# Support Localization and the Fisher Metric for off-the-grid Sparse Regularization

Clarice Poon1NicolasKeriven2Gabriel Peyré2<sup>1</sup>Centre for Mathematical Sciences<sup>2</sup>Département de Mathématiques et ApplicationsUniversity of CambridgeÉcole Normale SupérieureWilberforce Rd, Cambridge, United Kingdom45 rue d'Ulm, Paris, France

### Abstract

Sparse regularization is a central technique for both machine learning (to achieve supervised features selection or unsupervised mixture learning) and imaging sciences (to achieve super-resolution). Existing performance guaranties assume a separation of the spikes based on an ad-hoc (usually Euclidean) minimum distance condition, which ignores the geometry of the problem. In this article, we study the BLASSO (i.e. the off-the-grid version of  $\ell^1$  LASSO regularization) and show that the Fisher-Rao distance is the natural way to ensure and quantify support recovery, since it preserves the invariance of the problem under reparameterization. We prove that under mild regularity and curvature conditions, stable support identification is achieved even in the presence of randomized sub-sampled observations (which is the case in compressed sensing or learning scenario). On deconvolution problems, which are translation invariant, this generalizes to the multi-dimensional setting existing results of the literature. For more complex translation-varying problems, such as Laplace transform inversion, this gives the first geometry-aware guarantees for sparse recovery.

## 1 Introduction

### 1.1 Sparse Regularization

In this work, we consider the general problem of estimating an unknown Radon measure  $\mu_0 \in \mathcal{M}(\mathcal{X})$  defined over some metric space  $\mathcal{X}$  (for instance  $\mathcal{X} = \mathbb{R}^d$  for a possibly large d) from a few number m of randomized linear observations  $y \in \mathbb{C}^m$ , Let  $\Phi : \mathcal{M}(\mathcal{X}) \mapsto \mathbb{C}^m$  be defined by

$$\Phi \mu \stackrel{\text{def.}}{=} \frac{1}{\sqrt{m}} \left( \int_{\mathcal{X}} \varphi_{\omega_k}(x) \mathrm{d}\mu(x) \right)_{k=1}^m, \qquad (1.1)$$

where  $(\omega_1, \ldots, \omega_m)$  are identically and independently distributed according to some probability distribution  $\Lambda(\omega)$  on  $\omega \in \Omega$ , and for  $\omega \in \Omega$ ,  $\varphi_{\omega} : \mathcal{X} \to \mathbb{C}$  is a continuous function, denoted  $\varphi_{\omega} \in \mathscr{C}(\mathcal{X})$ . We further assume that  $\varphi_{\omega}(x)$  is normalized, that is

$$\mathbb{E}_{\omega}[|\varphi_{\omega}(x)|^2] = 1, \qquad \forall x \in \mathcal{X}.$$
(1.2)

The observations are  $y = \Phi \mu_0 + w$ , where  $w \in \mathbb{C}^m$  accounts for noise or modelling errors. Some representative examples of this setting include:

- Off-the-grid compressed sensing: off-the-grid compressed sensing, initially introduced in the special case of 1-D Fourier measurements on  $\mathcal{X} = \mathbb{T} = \mathbb{R}/\mathbb{Z}$  by (Tang et al., 2013), corresponds exactly to measurements of the form (1.1). This is a "continuous" analogous of the celebrated compressed sensing line of works (Candès et al., 2006; Donoho, 2006).
- Regression using a continuous dictionary: given a set of *m* training samples  $(\omega_k, y_k)_{k=1}^m$ , one wants to predicts the values  $y_k \in \mathbb{R}$  from the features  $\omega_k \in$  $\Omega$  using a continuous dictionary of functions  $\omega \mapsto$  $\varphi_{\omega}(x)$  (here  $x \in \mathcal{X}$  parameterizes the dictionary), as  $y_k \approx \int_{\mathcal{X}} \varphi_{\omega_k}(x) d\mu(x)$ . A typical example, studied for instance by Bach (2017) is the case of neural networks with a single hidden layer made of an infinite number of neurons, where  $\Omega = \mathcal{X} = \mathbb{R}^p$  and one uses ridge functions of the form  $\varphi_{\omega}(x) = \psi(\langle x, \omega \rangle)$ , for instance using the ReLu non-linearity  $\psi(u) = \max(u, 0)$ .
- Sketching mixtures: the goal is estimate a (hopefully sparse) mixture of density probability distributions on some domain  $\mathcal{T}$  of the form  $\xi(t) = \sum_{i} a_i \xi_{x_i}(t)$  where the  $(\xi_x)_{x \in \mathcal{X}}$  is a family of template densities,

and  $a_i \ge 0$ ,  $\sum_i a_i = 1$ . Introducing the measure  $\mu_0 = \sum_i a_i \delta_{x_i}$ , this mixture model is conveniently rewritten as  $\xi(t) = \int_{\mathcal{X}} \xi_x(t) d\mu_0(x)$ . The most studied example is the mixture of Gaussians, using (in 1-D for simplicity,  $\mathcal{T} = \mathbb{R}$ ) as  $\xi_x(t) \propto \sigma^{-1} e^{-\frac{(t-\tau)^2}{2\sigma^2}}$  where the parameter space is the mean and standard deviation  $x = (\tau, \sigma) \in \mathcal{X} = \mathbb{R} \times \mathbb{R}^+$ . In a typical machine learning scenario, one does not have direct access to  $\xi$  but rather to n i.i.d. samples  $(t_1, \ldots, t_n) \in \mathcal{T}^n$  drawn from  $\xi$ . Instead of recording this (possibly huge, specially when  $\mathcal{T}$  is high dimensional) set of data, following Gribonval et al. (2017), one computes "online" a small set  $y \in \mathbb{C}^m$  of m sketches against sketching functions  $\theta_{\omega}(t)$ , that is, for  $k = 1, \ldots, m$ ,

$$y_k \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{j=1}^n \theta_{\omega_k}(t_j) \approx \int_{\mathcal{T}} \theta_{\omega_k}(t) \xi(t) \mathrm{d}t.$$

These sketches exactly have the form (1.1) when defining the functions  $\varphi_{\omega}(x) \stackrel{\text{def.}}{=} \int_{\mathcal{T}} \theta_{\omega}(t)\xi_x(t)dt$ . A popular set of sketching functions, over  $\mathcal{T} = \mathbb{R}^d$  are Fourier atoms  $\theta_{\omega}(t) \stackrel{\text{def.}}{=} e^{i\langle\omega, t\rangle}$ , for which  $\varphi_{\cdot}(x)$  is the characteristic functions of  $\xi_x$ , which can generally be computed in closed form.

**BLASSO.** In all these applications, and many more, one is actually interested in recovering a discrete and *s*-sparse measure  $\mu_0$  of the form  $\mu_0 = \sum_{i=1}^s a_i \delta_{x_i}$  where  $(x_i, a_i) \in \mathcal{X} \times \mathbb{C}$ . An increasingly popular method to estimate such a sparse measure corresponds to solving a infinite-dimensional analogous of the Lasso regression problem

$$\min_{\mu \in \mathcal{M}(\mathcal{X})} \frac{1}{2} \left\| \Phi \mu - y \right\|_{2}^{2} + \lambda |\mu|(\mathcal{X}). \qquad (\mathcal{P}_{\lambda}(y))$$

Following De Castro and Gamboa (2012), we call this method the BLASSO (for Beurling-Lasso). Here  $|\mu|(\mathcal{X})$  is the so-called total variation of the measure  $\mu$ , and is defined as

$$|\mu|(\mathcal{X}) \stackrel{\text{\tiny def.}}{=} \sup \left\{ \operatorname{Re}\langle f, \mu \rangle \; ; \; f \in \mathscr{C}(\mathcal{X}), \left\| f \right\|_{\infty} \leqslant 1 \right\}.$$

Note that on unbounded  $\mathcal{X}$ , one needs to impose that f vanishes at infinity. If  $\mathcal{X} = \{x_i\}_i$  is a finite space, then this corresponds to the classical finite-dimensional Lasso problem (Tibshirani, 1996), because  $|\mu|(\mathcal{X}) = ||a||_1 \stackrel{\text{def.}}{=} \sum_i |a_i|$  where  $a_i = \mu(\{x_i\})$ . Similarly, if  $\mathcal{X}$  is possibly infinite but  $\mu = \sum_i a_i \delta_{x_i}$ , one also has that  $|\mu|(\mathcal{X}) = ||a||_1$ .

**Previous Works.** The BLASSO problem  $(\mathcal{P}_{\lambda}(y))$  was initially proposed by De Castro and Gamboa (2012), see also Bredies and Pikkarainen (2013). The first sharp analysis of the solution of this problem is

provided by Candès and Fernandez-Granda (2014) in the case of Fourier measurement on  $\mathbb{T}^d$ . They show that if the spikes are separated enough, then  $\mu_0$  is the unique solution of  $(\mathcal{P}_{\lambda}(y))$  when w = 0 and  $\lambda \to 0$ . Robustness to noise under this separation condition is addressed in (Candès and Fernandez-Granda, 2013; Fernandez-Granda, 2013: Azais et al., 2015). A refined stability results is detailed by Duval and Peyré (2015) which shows that conditions based on minimum separation imply support stability, which means that when ||w|| and  $||w|| / \lambda$  are small enough, then the solution of  $(\mathcal{P}_{\lambda}(y))$  has the same number of Diracs as  $\mu_0$ , and that both the amplitudes and positions of the spikes converges smoothly as  $w \to 0$ . These initial works have been extended by Tang et al. (2013) to the case of randomized compressive measurements of the form (1.1), when using Fourier sketching functions  $\varphi_{\omega}$ . In all these results, the separation condition are given for the Euclidean cases, which is an ad-hoc choice which does not take into account the geometry of the problem, and gives vastly sub-optimal theories for spatially varying operators (such as data-dependent kernels in supervised learning, Gaussian mixture estimation and Laplace transform in imaging, see Section 1.2).

While this is not the topic of the present paper, note that for positive spikes, the separation condition is in some cases not needed, see for instance (Schiebinger et al., 2015; Denoyelle et al., 2017). It is important to note that efficient algorithms have been developed to solve  $(\mathcal{P}_{\lambda}(y))$ , among which SDP relaxations for Fourier measurements (Candès and Fernandez-Granda, 2013) and Frank-Wolfe (also known as conditional gradient) schemes (Bredies and Pikkarainen, 2013; Boyd et al., 2017). Note also that while we focus here on variational convex approaches, alternative methods exist, in particular greedy algorithms (Gribonval et al., 2017) and (for Fourier measurements) Prony-type approaches (Schmidt, 1986; Roy and Kailath, 1989). To the best of our knowledge, their theoretical analysis in the presence of noise is more involved, see however (Liao and Fannjiang, 2016) for an analysis of robustness to noise when a minimum separation holds.

#### **1.2** The Fisher information metric

The empirial covariance operator is defined as  $\hat{K}(x,x') \stackrel{\text{def.}}{=} \frac{1}{m} \sum_{i} \overline{\varphi_{\omega_i}(x)} \varphi_{\omega_i}(x')$  and the deterministic limit as  $m \to +\infty$  is denoted K with

$$K(x, x') \stackrel{\text{def.}}{=} \int_{\Omega} \overline{\varphi_{\omega}(x)} \varphi_{\omega}(x') \mathrm{d}\Lambda(\omega).$$
(1.3)

Note that many covariance kernels can be written under the form (1.3). By Bochner's theorem, this includes all translation-invariant kernels, for which possible features are  $\varphi_{\omega}(x) = e^{i\omega^{\top}x}$ . The associated metric tensor is

$$\mathbf{H}_x \stackrel{\text{def.}}{=} \nabla_x \nabla_{x'} K(x, x) \in \mathbb{C}^{d \times d}. \tag{1.4}$$

Throughout, we assume that  $\mathbf{H}_x$  is positive definite for all  $x \in \mathcal{X}$ . Then,  $\mathbf{H}$  naturally induces a distance between points in our parameter space  $\mathcal{X}$ . Given a piecewise smooth curve  $\gamma : [0,1] \to \mathcal{X}$ , the length  $\ell_{\mathbf{H}}[\gamma]$ of  $\gamma$  is defined by  $\ell_{\mathbf{H}}[\gamma] \stackrel{\text{def.}}{=} \int_0^1 \sqrt{\langle \mathbf{H}_{\gamma(t)}\gamma'(t), \gamma'(t) \rangle} dt$ . Given two points  $x, x' \in \mathcal{X}$ , the distance from x to x', induced by  $\mathbf{H}$  is  $d_{\mathbf{H}}(x, x') \stackrel{\text{def.}}{=} \inf_{\gamma \in \mathcal{F}} \ell_{\mathbf{H}}[\gamma]$  where  $\mathcal{F}$  is the set of all piecewise smooth paths  $\gamma : [0, 1] \to \mathcal{X}$ with  $\gamma(0) = x$  and  $\gamma(1) = x'$ .

The metric **H** is closely linked to the Fisher information matrix (Fisher, 1925) associated with  $\Phi$ : since (1.2) holds,  $f(x,\omega) \stackrel{\text{def.}}{=} |\varphi_{\omega}(x)|^2$  can be interpreted as a probability density function for the random variable  $\omega$ conditional on parameter x, and the metric  $\mathbf{H}_x$  is equal (up to rescaling) to its Fisher information matrix, since

$$\int \nabla \left( \log f(x,\omega) \right) \nabla \left( \log f(x,\omega) \right)^{\top} f(x,\omega) d\Lambda(\omega)$$
$$= 4 \mathbb{E}_{\omega} \left[ \operatorname{Re} \left( \overline{\nabla \varphi_{\omega}(x)} \nabla \varphi_{\omega}(x)^{\top} \right) \right] = 4 \mathbf{H}_{x}.$$

The distance  $d_{\mathbf{H}}$  is called the "Fisher-Rao" geodesic distance (Rao, 1945) and is used extensively in information geometry for estimation and learning problems on parametric families of distributions (Amari and Nagaoka, 2007). The Fisher-Rao is the unique Riemannian metric on a statistical manifold (Cencov, 2000) and it is invariant to reparameterization, which matches the invariance of the BLASSO problem ( $\mathcal{P}_{\lambda}(y)$ ) to reparameterization of the space  $\mathcal{X}$ . Although  $d_{\mathbf{H}}$  has been used in conjunction with kernel methods (see for instance Burges (1999)), to the best of our knowledge, it is the first time this metric is put forward to analyze the performance of off-the-grid sparse recovery problems.

#### 1.2.1 Examples

We detail some popular learning and imaging examples.

The Fejér kernel One of the first seminal result of super-resolution with sparse regularization was given by Candès and Fernandez-Granda (2014) for this kernel, which corresponds to discrete Fourier measurements on the torus. We give a multi-dimensional generalization of this result here. Let  $f_c \in \mathbb{N}$ ,  $\mathcal{X} \in \mathbb{T}^d$ ,  $\Omega = \{\omega \in \mathbb{Z}^d ; \|\omega\|_{\infty} \leq f_c\}$ . Let  $\varphi_{\omega}(x) \stackrel{\text{def.}}{=} e^{i2\pi\omega^T x}$  and  $\Lambda(\omega) \propto \prod_{j=1}^d g(\omega_j)$  where  $g(j) = \frac{1}{f_c} \sum_{k=\max(j-f_c,-f_c)}^{\min(j+f_c,f_c)} (1-|k/f_c|)(1-|(j-k)/f_c|)$ . Note that this corresponds to sampling discrete Fourier frequencies. Then, the associated kernel is the Fejér kernel  $K(x, x') = \prod_{i=1}^d \kappa(x_i - x'_i)$ , where  $\kappa(x) \stackrel{\text{def.}}{=}$ 

 $\operatorname{sinc}_{f_c/2+1}^4(x)$  where  $\operatorname{sinc}_s(x) \stackrel{\text{def.}}{=} s^{-1} \sin(\pi s x) / \sin(\pi x)$ , which has a constant metric tensor  $\mathbf{H}_x = C_{f_c} \operatorname{Id}$  and  $d_{\mathbf{H}}(x, x') = \sqrt{C_{f_c}} \|x - x'\|_2$  is a scaled Euclidean metric (quotiented by the action of translation modulo 1 on  $\mathbb{T}^d$ ), where  $C_{f_c} = -\kappa''(0) = \frac{\pi^2 f_c(f_c+4)}{3}$ .

The Gaussian kernel Let  $\Sigma \in \mathbb{R}^{d \times d}$  be a positive semidefinite matrix,  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $\Omega = \mathbb{R}^d$ . Let  $\varphi_{\omega}(x) = e^{i\omega^\top x}$  and  $\Lambda(\omega) = \mathcal{N}(0, \Sigma^{-1})$ , the centered Gaussian distribution with covariance  $\Sigma^{-1}$ . This can be interpreted as sampling *continuous* Fourier frequencies. Then, the associated kernel is  $K(x, x') = e^{-\frac{1}{2} ||x-x'||_{\Sigma^{-1}}^2}$  where  $||x||_{\Sigma} = \sqrt{x^\top \Sigma x}$ , with constant metric  $\mathbf{H}_x = \Sigma^{-1}$ , and  $d_{\mathbf{H}}(x, x') = ||x - x'||_{\Sigma^{-1}}$ . In Section 3, we also detail how to exploit this kernel for Gaussian Mixture Model (GMM) estimation with the BLASSO.

The Laplace transform Let  $\bar{\alpha} = (\alpha_j) \in \mathbb{R}^d_+, \mathcal{X} \subseteq (0, +\infty)^d$  and  $\Omega = \mathbb{R}^d_+$ . A (sampled) Laplace transform is defined by setting  $\varphi_{\omega}(x) = \prod_{i=1}^d \sqrt{\frac{2(x_i + \alpha_i)}{\alpha_i}} e^{-\langle x, \omega \rangle}$ and  $\Lambda(\omega) = \prod_{j=1}^d (2\alpha_j) e^{-\langle 2\bar{\alpha}, \omega \rangle}$ . Then,  $K(x, x') = \prod_{i=1}^d \kappa(x_i + \alpha_i, x'_i + \alpha_i)$  where  $\kappa(a, b) = \frac{2\sqrt{ab}}{a+b}$ , with metric  $\mathbf{H}_x$  as the diagonal matrix with diagonal  $((2(x_i + \alpha_i))^{-2})_{i=1}^d$  and distance  $d_{\mathbf{H}}(x, x') = \sqrt{\sum_i \left| \log \left(\frac{x_i + \alpha_i}{x'_i + \alpha_i}\right) \right|^2}$ . We remark that this kernel, associated to the Laplace transform (which should not be confused with the translation-invariant Laplace kernel  $\exp(-||x - x'||)$ ) appears in some microscopy imaging technique, see for instance Boulanger et al. (2014). Unlike the previous examples, it is not translationinvariant, and therefore the metric  $\mathbf{H}_x$  is not constant. Our results show that the corresponding Fisher metric is the natural way to impose the separation condition in super-resolution.

#### 1.3 Contributions.

Our main contribution is Theorem 1, which states that if the sought after spikes positions  $X_0$  are sufficiently separated with respect to the Fisher distance  $d_{\mathbf{H}}$ , then the solution to  $(\mathcal{P}_{\lambda}(y))$  is support stable (that is, the solution of the BLASSO is formed of exactly *s* Diracs) provided that the number of random noisy measurements *m* is, up to log factors and under the assumption of random signs of the amplitudes  $a_0$ , linear in *s*, and the noise level ||w|| is less than 1/s. In the case of translation invariant kernels, this generalizes existing results to a large class of multi-dimensional kernels, and also provides for the first time a quantitative bounds on the impact of the noise and sub-sampling on the spikes positions and amplitudes errors. For non-translation kernels, this provides for the first time a meaningful support recovery guarantee, a typical example being the Laplace kernel (see Section 1.2).

### 2 Key concepts

Notation for derivatives. Given  $f \in \mathscr{C}^{\infty}(\mathscr{X})$ , by interpreting the  $r^{th}$  derivative as a multilinear map:  $\nabla^r f : (\mathbb{C}^d)^r \to \mathbb{C}$ , so given  $Q \stackrel{\text{def.}}{=} \{q_\ell\}_{\ell=1}^r \in (\mathbb{C}^d)^r$ ,

$$\nabla^r f[Q] = \sum_{i_1, \cdots, i_r} \partial_{i_1} \cdots \partial_{i_r} f(x) q_{1,i_1} \cdots q_{r,i_r}.$$

and we define the  $r^{th}$  normalized derivative of f as

$$\mathbf{D}_r [f](x)[Q] \stackrel{\text{def.}}{=} \nabla^r f(x)[\{\mathbf{H}_x^{-\frac{1}{2}}q_i\}_{i=1}^r]$$

with norm  $\|\mathbf{D}_r[f](x)\| \stackrel{\text{def.}}{=} \sup_{\forall \ell, \|q_\ell\| \leq 1} |\mathbf{D}_r[f](x)[Q]|$ . For  $i, j \in \{0, 1, 2\}$ , let  $K^{(ij)}(x, x')$  be a "bi"-multilinear map, defined for  $Q \in (\mathbb{C}^d)^i$  and  $V \in (\mathbb{C}^d)^j$  as

$$[Q]K^{(ij)}(x,x')[V] \stackrel{\text{def.}}{=} \mathbb{E}[\overline{\mathcal{D}_i[\varphi_\omega](x)[Q]}\mathcal{D}_j[\varphi_\omega](x')[V]]$$

and  $\|K^{(ij)}(x,x')\| \stackrel{\text{def.}}{=} \sup_{Q,V} \|[Q]K^{(ij)}(x,x')[V]\|$ where the supremum is defined over all  $Q \stackrel{\text{def.}}{=} \{q_\ell\}_{\ell=1}^i$ ,  $V \stackrel{\text{def.}}{=} \{v_\ell\}_{\ell=1}^j$  with  $\|q_\ell\| \leq 1$ ,  $\|v_\ell\| \leq 1$ . Note that  $D_2[f](x)$  and  $K^{(02)}(x,x')$  can also be interpreted as a matrix in  $\mathbb{C}^{d \times d}$ , and we have the normalization  $K^{(02)}(x,x) = -\text{Id}$  for all x.

### 2.1 Admissible kernel and separation

In previous studies on the recovery properties of  $(\mathcal{P}_{\lambda}(y))$  (Candès and Fernandez-Granda, 2014; Bhaskar et al., 2013; Bendory et al., 2016; Duval and Peyré, 2015; Fernandez-Granda, 2016), recovery bounds are attained in the context of K being *admissible* and a separation condition on the underlying positions  $\{x_j\}_j$ . Namely, given  $X = \{x_j\}_j$ , that  $\min_{i\neq j} d_{\mathbf{H}}(x_i, x_j)$  is sufficiently large with respect to the decay properties of K. For example, in the case where  $\Phi$  corresponds to Fourier sampling on a grid, up to frequency  $f_c$ , this separation condition is  $\min_{j\neq \ell} ||x_j - x_\ell||_2 \gtrsim 1/f_c$ . In fact, if  $\operatorname{sign}(a_j)$  can take arbitrary values in  $\{+1, -1\}$ , this separation condition is a necessary to ensure exact recovery for the BLASSO (Tang, 2015).

Following the aforementioned works, we introduce the notion of an admissible kernel.

**Definition 1.** A kernel K will be said admissible with respect to  $\mathcal{K} \stackrel{\text{def.}}{=} \{r_{\text{near}}, \Delta, \varepsilon_i, B_{ij}, s_{\max}\}, \text{ where } 0 < r_{\text{near}} < \Delta/4 \text{ is a neighborhood size, } \varepsilon_0 \in (0, 1), \ \varepsilon_2 \in (0, r_{\text{near}}^{-2}) \text{ are respectively a distance to } 1 \text{ and a curvature, } \Delta > 0 \text{ is a minimal separation, } B_{ij} > 0 \text{ for } i, j = 0, \dots, 2 \text{ are some constants and } s_{\max} \in \mathbb{N}^* \text{ is a maximal sparsity level, if}$ 

- 1. Uniform bounds: For  $(i, j) \in \{(0, 0), (1, 0)\},$   $\sup_{x,x' \in \mathcal{X}} \|K^{(ij)}(x, x')\| \leq B_{ij};$  for  $(i, j) \in \{(0, 2), (1, 1), (1, 2)\}$  and all x, x'such that  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$  or  $d_{\mathbf{H}}(x, x') > \Delta/4, \|K^{(ij)}(x, x')\| \leq B_{ij};$  and finally,  $\sup_{x \in \mathcal{X}} \|K^{(22)}(x, x)\| \leq B_{22}.$
- 2. Neighborhood of each point: For all  $x \in \mathcal{X}$ , K(x,x) = 1 and for all  $x, x' \in \mathcal{X}$  with  $d_{\mathbf{H}}(x,x') \leq r_{\text{near}}$ ,  $\operatorname{Re}\left(K^{(02)}(x,x')\right) \leq -\varepsilon_{2}\operatorname{Id}$  and  $\left\|\operatorname{Im}\left(K^{(02)}(x,x')\right)\right\| \leq c\varepsilon_{2}$ , where  $c \stackrel{\text{def.}}{=} \frac{1}{2}\sqrt{\frac{2-\varepsilon_{2}r_{\text{near}}^{2}}{\varepsilon_{2}r_{\text{near}}^{2}}}$ and for  $d_{\mathbf{H}}(x,x') \geq r_{\text{near}}$ ,  $|K(x,x')| \leq 1 - \varepsilon_{0}$ .
- 3. Separation: For  $d_{\mathbf{H}}(x, x') \ge \Delta/4$ , for all  $i, j \in \{0, \ldots, 2\}$  with  $i + j \le 3$ ,  $||K^{(ij)}(x, x')|| \le \frac{h}{s_{\max}}$ , where  $h \stackrel{\text{def.}}{=} \min_{i \in \{0,2\}} \left(\frac{\varepsilon_i}{32B_{1i}+32}, \frac{5\varepsilon_2}{16B_{12}+24}\right)$ .

Additionally, there exists  $C_{\mathbf{H}} \ge 0$  such that for  $d_{\mathbf{H}}(x, x_0) \leqslant r_{\text{near}}$ :  $\left\| \operatorname{Id} - \mathbf{H}_{x_0}^{-\frac{1}{2}} \mathbf{H}_x^{\frac{1}{2}} \right\| \leqslant C_{\mathbf{H}} d_{\mathbf{H}}(x, x_0).$ We also denote  $d_{\mathbf{H}}(X, X_0) = \sqrt{\sum_i d_{\mathbf{H}}(x_i, x_{0,i})^2}$  and  $B \stackrel{\text{def.}}{=} \sum_{i+j \leqslant 3} B_{ij}$  and  $\varepsilon \stackrel{\text{def.}}{=} \min\{\varepsilon_0, \varepsilon_2\}.$ 

Intuitively, these three conditions express the following facts: 1) the kernel and its derivatives are uniformly bounded, 2) near x = x', the kernel has negative curvature, and otherwise it is strictly less than 1, and 3) for x and x' sufficiently separated, the kernel and all its derivatives have a small value.

#### 2.2 Almost bounded random features

Ideally, we would like our features and its derivatives to be uniformly bounded for all  $\omega$ . However this may not be the case: think of  $e^{i\omega^{\top}x}$  where the support of the distribution  $\Lambda$  is not bounded. Hence our results will be dependent on the probability that the derivatives are greater than some value T decays sufficiently quickly as T increases. In the following, for  $r \in \{0, 1, 2, 3\}$ ,  $L_r(\omega) \stackrel{\text{def.}}{=} \sup_{x \in \mathcal{X}} \|D_r [\varphi_{\omega}](x)\|$ , and let  $F_r$  be such that  $\mathbb{P}_{\omega} (L_r(\omega) > t) \leq F_r(t)$ .

### 2.3 Key assumptions

Our main result will be valid under the following assumptions.

**I.** On the domain and limit kernel Let  $\mathcal{X}$  be a compact domain with radius  $R_{\mathcal{X}} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} d_{\mathbf{H}}(x,x')$ . Assume the kernel is admissible wrt  $\mathcal{K} \stackrel{\text{def.}}{=} \{r_{\text{near}}, \Delta, \varepsilon_i, B_{ij}, s_{\text{max}}\}.$ 

**II. Assumption on the underlying signal** For  $s \leq s_{\max}$ , let  $a_0 \in \mathbb{C}^s$  and let  $X_0 \stackrel{\text{def.}}{=} (x_{0,j})_{j=0}^s$  be such that  $d_{\mathbf{H}}(x_{0,i}, x_{0,j}) \geq \Delta$  for  $i \neq j$ . The underlying measure is assumed to be  $\mu_0 = \sum_{j=1}^s a_{0,j} \delta_{x_{0,j}}$ .

III. Assumption on the sampling complexity For  $\rho > 0$ , suppose that  $m \in \mathbb{N}$  and  $\{\overline{L}_i\}_{i=0}^3 \in \mathbb{R}^4_+$  are chosen such that

$$\sum_{j=0}^{3} F_j(\bar{L}_j) \leqslant \frac{\rho}{m}, \quad \text{and}$$

$$\max_{j=0}^{3} \{\bar{L}_j^2 \sum_{i=0}^{3} F_i(\bar{L}_i) + 2 \int_{\bar{L}_j}^{\infty} tF_j(t) dt\} \leqslant \frac{\varepsilon}{m},$$
(2.1)

and either one of the following hold:

$$m \gtrsim C \cdot s \cdot \log\left(N^d/\rho\right) \log\left(s/\rho\right),$$
 (2.2)

or 
$$m \gtrsim C \cdot s^{3/2} \cdot \log\left(N^d/\rho\right)$$
, (2.3)

where  $C \stackrel{\text{def.}}{=} \varepsilon^{-2} (\bar{L}_2^2 B_{11} + \bar{L}_1^2 B_{22} + (B_0 + B_2) \bar{L}_{01}^2), N \stackrel{\text{def.}}{=} \mathbb{L}_3 dR_{\mathcal{X}} (r_{\text{near}} \varepsilon)^{-1} \text{ and } \mathbb{L}_r = \max_{i=1}^r \bar{L}_i.$ 

Remark 1. Our main theorem presents support stability guarantees under the sampling complexity rate (2.2) if sign $(a_0) = (a_{0,i}/|a_{0,i}|)_{i=1}^s$  forms a Steinhaus sequence, that is, iid uniformly distributed on the complex unit circle. This assumption has been used before in compressed sensing (Candès and Romberg, 2007; Tang et al., 2013) to achieve this optimal complexity (see also Foucart and Rauhut (2013), Chap. 14). As noted in previous works, this random signs assumption is likely to be a proof artefact, however achieving optimal complexity without it may require more involved arguments (Candes and Plan, 2011). When the signs are arbitrary, we prove our results under (2.3). Although this  $s^{3/2}$  scaling is still sub-optimal in s, we remark it improves upon the previous theoretical rate of  $s^2$  (up to log factors) (Li and Chi, 2017).

**Remark 2.** The assumption on the choice of  $\bar{L}_r$  ensures that with high probability,  $D_r [\varphi_{\omega}](x)$  is uniformly bounded up to r = 3. Note also that, generally, the  $\{\bar{L}_r\}$  depend on m, through (2.1). However, in all our examples: either a)  $\sup_{x \in \mathcal{X}} \|D_r [\varphi_{\omega}](x)\|$  are already uniformly bounded, in which case  $\bar{L}_i$  can be chosen independently of  $\rho$  and m (for instance this is the case of the Fejér kernel); or b) the  $F_r(t)$  are exponentially decaying, in which case we can show that  $\bar{L}_r = \mathcal{O}(\log(m/\rho)^p)$  for some p > 0, which only incurs additional logarithmic terms on the bounds (2.2) and (2.3). This is the case of the Gaussian or Laplace transform kernel.

### 3 Main result

Our main theorem below states quantitative exact support recovery bounds under a minimum separation condition according to  $d_{\rm H}$ .

**Theorem 1.** Let  $\rho > 0$ , suppose that K is admissible, and that  $a_0$ ,  $X_0$ , m and  $\overline{L}_i$  satisfy

the assumptions of Section 2.3. Let  $\mathcal{D}_{\lambda_0,c_0} \stackrel{\text{def.}}{=} \{(\lambda, w) \in \mathbb{R}_+ \times \mathbb{C}^m ; \lambda \leq \lambda_0, \|w\| \leq c_0 \lambda\}$  where  $c_0 \sim \min\left(\frac{\varepsilon_0}{L_0}, \frac{\varepsilon_2}{L_2}\right)$  and  $\lambda_0 \sim D/s$  with

$$D \stackrel{\text{def.}}{=} \underline{a} \min \left( r_{\text{near}} \sqrt{s}, \ \frac{\varepsilon \sqrt{s}}{\mathbb{L}_2^2 \|a\|}, \ \frac{\varepsilon}{C_{\mathbf{H}}(B + \mathbb{L}_2^2)} \right)$$
(3.1)

where  $\underline{a} = \min\{|a_{0,i}|, |a_{0,i}|^{-1}\}$ . Suppose that either sign $(a_0)$  is a Steinhaus sequence and m satisfies (2.2) or sign $(a_0)$  is an arbitrary sign sequence and m satisfies (2.3). Then, with probability at least  $1 - \rho$ ,

- (i) for all  $v \stackrel{\text{def.}}{=} (\lambda, w) \in \mathcal{D}_{\lambda_0, c_0}$ ,  $(\mathcal{P}_{\lambda}(y))$  has a unique solution which consists of exactly s spikes. Moreover, up to a permutation of indices, the solution can be written as  $\sum_{i=1}^{s} a_i^v \delta_{x_i^v}$ , and  $\operatorname{sign}(a_i^v) = \operatorname{sign}(a_{0,i})$  for all  $i = 1, \ldots, s$
- (ii) The mapping  $v \in \mathcal{D}_{\lambda_0,c_0} \mapsto (a^v, X^v)$  is  $\mathscr{C}^1$  and we have the error bound

$$a^{v} - a_{0} \| + d_{\mathbf{H}}(X^{v}, X_{0}) \leqslant \frac{\sqrt{s}(\lambda + \|w\|)}{\min_{i} |a_{0,i}|}$$
 (3.2)

We detail below the values relating to the sampling complexity corresponding to each of the examples detailed in Section 1.2.1. The corresponding proofs can be found in Section F of the appendix.

**Discrete Fourier sampling** The Fejer kernel of order  $f_c \ge 128$  is admissible with  $\Delta = \mathcal{O}(\sqrt{d\sqrt[4]{s_{\text{max}}}}),$  $r_{\text{near}} = 1/(8\sqrt{2}), \ \varepsilon_0 = 0.00097, \ \varepsilon_2 = 0.941, \ B_{01} =$  $\mathcal{O}(d), B_{11} = B_{02} = B_{12} = \mathcal{O}(1) \text{ and } B_{22} = \mathcal{O}(d).$ Moreover,  $\bar{L}_r = \mathcal{O}(d^{r/2})$ . Hence, up to logarithmic terms, Thm. 1 is applicable with  $m = \mathcal{O}(sd^3)$  when the random signs assumption holds, and  $m = \mathcal{O}(s^{\frac{3}{2}}d^3)$ in the general case, with guaranteed support stability when  $\lambda = O(s^{-1}d^{-2}), \|w\| = O(s^{-1}d^{-3})$ . Note that our choice of  $\Delta$  imposes that  $\|x_i - x_j\|_2 \gtrsim \sqrt{d} s_{\max}^{1/4} / f_c$ whereas the previous result of Candès and Fernandez-Granda (2014) requires  $||x_i - x_j||_{\infty} \gtrsim C_d/f_c$  with no dependency in  $s_{\text{max}}$ , however, their proof would imply that the constant  $C_d$  grows exponentially in d. Since we are interested in having a general theory in arbitrary dimension, we have opted to present a polynomial dependency on  $s_{\max}$ .

**Continuous Gaussian Fourier sampling** In the appendix we prove that the kernel is admissible with  $\Delta = \mathcal{O}\left(\sqrt{\log s_{\max}}\right), r_{\text{near}} = 1/\sqrt{2}, \varepsilon_0 = 1 - e^{-\frac{1}{4}}, \varepsilon_2 = e^{-\frac{1}{4}}/2, B_{ij} = \mathcal{O}(1) \text{ for } i+j \leq 3, B_{22} = \mathcal{O}(d) \text{ and } \bar{L}_r = \left(d + \log\left(\frac{dm}{\rho}\right)^2\right)^{\frac{r}{2}}$  (as mentioned before, the dependence in *m* only incurs additional logarithmic factors in (2.2) and (2.3)). Hence, up to log factors, the

sample complexity and noise level for the application of Thm. 1 is the same as for the Fejér kernel.

Laplace sampling The associated kernel is admissible with  $\Delta = \mathcal{O}(d + \log(ds_{\max}))$ ,  $r_{\text{near}} = 0.2$ ,  $\varepsilon_0 = 0.005$ ,  $\varepsilon_2 = 1.52$ ,  $B_{ij} = \mathcal{O}(1)$  for  $i + j \leq 3$  and  $B_{22} = \mathcal{O}(d)$ . Define  $\bar{R}_{\mathcal{X}} = \left(1 + \frac{R_{\mathcal{X}}}{\min_i \alpha_i}\right)^d$  (where we recall that  $R_{\mathcal{X}}$  is the radius of  $\mathcal{X}$ ). Assuming for simplicity that all  $\alpha_j$  are distinct, we can set  $\bar{L}_r = \bar{R}_{\mathcal{X}}(R_{\mathcal{X}} + \|\alpha\|_{\infty})^r \left(\sqrt{d} + \max_i \frac{1}{\alpha_i} \log\left(\frac{d\beta_i m \bar{R}_{\mathcal{X}}}{\rho \alpha_i}\right)\right)^r$ Hence, choosing  $\alpha_i \sim d$ , we have that  $\bar{R}_{\mathcal{X}} = (1)$ and up to log factors, (2.2) is  $\mathcal{O}(sd^7)$  and (2.3) is  $\mathcal{O}(s^{3/2}d^7)$ , and support stability is guaranteed when  $\lambda = \mathcal{O}(s^{-1}d^{-3})$  and  $\|w\| = \mathcal{O}(s^{-1}d^{-5})$ . Note that despite the stronger dependency on d, for practical applications (microscopy), one is typically only interested in the low dimensional setting of d = 2, 3.

**Gaussian mixture learning** Consider *n* datapoints  $z_1, \ldots, z_n \in \mathbb{R}^d$  drawn *iid* from a mixture of Gaussians  $\sum_i a_{0,i} \mathcal{N}(x_{0,i}, \Sigma)$  with means  $x_{0,i} \in \mathcal{X} \subset \mathbb{R}^d$  and known covariance  $\Sigma$ , where  $\mathcal{X}$  is bounded. Consider the following procedure:

- draw  $\omega_j$  *iid* from  $\mathcal{N}(0, \Sigma^{-1}/d)$  (the 1/d normalization is necessary to avoid an exponential dependency in d later on)
- compute the generalized moments  $y = \frac{1}{\sqrt{m}} \sum_{i=1}^{n} (e^{i \langle \omega_j, x_i \rangle})_{j=1}^m$
- solve the BLASSO with features  $\varphi_{\omega}(x) = e^{i\langle\omega,x\rangle}e^{-\frac{1}{2}\|\omega\|_{\Sigma}^2}$ , to obtain a distribution  $\tilde{\mu}$

Then, as described in the introduction, we can interpret y as noisy Fourier measurements of  $\mu_0 = \sum_i a_{0,i} \delta_{x_{0,i}}$  in the space of means  $\mathcal{X}$ , where the "noise" w corresponds to using the empirical average over the  $z_i$  instead of a true integration. It is easily bounded with probability  $1-\rho$  by  $||w|| \leq \mathcal{O}\left(\sqrt{\frac{\log(1/\rho)}{n}}\right)$ , by a simple application of Hoeffding's inequality (Gribonval et al., 2017).

The associated kernel is the Gaussian kernel with covariance  $(2 + d)\Sigma$  and hence, our result states that, if  $||x_i - x_j||_{\Sigma^{-1}} \ge \sqrt{d \log s}$ , and the number of measurements and sample complexity satisfy, up to logarithmic terms,  $m = \mathcal{O}\left(s^{\frac{3}{2}}d^3\right)$ ,  $n = \mathcal{O}\left(s^2d^6/\min_i|a_{0,i}|^2\right)$  and  $\lambda_0 = \mathcal{O}\left(\frac{\min_i|a_{0,i}|}{\sqrt{sd^2}||a_0||_2}\right)$ , then, with probability  $1 - \rho$  on both samples  $z_j$  and frequencies  $\omega_j$ , the distribution  $\tilde{\mu}$  is formed of exactly *s* Diracs, and their positions and weights converge to the means and weights of the GMM. Let us give a few remarks on this result.

