$ :

$$\|\delta_t\|_h^2 = \rho(0) = 1 \quad (26)$$

$$\langle \delta_t, \delta'_{t,v} \rangle_h = -\rho'(0) = 0 \quad (27)$$

$$\|\delta'_{t,v}\|_h^2 = |\rho''(0)| \quad (28)$$

*Proof.* See Appendix A.2. □

From [19, Lemma 6.5], we have the following Lemma:

**Lemma 2.3.** Suppose for all two  $\epsilon$ -separated dipoles,  $\langle \pi_1, \pi_2 \rangle_h \leq \mu \|\pi_1\|_h \|\pi_2\|_h$  (mutual coherence). Then for  $k$ ,  $\epsilon$ -separated dipoles  $\pi_1, \dots, \pi_k$  such that  $\max_i \|\pi_i\|_h > 0$ , we have

$$1 - (k-1)\mu \leq \frac{\|\sum_{i=1,k} \pi_i\|_h^2}{\sum_{i=1,k} \|\pi_i\|_h^2} \leq 1 + (k-1)\mu. \quad (29)$$

We can generalize the previous result to generalized dipoles.

**Lemma 2.4.** Let two  $\epsilon$ -separated **generalized** dipoles  $\nu_1, \nu_2$ . Suppose for all two  $\epsilon$ -separated dipoles  $\pi_1, \pi_2$ ,  $\langle \pi_1, \pi_2 \rangle_h \leq \mu \|\pi_1\|_h \|\pi_2\|_h$  (mutual coherence). Then we have:

$$\langle \nu_1, \nu_2 \rangle_h \leq \mu \|\nu_1\|_h \|\nu_2\|_h \quad (30)$$

*Proof.* See Appendix A.2. □

A consequence of the previous result is the following Lemma:

**Lemma 2.5.** Suppose for all two  $\epsilon$ -separated generalized dipoles,  $\langle \nu_1, \nu_2 \rangle_h \leq \mu \|\nu_1\|_h \|\nu_2\|_h$  (mutual coherence). Then for  $k$   $\epsilon$ -separated generalized dipoles  $\nu_1, \dots, \nu_k$  such that  $\max_i \|\nu_i\|_h > 0$ , we have

$$1 - (k-1)\mu \leq \frac{\|\sum_{i=1,k} \nu_i\|_h^2}{\sum_{i=1,k} \|\nu_i\|_h^2} \leq 1 + (k-1)\mu. \quad (31)$$

*Proof.* See Appendix A.2. □

We are now able to define the Restricted Isometry Property (RIP). The secant set of the model set  $\Sigma$  is  $\Sigma - \Sigma := \{x - y : x \in \Sigma, y \in \Sigma\}$ .

**Definition 2.7** (RIP). *A has the RIP on  $\Sigma - \Sigma$  with respect to  $\|\cdot\|$  with constant  $\gamma$  if for all  $x \in \Sigma - \Sigma$ :*

$$(1 - \gamma)\|x\|^2 \leq \|Ax\|_2^2 \leq (1 + \gamma)\|Ax\|^2. \quad (32)$$

In the following we will suppose that  $A$  has RIP  $\gamma$  on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$  with respect to  $\|\cdot\|_h$ , i.e. for  $\sum_{r=1,k} a_r \delta_{t_r} - \sum_{r=1,k} b_r \delta_{s_r} \in \Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$ , we have

$$\begin{aligned} (1 - \gamma) \left\| \sum_{r=1,k} (a_r \delta_{t_r} - b_r \delta_{s_r}) \right\|_h^2 &\leq \left\| A \sum_{r=1,k} (a_r \delta_{t_r} - b_r \delta_{s_r}) \right\|_2^2 \\ &\leq (1 + \gamma) \left\| \sum_{r=1,k} a_r \delta_{t_r} - b_r \delta_{s_r} \right\|_h^2. \end{aligned} \quad (33)$$

From [19], with a Gaussian kernel  $h$  it is possible to build a random  $A$  with RIP constant  $\gamma$ . With this choice of  $A$ , the ideal minimization (2) yields a stable and robust estimation of  $x_0$  with respect to the  $\|\cdot\|_h$ .

In [8], stable recovery for  $\epsilon$ -separated Diracs is guaranteed on the Torus with the metric  $\|K_{hi} * \cdot\|_{L^1}$  where  $K_{hi} *$  is the convolution with a Fejér kernel. From [5, IV.A], this guarantees a lower RIP with respect to this metric. Indeed, the  $L^1$ -norm of trigonometric polynomials (on  $[0, 1]$ ) is lower bounded by their  $L^2$ -norm, i.e. there is an absolute constant  $D > 0$  depending on  $K_{hi}$  such that  $\|K_{hi} * \cdot\|_{L^1} \geq D \|K_{hi} * \cdot\|_{L^2}$  (see [26, p. 230]). Applying Lemma A.1 from the Annex on the Fejér kernel shows that there exists a kernel metric  $\|\cdot\|_{h_K}$  that lower bounds  $\|K_{hi} * \cdot\|_{L^1}$  for sums of Diracs. This guarantees the existence of a lower RIP with respect to a kernel metric for the conventional deterministic spike super-resolution setting.

The RIP on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$  implies a RIP on  $\epsilon$ -separated generalized dipoles.

**Lemma 2.6** (RIP on generalized dipoles). *Suppose  $A$  has the RIP on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$  with constant  $\gamma$ . Let  $(\nu_r)_{r=1,k}$ ,  $k$   $\epsilon$ -separated dipoles supported in  $\text{rint}\mathcal{B}_2(R)$ , we have*

$$(1 - \gamma) \left\| \sum_{r=1,k} \nu_r \right\|_h^2 \leq \left\| A \left( \sum_{r=1,k} \nu_r \right) \right\|_2^2 \leq (1 + \gamma) \left\| \sum_{r=1,k} \nu_r \right\|_h^2. \quad (34)$$

*Proof.* See Appendix A.2. □

Finally, we will need a last estimate. To state it, we need first to introduce the following definition:

**Definition 2.8.** *Let  $A$  such that the  $\alpha_l$  are in  $\mathcal{C}^2(\mathcal{B}_2(R))$ . We define*

$$D_{A,R} := \sup_{1 \leq l \leq m; v \in \mathbb{R}^d, w \in \mathbb{R}^d: \|v\|_2 = \|w\|_2 = 1; t \in \mathcal{B}_2(R)} |\alpha''_{l,v,w}(t)|. \quad (35)$$

*The constant  $D_{A,R}$  is finite, and it is thus a bound of the directional second derivatives of the  $\alpha_l$  over  $\mathcal{B}_2(R)$ .*

**Lemma 2.7.** *Let  $A$  such that the  $\alpha_l$  are in  $\mathcal{C}^2(\mathcal{B}_2(R))$ . Then, for any  $t \in \mathcal{B}_2(R)$ , with directions  $v_1, v_2$ , we have*

$$\|A\delta''_{t,v_1,v_2}\|_2 \leq \sqrt{m}D_{A,R}. \quad (36)$$

where  $D_{A,R}$  is defined in Equation (35).

*Proof.* See Appendix A.2. □

## 2.4 Control of the conditioning of the Hessian with the restricted isometry property

We can now give a lower (resp. upper) bound for the highest (resp. lowest) eigenvalues of the Hessian matrix  $H$  of  $g$  (computed in Proposition 2.2).

**Theorem 2.1** (Control of the Hessian). *Let  $\theta = (a_1, \dots, a_k, t_1, \dots, t_k) \in \Theta_{k,\epsilon}$  with  $t \in \text{rint}\mathcal{B}_2(R)$  and  $\theta^* \in \Theta_{k,\epsilon}$  a minimizer of (7). Suppose  $h$  follows Assumption 2.1. Let  $H$  the Hessian of  $g$  at  $\theta$ . Suppose  $A$  has RIP  $\gamma$  on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon^*}$ . We have*

$$\sup_{\|u\|_2=1} u^T H u \leq 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (a_r^2 |\rho''(0)|)_{r=1,l}) + \xi; \quad (37)$$

$$\inf_{\|u\|_2=1} u^T H u \geq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l}) - \xi \quad (38)$$

where  $\xi = 2(d+1) \max(\max_r |a_r| \sqrt{m}D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|})(\|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2)$ , the constant  $D_{A,R}$  is defined in (35) and  $e$  is the finite energy measurement noise.

