

# Off-the-grid learning of mixtures from a continuous dictionary

Cristina Butucea<sup>1,a</sup>, Jean-François Delmas<sup>2,b</sup> Anne Dutfoy<sup>3,d</sup> and Clément Hardy<sup>2,3,c</sup>

<sup>1</sup>CREST, ENSAE, IP Paris, France, <sup>a</sup>[cristina.butucea@ensae.fr](mailto:cristina.butucea@ensae.fr)

<sup>2</sup>CERMICS, École des Ponts, France, <sup>b</sup>[jean-francois.delmas@enpc.fr](mailto:jean-francois.delmas@enpc.fr); <sup>c</sup>[clement.hardy@enpc.fr](mailto:clement.hardy@enpc.fr)

<sup>3</sup>EDF R&D, Palaiseau, France, <sup>d</sup>[anne.dutfoy@edf.fr](mailto:anne.dutfoy@edf.fr)

**Abstract.** We consider a general non-linear model where the signal is a finite mixture of an unknown, possibly increasing, number of features issued from a continuous dictionary parameterized by a real non-linear parameter. The signal is observed with Gaussian (possibly correlated) noise in either a continuous or a discrete setup. We propose an off-the-grid optimization method, that is, a method which does not use any discretization scheme on the parameter space, to estimate both the non-linear parameters of the features and the linear parameters of the mixture.

We use recent results on the geometry of off-the-grid methods to give minimal separation on the true underlying non-linear parameters such that interpolating certificate functions can be constructed. Using also tail bounds for suprema of Gaussian processes we bound the prediction error with high probability. Assuming that the certificate functions can be constructed, our prediction error bound is up to log-factors similar to the rates attained by the Lasso predictor in the linear regression model. We also establish convergence rates that quantify with high probability the quality of estimation for both the linear and the non-linear parameters.

We develop in full details our main results for two applications: the Gaussian spike deconvolution and the scaled exponential model.

*MSC2020 subject classifications:* Primary 62G08; secondary 62G05

*Keywords:* Continuous dictionary, Interpolating certificates, Mixture model, Non-linear regression model, Off-the-grid methods, Sparse spike deconvolution

## 1. Introduction

### 1.1. Model and method

Assume we observe a random element  $y$  of an Hilbert space and we consider a signal-plus-noise structure for the observation  $y$ , where the noise is distributed according to a centered Gaussian process. The signal is modeled as a mixture model, by a linear combination of at most  $K$  features of the form  $\varphi(\theta)$  for some parameters  $\theta \in \Theta$ , where  $\Theta \subseteq \mathbb{R}$  is an interval of parameters and  $\varphi$  is a smooth function defined on  $\Theta$  and taking values in the Hilbert space. We denote by  $(\varphi(\theta), \theta \in \Theta)$  the continuous dictionary.

In order to capture a great variety of examples, we shall assume there exists a Hilbert space  $H_T$ , endowed with the scalar product  $\langle \cdot, \cdot \rangle_T$  and the norm  $\|\cdot\|_T$ , where  $T$  is a parameter belonging to  $\mathbb{N}$ , such that: the observed process  $y$  belongs to  $H_T$ ; for all  $\theta \in \Theta$ , the feature  $\varphi_T(\theta)$  (which may depend on  $T$ ) belongs to  $H_T$  and is non degenerate, *i.e.*  $\|\varphi_T(\theta)\|_T$  is finite and non zero; the noise process  $w_T$ , which might also depend on the parameter  $T$  is a centered Gaussian process belonging to  $H_T$ .

We consider the model with unknown parameters  $\beta^*$  in  $\mathbb{R}^K$  and  $\vartheta^*$  in  $\Theta^K$ :

$$(1) \quad y = \beta^* \Phi_T(\vartheta^*) + w_T \quad \text{in } H_T,$$

where the multivariate function  $\Phi_T$  is defined on  $\Theta^K$  by:

$$\Phi_T(\vartheta) = (\phi_T(\theta_1), \dots, \phi_T(\theta_K))^\top \quad \text{for } \vartheta = (\theta_1, \dots, \theta_K) \in \Theta^K$$

and the function  $\phi_T$  defined on  $\Theta$  is the normalized feature  $\varphi_T(\theta)$ :

$$(2) \quad \phi_T(\theta) = \frac{\varphi_T(\theta)}{\|\varphi_T(\theta)\|_T}.$$

We assume from now on that the unknown  $K$  dimensional vector  $\beta^*$  is sparse, *i.e* it has  $s$  non zero entries or, equivalently,  $\beta^* \in \mathcal{B}_0(s) = \{\beta \in \mathbb{R}^K, \|\beta\|_{\ell_0} = s\}$ , where  $\|\beta\|_{\ell_0}$  counts the number of non zero entries of the vector  $\beta$ . Let  $S^*$  be the support of  $\beta^*$ :

$$S^* = \text{Supp}(\beta^*) = \{k \in \{1, \dots, K\}, \beta_k^* \neq 0\},$$

and call  $s = \text{Card } S^*$  the sparsity parameter. We are interested in predicting observations and in recovering the unknown parameters. Let us denote in general by  $u_S$  the vector  $u$  in  $\mathbb{R}^K$  restricted to the coordinates in  $S$  for any non-empty set  $S \subseteq \{1, \dots, K\}$ . We estimate both the vector  $\beta_{S^*}^*$  with unknown  $s$  and the vector  $\vartheta_{S^*}^*$  with entries in some compact set  $\Theta_T$  containing the parameters of those functions from our continuous dictionary that appear in the mixture model. Note that when applying the same permutation on the coordinates of  $\beta^*$  and the coordinates of  $\vartheta^*$ , we obtain the same model. Thus, the vectors  $\beta^*$  and  $\vartheta^*$  are defined up to such a joint permutation. Moreover, we have  $\beta^* \Phi_T(\vartheta^*) = \beta_{S^*}^* \Phi_T(\vartheta^*)_{S^*}$ , where, by definition,  $\Phi_T(\vartheta^*)_{S^*} = \Phi_T(\vartheta_{S^*}^*)$ . Our model is linear and sparse in  $\beta^*$  but it is non-linear in  $\vartheta^*$ .

We make the following assumption on the noise process  $w_T$ , where the decay rate  $\Delta_T > 0$  controls the noise variance decay as the parameter  $T$  grows and  $\sigma > 0$  is the intrinsic noise level.

**Assumption 1.1** (Admissible noise). *Let  $T \in \mathbb{N}$ . The noise process  $w_T$  belongs to  $H_T$  a.s., and there exist a noise level  $\sigma > 0$  and a decay rate  $\Delta_T > 0$  such that for all  $f \in H_T$ , the random variable  $\langle f, w_T \rangle_T$  is a centered Gaussian random variable satisfying:*

$$(3) \quad \text{Var}(\langle f, w_T \rangle_T) \leq \sigma^2 \Delta_T \|f\|_T^2.$$

In our model, the parameter  $T$  may be understood as the amount of information that we have on the underlying signal.

In order to recover the sparse vector  $\beta^*$  as well as the associated parameters  $\vartheta_{S^*}^*$  (up to a permutation), we solve the following regularized optimization problem with a real tuning parameter  $\kappa > 0$ :

$$(4) \quad (\hat{\beta}, \hat{\vartheta}) \in \underset{\beta \in \mathbb{R}^K, \vartheta \in \Theta_T^K}{\text{argmin}} \quad \frac{1}{2} \|y - \beta \Phi_T(\vartheta)\|_T^2 + \kappa \|\beta\|_{\ell_1},$$

where the smooth function  $\Phi_T$  is defined on the set  $\Theta_T^K$ , with  $\Theta_T$  a compact interval. Therefore the existence of at least a solution is guaranteed. The functional that we minimize in this problem is composed of a data fidelity term and a penalty term. The penalty is expressed with a  $\ell_1$ -norm on the vector  $\beta = (\beta_1, \dots, \beta_K)$ , *i.e* the sum of the absolute values of its coordinates:  $\|\beta\|_{\ell_1} = \sum_{i=1}^K |\beta_i|$ . This penalization is similar to that of the Lasso problem (also referred to as Basis pursuit) introduced in [48] and extensively studied since then (see [13] for a comprehensive survey). The optimization of the non-linear parameters is not performed on the whole set of parameters  $\Theta$  but rather on a compact subset  $\Theta_T$  indexed by the parameter  $T$ . Indeed, it may be necessary to restrict the set of parameters, *e.g.* in a finite mixture model where we consider a location parameter we can only recover those parameters within the support of the observations.

In the more general Beurling Lasso (BLasso) framework, one can rewrite the problem (4) in a measure setting. The actual solution  $(\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_K), \hat{\vartheta} = (\hat{\vartheta}_1, \dots, \hat{\vartheta}_K))$  of (4) is then seen as the atomic measure  $\hat{\mu} = \sum_{k=1}^K \hat{\beta}_k \delta_{\hat{\vartheta}_k}$ , where the amplitudes and the locations of the Dirac masses correspond respectively to the linear coefficients in the mixture and the parameters of the features. The measure  $\hat{\mu}$  is also a solution of the BLasso problem when the latter admits atomic solutions composed of less than  $K$  atoms. This is in particular the case in the discrete-time model, with  $T$  design points, presented in Section 1.2.1 where  $K \geq T$  according to [11]. However, to the best of our knowledge, there are no such results when  $H_T$  is a general Hilbert space.

## 1.2. Examples

In this section we give examples of both discrete and continuous-time models that are covered by our general setup. We discuss how  $T$  indicates the amount of information that the data contain on the unknown underlying signal. Indeed, in the discrete case, the amount of information grows as the number  $T$  of the design points over which the process is observed increases, while the largest step-size decreases; in the continuous case, it grows as the decay rate  $\Delta_T$  of the noise variance decreases.

We emphasize the various structures of noise processes that are admissible by giving several examples of discrete or continuous-time noise processes that satisfy our assumptions. They are frequently used in discrete regression models or continuous models like the Gaussian white noise model, see [50] or [31].

### 1.2.1. Discrete-time models

Consider a real-valued process  $y$  observed over the points  $t_1 < \dots < t_T$  on  $[0, 1]$ , with  $T \in \mathbb{N}^*$ . Let  $H_T = L^2(\lambda_T)$  be the Hilbert space of real valued functions defined on  $[0, 1]$  and square integrable with respect to some probability measure  $\lambda_T$  on  $\{t_1, \dots, t_T\}$ . Let the noise  $w_T \in H_T$  be given by  $w_T(t) = \sum_{j=1}^T G_j \mathbf{1}_{\{t_j\}}(t)$ , where  $G_1, \dots, G_T$  are centered Gaussian random variables and  $\mathbf{1}_A$  denotes the indicator function of an arbitrary set  $A$ . Thus, the observations are:

$$(5) \quad y(t_j) = \sum_{k \in S^*} \beta_k^* \cdot \phi_T(\theta_k^*, t_j) + G_j, \quad j = 1, \dots, T.$$

The risk is measured by:

$$\|y - \beta \Phi_T(\vartheta)\|_T^2 = \sum_{j=1}^T \left( y(t_j) - \sum_{k=1}^K \beta_k \cdot \phi_T(\theta_k, t_j) \right)^2 \lambda_T(t_j).$$

Now, let  $\lambda_T = \Delta_T \sum_{j=1}^T \delta_{t_j}$ , where  $\delta_x$  denotes the Dirac mass at  $x$ . In the particular case where  $\Delta_T = 1/T$ , one can approximate the measure  $\lambda_T$  for  $T$  large by the Lebesgue measure on  $[0, 1]$ , say  $\text{Leb}$ . In various examples, it is also easier to compute the norms of the features and of their derivatives in the Hilbert space  $L^2(\text{Leb})$ . This amounts to seeing  $H_T$  as approximating Hilbert spaces of the fixed Hilbert space  $L^2(\text{Leb})$ .

Let us now see that, if the noise variables  $G_1, \dots, G_T$  are independent centered Gaussian random variables with variance  $\sigma^2$ , then Assumption 1.1 holds with an equality:

$$\text{Var}(\langle f, w_T \rangle_T) = \sigma^2 \Delta_T \|f\|_T^2.$$

If  $(G_1, \dots, G_T)$  is a centered Gaussian vector of dimension  $T$  and covariance matrix with each diagonal entry  $\sigma^2$ , then Assumption 1.1 holds with  $\Delta_T$  multiplied by the spectral radius  $\varrho_T \in [1, T]$  of the correlation matrix:

$$\text{Var}(\langle f, w_T \rangle_T) \leq \sigma^2 \Delta_T \varrho_T \|f\|_T^2.$$

### 1.2.2. Continuous-time models with truncated white noise or colored noise

Consider the set  $\mathcal{C} = \mathcal{C}([0, 1], \mathbb{R})$  of  $\mathbb{R}$ -valued continuous functions defined on  $[0, 1]$ , an orthonormal base  $(\psi_j, j \in \mathbb{N})$  of  $L^2 = L^2([0, 1], \text{Leb})$  of elements of  $\mathcal{C}$ , where  $\text{Leb}$  is the Lebesgue measure on  $[0, 1]$ . We simply denote by  $\langle \cdot, \cdot \rangle_{L^2}$  the corresponding scalar product. Let  $p = (p_j, j \in \mathbb{N})$  be a sequence of non-negative real numbers and set  $\text{Supp}(p) = \{j \in \mathbb{N} : p_j > 0\}$  its support. Let  $H_T$  be the completion of the vector space generated by the base  $(\psi_j, j \in \text{Supp}(p))$  (which is also the completion of  $\mathcal{C}$  if  $p$  is positive and bounded), with respect to the scalar product:

$$\langle f, g \rangle_T = \sum_{j \in \mathbb{N}} p_j \langle f, \psi_j \rangle_{L^2} \langle g, \psi_j \rangle_{L^2}.$$

Notice that the Hilbert space  $H_T$  does not depend on the parameter  $T$  unless  $p$  depends on  $T$ . Let us recall that if  $p \equiv 1$ , that is, the sequence  $p$  is constant equal to 1, then  $H_T = L^2$ . In this model we observe a continuous path:

$$(6) \quad y(t) = \sum_{k \in S^*} \beta_k^* \phi_T(\theta_k^*, t) + w_T(t), \quad t \in [0, 1].$$

The risk is measured by:

$$\|y - \beta \Phi_T(\theta)\|_T^2 = \sum_{j \in \mathbb{N}} p_j \left( \int_0^1 (y(t) - \beta \Phi_T(\theta, t)) \cdot \psi_j(t) dt \right)^2.$$

Let  $\xi = (\xi_j, j \in \mathbb{N})$  be a weight sequence of non-negative real numbers such that the sequence  $p \circ \xi := (p_j \xi_j, j \in \mathbb{N})$  is summable. Consider the noise  $w_T = \sum_{j \in \text{Supp}(p)} \sqrt{\xi_j} G_j \psi_j$ , where  $(G_j, j \in \mathbb{N})$  are independent centered Gaussian random variables with variance  $\sigma^2$ . Notice Assumption 1.1 holds as  $\|w_T\|_T^2 = \sum_{k \in \mathbb{N}} p_j \xi_j G_j^2$  is a.s. finite and, with  $\Delta_T = \sup_{\mathbb{N}} p \circ \xi$ :

$$\text{Var}(\langle f, w_T \rangle_T) = \sigma^2 \sum_{j \in \mathbb{N}} p_j^2 \xi_j \langle f, \psi_j \rangle_{L^2}^2 \leq \sigma^2 \Delta_T \|f\|_T^2.$$

Notice that the noise  $w_T$  does not depend on the parameter  $T$  unless  $p$  or  $\xi$  depends on  $T$ .