*On model selection.* Besides convexity (with respect to the distribution of means) of the BLASSO, which is

not the case of classical likelihood- or moments-based methods for learning GMM, the most striking feature of our approach is probably the support stability: with a sample complexity that is polynomial in s and d, the BLASSO yields *exactly* the right number of components for the GMM. Despite the huge literature on model selection for GMM, to our knowledge, this is one of the only result which is *non-asymptotic* in sample complexity, as opposed to many approaches (Roeder and Wasserman, 1997; Huang et al., 2013) which guarantee that the selected number of components approaches the correct one when the number of samples grows to infinity.

On separation condition. Our separation condition of  $\sqrt{d\log s}$  is, up to the logarithmic term, similar to the  $\sqrt{d}$  found in the seminal work by Dasgupta (1999). This was later improved by different methods (Dasgupta and Schulman, 2000; Vempala and Wang, 2004), until the most recent results on the topic (Moitra and Valianty, 2010) show that it is possible to learn a GMM with no separation condition, provided the sample complexity is exponential in s, which is a necessary condition (Moitra and Valianty, 2010). As mentioned in the introduction, similar results exist for the BLASSO: Denoyelle et al. (2017) showed that in one dimension, one can identify s positive spikes with no separation, provided the noise level is exponentially small with s. Hence learning GMM with the BLASSO and no separation condition may be feasible, which we leave for future work, however we note that the multi-dimensional case is still largely an open problem (Poon and Pevré, 2017).

On known covariance. An important path for future work is to handle arbitrary covariance. When the components all share the same mean and have diagonal covariance, the Fisher metric is related, up to a change of variables, to the Laplace transform kernel case treated earlier. When both means and covariance vary, in one dimension, the Fisher metric is related to the Poincaré half-plane metric (Costa et al., 2015). In the general case, it does not have a closed-form expression. We leave the treatment of these cases for future work.

### 4 Sketch of proof

#### 4.1 Background on dual certificates

Our approach to establishing that the solutions to  $(\mathcal{P}_{\lambda}(y))$  are support stable is via the study of the associated dual solutions in accordance to the framework introduced in Duval and Peyré (2015). We first recall some of their key ideas. In order to study the support stability properties of  $(\mathcal{P}_{\lambda}(y))$  in the small noise regime, we consider the limit problem as  $\lambda \to 0$  and  $||w|| \to 0$ ,

that is

Þ

$$\min_{\mu \in \mathcal{M}(\mathcal{X})} |\mu|(\mathcal{X}) \text{ subject to } \Phi \mu = y. \qquad (\mathcal{P}_0(y))$$

The dual of  $(\mathcal{P}_{\lambda}(y))$  and  $(\mathcal{P}_{0}(y))$  are

$$\min_{p} \left\{ \|y/\lambda - p\|_{2}^{2} ; \|\Phi^{*}p\|_{\infty} \leqslant 1 \right\} \qquad (\mathcal{D}_{\lambda}(y))$$

$$\max_{p} \left\{ \langle y, p \rangle ; \left\| \Phi^* p \right\|_{\infty} \leq 1 \right\}. \qquad (\mathcal{D}_0(y))$$

Any solution  $\mu_{\lambda}$  of  $(\mathcal{P}_{\lambda}(y))$  to related to the (unique) solution  $p_{\lambda}$  of  $(\mathcal{D}_{\lambda}(y))$  by  $-p_{\lambda} = \frac{1}{\lambda}(\Phi\mu_{\lambda} - y)$  and writing  $\eta_{\lambda} \stackrel{\text{def.}}{=} \Phi^* p_{\lambda}, \langle \eta_{\lambda}, \mu_{\lambda} \rangle = |\mu_{\lambda}|(\mathcal{X})$ . Note that  $\operatorname{Supp}(\mu_{\lambda}) \subseteq \{x \in \mathcal{X} ; |\Phi^* p_{\lambda}(x)| = 1\}$ , so  $\eta_{\lambda}$  "certifies" the support of  $\mu_{\lambda}$  and is often referred to as a *dual certificate*. Furthermore, by defining the minimal norm certificate  $\eta_0$  as  $\eta_0 \stackrel{\text{def.}}{=} \Phi^* p_0$  where

$$p_0 = \operatorname{argmin} \{ \|p\|_2 \; ; \; p \text{ is a solution to } (\mathcal{D}_0(y)) \}$$
(4.1)

one can show that  $p_{\lambda}$  converges as  $\lambda \to 0$  to  $p_0$  and hence  $\eta_{\lambda}$  converges to  $\eta_0 \stackrel{\text{def.}}{=} \Phi^* p_0$  in  $L^{\infty}$ . When  $\lambda$ and ||w|| are sufficiently small, solutions to  $(\mathcal{P}_{\lambda}(y))$  are support stable provided that  $\eta_0$  (called the minimal norm certificate) is *nondegenerate*, that is  $\eta_0(x_i) =$  $\operatorname{sign}(a_i)$  for  $i = 1, \ldots, s$  and  $\nabla^2 |\eta_0|^2(x_i)$  is negative definite. This is proven to be an almost sharp condition for support stability, since Duval and Peyré (2017) provided explicit examples where  $|\eta_0(x)| = 1$  for some  $x \notin \{x_i\}_i$  implies that  $(\mathcal{P}_{\lambda}(y))$  recovers more than sspikes under arbitrarily small noise.

**Pre-certificates** In practice, the minimal norm certificate is hard to compute and analyse due to the nonlinear  $\ell^{\infty}$  constraint in (4.1). So, one often introduces a proxy which can be computed in closed form by solving an linear system associated to the following least squares problem:  $\eta_X \stackrel{\text{def.}}{=} \Phi^* p$  where

$$p_X \stackrel{\text{def.}}{=} \operatorname{argmin}\{\|p\|_2 ; (\Phi^* p)(x_i) = \operatorname{sign}(a_i), \\ \nabla(\Phi^* p)(x_i) = 0\}.$$
(4.2)

Note that if  $\eta_X$  satisfies  $\|\eta_X\|_{\infty} \leq 1$ , then  $\eta_X = \eta_0$ .

**Computation of**  $\eta_X$  For  $x \in \mathcal{X}$ , let  $\varphi(x) \stackrel{\text{def.}}{=} \frac{1}{\sqrt{m}} (\varphi_{\omega_k}(x))_{k=1}^m$ . For  $X = \{x_i\}_{i=1}^s$  we define  $\Gamma_X$ :  $\mathbb{C}^{s(d+1)} \to \mathbb{C}^m$  as  $\Gamma_X([\alpha, \beta]) \stackrel{\text{def.}}{=} \sum_{i=1}^s \alpha_i \varphi(x_i) + \nabla \varphi(x_i)^\top \beta_i$  where  $\nabla \varphi \in \mathbb{C}^{m \times d}$ . Then, the minimizer of (4.2) is  $p_X = \Gamma_X^{*,\dagger} (\stackrel{(\text{sign}(a))}{\mathbf{0}_{sd}})$ . Furthermore, when  $\Gamma_X$ is full rank, we can write  $\hat{\eta}_X(x) \stackrel{\text{def.}}{=} \sum_i \hat{\alpha}_i \hat{K}(x_i, x) + \langle \hat{\beta}_i, \nabla_1 \hat{K}(x_i, x) \rangle$ , where  $\hat{\alpha}_i \in \mathbb{C}$ ,  $\hat{\beta}_i \in \mathbb{C}^d$  are such that  $(\hat{\beta}) = (\Gamma_X^* \Gamma_X)^{-1} (\stackrel{(\text{sign}(a))}{\mathbf{0}_{sd}})$ , and the hat notation refers to the fact that we are using sub-sampled measurements. The limit precertificate is defined as  $\eta_X(x) \stackrel{\text{def.}}{=} \sum_i \alpha_i K(x_i, x) + \langle \beta_i, \nabla_1 K(x_i, x) \rangle$ , where  $\binom{\alpha}{\beta} = (\mathbb{E}[\Gamma_X^* \Gamma_X])^{-1} \binom{\operatorname{sign}(a)}{0_{sd}}.$ 

The key to establishing our recovery results is to show that  $\hat{\eta}_X$  is nondegenerate. In this paper, we will actually prove a stronger notion of nondegeneracy:

**Definition 2.** Let  $a \in \mathbb{C}^s$ ,  $X = \{x_i\}_{i=1}^s \in \mathcal{X}^s$  for some  $s \in \mathbb{N}$ , and  $\varepsilon_0, \varepsilon_2, r > 0$ . We say that  $\eta \in \mathscr{C}^1(\mathcal{X})$ is  $(\varepsilon_0, \varepsilon_2)$ -nondegenerate with respect to a, X and r if for all i,  $\eta(x_i) = \operatorname{sign}(a_i)$ ,  $\nabla \eta(x_i) = 0$  and

$$\begin{aligned} \forall x \in \mathcal{X}^{\text{far}}, \ |\eta(x)| &\leq 1 - \varepsilon_0 \\ \forall x \in \mathcal{X}^{\text{near}}_i, |\eta(x)| &\leq 1 - \varepsilon_2 d_{\mathbf{H}}(x, x_j)^2 \end{aligned}$$

where  $\mathcal{X}_{j}^{\text{near}} \stackrel{\text{def.}}{=} \{x \in \mathcal{X} ; d_{\mathbf{H}}(x_{i}, x) \leq r\}$  and  $\mathcal{X}^{\text{far}} \stackrel{\text{def.}}{=} \mathcal{X} \setminus \bigcup_{j=1}^{s} \mathcal{X}_{j}^{\text{near}}.$ 

Our proof proceeds in three steps:

- 1. Show that under admissibility of the kernel and sufficient separation, the limit precertificate  $\eta_{X_0}$  is non-degenerate (see Theorem 2).
- 2. Show that this non-degeneracy transfers to  $\hat{\eta}_X$  when m is large enough and X is close to  $X_0$ . This is the purpose of Section 4.3.
- 3. As discussed, nondegeneracy of  $\hat{\eta}_{X_0}$  automatically guarantees support stability when  $(\lambda, w) \in \mathcal{D}_{\lambda_0, c_0}$ for  $\lambda_0$  and  $c_0$  sufficiently small. To conclude we simply need to quantify  $\lambda_0$  and  $c_0$ . This is the purpose of Section 4.4. In particular, given  $(\lambda, w)$ , we construct a candidate solution by means of (a quantitative version of) the Implicit Function Theorem, and show that it is indeed a true solution using the previous results.

#### 4.2 Non-degeneracy of the limit certificate

Our first result shows that the "limit precertificate"  $\eta_{X_0}$  is nondegenerate:

**Theorem 2.** Assume the kernel is admissible wrt  $\mathcal{K}$  (see Definition 1). Then, for  $s \leq s_{\max}$ , for all  $a = (a_j)_{j=1}^s \in \mathbb{C}^s$  and  $X = \{x_j\}_{j=1}^s \in \mathcal{X}^s$  such that  $d_{\mathbf{H}}(x_i, x_j) \geq \Delta$ , the function  $\eta_{X_0}$  is  $(\frac{\varepsilon_0}{2}, \frac{\varepsilon_2}{2})$ -nondegenerate with respect to a, X and  $r_{\operatorname{near}}$ .

The proof of this result can be found in Appendix B and is a generalization of the arguments of Candès and Fernandez-Granda (2014) (see also Bendory et al. (2016)). We remark that unlike previous works which focus on translation invariant kernels, the Fisher metric provides a natural way to understand the required separation between the points in X and thus open up the possibility of analysing more complex problems such as Laplace transform inversion.

#### 4.3 The randomized setting

For the remainder of this paper, we consider solutions of  $(\mathcal{P}_{\lambda}(y))$  given  $y = \Phi \mu_{a_0, X_0} + w$  for some fixed  $a_0 \in \mathbb{C}^s$  and  $X_0 \in \mathcal{X}^s$ . The following result shows that  $\hat{\eta}_X$  is nondegenerate for all X close to  $X_0$ :

**Theorem 3.** Let  $\rho > 0$ . Under the assumptions of Section 2.3, and assuming that either m satisfies (2.2) and sign( $a_0$ ) is a Steinhaus sequence, or m satisfies (2.3) and sign( $a_0$ ) is an arbitrary sign sequence, with probability at least  $1 - \rho$ : for all  $X \in \mathcal{X}^s$  such that

$$d_{\mathbf{H}}(X, X_0) \lesssim \min\left(r_{\text{near}}, \frac{\varepsilon_r}{C_{\mathbf{H}}\sqrt{s}\max\left(B, \bar{L}_{12}\bar{L}_r\right)}\right), \quad (4.3)$$

 $\Gamma_X$  is full rank and  $\hat{\eta}_X$  is  $(\varepsilon_0/8, \varepsilon_2/8)$ -nondegenerate with respect to  $a_0$ , X and  $r_{\text{near}}$ .

The proof of this result is given in Appendix D. We simply make a remark on the proof here: We first prove that  $\hat{\eta}_{X_0}$  is nondegenerate by bounding variations between  $\eta_{X_0}$  and  $\hat{\eta}_{X_0}$ . The proof of this fact is a generalization of the arguments in Tang et al. (2013) to the multidimensional and general operator case. We then exploit the fact the  $\varphi$  is smooth and hence,  $\Gamma_X^* \Gamma_X$ satisfies certain Lipschitz properties with respect to X, to bound the local variation between  $\hat{\eta}_X$  and  $\hat{\eta}_{X_0}$ .

#### 4.4 Quantitative support recovery

This final section concludes the proof of Theorem 1 by quantifying the regions for  $\lambda$  and ||w|| for which support stability is guaranteed.

Solution of the noisy BLASSO. Let  $\Phi_X : \mathbb{C}^s \to \mathbb{C}^m$  be defined by  $\Phi_X a = \sum_{i=1}^s a_i \varphi(x_i)$ . Recall that  $\mu_{a,X} = \sum_i a_i \delta_{x_i}$  is a solution to the BLASSO with  $y = \Phi \mu_{a_0,X_0} + w$  if and only if  $\hat{\eta}_{\lambda} = \Phi^* p_{\lambda}$ , with  $p_{\lambda} = \frac{1}{\lambda}(y - \Phi_X a)$ , satisfies  $\|\hat{\eta}_{\lambda}\|_{\infty} \leq 1$  and  $\hat{\eta}(x_j) = \operatorname{sign}(a_j)$ . In that case,  $p_{\lambda}$  is the *unique* solution to the dual of the BLASSO. Moreover, if  $|\hat{\eta}_{\lambda}(x)| < 1$  for  $x \neq x_i$  and  $\Phi_X$  is full rank (which follows by Theorem D.2), then  $\mu_{a,X}$  is also the unique solution of the primal.

**Construction of a solution** Following Denoyelle et al. (2017), we define the function  $f : \mathbb{C}^s \times \mathcal{X}^s \times \mathbb{R}_+ \times \mathbb{C}^m$  by

$$f(u,v) \stackrel{\text{def.}}{=} \Gamma_X^*(\Phi_X a - \Phi_{X_0} a_0 - w) + \lambda \begin{pmatrix} \operatorname{sign}(a_0) \\ 0_{sd} \end{pmatrix}$$

where u = (a, X) and  $v = (\lambda, w)$ . Observe that having f(u, v) = 0 ensures the existence of  $\hat{\eta}_{\lambda}$  defined as above that satisfies  $\hat{\eta}_{\lambda}(x_i) = \operatorname{sign}(a_{0,i})$  and  $\nabla \hat{\eta}_{\lambda}(x_i) = 0$ . We will use it to construct a nondegenerate solution to  $\mathcal{D}_{\lambda}(y)$  for small  $\lambda$  and ||w||. Now, f is continuously differentiable, with explicit forms of

 $\partial_v f(u,v)$  and  $\partial_u f(u,v)$  given in (E.1) and (E.2) in the appendix, and in particular, letting  $u_0 = (a_0, X_0)$ ,  $\partial_u f(u_0, 0) = \Gamma^*_{X_0} \Gamma_{X_0} J_a$ , where  $J_a$  is the diagonal matrix with  $\binom{1}{a} \otimes 1_d \in \mathbb{C}^{s(d+1)}$  along its diagonal and  $\Gamma_{X_0}$ is full rank (with probability at least  $1 - \rho$ ) by Theorem D.2. So,  $\partial_u f(u_0, 0)$  is invertible and  $f(u_0, 0) = 0$ . Hence, by the Implicit Function Theorem, there exists a neighbourhood V of 0 in  $\mathbb{C} \times \mathbb{C}^m$ , a neighbourhood U of  $u_0$  in  $\mathbb{C}^s \times \mathcal{X}^s$  and a Fréchet differentiable function  $g: V \to U$  such that for all  $(u, v) \in U \times V$ , f(u, v) = 0if and only if u = q(v). So, to establish support stability for  $(\mathcal{P}_{\lambda}(y))$ , we simply need to estimate the size of the neighbourhood V on which q is well defined, and given  $(\lambda, w) \in V$ , for  $(a, Z) = g((\lambda, w))$ , to check that the associated certificate  $\hat{\eta}_{\lambda,w} \stackrel{\text{def.}}{=} \Phi^* p_{\lambda,w}$  with  $p_{\lambda,w} \stackrel{\text{def.}}{=} \frac{1}{\lambda} (\Phi_X a - \Phi_{X_0} a_0 - w) \text{ is nondegenerate.}$ 

Indeed, one can prove (see Theorem E.1) that with probability at least  $1 - \rho$ , V contains the ball  $B_r(0)$  with radius  $r \sim \frac{1}{\sqrt{s}} \min\left(\frac{\min\{r_{\text{near}}, (C_{\mathbf{H}}B)^{-1}\}}{\min_i |a_{0,i}|}, \frac{1}{\bar{L}_{01}\bar{L}_{12}(1+||a_0||)}\right)$  and given any  $v \in B_r(0)$ , (a, X) = g(v) indeed satisfy the error bound (3.2).

Checking that the candidate solution is a true solution It remains to check that  $g(\lambda, w)$  defines a valid certificate and is non-degenerate (and hence,  $\sum_i a_i \delta_{x_i}$  is the unique solution to  $(\mathcal{P}_{\lambda}(y))$ ) provided that  $\lambda, w$  satisfy (3.1). Given  $(\lambda, w) \in V$ , let  $(a, X) = g((\lambda, w))$ . Define  $\hat{\eta}_{\lambda,w} \stackrel{\text{def.}}{=} \frac{1}{\lambda} \Phi^*(\Phi_X a - \Phi_{X_0} a_0 - w)$  and following Denoyelle et al. (2017), one can show that

$$\hat{\eta}_{\lambda,w} = \hat{\eta}_X + \varphi(\cdot)^\top \Pi_X \frac{w}{\lambda} + \frac{1}{\lambda} \varphi(\cdot)^\top \Pi_X \Phi_{X_0} a_0$$

where  $\Pi_X$  is the orthogonal projection onto  $\operatorname{Im}(\Gamma_X)^{\perp}$ .

Note that since we have the error bound (3.2), our choice of  $\lambda$  and ||w|| ensures that (4.3) holds and hence, Theorem D.2 implies that  $\hat{\eta}_X$  is nondegenerate with probability at least  $1 - \rho$ . To conclude, it is sufficient to show that the two remaining terms are sufficiently small, so that  $\hat{\eta}_{\lambda,w}$  remains non-degenerate. Under  $\bar{E}$ ,  $||D_r[\varphi_{\omega}](\cdot)|| \leq \bar{L}_r$ , and for any  $z \in \mathbb{C}^m$ ,  $||D_r[\varphi^\top z] \cdot || \leq \bar{L}_r ||z||$ . Therefore, since  $\Pi_X$  is a projection, we have  $||D_r[\varphi(\cdot)^\top \Pi_X \frac{w}{\lambda}]|| \lesssim \varepsilon_r$  when  $||w||/\lambda \lesssim \varepsilon_r/\bar{L}_r$ . Finally, since  $\Phi_{X_0}a_0 = \sum_{j=1}^s \varphi(x_{0,j})$ , by Taylor expansion of  $\varphi(x_{0,j})$  around  $x_j$  and applying  $\Pi_X$  (see Lemma E.1 for this computation), we have

$$\left\|\frac{1}{\lambda}\Pi_X\Gamma_{X_0}\begin{pmatrix}a_0\\0_{sd}\end{pmatrix}\right\| \leqslant \frac{\bar{L}_2}{\lambda} \|a_0\|_{\infty} d_H(X,X_0)^2.$$

Since g satisfies (3.2) our choice of  $\lambda_0 = \mathcal{O}(s^{-1})$  ensures that we can upper bound this

by  $\bar{L}_2 \|a_0\|_{\infty} \frac{s\left(\lambda + \|w\|^2/\lambda\right)}{\min|a_{0,i}|^2} \lesssim \varepsilon$  and consequently,  $\frac{1}{\lambda} \left\| \mathbf{D}_r \left[ \varphi(\cdot)^\top \Pi_X \Phi_{X_0} a_0 \right] \right\| \lesssim \varepsilon_r.$ 

### Acknowledgements

We would like to thank Ben Adcock for a helpful conversation regarding the stochastic gradient bounds. This work was partly funded by the CFM-ENS chair "Modèles et Sciences des données" and the European Research Council, NORIA project.

#### References

- M. Akkouchi. On the convolution of exponential distributions.
- S.-i. Amari and H. Nagaoka. Methods of information geometry, volume 191. American Mathematical Soc., 2007.
- J.-M. Azais, Y. De Castro, and F. Gamboa. Spike detection from inaccurate samplings. Applied and Computational Harmonic Analysis, 38(2):177–195, 2015.
- F. Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017.
- T. Bendory, S. Dekel, and A. Feuer. Robust recovery of stream of pulses using convex optimization. *Journal of mathematical analysis and applications*, 442(2):511–536, 2016.
- 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.
- J. Boulanger, C. Gueudry, D. Münch, B. Cinquin, P. Paul-Gilloteaux, S. Bardin, C. Guérin, F. Senger, L. Blanchoin, and J. Salamero. Fast high-resolution 3D total internal reflection fluorescence microscopy by incidence angle scanning and azimuthal averaging. *Proceedings of the National Academy of Sciences*, 111(48):17164–17169, 2014.
- N. Boyd, G. Schiebinger, and B. Recht. The alternating descent conditional gradient method for sparse inverse problems. SIAM Journal on Optimization, 27(2):616–639, 2017.
- K. Bredies and H. K. Pikkarainen. Inverse problems in spaces of measures. ESAIM: Control, Optimisation and Calculus of Variations, 19(1):190–218, 2013.
- C. J. Burges. Geometry and invariance in kernel based methods. 1999.
- E. Candès and J. Romberg. Sparsity and incoherence in compressive sampling. *Inverse Problems*, 23(3):969–985, 2007.
- E. J. Candès and C. Fernandez-Granda. Super-resolution from noisy data. *Journal of Fourier Analysis and Appli*cations, 19(6):1229–1254, 2013.
- 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.
- E. J. Candes and Y. Plan. A probabilistic and RIPless theory of compressed sensing. *IEEE Transactions on Information Theory*, 57(11):7235–7254, 2011.

- E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on information theory*, 52(2):489–509, 2006.
- N. N. Cencov. Statistical decision rules and optimal inference. Number 53. American Mathematical Soc., 2000.
- S. I. Costa, S. A. Santos, and J. E. Strapasson. Fisher information distance: A geometrical reading. In *Discrete Applied Mathematics*, volume 197, pages 59–69. Elsevier B.V., 2015.
- S. Dasgupta. Learning mixtures of Gaussians. In IEEE 51st Annual Symposium on Foundations of Computer Science, number May, 1999.
- S. Dasgupta and A. Gupta. An Elementary Proof of a Theorem of Johnson and Lindenstrauss. *Random Structures* and Algorithms, 22(1):60–65, 2003.
- S. Dasgupta and L. J. Schulman. A Two-Round Variant of EM for Gaussian Mixtures. Uncertainty in Artificial Intelligence, pages 152–159, 2000.
- Y. De Castro and F. Gamboa. Exact reconstruction using Beurling minimal extrapolation. *Journal of Mathematical Analysis and applications*, 395(1):336–354, 2012.
- Q. Denoyelle, V. Duval, and G. Peyré. Support recovery for sparse super-resolution of positive measures. to appear in Journal of Fourier Analysis and Applications, 2017.
- D. L. Donoho. Compressed sensing. IEEE Transactions on information theory, 52(4):1289–1306, 2006.
- V. Duval and G. Peyré. Exact support recovery for sparse spikes deconvolution. Foundations of Computational Mathematics, 15(5):1315–1355, 2015.
- V. Duval and G. Peyré. Sparse spikes super-resolution on thin grids I: the LASSO. *Inverse Problems*, 33(5):055008, 2017.
- C. Fernandez-Granda. Support detection in superresolution. Proc. Proceedings of the 10th International Conference on Sampling Theory and Applications, pages 145–148, 2013.
- C. Fernandez-Granda. Super-resolution of point sources via convex programming. *Information and Inference: A Journal of the IMA*, 5(3):251–303, 2016.
- R. A. Fisher. Theory of statistical estimation. In Mathematical Proceedings of the Cambridge Philosophical Society, volume 22, pages 700–725. Cambridge University Press, 1925.
- S. Foucart and H. Rauhut. A Mathematical Introduction to Compressive Sensing. Applied and Numerical Harmonic Analysis. Springer New York, NY, 2013.
- R. Gribonval, G. Blanchard, N. Keriven, and Y. Traonmilin. Compressive statistical learning with random feature moments. arXiv preprint arXiv:1706.07180, 2017.
- T. Huang, H. Peng, and K. Zhang. Model Selection for Gaussian Mixture Models. *Statistica Sinica*, pages 1–27, 2013.
- M. Ledoux and M. Talagrand. *Probability in Banach Spaces: isoperimetry and processes.* Springer Science & Business Media, 2013.
- Y. Li and Y. Chi. Stable separation and super-resolution of mixture models. *Applied and Computational Harmonic Analysis*, 2017.

- W. Liao and A. Fannjiang. MUSIC for single-snapshot spectral estimation: Stability and super-resolution. Applied and Computational Harmonic Analysis, 40(1):33–67, 2016.
- S. Minsker. On some extensions of bernstein's inequality for self-adjoint operators. *Statistics & Probability Letters*, 127:111–119, 2017.
- A. Moitra and G. Valianty. Settling the polynomial learnability of mixtures of Gaussians. Proceedings - Annual IEEE Symposium on Foundations of Computer Science, FOCS, pages 93–102, 2010.
- C. Poon and G. Peyré. Multi-dimensional Sparse Superresolution. pages 1–42, 2017.
- C. R. Rao. Information and the accuracy attainable in the estimation of statistical parameters. *Bull. Calcutta Math. Soc.*, 37:81–91, 1945.
- K. Roeder and L. Wasserman. Practical Bayesian density estimation using mixtures of normal. JASA, 92:894–902, 1997.
- R. Roy and T. Kailath. ESPRIT-estimation of signal parameters via rotational invariance techniques. *IEEE Transactions on acoustics, speech, and signal processing*, 37(7):984–995, 1989.
- G. Schiebinger, E. Robeva, and B. Recht. Superresolution without separation. arXiv preprint arXiv:1506.03144, 2015.
- R. Schmidt. Multiple emitter location and signal parameter estimation. *IEEE transactions on antennas and propagation*, 34(3):276–280, 1986.
- K. Sridharan. A Gentle Introduction to Concentration Inequalities. Technical report, 2002.
- G. Tang. Resolution limits for atomic decompositions via markov-bernstein type inequalities. In Sampling Theory and Applications (SampTA), 2015 International Conference on, pages 548–552. IEEE, 2015.
- 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.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. Journal of the Royal Statistical Society. Series B (Methodological), pages 267–288, 1996.
- J. A. Tropp. An introduction to matrix concentration inequalities. (December), 2015.
- S. Vempala and G. Wang. A spectral algorithm for learning mixture models. *Journal of Computer and System Sciences*, 68(4):841–860, 2004.

### A Notations.

In this section, we recall and introduce some notation which will be used throughout the appendix.

**Block norms.** By default,  $\|\cdot\|$  is the Euclidean norm for vector and spectral norm for matrices. For a vector  $x = [x_1, \ldots, x_s] \in \mathbb{C}^{sd}$  formed of s blocks  $x_i \in \mathbb{C}^d$ ,  $1 \leq i \leq s$ , we define the block norm

$$\|x\|_{\text{block}} \stackrel{\text{def.}}{=} \sup_{1 \leqslant i \leqslant s} \|x_i\|_2$$

For a vector  $q = [q_1, \ldots, q_s, Q_1, \ldots, Q_s] \in \mathbb{C}^{s(d+1)}$  decomposed such that  $q_i \in \mathbb{C}$  and  $Q_i \in \mathbb{C}^d$ , we define

$$||q||_{*,\infty} \stackrel{\text{def.}}{=} \max_{i=1}^{s} \{ |q_i|, ||Q_i|| \}.$$

Kernel The empirical kernel is defined as

$$\hat{K}(x,x') = \frac{1}{m} \sum_{k=1}^{m} \overline{\varphi_{\omega_k}(x)} \varphi_{\omega_k}(x')$$

and the limit kernel is  $K(x,x) \stackrel{\text{def.}}{=} \mathbb{E}_{\omega}[\overline{\varphi_{\omega}(x)}\varphi_{\omega}(x')]$ . The metric tensor associated to this kernel is

$$\mathbf{H}_x \stackrel{\text{def.}}{=} \mathbb{E}_{\omega}[\overline{\nabla \varphi_{\omega}(x)} \nabla \varphi_{\omega}(x)^{\top}]$$

Given an event E, we write  $K_E(x, x') \stackrel{\text{def.}}{=} \mathbb{E}_{\omega}[\hat{K}(x, x')|E]$  to denote the conditional expectation on E.

**Derivatives** Given  $f \in \mathscr{C}^{\infty}(\mathcal{X})$ , by interpreting the  $r^{th}$  derivative as a multilinear map:  $\nabla^r f : (\mathbb{C}^d)^r \to \mathbb{C}$ , so given  $Q \stackrel{\text{def.}}{=} \{q_\ell\}_{\ell=1}^r \in (\mathbb{C}^d)^r$ ,

$$\nabla^r f[Q] = \sum_{i_1, \cdots, i_r} \partial_{i_1} \cdots \partial_{i_r} f(x) q_{1,i_1} \cdots q_{r,i_r}.$$

and we define the  $r^{th}$  normalized derivative of f as

$$\mathcal{D}_r[f](x)[Q] \stackrel{\text{def.}}{=} \nabla^r f(x)[\{\mathbf{H}_x^{-\frac{1}{2}}q_i\}_{i=1}^r]$$

with norm  $\|\mathbf{D}_r[f](x)\| \stackrel{\text{def.}}{=} \sup_{\forall \ell, \|q_\ell\| \leq 1} |\mathbf{D}_r[f](x)[Q]|$ . We will sometimes make use the the multiarray interpretation:  $\mathbf{D}_0[f] = f$ ,  $\mathbf{D}_1[f](x) = \mathbf{H}_x^{-\frac{1}{2}} \nabla f(x) \in \mathbb{C}^d$ ,  $\mathbf{D}_2[f](x) = \mathbf{H}_x^{-\frac{1}{2}} \nabla^2 f(x) \mathbf{H}_x^{-\frac{1}{2}} \in \mathbb{C}^{d \times d}$ .

For a bivariate function  $K : \mathcal{X} \times \mathcal{X} \to \mathbb{C}$ ,  $\partial_{1,i}$  (resp.  $\partial_{2,i}$ ) designates the derivative with respect to the *i*<sup>th</sup> coordinate of the first variable (resp. second variable), and similarly  $\nabla_i$  and  $\nabla_i^2$  denote the gradient and Hessian on the *i*<sup>th</sup> coordinate respectively.

For  $i, j \in \{0, 1, 2\}$ , let  $K^{(ij)}(x, x')$  be a "bi"-multilinear map, defined for  $Q \in (\mathbb{C}^d)^i$  and  $V \in (\mathbb{C}^d)^j$  as

$$[Q]K^{(ij)}(x,x')[V] \stackrel{\text{def.}}{=} \mathbb{E}[\overline{\mathcal{D}_i[\varphi_\omega](x)[Q]}\mathcal{D}_j[\varphi_\omega](x')[V]]$$

and  $\|K^{(ij)}(x,x')\| \stackrel{\text{def.}}{=} \sup_{Q,V} \|[Q]K^{(ij)}(x,x')[V]\|$  where the supremum is defined over all  $Q \stackrel{\text{def.}}{=} \{q_\ell\}_{\ell=1}^i, V \stackrel{\text{def.}}{=} \{v_\ell\}_{\ell=1}^j$  with  $\|q_\ell\| \leq 1, \|v_\ell\| \leq 1$ .

When  $i + j \leq 2$ , an equivalent definition is  $K^{(ij)}(x, x') = \mathbb{E}[\overline{\mathcal{D}_i[\varphi_\omega]}(x)\mathcal{D}_j[\varphi_\omega](x')^\top]$ , and we note that  $K^{(00)} = K$ , and we have normalized so that  $\operatorname{Re}\left(K^{(11)}(x, x)\right) = -\operatorname{Re}\left(K^{(02)}(x, x)\right)$ . Finally, we will make use of the still equivalent definition:  $[q]K^{(12)}(x, x') = \mathbb{E}[\overline{q^\top \mathcal{D}_1}[\varphi_\omega](x)\mathcal{D}_2[\varphi_\omega](x')^\top] \in \mathbb{C}^{d \times d}$ .

 $\mathbf{Kernel\ constants}\quad \text{For for } i,j\in\{(0,0),(0,1)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ define } B_{ij} \stackrel{\text{def.}}{=} \sup_{x,x'\in\mathcal{X}} \left|K^{(ij)}(x,x')\right|, \text{ for } (i,j)\in\{(0,2),(1,2)\}, \text{ def.}$ 

$$B_{ij} \stackrel{\text{def.}}{=} \sup \left\{ \left\| K^{(ij)}(x,x') \right\| ; d_{\mathbf{H}}(x,x') \leqslant r_{\text{near}} \text{ or } d_{\mathbf{H}}(x,x') > \Delta/2 \right\}.$$

and define for i = 1, 2

$$B_{ii} \stackrel{\text{def.}}{=} \sup_{x \in \mathcal{X}} \left\| K^{(ii)}(x, x) \right\|$$

For convenience, we define

$$B_i \stackrel{\text{def.}}{=} B_{0i} + B_{1i} + 1, \quad B \stackrel{\text{def.}}{=} \sum_{\substack{i,j \in \{0,1,2\}\\i+j \leq 3}} B_{ij} + 1.$$
(A.1)

**Matrices and vectors** We will make use of the following vectors and matrices throughout: Given  $X \stackrel{\text{def.}}{=} \{x_j\}_{j=1}^s \in \mathcal{X}^s$  and  $a \in \mathbb{C}^s$  which are always clear from context, define the vector  $\gamma_X(\omega) \in \mathbb{C}^{s(d+1)}$  as

$$\gamma_X(\omega) \stackrel{\text{def.}}{=} \left( \left( \overline{\varphi_\omega(x_i)} \right)_{i=1}^s, \left( \overline{D_1[\varphi_\omega](x_i)}^\top \right)_{i=1}^s \right)^\top, \tag{A.2}$$

and

$$\Upsilon_X \stackrel{\text{def.}}{=} \mathbb{E}_{\omega}[\gamma(\omega)\gamma(\omega)^*] \in \mathbb{C}^{s(d+1)\times s(d+1)}$$
$$\mathbf{f}_X(x) \stackrel{\text{def.}}{=} \mathbb{E}_{\omega}[\gamma(\omega)\varphi_{\omega}(x)] \in \mathbb{C}^{s(d+1)}$$
$$\alpha \stackrel{\text{def.}}{=} \Upsilon_X^{-1}\mathbf{u}_s, \qquad \mathbf{u}_s = \begin{pmatrix} \operatorname{sign}(a) \\ 0_{sd} \end{pmatrix}.$$

Note that the diagonal of  $\Upsilon$  has only 1's. For  $\omega_1, \ldots, \omega_m$ , we denote their empirical versions as:

$$\hat{\Upsilon}_X \stackrel{\text{def.}}{=} \frac{1}{m} \sum_{k=1}^m \gamma(\omega_k) \gamma(\omega_k)^*,$$
$$\hat{\mathbf{f}}_X(x) \stackrel{\text{def.}}{=} \frac{1}{m} \sum_{k=1}^m \gamma(\omega_k) \varphi_{\omega_k}(x), \quad \hat{\alpha} \stackrel{\text{def.}}{=} \hat{\Upsilon}_X^{-1} \mathbf{u}_s$$

which will serve us to construct our certificate, using the properties of their respective limit version. We remark that  $\mathbf{G}_X^{-1/2}\Gamma_X^*\Gamma_X\mathbf{G}_X^{-1/2} = \hat{\Upsilon}_X$ , where  $\Gamma_X$  is defined in the main paper and

$$\mathbf{G}_X = \begin{pmatrix} \mathrm{Id}_s & & 0 \\ & \mathbf{H}_{x_1} & & \\ & & \ddots & \\ 0 & & & \mathbf{H}_{x_s} \end{pmatrix}$$

The vanishing derivative pre-certificate  $\hat{\eta}_X$  is  $\hat{\alpha}^{\top} \hat{\mathbf{f}}_X(\cdot)$  and the limit pre-certificate is  $\eta_X \stackrel{\text{def.}}{=} \alpha^{\top} \mathbf{f}_X(\cdot)$ . When the set of points X is clear from context, we will drop the subscript X and write instead  $\gamma$ ,  $\Upsilon$ ,  $\mathbf{f}$ ,  $\eta$ , and so on.

Metric induced distances Given  $X = (x_j)_{j=1}^s \in \mathcal{X}^s$  and  $X' = (x'_j)_{j=1}^s \in \mathcal{X}^s$ , denote  $d_{\mathbf{H}}(X, X') \stackrel{\text{def.}}{=} \sqrt{\sum_j d_{\mathbf{H}}(x_j, x'_j)^2}$ . Observe also that  $\mathbf{G}_X$  is positive definite for all X and induces a metric on  $\mathbb{R}^s \times \mathcal{X}^s$  so that given  $a, a' \in \mathbb{R}^s$  and  $X, X' \in \mathcal{X}^s$ ,

$$d_G((a, X), (a', X')) = \sqrt{\|a - a'\|_2^2} + d_{\mathbf{H}}(X, X')^2.$$

Stochastic gradient bounds For  $r \in \mathbb{N}$ ,

$$L_{r}(\omega) = \sup_{x \in \mathcal{X}} \left\| \mathbf{D}_{r} \left[ \varphi_{\omega} \right](x) \right\|,$$

and  $L_{ij}(\omega) \stackrel{\text{def.}}{=} \sqrt{L_i(\omega)^2 + L_j(\omega)^2}$ . For i = 0, 1, 2, 3, let  $F_i$  be such that

$$\mathbb{P}_{\omega}\left(L_{j}(\omega) > t\right) \leqslant F_{i}(t),$$

Throughout, for  $(\bar{L}_j)_{j=0}^3 \in \mathbb{R}^4_+$ , the event  $\bar{E}$  is defined as

$$\bar{E} \stackrel{\text{def.}}{=} \bigcap_{k=1}^{m} E_{\omega_k} \quad \text{where} \quad E_{\omega} \stackrel{\text{def.}}{=} \{L_j(\omega) \leq \bar{L}_j, \ \forall j = 0, 1, 2, 3\}.$$
(A.3)

## **B** Proof of Theorem 2

In this section, we consider the (limit) vanishing derivative pre-certificate

$$\eta(x) = \mathbf{u}^{\top} \Upsilon_X^{-1} \mathbf{f}_X(x).$$

Note that

$$D_{2}[\eta](x) = \sum_{i=1}^{s} \alpha_{1,i} K^{(02)}(x_{i}, x) + [\alpha_{2,i}] K^{(12)}(x_{i}, x)$$

where we have decomposed  $\alpha = [\alpha_{1,1}, \ldots, \alpha_{1,s}, \alpha_{2,1}, \ldots, \alpha_{2,s}] \in \mathbb{C}^{s(d+1)}$  where  $\alpha_{2,i} \in \mathbb{C}^d$ .