*Proof.* See Appendix A.3. □

**Remark 2.1.** *Notice that, in the noiseless case, (38) ensures in particular that  $g$  has a positive Hessian matrix in  $\theta^*$ . Moreover, if  $\min_r |a_r| > 0$ , there exists a neighbourhood of  $\theta^*$ , in which  $g$  remains convex. We will give an explicit size for this neighbourhood in the next section. Notice also that (37) gives an upper bound on the Lipschitz constant of the gradient of  $g$ . This implies the existence of a basin of attraction (see Definition 1.3) with a uniform bound for the step size.*

**Remark 2.2.** *With the method to choose  $A$  from [19, Lemma 6.5], for any  $\gamma$  and  $m \gtrsim k^2 d \text{polylog}(k, d) / \gamma^2$ , we can find  $A$  that has RIP with high probability with a kernel  $h_0$  having the right properties.*

We can control the conditioning of the Hessian matrix  $\kappa(H)$  at a global minimum as the term  $\|A\phi(\theta) - A\phi(\theta^*)\|_2$  vanishes in the control from Theorem 2.1. Particularly, in the noiseless case we have the following Corollary. The lower bound is useful to confirm the dependency on the ratio of amplitudes when it converges to  $+\infty$ . For this next result, we make the additional assumption that  $\min_r |a_r| > 0$ . In practice, this amounts to assuming that when estimating the Diracs, we do not over-estimate their number (which will often be the case, in particular in the presence of noise). When the number of Diracs is overestimated, the minimizers of (7) are points that are not isolated, the notion of basin of attraction would have to be generalized to a basin of attraction of a set of minimizers (when  $a_r = 0$ ,  $g(\theta)$  does not depend on  $t_r$ ), which is out of the scope of this article for clarity purpose.

**Corollary 2.1.** *Let  $x_0 = \sum_{r=1,k} a_r \delta_{t_r} \in \Sigma_{k,\epsilon} = \phi(\theta_0)$  and  $e = 0$ . Suppose  $h$  follows Assumption 2.1. Let  $H$  the Hessian of  $g$  at  $\theta_0$ . Suppose  $A$  has RIP  $\gamma$  on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$ , and that  $\min_r |a_r| > 0$ . We have*

$$\begin{aligned} \frac{(1 - \gamma) \max(1, (a_r^2 |\rho''(0)|)_{r=1,l})}{(1 + \gamma) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l})} &\leq \kappa(H) \\ &\leq \frac{(1 + \gamma)(1 + (k - 1)\mu) \max(1, (a_r^2 |\rho''(0)|)_{r=1,l})}{(1 - \gamma)(1 - (k - 1)\mu) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l})}. \end{aligned} \quad (39)$$

*Proof.* See Appendix A.3. □

It is easy to see that for a noise  $e$  with small enough energy (i.e. such that  $\xi$  is strictly lower than  $2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l})$ , if  $\min_r |a_r| > 0$ , then the Hessian at a global minimum is strictly positive. Of course, this may require a very small noise since the ratio of amplitudes at the global minimum can be large.

**Remark 2.3.** *We remark that for a same maximal ratio of amplitudes in  $\theta^*$ , a better conditioning bound is achieved when  $\max_{r=1,l} a_r^2 |\rho''(0)| \geq 1 \geq \min_{r=1,l} a_r^2 |\rho''(0)|$ . We attribute this to the fact that we estimate amplitudes and locations at the same time. The amplitudes must be appropriately scaled to match the variations of  $g$  with respect to locations. Intuitively, alternate descent with respect to amplitudes and locations might be better than the classical gradient descent for easily setting the descent step.*

**Remark 2.4.** *As  $g$  is  $\mathcal{C}^2$ , ensuring the strict positivity of the Hessian at the global minimum guarantees the existence of a basin of attraction as emphasized in Section 1.3. In the next Section, we give an explicit formulation of a basin of attraction.*

### 3 Explicit basin of attraction of the global minimum

Let  $\theta_1 \in \mathbb{R}^d$ . Can we guarantee, for some notion of distance  $d$ , that  $d(\theta_1, \theta^*) \leq C$  and  $\theta_1 \neq \theta^*$ , with  $C$  an explicit constant, implies  $\nabla g(\theta_1) \neq 0$ ? The following theorems show that it is in fact the case. With a strong RIP assumption, we can give an explicit basin of attraction of the global minimum for minimization (7) without separation constraints.

#### 3.1 Uniform control of the Hessian

In the noiseless case, a global minimum  $\theta^*$  of the constrained minimization of  $g$  over  $\Theta_{k,\epsilon}$  is also a global minimum of the unconstrained minimization because  $g(\theta^*) = 0$ . In the presence of noise, we can no longer guarantee that the minimizer of the constrained problem  $\theta^*$  is a global minimum of the unconstrained problem. However, the shape of the constraint guarantees that it is a local minimum (see next Lemma).

**Lemma 3.1.** *Suppose  $\theta^* = (a_1, \dots, a_k, t_1, \dots, t_k)$  is a result of constrained minimization (7) with  $t_i \in \text{rint}\mathcal{B}_2(R)$ . Then  $\theta^*$  is a local minimum of  $g$ .*

*Proof.* let  $\theta^* = (a_1, \dots, a_k, t_1, \dots, t_k)$ . As for all  $i \neq j$ ,  $\|t_i - t_j\|_\infty > \epsilon$ , there exists  $\eta > 0$  such that for all  $\theta = (b_1, \dots, b_k, s_1, \dots, s_k)$  such that  $\|s_i - t_i\|_\infty < \eta$ , we have  $\theta \in \Theta_{k,\epsilon}$ . Hence,  $\theta^* + B_\infty(\eta) \subset \Theta_{k,\epsilon}$ , and  $\theta^* \in \arg \min_{\theta \in \theta^* + B_\infty(\eta)} g(\theta)$ .  $\square$

Hence we can still calculate a basin of attraction of  $\theta^*$  (for the unconstrained minimization). The expression of the basin in the next Section is a direct consequence of the following Theorem that uniformly control the Hessian of  $g$  in an explicit neighbourhood of  $\theta^*$ .

**Theorem 3.1.** *Suppose  $A$  has RIP  $\gamma$  on  $\Sigma_{k,\frac{\epsilon}{2}} - \Sigma_{k,\frac{\epsilon}{2}}$  and that  $h$  follows Assumption 2.1 and has mutual coherence constant  $\mu$  on  $\frac{\epsilon}{2}$ -separated dipoles. Let  $\theta^* = (a_1, \dots, a_k, t_1, \dots, t_k) \in \Theta_{k,\epsilon}$  be a result of constrained minimization (7) such that  $t_i \in \text{rint}\mathcal{B}_2(R)$ . Suppose  $0 < |a_1| \leq |a_2| \dots \leq |a_k|$ . Let  $0 \leq \beta \leq \frac{\epsilon}{4}$  and*

$$\Lambda_{\theta^*,\beta} := \{\theta : \|\theta - \theta^*\|_2 < \beta\}. \quad (40)$$

If  $\theta \in \Lambda_{\theta^*,\beta}$ , then  $H$  the Hessian of  $g$  at  $\theta$  has the following bounds :

$$\sup_{\|u\|_2=1} u^T H u \leq 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (|a_k| + \beta)^2 |\rho''(0)|) + \xi; \quad (41)$$

$$\inf_{\|u\|_2=1} u^T H u \geq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (|a_1| - \beta)^2 |\rho''(0)|) - \xi \quad (42)$$

where  $\xi = 2(d + 1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\sup_{\theta \in \Lambda_{\theta^*,\beta}} \|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2)$ , the constant  $D_{A,R}$  is given in (35) and  $e$  is the finite energy measurement noise.

*Proof.* See Appendix A.4.  $\square$

**Remark 3.1.** *We observe that we require a stronger RIP than the usual one on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$  to guarantee that unconstrained minimization converges in the basin of attraction  $\Lambda_{\theta^*,\beta}$ .*

The set  $\Lambda_{\theta^*,\beta}$  is an open  $\ell^2$  ball centered on  $\theta^*$ . The choice of this set, besides its simplicity, is useful to guarantee the convergence of the gradient descent. We could guarantee the positivity of the Hessian on bigger sets with more complicated formulations. Guaranteeing that iterates of the gradient descent stay in such sets would become much harder then. When the separation constraint is added for the basin of attraction (we look for potential critical points in  $\Sigma_{k,\epsilon}$ ), we can provide better bounds. We will discuss what we could expect from constrained descent algorithms in Section 4.