The truncated white noise model corresponds to  $p \equiv 1$  and  $\xi = (\xi_k = \mathbf{1}_{\{j \leq T\}}, j \in \mathbb{N})$ . In this case  $\Delta_T = 1$  and  $\|w_T\|_T^2$  is a.s. of order  $\sigma^2 T$  by the strong law of large numbers. The white noise corresponds to the limit case  $T = +\infty$ , which does not satisfy the hypothesis as a.s. its  $L^2$ -norm is infinite. Let us mention that the bounds given in the main theorems in Section 2 rely on  $\|w_T\|_T$  being finite and not on its value.

Consider again  $p \equiv 1$ . Thanks to the Karhunen-Loève's decomposition, the scaled Brownian motion  $w_T = C_T B$ , with  $B$  the Brownian motion on  $[0, 1]$  and  $C_T$  a positive constant, corresponds to the orthonormal base functions  $\psi_k(t) = \sqrt{2} \sin((2k+1)\pi t/2)$  for  $t \in [0, 1]$  and the weights  $\xi_k = 4C_T^2/(2k+1)^2\pi^2$  for  $k \in \mathbb{N}$ , and  $\sigma^2 = 1$ . In this case, we have  $\langle f, w_T \rangle_T = C_T \int_0^1 f(s)B(s)ds$  for  $f \in L^2$  and Assumption 1.1 holds with  $\sigma^2 = 1$  and  $\Delta_T = \sup_{\mathbb{N}} p \circ \xi = 4C_T^2/\pi^2$ .

### 1.3. Previous work

The model (1) in the particular case where  $\vartheta^*$  is supposed given and the observations depend linearly on a vector  $\beta^*$  has long been studied in the literature. Assume for simplicity that  $H_T = \mathbb{R}^T$  is the  $T$ -dimensional Euclidean space, so that  $\Phi_T \in \mathbb{R}^{K \times T}$  is a matrix whose entries are known and can be either random or deterministic,  $y \in \mathbb{R}^T$  is an observed vector and  $w_T \in \mathbb{R}^T$  is a vector of noise (often assumed to be Gaussian). Even when  $K$  is larger than  $T$  the estimation of  $\beta^*$  is still consistent provided the vector  $\beta^*$  is sparse and a null space property is verified by the matrix  $\Phi_T$ , or some sufficient condition saying that the lines of  $\Phi_T$  are not too colinear (see [51] for a complete overview). The Lasso estimator [48] or the Dantzig selector [15] are efficient to perform such estimation and the quality of the estimation with respect to the dimension of the problem is now well known. The authors of [9] have given bounds for the prediction error for both estimators.

We consider here a highly non-linear extension of this model that consists in assuming that the matrix  $\Phi_T = \Phi_T(\vartheta^*)$  depends non-linearly on a parameter  $\vartheta^*$  to be estimated. In our model (1),  $\Phi_T$  is composed of  $K$  row vectors belonging to a parametric family or by  $K$  features belonging to a continuous dictionary and the observed data  $y$  may be either a vector or a function. This model has proven to be relevant in many fields such as microscopy, astronomy, spectroscopy, imaging or signal processing.

When the observation  $y$  belongs to a finite-dimensional Hilbert space and the dimension  $K$  is fixed and small compared to  $T$ , the model received attention several decades ago and gave rise to separable least square problems and resolution methods such as variable projection (see [33, 34]). These papers mainly provided numerical methods but let us mention the consistency result in [35] for non-linear regression models.

On the contrary, when  $K$  is arbitrarily large many problems remain open. One of the natural ideas to estimate the underlying parameters could be to discretize the parameter space  $\Theta$  and return to the study of a linear model. It would amount to considering a finite subfamily of  $(\varphi(\theta), \theta \in \Theta)$  as in [46] and deal with overcomplete dictionary learning techniques (also referred to as sparse coding, see [25, 40]). In this case, sparse estimators for linear models such as the Lasso are available. However, in sparse spike deconvolution where the family  $(\varphi(\theta), \theta \in \Theta)$  is a family of spikes parametrized by a location parameter, the authors of [27] have shown that in the presence of noise discretizing the space of parameters and solving a Lasso problem tends to produce clusters of spikes around the spikes one seeks to locate. That is why it is preferable to use off-the-grid methods. By off-the-grid, we mean that the methods employed do not use discretization schemes on the parameter set  $\Theta$ . In [26], the authors show that in presence of a small noise, the BLasso only induces a slight perturbation of the spikes locations and amplitudes and does not produce clusters. The BLasso was introduced in [23] and has been studied in many papers since then mostly by the compressed sensing and super-resolution communities (see [17], [5] among many others). It is basically an off-the-grid extension of the classical Lasso for continuous dictionary learning. The optimization problem is formulated as a convex minimization over the space of Radon measures. In the BLasso framework, the dimension  $K$  in (1) is infinite and the linear coefficients and non-linear parameters are encoded by an atomic measure made of weighted Dirac functions. By solving a minimization problem over Radon measures, the aim is to recover an atomic measure. It raises the question of whether such a solution exists. In [11] the question is answered by the affirmative when the observed data  $y$  belongs to a finite-dimensional Hilbert space  $H_T$ . When this is not the case, i.e.  $H_T$  is infinite dimensional, the question is open. In this paper, we avoid the problem by assuming a bound  $K$  on the number of functions in the mixture and restricting the space over which the BLasso is performed to the atomic measures with at most  $K$  atoms. The numerical methods used to solve the BLasso such as the Sliding Frank-Wolfe algorithm (see [24] and [14, 32] for applications in spectroscopy and imaging), also called the alternating descent conditional gradient method (see [10]), and the conic particle gradient descent (see [21]), seek a solution directly in the space of Dirac mixtures. Hence, our formulation (4) is closer to the way algorithms proceed. Let us mention that other methods such as Orthogonal Matching Pursuit (see [28]) exist to tackle the problem of sparse learning from a continuous dictionary. Typically, the case of sparse spike deconvolution where the dictionary consists of Gaussian functions continuously parametrized by a location parameter is not included.

The study of the regression over a continuous dictionary in the framework of the BLasso has been quite specific to the dictionary considered. The literature first focused on the dictionary of complex exponential functions parametrized by their frequency ( $\varphi(\theta) : t \mapsto e^{i2\pi\langle t, \theta \rangle}$ ,  $\theta \in \Theta$ ) where  $\Theta$  is the  $d$ -dimensional torus (see [18]). In [12], a bound is given for the prediction error for this dictionary. The proof extends a previous result obtained in [47] for atomic norm denoising. What is particularly interesting is that the rates obtained for the prediction error almost reach the minimax rates achievable for linear models (see [16, 42]) provided that the frequencies are sufficiently separated. The separation condition between the non-linear parameters to estimate is inherent to the BLasso unless we assume the positivity of the linear parameters as in [44].

For results on a wider range of dictionaries, let us highlight the work of [26] that gives recovery and robustness to noise results for spike deconvolution. Let us also mention the recent work of [8] that generalizes some exact recovery results for a broader family of dictionaries as well as the paper [7] that gives robustness to noise guarantees for a family of shifted functions ( $\varphi(\theta) = k(\cdot - \theta)$ ,  $\theta \in \Theta$ ) of a given specific function  $k$ . In a density model that is a mixture of shifted functions, [22] studies a modification of the BLasso by considering a weighted  $L^2$  prediction error.

The case of non-translation invariant families remained for long intractable without very pessimistic separation conditions. In [41] the authors set a natural geometric framework to analyse the estimation problem. The separation condition between the parameters appears naturally in terms of a metric. In their paper, the design over which the observation are made is distributed according to a probability distribution. Their main result shows that in presence of noise the BLasso recovers a measure close to the one to be estimated with respect to a Wasserstein metric.

#### 1.4. Contributions

This paper addresses the problem of learning sparse mixtures from a continuous dictionary for a wide variety of regression models within a common framework. Indeed, we tackle a wide range of possible dictionaries of sufficiently smooth features, observation schemes and Gaussian noises with various structures. The observations are supposed to belong to a Hilbert space  $H_T$ . Continuous observations over an interval of  $\mathbb{R}$  as well as discrete observations at given design points are therefore included in our framework. Furthermore, the Hilbert structure and the mild assumption we make on the noise, encompass a wide range of Gaussian noises. In particular, our framework allows to take into account the case of correlated Gaussian noise processes.

The main results of this paper gives a high-probability bound for the prediction error:

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_T,$$

where  $(\hat{\beta}, \hat{\vartheta})$  is the solution of the optimization problem (4). Contrary to the BLasso optimization program over a set of measures whose result can be a diffuse measure, our formulation of the optimization problem has always a solution belonging to a finite set of values. Our prediction error bound matches (up to logarithmic factors and with high probability) that obtained in the linear case, that is when  $\vartheta^*$  is known and does not need to be estimated. We also give high-probability bounds on some loss functions comparing the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  given by (4) to the parameters  $\beta^*$  and  $\vartheta^*$ , respectively. Our work extends results that were so far restricted to the specific case of a dictionary consisting of complex exponentials continuously parameterized by their frequencies (see [12, 47]). When the optimization problem produces a cluster of features to approximate an element of the mixture, we also show that there can be no compensation between the amplitudes of the features involved.

Following works in compressed sensing and super-resolution (see [17, 18] among others), our bounds rely on the existence of interpolating functions called “certificates” (see Assumptions 6.1 and 6.2) instead of relying on compatibility conditions or Restricted Eigenvalue conditions. We give in Section 7 sufficient conditions for the existence of certificates and an explicit way to construct such functions in the spirit of [41]. We show in this paper that such functions can be constructed provided the non-linear parameters belonging to  $\Theta$  are well separated with respect to a Riemannian metric  $\mathfrak{d}_T$  (defined in Section 4.1) associated to the kernel  $\mathcal{K}_T(\theta, \theta') = \langle \phi_T(\theta), \phi_T(\theta') \rangle_T$ . This minimal separation distance between the non-linear parameters needs to be rather large, comparable to  $s$ , in a general context. However, it can be significantly reduced to a constant order in more particular cases such as the sparse spike deconvolution, see Remark 8.2. The Riemannian metric appears naturally when it comes to tackle a wide variety of dictionaries. In addition, it leads to a lot of invariances in many quantities useful in the proofs. Typically, the Riemannian metrics  $\mathfrak{d}_T$  and  $\mathfrak{d}_T^h$  associated respectively to the kernel  $\mathcal{K}_T(\cdot, \cdot)$  and the warped kernel  $\mathcal{K}_T^h = \mathcal{K}_T(h(\cdot), h(\cdot))$  for some smooth enough diffeomorphism  $h$  are equal and we have  $\mathfrak{d}_T(\theta, \theta') = \mathfrak{d}_T^h(h^{-1}(\theta), h^{-1}(\theta'))$ .

Our statistical results rely on tail bounds for suprema of Gaussian processes: following [12], instead of using controls on  $\|w_T\|_T$  as in the seminal works [26, 41], we used bounds, based on the noise structure from Assumption 1.1, on quantities of the form  $\sup_{\Theta_T} \langle f(\theta), w_T \rangle_T$  for some  $H_T$ -valued functions  $f$  built from the dictionary  $(\varphi_T(\theta), \theta \in \Theta)$  and its derivative. This approach is relevant as for some models the quantity  $\|w_T\|_T$  may be very large, see for example the truncated white noise model from Section 1.2.2. We note that the nonlinear parameter  $\theta$  is univariate in our setup. Generalization to multivariate non-linear parameters is possible, but highly technical. Indeed, the construction of the certificates holds in the multivariate setting, but the exponential bounds for suprema of Gaussian fields are less precise concerning their dependence on the dimension.

We give next two applications of our results respectively to the Gaussian sparse spike deconvolution and to the Scaled exponential model also known as Laplace transform inversion. They illustrate how the stringent assumptions in all generality, become less restrictive in more precise setups. The full derivation of these examples can be found in Sections 8 and 9, respectively.