We aim to prove that  $\eta$  is nondegenerate if K is an admissible kernel. Our first lemma shows that nondegeneracy of  $\eta$  within each small neighbourhood of  $x_i$  can be established by controlling the real and imaginary parts of  $D_2[\eta]$  in each small region:

**Lemma B.1.** Let  $\varepsilon > 0$ . Let  $a_0 \neq 0$ ,  $x_0 \in \mathcal{X}$  and let  $\sigma \in \mathbb{C}$  be such that  $|\sigma| = 1$ . Suppose that  $\eta \in \mathscr{C}^2(\mathcal{X}; \mathbb{C})$  is such that  $\eta(x_0) = \sigma$ ,  $\nabla \eta(x_0) = 0$  and  $\operatorname{Re}(\overline{\sigma}D_2[\eta](x_0)) \prec -\varepsilon \operatorname{Id}$ . Then,  $\nabla^2 |\eta|^2(x_0) \prec -2\varepsilon \operatorname{Id}$ . If in addition, we have c, r > 0 with  $\varepsilon r < 1$  and  $c^2 \leq (1 - \varepsilon r^2)/(\varepsilon r^2)$  such that for all x such that  $d_{\mathbf{H}}(x, x_0) \leq r$ ,

$$\operatorname{Re}\left(\overline{\sigma}\mathrm{D}_{2}\left[\eta\right](x)\right)\prec-\varepsilon\mathrm{Id}\quad and\quad \left\|\operatorname{Im}\left(\overline{\sigma}\mathrm{D}_{2}\left[\eta\right](x)\right)\right\|\leqslant c\varepsilon,$$

then,  $|\eta(x)|^2 \leq 1 - \varepsilon^2 d_{\mathbf{H}}(x, x_0)^2$  for all x such that  $d_{\mathbf{H}}(x, x_0) \leq r$ .

*Proof.* The first claim follows immediately from the computation: by writing  $\eta = \eta_r(x) + i\eta_i(x)$  where  $\eta_i$  and  $\eta_r$  are real valued functions,

$$\frac{1}{2}D_{2}\left[\left|\eta\right|^{2}\right] = \operatorname{Re}\left(\overline{D_{1}\left[\eta\right]}D_{1}\left[\eta\right]^{\top} + D_{2}\left[\eta\right]\overline{\eta}\right),$$

and evaluation at  $x_0$  gives the required result.

Let  $\gamma: [0,1] \to \mathcal{X}$  be a piecewise smooth path such that  $\gamma(0) = x_0, \gamma(1) = x$ .

$$\eta(x) = \eta(x_0) + \int_0^1 (1-t) \langle \nabla^2 \eta(\gamma(t)) \gamma'(t), \gamma'(t) \rangle dt$$
  
=  $\eta(x_0) + \int_0^1 (1-t) \langle D_2[\eta](\gamma(t)) \mathbf{H}_{\gamma(t)}^{\frac{1}{2}} \gamma'(t), \mathbf{H}_{\gamma(t)}^{\frac{1}{2}} \gamma'(t) \rangle dt.$ 

So,

$$\operatorname{Re}\left(\overline{\operatorname{sign}(a_0)}\eta(x)\right) = 1 + \inf_{\gamma} \operatorname{Re}\left(\overline{\operatorname{sign}(a_0)} \int_0^1 (1-t) \langle \mathcal{D}_2\left[\eta\right](\gamma(t)) \mathbf{H}_{\gamma(t)}^{\frac{1}{2}} \gamma'(t), \ \mathbf{H}_{\gamma(t)}^{\frac{1}{2}} \gamma'(t) \rangle \mathrm{d}t\right) \leqslant 1 - \varepsilon d_{\mathbf{H}}(x, x')^2$$

if we minimise over all paths from x to  $x_0$ . Similarly,

$$\left\|\operatorname{Im}\left(\overline{\operatorname{sign}(a_0)}\eta(x)\right)\right\| \leq c\varepsilon d_{\mathbf{H}}(x,x_0)^2$$

Therefore,

$$\begin{aligned} \left|\eta(x)\right|^2 &\leqslant \left|1 - \varepsilon d_{\mathbf{H}}(x, x_0)^2\right|^2 + \left|\varepsilon \varepsilon d_{\mathbf{H}}(x, x_0)^2\right|^2 \\ &\leqslant 1 - 2\varepsilon d_{\mathbf{H}}(x, x_0)^2 + \varepsilon^2 d_{\mathbf{H}}(x, x_0)^4 + c^2 \varepsilon^2 d_{\mathbf{H}}(x, x_0)^4 \\ &= 1 - \varepsilon d_{\mathbf{H}}(x, x_0)^2 - \varepsilon d_{\mathbf{H}}(x, x_0)^2 \left(1 - \varepsilon d_{\mathbf{H}}(x, x_0)^2 \left(1 + c^2\right)\right) \leqslant 1 - \varepsilon d_{\mathbf{H}}(x, x_0)^2. \end{aligned}$$

Proof of Theorem 2. In order to show that  $\eta$  is  $(\varepsilon_0/2, \varepsilon_2/2)$ -nondegenerate, it is enough to show that

$$\forall x \in \mathcal{X}^{\text{far}}, \quad |\eta(x)| \leq 1 - \varepsilon_0/2 \tag{B.1}$$

$$\forall x \in \mathcal{X}^{\text{near}}, \quad \text{Re}\left(\overline{\text{sign}(a_j)} D_2\left[\eta\right](x)\right) \prec -\frac{\varepsilon_2}{2} \text{Id} \quad \text{and} \quad \left\|\text{Im}\left(\overline{\text{sign}(a_j)} D_2\left[\eta\right](x)\right)\right\| \leqslant \frac{p}{4}\varepsilon_2 \tag{B.2}$$

where  $p = \sqrt{\frac{1 - \varepsilon_2 r_{\text{near}}^2/2}{\varepsilon_2 r_{\text{near}}^2/2}}$ .

We first prove that the matrix  $\Upsilon$  is invertible. To this end, we write

$$\Upsilon = \begin{pmatrix} \Upsilon_0 & \Upsilon_1^\top \\ \Upsilon_1 & \Upsilon_2 \end{pmatrix} \tag{B.3}$$

where  $\Upsilon_0 \stackrel{\text{def.}}{=} (K(x_i, x_j))_{i,j=1}^s \in \mathbb{C}^{s \times s}, \ \Upsilon_1 \stackrel{\text{def.}}{=} (K^{(10)}(x_i, x_j))_{i,j=1}^s \in \mathbb{C}^{sd \times s}, \text{ and } \Upsilon_2 \stackrel{\text{def.}}{=} (K^{(11)}(x_i, x_j))_{i,j=1}^s \in \mathbb{C}^{sd \times sd}.$  By definition of  $K^{(ij)}$ ,  $\Upsilon$  (and also  $\Upsilon_0$  and  $\Upsilon_2$ ) has only 1's on its diagonal.

To prove the invertibility of  $\Upsilon$ , we use the Schur complement of  $\Upsilon$ , and in particular it suffices to prove that  $\Upsilon_2$  and the Schur complement  $\Upsilon_S \stackrel{\text{def.}}{=} \Upsilon_0 - \Upsilon_1 \Upsilon_2^{-1} \Upsilon_1^{\top}$  are both invertible. To show that  $\Upsilon_2$  is invertible, we define  $A_{ij} = K^{(11)}(x_i, x_j)$ . So  $\Upsilon_2$  has the form:

$$\Upsilon_2 = \begin{pmatrix} \mathrm{Id} & A_{12} & \dots & A_{1s} \\ A_{21} & \mathrm{Id} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ A_{s1} & \dots & \dots & \mathrm{Id} \end{pmatrix}$$

and by Lemma G.6, we have

$$\|\mathrm{Id} - \Upsilon_2\|_{\mathrm{block}} \leqslant \max_i \sum_j \|A_{ij}\| \leqslant 1/4.$$

Since  $\|\mathrm{Id} - \Upsilon_2\|_{\mathrm{block}} < 1$ ,  $\Upsilon_2$  is invertible, and we have  $\|\Upsilon_2^{-1}\|_{\mathrm{block}} \leq \frac{1}{1 - \|I - \Upsilon_2\|_{\mathrm{block}}} \leq \frac{4}{3}$ . Next, again with Lemma G.6, we can bound

$$\|I - \Upsilon_0\|_{\infty} = \max_i \sum_{j \neq i} |K(x_i, x_j)| \leq \frac{\varepsilon_0}{16}$$
$$\|\Upsilon_1\|_{\infty \to \text{block}} \leq \max_i \sum_j \left\|K^{(10)}(x_i, x_j)\right\| \leq h \quad \text{since } K^{(10)}(x, x) = 0$$
$$\Upsilon_1^\top\|_{\text{block} \to \infty} \leq \max_i \sum_j \left\|K^{(10)}(x_j, x_i)\right\| \leq h$$

Hence, we have

$$\|I - \Upsilon_S\|_{\infty} \leqslant \|I - \Upsilon_0\|_{\infty} + \|\Upsilon_1^{\top}\|_{\text{block}\to\infty} \|\Upsilon_2^{-1}\|_{\text{block}} \|\Upsilon_1\|_{\infty\to\text{block}} \leqslant \frac{\varepsilon_0}{16} + \frac{4}{3}h^2 \leqslant \frac{\varepsilon_0}{8}$$
(B.4)

since  $h \leq \frac{\varepsilon_0}{32}$ . Therefore the Schur complement of  $\Upsilon$  is invertible and so is  $\Upsilon$ .

**Expression of**  $\eta$ . By definition,  $\eta = \text{satisfies } \eta(x_i) = \text{sign}(a_i)$  and  $\nabla \eta(x_i) = 0$ . We divide:

$$\alpha = \Upsilon^{-1} \mathbf{u}_s = \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix}$$

where  $\alpha_1 \in \mathbb{C}^s$  and  $\alpha_2 \in \mathbb{C}^{sd}$ , and we denote  $\alpha_{2,i} \in \mathbb{C}^d$  blocks such that  $\alpha_2 = [\alpha_{2,1}, \ldots, \alpha_{2,s}]$ .

The Schur's complement of  $\Upsilon$  allows us to express  $\alpha_1$  and  $\alpha_2$  as

$$\alpha_1 = \Upsilon_S^{-1} \operatorname{sign}(a), \qquad \alpha_2 = -\Upsilon_2^{-1} \Upsilon_1 \Upsilon_S^{-1} \operatorname{sign}(a)$$
(B.5)

and therefore we can bound

$$\|\alpha_1\|_{\infty} \leqslant \frac{1}{1 - \varepsilon_0/8} \tag{B.6}$$

$$\|\alpha_2\|_{\text{block}} \leqslant \frac{8}{3}h \leqslant 4h \tag{B.7}$$

Moreover, we have

$$\|\alpha_1 - \operatorname{sign}(a)\|_{\infty} \leq \|I - \Upsilon_S^{-1}\|_{\infty} \leq \|\Upsilon_S^{-1}\|_{\infty} \|I - \Upsilon_S\|_{\infty} \leq \frac{1}{4}$$
(B.8)

**Non-degeneracy.** We can now prove that  $\eta$  is non-degenerate.

Let x be such that  $d_{\mathbf{H}}(x_i, x) \leq r_{\text{near}}$ . We need to prove that for all x such that  $d_{\mathbf{H}}(x, x_i) \leq r$ ,

$$\operatorname{Re}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\eta\right](x)\right) \prec -\frac{\varepsilon_{2}}{2}\operatorname{Id} \quad \text{and} \quad \left\|\operatorname{Im}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\eta\right](x)\right)\right\| \leq \frac{\varepsilon_{2}}{2}\sqrt{\frac{2-\varepsilon r_{\operatorname{near}}^{2}}{\varepsilon_{2}r_{\operatorname{near}}^{2}}},$$

Then, since  $r_{\text{near}} \leq \Delta/2$  and the  $x_i$ 's are  $\Delta$ -separated, for all  $j \neq i$  we have  $d_{\mathbf{H}}(x, x_j) \geq \Delta/2$ . Then, we have

$$\overline{\operatorname{sign}(a_i)} \mathcal{D}_2[\eta](x) = \overline{\operatorname{sign}(a_i)} \left[ \alpha_{1,i} K^{(02)}(x_i, x) + \sum_{j \neq i} \alpha_{1,j} K^{(02)}(x_j, x) + [\alpha_{2,i}] K^{(12)}(x_i, x) + \sum_{j \neq i} [\alpha_{2,j}] K^{(12)}(x_j, x) \right]$$

$$\operatorname{Re}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\eta\right](x)\right) \preccurlyeq (1 - \|\alpha_{1} - \operatorname{sign}(a)\|_{\infty})\operatorname{Re}\left(K^{(02)}(x_{i}, x)\right) + \|\alpha_{1}\|_{\infty}\sum_{j\neq i}\left\|K^{(02)}(x_{j}, x)\right\| \operatorname{Id} \\ + \left(\left\|K^{(12)}(x_{i}, x)\right\| + \sum_{j\neq i}\left\|K^{(12)}(x_{j}, x)\right\|\right) \|\alpha_{2}\|_{\operatorname{block}}\operatorname{Id} \\ \preccurlyeq \left(-\frac{3}{4}\varepsilon_{2} + \frac{1}{1 - \varepsilon_{0}/8}\frac{\varepsilon_{2}}{16} + 4h(B_{12} + 1)\right)\operatorname{Id} \preccurlyeq \varepsilon_{2}\left(-\frac{3}{4} + \frac{1}{4}\right)\operatorname{Id} \preccurlyeq -\frac{\varepsilon_{2}}{2}\operatorname{Id}$$

Taking the imaginary part, we have

$$\begin{split} \left\| \operatorname{Im}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\eta\right](x)\right) \right\| &\leqslant (1 + \|\alpha_{1} - \operatorname{sign}(a)\|) \left\| \operatorname{Im}\left(K^{(02)}(x_{i}, x)\right) \right\| + \|\alpha_{1}\|_{\infty} \sum_{j \neq i} \left\| K^{(02)}(x_{j}, x) \right\| \\ &+ \left( \left\| K^{(12)}(x_{i}, x) \right\| + \sum_{j \neq i} \left\| K^{(12)}(x_{j}, x) \right\| \right) \|\alpha_{2}\|_{\operatorname{block}} \\ &\leqslant \left( \frac{5c\varepsilon_{2}}{4} + \frac{1}{(1 - \varepsilon_{0}/8)}h + 4h(B_{12} + 1) \right) \leqslant \frac{5c\varepsilon_{2}}{4} + h\left(4B_{12} + 6\right) \leqslant \frac{\varepsilon_{2}}{2} \sqrt{\frac{2 - \varepsilon r_{\operatorname{near}}^{2}}{\varepsilon_{2}r_{\operatorname{near}}^{2}}}. \end{split}$$

So, by Lemma B.1, for each i = 1, ..., s,  $|\eta(x)| \leq 1 - \varepsilon_2/2d_{\mathbf{H}}(x, x_i)$  for all  $x \in \mathcal{X}$  such that  $d_{\mathbf{H}}(x, x_i) \leq r_{\text{near}}$ . Next, for any x such that  $d_{\mathbf{H}}(x, x_i) \geq r_{\text{near}}$  for all  $x_i$ 's, we can say that there exists (at most) one index i such that  $d_{\mathbf{H}}(x, x_i) \ge r_{\text{near}}$  and for all  $j \ne i$  we have  $d_{\mathbf{H}}(x, x_j) \ge \Delta/2$ . We have

$$\begin{aligned} |\eta(x)| &= \left| \alpha_{1,i} K(x_i, x) + \sum_{j \neq i} \alpha_{1,j} K(x_j, x) \right. \\ &+ K^{(10)}(x_i, x)^\top \alpha_{2,i} + \sum_{j \neq i} K^{(10)}(x_j, x)^\top \alpha_{2,j} \right| \\ &\leqslant \|\alpha_1\|_{\infty} \left( |K(x_i, x)| + \sum_{j \neq i} |K(x_j, x)| \right) \\ &+ \|\alpha_2\|_{\text{block}} \left( \left\| K^{(10)}(x_i, x) \right\| + \sum_{j \neq i} \left\| K^{(10)}(x_j, x) \right\| \right) \\ &\leqslant \frac{1 - \varepsilon_0 + \varepsilon_0 / 16}{1 - \varepsilon_0 / 8} + 4h(B_{10} + 1) \leqslant 1 - \frac{\varepsilon_0}{2} \,. \end{aligned}$$

**Remark B.1.** Assuming that the derivatives of the kernel decay like a function f(||x - x'||) when, there is always a separation  $\Delta \propto f^{-1}(1/(Cs_{\max})))$  such that the kernel is admissible. Ex: when  $f = x^{-p}$ , we have  $\Delta \propto s_{\max}^{1/p}$  (eg Cauchy). When  $f = e^{-x^p}$ , we have  $\Delta \propto \log^{1/p}(s_{\max})$  (eg Gaussian).

### C Preliminaries

In this section, we present some preliminary results which will be used for proving our main results. We assume that K is admissible, and given a set of points  $X \in \mathcal{X}^s$ , let  $\mathcal{X}_j^{\text{near}} \stackrel{\text{def.}}{=} \{x \in \mathcal{X} ; d_{\mathbf{H}}(x, x_j) \leq r_{\text{near}}\}, \mathcal{X}^{\text{near}} \stackrel{\text{def.}}{=} \bigcup_{j=1}^s \mathcal{X}_j^{\text{near}}$  and  $\mathcal{X}^{\text{far}} \stackrel{\text{def.}}{=} \mathcal{X} \setminus \mathcal{X}^{\text{near}}$ .

#### C.1 On the determistic kernel

For an admissible kernel, we have the following additional bounds that will be handy.

**Lemma C.1.** Assume K is an admissible kernel, let  $X \in \mathcal{X}^s$  be  $\Delta$ -separated points. Then we have the following:

(i) We have seen that  $\Upsilon$  is invertible. Additionally it satisfies

$$\|\operatorname{Id} - \Upsilon\| \leq \frac{1}{2} \quad and \quad \|\operatorname{Id} - \Upsilon\|_{*,\infty} \leq \frac{1}{2}.$$
 (C.1)

(ii) For any vector  $q \in \mathbb{C}^{s(d+1)}$  and any  $x \in \mathcal{X}^{\text{far}}$ , we have

$$\|\mathbf{f}(x)\| \leqslant B_0 \quad and \quad \left|q^{\top} \mathbf{f}(x)\right| \leqslant B_0 \|q\|_{*,\infty} \tag{C.2}$$

(iii) For any vector  $q \in \mathbb{C}^{s(d+1)}$  and any  $x \in \mathcal{X}^{\text{near}}$  we have the bound:

$$\left\| \mathbf{D}_{2} \left[ q^{\mathsf{T}} \mathbf{f}(.) \right] (x) \right\| \leq \left\| q \right\| B_{2} \quad and \quad \left\| \mathbf{D}_{2} \left[ q^{\mathsf{T}} \mathbf{f}(.) \right] (x) \right\| \leq \left\| q \right\|_{*,\infty} B_{2}$$
(C.3)

*Proof.* We bound the spectral norm of Id –  $\Upsilon$ . Define  $y \in \mathbb{C}^{s(d+1)}$  decomposed as  $y = [y_1, \ldots, y_s, Y_1, \ldots, Y_s]$ 

where  $Y_i \in \mathbb{R}^d$ , such that  $||y|| \leq 1$ . We have

$$\begin{split} \|(\mathrm{Id}-\Upsilon)y\|^{2} &= \sum_{i=1}^{s} \left| \sum_{j\neq i} K(x_{i},x_{j})y_{j} + \sum_{j=1}^{s} K^{(10)}(x_{i},x_{j})^{\top}Y_{j} \right|^{2} \\ &+ \left\| \sum_{j} y_{j}K^{(10)}(x_{i},x_{j}) + \sum_{j\neq i} K^{(11)}(x_{i},x_{j})Y_{j} \right\|^{2} \\ &\leq \sum_{i=1}^{s} \left( \sum_{j\neq i} |K(x_{i},x_{j})| |y_{j}| + \sum_{j=1}^{s} \left\| K^{(10)}(x_{i},x_{j}) \right\| \|Y_{j}\| \right)^{2} \\ &+ \left( \sum_{j} |y_{j}| \left\| K^{(10)}(x_{i},x_{j}) \right\| + \sum_{j\neq i} \left\| K^{(11)}(x_{i},x_{j}) \right\| \|Y_{j}\| \right)^{2} \\ &\leq \max_{d_{\mathbf{H}}(x,x') \geq \Delta} \left( |K(x,x')|, \left\| K^{(10)}(x,x') \right\|, \left\| K^{(11)}(x,x') \right\| \right)^{2} \sum_{i} 2 \left( \sum_{j} |y_{j}| + \|Y_{j}\| \right)^{2} \\ &\leq 4s^{2} \max_{d_{\mathbf{H}}(x,x') \geq \Delta} \left( |K(x,x')|, \left\| K^{(10)}(x,x') \right\|, \left\| K^{(11)}(x,x') \right\| \right)^{2} \end{split}$$

by Cauchy-Schwartz inequality and since  $K^{(10)}(x, x) = 0$  for all  $x \in \mathcal{X}$ . Since by hypothesis we have

$$\max_{d_{\mathbf{H}}(x,x') \ge \Delta} \left( |K(x,x')|, \left\| K^{(10)}(x,x') \right\|, \left\| K^{(11)}(x,x') \right\| \right) \le \frac{1}{4s_{\max}},$$

we obtain

$$\|\mathrm{Id} - \Upsilon\| \leqslant \frac{1}{2} \tag{C.4}$$

and we deduce (i). A near identical argument also yields  $\|\Upsilon - \mathrm{Id}\|_{*,\infty} \leq \frac{1}{4}$ . For (ii), let  $x \in \mathcal{X}^{\mathrm{far}}$ , then we have

$$\|\mathbf{f}(x)\| \leq \left(\sum_{i=1}^{s} |K(x_i, x)|^2 + \left\|K^{(10)}(x_i, x)\right\|^2\right)^{\frac{1}{2}}$$
$$\leq \left(B_{00}^2 + \frac{(s-1)\varepsilon_0^2}{(16s_{\max})^2} + B_{10}^2 + \frac{(s-1)}{s_{\max}^2}\right)^{\frac{1}{2}} \leq B_0$$

for which, similar to the proof above, we have used the fact that x is  $\Delta/2$ -separated from at least s-1 points  $x_i$ . Similarly, for any vector  $q = [q_1, \ldots, q_s, Q_1, \ldots, Q_s] \in \mathbb{C}^{s(d+1)}$  and any  $x \in \mathcal{X}^{\text{far}}$ , we have

$$\|q^{\top}\mathbf{f}(x)\| \leq \sum_{i=1}^{s} |q_{i}| |K(x_{i}, x)| + \|Q_{i}\| \|K^{(10)}(x_{i}, x)\|$$
  
 
$$\leq \|q\|_{*,\infty} \left(B_{00} + \frac{(s-1)\varepsilon_{0}}{32s_{\max}} + B_{10} + \frac{(s-1)\varepsilon_{0}}{32s_{\max}}\right) \leq B_{0} \|q\|_{*,\infty} .$$

For any  $x \in \mathcal{X}^{\text{near}}$  we have the bound:

$$\begin{aligned} \left\| \mathcal{D}_{2} \left[ q^{\top} \mathbf{f} \right] (x) \right\| &= \left\| \sum_{i=1}^{s} q_{i} K^{(02)}(x_{i}, x) + [Q_{i}] K^{(12)}(x_{i}, x) \right\| \\ &\leq \left\| q \right\| \left( \sum_{i=1}^{s} \left\| K^{(02)}(x_{i}, x) \right\|^{2} + \left\| K^{(12)}(x_{i}, x) \right\|^{2} \right)^{\frac{1}{2}} \\ &\leq \left\| q \right\| B_{2} \end{aligned}$$

and

$$\| \mathbf{D}_{2} \left[ q^{\top} \mathbf{f} \right] (x) \| = \left\| \sum_{i=1}^{s} q_{i} K^{(02)}(x_{i}, x) + [Q_{i}] K^{(12)}(x_{i}, x) \right\|$$
$$\leq \| q \|_{*,\infty} \left( \sum_{i=1}^{s} \left\| K^{(02)}(x_{i}, x) \right\| + \left\| K^{(12)}(x_{i}, x) \right\| \right)$$
$$\leq \| q \|_{*,\infty} B_{2}$$

		٦

### C.2 Lipschitz bounds

**Lemma C.2** (Local Lipschitz constant of  $\varphi_{\omega}$  and higher order derivatives). Suppose that  $\|D_j[\varphi_{\omega}](x)\| \leq \bar{L}_j$  for all  $x \in \mathcal{X}$ . For all x, x' with  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$ , we have

- (i)  $|\varphi_{\omega}(x) \varphi_{\omega}(x')| \leq \mathcal{L}_0 d_{\mathbf{H}}(x, x'),$
- (*ii*)  $\| \mathbf{D}_1 [\varphi_\omega] (x) \mathbf{D}_1 [\varphi_\omega] (x') \| \leq \mathcal{L}_1 d_{\mathbf{H}}(x, x'),$
- (*iii*)  $\| \mathbf{D}_2 [\varphi_{\omega}] (x) \mathbf{D}_2 [\varphi_{\omega}] (x') \| \leq \mathcal{L}_2 d_{\mathbf{H}}(x, x'),$

where  $\mathcal{L}_0 \stackrel{\text{def.}}{=} \bar{L}_1$ ,  $\mathcal{L}_1 \stackrel{\text{def.}}{=} \bar{L}_1 C_{\mathbf{H}} + \bar{L}_2 (1 + C_{\mathbf{H}} r_{\text{near}})$  and  $\mathcal{L}_2 \stackrel{\text{def.}}{=} \bar{L}_2 (C_{\mathbf{H}} + C_{\mathbf{H}}^2 r_{\text{near}} + 1) + \bar{L}_3 (1 + C_{\mathbf{H}} r_{\text{near}})^2$ . As a consequence, for all  $X = (x_j)$  and  $X' = (x'_j)$  such that  $d_{\mathbf{H}}(x_j, x'_j) \leq r_{\text{near}}$ , we have

$$\sup_{\|q\|=1} \left\| \mathbf{D}_r \left[ q^\top (\hat{\mathbf{f}}_X - \hat{\mathbf{f}}_{X'}) \right] (y) \right\| \leqslant \bar{L}_r \sqrt{\mathcal{L}_0^2 + \mathcal{L}_1^2} d_{\mathbf{H}}(X, X').$$

*Proof.* Let  $x, x' \in \mathcal{X}$  with  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$ . Recall that  $\left\|\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}} - \text{Id}\right\| \leq C_{\mathbf{H}}d_{\mathbf{H}}(x, x')$ , and so,  $\left\|\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}}\right\| \leq 1 + C_{\mathbf{H}}r_{\text{near}}$ .

Let  $p: [0,1] \to \mathcal{X}$  be a piecewise smooth path such that p(0) = x', p(1) = x. Then, by Taylor's theorem,

$$\varphi_{\omega}(x) - \varphi_{\omega}(x') = \int_{t=0}^{1} \langle \mathbf{H}_{p(t)}^{-\frac{1}{2}} \nabla \varphi_{\omega}(p(t)), \, \mathbf{H}_{p(t)}^{\frac{1}{2}} p'(t) \rangle \mathrm{d}t \leqslant \bar{L}_{1} \int_{0}^{1} \left\| \mathbf{H}_{p(t)}^{\frac{1}{2}} p'(t) \right\| \mathrm{d}t \tag{C.5}$$

so taking the minimum over all paths p yields  $|\varphi_{\omega}(x) - \varphi_{\omega}(x')| \leq \bar{L}_1 d_{\mathbf{H}}(x, x')$ . Given  $q \in \mathbb{R}^d$ , by Taylor's theorem,

$$D_{1}[\varphi_{\omega}](x)[q] = \nabla\varphi(x)[\mathbf{H}_{x}^{-\frac{1}{2}}q] = \nabla\varphi(x')[\mathbf{H}_{x}^{-\frac{1}{2}}q] + \int \nabla^{2}\varphi_{\omega}(p(t))[\mathbf{H}_{x}^{-\frac{1}{2}}q, p'(t)]dt$$

$$= D_{1}[\varphi_{\omega}](x')[q] + D_{1}[\varphi_{\omega}](x')[(\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}} - \mathrm{Id})q] + \int D_{2}[\varphi_{\omega}](p(t))[\mathbf{H}_{p(t)}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}}q, \mathbf{H}_{p(t)}^{\frac{1}{2}}p'(t)]dt$$
(C.6)

Therefore,

$$\left\| \mathbf{D}_{1} \left[ \varphi_{\omega} \right](x) - \mathbf{D}_{1} \left[ \varphi_{\omega} \right](x') \right\| \leqslant \bar{L}_{1} C_{\mathbf{H}} d_{\mathbf{H}}(x, x') + \bar{L}_{2} (1 + C_{\mathbf{H}} r_{\mathrm{near}}) d_{\mathbf{H}}(x, x').$$

Finally, for all  $q_1, q_2 \in \mathbb{R}^d$ , by Taylor's theorem

$$D_{2} [\varphi_{\omega}] (x)[q_{1},q_{2}] - D_{2} [\varphi_{\omega}] (x')[q_{1},q_{2}] = \nabla^{2} \varphi_{\omega}(x) [\mathbf{H}_{x}^{-\frac{1}{2}}q_{1},\mathbf{H}_{x}^{-\frac{1}{2}}q_{2}] - \nabla^{2} \varphi_{\omega}(x') [\mathbf{H}_{x'}^{-\frac{1}{2}}q_{1},\mathbf{H}_{x'}^{-\frac{1}{2}}q_{2}] = D_{2} [\varphi_{\omega}] (x') [\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}}q_{1}, (\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}} - \mathrm{Id})q_{2}] + D_{2} [\varphi_{\omega}] (x') [(\mathbf{H}_{x'}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}} - \mathrm{Id})q_{1}, q_{2}] + \int D_{3} [\varphi_{\omega}] (p(t)) [\mathbf{H}_{p(t)}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}}q_{1}, \mathbf{H}_{p(t)}^{\frac{1}{2}}\mathbf{H}_{x}^{-\frac{1}{2}}q_{2}, \mathbf{H}_{p(t)}^{\frac{1}{2}}p'(t)] dt.$$
(C.7)

Therefore,

$$\|\mathbf{D}_{2}[\varphi_{\omega}](x) - \mathbf{D}_{2}[\varphi_{\omega}](x')\| \leq \left(\bar{L}_{2}\left((1 + C_{\mathbf{H}}r_{\mathrm{near}})C_{\mathbf{H}} + 1\right) + \bar{L}_{3}(1 + C_{\mathbf{H}}r_{\mathrm{near}})^{2}\right) d_{\mathbf{H}}(x, x').$$

By applying these Lipschitz bounds, we obtain

$$\sup_{\|q\|=1} \left\| D_r \left[ q^{\top}(\hat{\mathbf{f}}_X - \hat{\mathbf{f}}_{X'}) \right](y) \right\|^2$$
  
$$\leqslant \sum_{j=1}^s \left\| \hat{K}^{(0r)}(x_j, y) - \hat{K}^{(0r)}(x'_j, y) \right\|^2 + \sum_{j=1}^s \left\| \hat{K}^{(1r)}(x_j, y) - \hat{K}^{(1r)}(x'_j, y) \right\|^2$$
  
$$\leqslant \sum_{j=1}^s \mathcal{L}_0^2 \bar{L}_r^2 d_{\mathbf{H}}(x_j, x'_j)^2 + \sum_{j=1}^s \mathcal{L}_1^2 \bar{L}_r^2 d_{\mathbf{H}}(x_j, x'_j)^2$$
  
$$= \left( \mathcal{L}_0^2 + \mathcal{L}_1^2 \right) \bar{L}_r^2 d_{\mathbf{H}}(X, X')^2$$

**Lemma C.3** (Local Lipschitz constant of  $\hat{K}^{(ij)}$ ). Let  $x_1, x_0 \in \mathcal{X}$ . Let  $i, j \in \{0, 1, 2\}$  with  $i + j \leq 3$ . Define

$$A_{ij} = \sup_{x} \left\| \hat{K}^{(ij)}(x, x_0) \right\|$$

where x ranges over  $d_{\mathbf{H}}(x, x_1) \leq r_{\text{near}}$ . Then, for all x such that  $d_{\mathbf{H}}(x, x_1) \leq r_{\text{near}}$ ,

$$\begin{aligned} \left\| \hat{K}^{(0j)}(x,x_0) - \hat{K}^{(0j)}(x_1,x_0) \right\| &\leq A_{1j} d_{\mathbf{H}}(x,x_1) \\ \left\| \hat{K}^{(1j)}(x,x_0) - \hat{K}^{(1j)}(x_1,x_0) \right\| &\leq \left( C_{\mathbf{H}} A_{1j} + (1 + C_{\mathbf{H}} r_{\text{near}}) A_{2j} \right) d_{\mathbf{H}}(x,x_1) \end{aligned}$$

The same results hold if we replace  $\hat{K}$  by K.

*Proof.* The Lipschitz bounds on  $\hat{K}^{ij}$  follow by combining

$$[q_1, \dots, q_i](\hat{K}^{(ij)}(x, x_0) - \hat{K}^{(ij)}(x_1, x_0))[v_1, \dots, v_j]$$
  
=  $\hat{\mathbb{E}} \operatorname{Re} \left( \overline{(\operatorname{D}_i [\varphi_{\omega}](x) - \operatorname{D}_i [\varphi_{\omega}](x_1))[q_1, \dots, q_i]} \operatorname{D}_j [\varphi_j](x_0)[v_1, \dots, v_j] \right)$ 

where  $\hat{\mathbb{E}}$  indicates either empirical expectation or true expectation with (C.5), (C.6) and (C.7).

### C.3 Probability bounds

In the proof of our main results, we will often assume that event  $\overline{E}$  (see (A.3)) holds since our assumptions in Section 2.3 imply that  $\mathbb{P}(\overline{E}^c) \leq \rho/m$ . The following lemma shows that our assumptions also imply that  $\mathbb{E}_{\omega}[L_i(\omega)^2 \mathbb{1}_{E_{\omega}^c}] \leq \frac{\varepsilon}{m}$ . and this is a condition which our proofs will often rely upon. **Lemma C.4.** The following holds.  $\mathbb{P}(E_{\omega}^{c}) \leq \sum_{i} F_{i}(\bar{L}_{i})$  and

$$\mathbb{E}_{\omega}[L_j(\omega)^2 \mathbf{1}_{E_{\omega}^c}] \leqslant 2 \int_{\bar{L}_j}^{\infty} tF_j(t) \mathrm{d}t + \bar{L}_j^2 \sum_i F_i(\bar{L}_i)$$

*Proof.* Let  $E_{\omega,j}$  be the event that  $L_r(\omega) \leq \bar{L}_r$ , so  $E_\omega = \bigcap_{j=0}^3 E_{\omega,j}$ . By the union bound,  $\mathbb{P}(E_\omega^c) \leq \sum_j \mathbb{P}(E_{\omega,j}^c) \leq \sum_i F_i(\bar{L}_i)$ .

For the second claim, observe that  $E_{\omega}^c = \bigcup_i E_{\omega,i}^c$  so that  $\mathbb{E}[L_j(\omega)^2 \mathbf{1}_{E_{\omega}^c}] \leq \sum_i \mathbb{E}[L_j(\omega)^2 \mathbf{1}_{E_{\omega,i}^c}]$  and we have

$$\mathbb{E}[L_j(\omega)^2 \mathbf{1}_{E_{\omega,i}^c}] = \int_0^\infty \mathbb{P}(L_j(\omega)^2 \mathbf{1}_{E_{\omega,i}^c} \ge t) \mathrm{d}t$$
$$= \int_0^\infty \mathbb{P}\left((L_j(\omega)^2 \ge t) \cap (L_i(\omega) \ge \bar{L}_i)\right) \mathrm{d}t$$
$$\leqslant \bar{L}_j^2 F_i(\bar{L}_i) + \int_{\bar{L}_j^2}^\infty F_j(\sqrt{t}) \mathrm{d}t = \bar{L}_j^2 F_i(\bar{L}_i) + 2\int_{\bar{L}_j}^\infty t F_j(t) \mathrm{d}t$$

where we have bounded  $\mathbb{P}\left((L_j(\omega)^2 \ge t) \cap (L_i(\omega) \ge \bar{L}_i)\right)$  by respectively  $\mathbb{P}(L_i(\omega) \ge \bar{L}_i) \le F_i(\bar{L}_i)$  in the first term and by  $\mathbb{P}(L_j(\omega)^2 \ge t) \le F_j(\sqrt{t})$  in the second term.  $\Box$ 

### C.3.1 Concentration inequalities

The following result is an adaption of the Matrix Bernstein inequality for dealing with conditional probabilities. **Lemma C.5** (Adapted unbounded Matrix Bernstein). Let  $A_j \in \mathbb{R}^{d_1 \times d_2}$  be a family of iid matrices for  $j = 1, \ldots, m$ . Let  $Z = \frac{1}{m} \sum_{j=1}^{m} A_j$  and let  $\overline{Z} = \mathbb{E}[Z]$ . Let  $t \in (0, 4 ||\mathbb{E}[A_1]||]$ . Let events  $E_j$  be independent events such that  $E_j \subseteq \{||A_j|| \leq L\}$  and let  $E = \cap_j E_j$ . Suppose that we have

$$\mathbb{P}(E_j^c) \leqslant \frac{t}{t+4 \left\|\mathbb{E}[A_1]\right\|} \quad and \quad \mathbb{E}[\left\|A_j\right\| \mathbf{1}_{E_j^c}] \leqslant \frac{t}{4}$$

Then a first consequence is that we have  $\mathbb{E}_E[Z] = \mathbb{E}_{E_j}[A_j]$  for all j and  $\|\mathbb{E}[Z] - \mathbb{E}_E[Z]\| \leq \frac{t}{2}$ .

Finally, assuming that

$$\sigma^{2} \stackrel{\text{def.}}{=} \max_{j} \{ \left\| \mathbb{E}_{E_{j}}[A_{j}A_{j}^{*}] \right\|, \left\| \mathbb{E}_{E_{j}}[A_{j}^{*}A_{j}] \right\| \} < \infty$$

we have

$$\mathbb{P}_E\left(\|Z - \mathbb{E}[Z]\| \ge t\right) \le (d_1 + d_2) \exp\left(-\frac{mt^2/4}{\sigma^2 + Lt/3}\right)$$

*Proof.* We first bound  $||\mathbb{E}[Z] - \mathbb{E}_E[Z]||$ . First observe that  $\mathbb{E}[Z] = \mathbb{E}_{E_1}[A_1]$  and  $\mathbb{E}_E Z = \mathbb{E}_{E_1}[A_1]$  since  $A_j$  are iid. Moreover,

$$\mathbb{E}[A_1] = \mathbb{E}[A_1 1_{E_1}] + \mathbb{E}[A_1 1_{E_1^c}] = \mathbb{E}[A_1 | E_1] \mathbb{P}(E_1) + \mathbb{E}[A_1 1_{E_1^c}]$$

Hence,

$$\begin{aligned} \|\mathbb{E}[A_1] - \mathbb{E}_{E_1}[A_1]\| &= \left\| (P(E_1) - 1)\mathbb{E}_{E_1}[A_1] + \mathbb{E}[A_1 1_{E_1^c}] \right\| \\ &\leq \mathbb{P}(E_1^c) \|\mathbb{E}[A_1]\| + P(E_1^c) \|\mathbb{E}[A_1] - \mathbb{E}_{E_1}[A_1]\| + \mathbb{E}[\|A_1\| 1_{E_1^c}]. \end{aligned}$$

Therefore,

$$\|\mathbb{E}[A_1] - \mathbb{E}_{E_1}[A_1]\| \leqslant \frac{P(E_1^c) \|\mathbb{E}[A_1]\| + \mathbb{E}[\|A_1\| \, \mathbf{1}_{E_1^c}]}{1 - \mathbb{P}(E_1^c)} \leqslant \frac{t}{2}$$

For the second statement,

$$\mathbb{P}_E(\|Z - \mathbb{E}[Z]\| \ge t) \le \mathbb{P}_E(\|Z - \mathbb{E}_E[Z]\| \ge t - \|\mathbb{E}[Z] - \mathbb{E}_E[Z]\|)$$
$$\le \mathbb{P}_E(\|Z - \mathbb{E}_E[Z]\| \ge t/2).$$

To conclude, we apply Bernstein's inequality (Lemma G.2) to  $Y_j = A_j - \mathbb{E}[A_j|E] = Y_j = A_j - \mathbb{E}[A_j|E_j]$  conditional to E. Observe that

$$0 \leq \mathbb{E}_E[Y_j Y_j^{\top}] \leq \mathbb{E}_E[A_j A_j^{\top}] - \mathbb{E}_E[A_j] \mathbb{E}_E[A_j]^{\top}] \leq \mathbb{E}_E[A_j A_j^{\top}],$$

which yields  $\left\|\mathbb{E}_{E}[Y_{j}Y_{j}^{\top}]\right\| \leq \left\|\mathbb{E}[A_{j}A_{j}^{\top}]\right\|$  and similarly,  $\left\|\mathbb{E}_{E}[Y_{j}^{\top}Y_{j}]\right\| \leq \left\|\mathbb{E}_{E}[A_{j}^{\top}A_{j}]\right\|$ . So by Bernstein's inequality

$$\mathbb{P}_E(\|Z - \mathbb{E}_E[Z]\| \ge t/2) \le 2(d_1 + d_2) \exp\left(-\frac{mt^2/4}{\sigma^2 + Lt/3}\right).$$

Corollary C.1. Let  $x, x' \in \mathcal{X}$ . If

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t+4 \left\| K^{(ij)}(x,x') \right\|} \quad and \quad \mathbb{E}[L_{ij}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{t}{4}$$

then  $\left\| K_{\bar{E}}^{(ij)}(x,x') - K^{(ij)}(x,x') \right\| \leq t/2.$ 

**Proposition C.1.** Let t > 0 and assume that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t+6} \quad and \quad \mathbb{E}[L_{01}(\omega)^{2}1_{E_{\omega}^{c}}] \leqslant \frac{t}{4s}$$

then  $\|\Upsilon - \Upsilon_{\bar{E}}\| \leq t/2$  and

$$\mathbb{P}_{\bar{E}}(\left\|\Upsilon - \hat{\Upsilon}\right\| \ge t) \le 4(d+1)s \exp\left(-\frac{mt^2/4}{s\bar{L}_{01}^2(3+t/3)}\right)$$

Consequently,

$$\mathbb{P}_{\bar{E}}(\left\|\Upsilon^{-1} - \hat{\Upsilon}^{-1}\right\| \ge t) \le 4(d+1)s \exp\left(-\frac{mt^2}{16s\bar{L}_{01}^2(3+2\tilde{t})}\right)$$

*Proof.* We apply Lemma C.5 to  $A_j = \gamma(\omega_j)\gamma(\omega_j)^*$  with the following observations:

• for each  $\omega$ ,

$$\left\|\gamma(\omega)\gamma(\omega)^*\right\| \leq \left\|\gamma(\omega)\right\|^2 \leq s \max_{x \in \mathcal{X}} \{\left\|\mathbf{D}_1\left[\varphi_{\omega}\right](x)\right\|^2 + \left|\varphi_{\omega}(x)\right|^2\}$$

so under event  $\overline{E}$ ,  $||A_j|| \leq s\overline{L}_{01}^2$ .