**Theorem 3.2.** *Suppose  $A$  has RIP  $\gamma$  on  $\Sigma_{k,\epsilon} - \Sigma_{k,\epsilon}$  and that  $h$  follows Assumption 2.1 and has mutual coherence constant  $\mu$  on  $\epsilon$ -separated dipoles. Suppose  $0 < |a_1| \leq |a_2| \dots \leq |a_k|$ . Let  $\theta^* = (a_1, \dots, a_k, t_1, \dots, t_k) \in \Theta_{k,\epsilon}$  be a result of constrained minimization (7) such that  $t_i \in \text{rint}\mathcal{B}_2(R)$ . Let  $\beta \geq 0$  and*

$$\Lambda_{\theta^*,\beta} := \{\theta : \|\theta - \theta^*\|_2 < \beta\}. \quad (43)$$

Then for  $\theta \in \Theta_{k,\epsilon} \cap \Lambda_{\theta^*,\beta}$ , then  $H$  the Hessian of  $g$  at  $\theta$  has the following bounds:

$$\sup_{\|u\|_2=1} u^T H u \leq 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (|a_k| + \beta)^2 |\rho''(0)|) + \xi; \quad (44)$$

$$\inf_{\|u\|_2=1} u^T H u \geq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (|a_1| - \beta)^2 |\rho''(0)|) - \xi \quad (45)$$

where  $\xi = 2(d + 1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\sup_{\theta \in \Lambda_{\theta^*,\beta}} \|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2)$ , the constant  $D_{A,R}$  is given in (35) and  $e$  is the finite energy measurement noise.

*Proof.* See Appendix A.4. □

### 3.2 Explicit basin of attraction in the noiseless and noisy case

With the help of this uniform control of the Hessian we give an explicit (yet suboptimal) basin of attraction.

**Corollary 3.1** (of Theorem 3.1, noiseless case). *Under the hypotheses of Theorem 3.1, let  $\theta^* \in \Theta_{k,\epsilon}$  be a result of constrained minimization (7). Let  $a^* = (a_1, a_2, \dots, a_k)$ . Take*

$$\beta_{max} := \min \left( c_h, \frac{|a_1|}{2}, C_1 C_2 \right)$$

where  $C_1 = \frac{(1-\gamma)(1-(k-1)\mu)}{(d+1)\sqrt{1+\gamma}\sqrt{1+(k-1)\mu}}$  and  $C_2 = \frac{\min(1, |a_1|^2 |\rho''(0)|/4)}{\max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \sqrt{1+2|\rho''(0)|} \|a^*\|_2^2}$ . Then the set  $\Lambda_{\theta^*,\beta_{max}}$  is a basin of attraction of  $\theta^*$ .

*Proof.* See Appendix A.4. □

The parameter  $\beta$  controls the distance between a parameter and the optimal parameter. When the RIP constant  $\gamma$  decreases (and generally as the number of measurement increases), the size of the basin of attraction increases. In both the context of regular Fourier sampling and random Fourier sampling, the constant  $D_{A,R}$  is bounded when  $m$  increases. When the mutual coherence constant  $\mu$  decreases, the basin of attraction also increases. The size of the basin also decreases as the ratio of amplitudes  $\frac{a_1}{a_k}$  decreases. We observe again that performing the descent with respect to amplitudes and positions at the same time yields pessimistic bounds for the basin of attraction. Finally, we note that the smaller  $\beta$  is, the smaller is the upper bound on the operator norm of the Hessian.

When the noise contaminating the measurements is small enough, we have similar results with a smaller basin of attraction.

**Corollary 3.2** (of Theorem 3.1, noisy case). *Under the hypotheses of Theorem 3.1, let  $\theta^* \in \Theta_{k,\epsilon}$  be a result of constrained minimization (7). Let  $a^* = (a_1, a_2, \dots, a_k)$ . Take*

$$\beta_{max} := \min \left( c_h, \frac{|a_1|}{2}, C_1 C_3 \right)$$



where  $C_1 = \frac{(1-\gamma)(1-(k-1)\mu)}{(d+1)\sqrt{1+\gamma}\sqrt{1+(k-1)\mu}}$  and  $C_3 = \frac{\min(1, |a_1|^2 |\rho''(0)|/4)}{\max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|})(1 + \sqrt{1+2|\rho''(0)|} \|a^*\|_2^2)}$

Suppose  $\|e\|_2 \leq \sqrt{1+\gamma} \sqrt{1+(k-1)\mu} \beta$ . Then the set  $\Lambda_{\theta^*, \beta_{max}}$  is a basin of attraction of  $\theta^*$ .

*Proof.* See Appendix A.4. □

## 4 Towards new descent algorithms for sparse spike estimation?

We have shown that, given an appropriate measurement operator for separated Diracs, a good initialization is sufficient to guarantee the success of a simple gradient descent. Such gradient descent is used in the practical setting of compressive statistical learning [21]. Our result on unconstrained minimization explains why the use of such gradient descent is valid in this setting. If we could guarantee additionally that by greedily estimating Diracs, we fall within the basin of attraction, we would have a full non-convex optimization technique with guarantees of convergence to a global minimum.

In other works [15, 16], it has been shown that discretization (on grids) of convex methods have a tendency to produce spurious spikes at Dirac locations. Our results seem to indicate that merging spikes that are close to each other when performing a gradient descent might break the barrier between continuous and discrete methods.

Theorem 3.2 brings another question as the Hessian of  $g$  is more easily controlled in  $\Theta_{k,\epsilon}$ . More generally, can we build a simple descent algorithm that stays in  $\Theta_{k,\epsilon}$  to get larger basins of attraction? Consider the problem for  $d = 1$  in the noiseless case for the sake of clarity. We want to use the following descent algorithm:

$$\theta_{i+1} = P_{\Theta_{k,\epsilon}}(\theta_i - \tau \nabla g(\theta_i)) \quad (46)$$

Where  $P_{\Theta_{k,\epsilon}}$  is a projection onto the separation constraint. Notice that since  $\Theta_{k,\epsilon}$  is not a convex set, we cannot easily define the orthogonal projection onto it (it may not even exist).

If we suppose that the gradient descent step decreases  $g$  (i.e.  $g(\theta_i - \tau \nabla g(\theta_i)) < g(\theta_i)$ ), is it possible to guarantee that applying projection step keeps decreasing  $g$ ? Consider:

$$P_{\Theta_{k,\epsilon}}(\theta) \in \arg \min_{\tilde{\theta} \in \Theta_{k,\epsilon}} \left| \|A\phi(\tilde{\theta}) - y\|_2 - \|A\phi(\theta) - y\|_2 \right| \quad (47)$$

First consider the following Lemma:

**Lemma 4.1.** *Let  $d = 1$ . Let  $\theta_0, \theta_1 \in \Theta_{k,\epsilon}$ . Let  $g(\theta) = \|A\phi(\theta) - A\phi(\theta_0)\|$ . Then for all  $\alpha$  such that  $0 = g(\theta_0) \leq \alpha \leq g(\theta_1)$ , there exists  $\theta^* \in \Theta_{k,\epsilon}$  such that  $g(\theta^*) = \alpha$ .*

*Proof.* See Appendix A.5. □

Lemma 4.1 essentially guarantees that it is possible to continuously map the interval  $[0, g(\theta_1)]$  by  $g$  with elements of  $\Theta_{k,\epsilon}$ . Hence, at a step  $i + 1$ , we have

$$|g(\theta_{i+1}) - g(\theta_i)| = |g(\theta_i - \tau \nabla g(\theta_i)) - g(\theta_i)|. \quad (48)$$

The projection  $P_{\Theta_{k,\epsilon}}$  defined by (47) is not easy to calculate (in fact, it is a similar optimization problem as the main problem). Other more "natural" projections on  $\Theta_{k,\epsilon}$  could be defined as :

$$P_{\Theta_{k,\epsilon}}(\theta) \in \phi^{-1}(\arg \inf_{x \in \Sigma_{k,\epsilon}} \|Ax - A\phi(\theta)\|_2) \quad (49)$$

or

$$P_{\Theta_{k,\epsilon}}(\theta) \in \phi^{-1}(\arg \inf_{x \in \Sigma_{k,\epsilon}} \|x - \phi(\theta)\|_h). \quad (50)$$

However they suffer from the same calculability drawback. This suggests to build a new family of heuristic algorithms of spike estimation where we propose heuristics to approach the projection of  $\hat{\theta}_{i+1}$  on  $\Theta_{k,\epsilon}$ . Recovery guarantees would be obtained by guaranteeing that the projection heuristic does not increase the value of  $g$  by too much compared to the gradient descent step.