#### 1.4.1. Gaussian sparse spike deconvolution, see Section 8.

Consider the discrete-time model (5) as described in section 1.2.1, where a process  $y$  is observed over a regular grid  $t_1 < \dots < t_T$  on the interval  $[a_T, b_T]$  with step size  $\Delta_T = (b_T - a_T)/T$ , where  $T \in \mathbb{N}^*$ ,  $b_T = -a_T = \sigma_0 \sqrt{\log(T)}$  and  $\sigma_0 > 0$  is some fixed scale factor. Assume the observations are corrupted by independent centered Gaussian random variables of variance  $\sigma^2$ .

The dictionary consists of Gaussian spikes that are continuously translated:

$$\left( \varphi(\theta) = \exp \left( -\frac{(\theta - \cdot)^2}{2\sigma_0^2} \right), \quad \theta \in \mathbb{R} \right).$$

This model can be viewed as a non-linear extension of the Gaussian sequence model, where the mean vector is a linear combination of shifted Gaussian spikes. We are interested in recovering the unknown shift parameters  $(\theta_k^*)_{1 \leq k \leq s}$  belonging to the compact set  $\Theta_T = [(1 - \epsilon)a_T, (1 - \epsilon)b_T] \subset [a_T, b_T]$ , where  $\epsilon$  is a given positive shrinkage, as well as the unknown linear parameters  $\beta^*$ .

We apply our main result, Theorem 2.1, which gives that: if the number of observations  $T$  is sufficiently large (depending on  $\sigma_0$ ,  $\epsilon$  and the sparsity  $s$ ) and if the shift parameters are separated, i.e. such that for all  $\ell \neq k$ ,  $|\theta_k^* - \theta_\ell^*| \gtrsim \sigma_0$ , the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in the minimization problem (4) using the regularization weight  $\kappa = \mathcal{C}\sigma\sqrt{\Delta_T \log(T)}$  achieve the following prediction error bound:

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_T \leq \mathcal{C}' \sigma \sqrt{s \frac{\log(T)}{T}},$$

with probability greater than  $1 - \mathcal{C}'' T^{-\gamma}$ , for some  $\gamma > 0$ , where  $\mathcal{C}/\sqrt{\gamma}$ ,  $\mathcal{C}'/\sqrt{\gamma}$  and  $(\sqrt{\gamma} \wedge 1)\mathcal{C}''$  are some universal constants and  $\|f\|_T = \frac{1}{\sqrt{T}} \sqrt{\sum_{j=1}^T f(t_j)^2}$ . See Remark 8.4 for details, with  $\gamma' = \gamma$  therein.

#### 1.4.2. Scaled exponential model, see Section 9.

Consider the continuous time model (6) where the real-valued process  $y$  is observed on  $\mathbb{R}_+$  and assume that this process is an element of the Hilbert space  $H_T = L^2(\mathbb{R}_+, \text{Leb})$  where  $\text{Leb}$  denotes here the Lebesgue measure over  $\mathbb{R}_+$ . We write  $H$  instead of  $H_T$  for the Hilbert space and we write  $\langle \cdot, \cdot \rangle$  its scalar product and  $\|\cdot\|$  its associated norm.

Let the noise process be a truncated white noise as in Section 1.2.2 such that  $w_T = \sum_{k=1}^T (1/\sqrt{T}) G_k \psi_k$ , where  $(G_k, k \in \mathbb{N})$  are independent centered Gaussian random variables with variance  $\sigma^2$  and  $(\psi_k, k \in \mathbb{N})$  denotes an orthonormal basis of  $H$ . We stress the fact that by the law of large numbers  $\|w_T\|^2$  tends almost surely to  $\sigma^2 > 0$ . Therefore the upper bounds from previous results on super-resolution and BLasso (see [26] or [41]) do not apply here, as they hold for noise processes having  $\|w_T\|$  tending to zero.

Let the dictionary consist of the exponential functions :

$$(\varphi(\theta) = \exp(-\theta \cdot), \quad \theta \in \mathbb{R}_+^*).$$

We aim at recovering the unknown scale parameters  $(\theta_k^*)_{1 \leq k \leq s}$  belonging to a compact set whose diameter may depend on  $T \in \mathbb{N}^*$ , say  $\Theta_T = [T^{-\gamma}, T^\gamma]$ , with  $\gamma > 0$ , as well as the unknown linear parameters  $\beta^*$ .

We apply our main result, Theorem 2.1, which gives that: if the scale parameters are separated, i.e. such that for all  $\ell \neq k$ ,  $|\log(\theta_k^*/\theta_\ell^*)| \gtrsim 1$ , the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in the minimization problem (4), using the regularization weight



$\kappa = \mathcal{C} \sigma \sqrt{\log(T)/T}$  achieve the following prediction bound:

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\| \leq \mathcal{C}' \sigma \sqrt{s \frac{\log(T)}{T}},$$

with probability larger than  $1 - \mathcal{C}'' T^{-\gamma} (1 \vee \sqrt{\gamma \log(T)})$ , where  $\mathcal{C}/\sqrt{\gamma}$ ,  $\mathcal{C}'/\sqrt{\gamma}$  and  $\mathcal{C}''$  are some universal constants. See Remark 9.4 for details, with  $\gamma' = \gamma$  therein.

## 2. Main Results

Recall that we consider the model (1) that we can write in an equivalent way as:

$$y = \sum_{j \in S^*} \beta_j^* \frac{\varphi_T(\theta_j^*)}{\|\varphi_T(\theta_j^*)\|_T} + w_T \quad \text{in } H_T,$$

with  $S^*$  the support of the vector  $\beta^*$ . The main theorem of this paper gives the behavior of the prediction error with respect to: the decay rate of the noise variance  $\Delta_T$ , the parameter  $T \in \mathbb{N}$ , the sparsity  $s \in \mathbb{N}^*$ , the upper bound on the number of components in the mixed signal  $K$  and the intrinsic noise level  $\sigma$ . We shall consider assumptions on the regularity of the dictionary  $\varphi_T$ , on the parameter space  $\Theta_T$  on which the optimization is performed and on the noise  $w_T$ . Using the features  $\varphi_T$  we build a kernel  $\mathcal{K}_T$  on the space of parameters  $\Theta$  and an associated Riemannian metric  $\mathfrak{d}_T$ , see Section 4, which is the intrinsic metric, rather than the usual Euclidean metric. More assumptions are necessary on the closeness of the kernel  $\mathcal{K}_T$  and its derivatives defined in (29) to a limit kernel  $\mathcal{K}_\infty$  and its derivatives.

The theorem is stated assuming the existence of certificate functions, see Assumptions 6.1 and 6.2. Sufficient conditions for their existence are given later in Section 7, in which Propositions 7.4 and 7.6 show that the limit kernel  $\mathcal{K}_\infty$  must be uniformly bounded and have concavity properties. In this case, the existence of certificates stands provided the underlying non-linear parameters to be estimated are sufficiently separated according to the Riemannian metric  $\mathfrak{d}_T$ , see Condition (iii) in Propositions 7.4 and 7.6.

In the following result the parameter set  $\Theta_T$  is a one dimensional compact interval. We note  $|\Theta_T|_{\mathfrak{d}_T}$  its length with respect to the Riemannian metric  $\mathfrak{d}_T$  on  $\Theta^2$  associated to the kernel  $\mathcal{K}_T$ .

**Theorem 2.1.** *Assume we observe the random element  $y$  of  $H_T$  under the regression model (1) with unknown parameters  $\beta^*$  and  $\vartheta^* = (\theta_1^*, \dots, \theta_K^*)$  a vector with entries in  $\Theta_T$ , a compact interval of  $\mathbb{R}$ , such that:*

- (i) **Admissible noise:** *The noise process  $w_T$  satisfies Assumption 1.1 for a noise level  $\sigma > 0$  and a decay rate for the noise variance  $\Delta_T > 0$ .*
- (ii) **Regularity of the dictionary  $\varphi_T$ :** *The dictionary function  $\varphi_T$  satisfies the smoothness conditions of Assumption 3.1. The function  $g_T$  defined in (14), satisfies the positivity condition of Assumption 3.2.*
- (iii) **Regularity of the limit kernel:** *The kernel  $\mathcal{K}_\infty$  and the functions  $g_\infty$  and  $h_\infty$ , defined on an interval  $\Theta_\infty \subset \Theta$ , see (16) and (33), satisfy the smoothness conditions of Assumption 5.1.*
- (iv) **Proximity to the limit kernel:** *The kernel  $\mathcal{K}_T$  defined from the dictionary, see (29), is sufficiently close to the limit kernel  $\mathcal{K}_\infty$  in the sense that Assumption 5.2 holds.*
- (v) **Existence of certificates:** *The set of unknown parameters  $\mathcal{Q}^* = \{\theta_k^*, k \in S^*\}$ , with  $S^* = \text{Supp}(\beta^*)$ , satisfies Assumptions 6.1 and 6.2 with the same  $r > 0$ .*

*Then, there exist finite positive constants  $\mathcal{C}_0, \mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$  depending on the kernel  $\mathcal{K}_\infty$  defined on  $\Theta_\infty$  and on  $r$  such that for any  $\tau > 1$  and a tuning parameter:*

$$\kappa \geq \mathcal{C}_1 \sigma \sqrt{\Delta_T \log \tau},$$

*we have the prediction error bound of the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in (4) given by:*

$$(7) \quad \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_T \leq \mathcal{C}_0 \sqrt{s} \kappa,$$

*with probability larger than  $1 - \mathcal{C}_2 \left( \frac{|\Theta_T|_{\mathfrak{d}_T}}{\tau \sqrt{\log \tau}} \vee \frac{1}{\tau} \right)$ . Moreover, with the same probability, the difference of the  $\ell_1$ -norms of  $\hat{\beta}$  and  $\beta^*$  is bounded by:*

$$(8) \quad \left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| \leq \mathcal{C}_3 \kappa s.$$

This result holds for both the continuous and discrete settings described in Section 1.2, covers a wide range of smooth dictionaries, and is proven under mild assumptions on the noise. We discuss in the next remark that the prediction error is, up to a logarithmic factor, almost optimal.

*Remark 2.2* (Comparison with the Lasso estimator). Let us consider the discrete-time model where the observation space is the Hilbert space  $H_T = \mathbb{R}^T$  endowed with the Euclidean norm  $\|\cdot\|_{\ell_2}$ . The observation  $y \in \mathbb{R}^T$  comes from the model (1) where the noise is a Gaussian vector with independent entries of variance  $\sigma^2$ . In this setting, the decay rate of the noise variance is fixed with  $\Delta_T = 1$ .

We first consider that the parameters  $\vartheta^*$  are known. In this case, the model becomes the classical high-dimensional regression model and the Lasso estimator  $\hat{\beta}_L$  can be used to estimate  $\beta^*$  under coherence assumptions on the finite dictionary made of the rows of the matrix  $\Phi^* = \Phi_T(\vartheta^*)$  (see [9]). The behavior of the Lasso estimator has been studied in the literature and its prediction risk tends to zero at the rate:

$$(9) \quad \frac{1}{T} \|(\hat{\beta}_L - \beta^*)\Phi^*\|_{\ell_2}^2 = \mathcal{O}\left(\frac{\sigma^2 s \log(K)}{T}\right)$$

with high probability, larger than  $1 - 1/K^\gamma$  for some positive constant  $\gamma > 0$ . Furthermore, in the case where  $\beta^*$  is an unknown  $s$ -sparse vector,  $\vartheta^*$  is known and  $\Phi^*$  verifies a coherence property, then the lower bounds of order  $\sigma^2 s \log(K/s)/T$  in expected value can be deduced from the more general bounds for group sparsity in [38] (see also [42]). The non-asymptotic prediction lower bounds for the prediction error given in [42] are:

$$\inf_{\hat{\beta}} \sup_{\beta^* \text{ } s\text{-sparse}} \mathbb{E} \left[ \frac{1}{T} \|(\hat{\beta} - \beta^*)\Phi^*\|_{\ell_2}^2 \right] \geq C \cdot \frac{\sigma^2 s \log(K/s)}{T},$$

where the infimum is taken over all the estimators  $\hat{\beta}$  (square integrable measurable functions of the observation  $y$ ) and for some constant  $C > 0$  free of  $s$  and  $T$ . When the parameters  $\vartheta^*$  are unknown, Theorem 2.1 gives an upper bound for the prediction risk which is, up to a logarithmic factor, almost the best rate we could achieve even knowing the non-linear parameters  $\vartheta^*$ . Consider the estimators in (4) where the Riemannian diameter of the set  $\Theta_T$  is bounded by a constant free of  $T$  (this is the case of Example 5.1 below). By squaring (7) and then dividing it by  $T$ , we obtain from Theorem 2.1 with  $\kappa = C_1 \sigma \sqrt{\Delta_T \log \tau}$  and  $\tau = T^\gamma$  for some given  $\gamma > 0$ , that with high probability, larger than  $1 - C/T^\gamma$ :

$$(10) \quad \frac{1}{T} \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{\ell_2}^2 = \mathcal{O}\left(\frac{\sigma^2 s \log(T)}{T}\right).$$

Let us mention that [47] also obtained a similar prediction error (10) for the specific dictionary given by the complex exponential functions  $(\varphi(\theta) : t \mapsto e^{i2\pi t\theta}, \theta \in \Theta = [0, 2\pi])$ ; notice that the proof therein uses the Parseval's identity for Fourier series as well as Markov-Bernstein type inequalities for trigonometric polynomials. Even if the structure of our proof is in the spirit of [47], our result is more general and does not rely on the convex setting of the BLasso approach.

*Remark 2.3* (Proximity to the limit kernel). We comment on Condition (iv) on the proximity of the kernels  $\mathcal{K}_T$  and  $\mathcal{K}_\infty$ , which also appears as Conditions (iv)-(v) in Proposition 7.4 (and similarly as Condition (iv) in Proposition 7.6).

In the examples of Sections 3.2.2 and 3.2.4 on translation or scaling model with a continuum of observations, the parameter  $T$  does not play any role in the definition of  $\mathcal{K}_T$ , so that one can take  $\mathcal{K}_\infty$  equal to  $\mathcal{K}_T$ . In this case, the proximity conditions on the kernels are trivially satisfied.