- By Lemma C.1,  $\|\mathbb{E}[A_j]\| = \|\Upsilon\| \leq 3/2$ ,
- We may set  $\sigma^2 = \bar{L}_{01}(3/2 + t/2)$  since

$$0 \leq \mathbb{E}_{\bar{E}}[A_1A_1^*] = \mathbb{E}_{\bar{E}}[A_1^*A_1] = \mathbb{E}_{\bar{E}}[\|\gamma(\omega_j)\|^2 \gamma(\omega_j)\gamma(\omega_j)^*] \leq \bar{L}_{01}(\|\mathbb{E}[A_j]\| + t/2) \mathrm{Id}.$$

The last claim is because  $\left\| \Upsilon - \hat{\Upsilon} \right\| \leq t$  implies that  $\|\Upsilon\| \leq 3/2 + t$ ,  $\|\Upsilon^{-1}\| \leq \frac{\|\Upsilon\|}{1 - \|\Upsilon^{-1}\|} \leq \frac{3}{2 - 4t}$  and  $\left\| \Upsilon^{-1} - \hat{\Upsilon}^{-1} \right\| \leq \|\Upsilon^{-1}\| \left\| \Upsilon - \hat{\Upsilon} \right\| \left\| \hat{\Upsilon}^{-1} \right\| \leq \frac{3t}{1 - 2t}$  and writing  $\tilde{t} = \frac{3t}{1 - 2t}$  is equivalent to  $t = \tilde{t}/(3 + 2\tilde{t})$ .  $\Box$ 

## Bounds on $\hat{\mathbf{f}}_X$ applied to a fixed vector

**Proposition C.2.** Let  $t \in (0,1)$ ,  $r \in \{0,2\}$ ,  $q \in \mathbb{C}^{s(d+1)}$  and  $y \in \mathcal{X}_r$ , where  $\mathcal{X}_0 \stackrel{\text{def.}}{=} \mathcal{X}$  and  $\mathcal{X}_2 \stackrel{\text{def.}}{=} \mathcal{X}^{\text{near}}$ . If

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t+4B_{r}} \quad and \quad \mathbb{E}[L_{01}(\omega)L_{r}(\omega)\mathbf{1}_{E_{\omega}^{c}}] \leqslant \frac{t}{4\sqrt{s}}$$

then

$$\mathbb{P}_{\bar{E}}\left(\left\|\mathbf{D}_{r}\left[\left(\hat{\mathbf{f}}_{X_{0}}-\mathbf{f}_{X_{0}}\right)^{\top}q\right](y)\right\| \ge t \left\|q\right\|\right) \le 2\tilde{d}\exp\left(\frac{-mt^{2}/4}{2\bar{L}_{r}^{2}+\bar{L}_{r}\bar{L}_{01}t/(3\sqrt{s})}\right)$$

where  $\tilde{d} = 1$  if r = 0 and  $\tilde{d} = d$  if r = 2. As a consequence, since  $\sqrt{2s} \|q\|_{*,\infty} \ge \|q\|_2$ , we have

$$\mathbb{P}_{E}\left(\left\|\mathbf{D}_{r}\left[\left(\mathbf{f}_{X_{0}}-\hat{\mathbf{f}}_{X_{0}}\right)^{\top}q\right](y)\right\| \ge t \left\|q\right\|_{*,\infty}\right) \le 2\tilde{d}\exp\left(\frac{-mt^{2}}{16s(\bar{L}_{r}^{2}+8\bar{L}_{r}\bar{L}_{01}t/(3\sqrt{2}))}\right)$$

provided that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 4\sqrt{2s}B_{r}} \quad and \quad \mathbb{E}[L_{01}(\omega)L_{r}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{t}{4\sqrt{2s}}.$$

*Proof.* Without loss of generality, assume that ||q|| = 1. First note that

$$\mathbf{D}_{r}\left[\left(\hat{\mathbf{f}}_{X_{0}}-\mathbf{f}_{X_{0}}\right)^{\top}q\right](y)=\frac{1}{m}\sum_{k=1}^{m}q^{\top}\gamma(\omega_{k})\mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](y)-\mathbb{E}[q^{\top}\gamma(\omega_{k})\mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](y)].$$

We first consider the case of r = 0. We apply Lemma C.5 to  $A_k \stackrel{\text{def.}}{=} q^\top \gamma(\omega_k) \varphi_{\omega_k}(y) \in \mathbb{C}$ : Note that  $|A_k| \leq \sqrt{s} L_{01}(\omega_k) L_0(\omega_k)$  and  $|\mathbb{E}[A_k]| \leq B_0$ .

- Under event  $E_{\omega_k}$ ,  $||A_k|| \leq \bar{L}_2 \bar{L}_{01} \sqrt{s} \stackrel{\text{def.}}{=} L$ .
- $\mathbb{E}_{\bar{E}} |A_k|^2 = \mathbb{E}_{\bar{E}} [\langle \gamma(\omega_k) \gamma(\omega_k)^* q, q \rangle |\varphi_{\omega_k}(y)|^2] \leq \bar{L}_0^2 ||\Upsilon_{\bar{E}}|| \leq (3/2 + t/2) \bar{L}_0^2 \leq 2\bar{L}_0^2 \stackrel{\text{def.}}{=} \sigma^2.$

For the case r = 2, we apply Lemma C.5 with  $A_k \stackrel{\text{def.}}{=} q^\top \gamma(\omega_k) D_2[\varphi_{\omega_k}](y) \in \mathbb{C}^{d \times d}$ . Then,  $||A_k|| \leq \sqrt{s} L_{01}(\omega_k) L_2(\omega_k)$ ,  $||\mathbb{E}[A_k]|| \leq B_2$ , under event  $E_{\omega_k}$ ,  $||A_k|| \leq \overline{L}_2 \overline{L}_{01} \sqrt{s} \stackrel{\text{def.}}{=} L$  and

$$\left\|\mathbb{E}_{\bar{E}}[A_{k}A_{k}^{*}]\right\| = \left\|\mathbb{E}_{\bar{E}}[A_{k}^{*}A_{k}]\right\| = \left\|\mathbb{E}_{\bar{E}}[D_{2}\left[\varphi_{\omega_{k}}\right](y)^{*}D_{2}\left[\varphi_{\omega_{k}}\right](y)\left|q^{\top}\gamma(\omega_{k})\right|^{2}\right]\right\| \leqslant \bar{L}_{2}^{2}\mathbb{E}_{\bar{E}}[\left|q^{\top}\gamma(\omega_{k})\right|^{2}] \leqslant 2\bar{L}_{2}^{2} \stackrel{\text{def.}}{=} \sigma^{2}.$$

Lemma C.6. Assume that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 6\sqrt{2s}} \quad and \quad \mathbb{E}[L_{01}(\omega)^{2}1_{\bar{E}^{c}}] \leqslant \frac{t}{4\sqrt{2s^{3/2}}}$$
  
Let  $q \in \mathbb{C}^{s(d+1)}$ . Then, for all  $t \geqslant \frac{2\sqrt{2s}\bar{L}_{01}\bar{L}_{1}}{m} + \sqrt{\frac{8s^{2}\bar{L}_{01}^{2}\bar{L}_{1}^{2}}{m^{2}} + \frac{144s\bar{L}_{1}^{2}}{m}}$ , we have for each  $x_{i} \in X$ ,  
 $\mathbb{P}_{E}\left(\left\|D_{1}\left[q^{\top}(\mathbf{f}_{X} - \hat{\mathbf{f}}_{X})\right](x_{i})\right\|_{2} > 2t \left\|q\right\|_{*,\infty}\right) \leqslant 28 \exp\left(-\frac{mt^{2}/(4s)}{2\bar{L}_{1}^{2} + \sqrt{2t}\bar{L}_{1}\bar{L}_{01}/3}\right)$ 

*Proof.* For each  $x_i \in X$ ,

$$\left\| \mathbb{D}_1 \left[ \left( \mathbb{E}_{\bar{E}}[q^\top \hat{\mathbf{f}}_X] - q^\top \mathbf{f}_X \right) \right] (x_i) \right\| \leq \|\Upsilon - \Upsilon_{\bar{E}}\| \|q\| \leq \frac{t}{\sqrt{2s}} \|q\|,$$

by Proposition C.1. For convenience, we drop the subscript X from  $\mathbf{f}_X$ . Fix  $i \in \{1, \ldots, s\}$ . Observe that

$$\mathbb{P}_{E}\left(\left\|\mathbf{D}_{1}\left[\boldsymbol{q}^{\top}(\mathbf{f}-\hat{\mathbf{f}})\right](\boldsymbol{x}_{i})\right\|_{2} > 2t \left\|\boldsymbol{q}\right\|_{*,\infty}\right) \leq \mathbb{P}_{E}\left(\left\|\mathbf{D}_{1}\left[\boldsymbol{q}^{\top}(\mathbf{f}-\hat{\mathbf{f}})\right](\boldsymbol{x}_{i})\right\|_{2} > \frac{2t}{\sqrt{2s}} \left\|\boldsymbol{q}\right\|_{2}\right)$$
$$\leq \mathbb{P}_{E}\left(\left\|\mathbf{D}_{1}\left[\boldsymbol{q}^{\top}(\mathbb{E}_{\bar{E}}[\hat{\mathbf{f}}]-\hat{\mathbf{f}})\right](\boldsymbol{x}_{i})\right\|_{2} > \frac{t}{\sqrt{2s}} \left\|\boldsymbol{q}\right\|_{2}\right)$$

The claim of this lemma follows by applying Lemma G.3: Let

$$Y_{k} = \mathbf{D}_{1} \left[ \varphi_{\omega_{k}} \right] (x_{i}) \gamma(\omega_{k})^{\top} q - \mathbb{E}_{\bar{E}} \mathbf{D}_{1} \left[ \varphi_{\omega_{k}} \right] (x_{i}) \gamma(\omega)^{\top} q \in \mathbb{C}^{d},$$

and observe that  $D_1\left[q^{\top}(\hat{\mathbf{f}} - \mathbb{E}_{\bar{E}}[\hat{\mathbf{f}}])\right](x_i) = \frac{1}{m}\sum_k Y_k$ . Without loss of generality, assume that  $||q||_2 = 1$ . We apply Lemma G.3. Observe that conditional on event E,

- $||Y_k||_2 \leq 2 ||q||_2 ||\gamma(\omega_k)||_2 ||D_1[\varphi_{\omega_k}](x_i)||_2 \leq 2\sqrt{s}\bar{L}_{01}\bar{L}_1.$
- $\mathbb{E}_E \|Y_k\|^2 \leq \mathbb{E}_E[|\gamma(\omega_k)^\top q|^2 \mathbf{D}_1[\varphi_{\omega_k}](x_i)\mathbf{D}_1[\varphi_{\omega_k}](x_i)^\top] \leq \bar{L}_1^2 \|\Upsilon_E\|$ . So,  $\sigma^2 \leq m\bar{L}_1^2 \|\Upsilon_E\| \leq m\bar{L}_1^2(t+\|\Upsilon\|) \leq m\bar{L}_1^2(t/2+3/2) \leq 2m\bar{L}_1^2$  (here we are talking about the  $\sigma^2$  in Lemma G.3).

Therefore, for all

$$t \ge \frac{2\sqrt{2}s\bar{L}_{01}\bar{L}_{1}}{m} + \sqrt{\frac{8s^{2}\bar{L}_{01}^{2}\bar{L}_{1}^{2}}{m^{2}}} + \frac{144s\bar{L}_{1}^{2}}{m}$$
$$\mathbb{P}\left(\left\|\frac{1}{m}\sum_{k=1}^{m}Y_{k}\right\|_{2} \ge \frac{t}{\sqrt{2s}}\right) \le 28\exp\left(-\frac{mt^{2}/(4s)}{2\bar{L}_{1}^{2} + \sqrt{2t}\bar{L}_{1}\bar{L}_{01}/3}\right)$$

**Proposition C.3** (Block norm bound on  $\hat{\Upsilon}$  applied to a fixed vector). Suppose that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 6\sqrt{s}(B_{0} + 1)} \quad and \quad \mathbb{E}[L_{01}(\omega)^{2}1_{\bar{E}^{c}}] \leqslant \frac{t}{4s^{3/2}(1 + 4B_{0})}$$

Then, for all

$$t \ge \left(\frac{4\sqrt{2}s\bar{L}_{01}\bar{L}_1}{m} + \sqrt{\frac{32s^2\bar{L}_{01}^2\bar{L}_1^2}{m^2} + \frac{576s\bar{L}_1^2}{m}}\right)$$

we have

$$\mathbb{P}_{E}\left(\left\|(\Upsilon-\hat{\Upsilon})q\right\|_{*,\infty} \ge t \left\|q\right\|_{*,\infty}\right) \le 32s \exp\left(-\frac{mt^{2}}{s\left(32\bar{L}_{1}^{2}+34t\bar{L}_{1}\bar{L}_{01}\right)}\right).$$
(C.8)

*Proof.* Let  $S_0 \stackrel{\text{def.}}{=} \{1, \ldots, s\}$  and  $S_j \stackrel{\text{def.}}{=} \{s + (j-1)d + 1, \ldots, s + jd\}$  for  $j = 1, \ldots, s$ . Observe that by the union bound

$$\mathbb{P}_{E}\left(\left\|(\Upsilon-\hat{\Upsilon})q\right\|_{*,\infty} \ge t \|q\|_{*,\infty}\right)$$

$$\leq \mathbb{P}_{E}\left(\left\|((\Upsilon-\hat{\Upsilon})q)_{S_{0}}\right\|_{\infty} \ge t \|q\|_{*,\infty}\right) + \sum_{j=1}^{s} \mathbb{P}_{E}\left(\left\|((\Upsilon-\hat{\Upsilon})q)_{S_{j}}\right\|_{2} \ge t \|q\|_{*,\infty}\right)$$

$$\leq \sum_{j=1}^{s} \mathbb{P}_{E}\left(\left\|((\Upsilon-\hat{\Upsilon})q)_{j}\right\| \ge t \|q\|_{*,\infty}\right) + \sum_{j=1}^{s} \mathbb{P}_{E}\left(\left\|((\Upsilon-\hat{\Upsilon})q)_{S_{j}}\right\|_{2} \ge t \|q\|_{*,\infty}\right).$$
(C.9)

To bound the first sum, observe that  $((\Upsilon - \hat{\Upsilon})q)_j = (\mathbf{f}(x_j) - \hat{\mathbf{f}}(x_j))^\top q$  and  $((\Upsilon - \hat{\Upsilon})q)_{S_j} = D_1 \left[ q^\top (\mathbf{f} - \hat{\mathbf{f}}) \right] (x_j)$ . So, the first sum can be bounded by applying Proposition C.2. The second sum can be bounded by applying Lemma C.6.

Norm bounds for  $\hat{\mathbf{f}}$  We will repeatedly make use of the following result on  $\hat{\mathbf{f}}_X$ . This result is due to concentration bounds on the kernel  $\hat{K}$  which are derived subsequently.

**Proposition C.4** (Bound on  $\hat{\mathbf{f}}_X$ ). Let  $X \in \mathcal{X}^s$ . Let  $\rho > 0$ . Assume that for all  $(i, j) \in \{(0, 0), (1, 0), (0, 2), (1, 2)\}$ ,

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 4\sqrt{s} \max\{B_{0}, B_{2}\}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{t}{4\sqrt{s}}$$

Then, given any  $y \in \mathcal{X}$ ,

$$\mathbb{P}_{\bar{E}}\left(\left\|\hat{\mathbf{f}}_{X}(y) - \mathbf{f}_{X}(y)\right\| \ge t\right) \le 4sd \exp\left(-\frac{mt^{2}/8}{3s\bar{L}_{01}^{2}}\right).$$
(C.10)

and given any  $y \in \mathcal{X}^{\text{near}}$ , writing  $\hat{\mathbf{f}}_X = (\hat{f}_j)_{j=1}^p$  and  $\mathbf{f}_X = (f_j)_{j=1}^p$  with p = s(d+1), we have

$$\mathbb{P}_{\bar{E}}\left(\sup_{\|q\|=1}\sqrt{\sum_{j=1}^{p}\left\|D_{2}\left[\hat{f}_{j}-f_{j}\right](y)q\right\|^{2}} > t\right) \leqslant s(3d+d^{2})\exp\left(-\frac{mt^{2}/8}{s(\bar{L}_{2}^{2}B_{11}+\bar{L}_{1}^{2}B_{22}+\bar{L}_{01}\bar{L}_{2})}\right).$$
(C.11)

*Proof.* Let  $i, j \in \mathbb{N}_0$  with  $i + j \leq 2$ . Let  $[s] \stackrel{\text{def.}}{=} \{1, \ldots, s\}$  and  $I \stackrel{\text{def.}}{=} \{(0, 0), (1, 0)\}$ , By Lemma C.7 and the union bound,

$$\mathbb{P}_{\bar{E}}\left(\exists (i,j) \in I, \exists \ell \in [s], \left\|\hat{K}^{(ij)}(x_{\ell},y) - K^{(ij)}(x_{\ell},y)\right\| \ge \frac{t}{\sqrt{s}}\right) \leqslant 4sd \exp\left(-\frac{mt^2/4}{3s\bar{L}_{01}^2}\right).$$
(C.12)

So, (C.10) follows because

$$\left\|\hat{\mathbf{f}}_{X}(y) - \mathbf{f}_{X}(y)\right\| \leq \sqrt{\sum_{i=1}^{s} \left|\hat{K}(x_{i}, y) - K(x_{i}, y)\right|^{2} + \left\|\hat{K}^{(10)}(x_{i}, y) - K^{(10)}(x_{i}, y)\right\|^{2}} \leq \sqrt{2}t$$

By Lemma C.7, Lemma C.9 and the union bound, letting  $I_2 \stackrel{\text{def.}}{=} \{(0,2), (1,2)\}$ , we have

$$\mathbb{P}_{\bar{E}}\left(\exists (i,j) \in I_2, \exists \ell \in [s], \left\|\hat{K}^{(ij)}(x_\ell, y) - K^{(ij)}(x_\ell, y)\right\| \ge \frac{t}{\sqrt{s}}\right) \le 2sd \exp\left(-\frac{mt^2/4}{2s(\bar{L}_2^2 + \bar{L}_0\bar{L}_2)}\right) + s(d+d^2)\exp\left(-\frac{mt^2/4}{s(\bar{L}_2^2B_{11} + \bar{L}_1^2B_{22} + \bar{L}_1\bar{L}_2)}\right).$$
(C.13)

and (C.11) follows since given  $q \in \mathbb{C}^d$ , ||q|| = 1, we have

$$\sum_{j=1}^{p} \left\| \mathcal{D}_2\left[ \hat{f}_j - f_j \right](y) q \right\|^2 \leqslant \sum_{j=1}^{s} \left( \left\| \hat{K}^{(02)}(x_j, y) - K^{(02)}(x_j, y) \right\|^2 + \left\| \hat{K}^{(12)}(x_j, y) - K^{(12)}(x_j, y) \right\|^2 \right) \leqslant 2t^2$$

**Lemma C.7** (Concentration on kernel). Let t > 0,  $x, x' \in \mathcal{X}$ . Let  $i, j \in \mathbb{N}_0$  with  $i + j \leq 2$ . Assume

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t+4 \left\| K^{(ij)}(x,x') \right\|}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{t}{4}$$

then

$$\mathbb{P}_{\bar{E}}\left(\left\|\hat{K}^{(ij)}(x,x') - K^{(ij)}(x,x')\right\| \ge t\right) \le 2d \exp\left(-\frac{mt^2}{\bar{L}_p^2(b_{ij}+1) + \bar{L}_i\bar{L}_jt/3}\right)$$

where  $p = \max(i, j)$  and  $b_{ij} = 1$  if  $\min(i, j) = 0$  and  $b_{ij} \stackrel{\text{def.}}{=} ||K^{(11)}(x, x')||$  otherwise.

Proof. It is an immediate application of Lemma C.5 with  $A_k = \operatorname{Re}\left(\overline{\operatorname{D}_i\left[\varphi_{\omega_k}\right](x)}\operatorname{D}_j\left[\varphi_{\omega_k}\right](x')^{\top}\right)$  for  $k = 1, \ldots, m$ . Note that  $A_k \in (\mathbb{R}^d)^{i+j}$  if  $(i,j) \in \{(0,0), (0,1), (1,0)\}$  and  $A_k \in \mathbb{R}^{d \times d}$  if  $\max(i,j) = 2$ . noting that under  $\overline{E}$ ,  $||A_k|| \leq \overline{L}_i \overline{L}_j$ . Next, we need to bound  $||\mathbb{E}_{\overline{E}}[A_k A_k^*]||$  and  $||\mathbb{E}_{\overline{E}}[A_k^* A_k]||$ . We present only the argument for (i, j) = (0, 2), since all the other cases are similar:

$$0 \leq \mathbb{E}_{\bar{E}} A_k A_k^* \leq \mathbb{E}_{\bar{E}} [\|\varphi_{\omega_k}(x')\|^2 \operatorname{D}_2 [\varphi_{\omega_k}](x) \operatorname{D}_2 [\varphi_{\omega}](x)^*] \leq \bar{L}_2^2 \mathbb{E}_{\bar{E}} \|\varphi_{\omega_k}(x')\|^2 \operatorname{Id} = \bar{L}_2^2 |K_{\bar{E}}(x',x')| \operatorname{Id} \leq (1+t/2)\bar{L}_2^2 \operatorname{Id}$$

so  $\|\mathbb{E}_{\bar{E}}A_kA_k^*\| \leq (1+t/2)\bar{L}_2^2$ . Similarly,  $\|\mathbb{E}_{\bar{E}}A_k^*A_k\| \leq (1+t/2)\bar{L}_2^2$  and

$$\|\mathbb{E}_{\bar{E}}A_k^*A_k\|, \|\mathbb{E}_{\bar{E}}A_kA_k^*\| \leq L_p^2(B_{qq}+t/2)$$

where  $p = \max(i, j)$  and  $q = \min(i, j)$ .

Applying a grid on  $\mathcal{X}^{\text{near}}$ , we get a uniform version.

**Lemma C.8.** Let  $i, j \in \mathbb{N}_0$  with  $i + j \leq 2$ , and assume that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 16B_{ij}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{t}{16}.$$

Then

$$\mathbb{P}_{\bar{E}}\left(\exists x, x' \in \mathcal{X}^{\text{near}}, \left\|\hat{K}^{(ij)}(x, x') - K^{(ij)}(x, x')\right\| \ge t\right)$$
$$\leqslant 2ds^2 \exp\left(-\frac{mt^2/16}{L_p^2(B_{qq}+1) + \bar{L}_i\bar{L}_jt/12} + 2d\log\left(\frac{4(\mathcal{L}_i\bar{L}_j + \bar{L}_i\mathcal{L}_j)}{t}\right)\right)$$

where  $p = \max(i, j)$  and  $q = \min(i, j)$  and  $\mathcal{L}_i, \mathcal{L}_j$  are as in Lemma C.2

*Proof.* We define a  $\delta$ -covering of  $\mathcal{X}^{\text{near}}$  for the metric  $d_{\mathbf{H}}$  with  $\delta = \min\left(r_{\text{near}}, \frac{t}{4(\mathcal{L}_i \bar{L}_j + \bar{L}_i \mathcal{L}_j)}\right)$  of size  $s\left(\frac{r_{\text{near}}}{\delta}\right)^d$ . Let this covering be denoted by  $\mathcal{X}^{\text{grid}}$ .

By the union bound and Lemma C.7,

$$\mathbb{P}_{\bar{E}}\left(\exists x, x' \in \mathcal{X}^{\text{grid}} \text{ s.t. } \left\| \hat{K}^{(ij)}(x, x') - K^{(ij)}(x, x') \right\| \ge t/4 \right) \le 2ds^2 \left(\frac{r_{\text{near}}}{\delta}\right)^{2d} \exp\left(-\frac{mt^2/16}{L_p^2(B_{qq}+1) + \bar{L}_i\bar{L}_jt/12}\right)$$

where  $p = \max(i, j)$  and  $q = \min(i, j)$ . This gives the required upper bound: Given any  $x, x' \in \mathcal{X}$ , let  $x_{\text{grid}}, x'_{\text{grid}} \in \mathcal{X}^{\text{grid}}$  be such that  $d_{\mathbf{H}}(x, x_{\text{grid}}), d_{\mathbf{H}}(x', x'_{\text{grid}}) \leq \delta$ . Then, under event  $\overline{E}$ , by Lemma C.2,

$$\left\|\hat{K}^{(ij)}(x,x') - \hat{K}^{(ij)}(x_{\text{grid}},x'_{\text{grid}})\right\| \leq (\mathcal{L}_i\bar{L}_j + \bar{L}_i\mathcal{L}_j)\delta \leq t/4$$

By Jensen's inequality and since  $\left\|K_{\bar{E}}^{(ij)}(x,x') - K^{(ij)}(x,x')\right\| \leq t/4$  for all x, x', we have

$$\left\| K^{(ij)}(x,x') - K^{(ij)}(x_{\text{grid}},x'_{\text{grid}}) \right\| \leq t/2.$$

We now derive analogous results for the kernel differentiated 3 times.

**Lemma C.9** (Concentration on order 3 kernel). Let  $x, x' \in \mathcal{X}^{\text{near}}$ . Assume that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{t}{t + 4 \max\{B_{12}, B_{22}\}}, \quad \mathbb{E}[(L_{1}(\omega)L_{2}(\omega) + L_{2}^{2}(\omega))1_{E_{\omega}^{c}}] \leqslant \frac{t}{4}$$

r	_	_	٦
L			
L			
L	_	_	_

For j = 1, ..., m, let  $a_i = (D_1 [\overline{\varphi_{\omega_j}}](x))_i \in \mathbb{C}$ ,  $D \stackrel{\text{def.}}{=} D_2 [\varphi_{\omega}](x') \in \mathbb{C}^{d \times d}$  and  $A_j \stackrel{\text{def.}}{=} (a_1 D \quad a_2 D \quad \cdots \quad a_d D)^\top \in \mathbb{C}^{d^2 \times d}$ (C.14)

Let  $Z \stackrel{\text{def.}}{=} \frac{1}{m} \sum_{j=1}^{m} (A_j - \mathbb{E}[A_j])$ . Then, given

$$g(x') \stackrel{\text{def.}}{=} (g_i(x'))_{i=1}^d \stackrel{\text{def.}}{=} \sum_{k=1}^m \left( \overline{\mathcal{D}_1 \left[\varphi_{\omega_k}\right](x)} \varphi_{\omega}(x') - \mathbb{E}[\overline{\mathcal{D}_1 \left[\varphi_{\omega_k}\right](x)} \varphi_{\omega}(x')] \right)$$
$$= \hat{K}^{(10)}(x, x') - K^{(10)}(x, x'),$$

(*i*)  $\sup_{q \in \mathbb{C}^d, \|q\| \leq 1} \sum_{i=1}^d \|D_2[g_i](x')q\|^2 = \|Z\|^2$ ,

(*ii*)  $\sup_{q \in \mathbb{C}^d, \|q\| \leq 1} \left\| \mathcal{D}_2\left[q^\top g\right](x') \right\| = \left\| \hat{K}^{(12)}(x, x') - K^{(12)}(x, x') \right\| \leq \|Z\|.$ 

and

$$\mathbb{P}_{\bar{E}}\left(\|Z\| \ge t\right) \leqslant (d+d^2) \exp\left(-\frac{mt^2/4}{\tilde{B} + \bar{L}_1 \bar{L}_2 t/3}\right)$$

where  $\tilde{B} \stackrel{\text{def.}}{=} \max\{\bar{L}_2^2(B_{11}+t/2), \bar{L}_1^2(B_{22}+t/2)\}.$ 

*Proof.* The claim (i) is simply by definition, since  $Zq = (D_2[g_i](x')q)_{i=1}^d \in \mathbb{C}^{d^2}$ . For (ii), the first equality is simply be definition, and for the inequality, observe that

$$\sup_{q \in \mathbb{C}^{d}, \|q\| \leq 1} \left\| \mathbf{D}_{2} \left[ q^{\top} g \right] (x') \right\| = \sup_{q \in \mathbb{C}^{d}, \|q\| \leq 1} \sup_{p \in \mathbb{C}^{d}, \|p\| \leq 1} \left\| \sum_{i=1}^{d} q_{i} \mathbf{D}_{2} \left[ g_{i} \right] (x') p \right\|$$
$$\leq \sup_{q \in \mathbb{C}^{d}, \|q\| \leq 1} \sup_{p \in \mathbb{C}^{d}, \|p\| \leq 1} \left\| q \right\| \sqrt{\sum_{i=1}^{d} \left\| \mathbf{D}_{2} \left[ g_{i} \right] (x') p \right\|^{2}} \leq \left\| Z \right\|.$$

Finally, the probability bound follows by applying Lemma C.5: First note that under  $\bar{E}$ ,  $||A_j|| \leq \bar{L}_1 \bar{L}_2$ . It remains to bound  $||\mathbb{E}_{\bar{E}}[A_j^*A_j]||$  and  $||\mathbb{E}_{\bar{E}}[A_jA_j^*]||$ :

$$\sup_{\|q\| \leq 1} \mathbb{E}_{\bar{E}} \langle A_j^* A_j q, q \rangle = \sup_{\|q\| \leq 1} \mathbb{E}_E \sum_{i=1}^d \left| (\mathbf{D}_1 \left[ \varphi_{\omega_j} \right] (x))_i \right|^2 \left\| \mathbf{D}_2 \left[ \varphi_{\omega} \right] (x') q \right\|^2$$
$$\leq \sup_{\|q_k\| \leq 1} \bar{L}_1^2 \mathbb{E}_{\bar{E}} \overline{\mathbf{D}_2 \left[ \varphi_{\omega} \right] (x') [q_1, q_2]} \mathbf{D}_2 \left[ \varphi_{\omega} \right] (x') [q_3, q_4$$
$$\leq \bar{L}_1^2 \left\| K_{\bar{E}}^{(22)} (x, x) \right\| \leq \bar{L}_1^2 (B_{22} + t/2).$$

Given  $p_i \in \mathbb{C}^d$  for  $i = 1, \ldots, d$  such that  $\sum_i ||p_i||^2 \leq 1$ , write  $P = (p_1 \quad p_2 \quad \cdots \quad p_d) \in \mathbb{C}^{d \times d}$  and  $\bar{p} = (p_1^\top \quad p_2^\top \quad \cdots \quad p_d^\top)^\top \in \mathbb{C}^{d^2}$ . Then,

$$\mathbb{E}_{E} \langle A_{j} A_{j}^{*} \bar{p}, \bar{p} \rangle = \mathbb{E}_{E} \left\| \sum_{i=1}^{d} (\mathcal{D}_{1} \left[ \varphi_{\omega_{j}} \right] (x))_{i} \mathcal{D}_{2} \left[ \varphi_{\omega_{j}} \right] (x') p_{i} \right\|^{2}$$

$$= \mathbb{E}_{E} \left\| \mathcal{D}_{2} \left[ \varphi_{\omega_{j}} \right] (x') P \mathcal{D}_{1} \left[ \varphi_{\omega_{j}} \right] (x) \right\|^{2}$$

$$\leq \bar{L}_{2}^{2} \mathbb{E}_{E} \sum_{i} \left| \sum_{k} p_{i,k} (\mathcal{D}_{1} \left[ \varphi_{\omega_{j}} \right] (x))_{k} \right|^{2}$$

$$= \bar{L}_{2}^{2} \sum_{i} \langle \hat{K}_{\bar{E}}^{(11)}(x, x) p_{i}, p_{i} \rangle \leq \bar{L}_{2}^{2} \left\| \hat{K}_{\bar{E}}^{(11)}(x, x) \right\|^{2} \sum_{i} \left\| p_{i} \right\|^{2} \leq \bar{L}_{2}^{2} (B_{11} + t/2).$$

Lemma C.10 (Uniform concentration on order 3 kernel). Assume

$$\mathbb{P}(E_{\omega}^{c}) \leq \frac{t}{t + 16 \max\{B_{12}, B_{22}\}}, \quad \mathbb{E}[L_{1}(\omega)L_{2}(\omega)1_{E_{\omega}^{c}}] \leq \frac{t}{16}$$

then

$$\mathbb{P}_{\bar{E}}\left(\exists x, x' \in \mathcal{X}^{\text{near}}, \ \left\|\hat{K}^{(12)}(x, x') - K^{(12)}(x, x')\right\| \ge t\right) \\ \leqslant s^{2}(d+d^{2}) \exp\left(-\frac{mt^{2}/16}{\bar{B} + \bar{L}_{1}\bar{L}_{2}t/6} + 2d\log\left(\frac{8(\mathcal{L}_{1}\bar{L}_{2} + \bar{L}_{2}\mathcal{L}_{2})}{t}\right)\right)$$

where  $\tilde{B} \stackrel{\text{def.}}{=} \max\{\bar{L}_2^2(B_{11}+t/2), \bar{L}_1^2(B_{22}+t/2)\}, \mathcal{L}_1, \mathcal{L}_2 \text{ are as in Lemma C.2.}$ 

Proof. Let  $\mathcal{X}^{\text{grid}}$  be a  $\delta$ -covering of  $\mathcal{X}^{\text{near}}$  for the metric  $d_{\mathbf{H}}$  with  $\delta = \min\left(r_{\text{near}}, \frac{t}{8(\mathcal{L}_1\bar{L}_2 + \mathcal{L}_2\bar{L}_2)}\right)$  of size at most  $s\left(\frac{8(\mathcal{L}_1\bar{L}_2 + \mathcal{L}_2\bar{L}_2)}{t}\right)^d$ . By Lemma C.9 and the union bound,

$$\mathbb{P}_{\bar{E}}\left(\exists x, x' \in \mathcal{X}^{\text{grid}}, \left\|\hat{K}^{(ij)}(x, x') - K^{(ij)}(x, x')\right\| \ge t/2\right)$$
  
$$\leqslant s^{2}(d+d^{2}) \left(\frac{8(\bar{L}_{1}\bar{L}_{2} + \bar{L}_{2}^{2})}{t}\right)^{2d} \exp\left(-\frac{mt^{2}/16}{\bar{L}_{2}^{2}(B_{11} + t/4) + \bar{L}_{1}\bar{L}_{2}t/6}\right) \stackrel{\text{def.}}{=} \rho.$$

Moreover, under event  $\overline{E}$ , given any  $x, x' \in \mathcal{X}^{\text{near}}$ , there exists grid points  $x_{\text{grid}}, x'_{\text{grid}}$  such that

$$d_{\mathbf{H}}(x, x_{\text{grid}}), d_{\mathbf{H}}(x', x'_{\text{grid}}) \leqslant \delta$$

and

$$\begin{split} \left\| \left( \hat{K}^{(12)}(x,x') - K^{(12)}(x,x') \right) \right\| &\leq \left\| \left( \hat{K}^{(12)}(x_{\text{grid}},x'_{\text{grid}}) - K^{(12)}(x_{\text{grid}},x'_{\text{grid}}) \right) \right\| \\ &+ \left\| \left( \hat{K}^{(12)}(x,x') - \hat{K}^{(12)}(x_{\text{grid}},x'_{\text{grid}}) \right) \right\| \\ &+ \left\| \left( K^{(12)}(x,x') - K^{(12)}(x_{\text{grid}},x'_{\text{grid}}) \right) \right\|, \end{split}$$

and by Lemma C.2, under event  $\bar{E}$ ,

$$\left\| \left( \hat{K}^{(12)}(x,x') - \hat{K}^{(12)}(x_{\text{grid}},x'_{\text{grid}}) \right) \right\| \leq (\mathcal{L}_1 \bar{L}_2 + \mathcal{L}_2 \bar{L}_2) \delta \leq t/8.$$

and by Jensen's inequality and since  $\left\|K^{(12)}(x,y)-K^{(12)}_{\bar{E}}(x,y)\right\|\leqslant t/8,$ 

$$\left\| \left( K^{(12)}(x,y) - K^{(12)}(x_{\text{grid}},y) \right) \right\| \leq 3t/8.$$

Therefore, conditional on  $\bar{E}$ ,  $\left\| \left( \hat{K}^{(12)}(x,y) - K^{(12)}(x,y) \right) \right\| < t$  with probability at least  $1 - \rho$ .

## D Proof of Theorem 3

In all the rest of the proofs we fix  $X_0 \in \mathcal{X}^s$  to be  $\Delta$ -separated points,  $a_0 \in \mathbb{C}^s$ , and let  $\mathbf{u} = (\operatorname{sign}(a_0), 0_{sd})$ . We denote  $\mathcal{X}_i^{\operatorname{near}} = \{x \in \mathcal{X} ; d_{\mathbf{H}}(x, x_{0,i}) \leq r_{\operatorname{near}}\}$  and  $\mathcal{X}^{\operatorname{near}} = \cup_i \mathcal{X}_i^{\operatorname{near}}$  and  $\mathcal{X}^{\operatorname{far}} = \mathcal{X} \setminus \mathcal{X}^{\operatorname{near}}$ .

Since K is an admissible kernel, from (B.2) and (B.1) in the proof of Theorem 2  $\eta_{X_0}$  satisfies

- (i) for all  $y \in \mathcal{X}^{\text{far}}$ ,  $|\eta_{X_0}(y)| \leq 1 \frac{1}{2}\varepsilon_0$ ,
- (ii) for all  $y \in \mathcal{X}^{\text{near}}(i)$ ,  $-\text{Re}\left(\text{sign}(a_i)\text{D}_2\left[\eta_{X_0}\right](y)\right) \succeq \frac{1}{2}\varepsilon_2 \text{Id}$  and  $\|\text{Im}\left(\text{sign}(a_i)\text{D}_2\left[\eta_{X_0}\right](y)\right)\| \leqslant \left(\frac{p}{2}\right)\frac{1}{2}\varepsilon_2$ .

 $p \stackrel{\text{def.}}{=} \sqrt{(1 - \varepsilon_2 r_{\text{near}}^2/2)/(\varepsilon_2 r_{\text{near}}^2/2)} \ge 1,$ 

since  $\varepsilon_2 r_{\text{near}}^2 \leq 1$  by assumption of K being admissible. We aim to show that, for X close to  $X_0$ ,  $\hat{\eta}_X$  is nondegenerate by showing that  $\|D_r[\hat{\eta}_X] - D_r[\eta_{X_0}]\| \leq c\varepsilon_r$  for some positive constant c sufficiently small.