## A Annex

### A.1 Proofs for Section 2.2

*Proof of Proposition 2.1.*

$$\begin{aligned} \frac{\partial g(\theta)}{\partial a_r} &= \frac{\partial}{\partial a_r} \sum_{l=1}^m \left| \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right|^2 \\ &= \sum_{l=1}^m 2\mathcal{R}e \left( \alpha_l(t_{r,j}) \left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right) \right) \\ &= 2\mathcal{R}e \langle A\delta_{t_r}, A\phi(\theta) - y \rangle. \end{aligned} \quad (51)$$

Similarly,

$$\begin{aligned} \frac{\partial g(\theta)}{\partial t_{r,j}} &= \frac{\partial}{\partial t_r} \sum_{l=1}^m \left| \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right|^2 \\ &= \sum_{l=1}^m 2\mathcal{R}e \left( a_r \partial_j \alpha_l(t_r) \left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right) \right) \\ &= -2a_r \mathcal{R}e \langle A\delta'_{t_r,j}, A\phi(\theta) - y \rangle. \end{aligned} \quad (52)$$

□

*Proof of Proposition 2.2.* For  $H_{1,r,s}$ ,

$$\begin{aligned} \frac{\partial^2 g(\theta)}{\partial a_r \partial a_s} &= \frac{\partial}{\partial a_s} \sum_{l=1}^m 2\mathcal{R}e \left( \alpha_l(t_r) \overline{\left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right)} \right) \\ &= \sum_{l=1}^m 2\mathcal{R}e \left( \alpha_l(t_r) \overline{\alpha_l(t_s)} \right). \end{aligned} \quad (53)$$

For  $H_{2,r,j_1,s,j_2}$ ,

$$\begin{aligned} \frac{\partial^2 g(\theta)}{\partial t_{r,j_1} \partial t_{s,j_2}} &= \frac{\partial}{\partial t_{s,j_1}} \sum_{l=1}^m 2\mathcal{R}e \left( a_r \partial_{j_1} \alpha_l(t_r) \overline{\left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right)} \right) \\ &= \sum_{l=1}^m 2\mathcal{R}e \left( a_r \partial_{j_1} \alpha_l(t_r) \overline{\left( a_s \partial_{j_2} \alpha_l(t_s) \right)} \right) \\ &\quad + \mathbf{1}(r=s) \sum_{l=1}^m 2\mathcal{R}e \left( a_r \partial_{j_2} \partial_{j_1} \alpha_l(t_r) \overline{\left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right)} \right). \end{aligned} \quad (54)$$

For  $H_{12,r,s,j}$

$$\begin{aligned} \frac{\partial^2 g(\theta)}{\partial a_r \partial t_{s,j}} &= \frac{\partial}{\partial t_{s,j}} \sum_{l=1}^m 2\mathcal{R}e(\alpha_l(t_r)) \overline{\left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right)} \\ &= \sum_{l=1}^m 2\mathcal{R}e \left( \alpha_l(t_r) \overline{\left( a_s \partial_j \alpha_l(t_s) \right)} \right) \\ &\quad + \mathbf{1}(r=s) \sum_{l=1}^m 2\mathcal{R}e \left( \partial_j \alpha_l(t_r) \overline{\left( \sum_{i=1}^k a_i \alpha_l(t_i) - y_l \right)} \right). \end{aligned} \quad (55)$$

□

## A.2 Proofs for Section 2.3

*Proof of Lemma 2.1.* First remark that a generalized dipole  $\nu = a\delta_t + b\delta'_{t,v}$  with  $\|v\|_2 = 1$  is the limit in the distributional sense of the dipoles  $\nu^\eta = a\delta_t - b\frac{\delta_t + \eta v - \delta_t}{\eta}$  when  $\eta \rightarrow 0$ .

Now let two generalized dipoles  $\nu_1 = a_1\delta_{t_1} + b_1\delta'_{t_1,v_1}$ ,  $\nu_2 = a_2\delta_{t_2} + b_2\delta'_{t_2,v_2}$ . The  $\nu_i$  are the limit (in the distributional sense) of a family of dipole  $\nu_i^{\eta_i}$  for  $\eta_i \rightarrow 0^+$ . Let  $f(t) = \rho(\|t\|_2)$ . We have

$$\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} f(t-s) d\nu_1^{\eta_1}(t) d\nu_2^{\eta_2}(s). \quad (56)$$

Remark that by construction  $g_{\eta_1}(s) := \int_{\mathbb{R}^d} f(t-s) d\nu_1^{\eta_1}(t) \rightarrow_{\eta_1 \rightarrow 0^+} g(s) := a_1 f(t_1 - s) + b_1 \langle \delta'_{t_1,v_1}, f(\cdot - s) \rangle < +\infty$  where  $g_{\eta_1}$  is in  $\mathcal{C}^2$  and  $g$  is in  $\mathcal{C}^1$  thanks to the assumption

on  $h$  and  $\rho$ . Hence by boundedness of the integrals and the dominated convergence theorem, for any  $\eta_2$ ,

$$\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h \rightarrow_{\eta_1 \rightarrow 0^+} \int_{\mathbb{R}^d} g(s) d\nu_2^{\eta_2}(s). \quad (57)$$

Moreover, by construction of  $\nu_2^{\eta_2}$ , and symmetry of  $f$  (i.e.  $f(t_1 - t_2) = f(t_2 - t_1)$ ),

$$\begin{aligned} \int_{\mathbb{R}^d} g(s) d\nu_2^{\eta_2}(s) &\rightarrow_{\eta_2 \rightarrow 0^+} a_1 a_2 f(t_1 - t_2) \\ &+ a_2 b_1 \langle \delta'_{t_1, v_1}, f(\cdot - t_2) \rangle + a_1 b_2 \langle \delta'_{t_2, v_2}, f(t_1 - \cdot) \rangle \\ &- b_1 b_2 \langle \delta'_{t_2, v_2}, f'_{v_1}(t_1 - \cdot) \rangle \\ &= a_1 a_2 f(t_1 - t_2) - a_2 b_1 f'_{v_1}(t_1 - t_2) - a_1 b_2 f'_{v_2}(t_2 - t_1) \\ &- b_1 b_2 f''_{v_1, v_2}(t_1 - t_2) \end{aligned} \quad (58)$$

We define  $\langle \nu_1, \nu_2 \rangle_h := a_1 a_2 f(t_1 - t_2) - a_2 b_1 f'_{v_1}(t_1 - t_2) - a_1 b_2 f'_{v_2}(t_2 - t_1) - b_1 b_2 f''_{v_1, v_2}(t_1 - t_2)$ . We just showed that

$$\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h \rightarrow_{\eta_1 \rightarrow 0^+, \eta_2 \rightarrow 0^+} \langle \nu_1, \nu_2 \rangle_h. \quad (59)$$

Note that the value of  $\langle \nu_1, \nu_2 \rangle_h$  only depends on  $\rho, \nu_1, \nu_2$ .  $\square$

*Proof of Lemma 2.2.* Using Lemma 2.1 with  $t_1 = t_2 = t$ ,  $b_1 = b_2 = 0$  and  $a_1 = a_2 = 1$  gives

$$\|\delta_t\|_h^2 = \rho(0). \quad (60)$$

Using Lemma 2.1 with  $t_1 = t_2 = t$ ,  $b_1 = a_2 = 0$  and  $a_1 = b_2 = 1$  gives

$$\langle \delta_t, \delta'_{t, v} \rangle_h := -f'_v(0) = -\lim_{\eta \rightarrow 0^+} \frac{\rho(\eta \|v\|) - \rho(0)}{\eta} = -\rho'(0) = 0. \quad (61)$$

Using Lemma 2.1 with  $t_1 = t_2 = t$ ,  $b_1 = b_2 = 1$  and  $a_1 = a_2 = 0$  gives

$$\|\delta'_{t, v}\|_h^2 := -f''_v(0) = |\rho''(0)|. \quad (62)$$

$\square$

*Proof of Lemma 2.4.* Using the construction from the proof of Lemma 2.1, let two  $\epsilon$ -separated generalized dipole  $\nu_1, \nu_2$ . The  $\nu_i$  are the limit (in the distributional sense) of a family of  $\epsilon$ -separated dipole  $\nu_i^{\eta_i}$  for  $\eta_i \rightarrow 0^+$ . With the hypothesis, we have

$$\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h \leq \mu \|\nu_1^{\eta_1}\|_h \|\nu_2^{\eta_2}\|_h. \quad (63)$$

Furthermore,

$$\langle \nu_1^{\eta_1}, \nu_2^{\eta_2} \rangle_h \rightarrow_{\eta_1 \rightarrow 0^+, \eta_2 \rightarrow 0^+} \langle \nu_1, \nu_2 \rangle_h. \quad (64)$$

Let  $\nu = a\delta_t + b\delta'_{t,v}$  with  $\|\nu\|_2 = 1$  and  $\nu^\eta = a\delta_t - b\frac{\delta_t + \eta v - \delta_t}{\eta} = \left(a + \frac{b}{\eta}\right)\delta_t - b\frac{\delta_t + \eta v}{\eta}$ . We have  $\|\nu\|_h^2 = a^2 + b^2|\rho''(0)|$  (with Lemma 2.2) and

$$\begin{aligned}\|\nu^\eta\|_h^2 &= \left(a + \frac{b}{\eta}\right)^2 + \left(\frac{b}{\eta}\right)^2 - 2\left(a + \frac{b}{\eta}\right)\frac{b}{\eta}\rho(\eta) \\ &= a^2 + 2\left(\frac{b}{\eta}\right)^2 + 2\frac{ab}{\eta} - 2\frac{ab}{\eta}\rho(\eta) - 2\left(\frac{b}{\eta}\right)^2\rho(\eta) \\ &= a^2 + 2\frac{ab}{\eta}(1 - \rho(\eta)) + 2\frac{b^2}{\eta^2}(1 - \rho(\eta)).\end{aligned}\tag{65}$$

But  $\frac{1-\rho(\eta)}{\eta} = \frac{\rho(0)-\rho(\eta)}{\eta} \rightarrow -\rho'(0)$  when  $\eta \rightarrow 0^+$ , and  $\rho'(0) = 0$ .