The example from Section 8 is devoted to the Gaussian sparse spike deconvolution, that is, to a mixture of Gaussian translation invariant features observed in a discrete regression model on a regular grid of size  $T$ . In this case, we built a family of models  $(H_T, \varphi_T, w_T, \Theta_T)$  with a dictionary  $\varphi_T$  which does not depend on  $T$  and such that the kernel  $\mathcal{K}_T$  and its derivatives converge to  $\mathcal{K}_\infty$  (and also  $\rho_T$  from (35) converges to 1). In this setting, the proximity condition of Theorem 2.1 holds for  $T$  large enough, say  $T$  larger than some  $T_0$  which depends on  $\mathcal{K}_\infty$ , see Assumption 5.2. The existence of the certificates, see Propositions 7.4 and 7.6, also requires a proximity criterion which is achieved for  $T$  large enough, say  $T$  larger than some  $T_1$  which depends on  $\mathcal{K}_\infty$  and is increasing with the sparsity parameter  $s$  (see for example Condition (v) in Proposition 7.4).

*Remark 2.4* (On the dimension  $K$ , the upper bound of the sparsity). We remark that neither the bound on the prediction error nor the probability on which the bound holds, depends on the upper bound  $K$  on the sparsity  $s$ . Therefore, the value of  $K$  can be taken arbitrarily large. It is not surprising that  $K$  does not have any impact on the bound since the optimisation problem (4) could be formulated without any bound on the sparsity. Indeed, the problem (4) can be embedded in an optimization problem over a space of measures following the literature on the BLasso introduced in [23]. See also Remark 2.6.



The next theorem gives bounds on the differences between the parameters  $\hat{\beta}$  given by the optimization problem (4) and the “true” parameters  $\beta^*$  for active features having their parameter  $\hat{\theta}_\ell$  close, with respect to the Riemannian metric  $\mathfrak{d}_T$ , to a parameter  $\theta_k^*$ , with  $k$  in  $S^*$ . For  $r > 0$  given by Assumptions 6.1 and 6.2, we define:

- The support of  $\hat{\beta}$  given by the optimization problem (4):  $\hat{S} = \text{Supp}(\hat{\beta}) = \{\ell : \hat{\beta}_\ell \neq 0\}$ .
- The near region  $\tilde{S}(r)$  given by:

$$\tilde{S}(r) = \bigcup_{k \in S^*} \tilde{S}_k(r) \quad \text{where} \quad \tilde{S}_k(r) = \left\{ \ell \in \hat{S} : \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*) \leq r \right\},$$

which corresponds to the set of indices  $\ell$  in the support of  $\hat{\beta}$  such that the corresponding parameter  $\hat{\theta}_\ell$  is close to one of the true parameters  $\theta_k^*$ , for some  $k \in S^*$ .

The set  $\hat{S} \setminus \tilde{S}(r)$  is also called the far region. Notice that the sets  $\tilde{S}_k(r)$  with  $k \in S^*$  are pairwise disjoint under Assumption 6.1, and that they can be empty. In what follows, we use the convention  $\sum_{\emptyset} = 0$ .

**Theorem 2.5.** *We consider the model in Theorem 2.1 and suppose that Assumptions (i)-(v) therein hold. Then, there exist finite positive constants  $C_1, C_2, C_3, C_4, C_5$  and  $C_6$  depending on  $\mathcal{K}_\infty$  defined on  $\Theta_\infty$  and on  $r$  such that for any  $\tau > 1$  and a tuning parameter:*

$$\kappa \geq C_1 \sigma \sqrt{\Delta_T \log \tau}$$

the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in (4) satisfy the following bounds with probability larger than  $1 - C_2 \left( \frac{|\Theta_T| \mathfrak{d}_T}{\tau \sqrt{\log \tau}} \vee \frac{1}{\tau} \right)$ :

$$(11) \quad \sum_{k \in S^*} \left| |\beta_k^*| - \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \right| \leq C_3 \kappa s, \quad \sum_{k \in S^*} \left| \beta_k^* - \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \right| \leq C_4 \kappa s \quad \text{and} \quad \left\| \hat{\beta}_{\tilde{S}(r)^c} \right\|_{\ell_1} \leq C_5 \kappa s,$$

$$(12) \quad \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} \left| \hat{\beta}_\ell \right| \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*)^2 \leq C_6 \kappa s,$$

where for a subset  $S$  of  $\mathcal{I} = \{1, \dots, K\}$ , the set  $S^c$  denotes the complementary set of  $S$  in  $\mathcal{I}$ , that is  $\mathcal{I} \setminus S$ .

Notice that each linear parameter  $\beta_k^*$  can be estimated by the sum of several linear coefficients  $\hat{\beta}_\ell$  with  $\ell \in \{1, \dots, K\}$ . The first two inequalities in (11) show that there can be no compensation between the estimators  $\hat{\beta}_\ell$  that approximate the same  $\beta_k^*$  with  $k \in S^*$ , meaning that there can be no large values of  $\hat{\beta}_\ell$  having different signs that sum up to a possibly small (in absolute value) true  $\beta_k^*$ . The second inequality in (11) gives the estimation rate of the linear parameters  $\beta_k^*$  with  $k \in S^*$ . The last bound in (11) basically means that when an estimation  $\hat{\theta}_\ell$  with  $\ell \in \{1, \dots, K\}$  is far from any parameter  $\theta_k^*$  with  $k \in S^*$ , that is at a distance greater than  $r$ , the associated parameters  $\hat{\beta}_\ell$  drop to zero if the tuning parameter  $\kappa$  is taken equal to its lower bound and the decay rate of the noise variance  $\Delta_T$  drops to zero. Therefore, the contribution of the parameters  $\hat{\theta}_\ell$  in the far region, that are not in  $\tilde{S}(r)$ , will drop to zero as well.

*Remark 2.6* (Again on the dimension  $K$ ). As in Theorem 2.1, we remark that neither the bounds nor the probability of the event on which the bounds hold depend on the upper bound  $K$  on the sparsity  $s$ .

If the optimization on  $\vartheta$  in (4) is performed over a subset of  $\Theta_T$  in which the coordinates of the considered vectors are at a distance greater than  $2r$  pairwise with respect to the Riemannian metric  $\mathfrak{d}_T$ , then the sets  $\tilde{S}_k(r)$  contain at most one element. However, by doing so, we introduce an upper bound on the dimension  $K$  whereas in Theorem 2.1 the dimension  $K$  can be arbitrarily large. Indeed,  $\Theta_T$  is a compact set and therefore contains a finite number of balls of size  $2r$ .

**Outline of the paper.** In Section 3, we give the definition of the kernel  $\mathcal{K}_T$  measuring the correlation between two elements in the continuous dictionary and we present the regularity assumptions on the function  $\varphi_T$ . Section 4 introduces the Riemannian geometry framework useful in our context. Section 5 defines the convergence (or closeness condition) of kernels  $\mathcal{K}_T$  towards a limit kernel  $\mathcal{K}_\infty$ . Then, we require properties on the limit kernel  $\mathcal{K}_\infty$  and propagate them to the kernels  $\mathcal{K}_T$  thanks to this convergence. In Section 6, we present the assumptions on the existence of the so-called certificate functions used to state Theorems 2.1 and 2.5. We give sufficient conditions for the existence of certificate functions in Section 7. The examples of Gaussian sparse spike deconvolution and of Scaled exponential family in our regression model is fully detailed in Section 8 and 9, respectively. Then, the Appendix A is dedicated to the proofs of Theorems 2.1 and 2.5. The proofs of existence and explicit constructions of the certificates are detailed in the Appendix B. Other intermediate results are proven in Appendix C.

### 3. Dictionary of features

We present in the next section the regularity assumptions on the features  $(\varphi_T(\theta), \theta \in \Theta)$  we shall consider and then give examples of families of features satisfying such assumptions.

#### 3.1. Assumptions on the regularity of the features

Let  $T \in \mathbb{N}$  be fixed. We consider the Hilbert space  $(H_T, \langle \cdot, \cdot \rangle_T)$  and the features  $(\varphi_T(\theta), \theta \in \Theta)$  which are elements of  $H_T$ . We shall consider the following regularity assumptions on the features.

**Assumption 3.1** (Smoothness of  $\varphi_T$ ). *We assume that the function  $\varphi_T : \Theta \rightarrow H_T$  is of class  $\mathcal{C}^3$  and  $\|\varphi_T(\theta)\|_T > 0$  on  $\Theta$ .*

Recall  $\phi_T = \varphi_T / \|\varphi_T\|_T$  from (2) and notice that  $\phi_T$ , and thus  $\Phi_T$ , are continuous functions. Under Assumption 3.1, elementary calculations using (123) give:

$$(13) \quad \partial_\theta \phi_T(\theta) = \frac{\partial_\theta \varphi_T(\theta)}{\|\varphi_T(\theta)\|_T} - \frac{\varphi_T(\theta) \langle \varphi_T(\theta), \partial_\theta \varphi_T(\theta) \rangle_T}{\|\varphi_T(\theta)\|_T^3},$$

and thus, we deduce that the function  $g_T : \Theta \mapsto \mathbb{R}_+$  defined by:

$$(14) \quad g_T(\theta) = \|\partial_\theta \phi_T(\theta)\|_T^2$$

is well defined and continuous.

We shall consider the following non-degeneracy assumption on the features.

**Assumption 3.2** (Positivity of  $g_T$ ). *Assumption 3.1 holds and we have  $g_T > 0$  on  $\Theta$ .*

Even if Assumption 3.2 requires Assumption 3.1, in the following we shall stress when Assumption 3.1 is in force.

The next lemma gives a sufficient condition on  $\varphi_T$  for Assumption 3.2 to hold.

**Lemma 3.1** (On the positivity of  $g_T$ ). *Suppose Assumption 3.1 holds. If the elements  $\varphi_T(\theta)$  and  $\partial_\theta \varphi_T(\theta)$  of  $H_T$  are linearly independent for all  $\theta \in \Theta$  and  $\|\partial_\theta \varphi_T(\theta)\|_T > 0$  for all  $\theta \in \Theta$ , then  $g_T$  is positive on  $\Theta$ .*

**Proof.** For simplicity, we remove the subscript  $T$ , and for example write simply  $\phi = \varphi / \|\varphi\|$ . Recall that by Assumption 3.1 we have  $\|\varphi(\theta)\| > 0$ . Assume there exists  $\theta \in \Theta$  such that  $g(\theta) = 0$ , that is  $\partial_\theta \phi(\theta) = 0$ . Since  $\|\varphi(\theta)\| > 0$ , we deduce from (13) that  $\partial_\theta \varphi(\theta) \|\varphi(\theta)\|^2 - \varphi(\theta) \langle \varphi(\theta), \partial_\theta \varphi(\theta) \rangle = 0$ . Then use that by assumption  $\partial_\theta \varphi(\theta) \neq 0$  and  $\|\varphi(\theta)\| > 0$ , to get that  $\varphi(\theta)$  and  $\partial_\theta \varphi(\theta)$  are linearly dependent. In conclusion, we get that if  $\varphi(\theta)$  and  $\partial_\theta \varphi(\theta)$  are linearly independent, then  $g(\theta) > 0$ .  $\square$

#### 3.2. Examples of regular features

The aim of this section of examples is to stress that a large variety of dictionaries of features and type of parameters verify Assumptions 3.1 and 3.2.

##### 3.2.1. Translation discrete-time model

Let  $t_1 < \dots < t_T$  be a grid on  $\mathbb{R}$  of size  $T \in \mathbb{N}$ ,  $\lambda_T$  an atomic measure whose support is the grid, and  $H_T = L^2(\lambda_T)$ . Consider the translation invariant dictionary:

$$(15) \quad (\varphi_T(\theta) = k(\cdot - \theta), \theta \in \Theta),$$

with  $\Theta = \mathbb{R}$  and  $k$  is a real-valued  $\mathcal{C}^3$  function defined on  $\mathbb{R}$ . Notice the dictionary does not depend on  $T$ . We now consider usual choices for the function  $k$ .

For the Gaussian function  $k(t) = e^{-t^2/2}$  and the Cauchy function  $k(t) = 1/(1+t^2)$ , we get that Assumption 3.1 holds and, using Lemma 3.1 that Assumption 3.2 is also satisfied provided respectively  $T \geq 2$  and  $T \geq 3$ .

For the Shannon scaling function  $k(t) = \text{sinc}(t) = \sin(\pi t)/(\pi t)$ , Assumption 3.1 holds provided that  $\lambda_T((a + \mathbb{Z})^c) > 0$  for all  $a \in \mathbb{R}$ , that is the grid is not a subset of  $a + \mathbb{Z}^*$  for some  $a \in \mathbb{R}$ . There is no easy way to write conditions on the grid, based on the use of Lemma 3.1, for Assumption 3.2 to hold (let us mention that  $T \geq 2$  and  $\min_{1 \leq i \leq T-1} (t_{i+1} - t_i) < 1/2$  is a sufficient condition for Assumption 3.2 to hold).

Eventually notice that the Laplace function  $k(t) = e^{-|t|}$  is not smooth enough for Assumption 3.1 to hold.

### 3.2.2. Translation model with a continuum of observations

Let  $T \in \mathbb{N}$  (which does not play a role here) and  $H_T = L^2(\text{Leb})$ , where  $\text{Leb}$  is the Lebesgue measure on  $\mathbb{R}$ . In this framework, the observation  $y$  defined in (1) is a continuum of observations. Consider the translation invariant dictionaries from Section 3.2.1, where  $k$  is either the Gaussian, the Cauchy or the Shannon scaling function. Notice that the Hilbert space and the dictionary do not depend on  $T$ . Then, it is easy to check that Assumptions 3.1 and 3.2 hold.