#### **D.1** Nondegeneracy of $\hat{\eta}_{X_0}$

We first establish the nondegeneracy of  $\hat{\eta}_{X_0}$ , our proof can be seen as a generalisation of the techniques in Tang et al. (2013) to the multidimensional setting with general sampling operators:

**Theorem D.1.** Let  $\rho > 0$  and assume that the assumptions in Section 2.3 hold. Assume also that either (a) or (b) holds:

(a)  $sign(a_0)$  is a Steinhaus sequence and

$$m \gtrsim C \cdot s \cdot \log\left(\frac{N^d}{\rho}\right) \log\left(\frac{s}{\rho}\right)$$

(b)  $sign(a_0)$  is an arbitrary sequence from the complex unit circle, and

$$m\gtrsim C\cdot s^{3/2}\cdot \log\left(\frac{N^d}{\rho}\right)$$

where C, N are defined in the main paper. Then with probability at least  $1 - \rho$ , the following hold: For all  $y \in \mathcal{X}^{far}$ ,  $|\hat{\eta}_{X_0}(y)| \leq 1 - \frac{7}{16}\varepsilon_0$ , and for all  $y \in \mathcal{X}^{near}(i)$ ,  $-\operatorname{Re}(\operatorname{sign}(a_i)\operatorname{D}_2[\hat{\eta}_{X_0}](y)) \approx \frac{7}{16}\varepsilon_2\operatorname{Id}$  and  $||\operatorname{Im}(\operatorname{sign}(a_i)\operatorname{D}_2[\hat{\eta}_{X_0}](y))|| \leq (\frac{p}{2} + \frac{p}{8})\frac{1}{2}\varepsilon_2$  and hence,  $\hat{\eta}_{X_0}$  is  $(\frac{7}{16}\varepsilon_0, \frac{7}{16}\varepsilon_2)$ -nondegenerate.

*Proof.* Note that

$$\frac{8}{7}\left(\frac{p}{2}+\frac{p}{8}\right) = \frac{5}{8}p < \sqrt{\frac{1-7\varepsilon_2 r_{\mathrm{near}}^2/16}{7\varepsilon_2 r_{\mathrm{near}}^2/16}}$$

so  $\hat{\eta}_{X_0}$  is  $(\frac{7}{16}\varepsilon_0, \frac{7}{16}\varepsilon_2)$ -nondegenerate by Lemma B.1

Let  $c \stackrel{\text{def.}}{=} 1/32$ . Observe that by assumption and Lemma C.4,  $\mathbb{P}(\bar{E}) \leq \rho/2$ . Therefore, it is sufficient to prove that conditional on  $\bar{E}$ , with probability at least  $1 - \delta$  with  $\delta \stackrel{\text{def.}}{=} \rho/2$ ,  $\hat{\eta}_{X_0}$  is nondegenerate.

We will repeatedly use the fact that our assumptions (by Lemma C.4) also imply that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{\varepsilon}{m}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{\varepsilon}{m}$$

for all  $(i, j) \in \{(0, 0), (1, 0), (0, 2), (1, 2)\},\$ 

### Step I: Proving nondegeneracy on a finite grid.

Let  $\mathcal{X}_{grid}^{far} \subset \mathcal{X}^{far}$  and  $\mathcal{X}_{grid}^{far} \subset \mathcal{X}^{near}$ , be finite point sets. Let

$$Q_r(y) \stackrel{\text{def.}}{=} \| \mathbf{D}_r [\hat{\eta}_{X_0}] (y) - \mathbf{D}_r [\eta_{X_0}] (y) \|, \qquad r = 0, 2.$$

We first prove that conditional on  $\overline{E}$ , with probability at least  $1 - \delta$  where  $\delta \stackrel{\text{def.}}{=} \rho/2$ , that  $Q_0(y) \leq c\varepsilon_0$  for all  $y \in \mathcal{X}_{\text{grid}}^{\text{far}}$  and  $Q_2(y) \leq c\varepsilon_2$  for all  $y \in \mathcal{X}_{\text{grid}}^{\text{far}}$ .

Let us first recall some facts which were proven in the previous section: Let  $a, t \in (0, 1)$  and write  $\mathbf{f} = (\bar{f}_j)_{j=1}^{s(d+1)}$ and  $\hat{\mathbf{f}} = (f_j)_{j=1}^{s(d+1)}$ . Let  $q_0 \stackrel{\text{def.}}{=} \Upsilon^{-1} \mathbf{u}$ , so  $||q_0|| \leq 2\sqrt{s}$ . Let F be the event that (a)  $\left\| \Upsilon^{-1} - \hat{\Upsilon}^{-1} \right\| \leq t$ , (b)  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left\| \hat{\mathbf{f}}_{X_0}(y) - \mathbf{f}_{X_0}(y) \right\| \leq a\varepsilon_0$ , (c)  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{near}}, \sup_{q \in \mathbb{C}^d, \|q\|=1} \sqrt{\sum_{j=1}^p \left\| D_2 \left[ f_j - \bar{f}_j \right](y) q \right\|^2} \leq a\varepsilon_2$ ,

Let G be the event that

(d)  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left| (\hat{\mathbf{f}}_{X_0}(y) - \mathbf{f}_{X_0}(y))^\top q_0 \right| \leq 2a\varepsilon_0$ (e)  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{near}}, \left\| \mathbf{D}_2 \left[ (\hat{\mathbf{f}}_{X_0} - \mathbf{f}_{X_0})^\top q_0 \right](y) \right\| \leq 2a\varepsilon_2$ 

then provided that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{u}{u + \max\{4\sqrt{s}B_{ij}, 6\}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{u}{4s}$$
(D.1)

where  $u = \min\{a\varepsilon_i, t\}$ , we have

$$\begin{split} \mathbb{P}_{\bar{E}}(F^{c}) \leqslant &4(d+1)s \exp\left(-\frac{mt^{2}}{16s\bar{L}_{01}^{2}(3+2t)}\right) \\ &+ 4sd \left|\mathcal{X}_{\text{grid}}^{\text{far}}\right| \exp\left(-\frac{m(a\varepsilon_{0})^{2}/8}{s(\bar{L}_{01}^{2}(B_{11}+1)+\bar{L}_{01}^{2})}\right) \\ &+ s(3d+d^{2}) \left|\mathcal{X}_{\text{grid}}^{\text{near}}\right| \exp\left(-\frac{m(a\varepsilon_{2})^{2}/8}{s(\bar{L}_{2}^{2}B_{11}+\bar{L}_{1}^{2}B_{22})+\bar{L}_{01}\bar{L}_{2})}\right) \end{split} \tag{D.2}$$
$$\\ \mathbb{P}_{\bar{E}}(G^{c}) \leqslant 2 \left|\mathcal{X}_{\text{grid}}^{\text{far}}\right| \exp\left(-\frac{ma^{2}\varepsilon_{0}^{2}}{s(8\bar{L}_{0}^{2}+\frac{4}{3}\bar{L}_{0}\bar{L}_{01}a\varepsilon_{0})}\right) \\ &+ 2d \left|\mathcal{X}_{\text{grid}}^{\text{near}}\right| \exp\left(-\frac{ma^{2}\varepsilon_{2}^{2}}{s(8\bar{L}_{2}^{2}+\frac{4}{3}\bar{L}_{2}\bar{L}_{01}a\varepsilon_{2})}\right), \end{split}$$

where for  $\mathbb{P}_{\bar{E}}(F^c)$ , the first term on the right is due to Proposition C.1, the second and third are due to Proposition C.4 while the bound on  $\mathbb{P}_{\bar{E}}(G^c)$  is due to Proposition C.2 (noting that, since this probability bound over the  $\omega_j$  is valid for all fixed **u**, and the  $\omega_j$  and the signs are independent, it is valid with the same probability over both  $\omega_j$  and **u**).

Observe that

$$\| \mathbf{D}_{j} [\hat{\eta}_{X_{0}}] (y) - \mathbf{D}_{j} [\eta_{X_{0}}] (y) \| = \left\| \mathbf{D}_{j} \left[ (\hat{\alpha}_{X_{0}} - \alpha_{X_{0}})^{\top} \hat{\mathbf{f}}_{X_{0}} \right] (y) + \mathbf{D}_{j} \left[ \alpha_{X_{0}}^{\top} (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) \right] (y) \right\|$$

$$\leq \left\| \mathbf{D}_{j} \left[ \mathbf{u}^{\top} \left( (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \hat{\mathbf{f}}_{X_{0}} + \Upsilon^{-1} (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) \right) \right] (y) \right\|$$
(D.3)

### Step I (a): Random signs

We first bound (D.3) in the case where **u** is a Steinhaus sequence.

Let  $\beta_1(y) \stackrel{\text{def.}}{=} (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \hat{\mathbf{f}}_{X_0}(y)$  and  $\beta_2(y) \stackrel{\text{def.}}{=} \Upsilon^{-1}(\hat{\mathbf{f}}_{X_0}(y) - \mathbf{f}_{X_0}(y))$ . Then, event F implies that  $\|\beta_1(y)\| \leq t(B_0 + a\varepsilon_0)$  for all  $y \in \mathcal{X}_{\text{grid}}^{\text{far}}$ , and event G implies that  $|\mathbf{u}^\top \beta_2(y)| \leq 2a\varepsilon_0$ . So,

$$\begin{aligned} \mathbb{P}_{\bar{E}}\left(\left|\exists y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \mathbf{u}^{\top}(\beta_{1}+\beta_{2})(y)\right| > c\varepsilon_{0}\right) \\ &\leqslant \mathbb{P}_{F\cap\bar{E}}\left(\exists y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left|\mathbf{u}^{\top}\beta_{1}(y)\right| > \frac{c}{2}\varepsilon_{0}\right)\mathbb{P}_{\bar{E}}(F) + \mathbb{P}_{\bar{E}}\left(F^{c}\right) \\ &+ \mathbb{P}_{G\cap\bar{E}}\left(\exists y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left|\mathbf{u}^{\top}\beta_{2}(y)\right| > \frac{c}{2}\varepsilon_{0}\right)\mathbb{P}_{\bar{E}}(G) + \mathbb{P}_{\bar{E}}\left(G^{c}\right) \\ &\leqslant \mathbb{P}_{F\cap\bar{E}}\left(\exists y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left|\mathbf{u}^{\top}\beta_{1}\right| > \frac{c}{2}\varepsilon_{0}\right) + \mathbb{P}_{\bar{E}}\left(F^{c}\right) + \mathbb{P}_{\bar{E}}\left(G^{c}\right) \\ &\leqslant 4\left|\mathcal{X}_{\text{grid}}^{\text{far}}\right| e^{-\frac{(c/4)^{2}\varepsilon_{0}^{2}}{8t^{2}(B_{0}+a\varepsilon_{0})^{2}}} + \mathbb{P}_{\bar{E}}(F^{c}) + \mathbb{P}_{\bar{E}}\left(G^{c}\right). \end{aligned}$$
(D.4)

where we set a = c/4 for the second inequality and the last inequality follows from Lemma G.4 and because **u** consists if random signs.

Now consider  $Q_2(y) = \mathcal{D}_2\left[\mathbf{u}^\top\beta\right](y)$ . Under event G,  $\left\|\mathcal{D}_2\left[\mathbf{u}^\top\beta_2\right](y)\right\| \leq \frac{c}{2}\varepsilon_2$ . Writing  $M = (\hat{\Upsilon}^{-1} - \Upsilon^{-1})$ , we have

$$D_{2}\left[\mathbf{u}^{\top}\beta_{1}\right](y) = D_{2}\left[\mathbf{u}^{\top}\left(M\hat{\mathbf{f}}_{X_{0}}\right)\right](y) = \sum_{\ell=1}^{p} \mathbf{u}_{\ell}\left(\sum_{j=1}^{p} M_{\ell j} D_{2}\left[f_{j}\right](y)\right).$$
(D.5)

We aim to bound (D.5) by applying the Matrix Hoeffding's inequality (Corollary G.1): let

$$Y_{\ell} \stackrel{\text{def.}}{=} \operatorname{Re}\left(\sum_{j=1}^{p} M_{\ell j} \mathcal{D}_{2}\left[f_{j}\right](y)\right) \in \mathbb{R}^{d \times d}$$

which is a symmetric matrix. Note that

$$\left\|\sum_{\ell=1}^{p} Y_{\ell}^{2}\right\| = \sup_{q \in \mathbb{R}^{d}, \|q\|=1} \sum_{\ell=1}^{p} \langle Y_{\ell}^{2}q, q \rangle = \sup_{q \in \mathbb{R}^{d}, \|q\|=1} \sum_{\ell=1}^{d} \|Y_{\ell}q\|^{2} \leqslant \sup_{q \in \mathbb{R}^{d}, \|q\|=1} \left\|\sum_{j=1}^{p} M_{\ell,j}(\mathcal{D}_{2}\left[f_{j}\right](y)q)\right\|^{2}.$$

Then, for a vector q of unit norm, let  $V_{j,n} \stackrel{\text{def.}}{=} (D_2[f_j](y)q)_n$  for  $j = 1, \ldots, p$  and  $n = 1, \ldots, d$ , then

$$\sum_{\ell=1}^{p} \left\| \sum_{j=1}^{p} M_{\ell,j}(\mathbf{D}_{2}[f_{j}](y)q) \right\|^{2} = \sum_{\ell=1}^{p} \sum_{n=1}^{d} \left| \sum_{j=1}^{p} M_{\ell,j} V_{j,n} \right|^{2} = \sum_{n=1}^{d} \|MV_{\cdot,n}\|^{2} \leq \|M\|^{2} \sum_{n=1}^{d} \|V_{\cdot,n}\|^{2}$$
$$= \|M\|^{2} \sum_{n=1}^{d} \sum_{j=1}^{p} |V_{j,n}|^{2} = \|M\|^{2} \sum_{j=1}^{p} \|\mathbf{D}_{2}[f_{j}](y)q\|^{2}.$$

Under event F, we have  $\|M\|^2 \sum_{j=1}^p \|D_2[f_j](y)q\|^2 \leq t^2 (B_2 + a\varepsilon_2)^2$ . Then,

$$\mathbb{P}_{F \cap \bar{E}}\left(\left\| \mathbf{D}_2\left[\mathbf{u}^\top \operatorname{Re}\left(M\hat{\mathbf{f}}_{X_0}\right)\right](y)\right\| \ge \frac{c\varepsilon_2}{\sqrt{2}}\right) \leqslant 2d \exp\left(-\frac{(c/2)^2\varepsilon_2^2}{4t^2(B_2 + a\varepsilon_2)^2}\right).$$

By repeating this argument for the imaginary part, we obtain

$$\mathbb{P}_{F \cap \bar{E}}\left(\left\| \mathcal{D}_2\left[\mathbf{u}^\top \operatorname{Im}\left(M\hat{\mathbf{f}}_{X_0}\right)\right](y)\right\| \ge \frac{c\varepsilon_2}{\sqrt{2}}\right) \le 2d \exp\left(-\frac{(c/2)^2\varepsilon_2^2}{4t^2(B_2 + a\varepsilon_2)^2}\right).$$

So,

$$\mathbb{P}_{\bar{E}}\left(\exists y \in \mathcal{X}_{\text{grid}}^{\text{near}}, \left\| \mathbf{D}_{2} \left[ \mathbf{u}^{\top} \beta(y) \right] \right\| > c\varepsilon_{2} \right) \\
\leqslant \mathbb{P}_{F \cap \bar{E}} \left( \exists y \in \mathcal{X}_{\text{grid}}^{\text{near}}, \left\| \mathbf{D}_{2} \left[ \mathbf{u}^{\top} \operatorname{Re} \left( M \hat{\mathbf{f}}_{X_{0}} \right) \right](y) \right\| \ge \frac{c}{2} \varepsilon_{2} \right) + \mathbb{P}_{\bar{E}}(F^{c}) + \mathbb{P}_{\bar{E}}(G^{c}) \\
\leqslant 4d \left| \mathcal{X}_{\text{grid}}^{\text{near}} \right| \exp \left( -\frac{(c/2)^{2} \varepsilon_{2}^{2}}{4t^{2} (B_{2} + a\varepsilon_{2})^{2}} \right) + \mathbb{P}_{\bar{E}}(F^{c}) + \mathbb{P}_{\bar{E}}(G^{c}).$$
(D.6)

Therefore,

$$1 - \mathbb{P}\left(Q_0(y_0) \leqslant c\varepsilon_0 \text{ and } Q_2(y_2) \leqslant c\varepsilon_2, \forall y_0 \in \mathcal{X}_{\text{grid}}^{\text{far}}, \forall y_2 \in \mathcal{X}_{\text{grid}}^{\text{near}}\right)$$
  
$$\leqslant 4 \left|\mathcal{X}_{\text{grid}}^{\text{far}}\right| \exp\left(-\frac{(c/2)^2 \varepsilon_0^2}{32t^2(B_0 + a\varepsilon_0)^2}\right) + 4d \left|\mathcal{X}_{\text{grid}}^{\text{near}}\right| \exp\left(-\frac{(c/2)^2 \varepsilon_2^2}{16t^2(B_2 + a\varepsilon_2)^2}\right) + 2\mathbb{P}_{\bar{E}}(F^c) + 2\mathbb{P}_{\bar{E}}(G^c).$$

The first 2 terms are each bounded by  $\delta/7$  by setting t such that

$$\frac{1}{t^2} = 2^{13} \log\left(\frac{112\bar{N}d}{\delta}\right) \frac{\left(\bar{B}+1\right)}{c^2 \varepsilon^2}$$

where  $\bar{B} \stackrel{\text{def.}}{=} \max\{B_0, B_2\}$ ,  $\varepsilon \stackrel{\text{def.}}{=} \min\{\varepsilon_0, \varepsilon_2\}$  and  $\bar{N} = \max(|\mathcal{X}_{\text{grid}}^{\text{near}}|, |\mathcal{X}_{\text{grid}}^{\text{far}}|)$ . The first term of (D.2) is bounded by  $\delta/7$  if

$$m \ge \frac{1}{t^2} \log\left(\frac{28(d+1)s}{\delta}\right) 64s\bar{L}_{01}^2 = s\bar{L}_{01}^2 \frac{2^{19}\left(\bar{B}+1\right)}{c^2\varepsilon^2} \log\left(\frac{112\bar{N}d}{\delta}\right) \log\left(\frac{28(d+1)s}{\delta}\right)$$

and the last 4 terms of (D.2) are each bounded by  $\delta/7$  provided that

$$m \gtrsim \log\left(\frac{28(s+d)d\bar{N}}{\delta}\right) \frac{16s(\bar{L}_2^2 B_{11} + \bar{L}_1^2 B_{22} + \bar{L}_{01}\bar{L}_2)}{c^2 \varepsilon^2}$$

So, to summarise, recalling that  $\delta = \rho/2$ ,  $\hat{\eta}_{X_0}$  is nondegenerate on  $\mathcal{X}_{\text{grid}}^{\text{near}}$  and  $\mathcal{X}_{\text{grid}}^{\text{far}}$  with probability at least  $1 - \delta$  (conditional on  $\bar{E}$ ) provided that

$$m \gtrsim \log\left(\frac{sdN}{\rho}\right) \log\left(\frac{sd}{\rho}\right) \frac{s(\bar{L}_{2}^{2}B_{11} + \bar{L}_{1}^{2}B_{22} + \bar{B}\bar{L}_{01}^{2} + \bar{L}_{01}\bar{L}_{2})}{\varepsilon^{2}}$$

and

$$\mathbb{P}(E_{\omega}^{c}) \lesssim \frac{\varepsilon}{\bar{B}^{3/2}\sqrt{s}\sqrt{\log(\bar{N}d/\rho)}} \quad \text{and} \quad , \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \lesssim \frac{\varepsilon}{4s\sqrt{B}\sqrt{\log(\bar{N}d/\rho)}}$$

Step I (b): Deterministic signs Assume now that u consists of arbitrary signs. We will show that (D.3) can be bounded by  $c\varepsilon$  when m is chosen as in condition (b) of this theorem. Let F' be the event that

(a') 
$$\left\| \Upsilon - \hat{\Upsilon} \right\| \leq \frac{t}{s^{1/4}} \text{ and } \left\| \Upsilon^{-1} - \hat{\Upsilon}^{-1} \right\| \leq \frac{t}{s^{1/4}}$$
  
(b')  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{far}}, \left\| (\hat{\mathbf{f}}_{X_0}(y) - \mathbf{f}_{X_0}(y)) \right\| \leq \frac{a\varepsilon_0}{s^{1/4}}$   
(c')  $\forall y \in \mathcal{X}_{\text{grid}}^{\text{near}}, \sup_{\|q\|=1} \left\| D_2 \left[ (\hat{\mathbf{f}}_{X_0} - \mathbf{f}_{X_0})^\top q \right] (y) \right\| \leq \frac{a\varepsilon_2}{s^{1/4}}$   
(f)  $\left\| (\Upsilon - \hat{\Upsilon}) \Upsilon^{-1} \mathbf{u} \right\|_{*,\infty} \leq a\varepsilon \left\| \Upsilon^{-1} \mathbf{u} \right\|_{*,\infty} \leq 2a\varepsilon.$ 

Then, provided that

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{u}{u + 6s(B_0 + B_2)} \quad \text{and} \quad \mathbb{E}[L_{01}(\omega)^2 \mathbb{1}_{\bar{E}^{c}}] \leqslant \frac{u}{4\bar{B}s^{3/2}},$$

with  $u = \min\{a\varepsilon_i, t\}$  as before, we have

$$\begin{split} \mathbb{P}_{\bar{E}}((F')^c) \leqslant & 4(d+1)s \exp\left(-\frac{mt^2}{16s^{3/2}\bar{L}_{01}^2(3+2t)}\right) \\ &+ 4sd \left|\mathcal{X}_{\text{grid}}^{\text{far}}\right| \exp\left(-\frac{m(a\varepsilon_0)^2/8}{s^{3/2}(\bar{L}_{01}^2(B_{11}+1)+\bar{L}_{01}^2)}\right) \\ &+ s(3d+d^2) \left|\mathcal{X}_{\text{grid}}^{\text{near}}\right| \exp\left(-\frac{m(a\varepsilon_2)^2/8}{s^{3/2}(\bar{L}_2^2B_{11}+\bar{L}_1^2B_{22}+\bar{L}_{01}\bar{L}_2)}\right) \\ &+ 32s \exp\left(-\frac{m4a^2\varepsilon^2}{s\left(32L_1^2+68a\varepsilon L_1\bar{L}_{01}\right)}\right). \end{split}$$

where the first bound is from Proposition C.1, the second and third are from Proposition C.4 and the final bound is due to Proposition C.3.

To bound (D.3), we first observe that if event G holds, then just as observed previously,  $|D_r[\mathbf{u}^\top\beta_2](y)| \leq 2a\varepsilon_r$ . To bound  $|\mathbf{u}^\top\beta_1(y)|$ , observe that

$$\begin{aligned} \mathbf{u}^{\top} \beta_{1}(y) &= \mathbf{u}^{\top} (\Upsilon^{-1} - \hat{\Upsilon}^{-1}) (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) + \mathbf{u}^{\top} (\Upsilon^{-1} - \hat{\Upsilon}^{-1}) \mathbf{f}_{X_{0}} \\ &= \mathbf{u}^{\top} (\Upsilon^{-1} - \hat{\Upsilon}^{-1}) (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) + \mathbf{u}^{\top} \Upsilon^{-1} (\hat{\Upsilon} - \Upsilon) \hat{\Upsilon}^{-1} \mathbf{f}_{X_{0}} \\ &= \mathbf{u}^{\top} (\Upsilon^{-1} - \hat{\Upsilon}^{-1}) (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) + \mathbf{u}^{\top} \Upsilon^{-1} (\hat{\Upsilon} - \Upsilon) (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \mathbf{f}_{X_{0}} + \mathbf{u}^{\top} \Upsilon^{-1} (\hat{\Upsilon} - \Upsilon) \Upsilon^{-1} \mathbf{f}_{X_{0}} \end{aligned}$$

Under event F',

• 
$$\left| \mathbf{u}^{\top} (\Upsilon^{-1} - \hat{\Upsilon}^{-1}) (\hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}}) \right| \leq \sqrt{s} \left\| \Upsilon^{-1} - \hat{\Upsilon}^{-1} \right\| \left\| \hat{\mathbf{f}}_{X_{0}} - \mathbf{f}_{X_{0}} \right\| \leq ta\varepsilon$$
  
•  $\left| \mathbf{u}^{\top} \Upsilon^{-1} (\hat{\Upsilon} - \Upsilon) (\hat{\Upsilon}^{-1} - \Upsilon^{-1}) \mathbf{f}_{X_{0}} \right| \leq \sqrt{s} \cdot 2 \cdot \left\| \hat{\Upsilon} - \Upsilon \right\| \left\| \hat{\Upsilon}^{-1} - \Upsilon^{-1} \right\| B_{0} \leq 2t^{2} B_{0}$   
•  $\left\| \Upsilon^{-1} (\hat{\Upsilon} - \Upsilon) \Upsilon^{-1} \mathbf{u} \right\|_{*,\infty} \leq \left\| \Upsilon^{-1} \right\|_{*,\infty} \left\| (\hat{\Upsilon} - \Upsilon) \Upsilon^{-1} \mathbf{u} \right\|_{*,\infty} \leq 4a\varepsilon.$ 

Finally, given any vector q such that  $||q||_{*,\infty} \leq 4a\varepsilon$ , we have  $|q^{\top}\mathbf{f}_{X_0}| \leq 4a\varepsilon B_0$ . Therefore,

$$\left|\mathbf{u}^{\top}\beta_{1}(y)\right| \leqslant ta + 2t^{2} + 4a\varepsilon B_{0}$$

and in a similar manner, we can show that the same upper bound holds for  $\|D_2[\mathbf{u}^\top\beta_1](y)\|$ . Therefore,

$$\left\| \mathbf{D}_r \left[ \mathbf{u}^\top \beta \right] (y) \right\| \leqslant c \varepsilon_r \tag{D.7}$$

if both F' and G hold, so conditional on  $\bar{E}$ , (D.7) holds with probability at least  $1 - \delta$  provided that

$$m \gtrsim s^{3/2} \cdot \frac{(\bar{L}_2^2 B_{11} + \bar{L}_1^2 B_{22} + \bar{B}\bar{L}_{01}^2 + \bar{L}_{01}\bar{L}_2)}{\varepsilon^2} \cdot \log\left(\frac{\bar{N}ds}{\rho}\right)$$

and

$$\mathbb{P}(E_{\omega}^{c}) \lesssim \frac{\varepsilon}{\bar{B}^{3/2} s \sqrt{\log(\bar{N}d/\rho)}} \quad \text{and} \quad , \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \lesssim \frac{\varepsilon}{s^{3/2} \sqrt{\bar{B}} \sqrt{\log(\bar{N}d/\rho)}}$$

Step II: Extending to the entire space To prove that  $\hat{\eta}_{X_0}$  is nondegenerate on the entire space  $\mathcal{X}$ , we first show that  $\hat{\eta}_{X_0}$  is locally Lipschitz (and hence determine how fine our grids  $\mathcal{X}_{\text{grid}}^{\text{near}}$ ,  $\mathcal{X}_{\text{grid}}^{\text{far}}$  need to be): for  $x, x' \in \mathcal{X}$  with  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$ ,

$$\|\mathbf{D}_{r}\left[\hat{\eta}_{X_{0}}\right](x) - \mathbf{D}_{r}\left[\hat{\eta}_{X_{0}}\right](x')\| = \left\|\frac{1}{m}\sum_{k=1}^{m}\mathbf{D}_{r}\left[\operatorname{Re}\left((\hat{\Upsilon}_{X}^{-1}\mathbf{u})^{\top}\gamma(\omega_{k})\varphi_{\omega_{k}}\right)\right](x)$$

$$- \mathbf{D}_{r}\left[\operatorname{Re}\left((\hat{\Upsilon}_{X}^{-1}\mathbf{u})^{\top}\gamma(\omega_{k})\varphi_{\omega_{k}}\right)\right](x')\right\|$$

$$= \left\|\frac{1}{m}\sum_{j=1}^{m}\operatorname{Re}\left(\left((\hat{\Upsilon}_{X}^{-1}\mathbf{u})^{\top}\gamma(\omega_{k})\right)\cdot\left(\mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](x) - \mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](x')\right)\right)\right\|$$

$$\leq \left\|\hat{\Upsilon}_{X}^{-1}\right\|\|\mathbf{u}\|\sqrt{s}\bar{L}_{01}\|\mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](x) - \mathbf{D}_{r}\left[\varphi_{\omega_{k}}\right](x')\right\|$$

$$(\mathbf{D}.9)$$

$$\leq 4s\bar{L}_{01}d_{\mathbf{H}}(x,x')\mathcal{L}_{r} \leq c\varepsilon_{r}.$$

$$(\mathbf{D}.10)$$

where we have applied Lemma C.2 to obtain the last line.

Choosing  $\mathcal{X}_{\text{grid}}^{\text{far}}$  to be a  $\delta_0 \stackrel{\text{def.}}{=} \frac{c\varepsilon_0}{4\mathcal{L}_0 \bar{L}_{01s}}$ -covering of  $\mathcal{X}^{\text{near}}$  (of size at most  $\mathcal{O}(R_{\mathcal{X}}/\delta_0)$ ),  $\mathcal{X}_{\text{grid}}^{\text{far}}$  to be a  $\delta_2 \stackrel{\text{def.}}{=} \frac{c\varepsilon_2}{4\mathcal{L}_2 \bar{L}_{01s}}$ -covering of  $\mathcal{X}^{\text{far}}$  (of size at most  $\mathcal{O}(R_{\mathcal{X}}/\delta_2)$ ). Then for any  $x \in \mathcal{X}^{\text{near}}$  and  $x' \in \mathcal{X}_{\text{grid}}^{\text{near}}$  such that  $d_{\mathbf{H}}(x, x') \leq \delta_0$ ,

$$|\hat{\eta}_{X_0}(x)| \leq |\hat{\eta}_{X_0}(x')| + |\hat{\eta}_{X_0}(x) - \hat{\eta}_{X_0}(x')| \leq 1 - \varepsilon_0 + 2c\varepsilon_0.$$

and given any  $x \in \mathcal{X}^{\text{far}}$ , let  $x' \in \mathcal{X}^{\text{far}}_{\text{grid}}$  be such that  $d_{\mathbf{H}}(x, x') \leq \delta_2$ , so

$$\operatorname{Re}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\hat{\eta}_{X_{0}}\right](x)\right) \preceq \operatorname{Re}\left(\overline{\operatorname{sign}(a_{i})}\operatorname{D}_{2}\left[\hat{\eta}_{X_{0}}\right](x')\right) + \left\|\operatorname{D}_{2}\left[\hat{\eta}_{X}\right](x) - \operatorname{D}_{2}\left[\hat{\eta}_{X}\right](x')\right\|\operatorname{Id} \preceq (-\varepsilon_{2} + 2c\varepsilon_{2})\operatorname{Id},$$

and

$$\left\| \operatorname{Im}\left(\overline{\operatorname{sign}(a_i)} \operatorname{D}_2\left[\hat{\eta}_{X_0}\right](x)\right) \right\| \leqslant \left\| \operatorname{Im}\left(\overline{\operatorname{sign}(a_i)} \operatorname{D}_2\left[\hat{\eta}_{X_0}\right](x')\right) \right\| + c\varepsilon_2 \leqslant (c_2 + c)\varepsilon_2.$$

### **D.2** Nondegeneracy transfer to $\hat{\eta}_X$ .

We are now ready to prove Theorem 3, which we restate below for clarity.

**Theorem D.2.** Under the assumptions of Theorem D.1, the following holds with probability at least  $1 - \rho$ : for all X such that

$$d_{\mathbf{H}}(X, X_0) \lesssim \min\left(r_{\text{near}}, \varepsilon_r (C_{\mathbf{H}} B \sqrt{s})^{-1}, \varepsilon_r (C_{\mathbf{H}} \bar{L}_{12} \bar{L}_r \sqrt{s})^{-1}\right), \tag{D.11}$$

we have

(i) for all 
$$y \in \mathcal{X}^{far}$$
,  $|\hat{\eta}_X(y)| \leq 1 - \frac{13}{32}\varepsilon_0$   
(ii) for all  $y \in \mathcal{X}^{near}(i)$ ,  $-\operatorname{Re}\left(\overline{\operatorname{sign}(a_i)}\operatorname{D}_2\left[\hat{\eta}_X\right](y)\right) \succeq \frac{13\varepsilon_2}{32}\operatorname{Id}$  and  $\left\|\operatorname{Im}\left(\overline{\operatorname{sign}(a_i)}\operatorname{D}_2\left[\hat{\eta}_X\right](y)\right)\right\| \leq (\frac{p}{2} + \frac{3p}{16})\frac{1}{2}\varepsilon_2$ 

Hence,  $\hat{\eta}_X$  is  $(\frac{13}{32}\varepsilon_0, \frac{13}{32}\varepsilon_2)$ -nondegenerate.

The proof essentially exploits the fact that  $\hat{\Upsilon}_X$ ,  $\hat{\mathbf{f}}_X$  are locally Lipschitz in X with respect to the metric  $d_{\mathbf{H}}$ , and consequently nondegeneracy of  $\hat{\eta}_{X_0}$  implies nondegeneracy of  $\hat{\eta}_X$  whenever  $d_{\mathbf{H}}(X, X_0)$  is sufficiently small.

### D.2.1 Proof of Theorem D.2

We begin with a lemma which shows that  $\hat{\Upsilon}_X$  is locally Lipschitz in X. Lemma D.1 (Lipschitz bound of  $\hat{\Upsilon}_X$ ). Let  $X_0 \in \mathcal{X}^s$  be  $\Delta$ -separated points. Assume that for all  $i + j \leq 3$ 

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{1}{1 + 16\sqrt{s}B_{ij}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{1}{16\sqrt{s}}$$

for all i, j = 0, ..., 2. Let  $\rho > 0$  and

$$m \gtrsim s(\bar{L}_2^2 B_{11} + \bar{L}_1^2 B_{22} + \bar{L}_{01} \bar{L}_2) \left( \log\left(\frac{sd}{\rho}\right) + d\log\left(sC_{\mathbf{H}} \max_{i=0}^3 \bar{L}_i\right) \right)$$

Then, conditional on event  $\overline{E}$ , with probability at least  $1 - \rho$ , the following hold:

• (i) for all X such that  $d_{\mathbf{H}}(x_i, x_{0,i}) \leq r_{\text{near}}$ , we have

$$\left\| \hat{\Upsilon}_X - \hat{\Upsilon}_{X_0} \right\| \lesssim C_{\mathbf{H}} B d_{\mathbf{H}}(X, X_0) \,.$$

• (ii) for all X such that  $d_{\mathbf{H}}(X, X_0) \lesssim \min\left(r_{\text{near}}, \frac{1}{C_{\mathbf{H}}B}\right)$ , we have  $\left\|\operatorname{Id} - \hat{\Upsilon}_X\right\| \leq \frac{3}{4}$  and  $\left\|\mathbf{G}_X^{-\frac{1}{2}}\Gamma_X^*\right\| \lesssim 1$ .

*Proof.* By Lemma C.8 and Lemma C.10, with probability at least  $1 - \rho$  conditonal on  $\overline{E}$ , for all  $(i, j) \in \{(0,0), (0,1), (1,1), (1,2)\}$  and all  $x, y \in \mathcal{X}^{\text{near}}$ ,

$$\left\| \hat{K}^{(ij)}(x,y) \right\| \leqslant \left\| K^{(ij)}(x,y) \right\| + \frac{1}{\sqrt{s}}$$

note that this also holds for  $\hat{K}^{(ji)}(x,y)$  since  $\hat{K}^{(ij)}(x,y) = \overline{\hat{K}^{(ij)}(y,x)}$ .

In particular, for all x, x' such that  $d_{\mathbf{H}}(x, x') \ge \Delta/4$ , we have  $\left\|\hat{K}^{(ij)}(x, x')\right\| \le \frac{2}{\sqrt{s}}$ . Take any X such that  $d_{\mathbf{H}}(x_i, x_{0,i}) \le r_{\text{near}}$ , we have that both  $x_i, x_{0,i}$  are at least  $\Delta/4$ -separated from  $x_j$  and  $x_{0,j}$ . Therefore, for  $k, \ell \in \{0, 1\}$ , using Lemma C.3:

$$\left\| \hat{K}^{(k\ell)}(x_i, x_j) - \hat{K}^{(k\ell)}(x_{i,0}, x_{j,0}) \right\| \lesssim \frac{C_{\mathbf{H}}}{\sqrt{s}} \sqrt{d_{\mathbf{H}}(x_i, x_{0,i})^2 + d_{\mathbf{H}}(x_j, x_{0,j})^2}$$

$$\left\| \hat{K}^{(k\ell)}(x_i, x_i) - \hat{K}^{(k\ell)}(x_{i,0}, x_{i,0}) \right\| \lesssim C_{\mathbf{H}} \left( B_{k+1,\ell} + B_{k,\ell+1} \right) d_{\mathbf{H}}(x_i, x_{0,i})$$
(D.12)

and therefore by Lemma G.6:

$$\begin{split} \left\| \hat{\Upsilon}_{X} - \hat{\Upsilon}_{X_{0}} \right\|^{2} &\leqslant \sum_{i,j=1}^{s} \sum_{k,\ell=0}^{1} \left\| \hat{K}^{(k\ell)}(x_{i},x_{j}) - \hat{K}^{(k\ell)}(x_{0,i},x_{0,j}) \right\|^{2} \\ &\leqslant 2 \sum_{i,j=1}^{s} \sum_{k,\ell=0}^{1} \left\| \hat{K}^{(k\ell)}(x_{i},x_{j}) - \hat{K}^{(k\ell)}(x_{0,i},x_{j}) \right\|^{2} + \left\| \hat{K}^{(\ell k)}(x_{j},x_{0,i}) - \hat{K}^{(\ell k)}(x_{0,j},x_{0,i}) \right\|^{2} \\ &\lesssim C_{\mathbf{H}}^{2} \left( \sum_{\substack{k,l \in \{0,1,2\}\\k+\ell \leqslant 3}} B_{k\ell} \right)^{2} \sum_{i} d_{\mathbf{H}}(x_{i},x_{0,i})^{2} + \frac{1}{s} \sum_{j \neq i} d_{\mathbf{H}}(x_{j},x_{0,j})^{2} \end{split}$$

which yields the desired result.

For the second statement, using Proposition C.1,  $\mathbb{P}_{\bar{E}}(\|\hat{\Upsilon}_{X_0} - \Upsilon_{X_0}\| > \frac{1}{8}) \leq \rho$ , so conditional on  $\bar{E}$ , we have with probability  $1 - \rho$ ,  $\|\hat{\Upsilon}_X - \hat{\Upsilon}_{X_0}\| \leq \frac{1}{8}$  and the claim follows since  $\|\mathrm{Id} - \Upsilon_{X_0}\| \leq \frac{1}{2}$  (due to Lemma C.1) implies that  $\|\mathrm{Id} - \hat{\Upsilon}_X\| \leq \frac{3}{4}$  and

$$\left\| \hat{\Upsilon}_X \right\| \leqslant 7/4 \text{ and } \left\| \mathbf{G}_X^{-\frac{1}{2}} \Gamma_X^* \right\| = \sqrt{\left\| \hat{\Upsilon}_X \right\|} \lesssim \sqrt{7}/2.$$

Proof of Theorem D.2. Since  $\hat{\eta}_{X_0}$  is nondegenerate with probability at least  $1 - \rho$ , the conclusion follows if we prove that for all  $x \in \mathcal{X}^{\text{far}}$  and all  $y \in \mathcal{X}^{\text{near}}$ ,

$$\|\mathbf{D}_{2}\left[\hat{\eta}_{X} - \hat{\eta}_{X_{0}}\right](x)\| \leqslant \varepsilon_{0}/32 \quad \text{and} \quad \|\mathbf{D}_{2}\left[\hat{\eta}_{X} - \hat{\eta}_{X_{0}}\right](y)\| \leqslant p\varepsilon_{2}/32 \tag{D.13}$$

with probability at least  $1 - \rho$ . We first write

$$\hat{\eta}_X(y) - \hat{\eta}_{X_0}(y) = \hat{\alpha}_X^\top (\hat{\mathbf{f}}_X - \hat{\mathbf{f}}_{X_0}) + (\hat{\alpha}_X - \hat{\alpha}_{X_0})^\top \hat{\mathbf{f}}_{X_0}(y).$$

Conditional on  $\bar{E}$ , with probability at least  $1 - \rho/2$ , we have by Lemma D.1 (note that our assumptions imply the assumptions of Lemma D.1),  $\|\Upsilon_X - \Upsilon_{X_0}\| \lesssim C_{\mathbf{H}}Bd_{\mathbf{H}}(X, X_0)$  and  $\|\Upsilon_X^{-1}\| \leq 4$ . So,

$$\left\| \mathbf{D}_r \left[ \left( \hat{\alpha}_X - \hat{\alpha}_{X_0} \right)^\top \hat{\mathbf{f}}_{X_0} \right] (y) \right\| \leqslant \sqrt{s} \left\| \Upsilon_X^{-1} - \Upsilon_{X_0}^{-1} \right\| \leqslant 8\sqrt{s} \left\| \hat{\Upsilon}_X - \hat{\Upsilon}_{X_0} \right\| \lesssim \sqrt{s} C_{\mathbf{H}} B d_{\mathbf{H}}(X, X_0).$$

By Lemma C.2, if  $\overline{E}$  occurs, then

$$\left\| \mathbf{D}_r \left[ \hat{\alpha}_X^\top (\hat{\mathbf{f}}_X - \hat{\mathbf{f}}_{X_0}) \right] (y) \right\| \leqslant C_r \left\| \hat{\alpha}_X \right\| d_{\mathbf{H}}(X, X_0) \leqslant C_r \left\| \hat{\Upsilon}_X^{-1} \right\| \sqrt{s} d_{\mathbf{H}}(X, X_0) \leqslant 4C_r \sqrt{s} d_{\mathbf{H}}(X, X_0),$$

where  $C_r \leq (1 + C_{\mathbf{H}}) \bar{L}_r \bar{L}_{12}$ . Finally, since  $\mathbb{P}(\bar{E}^c) \leq \rho/2$ , we have with probability at least  $1 - \rho$ , for all  $y \in \mathcal{X}$ , (D.13) holds provided that (D.11) holds. Combining with the nondegeneracy of  $\hat{\eta}_{X_0}$ , the conclusion follows with probability  $1 - 2\rho$ .