Moreover,  $\rho(\eta) = h(0) + \eta\rho'(0) + \frac{\eta^2}{2}\rho''(0) + o(\eta^2) = 1 - \frac{\eta^2}{2}|\rho''(0)| + o(\eta^2)$ . Hence  $\frac{1-\rho(\eta)}{\eta^2} \rightarrow_{\eta \rightarrow 0^+} \frac{1}{2}|\rho''(0)|$ .

We thus deduce that  $\|\nu^\eta\|_h^2 \rightarrow a^2 + b^2|\rho''(0)| = \|\nu\|_h$  when  $\eta \rightarrow 0^+$ .

Hence, with such choice of  $\nu_1^{\eta_1}, \nu_2^{\eta_2}$ , we can take the limit  $\eta_1, \eta_2 \rightarrow 0$  in Equation (63) to get the result.  $\square$

*Proof of Lemma 2.5.* Using Lemma 2.4, and the same proof as in Lemma 2.3, we get the result.  $\square$

*Proof of Lemma 2.6.* Let  $\nu_r = a_r\delta_{t_r} + b_r\delta'_{t_r,v}$  the  $\epsilon$ -separated generalized dipoles. Similarly to Lemma 2.4, take  $\nu_r^\eta = (a_r + \frac{b_r}{\eta})\delta_{t_r} - b_r\frac{\delta_{t_r} + \eta v}{\eta}$ . For sufficiently small  $\eta$  the  $\nu_r^\eta$  are  $\epsilon$ -separated dipoles, hence  $\sum \nu_r^\eta \in \Sigma - \Sigma$  and

$$(1 - \gamma) \left\| \sum_{r=1,k} \nu_r^\eta \right\|_h^2 \leq \left\| A \left( \sum_{r=1,k} \nu_r^\eta \right) \right\|_2^2 \leq (1 + \gamma) \left\| \sum_{r=1,k} \nu_r^\eta \right\|_h^2.\tag{66}$$

Now remark that  $g_1(\eta) = \|\sum_{r=1,k} \nu_r^\eta\|_h^2$  and  $g_2(\eta) = \|A(\sum_{r=1,k} \nu_r^\eta)\|_2^2$  are continuous functions of  $\eta$  that converge to  $\|\sum_{r=1,k} (a_r\delta_{t_r} + b_r\delta'_{t_r,v})\|_h^2$  and  $\|A(\sum_{r=1,k} (a_r\delta_{t_r} + b_r\delta'_{t_r,v}))\|_2^2$  when  $\eta \rightarrow 0$ :

- For  $g_1$ , use the same proof as in Lemma 2.4 with the linearity of the limit.

- For  $g_2$ :

$$\begin{aligned}
g_2(\eta) &= \sum_{l=1,m} \left| \sum_{r=1,k} \int \alpha_l(t) (a_r d\delta_{t_r}(t) - \frac{b_r}{\eta} (d\delta_{t_r+\eta v}(t) - d\delta_{t_r}(t))) \right|^2 \\
&= \sum_{l=1,m} \left| \sum_{r=1,k} \left( \alpha_l(t_r) a_r - \frac{b_r}{\eta} (\alpha_l(t_r + \eta v) - \alpha_l(t_r)) \right) \right|^2 \\
&\xrightarrow{\eta \rightarrow 0^+} \sum_{l=1,m} \left| \sum_{r=1,k} (\alpha_l(t_r) a_r - b_r (\alpha_l)'_v(t_r)) \right|^2 \\
&= \left\| A \left( \sum_{r=1,k} a_r \delta_{t_r} + b_r \delta'_{t_r, v} \right) \right\|_2^2.
\end{aligned} \tag{67}$$

Taking the limit of Equation (66) for  $\eta \rightarrow 0$  yields the result.  $\square$

*Proof of Lemma 2.7.* We have

$$\begin{aligned}
\|A\delta''_{t, v_1, v_2}\|_2^2 &= \sum_{l=1,m} |(A\delta''_{t, v_1, v_2})_l|^2 \\
&= \sum_{l=1,m} |\alpha''_{l, v_1, v_2}(t)|^2 \\
&\leq m \sup_{l=1,m; t \in \mathcal{B}_2(R)} |\alpha''_{l, v_1, v_2}(t)|^2 \leq m D_{A,R}^2
\end{aligned} \tag{68}$$

where  $D_{A,R}$  is given in (35), i.e.  $D_{A,R}$  is the supremum of directional second derivatives of the  $\alpha_l$  over  $\mathcal{B}_2(R)$ . We have  $D_{A,R} < +\infty$  because the  $\alpha_l$  are supposed to be in  $\mathcal{C}^2(\mathcal{B}_2(R))$ .  $\square$

**Lemma A.1.** *Let  $K$  be a symmetrical convolution kernel in  $\mathcal{C}^2$  and  $h_K : (t_1, t_2) \rightarrow h_K(t_1, t_2) = [K * K](t_1 - t_2)$  (the convolution of  $K$  by itself) then for any  $x \in \Sigma_{k, \epsilon} - \Sigma_{k, \epsilon}$ , we have*

$$\|x\|_{h_K}^2 = \|K * x\|_{L^2}^2. \tag{69}$$

*Proof of Lemma A.1.* Write  $x = \sum a_i \delta_{t_i}$  and use the symmetry of  $K$ :

$$\begin{aligned}
\|K * x\|_{L^2}^2 &= \int \left| \sum a_i K(t - t_i) \right|^2 dt = \sum_{i,j} a_i a_j \int K(t - t_i) K(t - t_j) dt \\
&= \sum_{i,j} a_i a_j \int K(t) K(t + t_i - t_j) dt \\
&= \sum_{i,j} a_i a_j [K * K](t_i - t_j) = \|x\|_{h_K}^2.
\end{aligned} \tag{70}$$

□

### A.3 Proofs for Section 2.4

We will use the following Lemma on directional derivatives of Diracs.

**Lemma A.2.** *Let  $u, t_0 \in \mathbb{R}^d$ . Suppose  $u \neq 0$ . Then,  $\sum_{i=1,d} u_i \delta'_{t_0, i} = \|u\|_2 \delta'_{t_0, \frac{u}{\|u\|_2}}$ .*

*Proof.* Let  $f$  a function in  $\mathcal{C}^2(\mathbb{R}^d)$ , we have

$$\int_{t \in \mathbb{R}^d} f(t) \sum_{i=1,d} u_i d\delta'_{t_0, i}(t) = - \sum_{i=1,d} u_i \partial_i f(t_0) = - \langle u_i, \nabla f(t_0) \rangle = - \|u\|_2 f'_{\frac{u}{\|u\|_2}}(t_0).$$

Hence,  $\sum_{i=1,d} u_i \delta'_{t_0, i} = \|u\|_2 \delta'_{t_0, \frac{u}{\|u\|_2}}$  □

To prove Theorem 2.1, we control first the eigenvalues of  $G$  in the decomposition  $H = G + F$ .

**Lemma A.3.** *Suppose  $h$  follows Assumption 2.1. Let  $\theta = (a_1, \dots, a_k, t_1, \dots, t_k) \in \Theta_{k, \epsilon}$  with  $t \in \text{rint}\mathcal{B}_2(R)$ . Let  $H$  the Hessian of  $g$  at  $\theta$ . Suppose  $A$  has RIP  $\gamma$  on  $\Sigma_{k, \epsilon} - \Sigma_{k, \epsilon}$ . We have*

$$\sup_{\|u\|_2=1} u^T G u \leq 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (a_r^2 |\rho''(0)|)_{r=1,l}); \tag{71}$$

$$\inf_{\|u\|_2=1} u^T G u \geq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l}). \tag{72}$$

where  $G$  is defined in Proposition 2.2.

*Proof.* Let  $u \in \mathbb{R}^{k(d+1)}$  such that  $\|u\|_2 = 1$ . We index  $u$  as follows:  $u_r \in \mathbb{R}$  for  $r = 1, k$ .