We see that this model, which can be seen as a continuous approximation (or limit) of the discrete models from Section 3.2.1 when  $T$  therein is large, is easier to handle than the corresponding discrete models.

### 3.2.3. Translation model with a varying scaling parameter

Let  $T \in \mathbb{N}$ ,  $H_T = L^2(\text{Leb})$ , where  $\text{Leb}$  is the Lebesgue measure on  $\mathbb{R}$ , and consider the translation invariant dictionary scaled by  $\bar{\sigma}_T > 0$  given by:

$$(\varphi_T(\theta) = k(\bar{\sigma}_T^{-1}(\cdot - \theta)), \theta \in \Theta),$$

with  $\Theta = \mathbb{R}$  and  $k$  is a real-valued  $\mathcal{C}^3$  function defined on  $\mathbb{R}$ . Contrary to Section 3.2.2, the features depend on  $T$ . Suppose that  $k$  is the Shannon scaling function (see Section 3.2.1) and consider the vector sub-space  $V_T$  given by the closure in  $H_T$  of the vector space spanned by the dictionary. According to [39, Theorem 3.5], the set  $V_T$  is the subset of  $H_T$  of all functions whose Fourier transform support is a subset of  $[-\pi/\bar{\sigma}_T, \pi/\bar{\sigma}_T]$ . Suppose that the sequence  $(\bar{\sigma}_T, T \in \mathbb{N})$  is decreasing to 0. Then the sequence  $(V_T, T \in \mathbb{N})$  is increasing and  $\bigcup_{T \in \mathbb{N}} V_T = H_T$ . This model provides an example of translation models with possibly varying, but known, scaling parameter  $\bar{\sigma}_T$ .

### 3.2.4. Scaling exponential model

Let  $T \in \mathbb{N}$  (which does not play a role here),  $H_T = L^2(\text{Leb})$ , where  $\text{Leb}$  is the Lebesgue measure on  $\mathbb{R}_+$ , and consider the scale invariant dictionary given by:

$$(\varphi_T(\theta) = k(\theta \cdot), \theta \in \Theta),$$

with  $\Theta = \mathbb{R}_+^*$  and the exponential function  $k : t \mapsto e^{-t}$ . This dictionary is used for example in fluorescence microscopy (see [24]). Clearly Assumption 3.1 holds as well as Assumption 3.2 as  $g_T(\theta) = 1/(4\theta^2)$ .

## 4. A Riemannian metric on the set of parameters

### 4.1. On the Riemannian metric in dimension one

Recall  $\Theta$  is an interval of  $\mathbb{R}$ . We call kernel a real-valued function defined on  $\Theta^2$ . Let  $\mathcal{K}$  be a symmetric kernel of class  $\mathcal{C}^2$  such that the function  $g_{\mathcal{K}}$  defined on  $\Theta$  by:

$$(16) \quad g_{\mathcal{K}}(\theta) = \partial_{x,y}^2 \mathcal{K}(\theta, \theta)$$

is positive and locally bounded, where  $\partial_x$  (resp.  $\partial_y$ ) denotes the usual derivative with respect to the first (resp. second) variable. Following [41], we define an intrinsic Riemannian metric, denoted  $\mathfrak{d}_{\mathcal{K}}$ , on the parameter set  $\Theta$  using the function  $g_{\mathcal{K}}$ . One of the motivations to use the Riemannian metric is to work with intrinsic quantities related to the parameters which are invariant by reparametrization, such as the diameter of (subsets of)  $\Theta$ . Since  $\Theta$  is one-dimensional and connected, the Riemannian metric  $\mathfrak{d}_{\mathcal{K}}(\theta, \theta')$  between  $\theta, \theta' \in \Theta$  reduces to:

$$(17) \quad \mathfrak{d}_{\mathcal{K}}(\theta, \theta') = |G_{\mathcal{K}}(\theta) - G_{\mathcal{K}}(\theta')|,$$

where  $G_{\mathcal{K}}$  is a primitive of  $\sqrt{g_{\mathcal{K}}}$ .

*Remark 4.1.* We refer to [37] and [43] for a general presentation on Riemannian manifolds, and we give an immediate application in dimension one which entails in particular (17). Let  $\Theta$  be a manifold (of dimension one). A path  $\gamma : [0, 1] \rightarrow \Theta$  is an admissible path if it is continuous, piecewise continuously differentiable with non-vanishing derivative. Its length is given by  $\mathcal{L}_{\mathcal{K}}(\gamma) = \int_0^1 |\dot{\gamma}_s| \sqrt{g_{\mathcal{K}}(\gamma_s)} ds$ , where  $|\dot{\gamma}_s|$  is seen as the norm of the vector  $\dot{\gamma}_s$  in the tangent space, and the scalar product on the tangent space at  $\theta \in \Theta$  is given by  $(u, v) \mapsto \langle u, g_{\mathcal{K}}(\theta)v \rangle$  with  $\langle \cdot, \cdot \rangle$  the usual Euclidean scalar product. (In our case, the tangent vector space is  $\mathbb{R}$  and the Euclidean scalar product reduces to the usual product). The Riemannian metric  $\mathfrak{d}_{\mathcal{K}}$  between  $\theta, \theta'$  in  $\Theta$  is then defined by:

$$(18) \quad \mathfrak{d}_{\mathcal{K}}(\theta, \theta') = \inf_{\gamma} \mathcal{L}_{\mathcal{K}}(\gamma),$$

where the infimum is taken over the admissible paths  $\gamma$  such that  $\gamma_0 = \theta$  and  $\gamma_1 = \theta'$ . It is not hard to see that  $\gamma$  is a minimizing path, that is,  $\mathfrak{d}_{\mathcal{K}}(\theta, \theta') = \mathcal{L}_{\mathcal{K}}(\gamma)$ , if and only if  $\gamma$  is monotone (and thus  $\gamma_s \in [\theta \wedge \theta', \theta \vee \theta']$  for all  $s \in [0, 1]$ ). This is equivalent to say that the sign of  $\dot{\gamma}_s$  is constant. Assume that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^1$ . The path  $\gamma$  is a geodesic if it is smooth with zero acceleration, that is, in dimension one for all  $s \in (0, 1)$ :

$$(19) \quad \ddot{\gamma}_s + \frac{1}{2} \frac{g'_{\mathcal{K}}(\gamma_s)}{g_{\mathcal{K}}(\gamma_s)} \dot{\gamma}_s^2 = 0.$$

This is equivalent to  $s \mapsto \dot{\gamma}_s \sqrt{g_{\mathcal{K}}(\gamma_s)}$  being constant, which implies that the geodesic is a minimizing path.

We now derive the equation of the geodesic path when  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^1$ . Recall  $G_{\mathcal{K}}$  denotes the primitive of  $\sqrt{g_{\mathcal{K}}}$ . It is continuous increasing and thus induces a one-to-one map from  $\Theta$  to its image. Set  $a = G_{\mathcal{K}}(\theta)$  and  $b = G_{\mathcal{K}}(\theta') - G_{\mathcal{K}}(\theta)$ , so that the path  $\gamma : [0, 1] \rightarrow \Theta$  defined by  $\gamma_s = G_{\mathcal{K}}^{-1}(a + bs)$  is a geodesic and minimizing path from  $\theta$  to  $\theta'$  with  $\mathcal{L}_{\mathcal{K}}(\gamma) = \mathfrak{d}_{\mathcal{K}}(\theta, \theta')$ .

Following [41], we introduce the covariant derivatives, see [2, Sections 3.6 and 5.6], which have elementary expressions as the set of parameters  $\Theta$  is one-dimensional. For a smooth function  $f$  defined on  $\Theta$  and taking values in an Hilbert space, say  $H$ , the covariant derivative  $D_{i;\mathcal{K}}[f]$  of order  $i \in \mathbb{N}$  is defined recursively by  $D_{0;\mathcal{K}}[f] = f$  and for  $i \in \mathbb{N}$ , assuming that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^i$ , and  $\theta \in \Theta$ :

$$(20) \quad D_{i+1;\mathcal{K}}[f](\theta) = g_{\mathcal{K}}(\theta)^{\frac{i}{2}} \partial_{\theta} \left( \frac{D_{i;\mathcal{K}}[f](\theta)}{g_{\mathcal{K}}(\theta)^{\frac{i}{2}}} \right).$$

In particular, we have for  $f \in \mathcal{C}^2(\Theta, H)$  (and assuming that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^1$  for the last equality) that:

$$(21) \quad D_{0;\mathcal{K}}[f] = f, \quad D_{1;\mathcal{K}}[f] = \partial_{\theta} f, \quad D_{2;\mathcal{K}}[f] = \partial_{\theta}^2 f - \frac{1}{2} \frac{g'_{\mathcal{K}}}{g_{\mathcal{K}}} \partial_{\theta} f.$$

We shall also consider the following modification of the covariant derivative, for  $i \in \mathbb{N}$ :

$$(22) \quad \tilde{D}_{i;\mathcal{K}}[f](\theta) = g_{\mathcal{K}}(\theta)^{-i/2} D_{i;\mathcal{K}}[f](\theta).$$

We have  $\tilde{D}_{0;\mathcal{K}}[f] = f$ , and we deduce from (20) that for  $i \in \mathbb{N}^*$ , assuming that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^i$ :

$$(23) \quad \tilde{D}_{i;\mathcal{K}} = \tilde{D}_{1;\mathcal{K}} \circ \tilde{D}_{i-1;\mathcal{K}} = \left( \tilde{D}_{1;\mathcal{K}} \right)^i,$$

so that  $\tilde{D}_{1;\mathcal{K}}$  can be seen as a derivative operator.

We now give an elementary variant of the Taylor-Lagrange expansion using the previously defined Riemannian metric and covariant derivatives. Its proof can be found in the Appendix, Section C.4.

**Lemma 4.2.** *Assume  $g_{\mathcal{K}}$  is positive and of class  $\mathcal{C}^1$ . Let  $f$  be a function defined on  $\Theta$  taking values in an Hilbert space of class  $\mathcal{C}^2$ . Setting  $f^{[i]} = \tilde{D}_{i;\mathcal{K}}[f]$  for  $i \in \{1, 2\}$ , we have that for all  $\theta, \theta_0 \in \Theta$ :*

$$(24) \quad f(\theta) = f(\theta_0) + \text{sign}(\theta - \theta_0) \mathfrak{d}_{\mathcal{K}}(\theta, \theta_0) f^{[1]}(\theta_0) + \mathfrak{d}_{\mathcal{K}}(\theta, \theta_0)^2 \int_0^1 (1-t) f^{[2]}(\gamma_t) dt,$$

where  $\gamma$  is a geodesic path such that  $\gamma_0 = \theta_0$ ,  $\gamma_1 = \theta$  (and thus  $\mathfrak{d}_{\mathcal{K}}(\theta, \theta_0) = \mathcal{L}_{\mathcal{K}}(\gamma)$ ).

For a real-valued function  $F$  defined on  $\Theta^2$ , we say that  $F$  is of class  $\mathcal{C}^{0,0}$  on  $\Theta^2$  if it is continuous on  $\Theta^2$ , and of class  $\mathcal{C}^{i,j}$  on  $\Theta^2$ , with  $i, j \in \mathbb{N}$ , as soon as:  $F$  is of class  $\mathcal{C}^{0,0}$ , and if  $i \geq 1$  then the function  $\theta \mapsto F(\theta, \theta')$  is of class  $\mathcal{C}^i$  on  $\Theta$  and its derivative  $\partial_x F$  is of class  $\mathcal{C}^{i-1,j}$  on  $\Theta^2$ , and if  $j \geq 1$  the function  $\theta' \mapsto F(\theta, \theta')$  is of class  $\mathcal{C}^j$  on  $\Theta$  and its derivative  $\partial_y F$  is of class  $\mathcal{C}^{i,j-1}$  on  $\Theta^2$ . For a real-valued symmetric function  $F$  defined on  $\Theta^2$  of class  $\mathcal{C}^{i,j}$ , we define the covariant derivatives  $D_{i,j;\mathcal{K}}[F]$  of order  $(i, j) \in \mathbb{N}^2$  recursively by  $D_{0,0;\mathcal{K}}[F] = F$  and for  $i, j \in \mathbb{N}$ , assuming that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^{\max(i,j)}$ , and  $\theta, \theta' \in \Theta$ :

$$(25) \quad D_{i+1,j;\mathcal{K}}[F](\theta, \theta') = g_{\mathcal{K}}(\theta)^{\frac{i}{2}} \partial_{\theta} \left( \frac{D_{i,j;\mathcal{K}}[F](\theta, \theta')}{g_{\mathcal{K}}(\theta)^{\frac{i}{2}}} \right) \quad \text{and} \quad D_{i,j;\mathcal{K}}[F](\theta, \theta') = D_{j,i;\mathcal{K}}[F](\theta', \theta).$$

In particular, we have  $D_{0,0;\mathcal{K}}[F] = F$ ,  $D_{1,0;\mathcal{K}} = \partial_x F$ ,  $D_{0,1;\mathcal{K}} = \partial_y F$  and  $D_{1,1;\mathcal{K}} = \partial_{xy}^2 F$ . We shall also consider the following modification of the covariant derivative, for  $i, j \in \mathbb{N}$ :

$$(26) \quad \tilde{D}_{i,j;\mathcal{K}}[F](\theta, \theta') = \frac{D_{i,j;\mathcal{K}}[F](\theta, \theta')}{g_{\mathcal{K}}(\theta)^{i/2} g_{\mathcal{K}}(\theta')^{j/2}}.$$

We have  $\tilde{D}_{1,0;\mathcal{K}} \circ \tilde{D}_{0,1;\mathcal{K}} = \tilde{D}_{0,1;\mathcal{K}} \circ \tilde{D}_{1,0;\mathcal{K}}$  and for  $i, j \in \mathbb{N}$ , assuming that  $g_{\mathcal{K}}$  is of class  $\mathcal{C}^{\max(i,j)}$ :

$$\tilde{D}_{i,j;\mathcal{K}} = \left( \tilde{D}_{1,0;\mathcal{K}} \right)^i \circ \left( \tilde{D}_{0,1;\mathcal{K}} \right)^j.$$