## E Supplementary results to the proof Theorem 1

Recall that in the proof of Theorem 1, we defined the function  $f: \mathbb{C}^s \times \mathcal{X}^s \times \mathbb{R}_+ \times \mathbb{C}^m$  by

$$f(u,v) \stackrel{\text{def.}}{=} \Gamma_X^*(\Phi_X a - \Phi_{X_0} a_0 - w) + \lambda \begin{pmatrix} \operatorname{sign}(a_0) \\ 0_{sd} \end{pmatrix}$$

where u = (a, X) and  $v = (\lambda, w)$ . This function f is differentiable with

$$\partial_v f(u,v) = \left( \begin{pmatrix} \operatorname{sign}(a_0) \\ 0_{sd} \end{pmatrix}, \ -\Gamma_X^* \right) \in \mathbb{C}^{s(d+1) \times m}, \tag{E.1}$$

and  $\partial_u f(u, v)$  is

$$\Gamma_X^* \Gamma_X J_a + \begin{pmatrix} 0_{1\times s} & A_{11} & 0 & \cdots & 0\\ 0_{1\times s} & 0 & A_{12} & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & \vdots\\ 0_{1\times s} & 0 & 0 & \cdots & A_{1s}\\ 0_{d\times s} & A_{21} & 0 & \cdots & 0\\ 0_{d\times s} & 0 & A_{22} & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & \vdots\\ 0_{d\times s} & 0 & 0 & \cdots & A_{2s} \end{pmatrix}$$
(E.2)

where  $A_{1j} \stackrel{\text{def.}}{=} \nabla_x \langle \varphi(x_j), z \rangle^\top$ ,  $A_{2j} \stackrel{\text{def.}}{=} \nabla_x^2 \langle \varphi(x_j), z \rangle$ ,  $z \stackrel{\text{def.}}{=} (\Phi_X a - \Phi_{X_0} a_0 - w)$  and  $J_a \in \mathbb{R}^{s(d+1) \times s(d+1)}$  is a the diagonal matrix:

$$J_a = \begin{pmatrix} \mathrm{Id}_{s \times s} & & 0 \\ & a_1 \mathrm{Id}_{d \times d} & & \\ & & \ddots & \\ 0 & & & a_s \mathrm{Id}_{d \times d} \end{pmatrix}$$

Letting  $u_0 = (a_0, X_0)$  and  $v_0 = (0, 0)$ ,  $\partial_u f(u_0, v_0) = \Gamma^*_{X_0} \Gamma_{X_0} J_a$  is invertible and  $f(u_0, v_0) = 0$ . Hence, by the Implicit Function Theorem, there exists a neighbourhood V of  $v_0$  in  $\mathbb{C} \times \mathbb{C}^m$ , a neighbourhood U of  $u_0$  in  $\mathbb{C}^s \times \mathcal{X}^s$ and a Fréchet differentiable function  $g: V \to U$  such that for all  $(u, v) \in U \times V$ , f(u, v) = 0 if and only if u = g(v). To conclude, we simply need to bound the size of the region on which g is well defined, and to bound the error between g(v) and g(0). Let us first remark that our assumptions imply that  $\mathbb{P}(\bar{E}^c) \leq \rho/2$  and

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{1}{1 + 16\sqrt{s}B_{ij}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{1}{16\sqrt{s}}, \tag{E.3}$$

for all i, j = 0, ..., 2. Therefore, it is sufficient to prove the existence of g conditional on event  $\overline{E}$ :

**Theorem E.1.** Assume that for all  $i + j \leq 3$ 

$$\mathbb{P}(E_{\omega}^{c}) \leqslant \frac{1}{1 + 16\sqrt{s}B_{ij}}, \quad \mathbb{E}[L_{i}(\omega)L_{j}(\omega)1_{E_{\omega}^{c}}] \leqslant \frac{1}{16\sqrt{s}}$$

for all i, j = 0, ..., 2. Let  $\rho > 0$  and suppose that

$$m \gtrsim s(\bar{L}_2^2 B_{11} + \bar{L}_1^2 B_{22} + \bar{L}_{01} \bar{L}_2) \left( \log\left(\frac{sd}{\rho}\right) + d\log\left(sC_{\mathbf{H}} \mathbb{L}_3\right) \right)$$

where  $\mathbb{L}_r \stackrel{\text{def.}}{=} \max_{i \leq r} L_r$ . Then, conditional on event  $\overline{E}$ , with probability at least  $1 - \rho$ : there exists a  $\mathscr{C}^1$  function g such that, for all  $v = (\lambda, w)$  such that  $||v|| \leq r$  with r satisfying

$$r = \mathcal{O}\left(\frac{1}{\sqrt{s}}\min\left(\frac{\min\{r_{\text{near}}, (C_{\mathbf{H}}B)^{-1}\}}{\min_{i}|a_{0,i}|}, \frac{1}{\bar{L}_{01}\bar{L}_{12}(1+||a_{0}||)}, \right)\right)$$
(E.4)

we have f(g(v), v) = 0 and  $g(0) = u_0$ . Furthermore, given  $(\lambda, w)$  in this ball,  $(a, X) \stackrel{\text{def.}}{=} g((\lambda, w))$  satisfies

$$||a - a_0|| + d_{\mathbf{H}}(X, X_0) \leqslant \frac{\sqrt{s}(\lambda + ||w||)}{\min_i |a_{0,i}|}.$$
(E.5)

We begin with some preliminary results before presenting the proof of this theorem in Section E.2.

#### E.1 Preliminary results

**Theorem E.2** (Quantitative implicit function theorem, adapted from Denoyelle et al. (2017)). Let  $F : \mathcal{H} \times \mathcal{Y} \to \mathbb{C}^n$ be a differentiable mapping where  $\mathcal{H}$  is a Hilbert space,  $\mathcal{Y} \subseteq \mathbb{C}^s \times \mathbb{R}^{sd}$ , n = s(d+1),  $\|\cdot\|$  be a norm on  $\mathcal{H}$ . For each  $y \in \mathcal{Y}$ , suppose that there exists a positive definite matrix  $\mathbf{G}_y$ , and let  $d_G$  be the associated metric. Assume that  $F(x_0, y_0) = 0$ , and that for  $x \in \mathcal{B}_{\|\cdot\|}(x_0, r_1), y \in \mathcal{B}_{d_G}(y_0, r_2), \partial_y F(x, y)$  is invertible and we have

$$\left\|\mathbf{G}_{y}^{-\frac{1}{2}}\partial_{x}F(x,y)\right\| \leqslant D_{1} \quad and \quad \left\|\mathbf{G}_{y}^{\frac{1}{2}}\partial_{y}F(x,y)^{-1}\mathbf{G}_{x}^{\frac{1}{2}}\right\| \leqslant D_{2}.$$

Then, defining  $R = \min\left(\frac{r_2}{D_1D_2}, r_1\right)$ , there exists a unique Fréchet differentiable mapping  $g : \mathcal{B}_{\|\cdot\|}(x_0, R) \to \mathcal{B}_{d_G}(y_0, r_2)$  such that  $g(x_0) = y_0$  and for all  $x \in \mathcal{B}_{\|\cdot\|}(x_0, R)$ , F(x, g(x)) = 0, and furthermore

$$dg(x) = -(\partial_y F(x, g(x)))^{-1} \partial_x F(x, g(x))$$

and consequently  $\left\|\mathbf{G}_{g(x)}^{\frac{1}{2}}\mathrm{d}g(x)\right\| \leq D_1 D_2.$ 

*Proof.* Let  $V^* = \bigcup_{V \in \mathcal{V}} V$ , where  $\mathcal{V}$  is the collection of all open sets  $V \in \mathbb{R}^m$  such that

- 1.  $x_0 \in V$ ,
- 2. V is star-shaped with respect to  $x_0$ ,

3. 
$$V \subset \mathcal{B}_{\parallel \cdot \parallel}(x_0, r_1),$$

4. there exists a  $\mathcal{C}^1$  function  $g: V \to \mathcal{B}_{d_G}(y_0, r_2)$  such that  $g(x_0) = y_0$  and F(x, g(x)) = 0 for all  $x \in V$ .

Observe that  $\mathcal{V}$  is non-empty by the (classical) Implicit Function Theorem. Moreover,  $\mathcal{V}$  is stable by union: indeed, all conditions expect the last one are easy to check. Now, let  $V, \tilde{V} \in \mathcal{V}$  and  $g, \tilde{g}$  be corresponding functions. The set  $\overline{V} = \{x \in V \cap \tilde{V}, g(x) = \tilde{g}(x)\}$  is non-empty (it contains  $x_0$ ), and closed in  $V \cap \tilde{V}$ . Moreover, it is open: for any  $x \in \overline{V}$ , by our assumptions  $\partial_y F(x, g(x))$  is invertible and the Implicit Function theorem applies at (x, g(x)), and by the uniqueness of the mapping resulting from it we obtain an open set around x in which g and  $\tilde{g}$  coincide. Hence  $\overline{V}$  is both closed and open in  $V \cap \tilde{V}$ , and by the connectedness of it  $\overline{V} = V \cap \tilde{V}$ . Therefore, there exists a function g' defined on  $V \cup \tilde{V}$  that satisfies condition 4. above (it is defined as g on V and  $\tilde{g}$  on  $\tilde{V}$ , which is well-posed for their intersection), and  $\mathcal{V}$  is indeed stable by union.

Hence  $V^* \in \mathcal{V}$ , let  $g^*$  be its corresponding function. It is unique by the arguments above, satisfies  $F(x, g^*(x)) = 0$ and

$$\begin{aligned} \mathbf{G}_{g^{*}(x)}^{\frac{1}{2}} \mathrm{d}g^{*}(x) &= -\mathbf{G}_{g^{*}(x)}^{\frac{1}{2}} (\partial_{y} F(x, g^{*}(x)))^{-1} \partial_{x} F(x, g^{*}(x)) \\ &= -(\mathbf{G}_{g^{*}(x)}^{-\frac{1}{2}} \partial_{y} F(x, g^{*}(x)) \mathbf{G}_{g^{*}(x)}^{-\frac{1}{2}})^{-1} \mathbf{G}_{g^{*}(x)}^{-\frac{1}{2}} \partial_{x} F(x, g^{*}(x)) \end{aligned}$$

for all  $x \in V^*$ . Note that by our assumptions  $\left\| \mathbf{G}_{g^*(x)}^{\frac{1}{2}} \mathrm{d}g^*(x) \right\| \leq D_1 D_2$ .

We finish the proof by showing that  $V^*$  contains a ball of radius  $r_2/(D_1D_2)$ . Let  $x \in \mathbb{R}^m$  with ||x|| = 1,  $R_x = \sup\{R, x_0 + Rx \in V^*\}$ , and  $x^* = x_0 + R_x x \in \partial V^*$ . Clearly  $0 < R_x \leq r_1$  since  $V^*$  is open, assume  $R_x < r_1$ . Our goal is to show that in that case  $R_x \geq \frac{r_1}{D_1D_2}$ . Since  $dg^*$  is bounded,  $g^*$  is uniformly continuous on  $V^*$  and it can be extended on  $\partial V^*$ , and by continuity  $F(x^*, g^*(x^*)) = 0$ . By contradiction, if  $g^*(x^*) \in \mathcal{B}_{d_G}(y_0, r_2)$ , by our assumptions we can apply the Implicit Function Theorem at  $(x^*, g^*(x^*))$ , and therefore extend  $g^*$  on an open set V that is not included in V<sup>\*</sup> such that  $V \cup V^* \in \mathcal{V}$ , which contradicts the maximality of V<sup>\*</sup>. Hence  $d_G(g^*(x^*), y_0) = r_2$ . Let  $\gamma : [0, 1] \to \mathcal{Y}$  be defined by  $\gamma(t) \stackrel{\text{def.}}{=} g^*(x^* + t(x_0 - x^*))$ , so  $\gamma'(t) = \mathrm{d}g^*(\gamma(t))(x_0 - x^*)$ . Then,

$$r_{2} = d_{G}(g^{*}(x^{*}), g^{*}(x_{0})) \leqslant \sqrt{\int_{0}^{1} \langle \mathbf{G}_{g^{*}(\gamma(t))} \gamma'(t), \gamma'(t) \rangle \mathrm{d}t}$$
$$= \sqrt{\int_{0}^{1} \left\| \mathbf{G}_{g^{*}(\gamma(t))}^{\frac{1}{2}} \mathrm{d}g^{*}(\gamma(t))(x_{0} - x^{*}) \right\|^{2} \mathrm{d}t} \leqslant D_{1} D_{2} R_{x}.$$

**Lemma E.1.** Assume that event  $\overline{E}$  occurs. Then, for all X such that  $d_{\mathbf{H}}(x_i, x_{0,i}) \leq r_{\text{near}}$ ,

$$\|\Pi_X \Gamma_{X_0} a\| \lesssim \begin{cases} \bar{L}_2 \, \|a\|_1 \max_i d_{\mathbf{H}}(x_i, x_{0,i})^2 \\ \bar{L}_2 \, \|a\|_\infty \, d_{\mathbf{H}}(X, X_0)^2 \end{cases}$$

*Proof.* Recall that  $\operatorname{Im}(\Gamma_X) = \{\varphi(x_i), J_{\varphi}(x_i)\}_i$ , and  $\Pi_X$  is a projector on  $\operatorname{Im}(\Gamma_X)^{\perp}$ . Also note that for  $d_{\mathbf{H}}(x_i, x_{0,i}) \leq r_{\operatorname{near}}$ , we have  $\left\| \mathbf{H}_{x_{0,i}}^{-\frac{1}{2}} \mathbf{H}_{x_i}^{\frac{1}{2}} \right\| \leq 1$ , and therefore under  $\bar{E}$ :

$$\left\|\mathbf{H}_{x_{0,i}}^{-\frac{1}{2}}\nabla^{2}\varphi_{\omega_{j}}(x_{i})\mathbf{H}_{x_{0,i}}^{-\frac{1}{2}}\right\| \lesssim \left\|\mathbf{D}_{2}\left[\varphi_{\omega_{j}}\right](x_{i})\right\| \leqslant \bar{L}_{2}$$

Let  $\gamma_i : [0,1] \to \mathcal{X}$  be any piecewise smooth curve such that  $\gamma_i(1) = x_{0,i}$  and  $\gamma_i(0) = x_i$ . Then, by Taylor expanding  $\varphi(\gamma_i(t))$  about t = 0, we obtain

$$\varphi(x_{0,i}) = \varphi(x_i) + \langle \nabla \varphi(x_i), \gamma'_i(0) \rangle + \int_0^1 \frac{1}{2} \langle \nabla^2 \varphi(\gamma_i(t)) \gamma'_i(t), \gamma'_i(t) \rangle \mathrm{d}t.$$

Therefore,

$$\Pi_X \Gamma_{X_0} a = \Pi_X \left( \sum_{i=1}^s a_i \varphi(x_{0,i}) \right) = \Pi_X \left( \sum_{i=1}^s \frac{a_i}{2} \int_0^1 \langle \nabla^2 \varphi(\gamma_i(t)) \gamma_i'(t), \gamma_i'(t) \rangle \mathrm{d}t \right)$$

Taking the norm implies

$$\|\Pi_X \Gamma_{X_0} a\| \leqslant \sum_{i=1}^s \frac{|a_i|}{2} \int_0^1 \bar{L}_2 \|\mathbf{H}_{\gamma_i(t)} \gamma_i'(t)\|^2 dt$$

and taking the infimum over all paths  $\gamma_i$  yields

$$\|\Pi_X \Gamma_{X_0} a\| \leq \bar{L}_2 \sum_i |a_i| d_{\mathbf{H}}(x_i, x_{0,i})^2.$$

### E.2 Proof of Theorem E.1

Our goal is to apply Theorem E.2. Let u = (a, X),  $u_0 = (a_0, X_0)$ ,  $v = (\lambda, w)$  and  $v_0 = (0, 0)$ . We must control  $\left\| \mathbf{G}_X^{-\frac{1}{2}} \partial_v f(u, v) \right\|$  and  $\left\| \mathbf{G}_X^{\frac{1}{2}} \partial_u f(u, v)^{-1} \mathbf{G}_X^{\frac{1}{2}} \right\|$  for (u, v) sufficiently close to  $(u_0, v_0)$ . Using Lemma D.1, conditional on event  $\overline{E}$ , with probability  $1 - \rho$  we have

$$\left\|\mathbf{G}_{X}^{-\frac{1}{2}}\partial_{v}f(u,v)\right\| \leqslant \left\|\mathbf{u}\right\| + \left\|\mathbf{G}_{X}^{-\frac{1}{2}}\Gamma_{X}\right\| \lesssim \sqrt{s}$$

To control  $\left\|\mathbf{G}_X^{\frac{1}{2}}\partial_u f(u,v)^{-1}\mathbf{G}_X^{\frac{1}{2}}\right\|$ , first observe that

$$\mathbf{G}_X^{-1/2} \partial_u f(u, v) \mathbf{G}_X^{-1/2} = \left( \mathbf{G}_X^{-1/2} \Gamma_X^* \Gamma_X \mathbf{G}_X^{-1/2} + M(u, v) \right) J_a$$

where

$$M(u,v) \stackrel{\text{def.}}{=} \begin{pmatrix} 0_{1\times s} & \frac{1}{a_1} \left( \mathbf{H}_{x_1}^{-\frac{1}{2}} \nabla[\langle \varphi, z \rangle](x_1) \right)^\top & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0_{1\times s} & 0 & \cdots & \frac{1}{a_s} \left( \mathbf{H}_{x_s}^{-\frac{1}{2}} \nabla[\langle \varphi, z \rangle](x_s) \right)^\top \\ 0_{d\times s} & \frac{1}{a_1} \mathbf{H}_{x_1}^{-\frac{1}{2}} \nabla^2[\langle \varphi, z \rangle](x_1) \mathbf{H}_{x_{0,1}}^{-\frac{1}{2}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0_{d\times s} & 0 & \cdots & \frac{1}{a_s} \mathbf{H}_{x_s}^{-\frac{1}{2}} \nabla^2[\langle \varphi, z \rangle](x_s) \mathbf{H}_{x_{0,s}}^{-\frac{1}{2}} \end{pmatrix},$$
(E.6)

where  $z = (\Phi_X a - \Phi_{X_0} a_0 - w)$ . Now, let us study the invertibility of  $\mathbf{G}_X^{-\frac{1}{2}} \Gamma_X^* \Gamma_X \mathbf{G}_X^{-\frac{1}{2}} + M(u, v)$  and bound the norm of its inverse.

**Lemma E.2** (Bound on M(u, v)). Let u = (a, X),  $v = (\lambda, w)$  and let M(u, v) be as defined in (E.6). Assume that  $\overline{E}$  occurs and given  $\varepsilon > 0$ , let  $c_{\varepsilon} \stackrel{\text{def.}}{=} \frac{\varepsilon \min_{i}|a_{0,i}|}{2L_{12}}$ . Then, for all  $X \in \mathcal{X}^{s}$ ,  $a \in \mathbb{R}^{s}$  and  $w \in \mathbb{C}^{m}$  such that

$$\|a - a_0\| \leq \frac{c_{\varepsilon}}{3\bar{L}_0}, \quad \|w\| \leq c_{\varepsilon}/3 \quad and \quad d_{\mathbf{H}}(X, X_0) \leq \min\left(r_{\mathrm{near}}, \frac{c_{\varepsilon}}{3\bar{L}_1 \|a_0\|}\right),$$

we have

$$\|M(u,v)\|\leqslant \varepsilon \quad and \quad \|M(u,v)\|_{*,\infty}\leqslant \varepsilon$$

*Proof.* First note that for  $r \in \mathbb{N}_0$ ,

$$\left\| \mathbf{D}_r \left[ \boldsymbol{\varphi}^\top \boldsymbol{z} \right] (\boldsymbol{x}_i) \right\| \leqslant \frac{1}{\sqrt{m}} \sum_{j=1}^m \left\| \boldsymbol{z}_j \mathbf{D}_r \left[ \boldsymbol{\varphi}_{\omega_j} \right] (\boldsymbol{x}_i) \right\| \leqslant L_r \left\| \boldsymbol{z} \right\|$$

Now, for  $\bar{q} = [q_1, \ldots, q_s, Q_1, \ldots, Q_s] \in \mathbb{C}^{s(d+1)}$ , where  $q_i \in \mathbb{C}$  and  $Q_i \in \mathbb{C}^d$ , and  $\|\bar{q}\| = 1$ , we have

$$\begin{split} \|M(u,v)\bar{q}\|^{2} &= \sum_{i=1}^{s} \left| \frac{1}{a_{i}} \left( \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \right)^{\top} Q_{i} \right|^{2} + \left\| \frac{1}{a_{i}} \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla^{2}[\varphi^{\top} z](x_{i}) \mathbf{H}_{x_{i}}^{-\frac{1}{2}} Q_{i} \right\|^{2} \\ &\leqslant \frac{4}{\min_{i} |a_{0,i}|^{2}} \|q\|^{2} \max_{i} \left( \left\| \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \right\|^{2} + \left\| \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla^{2}[\varphi^{\top} z](x_{i}) \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \right\|^{2} \right) \\ &= \frac{4}{\min_{i} |a_{0,i}|^{2}} \max_{i} \left( \left\| \mathbf{D}_{1} \left[ \varphi^{\top} z \right](x_{i}) \right\|^{2} + \left\| \mathbf{D}_{2} \left[ \varphi^{\top} z \right](x_{i}) \right\|^{2} \right) \\ &\leqslant \frac{4}{\min_{i} |a_{0,i}|^{2}} (\bar{L}_{1}^{2} + \bar{L}_{2}^{2}) \|z\|^{2} \end{split}$$

where we have used the fact that  $\min_{i} |a_{i}| \ge \min_{i} |a_{0,i}|/2$ . If  $\|\bar{q}\|_{*,\infty} = 1$ , then

$$\|M(u,v)\bar{q}\|_{*,\infty} = \max_{i} \{ \left\| \left( \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \right)^{\top} Q_{i} \right\|, \left\| \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \mathbf{H}_{x_{i}}^{-\frac{1}{2}} Q_{i} \right\|^{2} \} \\ \leq \max_{i} \{ \left\| \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \right\|, \left\| \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \nabla[\varphi^{\top} z](x_{i}) \mathbf{H}_{x_{i}}^{-\frac{1}{2}} \right\|^{2} \}$$

and the same bound holds.

Now it remains to bound ||z||. Writing  $\varphi(x) \stackrel{\text{def.}}{=} (\varphi_{\omega_k}(x))_{k=1}^m$ , we have

$$\begin{aligned} \|z\| &= \left\| \sum_{i} (a_{i}\varphi(x_{i}) - a_{0,i}\varphi(x_{0,i})) - w \right\| \\ &\leq \bar{L}_{0} \|a - a_{0}\| + \|a_{0}\| \max_{k} \sqrt{\sum_{i} |\varphi_{\omega_{k}}(x_{i}) - \varphi_{\omega_{k}}(x_{0,i})|^{2}} + \|w\| \\ &\leq \bar{L}_{0} \|a - a_{0}\| + \|a_{0}\| \bar{L}_{1} d_{\mathbf{H}}(X, X_{0}) + \|w\| \end{aligned}$$

where the last inequality follows from Lemma C.2.

The bound on ||M(u, v)|| from Lemma E.2 allows us to conclude that under event  $\overline{E}$ , taking

$$c \stackrel{\text{def.}}{=} \frac{\min_i |a_{0,i}|}{16\bar{L}_{12}} \tag{E.7}$$

for all  $X \in \mathcal{X}^s$ ,  $a \in \mathbb{R}^s$  and  $w \in \mathbb{C}^m$  such that

$$||a - a_0|| \leq \frac{c}{3\bar{L}_0}, ||w|| \leq c/3 \text{ and } d_{\mathbf{H}}(X, X_0) \leq \min\left(r_{\text{near}}, \frac{c}{3\bar{L}_1 ||a_0||}\right)$$

we have  $||M(u,v)|| \leq \frac{1}{8}$ . Combining this with Lemma D.1 gives

$$\left\| \mathrm{Id} - (\mathbf{G}_X^{-\frac{1}{2}} \Gamma_X^* \Gamma_X \mathbf{G}_X^{-\frac{1}{2}} + M(u, v)) \right\| \leq \left\| \mathrm{Id} - \mathbf{G}_X^{-\frac{1}{2}} \Gamma_X^* \Gamma_X \mathbf{G}_X^{-\frac{1}{2}} \right\| + \|M(u, v)\| < \frac{7}{8}$$

and therefore it is invertible and

$$\left\| (\mathbf{G}_{X}^{-\frac{1}{2}} \Gamma_{X}^{*} \Gamma_{X} \mathbf{G}_{X}^{-\frac{1}{2}} + M(u, v))^{-1} \right\| \leq \frac{1}{1 - \left\| \mathrm{Id} - (\mathbf{G}_{X}^{-\frac{1}{2}} \Gamma_{X}^{*} \Gamma_{X} \mathbf{G}_{X}^{-\frac{1}{2}} + M(u, v)) \right\|} = \mathcal{O}\left(1\right).$$

In this case,  $\partial_u f(u, v)$  is invertible, and we have

$$\left\| (\mathbf{G}_X^{-\frac{1}{2}} \partial_u f(u, v) \mathbf{G}_X^{-\frac{1}{2}})^{-1} \right\| = \left\| J_a^{-1} (\mathbf{G}_X^{-\frac{1}{2}} \Gamma_X^* \Gamma_X \mathbf{G}_X^{-\frac{1}{2}} + M(u, v))^{-1} \right\| \lesssim \frac{1}{\min_i |a_{0,i}|}$$

since  $||a - a_0|| \lesssim \min_i |a_{0,i}|$  by assumption.

Therefore we can apply Theorem E.2 with (recalling the definition of c in (E.7))

$$r_1 = c, \ D_1 = \mathcal{O}\left(\sqrt{s}\right), \ r_2 = \mathcal{O}\left(\min\left(r_{\text{near}}, \ \frac{c}{\tilde{L}_1 \| a_0 \|}, \frac{c}{\tilde{L}_0}, \frac{1}{C_{\mathbf{H}}B}\right)\right), \ D_2 = \mathcal{O}\left(\frac{1}{\min_i | a_{0,i} |}\right)$$

with  $B = \sum_{i+j \leqslant 3} B_{ij}$ , we obtain that g(v) is defined for  $v \in V \stackrel{\text{def.}}{=} \mathcal{B}_{\|\cdot\|_2}(0,r)$  with

$$r \stackrel{\text{def.}}{=} \min\left(\frac{r_2}{D_1 D_2}, r_1\right) = \frac{r_2}{D_1 D_2} = \mathcal{O}\left(\min\left(\frac{r_{\text{near}}}{\sqrt{s}\min_i |a_{0,i}|}, \frac{1}{\sqrt{s}\bar{L}_1\bar{L}_{12}} \|a_0\|, \frac{1}{\sqrt{s}\bar{L}_{12}\bar{L}_0}, \frac{1}{\sqrt{s}\min_i |a_{0,i}|C_{\mathbf{H}}B}\right)\right)$$

such that g is  $C^1$ , f(g(v), v) = 0,  $g(v_0) = u_0$ , where we recall that  $u_0 = (a_0, X_0)$  and  $v_0 = (0, 0)$ . Finally, from Theorem E.2 we also have that

$$\|\mathbf{G}_X \mathrm{d}g(v)\| \leqslant D_1 D_2 \lesssim \frac{\sqrt{s}}{\min_i |a_{0,i}|}$$

and by defining  $\gamma(t) = g(v_0 + t(v - v_0))$  for  $t \in [0, 1]$ , we have the following error bound between u = g(v) and

 $u_0 = g(v_0)$ :

$$d_G(u, u_0) = \sqrt{\|a - a_0\|_2^2 + d_{\mathbf{H}}(X, X_0)^2} \leqslant \sqrt{\int_0^1 \langle \mathbf{G}_{\gamma(t)} \gamma'(t), \gamma'(t) \rangle \mathrm{d}t}$$
$$= \sqrt{\int_0^1 \langle \mathbf{G}_{\gamma(t)} \mathrm{d}g(tv) v, \, \mathrm{d}g(tv) v \rangle \mathrm{d}t}$$
$$\leqslant \frac{\sqrt{s}}{\min_i |a_{0,i}|} \|v\|.$$

### **F** Examples

#### F.1 Fejér kernel

Let  $f \in \mathbb{N}$  and  $\mathcal{X} \in \mathbb{T}^d$  the *d*-dimensional torus. We consider the Fejér kernel

$$K(x, x') = \prod_{i=1}^{d} \kappa(x_i - x'_i)$$

where  $\kappa(x) \stackrel{\text{def.}}{=} \left( \frac{\sin\left(\left(\frac{f}{2}+1\right)\pi x\right)}{\left(\frac{f}{2}+1\right)\sin(\pi x)} \right)^4$ , with constant metric tensor

$$\mathbf{H}_{x} = C_{f} \mathrm{Id}$$
 and  $d_{\mathbf{H}}(x, x') = C_{f}^{-\frac{1}{2}} \|x - x'\|_{2}$ .

where  $C_f \stackrel{\text{def.}}{=} -\kappa''(0) = \frac{\pi^2}{3}f(f+4) \sim f^2$ . Note that  $K^{(ij)} = C_f^{-(i+j)/2} \nabla_1^i \nabla_2^j K$  and since the metric is constant, we can set  $C_{\mathbf{H}} \stackrel{\text{def.}}{=} 0$ .

### F.1.1 Discrete Fourier sampling

A random feature expansion associated with the Fejér kernel is obtained by choosing  $\Omega = \{\omega \in \mathbb{Z}^d ; \|\omega\|_{\infty} \leq f\}, \varphi_{\omega}(x) \stackrel{\text{def.}}{=} e^{i2\pi\omega^{\top}x}$ , and  $\Lambda(\omega) = \prod_{j=1}^d g(\omega_j)$  where  $g(j) = \frac{1}{f} \sum_{k=\max(j-f,-f)}^{\min(j+f,f)} (1 - |k/f|)(1 - |(j-k)/f|)$ . Note that this corresponds to sampling discrete Fourier frequencies. In this case, the derivatives of the random features are uniformly bounded with  $\|\nabla^j \varphi_{\omega}(x)\| = \|\omega\|^j = \mathcal{O}(C_f^{j/2} d^{j/2})$ . So, we can set  $\bar{L}_i = \mathcal{O}(d^{i/2})$ .

### F.1.2 Admissibility of the kernel

**Theorem F.1.** Suppose that  $f \ge 128$ . Then, K is an admissible kernel with  $r_{\text{near}} = 1/(8\sqrt{2})$ ,  $\varepsilon_2 = 0.941$ ,  $\varepsilon_0 = 0.00097$ ,  $h = \mathcal{O}(d^{-1/2})$  and  $\Delta = \mathcal{O}(d^{1/2}s_{\text{max}}^{1/4})$ ,  $B_{00} = B_{11} = B_{20} = \mathcal{O}(1)$ ,  $B_{01} = \mathcal{O}(d^{1/2})$  and  $B_{22} = \mathcal{O}(d)$ .

The remainder of this section is dedicated to proving this theorem. The uniform bounds on  $B_{ij}$  are due to Lemma F.4 (uniform bounds), and the bound on  $\Delta$  and h are due to Lemma F.3. From Lemma F.1, we see that by setting  $r_{\text{near}} \stackrel{\text{def.}}{=} \frac{1}{8\sqrt{2}}$ , for all  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$ ,  $K^{(20)}(x, x') \prec -\varepsilon_2 \text{Id}$  with  $\varepsilon_2 = (1 - 6r_{\text{near}}^2)(1 - r_{\text{near}}^2/(2 - r_{\text{near}}^2) - r_{\text{near}}^2) \geq 0.941$ . Finally, from Lemma F.2, we have that for for all  $d_{\mathbf{H}}(x, x') \geq r_{\text{near}}$ ,  $|K| \leq 1 - 1/(8^3 \cdot 2)$ , so we can set  $\varepsilon_0 \stackrel{\text{def.}}{=} 0.00097$ .

Before proving these lemmas, we first summarise in Section F.1.3 some key properties of the univariate Fejér kernel  $\kappa$  when  $f \ge 128$  which were derived in Candès and Fernandez-Granda (2014).

For notational convenience, write  $t_i \stackrel{\text{def.}}{=} x_i - x'_i$ ,  $\kappa_i \stackrel{\text{def.}}{=} \kappa(t_i)$ ,  $\kappa'_i \stackrel{\text{def.}}{=} \kappa'(t_i)$ , and so on. Let

$$K_i \stackrel{\text{def.}}{=} \prod_{\substack{k=1\\k\neq i}}^d \kappa_k, \quad K_{ij} \stackrel{\text{def.}}{=} \prod_{\substack{k=1\\k\neq i,j}}^d \kappa_k \quad \text{and} \qquad K_{ij\ell} \stackrel{\text{def.}}{=} \prod_{\substack{k=1\\k\neq i,j,\ell}}^d \kappa_k.$$

With this, we have:

$$\partial_{1,i}K(x,x') = \kappa'_i K_i$$
  
$$\partial_{1,i}\partial_{2,i}K(x,x') = -\kappa'_i K_i, \quad \text{and} \quad \forall i \neq j, \ \partial_{1,i}\partial_{2,j}K(x,x') = -\kappa'_i \kappa'_j K_{ij}.$$

Where convenient, we sometimes write  $K(t) = K(x - x') \stackrel{\text{def.}}{=} K(x, x')$ .