$u_r \in \mathbb{R}^d$  for  $r = k+1, 2k$  (it follows the indexing of  $H$  and  $G$  we used). Remark that

$$\begin{aligned}
u^T G u &= \sum_{r,s=1,k} u_r u_s G_{1,r,s} + \sum_{r=k+1,2k;j_1=1,d;s=k+1,2k;j_2=1,d} u_{r,j_1} u_{s,j_2} G_{2,r,j_1,s,j_2} \\
&+ \sum_{r=1,k;s=k+1,2k;j=1,d} u_r u_s G_{12,r,s,j} + \sum_{r=k+1,2k;j=1,d;s=1,k} u_{r,j} u_s G_{21,r,j,s} \\
&= 2 \sum_{r,s=1,k} \operatorname{Re} \langle A u_r \delta_{t_r}, A u_s \delta_{t_s} \rangle \\
&+ 2 \sum_{r=k+1,2k;j_1=1,d;s=k+1,2k;j_2=1,d} \operatorname{Re} \langle A u_{r,j_1} a_{r-k} \delta'_{t_{r-k},j_1}, A u_{s,j_2} a_{s-k} \delta'_{t_{s-k},j_2} \rangle \\
&- 2 \sum_{r=1,k;s=k+1,2k;j=1,d} \operatorname{Re} \langle A u_r \delta_{t_r}, A u_{s,j} a_{s-k} \delta'_{t_{s-k},j} \rangle \\
&- 2 \sum_{r=k+1,2k;j=1,d;s=1,k} \operatorname{Re} \langle A u_{r,j} a_{r-k} \delta'_{t_{r-k},j}, A u_s \delta_{t_s} \rangle
\end{aligned} \tag{73}$$

Thus we have

$$\begin{aligned}
u^T G u &= 2 \left\| A \sum_{r=1,k} u_r \delta_{t_r} \right\|_2^2 + 2 \left\| A \sum_{r=k+1,2k;j=1,d} u_{r,j} a_{r-k} \delta'_{t_{r-k},j} \right\|_2^2 \\
&- 2 \operatorname{Re} \left\langle A \sum_{r=1,k} u_r \delta_{t_r}, A \sum_{r=k+1,2k;j=1,d} u_{r,j} a_{r-k} \delta'_{t_{r-k},j} \right\rangle \\
&- 2 \operatorname{Re} \left\langle A \sum_{r=k+1,2k;j=1,d} u_{r,j} a_{r-k} \delta'_{t_{r-k},j}, A \sum_{r=1,k} u_r \delta_{t_r} \right\rangle \\
&= 2 \left\| A \left( \sum_{r=1,k} u_r \delta_{t_r} - \sum_{r=k+1,2k;j=1,d} u_{r,j} a_{r-k} \delta'_{t_{r-k},j} \right) \right\|_2^2 \\
&= 2 \left\| A \left( \sum_{r=1,k} \left( u_r \delta_{t_r} - a_r \sum_{j=1,d} u_{r+k,j} \delta'_{t_r,j} \right) \right) \right\|_2^2.
\end{aligned} \tag{74}$$

Using Lemma A.2, we have  $\sum_{j=1,d} w_j \delta'_{t_r,j} = \|w\|_2 \delta'_{t_r, \frac{w}{\|w\|_2}}$  and

$$u^T G u = 2 \left\| A \sum_{r=1,k} \left( u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}} \right) \right\|_2^2. \tag{75}$$

We use the lower RIP in Lemma 2.6,

$$u^T G u \geq 2(1 - \gamma) \left\| \sum_{r=1,k} \left( u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}} \right) \right\|_h^2. \tag{76}$$



Then the hypothesis on  $\|\cdot\|_h$  and Lemma 2.5 yields

$$\begin{aligned} & \left\| \sum_{r=1,k} (u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}) \right\|_h^2 \\ & \geq (1 - (k-1)\mu) \sum_{r=1,k} \|u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}\|_h^2 \end{aligned} \quad (77)$$

and

$$\begin{aligned} u^T G u & \geq 2(1-\gamma)(1-(k-1)\mu) \sum_{r=1,k} \|u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}\|_h^2 \\ & \geq 2(1-\gamma)(1-(k-1)\mu) \sum_{r=1,k} \left( |u_r|^2 - 2a_r u_r \|u_{k+r}\|_2 \langle \delta_{t_r}, \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}} \rangle_h \right. \\ & \quad \left. + a_r^2 \|u_{k+r}\|_2^2 \|\delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}\|_h^2 \right). \end{aligned} \quad (78)$$

Then using Lemma 2.2:

$$\begin{aligned} u^T G u & \geq 2(1-\gamma)(1-(k-1)\mu) \sum_{r=1,k} (|u_r|^2 + a_r^2 \|u_{k+r}\|_2^2 |\rho''(0)|) \\ & \geq 2(1-\gamma)(1-(k-1)\mu) \inf_{\|u\|_2=1} \sum_{r=1,k} (|u_r|^2 + \|u_{k+r}\|_2^2 a_r^2 |\rho''(0)|). \\ & = 2(1-\gamma)(1-(k-1)\mu) \min(1, (a_r^2 |\rho''(0)|)_{r=1,l}). \end{aligned} \quad (79)$$

Similarly, using the upper RIP in Lemma 2.6:

$$u^T G u \leq 2(1+\gamma) \left\| \sum_{r=1,k} (u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}) \right\|_h^2. \quad (80)$$

Then the hypothesis on  $\|\cdot\|_h$  yields (Lemma 2.5)

$$\begin{aligned} & \left\| \sum_{r=1,k} (u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}) \right\|_h^2 \\ & \leq (1 + (k-1)\mu) \sum_{r=1,k} \|u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}\|_h^2 \end{aligned} \quad (81)$$

and

$$u^T G u \leq 2(1+\gamma)(1+(k-1)\mu) \sum_{r=1,k} \|u_r \delta_{t_r} - a_r \|u_{r+k}\|_2 \delta'_{t_r, \frac{u_{r+k}}{\|u_{r+k}\|_2}}\|_h^2. \quad (82)$$

Then using Lemma 2.2:

$$\begin{aligned}
u^T G u &\leq 2(1 + \gamma)(1 + (k - 1)\mu) \sum_{r=1, k} (|u_r|^2 + a_r^2 \|u_{k+r}\|_2^2 |\rho''(0)|) \\
&\leq 2(1 + \gamma)(1 + (k - 1)\mu) \sup_{\|u\|_2=1} \sum_{r=1, k} (|u_r|^2 + \|u_{k+r}\|_2^2 a_r^2 |\rho''(0)|) \quad (83) \\
&= 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (a_r^2 |\rho''(0)|)_{r=1, l}).
\end{aligned}$$

□

*Proof of Theorem 2.1.* Let  $\theta^*$  a minimizer of (7). Consider  $H$  the Hessian of  $g$  at  $\theta$ . We recall that  $H = G + F$  (see Proposition 2.2). Using Lemma A.3, we just need to bound the operator norm of  $F$  and then to combine it with the bounds on the eigenvalues of  $G$  to get bounds on eigenvalues of  $H = G + F$ .

We use Lemma 2.7, the Cauchy-Schwartz and triangle inequalities. We have  $\|A\delta''_{t_r, j_1, j_2}\|_2 \leq \sqrt{m}D_{A, R}$  and

$$\begin{aligned}
|F_{2, r, j_1, s, j_2}| &\leq \mathbf{1}(r = s) 2|a_r| \|A\delta''_{t_r, j_1, j_2}\|_2 \|A\phi(\theta) - y\|_2 \\
&\leq \mathbf{1}(r = s) 2|a_r| \sqrt{m}D_{A, R} \|A\phi(\theta) - A\phi(\theta^*) + A\phi(\theta^*) - y\|_2. \quad (84) \\
&\leq \mathbf{1}(r = s) 2|a_r| \sqrt{m}D_{A, R} (\|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2).
\end{aligned}$$

Similarly, with Lemma 2.6,

$$\begin{aligned}
F_{12, r, s, j} &\leq \mathbf{1}(r = s) 2\sqrt{1 + \gamma} \|\delta'_{t_r, j}\|_h \|A\phi(\theta) - y\|_2 \\
&\leq \mathbf{1}(r = s) 2\sqrt{1 + \gamma} \sqrt{|\rho''(0)|} (\|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2). \quad (85)
\end{aligned}$$

Let  $\|\cdot\|_{op}$  be the  $\ell^2$  operator norm of a matrix. With Gerschgorin circle theorem [18], we have

$$\|F\|_{op} \leq \max_l \|F_{l, \cdot}\|_1 \quad (86)$$

where  $F_{l, \cdot}$  is the  $l$ -th row of  $F$ . We get

$$\begin{aligned}
\|F\|_{op} &\leq \max(d \max_{r, s, j} |F_{12, r, s, j}|, \max_{r, s, j} |F_{12, r, s, j}| + d \max_{r, j_1, s, j_2} |F_{2, r, j_1, s, j_2}|) \\
&\leq (d + 1) \max(\max_{r, s, j} |F_{12, r, s, j}|, \max_{r, j_1, s, j_2} |F_{2, r, j_1, s, j_2}|) \\
&\leq 2(d + 1) \max_r (|a_r| \sqrt{m}D_{A, R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2). \quad (87)
\end{aligned}$$

Hence, using Weyl's perturbation inequalities on  $H = G + F$ , i.e.  $\lambda_{min}(H) \geq \lambda_{min}(G) - \lambda_{max}(F)$  and  $\lambda_{max}(H) \leq \lambda_{max}(G) + \lambda_{max}(F)$ , we get the result.