For  $i, j \in \mathbb{N}$ , if  $\mathcal{K}$  is of class  $\mathcal{C}^{i \vee 1, j \vee 1}$ , we consider the real-valued function defined on  $\Theta^2$  by:

$$(27) \quad \mathcal{K}^{[i,j]} = \tilde{D}_{i,j;\mathcal{K}}[\mathcal{K}].$$

In particular, since  $\mathcal{K}$  is of class  $\mathcal{C}^2$ , we have:

$$(28) \quad \mathcal{K}^{[0,0]} = \mathcal{K} \quad \text{and} \quad \mathcal{K}^{[1,1]}(\theta, \theta) = 1.$$

#### 4.2. The kernel and the Riemannian metric associated to the dictionary of features

Let  $T \in \mathbb{N}$  be fixed and assume that Assumption 3.2 holds. We define the kernel  $\mathcal{K}_T$  on  $\Theta^2$  by:

$$(29) \quad \mathcal{K}_T(\theta, \theta') = \langle \phi_T(\theta), \phi_T(\theta') \rangle_T = \frac{\langle \varphi_T(\theta), \varphi_T(\theta') \rangle_T}{\|\varphi_T(\theta)\|_T \|\varphi_T(\theta')\|_T},$$

where we recall that  $\phi_T = \varphi_T / \|\varphi_T\|_T$ . When considering the kernel  $\mathcal{K}_T$ , we shall write  $g_T$  for  $g_{\mathcal{K}_T}$ , and similarly we shall use the notations  $\tilde{D}_{i,T}$  and  $\tilde{D}_{i,j;T}$  instead of  $\tilde{D}_{i;\mathcal{K}_T}$  and  $\tilde{D}_{i,j;\mathcal{K}_T}$ . Recall the derivatives of the kernel  $\mathcal{K}_T$  defined by (27). The next lemma insures in particular that the two definitions of  $g_T$  given by (14) and (16) are consistent, that is:

$$(30) \quad g_T(\theta) = \partial_{xy}^2 \mathcal{K}_T(\theta, \theta) = \|\partial_\theta \phi_T(\theta)\|_T^2.$$

The proof of the next lemma can be found in the Appendix, Section C.4.

**Lemma 4.3.** *Let  $T \in \mathbb{N}$  be fixed and assume that Assumptions 3.1 and 3.2 hold. Then, the symmetric kernel  $\mathcal{K}_T$  is of class  $\mathcal{C}^{3,3}$  on  $\Theta^2$  and for  $i, j \in \{0, \dots, 3\}$  and  $\theta, \theta' \in \Theta$ , we have:*

$$(31) \quad \mathcal{K}_T^{[i,j]}(\theta, \theta') = \langle \tilde{D}_{i,T}[\phi_T](\theta), \tilde{D}_{j,T}[\phi_T](\theta') \rangle_T.$$

We also have:

$$(32) \quad \sup_{\Theta^2} |\mathcal{K}_T^{[0,0]}| \leq 1, \quad \mathcal{K}_T^{[0,0]}(\theta, \theta) = 1, \quad \mathcal{K}_T^{[1,0]}(\theta, \theta) = 0, \quad \mathcal{K}_T^{[2,0]}(\theta, \theta) = -1 \quad \text{and} \quad \mathcal{K}_T^{[2,1]}(\theta, \theta) = 0.$$

### 5. Approximating the kernel associated to the dictionary

In the section we detail the assumptions guaranteeing the approximation of the kernel  $\mathcal{K}_T$  (which is usually difficult to compute) by a kernel  $\mathcal{K}_\infty$  (which is easier to handle). Both kernels are defined on  $\Theta^2$ , however, we shall qualify the approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_\infty$  and properties of  $\mathcal{K}_\infty$  on subsets of  $\Theta$ , respectively  $\Theta_T$  (which will be a compact interval) and  $\Theta_\infty$  (which will be an interval possibly unbounded). We use notations from Section 4 and recall the definition of  $g_{\mathcal{K}}$ , resp.  $\mathcal{K}^{[i,j]}$ , given in (16), resp. in (27). Assuming the kernel  $\mathcal{K}$  is of class  $\mathcal{C}^{3,3}$  and using the notation (27), we also set for  $\theta \in \Theta$ :

$$(33) \quad h_{\mathcal{K}}(\theta) = \mathcal{K}^{[3,3]}(\theta, \theta).$$

For simplicity, for an expression  $A$  we write  $A_*$  for  $A_{\mathcal{K}_*}$  where  $*$  is equal to  $T$  or  $\infty$ . We first give a regularity assumption on the kernel  $\mathcal{K}_\infty$ .

**Assumption 5.1** (Properties of the asymptotic kernel  $\mathcal{K}_\infty$  and function  $h_\infty$ ). *The symmetric kernel  $\mathcal{K}_\infty$  defined on  $\Theta^2$  is of class  $\mathcal{C}^{3,3}$ , the function  $g_\infty$  defined by (16) on  $\Theta$  is positive and locally bounded (as well as of class  $\mathcal{C}^2$ ), and we have  $\mathcal{K}_\infty(\theta, \theta) = -\mathcal{K}_\infty^{[2,0]}(\theta, \theta) = 1$  for  $\theta \in \Theta$ . The set  $\Theta_\infty \subseteq \Theta$  is an interval and we have:*

$$(34) \quad L_3 := \sup_{\Theta_\infty} h_\infty < +\infty, \quad \text{and} \quad L_{i,j} := \sup_{\Theta_\infty^2} |\mathcal{K}_\infty^{[i,j]}| < +\infty \quad \text{for all } i, j \in \{0, 1, 2\}.$$

Since  $\Theta_T$  is compact, under Assumptions 3.2 and 5.1, we deduce that the constant  $\rho_T$  below is positive and finite, where:

$$(35) \quad \rho_T = \max \left( \sup_{\Theta_T} \sqrt{\frac{g_T}{g_\infty}}, \sup_{\Theta_T} \sqrt{\frac{g_\infty}{g_T}} \right).$$

From the definition of the Riemannian metric given in (17) (see also (18)), we readily deduce that the metrics  $\mathfrak{d}_T$  and  $\mathfrak{d}_\infty$  are then strongly equivalent on  $\Theta_T$ ; more precisely we have that on  $\Theta_T^2$ :

$$(36) \quad \frac{1}{\rho_T} \mathfrak{d}_\infty \leq \mathfrak{d}_T \leq \rho_T \mathfrak{d}_\infty.$$

We then give an assumption on the quality of approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_\infty$ . We set:

$$(37) \quad \mathcal{V}_T = \max(\mathcal{V}_T^{(1)}, \mathcal{V}_T^{(2)}) \quad \text{with} \quad \mathcal{V}_T^{(1)} = \max_{i,j \in \{0,1,2\}} \sup_{\Theta_T^2} |\mathcal{K}_T^{[i,j]} - \mathcal{K}_\infty^{[i,j]}| \quad \text{and} \quad \mathcal{V}_T^{(2)} = \sup_{\Theta_T} |h_T - h_\infty|.$$

Let us recall that Assumption 3.2 implies regularity conditions on  $\mathcal{K}_T$ , see Lemma 4.3.

**Assumption 5.2** (Quality of the approximation). *Let  $T \in \mathbb{N}$  be fixed. Assumptions 3.2 and 5.1 hold, the interval  $\Theta_T \subset \Theta_\infty$  is a compact interval, and we have:*

$$\mathcal{V}_T \leq L_{2,2} \wedge L_3.$$

Notice that if Assumption 3.2 holds, then Assumptions 5.1 and 5.2 hold trivially when one takes  $\mathcal{K}_\infty = \mathcal{K}_T$  and  $\Theta_\infty = \Theta_T$ ; notice also that  $\rho_T = 1$  in this case. In the next example, the sequence of kernels  $(\mathcal{K}_T, T \in \mathbb{N})$  converges to the kernel  $\mathcal{K}_\infty$  as  $T$  goes to infinity, so that Assumption 5.2 holds for  $T$  large enough.

*Example 5.1.* We consider the discrete-time example from Section 1.2.1. We assume that the process  $y$  is a function defined on  $[0, 1]$  which, for  $T \in \mathbb{N}^*$  is observed through the regular grid  $\{t_{j,T} = j/T : 1 \leq j \leq T\}$ . The process  $y$  is seen as an element of the Hilbert space  $H_T = L^2(\lambda_T)$ , with the probability measure  $\lambda_T = \Delta_T \sum_{j=1}^T \delta_{t_{j,T}}$  on  $[0, 1]$  with  $\Delta_T = 1/T$ . Let  $\Theta$  be a compact interval of  $\mathbb{R}$  and set  $\Theta_T = \Theta_\infty = \Theta$ . Consider a dictionary  $(\varphi(\theta), \theta \in \Theta)$  independent of  $T$ , that is,  $\varphi_T = \varphi$  for all  $T \in \mathbb{N}^*$ , and assume that the function  $(\theta, t) \mapsto \varphi(\theta)(t)$  is defined on  $\Theta \times [0, 1]$  and of class  $\mathcal{C}^{3,0}$ . Assume that the dictionary satisfies the regularity assumptions of Assumption 3.2.

Let  $\text{Leb}$  be the Lebesgue measure on  $[0, 1]$ , so that  $(\lambda_T, T \in \mathbb{N}^*)$  converges weakly to  $\text{Leb}$ . Then, define the kernel  $\mathcal{K}_\infty$  by (29) with  $\varphi_T$  replaced by  $\varphi$  (as the dictionary does not depend on  $T$ ) and the scalar product  $\langle \cdot, \cdot \rangle_T$  by the usual scalar product on  $L^2(\text{Leb})$ . Thanks to Lemma 4.3, we deduce that Assumption 5.1 on the properties of  $\mathcal{K}_\infty$  is satisfied. Using the weak convergence of  $(\lambda_T, T \in \mathbb{N}^*)$  to  $\text{Leb}$ , we deduce that  $\lim_{T \rightarrow \infty} \partial_x^i \partial_y^j \mathcal{K}_T = \partial_x^i \partial_y^j \mathcal{K}_\infty$  uniformly on  $[0, 1]^2$  for all  $i, j \in \{0, \dots, 3\}$ . This implies that:

$$\lim_{T \rightarrow \infty} \mathcal{V}_T = 0 \quad \text{and} \quad \lim_{T \rightarrow \infty} \rho_T = 1.$$

Thus Assumption 5.2 holds for  $T$  large enough.

## 6. Certificates

In this section, we make assumptions on the existence of functions from  $\Theta$  to  $\mathbb{R}$  called certificates. These functions have interpolation properties that are corner stones in the proof of Theorem 2.1. The term “certificate” is inherited from the compressed sensing field where such functions were used to get rid of the Restricted Isometry Property condition (RIP) for exact reconstruction of signals (see [20] for details on the RIP condition). In [19], the authors showed that it is possible to reconstruct exactly a sparse signal from the observations of a finite number of Fourier coefficients by exhibiting a dual certificate. Many papers have followed this line of research since then (see e.g. [17, 18, 26]).



In sparse linear models the bounds for prediction error are proved using RIP, Restricted Eigenvalue or compatibility conditions (see [9, 51]). Among these assumptions, the compatibility conditions are the less restrictive. Indeed, the authors of [52] have shown that it is implied by both the RIP and the Restricted Eigenvalue. However, in many contexts even the weaker condition fails to hold. Typically the compatibility condition fails to hold in the context of super-resolution which aims at extracting the frequencies and amplitudes of a linear combination of complex exponentials from a small number of noisy time samples (see [12]).

In the papers [12] and [47], the authors achieve nearly optimal rates for the prediction error in the super-resolution framework using certificate functions. Their method and proof are however quite specific to complex exponentials and their certificates are trigonometric polynomials. The insightful paper of [26] builds certificates in a quite general setting for a one dimensional parameter set  $\Theta$ . In [22], the authors exhibit certificate functions to deal with more general probability density models where  $\Theta$  is multidimensional. However they are restricted to translation invariant dictionaries (15). The most general framework has been introduced in [41] where the Riemannian geometry is key to build in a natural way the so-called certificate functions. In fact a separation distance between the parameters to estimate is needed to build certificates and the Euclidean metric yields overly pessimistic minimum separation condition. In what follows we introduce new certificates, called derivative certificates, in order to control the prediction error.

We consider the following assumption in the spirit of [41]. We consider the setting where  $T$  may be finite. Let  $T \in \mathbb{N}$ ,  $H_T$  be an Hilbert space and  $(\varphi_T(\theta), \theta \in \Theta)$  a dictionary satisfying Assumptions 3.1 and 3.2, so that the kernel  $\mathcal{K}_T$  is of class  $\mathcal{C}^{3,3}$  on  $\Theta^2$ . Recall the Riemannian metric  $\mathfrak{d}_{\mathcal{K}_T}$  associated to  $\mathcal{K}_T$ , which we simply denote by  $\mathfrak{d}_T$ . We define the closed ball centered at  $\theta \in \Theta_T$  with radius  $r$  by:

$$\mathcal{B}_T(\theta, r) = \{\theta' \in \Theta_T, \mathfrak{d}_T(\theta, \theta') \leq r\} \subseteq \Theta_T.$$

Let  $\mathcal{Q}^*$  be a subset of  $\Theta_T$  of cardinal  $s$ . For  $r > 0$ , the near region of  $\mathcal{Q}^*$  is the union of balls  $\bigcup_{\theta^* \in \mathcal{Q}^*} \mathcal{B}_T(\theta^*, r)$  and its far region is the complementary of the near region in  $\Theta_T$ :  $\Theta_T \setminus \bigcup_{\theta^* \in \mathcal{Q}^*} \mathcal{B}_T(\theta^*, r)$ . Sufficient conditions for the next assumption to hold are given in Section 7.