### **F.1.3** Properties of $\kappa$

From (Candès and Fernandez-Granda, 2014, Equations (2.20)-(2.24) and (2.29)), for all  $t \in [-1/2, 1/2]$  and  $\ell = 0, 1, 2, 3$ :

$$1 - \frac{C_f}{2}t^2 \leqslant \kappa(t) \leqslant 1 - \frac{C_f}{2}t^2 + 8\left(\frac{1+2/f}{1+2/(2+f)}\right)^2 C_f^2 t^4 \leqslant 1 - \frac{C_f}{2}t^2 + 8C_f^2 t^4$$
$$|\kappa'(t)| \leqslant C_f t, \quad |\kappa''(t)| \leqslant C_f, \quad |\kappa'''(t)| \leqslant 3\left(\frac{1+2/f}{1+2/(2+f)}\right)^2 C_f^2 t \leqslant 12C_f^2 t \qquad (F.1)$$
$$\kappa'' \leqslant -C_f + \frac{3}{2}\left(\frac{1+2/f}{1+2/(2+f)}\right)^2 C_f^2 t^2 \leqslant -C_f + 6C_f^2 t^2.$$

By (Candès and Fernandez-Granda, 2014, Lemma 2.6),

$$\left|\kappa^{(\ell)}(t)\right| \leqslant \begin{cases} \frac{\pi^{\ell} H_{\ell}(t)}{(f+2)^{4-\ell}t^{4}}, & t \in [\frac{1}{2f}, \frac{\sqrt{2}}{\pi}]\\ \frac{\pi^{\ell} H_{\ell}^{\infty}}{(f+2)^{4-\ell}t^{4}}, & t \in [\frac{\sqrt{2}}{\pi}, \frac{1}{2}), \end{cases}$$

where  $H_0^{\infty} \stackrel{\text{def.}}{=} 1$ ,  $H_1^{\infty} \stackrel{\text{def.}}{=} 4$ ,  $H_2^{\infty} \stackrel{\text{def.}}{=} 18$  and  $H_3^{\infty} \stackrel{\text{def.}}{=} 77$ , and  $H_{\ell}(t) \stackrel{\text{def.}}{=} \alpha^4(t)\beta_{\ell}(t)$ , with

$$\alpha(t) \stackrel{\text{def.}}{=} \frac{2}{\pi(1 - \frac{\pi^2 t^2}{6})}, \quad \bar{\beta}(t) \stackrel{\text{def.}}{=} \frac{\alpha(t)}{ft} = \frac{2}{ft\pi(1 - \pi^2 t^2/6)}$$

and  $\beta_0(t) \stackrel{\text{def.}}{=} 1$ ,  $\beta_1(t) \stackrel{\text{def.}}{=} 2 + 2\bar{\beta}(t)$ ,  $\beta_2 \stackrel{\text{def.}}{=} 4 + 7\bar{\beta}(t) + 6\bar{\beta}(t)^2$  and  $\beta_3(t) \stackrel{\text{def.}}{=} 8 + 24\bar{\beta} + 30\bar{\beta}(t)^2 + 15\bar{\beta}(t)^3$ . Let us first remark that  $\bar{\beta}$  is decreasing on  $I \stackrel{\text{def.}}{=} [\frac{1}{2f}, \frac{\sqrt{2}}{\pi}]$ , so  $|\bar{\beta}(t)| \leq |\bar{\beta}(1/(2f))| \approx 1.2733$ , and  $a(t) \leq a(\sqrt{2}/\pi) = \frac{3}{\pi}$  on I. Therefore, on I,  $H_0(t) \leq \frac{3}{\pi}$ ,  $H_1(t) \leq 3.79$ ,  $H_2(t) \leq 18.83$  and  $H_3(t) \leq 98.26$ , and we can conclude that on  $[\frac{1}{2f}, \frac{1}{2}]$ , we have

$$\left|\kappa^{(\ell)}(t)\right| \leqslant \frac{\pi^{\ell} H_{\ell}^{\infty}}{(f+2)^{4-\ell} t^4}$$

where  $\bar{H}_0^{\infty} = 1$ ,  $\bar{H}_1^{\infty} \stackrel{\text{def.}}{=} 4$ ,  $\bar{H}_2^{\infty} \stackrel{\text{def.}}{=} 19$ ,  $\bar{H}_3^{\infty} \stackrel{\text{def.}}{=} 99$ . Combining with (F.1), we have  $\|\kappa^{(\ell)}\|_{\infty} \leq \kappa_{\ell}^{\infty}$  where  $\kappa_0^{\infty} \stackrel{\text{def.}}{=} 1$ ,  $\kappa_2^{\infty} \stackrel{\text{def.}}{=} C_f$ ,

$$\kappa_1^{\infty} \stackrel{\text{def.}}{=} \sqrt{C_f} \max\left(\frac{2\pi^4}{(\frac{1}{2} + \frac{1}{f})^3} \frac{f}{\sqrt{C_f}}, \frac{\sqrt{C_f}}{2f}\right) = \mathcal{O}(\sqrt{C_f})$$
$$\kappa_3^{\infty} \stackrel{\text{def.}}{=} (C_f)^{3/2} \max\left(\frac{99\pi^3}{(\frac{1}{2} + \frac{1}{f})} \left(\frac{2f}{\sqrt{C_f}}\right)^4, \frac{6\sqrt{C_f}}{f}\right) = \mathcal{O}((C_f)^{3/2}).$$

Finally, given  $p \in (0, 1)$ ,

$$(f+2)^4 t^4 \ge (1+p(f+2)^2 t^2)^2, \qquad \forall t \ge \frac{1}{\sqrt{(1-p)}(f+2)}.$$

Choosing  $p = \frac{1}{2}$  and using  $(f+2)^2 = (\frac{3}{\pi^2}C_f + 4) \ge \frac{3}{\pi^2}C_f$ , we have

$$\left|\kappa^{(\ell)}(t)\right| \leqslant \frac{\kappa_{\ell}^{\infty}}{(1+\frac{3}{2\pi^2}C_f t^2)^2}, \qquad \forall \ t^2 \geqslant \frac{2\pi^2}{3C_f}, \tag{F.2}$$

## **F.1.4** Bounds in neighbourhood of x' = x

**Lemma F.1.** Suppose that  $C_f ||t||_2^2 \leq c$  with c > 0 such that

$$\varepsilon \stackrel{\text{def.}}{=} (1 - 6c) \left( 1 - \frac{c}{2 - c} \right) - c > 0$$

Then,  $\hat{K}^{02}(t) \preceq -\varepsilon \mathrm{Id}.$ 

*Proof.* We need to show that  $\lambda_{\min}(-K^{(02)}(t)) \ge b$ . Let  $q \in \mathbb{R}^d$ , and note that

$$-\langle \nabla_2^2 Kq, q \rangle = -\sum_i \left( q_i \kappa_i'' K_i - \kappa_i' \sum_{j \neq i} q_j \kappa_j' K_{ij} \right) q_i$$
  
$$= -\left( \sum_i q_i^2 \kappa_i'' K_i - \sum_i q_i \kappa_i \sum_{j \neq i} q_j \kappa_j K_{ij} \right)$$
  
$$\geqslant \|q\|^2 \left( -\max_i \{\kappa_i'' K_i\} - \sum_j |\kappa_j'|^2 \right).$$
 (F.3)

We first consider  $\kappa_i'' K_i$ :

$$\kappa_i'' \leqslant -C_f + 6C_f^2 t_i^2,$$

$$K_i \geqslant \prod_{j \neq i} \left( 1 - \frac{C_f}{2} t_i^2 \right) \geqslant 1 - \frac{C_f}{2} \|t\|_2^2 - \left( \frac{C_f}{2} \|t\|_2^2 \right)^3 - \left( \frac{C_f}{2} \|t\|_2^2 \right)^5 - \cdots$$

$$\geqslant 1 - \frac{C_f \|t\|_2^2}{2(1 - \frac{C_f}{2} \|t\|_2^2)}.$$

and hence,

$$\kappa_i'' K_i \leqslant \left( -C_f + 6C_f^2 \left\| t \right\|_2^2 \right) \left( 1 - \frac{C_f \left\| t \right\|_2^2}{2(1 - \frac{C_f}{2} \left\| t \right\|_2^2)} \right)$$

For the second term,

$$\sum_{j} \left| \kappa_{j}^{\prime} \right|^{2} \leqslant C_{f}^{2} \left\| t \right\|_{2}^{2}.$$

Therefore,

$$\lambda_{\min}(-K^{(02)}(t)) \ge \left(1 - 6C_f \|t\|_2^2\right) \left(1 - \frac{C_f \|t\|_2^2}{2(1 - \frac{C_f}{2} \|t\|_2^2)}\right) - C_f \|t\|_2^2$$

**Lemma F.2.** Assume that  $\frac{1}{8\sqrt{C_f}} \ge ||t||_2$  Then,

$$K(t) \leq 1 - \frac{C_f}{4} \|t\|_2^2 + 16C_f^2 \|t\|_2^4$$

Consequently, for all

$$0 < c \leqslant \frac{1}{8\sqrt{2C_f}},$$

and all t such that  $\|t\|_2 \ge c$ ,

$$|K(t)| \leqslant 1 - \frac{C_f}{8}c^2.$$

*Proof.* First note that

$$|\kappa(u)| \leq 1 - \frac{C_f}{2}u^2 + 32C_f^2u^4 = 1 - u^2g(u)$$

where

$$g(u) \stackrel{\text{def.}}{=} C_f\left(\frac{1}{2} - 32C_f u^2\right),$$

and note that  $g(u) \in (0, \frac{C_f}{2})$  for  $u \in (0, 1/(8\sqrt{C_f}))$ . So, writing  $t = (t_i)_{i=1}^d$  and  $g_j \stackrel{\text{def.}}{=} g(t_j)$ , we have

$$K(t) = \prod_{j=1}^{d} \kappa(t_i) \leq \prod_{j=1}^{d} \left(1 - t_j^2 \cdot g(t_j)\right)$$
  
=  $1 - \sum_{j=1}^{d} t_j^2 g_j + \sum_{j \neq k} t_j^2 t_k^2 g_j g_k - \sum_{j \neq k \neq \ell} t_j^2 t_k^2 t_\ell^2 g_j g_k g_\ell + \cdots$ 

Note that

$$\begin{split} &-\sum_{j\neq k\neq \ell} t_j^2 t_k^2 t_\ell^2 \cdot g_j g_k g_\ell + \sum_{j\neq k\neq \ell\neq n} t_j^2 t_k^2 t_\ell^2 t_n^2 \cdot g_j g_k g_\ell g_n \\ &\leqslant -\sum_{j\neq k\neq \ell} t_j^2 t_k^2 t_\ell^2 \cdot g_j g_k g_\ell + \left(\sum_{j\neq k\neq \ell} t_j^2 t_k^2 t_\ell^2 \cdot g_j g_k g_\ell\right) \left(\sum_n t_n^2 g_n\right) \\ &\leqslant -\sum_{j\neq k\neq \ell} t_j^2 t_k^2 t_\ell^2 \cdot g_j g_k g_\ell \left(1 - \frac{C_f}{2} \left\|t\right\|_2^2\right) < 0 \end{split}$$

since  $\left(1 - \frac{C_f}{2} \|t\|_2^2\right) > 0$ . Also,

$$\sum_{j=1}^{d} t_j^2 g_j \leqslant \frac{C_f}{2} \sum_{j=1}^{d} t_j^2 < 1,$$

by assumption. So,

$$\begin{split} K(t) &\leqslant 1 - \sum_{j=1}^{d} t_{j}^{2} g_{j} + \sum_{j \neq k} t_{j}^{2} t_{k}^{2} g_{j} g_{k} \\ &\leqslant 1 - \sum_{j=1}^{d} t_{j}^{2} g_{j} + \frac{1}{2} \left( \sum_{j} t_{j}^{2} g_{j} \right)^{2} \leqslant 1 - \frac{1}{2} \sum_{j=1}^{d} t_{j}^{2} g_{j} \\ &\leqslant 1 - \frac{C_{f}}{2} \left( \frac{1}{2} \sum_{j=1}^{d} t_{j}^{2} - 32C_{f} \sum_{j=1}^{d} t_{j}^{4} \right) \leqslant 1 - \frac{C_{f}}{4} \|t\|_{2}^{2} + 16C_{f}^{2} \|t\|_{2}^{4}. \end{split}$$

Finally, observe that the function

$$q(z) \stackrel{\text{def.}}{=} \frac{C_f}{4} z^2 - 16C_f^2 z^4$$

is positive and increasing on the interval  $[0, \frac{1}{8\sqrt{2C_f}}]$ . So, for t satisfing

$$c \leqslant \|t\|_2 \leqslant \frac{1}{8\sqrt{2C_f}},\tag{F.4}$$

we have  $|K(t)| \leq 1 - q(c) \leq 1 - \frac{C_f}{8}c^2$ . Finally, since |K(t)| is decreasing as t increases, we trivially have that  $|K(t)| \leq 1 - q(c)$  for all t with  $||t||_2 \geq c$ .

### F.1.5 Bounds under separation

**Lemma F.3.** Let  $i, j \in \{0, 1, 2\}$  with  $i + j \leq 3$ . Let  $\bar{A} \ge \sqrt{\frac{4\pi^2}{3}}$  and  $||t||_2 \ge \bar{A}\sqrt{ds_{\max}^{1/4}}/\sqrt{C_f}$ . Then, we have  $||K^{(ij)}(t)|| \le d^{\frac{i+j-4}{2}}(\bar{A}^4s_{\max})^{-1}$ .

Proof. Write  $t = (t_j)_{j=1}^d$ . To bound  $K(t) = \prod_{j=1}^d \kappa(a_j)$ , we want to make use of the form (F.2). We can do this for each  $t_j$  such that  $|t_j| \ge \sqrt{\frac{2\pi^2}{3C_f}}$ . Note that there exists at least one such  $t_j$  since  $||t||_{\infty} \ge ||t||_2 / \sqrt{d} \ge \overline{As_{\max}^{1/4}} / \sqrt{C_f} \ge \sqrt{\frac{2\pi^2}{3C_f}}$ . If  $\{|t_j|\}_{j=1}^k \subset [0, \sqrt{\frac{2\pi^2}{3C_f}})$  for  $k \le d-1$ , then

$$k\frac{2\pi^2}{3C_f} + \sum_{j=k+1}^d t_j^2 \ge ||t||_2^2 \ge \frac{\bar{A}^2 ds_{\max}^{1/2}}{C_f},$$

which implies that  $\sum_{j=k+1}^{d} t_j^2 \ge \frac{1}{C_f} \left( \bar{A}^2 ds_{\max}^{1/2} - \frac{2\pi^2 (d-1)}{3} \right) \ge \frac{\bar{A}^2 ds_{\max}^{1/2}}{2C_f}$ , by our assumptions on  $\bar{A}$ . Therefore, we may assume that we have some  $d \ge p \ge 1$  such that  $\{b_j\}_{j=1}^p \subseteq \{t_j\}$  with  $|b_j| \ge \sqrt{\frac{2\pi^2}{3C_f}}$  and  $||b||_2 \ge \frac{\bar{A}\sqrt{d}\sqrt{4s_{\max}}}{\sqrt{2C_f}}$ . Observe that

$$\prod_{j=1}^{p} \left(1 + \frac{3C_f}{2\pi^2} b_j^2\right) \ge 1 + \frac{3C_f}{2\pi^2} \sum_{j=1}^{p} b_j^2 = 1 + \frac{3C_f}{2\pi^2} \left\|b\right\|_2^2 \ge 1 + \frac{3}{4\pi^2} \bar{A}^2 d\sqrt{s_{\max}}.$$

So, by applying the fact that  $|\kappa| \leq 1$ ,  $\kappa_0^{\infty} = 1$  and (F.2), we have

$$|K(t)| \leqslant \prod_{j=1}^{p} |\kappa(b_j)| \leqslant \prod_{j=1}^{p} \frac{1}{\left(1 + \frac{3C_f}{2\pi^2} b_j^2\right)^2} \leqslant \frac{1}{\left(1 + \frac{3}{4\pi^2} \bar{A}^2 d\sqrt{s_{\max}}\right)^2}$$

For  $|\kappa'_i K_i|$ , if  $i \notin \left\{ j \ ; \ |t_j| > \sqrt{\frac{2\pi^2}{3C_f}} \right\}$ , then

$$|\kappa_i' K_i| \leq \|\kappa_i'\|_{\infty} \prod_{j=1}^p |\kappa(b_j)| \leq \frac{\|\kappa_i'\|_{\infty}}{\left(1 + \frac{3}{4\pi^2} \bar{A}^2 d\sqrt{s_{\max}}\right)^2}$$

and otherwise, we have  $|\kappa'_i K_i| \leq |\kappa'(t_i)| \prod_{j \neq i} |\kappa(b_j)| \leq \frac{\kappa_1^{\infty}}{\left(1 + \frac{3}{4\pi^2} \bar{A}^2 d\sqrt{s_{\max}}\right)^2}$ . In a similar manner, writing  $V \stackrel{\text{def.}}{=} \left(1 + \frac{3}{4\pi^2} \bar{A}^2 d\sqrt{s_{\max}}\right)^{-2}$ , we can deduce that

$$\begin{aligned} |\kappa_i'K_i| &\leqslant \kappa_1^{\max}V, \qquad |\kappa_i''K_i| &\leqslant \kappa_2^{\max}V, \qquad \left|\kappa_i'\kappa_j'K_{ij}\right|^2 &\leqslant (\kappa_1^{\max})^2V\\ |\kappa_i'''K_i|^3 &\leqslant \kappa_3^{\max}V, \qquad \left|\kappa_i''\kappa_j'K_{ij}\right|^3 &\leqslant \kappa_2^{\max}\kappa_1^{\max}V, \qquad \left|\kappa_i'\kappa_j'\kappa_\ell'K_{ij\ell}\right| &\leqslant (\kappa_1^{\max})^3V \end{aligned}$$

Therefore,

$$\left\|K^{(10)}\right\| = \frac{1}{\sqrt{C_f}} \left\|\nabla_1 K\right\| \leqslant \frac{1}{\sqrt{C_f}} \sqrt{\sum_{j=1}^d \left|\kappa_j' K_j\right|^2} \leqslant \frac{\kappa_1^\infty}{\sqrt{C_f}} V\sqrt{d} \lesssim \frac{1}{\bar{A}^4 d^{3/2} s_{\max}}$$

Using Gershgorin theorem, we have

$$\left\|\nabla_{2}^{2}K(x,x')\right\| \leqslant \max_{1\leqslant i\leqslant d} \left\{\left|\kappa_{i}''K_{i}\right| + \left|\kappa_{i}'\right|\sum_{j\neq i}\left|\kappa_{j}'\right|\left|K_{ij}\right|\right\}$$

and hence,

$$\begin{aligned} \left\| K^{(02)} \right\| &= \frac{1}{C_f} \left\| \nabla_2^2 K \right\| \leqslant \frac{1}{C_f} \max_{i=1}^d \{ |\kappa_i'' K_i| + |\kappa_i'| \sum_{j \neq i} \left| \kappa_j' K_{ij} \right| \} \\ &\leqslant \frac{1}{C_f} V \left( \kappa_2^{\max} + (\kappa_1^{\max})^2 (d-1) \right) \leqslant \frac{\max\{\kappa_2^{\infty}, (\kappa_1^{\infty})^2\}}{C_f} V d \lesssim \frac{1}{\bar{A}^4 ds_{\max}}. \end{aligned}$$

Note also that  $\left\|K^{(11)}\right\| = \left\|K^{(02)}\right\|$ . Finally, since

$$\begin{aligned} \left\| \partial_{1,i} \nabla_{2}^{2} K(x,x') \right\| \leqslant \max \left\{ \left| \kappa_{i}^{\prime \prime \prime} K_{i} \right| + \left| \kappa_{i}^{\prime \prime} \right| \sum_{j \neq i} \left| \kappa_{j}^{\prime} \right| \left| K_{ij} \right|, \\ \max_{j \neq i} \left\{ \left| \kappa_{j}^{\prime \prime} \kappa_{i}^{\prime} K_{ij} \right| + \left| \kappa_{j}^{\prime} \kappa_{i}^{\prime \prime} K_{ij} \right| + \left| \kappa_{j}^{\prime} \right| \sum_{l \neq i,j} \left| \kappa_{l}^{\prime} \right| \left| K_{ij\ell} \right| \right\} \right\}, \end{aligned}$$

we have

$$\begin{split} \left\| K^{(12)} \right\| &= \frac{1}{C_f^{3/2}} \left\| \nabla_1 \nabla_2^2 K \right\| \\ &\leqslant \frac{1}{C_f^{3/2}} \sqrt{d} V \max\left( \kappa_3^{\max} + \kappa_2^{\max} \kappa_1^{\max} (d-1), 2\kappa_2^{\max} \kappa_1^{\infty} + (d-1)(\kappa_1^{\infty})^3 \right) \\ &\leqslant d^{3/2} \max\{ \kappa_3^{\infty}, \kappa_1^{\infty} \kappa_2^{\infty}, (\kappa_1^{\infty})^3 \} \frac{1}{C_f^{3/2}} V \lesssim \frac{1}{\bar{A}^4 d^{1/2} s_{\max}} \end{split}$$

_

### F.1.6 Uniform bounds

**Lemma F.4.** If  $r_{\text{near}} \sim 1/\sqrt{C_f}$ , then  $B_0 = \mathcal{O}(1)$ ,  $B_{01} = \mathcal{O}(\sqrt{d})$ ,  $B_{02} = B_{12} = B_{11} = \mathcal{O}(1)$  and  $B_{22} = \mathcal{O}(d)$ .

*Proof.* We have  $|K| \leq 1$ , and

$$\left\|\nabla K\right\|^{2} \leqslant \sum_{i} |\kappa_{i}|^{2} |K_{i}|^{2} \leqslant d(\kappa_{1}^{\infty})^{2} \lesssim C_{f} d,$$

so  $B_{01} = \mathcal{O}(\sqrt{d})$ . From (F.3), for all ||q|| = 1,

$$\langle \nabla_2^2 K(t)q, q \rangle \leq \max_i |\kappa_i''| ||q||_2^2 + ||q||_2^2 \sum_i |\kappa_i|^2 \leq C_f + C_f^2 ||t||^2 = \mathcal{O}(C_f),$$

for  $||t|| \leq 1/\sqrt{C_f}$ . So, since  $r_{\text{near}} \leq 2/\sqrt{C_f}$ ,  $||K^{02}(t)|| \leq 2 \stackrel{\text{def.}}{=} B_{02}$ . The norm bound for  $K^{11}$  is the same.

$$\begin{split} \left\| K^{(12)} \right\| &= \sup_{\|q\| = \|p\| = 1} \frac{1}{C_f^{3/2}} \left( \sum_k \sum_{k \neq i} \partial_{1,i} \left( \partial_{2,k}^2 K p_i q_k^2 + \partial_{1,i} \partial_{2,i} \partial_{2,k} K p_i q_i q_k \right) \\ &+ \sum_i \sum_k \sum_j \partial_{1,i} \partial_{2,j} \partial_{2,k} p_i p_j p_k + \sum_i \sum_{j \neq i} \partial_{1,i} \partial_{2,i} \partial_{2,j} K p_i q_i q_j + \sum_i \partial_{1,i} \partial_{2,j}^2 K p_i q_i^2 \right) \\ &= \sup_{\|q\| = \|p\| = 1} \frac{1}{C_f^{3/2}} \left( \sum_k \sum_{k \neq i} \kappa'_i \kappa''_k K_{ik} p_i q_k^2 + \kappa''_i \kappa'_k K_{ik} p_i q_i q_k \right) \\ &+ \sum_i \sum_k \sum_j \kappa'_i \kappa'_k \kappa'_j K_{ijk} p_i p_j p_k + \sum_i \sum_{j \neq i} \kappa''_i \kappa''_j K_{ij} p_i q_i q_j + \sum_i \kappa'_i \kappa''_j K_{ij} p_i q_i^2 \right) \\ &\leqslant \frac{1}{C_f^{3/2}} \left( 3 \|\kappa''\|_{\infty} \sqrt{\sum_i |\kappa'_k|^2} + \left( \sum_i |\kappa'_k|^2 \right)^{3/2} + \|\kappa'\|_{\infty} \|\kappa''\|_{\infty} \right) \\ &\leqslant \frac{1}{C_f^{3/2}} \left( 3 C_f^2 \|t\| + C_f^3 \|t\|^3 + \mathcal{O}(C_f^{3/2}) \right) = \mathcal{O}(1) \end{split}$$

for  $||t|| \leq 1/C_f^{1/2}$ .

We finally consider  $K^{(22)}(x, x)$ : for ||p|| = 1,

$$\begin{split} \sum_{i} \sum_{k} \sum_{j} \partial_{1,k} \partial_{1,i} \partial_{2,j} \partial_{2,i} K p_{j} p_{k} &= \sum_{i} \sum_{k \neq i} \kappa_{i}^{\prime\prime} \kappa_{k}^{\prime\prime} p_{j}^{2} K_{ik} + \sum_{i} \sum_{k \neq i} \kappa_{i}^{\prime\prime\prime} \kappa_{k}^{\prime} p_{i} p_{k} K_{ik} \\ &+ \sum_{i} \sum_{k} \sum_{j} \kappa_{i}^{\prime\prime} \kappa_{j}^{\prime} \kappa_{j}^{\prime} \kappa_{k}^{\prime} K_{ijk} p_{j} p_{k} + \sum_{i} \sum_{j} \kappa_{i}^{\prime\prime\prime\prime} \kappa_{j}^{\prime} p_{j} p_{i} K_{ij} + \sum_{i} \kappa_{i}^{\prime\prime\prime\prime\prime} p_{i}^{2} K_{i} \\ &= \sum_{i} \sum_{k \neq i} \kappa_{i}^{\prime\prime} \kappa_{k}^{\prime\prime} p_{j}^{2} K_{ik} + \sum_{i} \kappa_{i}^{\prime\prime\prime\prime\prime} p_{i}^{2} \\ &= d\mathcal{O}(C_{f}^{2}) \end{split}$$

since  $\kappa'(0) = \kappa'''(0) = 0$  and  $|\kappa''(0)| = \mathcal{O}(C_f), |\kappa''''(0)| = \mathcal{O}(C_f^2)$ . So,  $B_{22} = \mathcal{O}(d)$ .

### F.2 The Gaussian kernel

We consider the Gaussian kernel  $K(x, x') = \exp\left(-\frac{1}{2} \|x - x'\|_{\Sigma^{-1}}^2\right)$  in  $\mathbb{R}^d$ . Note that K is translation invariant, so that  $\mathbf{H}_x$  will be constant and equal to  $-\nabla^2 K(x, x)$ . For simplicity define t = x - x',  $\hat{K}_{\Sigma}(t) = \exp\left(-\frac{1}{2} \|t\|_{\Sigma^{-1}}^2\right)$  and for  $u \in \mathbb{R}$ ,  $\kappa(u) = \exp\left(-\frac{1}{2}u^2\right)$ . Denote by  $\{e_i\}$  the canonical basis of  $\mathbb{R}^d$ , and by  $f_i = \Sigma^{-1}e_i$  the  $i^{th}$  row of  $\Sigma^{-1}$ . We have the following:

$$\begin{aligned} \nabla \hat{K}_{\Sigma}(t) &= -\Sigma^{-1} t \hat{K}_{\Sigma}(t) \\ \nabla^2 \hat{K}_{\Sigma}(t) &= \left( -\Sigma^{-1} + \Sigma^{-1} t t^{\top} \Sigma^{-1} \right) \hat{K}_{\Sigma}(t) \\ \partial_{1,i} \nabla^2 \hat{K}_{\Sigma}(t) &= \left( \Sigma^{-1} t f_i^{\top} + f_i t^{\top} \Sigma^{-1} - (-\Sigma^{-1} + \Sigma^{-1} t t^{\top} \Sigma^{-1}) (t^{\top} f_i) \right) \hat{K}_{\Sigma}(t) \end{aligned}$$

Hence we have  $\mathbf{H}_x = -\nabla^2 \hat{K}_{\Sigma}(0) = \Sigma^{-1}$ , and, defining  $d_{\mathbf{H}}(x, x') = ||x - x'||_{\Sigma^{-1}} = ||\Sigma^{-\frac{1}{2}}(x - x')||$ , we have  $C_{\hat{K}} = 1, C_{\mathbf{H}} = 0$  (that is, the metric tensor of the kernel is constant, and  $d_{\mathbf{H}}$  is defined as the corresponding normalized norm).

Then, we have

$$\begin{aligned} \left\| K^{(10)}(x,x') \right\| &= \left\| K^{(01)}(x,x') \right\| = d_{\mathbf{H}}(x,x')\kappa(d_{\mathbf{H}}(x,x')) \\ \left\| K^{(02)}(x,x') \right\| &= \left\| K^{(11)}(x,x') \right\| \leqslant (d_{\mathbf{H}}(x,x')^2 + 1)\kappa(d_{\mathbf{H}}(x,x')) \\ K^{(02)}(x,x') &\leqslant (d_{\mathbf{H}}(x,x')^2 - 1)\kappa(d_{\mathbf{H}}(x,x')) \end{aligned}$$

and for  $q \in \mathbb{R}^d$  with ||q|| = 1, since

$$\sum_{i} (\Sigma^{\frac{1}{2}} \nabla \varphi_{\omega})_{i} q_{i} = \nabla \varphi_{\omega}^{\top} (\Sigma^{\frac{1}{2}} q) = \sum_{i} \partial_{i} \varphi_{\omega} (q^{\top} \Sigma^{\frac{1}{2}} e_{i})$$

we can write

$$K^{(12)}(x,x')q = \sum_{i=1}^{d} (q^{\top} \Sigma^{\frac{1}{2}} e_i) \Sigma^{\frac{1}{2}} \partial_{1,i} \nabla^2 \hat{K}_{\Sigma}(t) \Sigma^{\frac{1}{2}}$$

Thus we examine each term in  $\partial_{1,i} \nabla^2 \hat{K}_{\Sigma}$ . We have

$$\sum_{i} (q^{\top} \Sigma^{\frac{1}{2}} e_{i}) \Sigma^{\frac{1}{2}} \Sigma^{-1} t f_{i}^{\top} \Sigma^{\frac{1}{2}} = \Sigma^{-\frac{1}{2}} t \left( \sum_{i} q^{\top} \Sigma^{\frac{1}{2}} e_{i} e_{i}^{\top} \Sigma^{-\frac{1}{2}} \right) = \Sigma^{-\frac{1}{2}} t q^{\top}$$

and similarly  $\sum_i (q^\top \Sigma^{\frac{1}{2}} e_i) \Sigma^{\frac{1}{2}} f_i t^\top \Sigma^{-1} \Sigma^{\frac{1}{2}} = q t^\top \Sigma^{\frac{1}{2}}.$  Then

$$\sum_{i} (q^{\top} \Sigma^{\frac{1}{2}} e_{i}) (t^{\top} \Sigma^{-1} e_{i}) \Sigma^{\frac{1}{2}} \Sigma^{-1} \Sigma^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} q = (t^{\top} \Sigma^{\frac{1}{2}} q) \operatorname{Id}_{i} U^{\top} (t^{\top} \Sigma^{-1} e_{i}) \Sigma^{\frac{1}{2}} Z^{-1} \Sigma^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} = t^{\top} \Sigma^{-1} (\sum_{i} e_{i} e_{i}^{\top}) \Sigma^{\frac{1}{2}} Q^{-1} Z^{\frac{1}{2}} Q$$

and similarly  $\sum_i \sum_i (q^\top \Sigma^{\frac{1}{2}} e_i)(t^\top \Sigma^{-1} e_i) \Sigma^{\frac{1}{2}} \Sigma^{-1} t t^\top \Sigma^{-1} \Sigma^{\frac{1}{2}} = (t^\top \Sigma^{\frac{1}{2}} q) \Sigma^{-\frac{1}{2}} t t^\top \Sigma^{-\frac{1}{2}}$ . Hence at the end of the day

$$\left\| K^{(12)}(x,x') \right\| \le (3d_{\mathbf{H}}(x,x') + d_{\mathbf{H}}(x,x')^3)\kappa(d_{\mathbf{H}}(x,x'))$$

and this bound is automatically valid for  $K^{(21)}$  as well.

Finally, note that

$$\left\| K^{(22)}(x,x) \right\| = \sup_{\|p\| \leq 1} \langle \Sigma^{1/2} \nabla_2 \nabla_2 \cdot \left( \Sigma^{1/2} K^{(2,0)}(x,x) p \right), p \rangle$$

where  $\nabla_2$  is the divergence operator on the 2nd variable, and one can show that  $\|K^{(22)}(x,x)\| = (d+1)$ .

We are then going to use the fact that for any  $q \ge 1$  the function  $f(r) = r^q e^{-\frac{1}{2}r^2}$  defined on  $\mathbb{R}_+$  is increasing on  $[0, \sqrt{q}]$  and decreasing after, and its maximum value is  $f(\sqrt{q}) = \left(\frac{q}{e}\right)^{q/2}$ . Furthermore, it is easy to see that we have  $f(r) = r^q e^{-r^2/2} \le \left(\frac{2q}{2}\right)^{\frac{q}{2}} e^{-r^2/4}$  and therefore  $f(r) \le \varepsilon$  if  $r \ge 2 \left(\log\left(\frac{1}{\varepsilon}\right) + \frac{q}{2}\log\left(\frac{2q}{e}\right)\right)$ .

We define  $r_{\text{near}} = 1/\sqrt{2}$  and  $\Delta = C_1 \sqrt{\log(s_{\text{max}})} + C_2$  for some  $C_1$  and  $C_2$ .

1. Global Bounds. From what preceeds, we have

$$\left\|K^{(10)}\right\| \leqslant \frac{1}{\sqrt{e}}, \quad \left\|K^{(02)}\right\| \leqslant \frac{2}{e} + 1, \quad \left\|K^{(12)}\right\| \leqslant \frac{3}{\sqrt{e}} + \left(\frac{3}{e}\right)^{\frac{3}{2}}$$

and note that  $||K^{(11)}|| = ||K^{(02)}||$ , so for all  $i + j \leq 3$   $B_{ij} = \mathcal{O}(1)$ .

2. Near 0 For  $d_{\mathbf{H}}(x, x') \leq r_{\text{near}}$ , we have

$$K^{(02)} \preccurlyeq -\frac{e^{-\frac{1}{4}}}{2} \mathrm{Id}$$

and for  $d_{\mathbf{H}}(x, x') \ge \frac{1}{2}$ ,

$$|K| \leqslant e^{-\frac{1}{4}} = 1 - (1 - e^{-\frac{1}{4}})$$

and  $||K^{(22)}(x,x)|| = d + 1$ , so we have also  $\varepsilon_i = \mathcal{O}(1)$ , so  $B_i = B_{0i} + B_{1i} + 1 = \mathcal{O}(1)$  and  $B_{22} = d + 1$ .

3. Separation. Since  $\varepsilon_i = \mathcal{O}(1)$  and  $B_{ij} = \mathcal{O}(1)$ , every condition  $\|K^{(ij)}\| \lesssim \frac{1}{s_{\max}}$  is satisfied if  $\Delta \ge C_1 \sqrt{\log(s_{\max})} + C_2$  for some constant  $C_1$  and  $C_2$ .

### F.2.1 Fourier measurements with Gaussian frequencies

The random feature expansion for K is  $\varphi_{\omega}(x) = e^{i\omega^{\top}x}$  and  $\Lambda = \mathcal{N}(0, \Sigma^{-1})$ . We have immediately  $L_0 = 1$ . For  $j \ge 1$ , we have  $D_j [\varphi_{\omega}](x)[q_1, \ldots, q_j] = \left(\prod_i \omega^{\top}(\Sigma^{\frac{1}{2}}q_i)\right) \varphi_{\omega}(x)$  and therefore

$$\|\mathbf{D}_{j} [\varphi_{\omega}]\| \leqslant \|\omega\|_{\Sigma}^{j}$$

Now, we use  $\|\omega\|_{\Sigma}^{j} = (\|\Sigma^{\frac{1}{2}}\omega\|^{2})^{\frac{j}{2}} = W^{\frac{j}{2}}$  where W is a  $\chi^{2}$  variable with d degrees of freedom. Then, we use the following Chernoff bound (Dasgupta and Gupta, 2003): for  $x \ge d$ , we have

$$\mathbb{P}(W \ge x) \leqslant \left(\frac{ex}{d}e^{-\frac{x}{d}}\right)^{\frac{d}{2}} \leqslant \left(e\left(\sqrt{\frac{x}{d}}\right)^2 e^{-\frac{1}{2}\cdot\left(\sqrt{\frac{x}{d}}\right)^2} e^{-\frac{x}{2d}}\right)^{\frac{d}{2}} \leqslant 2^{\frac{d}{2}}e^{-\frac{x}{4}}$$

by using  $x^2 e^{-\frac{x^2}{2}} \leq \frac{2}{e}$ .

Hence we can define the  $F_j$  such that, for all  $t \ge d^{j/2}$ ,  $\mathbb{P}(L_j(\omega) \ge t) \le F_j(t) = 2^{\frac{d}{2}} \exp\left(-\frac{t^2_j}{4}\right)$ , and  $F_j(\bar{L}_j)$  is smaller than some  $\delta$  if  $\bar{L}_j \propto \left(d + \log \frac{1}{\delta}\right)^{\frac{j}{2}}$ . Then we must choose the  $L_j$  such that  $\int_{\bar{L}_j} tF_j(t)dt$  is bounded by some  $\delta$ . Taking  $L_j \ge d^{j/2}$  in any case, we have

$$\begin{split} \int_{\bar{L}_j} tF_j(t) \mathrm{d}t &= 2^{\frac{d}{2}} \int_{\bar{L}_j} t \exp\left(-\frac{t^{\frac{2}{j}}}{4}\right) \mathrm{d}t = 2^{\frac{d}{2}} \int_{\bar{L}_j^{\frac{2}{j}}} (j/2) t^{j-1} \exp\left(-\frac{t}{4}\right) \mathrm{d}t \\ &= 2^{\frac{d}{2}} (j/2) \int_{\bar{L}_j^{\frac{2}{j}}} \left(t^{j-1} \exp\left(-\frac{t}{8}\right)\right) \exp\left(-\frac{t}{8}\right) \mathrm{d}t \leqslant 2^{\frac{d}{2}} (j/2) \left(\frac{8(j-1)}{e}\right)^{j-1} \int_{\bar{L}_j^{\frac{2}{j}}} \exp\left(-\frac{t}{8}\right) \mathrm{d}t \\ &= 2^{\frac{d}{2}} j \left(\frac{8(j-1)}{e}\right)^{j-1} 8 \exp\left(-\bar{L}_j^{\frac{2}{j}}/8\right) \end{split}$$

Hence this quantity is bounded by  $\delta$  if  $\bar{L}_j \propto \left(d + \log\left(\frac{1}{\delta}\right)\right)^{\frac{j}{2}}$ . Then we have  $\bar{L}_j^2 F_i(\bar{L}_i) = \bar{L}_j^2 2^{\frac{d}{2}} \exp\left(-\frac{\bar{L}_i^2}{4}\right)$  which is also bounded by  $\delta$  if  $\bar{L}_j \propto \left(d + \left(\log\frac{d}{\delta}\right)^2\right)^{\frac{j}{2}}$ . At the end of the day, our assumptions are satisfied for

$$\bar{L}_j \propto \left(d + \left(\log \frac{dm}{\rho}\right)^2\right)^{\frac{1}{2}}$$

#### F.2.2 Gaussian mixture model learning

We apply the mixture model framework with the base distribution:

$$P_{\theta} = \mathcal{N}(\theta, \Sigma)$$

The random features on the data space are  $\varphi'_{\omega}(x) = C e^{i\omega^{\top}x}$  with Gaussian distribution  $\omega \sim \Lambda = \mathcal{N}(0, A)$  for some constant C and matrix A. Then, the features on the parameter space are  $\varphi_{\omega}(\theta) = \mathbb{E}_{x \sim P_{\theta}} \varphi'_{\omega}(x) = C e^{i\omega^{\top}\theta} e^{-\frac{1}{2} \|\omega\|_{\Sigma}^2}$ 

(that is, the characteristic function of Gaussians). Then, it is possible to show (Gribonval et al., 2017) that the kernel is

$$K(\theta, \theta') = C^2 \frac{\left|A^{-1}\right|^{\frac{1}{2}}}{\left|2\Sigma + A^{-1}\right|^{\frac{1}{2}}} e^{-\frac{1}{2}\left\|\theta - \theta'\right\|^2_{(2\Sigma + A^{-1})^{-1}}}$$

Hence we choose  $A = c\Sigma^{-1}$ ,  $C = (1+2c)^{\frac{d}{4}}$ , and we come back to the previous case  $K(\theta, \theta') = e^{-\frac{1}{2} \|\theta - \theta'\|_{\tilde{\Sigma}^{-1}}^2}$  with covariance  $\tilde{\Sigma} = (2+1/c)\Sigma$ . Hence  $\varepsilon_i = \mathcal{O}(1)$ ,  $B_{ij} = \mathcal{O}(1)$ ,  $d_{\mathbf{H}}(\theta, \theta') = \|\theta - \theta'\|_{\tilde{\Sigma}^{-1}} = \frac{1}{\sqrt{2+1/c}} \|\theta - \theta'\|_{\Sigma^{-1}}$ .

Admissible features. Unlike the previous case, the features are directly bounded and Lipschitz. We have

$$\begin{aligned} |\varphi_{\omega}(\theta)| &\leq C \stackrel{\text{def.}}{=} L_{0}, \\ \|\mathcal{D}_{j} \left[\varphi_{\omega}(\theta)\right]\| &= C \left\|\tilde{\Sigma}^{\frac{1}{2}}\omega\right\|^{j} e^{-\frac{\|\omega\|_{\Sigma}^{2}}{2}} = C \left(2 + 1/c\right)^{\frac{j}{2}} \left\|\Sigma^{\frac{1}{2}}\omega\right\|^{j} e^{-\frac{\|\omega\|_{\Sigma}^{2}}{2}} &\leq C \left(2 + 1/c\right)^{\frac{j}{2}} \left(\frac{j}{e}\right)^{\frac{j}{2}} \stackrel{\text{def.}}{=} L_{j}. \end{aligned}$$

Hence all constants  $L_j$  are in  $\mathcal{O}\left(C(2+1/c)^{\frac{j}{2}}\right)$  by choosing  $c = \frac{1}{d}$  they are in  $\mathcal{O}\left(d^{\frac{j}{2}}\right)$ .

### F.3 The Laplace transform kernel

Let  $\alpha \in \mathbb{R}^d_+$  and let  $\mathcal{X} \subset \mathbb{R}^d_+$  be a compact domain. Define for  $x \in \mathcal{X}$  and  $\omega \in \mathbb{R}^d_+$ ,

$$\varphi_{\omega}(x) \stackrel{\text{def.}}{=} \exp(-\langle x, \omega \rangle) \prod_{i=1}^{d} \sqrt{\frac{(x_i + \alpha_i)}{\alpha_i}} \quad \text{and} \quad \Lambda(\omega) \stackrel{\text{def.}}{=} \exp(-\langle 2\alpha, \omega \rangle) \prod_{i=1}^{d} (2\alpha_i),$$

The associated kernel is  $K(x, x') = \prod_{i=1}^{d} \kappa(x_i + \alpha_i, x'_i + \alpha_i)$  where  $\kappa$  is the 1D Laplace kernel

$$\kappa(u,v) \stackrel{\text{def.}}{=} 2 \frac{\sqrt{uv}}{(u+v)}$$

A direct computation shows that  $\mathbf{H}_x \in \mathbb{R}^{d \times d}$  is the diagonal matrix with  $(h_{x_i+\alpha_i})_{i=1}^d$  where  $h_x \stackrel{\text{def.}}{=} \partial_x \partial_{x'} \kappa(x, x) = (2x)^{-2}$ . Note that

$$d_{\kappa}(s,t) = \int_{\min\{s,t\}}^{\max\{s,t\}} (2x+2\alpha)^{-1} \mathrm{d}x = \left|\log\left(\frac{t+\alpha}{s+\alpha}\right)\right|$$
(F.5)

and so,  $d_{\mathbf{H}}(x, x') = \sqrt{\sum_{i=1}^{d} \left| \log \left( \frac{x_i + \alpha_i}{x'_i + \alpha_i} \right) \right|^2}.$ 

We have the following results concerning the boundedness of  $\|D_j[\varphi_{\omega}]\|$  and the admissibility of K:

**Theorem F.2** (Stochastic gradient bounds). Assume that the  $\alpha_i$ 's are all distinct. Then,  $\bar{L}_0(\omega) \leq \bar{L}_0 \stackrel{\text{def.}}{=} \left(1 + \frac{R_{\mathcal{X}}}{\min_i \alpha_i}\right)^d$  and for j = 1, 2, 3,

$$\mathbb{P}(L_j(\omega) \ge t) \le F_j(t) \stackrel{\text{def.}}{=} \sum_{i=1}^d \beta_i \exp\left(-\alpha_i \left(\frac{1}{2(R_{\mathcal{X}} + \|\alpha\|_{\infty})} \left(\frac{t}{\bar{L}_0}\right)^{1/j} - \sqrt{d}\right)\right)$$

and we have that  $\sum_i F_j(\bar{L}_j) \leqslant \delta$  and  $\bar{L}_j^2 \sum_i F_i(\bar{L}_i) + 2 \int_{\bar{L}_j}^{\infty} tF_j(t) dt \leqslant \delta$  provided that

$$\bar{L}_j \propto \bar{L}_0 (R_{\mathcal{X}} + \|\alpha\|_{\infty})^j \left(\sqrt{d} + \max_i \frac{1}{\alpha_i} \log\left(\frac{d\beta_i \bar{L}_0 (R_{\mathcal{X}} + \|\alpha\|_{\infty})}{\delta\alpha_i}\right)\right)^j.$$

where  $\beta_i = \prod_{j \neq i} \frac{\alpha_j}{\alpha_j - \alpha_i}$ . Note that  $\alpha_i \sim d$  implies that  $\bar{L}_0 \sim (1 + R_{\chi}/d)^d \sim e^{R_{\chi}}$ .

**Theorem F.3** (Admissiblity of K). The Laplace transform kernel K is admissible with  $r_{\text{near}} = 0.2$ ,  $C_{\text{H}} = 1.25$ ,

 $\varepsilon_0 = 0.005, \ \varepsilon_2 = 1.52. \ \text{For all } i+j \leqslant 3, \ B_{ij} = \mathcal{O}(1), \ B_{22} = \mathcal{O}(d), \ \Delta = \mathcal{O}(d+\log\left(d^{3/2}s_{\max}\right)) \ and \ h = \mathcal{O}(1).$ 

The first result Theorem F.2 is proved in Section F.3.1 and the second result, Theorem F.4 is a direct consequence of Theorem F.4 and Lemma F.5 in Section F.3.2.