□

*Proof of Corollary 2.1.* First, observe that at  $\theta_0$ ,  $F = 0$ .

The upper bound is a direct consequence of Theorem A.3.

We show the result in the case  $\max(1, (a_r^2|\rho''(0)|)_{r=1,l}) \neq 1$  and  $\min(1, (a_r^2|\rho''(0)|)_{r=1,l}) \neq 1$  (the proof is similar in the other case). For the lower bound let  $v \in \mathbb{R}^{k(d+1)}$  and  $i_0 = \arg \max_{r=1,l} (a_r^2|\rho''(0)|)$ , set  $\|v_{i_0}\|_2 = 1$  and  $v_j = 0$  for  $j \neq i_0$ . With Equation (75), we have

$$\sup_{\|u\|_2=1} u^T H u \geq v^T H v \geq 2(1 - \gamma) \max(1, (a_r^2|\rho''(0)|)_{r=1,l}). \quad (88)$$

Similarly, let  $v \in \mathbb{R}^{k(d+1)}$  and  $i_0 = \arg \min((a_r^2|\rho''(0)|)_{r=1,l})$ ,  $\|v_{i_0}\|_2 = 1$  and  $v_j = 0$  for  $j \neq i_0$ .

$$\inf_{\|u\|_2=1} u^T H u \leq 2(1 + \gamma) \min(1, (a_r^2|\rho''(0)|)_{r=1,l}). \quad (89)$$

□

#### A.4 Proofs for Section 3

*Proof of Theorem 3.1.* Let  $\theta^* = (a_1, \dots, a_k, t_1, \dots, t_k) \in \Theta_{k,\epsilon}$  the global minimum of  $g$  and  $\theta = (b_1, \dots, b_k, s_1, \dots, s_k) \in \Lambda_{\theta^*,\beta}$ .

First notice that  $\|\theta - \theta^*\|_2^2 \leq \beta^2$  implies that for any  $j$ , we have  $|a_j - b_j|^2 \leq \beta^2$  and

$$|a_1| - \beta \leq |a_j| - \beta \leq |b_j| \leq |a_j| + \beta \leq |a_k| + \beta. \quad (90)$$

We also have  $\|s_j - t_j\|_2 < \beta \leq \frac{\epsilon}{4}$ . Hence for  $i \neq j$  we have  $\|s_i - s_j\|_2 = \|s_i - t_i + t_i - t_j + t_j - s_j\|_2 \geq \|t_i - t_j\|_2 - \|t_i - s_i\|_2 - \|t_j - s_j\|_2 > \epsilon - 2\epsilon/4 = \epsilon/2$  and  $\phi(\theta) \in \Sigma_{k,\frac{\epsilon}{2}}$ .

We use Theorem 2.1 to get the bound on the min and max eigenvalues of the Hessian.

We can then plug Inequality (90) into the one of Theorem 2.1.

Finally we notice the fact that  $\sup_{\theta \in \Lambda_{\theta^*,\beta}} \|A\phi(\theta) - A\phi(\theta^*)\|_2$  exists because  $\Lambda_{\theta^*,\beta}$  is bounded.

□

*Proof of Theorem 3.2.* This is a direct consequence of Theorem 2.1. The proof follows the same lines as the one of Theorem 3.1. □

*Proof of Corollary 3.1.* The set  $\Lambda = \Lambda_{\theta^*,\beta}$  is an open set where the Hessian of  $g$  at  $\Lambda$  is positive as long as  $\xi \leq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (|a_1| - \beta)^2|\rho''(0)|)$  with Theorem 3.1.

In this case  $g$  is convex on  $\Lambda$ . Theorem 3.1 also gives a uniform bound for the operator norm of the Hessian:  $\|H\|_{op} \leq 2(1 + \gamma)(1 + (k - 1)\mu) \max(1, (|a_k| + \beta)^2|\rho''(0)|) + \xi$  and  $g$  has Lipschitz gradient. Moreover the gradient descent on Lipschitz smooth convex functions guarantees that  $\|\theta_n - \theta^*\|_2$  decreases, where  $\theta_n$  are the iterates of the gradient descent (this is proved using direct consequences of the nonexpensiveness of  $\langle \tau \nabla g(\theta), \cdot \rangle$  [1, Proposition 4.2 (iv)]). Hence the iterates  $\theta_n$  stay in  $\Lambda$  and we deduce from Corollary 1.1 that  $\Lambda$  is a basin of attraction.

Hence we just need to show that  $\xi \leq 2(1-\gamma)(1-(k-1)\mu) \min(1, (|a_1|-\beta)^2|\rho''(0)|)$ . Let  $\theta \in \Lambda$ , we have, with the RIP hypothesis,

$$\begin{aligned} \xi(\theta) &:= 2(d+1) \max(\max_r |a_r| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \|A\phi(\theta) - A\phi(\theta^*)\|_2 \\ &\leq 2(d+1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \sqrt{1+\gamma} \|\phi(\theta) - \phi(\theta^*)\|_h \\ &\leq 2(d+1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \sqrt{1+\gamma} \sqrt{1+(k-1)\mu} \sqrt{\sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2} \end{aligned} \quad (91)$$

where we wrote  $\theta^* = \sum_i a_i \delta_{t_i}$  and  $\theta = \sum_i b_i \delta_{s_i}$  such that  $|s_i - t_i| \leq \epsilon/4$ .

We now bound the term  $\sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2$ :

$$\begin{aligned} \sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2 &= \sum_i a_i^2 + b_i^2 - 2a_i b_i \rho(\|s_i - t_i\|_2) \\ &= \sum_i \rho(\|s_i - t_i\|_2) |a_i - b_i|^2 + (1 - \rho(\|s_i - t_i\|_2)) \sum_i a_i^2 + b_i^2 \end{aligned} \quad (92)$$

Using the hypothesis that  $\|\theta - \theta^*\|^2 \leq \beta^2$  and  $\beta \leq |a_1|/2$ , we have  $|b_i| \leq |a_i| + \beta \leq \frac{3}{2}|a_i|$ . With the assumption on  $h$  (and  $\rho$ ),

$$\begin{aligned} \sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2 &\leq \beta^2 + \frac{|\rho''(0)|}{2} \beta^2 \frac{13}{4} \|a^*\|_2^2 \\ &\leq \beta^2 + \frac{|\rho''(0)|}{2} \beta^2 4 \|a^*\|_2^2 \\ &\leq \beta^2 (1 + 2|\rho''(0)| \|a^*\|_2^2) \end{aligned} \quad (93)$$

where  $a^* = (a_1, \dots, a_k)$ . The fact that  $\beta \leq |a_1|/2$  implies

$$\begin{aligned} &\frac{\xi(\theta)}{\min(1, (|a_1|-\beta)^2|\rho''(0)|)} \\ &\leq \frac{2(d+1) \sqrt{1+\gamma} \sqrt{1+(k-1)\mu} \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \sqrt{1+2|\rho''(0)|} \|a^*\|_2^2 \beta}{\min(1, |a_1|^2 |\rho''(0)|/4)} \end{aligned} \quad (94)$$

Hence using the hypothesis that

$$\beta \leq \frac{(1-\gamma)(1-(k-1)\mu) \min(1, |a_1|^2 |\rho''(0)|/4)}{(d+1) \sqrt{1+\gamma} \sqrt{1+(k-1)\mu} \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1+\gamma} \sqrt{|\rho''(0)|}) \sqrt{1+2|\rho''(0)|} \|a^*\|_2^2}$$

we have

$$\xi(\theta) \leq 2(1-\gamma)(1-(k-1)\mu) \min(1, (|a_1|(1-\beta))^2 |\rho''(0)|). \quad (95)$$

□

*Proof of Corollary 3.2.* Following the same argument as Corollary 3.2, we just need to show that  $\xi \leq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (|a_1| - \beta)^2 |\rho''(0)|)$ . Let  $\theta \in \Lambda_{\theta^*, \beta}$ , we have, with the RIP hypothesis,

$$\begin{aligned} \xi(\theta) &:= 2(d + 1) \max(\max_r |a_r| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\|A\phi(\theta) - A\phi(\theta^*)\|_2 + \|e\|_2) \\ &\leq 2(d + 1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\sqrt{1 + \gamma} \|\phi(\theta) - \phi(\theta^*)\|_h + \|e\|_2) \\ &\leq 2(d + 1) \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (\sqrt{1 + \gamma} \sqrt{1 + (k - 1)\mu} \sqrt{\sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2} + \|e\|_2) \end{aligned} \quad (96)$$

where we wrote  $\theta^* = \sum_i a_i \delta_{t_i}$  and  $\theta = \sum_i b_i \delta_{s_i}$  such that  $|s_i - t_i| \leq \epsilon/4$ .