**Assumption 6.1** (Interpolating certificate). *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}^*$ ,  $r > 0$  and  $\mathcal{Q}^*$  be a subset of  $\Theta_T$  of cardinal  $s$ . Suppose Assumptions 3.1 and 3.2 on the dictionary  $(\varphi_T(\theta), \theta \in \Theta)$ , and Assumption 5.1 on the kernel  $\mathcal{K}_\infty$ , defined on  $\Theta^2$ , hold. Suppose that  $\mathfrak{d}_T(\theta, \theta') > 2r$  for all  $\theta, \theta' \in \mathcal{Q}^* \subset \Theta_T$ , and that there exist finite positive constants  $C_N, C'_N, C_F, C_B$ , with  $C_F < 1$ , depending on  $r$  and  $\mathcal{K}_\infty$  such that for any application  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  there exists an element  $p \in H_T$  satisfying:*

- (i) *For all  $\theta^* \in \mathcal{Q}^*$  and  $\theta \in \mathcal{B}_T(\theta^*, r)$ , we have  $|\langle \phi_T(\theta), p \rangle_T| \leq 1 - C_N \mathfrak{d}_T(\theta^*, \theta)^2$ .*
- (ii) *For all  $\theta^* \in \mathcal{Q}^*$  and  $\theta \in \mathcal{B}_T(\theta^*, r)$ , we have  $|\langle \phi_T(\theta), p \rangle_T - v(\theta^*)| \leq C'_N \mathfrak{d}_T(\theta^*, \theta)^2$ .*
- (iii) *For all  $\theta$  in  $\Theta_T$  and  $\theta \notin \bigcup_{\theta^* \in \mathcal{Q}^*} \mathcal{B}_T(\theta^*, r)$  (far region), we have  $|\langle \phi_T(\theta), p \rangle_T| \leq 1 - C_F$ .*
- (iv) *We have  $\|p\|_T \leq C_B \sqrt{s}$ .*

The function  $\eta : \theta \mapsto \langle \phi_T(\theta), p \rangle_T$  is the so-called “interpolating certificate” of the function  $v$ , as thanks to (ii) with  $\theta = \theta^*$ , the function  $\eta$  coincides with the function  $v$  on  $\mathcal{Q}^*$ . In addition, the interpolating certificate is required to have curvature properties in the near region and to be bounded by a constant strictly inferior to 1 in the far region. When  $r$  is sufficiently small (that is,  $r \leq \sqrt{2/(C_N + C'_N)}$ ) Conditions (i) and (ii) are equivalent to the fact that the function  $\eta$  is in-between two quadratic functions in the near region of  $\mathcal{Q}^*$ : for all  $\theta^* \in \mathcal{Q}^*$  such that  $v(\theta^*) = 1$  (resp.  $v(\theta^*) = -1$ ) and  $\theta \in \mathcal{B}_T(\theta^*, r)$ , we have  $1 - C'_N \mathfrak{d}_T(\theta^*, \theta)^2 \leq \eta(\theta) \leq 1 - C_N \mathfrak{d}_T(\theta^*, \theta)^2$  (resp.  $-1 + C_N \mathfrak{d}_T(\theta^*, \theta)^2 \leq \eta(\theta) \leq -1 + C'_N \mathfrak{d}_T(\theta^*, \theta)^2$ ).

Sufficient conditions for the next assumption to hold are also given in Section 7.

**Assumption 6.2** (Interpolating derivative certificate). *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}^*$ ,  $r > 0$  and  $\mathcal{Q}^*$  be a subset of  $\Theta_T$  of cardinal  $s$ . Suppose Assumptions 3.1 and 3.2 on the dictionary  $(\varphi_T(\theta), \theta \in \Theta)$ , and Assumption 5.1 on the kernel  $\mathcal{K}_\infty$ , defined on  $\Theta^2$ , hold. Assume that  $\mathfrak{d}_T(\theta, \theta') > 2r$  for all  $\theta, \theta' \in \mathcal{Q}^* \subset \Theta_T$  and that there exist finite positive constants  $c_N, c_F, c_B$  depending on  $r$  and  $\mathcal{K}_\infty$ , such that for any application  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  there exists an element  $q \in H_T$  satisfying:*

- (i) *For all  $\theta^* \in \mathcal{Q}^*$  and  $\theta \in \mathcal{B}_T(\theta^*, r)$ , we have:*

$$|\langle \phi_T(\theta), q \rangle_T - v(\theta^*) \operatorname{sign}(\theta - \theta^*) \mathfrak{d}_T(\theta, \theta^*)| \leq c_N \mathfrak{d}_T(\theta^*, \theta)^2.$$

- (ii) *For all  $\theta$  in  $\Theta_T$  and  $\theta \notin \bigcup_{\theta^* \in \mathcal{Q}^*} \mathcal{B}_T(\theta^*, r)$  (far region), we have  $|\langle \phi_T(\theta), q \rangle_T| \leq c_F$ .*

$$(iii) \quad \|q\|_T \leq c_B \sqrt{s}.$$

The function  $\theta \mapsto \langle \phi_T(\theta), q \rangle_T$  will be called an “interpolating derivative certificate” as it vanishes on  $Q^*$ . In addition, this function is required to decrease similarly to the function  $\mathfrak{d}_T(\cdot, \theta^*)$  near  $\theta^*$  and to be bounded in the far region of  $Q^*$ .

## 7. Sufficient conditions for the existence of certificates

In this section, we prove the existence of the certificate functions of Assumptions 6.1 and 6.2 provided that the parameters to be estimated are sufficiently separated in terms of the Riemannian metric. According to [45], the separation condition cannot be avoided to build certificate functions in general. It is however possible to remove this separation condition in some particular cases, see [44] for models with positive amplitudes.

In order to find sufficient conditions for the existence of the interpolating certificate functions of Assumption 6.1, we extend the construction from [41] to a non asymptotic setting. For the existence of the interpolating derivative certificate functions of Assumption 6.2, we generalize the proof of [17, Lemma 2.7] dedicated to the dictionary of complex exponential functions. The proofs for the existence of certificates given in Section B require boundedness and local concavity properties of the kernel  $\mathcal{K}_T$ . For practical application, they are deduced from the boundedness and local concavity properties of the kernel  $\mathcal{K}_\infty$  and the quality of approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_\infty$  discussed in Section 5.

### 7.1. Boundedness and local concavity of the kernel $\mathcal{K}_T$

In this work, we shall consider bounded kernels locally concave on the diagonal. More precisely, for  $T \in \bar{\mathbb{N}} = \mathbb{N} \cup \{\infty\}$  and  $r > 0$ , we define:

$$(38) \quad \varepsilon_T(r) = 1 - \sup \{ |\mathcal{K}_T(\theta, \theta')|; \quad \theta, \theta' \in \Theta_T \text{ such that } \mathfrak{d}_T(\theta', \theta) \geq r \},$$

$$(39) \quad \nu_T(r) = -\sup \left\{ \mathcal{K}_T^{[0,2]}(\theta, \theta'); \quad \theta, \theta' \in \Theta_T \text{ such that } \mathfrak{d}_T(\theta', \theta) \leq r \right\}.$$

The fact that  $\varepsilon_T(r)$  and  $\nu_T(r)$  are positive depends on the function  $\varphi_T$ , the space  $H_T$  and the set  $\Theta_T$ . Let us mention that in many examples the positiveness of  $\varepsilon_\infty(r)$  and  $\nu_\infty(r)$  is easy to check whereas the positiveness of  $\varepsilon_T(r)$  and  $\nu_T(r)$  might be more difficult to prove.

Notice that (32) for  $T \in \mathbb{N}$  and Assumption 5.1 for  $T = \infty$ , and the continuity of  $\mathcal{K}_T$  and  $\mathcal{K}_T^{[0,2]}$  give that:

$$(40) \quad \lim_{r \rightarrow 0^+} \varepsilon_T(r) = 0 \quad \text{and} \quad \lim_{r \rightarrow 0^+} \nu_T(r) = 1.$$

Recall  $\rho_T$  and  $\mathcal{V}_T$  defined in (35) and (37). The next lemmas state that if  $\varepsilon_\infty(r/\rho_T)$  (resp.  $\nu_\infty(r\rho_T)$ ) is positive and if the approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_\infty$  is good, i.e.  $\mathcal{V}_T$  is small, then  $\varepsilon_T(r)$  (resp.  $\nu_T(r)$ ) is also positive.

**Lemma 7.1.** *Let  $T \in \mathbb{N}$ . Suppose Assumptions 3.1, 3.2 and 5.1 hold. Then we have for  $r > 0$ :*

$$\varepsilon_T(r) \geq \varepsilon_\infty(r/\rho_T) - \mathcal{V}_T \quad \text{and} \quad \nu_T(r) \geq \nu_\infty(r\rho_T) - \mathcal{V}_T.$$

**Proof.** As Assumptions 3.2 and 5.1 hold, recall that  $\mathfrak{d}_\infty/\rho_T \leq \mathfrak{d}_T \leq \rho_T \mathfrak{d}_\infty$  on  $\Theta_T^2$ , see (36).

Let  $\theta, \theta' \in \Theta_T$  such that  $\mathfrak{d}_T(\theta', \theta) \geq r$ . We have  $\mathfrak{d}_\infty(\theta', \theta) \geq r/\rho_T$ . We get from the definition of  $\mathcal{V}_T$  that:

$$|\mathcal{K}_T(\theta, \theta')| \leq |\mathcal{K}_\infty(\theta, \theta')| + \mathcal{V}_T \leq 1 - \varepsilon_\infty(r/\rho_T) + \mathcal{V}_T.$$

Then, use (38) to get  $\varepsilon_T(r) \geq \varepsilon_\infty(r/\rho_T) - \mathcal{V}_T$ . We also have  $\mathfrak{d}_\infty(\theta', \theta) \leq r\rho_T$ . We deduce that:

$$-\mathcal{K}_T^{[0,2]}(\theta, \theta') \geq -\mathcal{K}_\infty^{[0,2]}(\theta, \theta') - \mathcal{V}_T \geq \nu_\infty(r\rho_T) - \mathcal{V}_T.$$

Finally, using (39), we obtain  $\nu_T(r) \geq \nu_\infty(r\rho_T) - \mathcal{V}_T$ . □

When we require in addition of the assumptions of Lemma 7.1 that  $\varepsilon_\infty(r/\rho_T) \wedge \nu_\infty(r\rho_T) > \mathcal{V}_T \geq 0$ , then we have  $\varepsilon_T(r) > 0$  and  $\nu_T(r) > 0$ .

## 7.2. Separation conditions for the non-linear parameters

In what follows, we measure the interferences (or the overlap) between the features in the mixture through a quantity  $\delta_T$  introduced in [41] and defined below. Let  $T \in \mathbb{N}$ ,  $\delta > 0$  and  $s \in \mathbb{N}^*$ . We define the set  $\Theta_{T,\delta}^s \subset \Theta_T^s$  of vector of parameters of dimension  $s \in \mathbb{N}^*$  and separation  $\delta > 0$  as:

$$\Theta_{T,\delta}^s = \left\{ (\theta_1, \dots, \theta_s) \in \Theta_T^s : \mathfrak{d}_T(\theta_\ell, \theta_k) > \delta \text{ for all distinct } k, \ell \in \{1, \dots, s\} \right\}.$$

Using the convention  $\inf \emptyset = +\infty$ , we set for  $u > 0$ :

$$(41) \quad \delta_T(u, s) = \inf \left\{ \delta > 0 : \max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s |\mathcal{K}_T^{[i,j]}(\theta_\ell, \theta_k)| \leq u \right. \\ \left. \text{for all } (i, j) \in \{0, 1\} \times \{0, 1, 2\} \text{ and } (\theta_1, \dots, \theta_s) \in \Theta_{T,\delta}^s \right\}.$$

The quantity  $\delta_T(u, s)$  is the minimum distance (with respect to the Riemannian metric  $\mathfrak{d}_T$ ) between  $s$  parameters so that the coherence of the associated dictionary is bounded by  $u$ . The notion of coherence between the features in the definition of  $\delta_T(u, s)$  is quite similar to the one used in compressed sensing (see [30, Section 5]). A standard problem in compressed sensing is to retrieve the vector  $\beta^*$  when the multivariate function  $\Phi_T(\vartheta^*)$  is known in the discrete setting of Section 1.2.1. In this framework, the matrix  $\Phi_T(\vartheta^*)$ , whose rows correspond to the  $K$  discretized functions in the dictionary, is known. The coherence is defined as  $\max_{1 \leq k \neq \ell \leq K} |\mathcal{K}_T(\theta_k^*, \theta_\ell^*)|$ . Usually, the smaller the coherence, the easier it is to retrieve the parameter  $\beta^*$ . The Babel function, introduced in [49], is even closer to our measure of overlap. We refer to [41] for a discussion on this function.