#### F.3.1 Stochastic gradient bounds

Proof of Theorem F.2. Let  $V \stackrel{\text{def.}}{=} (1 - 2(x_i + \alpha_i)\omega_i)_{i=1}^d \in \mathbb{R}^d$ . Then,

$$\begin{aligned} \|V\| &= \sqrt{\sum_{i} (1 - 2(x_i + \alpha_i)\omega_i)^2} \\ &\leqslant \sqrt{\sum_{i} 1 + 4(x_i + \alpha_i)^2 \omega_i^2} \leqslant \sqrt{d + 4(R_{\mathcal{X}} + \|\alpha\|_{\infty})^2 \|w\|^2} \\ &\leqslant \sqrt{d} + 2(R_{\mathcal{X}} + \|\alpha\|_{\infty}) \|w\| \end{aligned}$$

We have the following bounds:

$$\begin{aligned} |\varphi_{\omega}(x)| &\leqslant \prod_{i=1}^{d} \sqrt{1 + \frac{x_i}{\alpha_i}} \leqslant \left(1 + \frac{R_{\mathcal{X}}}{\min_i \alpha_i}\right)^d \stackrel{\text{def.}}{=} \bar{L}_0, \\ \mathrm{D}_1\left[\varphi_{\omega}\right](x) &= \varphi_{\omega}(x)V \implies \|\mathrm{D}_1\left[\varphi_{\omega}\right](x)\| \leqslant \bar{L}_0 \|V\| \\ \mathrm{D}_2\left[\varphi_{\omega}\right](x) &= \varphi_{\omega}(x)(VV^{\top} - 2\mathrm{Id}) \implies \|\mathrm{D}_2\left[\varphi_{\omega}\right](x)\| \leqslant \bar{L}_0 \min\{\|V\|^2, 2\}. \end{aligned}$$

and given  $u, q \in \mathbb{R}^d$ ,

$$D_{3}[\varphi_{\omega}](x)[q,q,u] = \varphi_{\omega}(x) \left( \langle u, V \rangle \langle q, V \rangle^{2} - 2 \|q\|^{2} - 4 \langle u, q \rangle \langle q, V \rangle + 8 \sum_{i} q_{i}^{2} u_{i} \right),$$

 $\mathbf{so}$ 

$$\|\mathbf{D}_{3}[\varphi_{\omega}](x)\| \leq |\varphi_{\omega}(x)| \left(\|V\|^{3} + 10 + 4 \|V\|\right) \leq \bar{L}_{0}5(\|V\|^{3} + 3),$$

And therefore, in general,

$$\|\mathbf{D}_{j}\left[\varphi_{\omega}\right](x)\| \leqslant L_{j}(\omega) \stackrel{\text{def.}}{=} \bar{R}_{\mathcal{X}}^{j+1} \left(\sqrt{d} + \|\omega\|\right)^{j}$$
$$\|\mathbf{D}_{j}\left[\varphi_{\omega}\right](x)\| \lesssim L_{j}(\omega) \stackrel{\text{def.}}{=} \bar{L}_{0} \left(\sqrt{d} + 2(R_{\mathcal{X}} + \|\alpha\|_{\infty}) \|w\|\right)^{j}$$

Assuming for simplicity that all  $\alpha_j$  are distinct, we have Akkouchi:

$$\mathbb{P}(\|w\| \ge t) \le \mathbb{P}(\|\omega\|_1 \ge t) = \sum_{i=1}^d \beta_i e^{-\alpha_i t}$$

where  $\beta_i = \prod_{j \neq i} \frac{\alpha_j}{\alpha_j - \alpha_i}$ , using the fact that  $\|\omega\|_1$  is a sum of independent exponential random variable. Hence, for all  $1 \leq j \leq 3$  and  $t \geq d^{\frac{j}{2}}$  we have

$$\mathbb{P}(L_j(\omega) \ge t) \le \mathbb{P}\left(\|w\| \ge \frac{1}{2(R_{\mathcal{X}} + \|\alpha\|_{\infty})} \left(\frac{t}{\bar{L}_0}\right)^{1/j} - \sqrt{d}\right)$$
$$\le F_j(t) \stackrel{\text{def.}}{=} \sum_{i=1}^d \beta_i \exp\left(-\alpha_i \left(\frac{1}{2(R_{\mathcal{X}} + \|\alpha\|_{\infty})} \left(\frac{t}{\bar{L}_0}\right)^{1/j} - \sqrt{d}\right)\right) \le \delta$$

and  $F_j(\bar{L}_j) \leqslant \delta$  if

$$\bar{L}_j \ge \bar{L}_0 \left( 2^j (R_{\mathcal{X}} + \|\alpha\|_{\infty})^j \left( \sqrt{d} + \max_i \frac{1}{\alpha_i} \log\left(\frac{d\beta_i}{\delta}\right) \right)^j \right)$$

Next, in a similar manner to the Gaussian case, we compute

$$\begin{split} \int_{\bar{L}_{j}} tF_{j}(t) \mathrm{d}t &= \sum_{i=1}^{d} \beta_{i} \int_{\bar{L}_{j}} t \exp\left(-\alpha_{i} \left(\frac{1}{2(R_{\mathcal{X}} + \|\alpha\|_{\infty})} \left(\frac{t}{\bar{L}_{0}}\right)^{1/j} - \sqrt{d}\right)\right) \mathrm{d}t \\ &= \bar{L}_{0}^{2} j \sum_{i=1}^{d} e^{\alpha_{i}\sqrt{d}} \beta_{i} \int_{(\bar{L}_{j}/\bar{L}_{0})^{1/j}} \exp\left(\frac{-\alpha_{i} u}{2(R_{\mathcal{X}} + \|\alpha\|_{\infty})}\right) u^{2j-1} \mathrm{d}u \\ &\leqslant \left(\frac{(2j-1)4(R_{\mathcal{X}} + \|\alpha\|_{\infty})}{e\alpha_{i}}\right)^{2j-1} \bar{L}_{0}^{2} j \sum_{i=1}^{d} e^{\alpha_{i}\sqrt{d}} \beta_{i} \int_{(\bar{L}_{j}/\bar{L}_{0})^{1/j}} \exp\left(\frac{-\alpha_{i} u}{4(R_{\mathcal{X}} + \|\alpha\|_{\infty})}\right) \mathrm{d}u \\ &\leqslant \left(\frac{4(R_{\mathcal{X}} + \|\alpha\|_{\infty})}{\alpha_{i}}\right)^{2j} \left(\frac{2j-1}{e}\right)^{2j-1} \bar{L}_{0}^{2} j \sum_{i=1}^{d} e^{\alpha_{i}\sqrt{d}} \beta_{i} \exp\left(\frac{-\alpha_{i}(\bar{L}_{j}/\bar{L}_{0})^{1/j}}{4(R_{\mathcal{X}} + \|\alpha\|_{\infty})}\right) \leqslant \delta \end{split}$$

if for all  $i = 1, \ldots, d$ ,

$$\frac{4(R_{\mathcal{X}} + \|\alpha\|_{\infty})}{\alpha_i} \left(2j \log\left(\frac{4(2j-1)(R_{\mathcal{X}} + \|\alpha\|_{\infty})}{e\alpha_i}\right) + \log(\bar{L}_0^2 j) + \alpha_i \sqrt{d} + \log\left(\frac{d\beta_i}{\delta}\right)\right) \leqslant \left(\frac{\bar{L}_j}{\bar{L}_0}\right)^{1/j}$$

that is,

$$\bar{L}_j \gtrsim \bar{L}_0 \left( 2^j (R_{\mathcal{X}} + \|\alpha\|_{\infty})^j \left( \sqrt{d} + \max_i \frac{1}{\alpha_i} \log\left(\frac{d\beta_i}{\delta}\right) \right)^j \right).$$

It remains to bound  $\bar{L}_j F_\ell(\bar{L}_\ell)$  with  $\ell, j \in \{0, 1, 2, 3\}$ : Let  $\bar{L}_\ell \ge \bar{L}_0 M^\ell$  for some M to be determined. Then,

$$\begin{split} \bar{L}_{j}F_{\ell}(\bar{L}_{\ell}) &\leqslant \bar{L}_{0}M^{j}\sum_{i=1}^{d}\beta_{i}\exp\left(\frac{-\alpha_{i}}{2(R_{\mathcal{X}}+\|\alpha\|_{\infty})}M+\alpha_{i}\sqrt{d}\right) \\ &= \bar{L}_{0}\sum_{i=1}^{d}\beta_{i}M^{j}\exp\left(\frac{-\alpha_{i}}{4(R_{\mathcal{X}}+\|\alpha\|_{\infty})}M\right)\exp\left(\frac{-\alpha_{i}}{4(R_{\mathcal{X}}+\|\alpha\|_{\infty})}M\right)e^{\alpha_{i}\sqrt{d}} \\ &\leqslant \bar{L}_{0}e^{-j}\sum_{i=1}^{d}\left(\frac{4j(R_{\mathcal{X}}+\|\alpha\|_{\infty})}{\alpha_{i}}\right)^{j}\beta_{i}\exp\left(\frac{-\alpha_{i}}{4(R_{\mathcal{X}}+\|\alpha\|_{\infty})}M\right)e^{\alpha_{i}\sqrt{d}} \\ &\leqslant \bar{L}_{0}e^{-3}\sum_{i=1}^{d}\left(\frac{12(R_{\mathcal{X}}+\|\alpha\|_{\infty})}{\alpha_{i}}\right)^{3}\beta_{i}\exp\left(\frac{-\alpha_{i}}{4(R_{\mathcal{X}}+\|\alpha\|_{\infty})}M\right)e^{\alpha_{i}\sqrt{d}} \leqslant \delta \end{split}$$

if for each  $i = 1, \ldots, d$ 

$$M \ge 4(R_{\mathcal{X}} + \|\alpha\|_{\infty}) \left(\sqrt{d} + \max_{i} \frac{1}{\alpha_{i}} \log\left(\frac{\bar{L}_{0}d\beta_{i}}{\delta e^{3}} \left(\frac{12(R_{\mathcal{X}} + \|\alpha\|_{\infty})}{\alpha_{i}}\right)^{3}\right)\right).$$

Therefore, similar to the Gaussian case, the conclusion follows for  $\bar{L}_0 = \left(1 + \frac{R_{\chi}}{\min_i \alpha_i}\right)^d$ , and for j = 1, 2, 3,

$$\bar{L}_j \propto \bar{L}_0 (R_{\mathcal{X}} + \|\alpha\|_{\infty})^j \left(\sqrt{d} + \max_i \frac{1}{\alpha_i} \log\left(\frac{d\beta_i \bar{L}_0 (R_{\mathcal{X}} + \|\alpha\|_{\infty})}{\delta\alpha_i}\right)\right)^j.$$

### F.3.2 Admissibility of the kernel

Metric variation We have the following lemma on the variation of the Fisher metric:

Lemma F.5. Suppose that  $d_{\mathbf{H}}(x, x') \leq c$ , then  $\left\| \operatorname{Id} - \mathbf{H}_{x'}^{1/2} \mathbf{H}_{x'} \right\| \leq (1 + ce^{c}) d_{\mathbf{H}}(x, x')$ .

*Proof.* Note that  $|1 - |(x_i + \alpha_i)/(x'_i + \alpha_i)|| \leq \max\{e^{d_{\kappa}(x_i, x'_i)} - 1, 1 - e^{-d_{\kappa}(x_i, x'_i)}\} \leq d_{\kappa}(x_i, x'_i)(1 + ce^c)$  for all  $d_{\kappa}(x_i, x'_i) \leq c$ . Therefore,

$$\|\mathrm{Id} - \mathbf{H}_{x}\mathbf{H}_{x'}\|^{2} = \sum_{i} |1 - |(x_{i} + \alpha_{i})/(x_{i}' + \alpha_{i})||^{2} \leq (1 + ce^{c})d_{\mathbf{H}}(x, x')$$

provided that  $d_{\mathbf{H}}(x, x') \leq c$ .

Admissibility of the kernel The following theorem provides bounds for K and its normalised derivatives. Theorem F.4. 1.  $|K(x,x')| \leq \min\{2^d e^{-\frac{1}{2}d_{\mathbf{H}}(x,x')}, \frac{8}{8+d_{\mathbf{H}}(x,x')^2}\}$ .

- 2.  $||K^{(10)}(x, x')|| \le \min\{2\sqrt{d} |K|, \sqrt{2}\}.$
- 3.  $\|K^{(11)}\| \leq \min\{9d |K|, 8\}$
- 4.  $||K^{(20)}|| \leq \min\{10d |K|, 8\}$  and  $\lambda_{\min}(-K^{(20)}) \geq (2 12d_{\mathbf{H}}(x, x')^2) K$ .
- 5.  $||K^{(12)}|| \leq \min\{66 |K| d^{3/2}, 16\sqrt{d} + 49\}$  and  $||K^{(12)}(x, x')|| \leq 34$  if  $d_{\mathbf{H}}(x, x') \leq 1$ .

6. 
$$||K^{(22)}|| \leq 16d + 9$$

In particular, for  $d_{\mathbf{H}}(x, x') \ge 2d\log(2) + 2\log\left(\frac{52d^{3/2}s_{\max}}{h}\right)$ , we have  $\left\|K^{(ij)}(x, x')\right\| \le \frac{h}{s_{\max}}$ .

To prove this result, we first present some bounds for the univariate Laplace kernel in Section F.3.3 before applying these bounds in Section F.3.4.

### F.3.3 1D Laplace kernel

In the following  $\kappa^{(ij)}(x,x') \stackrel{\text{def.}}{=} h_x^{-i/2} h_{x'}^{-j/2} \partial_x^i \partial_{x'}^j \kappa(x,x')$ . Lemma F.6. We have

(i) 
$$\kappa(x, x') = \operatorname{sech}\left(\frac{d_{\kappa}(x, x')}{2}\right) \leq 2e^{-\frac{1}{2}d_{\kappa}(x, x')},$$
  
(ii)  $|\kappa^{(10)}(x, x')| = 2 \left| \tanh\left(\frac{d_{\kappa}(x, x')}{2}\right)\kappa(x, x') \right|, and \left|\kappa^{(10)}\right| \leq 2 |\kappa|.$   
(iii)  $|\kappa^{(11)}| \leq 4 |\kappa|^3 + 4 |\kappa|$ 

(iv) 
$$\left|\kappa^{(20)}\right| \leq 6 \left|\kappa\right| and -\kappa^{(20)} \geq 2\kappa(x, x') \left(1 - 2 \tanh\left(\frac{d_{\kappa}(x, x')}{2}\right)\right).$$

(v) 
$$\left|\kappa^{(12)}\right| \leq 49 \left|\kappa\right|$$
.

(vi) 
$$\kappa^{(22)}(x,x) = 9$$
 for all x.

*Proof.* We first state the partial derivatives of  $\kappa$ :

$$\begin{split} \kappa(x,x') &= \frac{2\sqrt{xx'}}{x+x'},\\ \partial_x \kappa(x,x') &= \frac{x'(x'-x)}{\sqrt{xx'}(x+x')^2}\\ \partial_x \partial_{x'} \kappa(x,x') &= \frac{-x^2 + 6xx' - (x')^2}{2\sqrt{xx'}(x+x')^3}\\ \partial_x^2 \kappa(x,x') &= -\frac{(x')^2 \left((x+x')^2 + 4x(x'-x)\right)}{2 \left(xx'\right)^{3/2} (x+x')^3}\\ &= -\frac{(x')^2}{2 \left(xx'\right)^{3/2} (x+x')} - \frac{2x'(x'-x)}{(xx')^{1/2} (x+x')^3}\\ \partial_x \partial_{x'}^2 \kappa(x,x') &= \frac{x^3 + 13x^2x' - 33x(x')^2 + 3(x')^3}{4x'(xx')^{1/2} (x+x')^4}\\ \partial_x^2 \partial_{x'}^2 \kappa(x,x') &= -\frac{3x^4 + 60x^3x' - 270x^2(x')^2 + 60x(x')^3 + 3(x')^4}{8xx'(xx')^{1/2} (x+x')^5} \end{split}$$

(i)

$$\kappa(x,x') = 2\left(\sqrt{\frac{x}{x'}} + \sqrt{\frac{x'}{x}}\right)^{-1} = \frac{2}{e^{-\frac{d_{\kappa}(x,x')}{2}} + e^{\frac{d_{\kappa}(x,x')}{2}}} = \frac{1}{\cosh(\frac{d_{\kappa}(x,x')}{2})} \leqslant 2e^{-\frac{1}{2}d_{\kappa}(x,x')},$$

(ii) We have, assuming that x > x',

$$\begin{aligned} \kappa^{(10)}(x,x') &= 2x\partial_x \kappa(x,x') = 2\frac{x'-x}{x+x'}\kappa(x,x') \\ &= 2\left(\frac{1}{\frac{x}{x'}+1} - \frac{1}{1+\frac{x'}{x}}\right)\kappa(x,x') \\ &= 2\left(\frac{1}{1+\exp(d_\kappa(x,x'))} - \frac{1}{1+\exp(-d_\kappa(x,x'))}\right) \\ &= 2\left(\frac{\exp(-d_\kappa(x,x')) - \exp(d_\kappa(x,x'))}{2+\exp(d_\kappa(x,x')) + \exp(d_\kappa(x,x'))}\right) \\ &= \frac{-2\sinh(d_\kappa(x,x'))}{1+\cosh(d_\kappa(x,x'))}\kappa(x,x') \\ &= -2\tanh(d_\kappa(x,x')/2)\kappa(x,x'), \end{aligned}$$

(iii)

$$\kappa^{(11)} = 4xx'\partial_{x'}\partial_{x}\kappa(x,x') = 4xx'\frac{4xx' - (x - x')^{2}}{2\sqrt{xx'}(x + x')^{3}}$$
$$= 4\kappa(x,x')^{3} - \frac{4(x - x')^{2}}{(x + x')^{2}}\kappa(x,x')$$
$$= \kappa(x,x')\left(4\kappa(x,x')^{2} - 4\tanh^{2}(d_{\kappa}(x,x')/2)\right)$$

so  $|\kappa^{(11)}| \leq 4 |\kappa|^3 + 4 |\kappa|$ .

(iv)

$$\kappa^{(20)} = 4x^2 \partial_x^2 \kappa(x, x') = -\frac{4(xx')^{1/2} \left((x+x')^2 + 4x(x'-x)\right)}{2(x+x')^3}$$
$$= -2\kappa(x, x') \left(1 + \frac{2x(x'-x)}{(x+x')^2}\right)$$

so  $|\kappa^{20}| \leq 6 |\kappa|$ . Also,

$$-\kappa^{(20)} \ge 2\kappa(x, x') \left(1 - 2 \tanh(d_{\kappa}(x, x')/2)\right)$$

(v)

$$\kappa^{(12)} = 2x(2x')^2 \partial_x \partial_{x'}^2 \kappa(x, x')$$
  
=  $\kappa(x, x') \left( 1 + \frac{2v(5u^2 - 18uv + v^2)}{(u+v)^3} \right)$ 

so  $|\kappa^{(12)}| \leq 49 |\kappa|$ . (vi)

$$\kappa^{(22)} = 16(xx')^2 \partial_x^2 \partial_{x'}^2 \kappa(x,x')$$
  
=  $-3 - \frac{48xx'(x^2 - 6xx' + (x')^2)}{(x+x')^4}$ 

and  $\kappa^{(22)}(x, x) = 9$ .

## F.3.4 Proof of Theorem F.4

Let  $d_{\ell} \stackrel{\text{def.}}{=} d_{\kappa}(x_{\ell} + \alpha_{\ell}, x'_{\ell} + \alpha_{\ell})$  and note that  $d_{\mathbf{H}}(x, x') = \sqrt{\sum_{\ell} d_{\ell}^2}$ . Define  $g = \left(2 \tanh(\frac{d_{\ell}}{2})\right)_{\ell=1}^d$ . We first prove that

- (i)  $|K(x,x')| \leq \prod_{\ell=1}^{d} \operatorname{sech}(d_{\ell}/2) \leq \prod_{\ell=1}^{d} \frac{1}{1+d_{\ell}^2/8} \leq \frac{1}{1+\frac{1}{8}d_{\mathbf{H}}(x,x')^2}$
- (ii)  $||K^{(10)}(x, x')|| \leq ||g||_2 |K|.$

(iii) 
$$||K^{(11)}|| \leq |K| (||g||_2^2 + 5)$$

(iv) 
$$||K^{(20)}|| \leq |K| (||g||_2^2 + 6)$$
 and  $\lambda_{\min} (K^{(20)}) \geq K (2 - 3 ||g||_2^2)$ .

(v) 
$$||K^{(12)}|| \leq |K| (||g||_2^3 + 16 ||g||_2 + 49)$$

(vi) 
$$||K^{(22)}|| \le 16d + 9.$$

The result would then follow because

- $\operatorname{sech}(x) \leq 2e^{-x}$  and  $\operatorname{sech}(x) \leq (1 + x^2/2)^{-1}$ .
- $|\operatorname{tanh}(x)| \leq \min\{x, 1\}$ , so  $||g|| \leq \min\{d_{\mathbf{H}}(x, x'), 2\sqrt{d}\}$ ,

For example,  $\|K^{(12)}\| \leq \frac{1}{1+\frac{1}{8}d_{\mathbf{H}}(x,x')^2} \left( d_{\mathbf{H}}(x,x')^3 + 16d_{\mathbf{H}}(x,x') + 24 \right) \leq 8d_{\mathbf{H}}(x,x') + \frac{\sqrt{8}}{2} + 24 \leq 34$  when  $d_{\mathbf{H}}(x,x') \leq 1$ .

In the following, we write  $\kappa_{\ell}^{(ij)} \stackrel{\text{def.}}{=} \kappa^{(ij)} (x_{\ell} + \alpha_{\ell}, x'_{\ell} + \alpha_{\ell})$  and  $\kappa_{\ell} \stackrel{\text{def.}}{=} \kappa_{\ell}^{(00)}$  and  $K_i \stackrel{\text{def.}}{=} \prod_{j \neq i} \kappa_j$ . Moreover, we will make use of the inequalities for  $\kappa^{(ij)}$  derived in Lemma F.6.

L		J

(i) We have

$$|K(x,x')| \leq \prod_{\ell=1}^{d} \operatorname{sech}(d_{\ell}) \leq \prod_{\ell=1}^{d} \left(1 + \frac{d_{\ell}^2}{2}\right)^{-1} \leq \frac{1}{1 + d_{\mathbf{H}}(x,x')^2}.$$

(ii)

$$K^{(10)}(x,x') = \left(\kappa_{\ell}^{(10)}K_{\ell}\right)_{\ell=1}^{d} \implies \left\|K^{(10)}(x,x')\right\| \le \|g\|_{2} \|K\|_{2}$$

(iii) For  $i \neq j$ 

$$\left|K_{ij}^{(11)}\right| = \left|\kappa_i^{(10)}\kappa_j^{(01)}K_{ij}\right| \leqslant 4 \tanh\left(\frac{\mathbf{d}_i}{2}\right) \tanh\left(\frac{\mathbf{d}_j}{2}\right) |K|,$$

and  $\left|K_{ii}^{(11)}\right| = \left|\kappa_i^{(11)}K_i\right| \leqslant 5 |K|$ . So, given  $p \in \mathbb{R}^d$  of unit norm,

$$\langle K^{(11)}p, p \rangle = \sum_{i=1}^{d} \sum_{j \neq i} \kappa_i^{(10)} \kappa_j^{(01)} K_{ij} p_i p_j + \sum_{i=1}^{d} p_i^2 \kappa_i^{(11)} K_i$$

$$\leq |K| \left( \sum_{i=1}^{d} \sum_{j \neq i} 4 \tanh(d_i/2) \tanh(d_j/2) p_i p_j + 5 \sum_{i=1}^{d} p_i^2 \right)$$

$$\leq |K| \left( ||g||_2^2 + 5 \right)$$

(iv) For  $i \neq j$ ,  $K_{ij}^{(20)} = \kappa_i^{(10)} \kappa_j^{(10)} K_{ij}$ , and  $\left| K_{ii}^{(20)} \right| = \left| \kappa_i^{(20)} K_i \right| \le 6 |K|$  and  $-K_{ii}^{(20)} \ge 2K \left( 1 - 2 \tanh\left(\frac{\mathrm{d}_i}{2}\right) \right)$ .

$$\langle K^{(20)}p, p \rangle = \sum_{i=1}^{d} \sum_{j \neq i} \kappa_i^{(10)} \kappa_j^{(10)} K_{ij} p_i p_j + \sum_{i=1}^{d} p_i^2 \kappa_i^{(20)} K_i$$
  
 
$$\leqslant |K| \left( \sum_{i=1}^{d} \sum_{j \neq i} 4 \tanh(d_i/2) \tanh(d_j/2) p_i p_j + 6 \sum_{i=1}^{d} p_i^2 \right)$$
  
 
$$\leqslant |K| \left( ||g||_2^2 + 6 \right),$$

and

$$\langle -K^{(20)}p, p \rangle \ge K \left( 2 - 2 \|g\|_{\infty} - \|g\|_{2}^{2} \right)$$

(v) For  $i, j, \ell$  all distinct,

$$K_{ij\ell}^{(12)} = \kappa_i^{(10)} \kappa_j^{(01)} \kappa_\ell^{(01)} K_{ij\ell} \leqslant 8 \tanh\left(\frac{\mathbf{d}_i}{2}\right) \tanh\left(\frac{\mathbf{d}_j}{2}\right) \tanh\left(\frac{\mathbf{d}_\ell}{2}\right) K,$$

for all  $i, \ell$ ,

$$K_{ii\ell}^{(12)} = 8\kappa_i^{(11)}\kappa_\ell^{(01)}K_{i\ell} \leqslant 10\tanh\left(\frac{d_\ell}{2}\right)K$$
$$K_{iji}^{(12)} = \kappa_i^{(11)}\kappa_j^{(01)}K_{ij} \leqslant 10\tanh\left(\frac{d_j}{2}\right)K,$$

$$\begin{split} K_{ijj}^{(12)} &= \kappa_i^{(10)} \kappa_\ell^{(02)} K_{ij} \leqslant 12 \tanh\left(\frac{d_i}{2}\right) K, \text{ and } K_{iii}^{(12)} = \kappa_i^{(12)} K_i \leqslant 26 K. \text{ So, for } p, q \in \mathbb{R}^d \text{ of unit norm,} \\ &\sum_i \sum_j \sum_\ell K_{ij\ell}^{(12)} p_j p_\ell q_i = \sum_i \left( \sum_{j \neq i} \sum_\ell K_{ij\ell}^{(12)} p_j p_\ell q_i + \sum_\ell K_{ii\ell}^{(12)} p_i p_\ell q_i \right) \\ &= \sum_i \sum_{j \neq i} \left( \sum_{\ell \notin \{i,j\}} K_{ij\ell}^{(12)} p_j p_\ell q_i + K_{iji}^{(12)} p_j p_i q_i + K_{ijj}^{(12)} p_j^2 q_i \right) \\ &+ \sum_i \sum_{\ell \neq i} K_{ii\ell}^{(12)} p_i p_\ell q_i + \sum_i K_{iii}^{(12)} p_i^2 q_i \\ &\leqslant |K| \left( ||g||_2^3 + 16 ||g||_2 + 49 \right). \end{split}$$

(vi)

$$\begin{split} \left\| K^{(22)}(x,x) \right\| &= \sup_{\|p\|=1} \mathbb{E}[\langle \mathbf{H}_x^{-1/2} \nabla^2 \varphi_{\omega}(x) \mathbf{H}_x^{-1/2} p, \, \mathbf{H}_x^{-1/2} \nabla^2 \varphi_{\omega}(x) \mathbf{H}_x^{-1/2} p \rangle] \\ &\leqslant \sup_{\|p\|=1} \sum_i \sum_{k \neq i} \kappa_i^{(11)} \kappa_k^{(11)} p_i^2 + \sum_i \sum_{k \neq i} \kappa_i^{(12)} \kappa_k^{(10)} p_i p_k + \sum_i \sum_{k \neq i} \sum_{j \notin \{i,k\}} \kappa_i^{(11)} \kappa_k^{(10)} \kappa_j^{(01)} p_k p_j \\ &+ \sum_i \sum_{j \neq i} \kappa_i^{(21)} \kappa_j^{(01)} p_j p_i + \sum_i \kappa_i^{(22)} p_i^2 \\ &= \sup_{\|p\|=1} \sum_i \sum_{k \neq i} \kappa_i^{(11)} \kappa_k^{(11)} p_i^2 + \sum_i \kappa_i^{(22)} p_i^2 \\ &\leqslant d \left\| \kappa^{(11)} \right\|_{\infty} + \left\| \kappa^{(22)} \right\|_{\infty} \leqslant 16d + \left\| \kappa^{(22)} \right\|_{\infty}. \end{split}$$

since  $\kappa^{(10)}(x,x) = \kappa^{(01)}(x,x) = 0$ , and  $\kappa^{(11)}(x,x) = 4$  from the proof of (iii) in Lemma F.6.

### G Tools

### G.1 Probability tools

**Lemma G.1** (Bernstein's inequality (Sridharan (2002), Thm. 6)). Let  $x_1, \ldots, x_n \in \mathbb{C}$  be i.i.d. bounded random variables such that  $\mathbb{E}x_i = 0$ ,  $|x_i| \leq M$  and  $Var(x_i) \stackrel{\text{def.}}{=} \mathbb{E}[|x_i|^2] \leq \sigma^2$  for all *i*'s.

Then for all t > 0 we have

$$\mathcal{X}\left(\frac{1}{n}\sum_{i=1}^{n}x_{i} \ge t\right) \le 4\exp\left(-\frac{nt^{2}/4}{\sigma^{2} + Mt/(3\sqrt{2})}\right).$$
(G.1)

**Lemma G.2** (Matrix Bernstein (Tropp (2015), Theorem 6.1.1)). Let  $Y_1, ..., Y_m \in \mathbb{C}^{d_1, d_2}$  be complex random matrices with

$$\mathbb{E}Y_j = 0, \quad \|Y_j\| \leq L, \quad v(Y_j) := \max(\left\|\mathbb{E}Y_jY_j^*\right\|, \left\|\mathbb{E}Y_j^*Y_j\right\|) \leq M$$

for each index  $1 \leq j \leq m$ . Introduce the random matrix

$$Z = \frac{1}{m} \sum_{j} Y_j$$

Then

$$\mathbb{P}(\|Z\| \ge t) \le 2(d_1 + d_2)e^{-\frac{mt^2/2}{M + Lt/3}}$$
(G.2)

Lemma G.3 (Vector Bernstein for complex vectors Minsker (2017)). Let  $Y_1, \ldots, Y_M \in \mathbb{C}^d$  be a sequence of

independent random vectors such that  $\mathbb{E}[Y_i] = 0$ ,  $||Y_i||_2 \leq K$  for i = 1, ..., M and set

$$\sigma^2 \stackrel{\text{def.}}{=} \sum_{i=1}^M \mathbb{E} \left\| Y_i \right\|_2^2$$

Then, for all  $t \ge (K + \sqrt{K^2 + 36\sigma^2})/M$ ,

$$\mathbb{P}\left(\left\|\frac{1}{M}\sum_{i=1}^{M}Y_{i}\right\|_{2} \ge t\right) \leqslant 28\exp\left(-\frac{Mt^{2}/2}{\sigma^{2}/M + tK/3}\right)$$

**Lemma G.4** (Hoeffding's inequality ((Tang et al., 2013), Lemma G.1)). Let the components of  $u \in C^k$  be drawn *i.i.d.* from a symmetric distribution on the complex unit circle or 0, consider a vector  $w \in \mathbb{C}^k$ . Then, with probability at least  $1 - \rho$ , we have

$$\mathbb{P}\left(\left|\langle u, w \rangle\right| \ge t\right) \leqslant 4e^{-\frac{t^2}{4\|w\|^2}} \tag{G.3}$$

**Lemma G.5.** (Tropp, 2015, Theorem 4.1.1) Let the components of  $u \in \mathbb{R}^k$  be a Rademacher sequence and let  $Y_1, \ldots, Y_M \in \mathbb{C}^{d \times d}$  be self-adjoint matrices. Set  $\sigma^2 \stackrel{\text{def.}}{=} \left\| \sum_{\ell=1}^M Y_\ell^2 \right\|$ . Then, for t > 0,

$$\mathbb{P}\left(\left\|\sum_{\ell=1}^{M} u_{\ell} Y_{\ell}\right\| \ge t\right) \le 2d \exp\left(-\frac{t^2}{2\sigma^2}\right).$$
(G.4)

We were only able to find a reference for this result in the case where u is a Rademacher sequence, however, by the contraction principle (see (Ledoux and Talagrand, 2013, Theorem 4.4)), a similar statement is true for Steinhaus sequences (we write only for the case of real symmetric matrices because this is all we require in this paper, but of course, the same argument extends to complex self-adjoint matrices):

**Corollary G.1.** Let the components of  $u \in \mathbb{C}^k$  i.i.d. from a symmetric distribution on the complex unit circle or 0 and let  $B_1, \ldots, B_M \in \mathbb{R}^{d \times d}$  be symmetric matrices. Set  $\sigma^2 \stackrel{\text{def.}}{=} \left\| \sum_{\ell=1}^M B_\ell^2 \right\|$ . Then, for t > 0,

$$\mathbb{P}\left(\left\|\sum_{\ell=1}^{M} u_{\ell} B_{\ell}\right\| \ge t\right) \le 4d \exp\left(-\frac{t^2}{4\sigma^2}\right).$$
(G.5)

*Proof.* By the union bound,

$$\mathbb{P}\left(\left\|\sum_{\ell=1}^{M} u_{\ell} B_{\ell}\right\| \ge t\right) \leqslant \mathbb{P}\left(\left\|\sum_{\ell=1}^{M} \operatorname{Re}\left(u_{\ell}\right) B_{\ell}\right\| \ge \frac{t}{\sqrt{2}}\right) + \mathbb{P}\left(\left\|\sum_{\ell=1}^{M} \operatorname{Im}\left(u_{\ell}\right) B_{\ell}\right\| \ge \frac{t}{\sqrt{2}}\right).$$

By the contraction principle (Ledoux and Talagrand, 2013, Theorem 4.4),

$$\mathbb{P}\left(\left\|\sum_{\ell=1}^{M} \operatorname{Re}\left(u_{\ell}\right)B_{\ell}\right\| \geq \frac{t}{\sqrt{2}}\right) \leq \mathbb{P}\left(\left\|\sum_{\ell=1}^{M} \xi_{\ell}B_{\ell}\right\| \geq \frac{t}{\sqrt{2}}\right)$$

where  $\xi$  is a Rademacher sequence, and the same argument applies to the case of  $\operatorname{Im}(u_{\ell})$ . Therefore by Lemma G.5, we have  $\mathbb{P}\left(\left\|\sum_{\ell=1}^{M} u_{\ell} B_{\ell}\right\| \ge t\right) \le 4d \exp\left(-\frac{t^2}{4\sigma^2}\right)$ .

#### G.2 Linear algebra tools

The following simple lemma will be handy.

**Lemma G.6.** For  $1 \leq i, j \leq s$ , take any scalars  $a_{ij} \in \mathbb{R}$ , vectors  $Q_{ij}, R_{ij} \in \mathbb{R}^d$  and square matrices  $A_{ij} \in \mathbb{R}^{d \times d}$ .

1. Let  $M \in \mathbb{R}^{sd \times sd}$  be a matrix formed by blocks :

$$M = \begin{pmatrix} A_{11} & \dots & A_{1s} \\ \vdots & \ddots & \vdots \\ A_{s1} & \dots & A_{ss} \end{pmatrix}$$

 $Then \ we \ have$ 

$$\|M\|_{\text{block}} = \sup_{\|x\|_{\text{block}}=1} \|Mx\|_{\text{block}} \leqslant \max_{1 \leqslant i \leqslant s} \sum_{j=1}^{s} \|A_{ij}\|$$
(G.6)

Now, let  $P \in \mathbb{R}^{sd \times s}$  be a rectangular matrix formed by stacking vectors  $Q_{ij} \in \mathbb{R}^d$ :

$$M = \begin{pmatrix} Q_{11} & \dots & Q_{1s} \\ \vdots & \ddots & \vdots \\ Q_{s1} & \dots & Q_{ss} \end{pmatrix}$$

Then,

$$\|M\|_{\infty \to \text{block}} \leqslant \max_{1 \leqslant i \leqslant s} \sum_{j=1}^{s} \|Q_{ij}\|_{2}, \quad \|M^{\top}\|_{\text{block} \to \infty} \leqslant \max_{1 \leqslant i \leqslant s} \sum_{j=1}^{s} \|Q_{ji}\|_{2} \tag{G.7}$$

2. Consider  $A \in \mathbb{R}^{s(d+1) \times s(d+1)}$  decomposed as

$$M = \begin{pmatrix} a_{11} & \dots & a_{1s} & Q_{11}^{\top} & \dots & Q_{1s}^{\top} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ a_{s1} & \dots & a_{ss} & Q_{s1}^{\top} & \dots & Q_{ss}^{\top} \\ R_{11} & \dots & R_{1s} & A_{11} & \dots & A_{1s} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ R_{s1} & \dots & R_{ss} & A_{s1} & \dots & A_{ss} \end{pmatrix}$$

Then,

$$\|M\| \leqslant \sqrt{\sum_{i,j} a_{ij}^2 + \|Q_{ij}\|^2 + \|R_{ij}\|^2 + \|A_{ij}\|^2},$$
$$\|M\|_{\text{Block}} \leqslant \max_i \{\sum_j |a_{ij}| + \|Q_{ij}\|, \sum_j \|R_{ij}\| + \|A_{ij}\|\}$$

*Proof.* The proof is simple linear algebra.

1. Let x be a vector with  $||x||_{\text{block}} \leq 1$  decomposed into blocks  $x = [x_1, \ldots, x_s]$  with  $x_i \in \mathbb{R}^d$ , we have

$$\|Mx\|_{\text{block}}^2 = \max_{1 \le i \le s} \left\| \sum_{j=1}^s A_{ij} x_j \right\| \le \max_i \sum_j \|A_{ij}\| \|x_j\| \le \max_i \sum_j \|A_{ij}\|$$

2. Similarly,

$$\left\|M^{\top}x\right\|_{\infty} = \max_{1 \leqslant i \leqslant s} \left\|\sum_{j=1}^{s} Q_{ji}^{\top}x_{j}\right\| \leqslant \max_{i} \sum_{j} \left\|Q_{ji}\right\| \left\|x_{j}\right\| \leqslant \max_{i} \sum_{j} \left\|Q_{ji}\right\|$$

Then, taking  $x\in \mathbb{R}^s$  such that  $\|x\|_\infty\leqslant 1,$  we have

$$\|Mx\|_{\text{block}} = \max_{1 \le i \le s} \left\| \sum_{j=1}^{s} x_j Q_{ij} \right\| \le \max_i \sum_j \|Q_{ij}\|$$

3. Taking  $x = [x_1, ..., x_s, X_1, ..., X_s] \in \mathbb{R}^{s(d+1)}$  with ||x|| = 1, we have

$$\|Mx\|^{2} = \sum_{i=1}^{s} \left( \sum_{j=1}^{s} a_{ij}x_{j} + Q_{ij}^{\top}X_{j} \right)^{2} + \left\| \sum_{j=1}^{s} R_{ij}x_{j} + A_{ij}X_{j} \right\|^{2}$$
$$\leq \sum_{i=1}^{s} \left( \|x\| \sqrt{\sum_{j=1}^{s} a_{ij}^{2} + \|Q_{ij}\|^{2}} \right)^{2} + \left( \|x\| \sqrt{\sum_{j=1}^{s} \|R_{ij}\|^{2} + \|A_{ij}\|^{2}} \right)^{2}$$
$$\leq \sum_{i,j} a_{ij}^{2} + \|Q_{ij}\|^{2} + \|R_{ij}\|^{2} + \|A_{ij}\|^{2}$$

Now, if  $||x||_{\text{Block}} = 1$ , we have

$$\|Mx\|_{\text{Block}} = \max_{i} \left( \left| \sum_{j=1}^{s} a_{ij} x_{j} + Q_{ij}^{\top} X_{j} \right|, \left\| \sum_{j=1}^{s} R_{ij} x_{j} + A_{ij} X_{j} \right\| \right)$$
$$\leq \max_{i} \left( \sum_{j=1}^{s} |a_{ij}| + \|Q_{ij}\|, \sum_{j=1}^{s} \|R_{ij} x_{j} + A_{ij} X_{j}\| \right)$$

-	_	_	۰.
-			