Similarly as in Corollary 3.2, we bound the term  $\sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2$ :

$$\sum_i \|a_i \delta_{t_i} - b_i \delta_{s_i}\|_h^2 \leq \beta^2 (1 + 2|\rho''(0)| \|a^*\|_2^2) \quad (97)$$

The fact that  $\beta \leq |a_1|/2$  and  $\|e\|_2 \leq \sqrt{1 + \gamma} \sqrt{1 + (k - 1)\mu} \beta$  implies

$$\begin{aligned} &\frac{\xi(\theta)}{\min(1, (|a_1| - \beta)^2 |\rho''(0)|)} \\ &\leq \frac{2(d + 1) \sqrt{1 + \gamma} \sqrt{1 + (k - 1)\mu} \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (1 + \sqrt{1 + 2|\rho''(0)| \|a^*\|_2^2}) \beta}{\min(1, |a_1|^2 |\rho''(0)|/4)} \end{aligned} \quad (98)$$

Hence using the hypothesis that

$$\beta \leq \frac{(1 - \gamma)(1 - (k - 1)\mu) \min(1, |a_1|^2 |\rho''(0)|/4)}{(d + 1) \sqrt{1 + \gamma} \sqrt{1 + (k - 1)\mu} \max(|a_k| \sqrt{m} D_{A,R}, \sqrt{1 + \gamma} \sqrt{|\rho''(0)|}) (1 + \sqrt{1 + 2|\rho''(0)| \|a^*\|_2^2})},$$

we have

$$\xi(\theta) \leq 2(1 - \gamma)(1 - (k - 1)\mu) \min(1, (|a_1|(1 - \beta))^2 |\rho''(0)|). \quad (99)$$

□

## A.5 Proofs for Section 4

*Proof of Lemma 4.1.* Remark that  $g(\theta)$  does not depend on the ordering of the positions. Reorder  $\theta_0 = (a, t)$  and  $\theta_1 = (b, s)$  such that  $t_1 < t_2 \dots < t_k$  and  $s_1 < s_2 \dots < s_k$ . Consider the function  $g_1(\lambda) = g(\theta_\lambda)$  with  $\theta_\lambda = (1 - \lambda)\theta_0 + \lambda\theta_1$ . Remark that  $g_1$  is a continuous function of  $\lambda$  taking values  $g_1(0) = g(\theta_0)$  and  $g_1(1) = g(\theta_1)$ . Hence, with the

intermediate value theorem, there is  $\lambda$  such that  $g(\theta_\lambda) = g_1(\lambda) = \alpha$ . Moreover, denoting  $\theta_\lambda = (a_\lambda, t_\lambda)$ , we have, using the sorting of  $t$  and  $s$ , for  $1 \leq i < k$ ,

$$\begin{aligned} |t_{\lambda, i+1} - t_{\lambda, i}| &= |(1 - \lambda)t_{i+1} + \lambda s_{i+1} - (1 - \lambda)t_i - \lambda s_i| \\ &= (1 - \lambda)|t_{i+1} - t_i| + \lambda|s_{i+1} - s_i| > (1 - \lambda)\epsilon + \lambda\epsilon = \epsilon. \end{aligned} \tag{100}$$

Hence  $\theta_\lambda \in \Theta_{k, \epsilon}$ . □

## References

- [1] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*, volume 408. Springer, 2011.
- [2] B. N. Bhaskar, G. Tang, and B. Recht. Atomic norm denoising with applications to line spectral estimation. *IEEE Transactions on Signal Processing*, 61(23):5987–5999, 2013.
- [3] S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016.
- [4] T. Blumensath. Sampling and reconstructing signals from a union of linear subspaces. *IEEE Transactions on Information Theory*, 57(7):4660–4671, 2011.
- [5] A. Bourrier, M. Davies, T. Peleg, P. Perez, and R. Gribonval. Fundamental performance limits for ideal decoders in high-dimensional linear inverse problems. *Information Theory, IEEE Transactions on*, 60(12):7928–7946, 2014.
- [6] C. Boyer, Y. De Castro, and J. Salmon. Adapting to unknown noise level in sparse deconvolution. *Information and Inference: A Journal of the IMA*, 6(3):310–348, 2017.
- [7] V. Cambareri and L. Jacques. Through the haze: a non-convex approach to blind gain calibration for linear random sensing models. *Information and Inference: A Journal of the IMA*, 2018.
- [8] E. J. Candès and C. Fernandez-Granda. Super-resolution from noisy data. *Journal of Fourier Analysis and Applications*, 19(6):1229–1254, 2013.
- [9] E. J. Candès and C. Fernandez-Granda. Towards a mathematical theory of super-resolution. *Communications on Pure and Applied Mathematics*, 67(6):906–956, 2014.
- [10] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *arXiv preprint arXiv:1805.09545*, 2018.

- [11] P. G. Ciarlet, B. Miara, and J.-M. Thomas. *Introduction to numerical linear algebra and optimisation*. Cambridge University Press, 1989.
- [12] Y. De Castro, F. Gamboa, D. Henrion, and J.-B. Lasserre. Exact solutions to super resolution on semi-algebraic domains in higher dimensions. *arXiv preprint arXiv:1502.02436*, 2015.
- [13] V. Duval, P. Catala, and G. Peyré. A low-rank approach to off-the-grid sparse deconvolution. *preprint*, 2017.
- [14] V. Duval and G. Peyré. Exact support recovery for sparse spikes deconvolution. *Foundations of Computational Mathematics*, 15(5):1315–1355, 2015.
- [15] V. Duval and G. Peyré. Sparse regularization on thin grids i: the lasso. *Inverse Problems*, 33(5):055008, 2017.
- [16] V. Duval and G. Peyré. Sparse spikes super-resolution on thin grids ii: the continuous basis pursuit. *Inverse Problems*, 33(9):095008, 2017.
- [17] M. Golbabaee and M. E. Davies. Inexact gradient projection and fast data driven compressed sensing. *IEEE Transactions on Information Theory*, 2018.
- [18] G. H. Golub and C. F. Van Loan. *Matrix computations*, volume 3. JHU press, 2012.
- [19] R. Gribonval, G. Blanchard, N. Keriven, and Y. Traonmilin. Compressive Statistical Learning with Random Feature Moments. *Preprint*, 2017.
- [20] N. Keriven, A. Bourrier, R. Gribonval, and P. Pérez. Sketching for Large-Scale Learning of Mixture Models. *Preprint*, 2016.
- [21] N. Keriven, N. Tremblay, Y. Traonmilin, and R. Gribonval. Compressive k-means. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, pages 6369–6373. IEEE, 2017.
- [22] S. Ling and T. Strohmer. Regularized gradient descent: a non-convex recipe for fast joint blind deconvolution and demixing. *Information and Inference: A Journal of the IMA*, 2017.
- [23] C. Poon, N. Keriven, and G. Peyré. Support localization and the fisher metric for off-the-grid sparse regularization. *arXiv preprint arXiv:1810.03340*, 2018.
- [24] B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. R. Lanckriet. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 11(Apr):1517–1561, 2010.
- [25] G. Tang, B. N. Bhaskar, P. Shah, and B. Recht. Compressed sensing off the grid. *IEEE transactions on information theory*, 59(11):7465–7490, 2013.

- [26] A. F. Timan. *Theory of approximation of functions of a real variable*. International Series of Monographs on Pure and Applied Mathematics, vol. 34, Pergamon Press, Oxford, 1963.
- [27] Y. Traonmilin, N. Keriven, R. Gribonval, and G. Blanchard. Spikes super-resolution with random fourier sampling. In *SPARS workshop 2017*, 2017.
- [28] I. Waldspurger. Phase retrieval with random gaussian sensing vectors by alternating projections. *IEEE Transactions on Information Theory*, 2018.
- [29] T. Zhao, Z. Wang, and H. Liu. A nonconvex optimization framework for low rank matrix estimation. In *Advances in Neural Information Processing Systems*, pages 559–567, 2015.