*Remark 7.2* (Rewriting the separation condition with operator norm). We shall stress that the definition of  $\delta_T$  in (41) is related to the operator norm  $\|\cdot\|_{\text{op}}$  associated to the  $\ell_\infty$  norm on  $\mathbb{R}^s$ . We restate (41) using this operator norm  $\|\cdot\|_{\text{op}}$ , and leave the interested reader to check that another choice of operator norm does not improve the bounds on the certificates. Let us define for  $i, j = 0, 1, 2$  (assuming the kernel  $\mathcal{K}_T$  is smooth enough) and  $\vartheta = (\theta_1, \dots, \theta_s) \in \Theta_T^s$  the  $s \times s$  matrix:

$$(42) \quad \mathcal{K}_T^{[i,j]}(\vartheta) = \left( \mathcal{K}_T^{[i,j]}(\theta_k, \theta_\ell) \right)_{1 \leq k, \ell \leq s}.$$

Let  $I$  be the identity matrix of size  $s \times s$ . For  $i = 0$  or  $i = 1$ , since the diagonal coefficients of  $\mathcal{K}_T^{[i,i]}(\vartheta)$  are equal to 1, see (28), we get:

$$\left\| I - \mathcal{K}_T^{[i,i]}(\vartheta) \right\|_{\text{op}} = \max_{1 \leq k \leq s} \sum_{\ell \neq k} |\mathcal{K}_T^{[i,i]}(\theta_k, \theta_\ell)|.$$

Since the diagonal coefficients of  $\mathcal{K}_T^{[1,0]}(\vartheta)$ ,  $\mathcal{K}_T^{[0,1]}(\vartheta)$  and  $\mathcal{K}_T^{[1,2]}(\vartheta)$  are zero, see (32), we also get:

$$\left\| \mathcal{K}_T^{[1,0]}(\vartheta) \right\|_{\text{op}} = \max_{1 \leq k \leq s} \sum_{\ell \neq k} |\mathcal{K}_T^{[1,0]}(\theta_k, \theta_\ell)| \quad \text{and} \quad \left\| \mathcal{K}_T^{[1,2]}(\vartheta) \right\|_{\text{op}} = \max_{1 \leq k \leq s} \sum_{\ell \neq k} |\mathcal{K}_T^{[1,2]}(\theta_k, \theta_\ell)|$$

and by symmetry, with  $\|\cdot\|_{\text{op}}^*$  for the operator norm associated to the  $\ell_1$  norm:

$$\left\| \mathcal{K}_T^{[0,1]}(\vartheta) \right\|_{\text{op}} = \left\| \mathcal{K}_T^{[1,0]\top}(\vartheta) \right\|_{\text{op}} = \left\| \mathcal{K}_T^{[1,0]}(\vartheta) \right\|_{\text{op}}^* = \max_{1 \leq \ell \leq s} \sum_{k \neq \ell} |\mathcal{K}_T^{[1,0]}(\theta_k, \theta_\ell)| = \max_{1 \leq k \leq s} \sum_{\ell \neq k} |\mathcal{K}_T^{[0,1]}(\theta_k, \theta_\ell)|.$$

Since the diagonal coefficients of  $\mathcal{K}_T^{[2,0]}(\vartheta)$  are equal to -1, see (32), we also get:

$$\left\| I + \mathcal{K}_T^{[2,0]}(\vartheta) \right\|_{\text{op}} = \max_{1 \leq k \leq s} \sum_{\ell \neq k} |\mathcal{K}_T^{[2,0]}(\theta_k, \theta_\ell)|.$$

Thus, we have:

$$(43) \quad \delta_T(u, s) = \inf \left\{ \delta > 0 : A_{T,\ell_\infty}(\vartheta) \leq u, \vartheta \in \Theta_{T,\delta}^s \right\},$$

where:

$$(44) \quad A_{T, \ell_\infty}(\vartheta) = \max \left( \left\| I - \mathcal{K}_T^{[0,0]}(\vartheta) \right\|_{\text{op}}, \left\| I - \mathcal{K}_T^{[1,1]}(\vartheta) \right\|_{\text{op}}, \left\| I + \mathcal{K}_T^{[2,0]}(\vartheta) \right\|_{\text{op}}, \left\| \mathcal{K}_T^{[1,0]}(\vartheta) \right\|_{\text{op}}, \left\| \mathcal{K}_T^{[0,1]}(\vartheta) \right\|_{\text{op}}, \left\| \mathcal{K}_T^{[1,2]}(\vartheta) \right\|_{\text{op}} \right).$$

Lemma 7.3 below enables us to compare the separation distance at  $T$  fixed and at the limit case where  $T = +\infty$ . Recall that the constant  $\rho_T$  is defined in (35).

**Lemma 7.3.** *Let  $T \in \bar{\mathbb{N}}$  and  $s \in \mathbb{N}^*$ . Suppose Assumptions 3.1, 3.2 and 5.1 hold. Then, for  $u > 0$  and with:*

$$u_T(s) = u + (s-1)\mathcal{V}_T,$$

*we have:*

$$\delta_T(u_T(s), s) \leq \rho_T \delta_\infty(u, s) \quad \text{and} \quad \Theta_{T, \rho_T \delta_\infty(u, s)}^s \subseteq \Theta_{T, \delta_T(u_T(s), s)}^s.$$

**Proof.** Since Assumptions 3.2 and 5.1 hold, we have from (36) that  $\mathfrak{d}_T \leq \rho_T \mathfrak{d}_\infty$  on  $\Theta_T^2$ . Hence for any  $\delta > 0$ , we have the inclusion  $\Theta_{T, \rho_T \delta}^s \subseteq \Theta_{\infty, \delta}^s$ . In particular, we have for  $u > 0$  that  $\Theta_{T, \rho_T \delta_\infty(u, s)}^s \subseteq \Theta_{\infty, \delta_\infty(u, s)}^s$ . Using the triangle inequality and the definition of  $\mathcal{V}_T$  in (37), we have that for  $(i, j) \in \{0, 1\} \times \{0, 1, 2\}$  and  $(\theta_1, \dots, \theta_s) \in \Theta_T^s$ :

$$\sum_{k=1, k \neq \ell}^s |\mathcal{K}_T^{[i,j]}(\theta_\ell, \theta_k)| \leq \sum_{k=1, k \neq \ell}^s \left( |\mathcal{K}_\infty^{[i,j]}(\theta_\ell, \theta_k)| + \mathcal{V}_T \right).$$

Then, the inclusion  $\Theta_{T, \rho_T \delta_\infty(u, s)}^s \subseteq \Theta_{\infty, \delta_\infty(u, s)}^s$  gives that for all  $(i, j) \in \{0, 1\} \times \{0, 1, 2\}$  and  $(\theta_1, \dots, \theta_s) \in \Theta_{T, \rho_T \delta_\infty(u, s)}^s$ :

$$\sum_{k=1, k \neq \ell}^s |\mathcal{K}_T^{[i,j]}(\theta_\ell, \theta_k)| \leq u + (s-1)\mathcal{V}_T.$$

With  $u_T(s) = u + (s-1)\mathcal{V}_T$ , we deduce that  $\delta_T(u_T(s), s) \leq \rho_T \delta_\infty(u, s)$ , which proves the inclusion  $\Theta_{T, \rho_T \delta_\infty(u, s)}^s \subseteq \Theta_{T, \delta_T(u_T(s), s)}^s$ .  $\square$

### 7.3. The interpolating certificates

We define quantities which depend on  $\mathcal{K}_\infty$ ,  $\Theta_\infty$  and on real parameters  $r > 0$  and  $\rho \geq 1$ :

$$(45) \quad \begin{aligned} H_\infty^{(1)}(r, \rho) &= \frac{1}{2} \wedge L_{2,0} \wedge L_{2,1} \wedge \frac{\nu_\infty(\rho r)}{10} \wedge \frac{\varepsilon_\infty(r/\rho)}{10}, \\ H_\infty^{(2)}(r, \rho) &= \frac{1}{6} \wedge \frac{8\varepsilon_\infty(r/\rho)}{10(5 + 2L_{1,0})} \wedge \frac{8\nu_\infty(\rho r)}{9(2L_{2,0} + 2L_{2,1} + 4)}, \end{aligned}$$

where the constants involved are defined in (34). By recalling the behaviors of  $\varepsilon_\infty(r)$  and  $\nu_\infty(r)$  when  $r$  goes down to zero from (40), we have for  $\rho \geq 1$ :

$$\lim_{r \rightarrow 0^+} H_\infty^{(1)}(r, \rho) = 0 \quad \text{and} \quad \lim_{r \rightarrow 0^+} H_\infty^{(2)}(r, \rho) = 0.$$

We state the first main result of this section whose proof is given in Section B.

**Proposition 7.4** (Interpolating certificate). *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}^*$ ,  $\rho \geq 1$  and  $r > 0$ . We assume that:*

- (i) **Regularity of the dictionary  $\varphi_T$ :** Assumptions 3.1 and 3.2 hold.
- (ii) **Regularity of the limit kernel  $\mathcal{K}_\infty$ :** Assumption 5.1 holds, we have  $r \in (0, 1/\sqrt{2L_{2,0}})$ , and also  $\varepsilon_\infty(r/\rho) > 0$  and  $\nu_\infty(\rho r) > 0$ .
- (iii) **Separation of the non-linear parameters:** There exists  $u_\infty \in (0, H_\infty^{(2)}(r, \rho))$  such that:

$$\delta_\infty(u_\infty, s) < +\infty.$$

- (iv) **Closeness of the metrics  $\mathfrak{d}_T$  and  $\mathfrak{d}_\infty$** : We have  $\rho_T \leq \rho$ .  
(v) **Proximity of the kernels  $\mathcal{K}_T$  and  $\mathcal{K}_\infty$** :

$$\mathcal{V}_T \leq H_\infty^{(1)}(r, \rho) \quad \text{and} \quad (s-1)\mathcal{V}_T \leq H_\infty^{(2)}(r, \rho) - u_\infty.$$

Then, with the positive constants:

$$(46) \quad C_N = C_N(r) = \frac{\nu_\infty(\rho r)}{180}, \quad C'_N = \frac{5}{8}L_{2,0} + \frac{1}{8}L_{2,1} + \frac{1}{2}, \quad C_B = 2 \quad \text{and} \quad C_F = C_F(r) = \frac{\varepsilon_\infty(r/\rho)}{10} \leq 1,$$

Assumption 6.1 holds (with the same  $r$ ) for any subset  $\mathcal{Q}^* = \{\theta_i^*, 1 \leq i \leq s\}$  such that for all  $\theta \neq \theta' \in \mathcal{Q}^*$ :

$$\mathfrak{d}_T(\theta, \theta') > 2 \max(r, \rho_T \delta_\infty(u_\infty, s)).$$

Note that (i) concerns the dictionary  $\varphi_T$ , (ii) and (iii) the limit kernel  $\mathcal{K}_\infty$  and the set of parameters, and (iv) and (v) the regime for the parameters  $s$  and  $T$ .

**Remark 7.5** (On the assumptions of Proposition 7.4 when  $\mathcal{K}_\infty = \mathcal{K}_T$ ). In the setting where the limit kernel and the approximating kernel are equal, some assumptions in the proposition become less restrictive, without any changes to the proofs. If  $\mathcal{K}_\infty$  is chosen equal to  $\mathcal{K}_T$ , then  $\mathcal{V}_T = 0$  and  $\rho_T = 1$ , and also (iv) and (v) hold and  $\rho$  can be chosen equal to 1 and  $u_\infty$  can be chosen equal to  $H_\infty^{(2)}(r, 1)$ .

We now give the second main result of this section whose proof is given in Section B.2.

**Proposition 7.6** (Interpolating derivative certificate). *Let  $T \in \mathbb{N}$  and  $s \in \mathbb{N}^*$ . We assume that:*

- (i) **Regularity of the dictionary  $\varphi_T$** : Assumptions 3.1 and 3.2 hold.  
(ii) **Regularity of the limit kernel  $\mathcal{K}_\infty$** : Assumption 5.1 holds.  
(iii) **Separation of the non-linear parameters**: There exists  $u'_\infty \in (0, 1/6)$ , such that:

$$\delta_\infty(u'_\infty, s) < +\infty.$$

- (iv) **Proximity of the kernels  $\mathcal{K}_T$  and  $\mathcal{K}_\infty$** : We have:

$$\mathcal{V}_T \leq 1 \quad \text{and} \quad (s-1)\mathcal{V}_T + u'_\infty \leq 1/6.$$

Then, with the positive constants:

$$(47) \quad c_N = \frac{1}{8}L_{2,0} + \frac{5}{8}L_{2,1} + \frac{7}{8}, \quad c_B = 2 \quad \text{and} \quad c_F = \frac{5}{4}L_{1,0} + \frac{7}{4},$$

Assumption 6.2 holds for any  $r > 0$  and any subset  $\mathcal{Q}^* = \{\theta_i^*, 1 \leq i \leq s\}$  such that for all  $\theta \neq \theta' \in \mathcal{Q}^*$ :

$$\mathfrak{d}_T(\theta, \theta') > 2 \max(r, \rho_T \delta_\infty(u'_\infty, s)).$$

Let us briefly indicate how the certificates are constructed in Section B using the features of the dictionary. Let  $\alpha = (\alpha_1, \dots, \alpha_s)$  and  $\xi = (\xi_1, \dots, \xi_s)$  be elements of  $\mathbb{R}^s$ . Let  $p_{\alpha, \xi} \in H_T$  be defined by:

$$p_{\alpha, \xi} = \sum_{k=1}^s \alpha_k \phi_T(\theta_k^*) + \sum_{k=1}^s \xi_k \phi_T^{[1]}(\theta_k^*),$$

where  $\phi_T^{[1]}$  denotes the derivative  $\tilde{D}_{1;T}[\phi_T]$ . Using (31) in Lemma 4.3, set the interpolating real-valued function  $\eta_{\alpha, \xi}$  defined on  $\Theta$  by:

$$\eta_{\alpha, \xi}(\theta) := \langle \phi_T(\theta), p_{\alpha, \xi} \rangle_T = \sum_{k=1}^s \alpha_k \mathcal{K}_T(\theta, \theta_k^*) + \sum_{k=1}^s \xi_k \mathcal{K}_T^{[0,1]}(\theta, \theta_k^*).$$

By Assumption 3.2 on the regularity of  $\varphi_T$  and the positivity of  $g_T$  and Lemma 4.3, we get that the function  $\eta_{\alpha, \xi}$  is of class  $\mathcal{C}^3$  on  $\Theta$ , and using (23), we get that:

$$\eta_{\alpha, \xi}^{[1]} := \tilde{D}_{1;T}[\eta_{\alpha, \xi}](\theta) = \sum_{k=1}^s \alpha_k \mathcal{K}_T^{[1,0]}(\theta, \theta_k^*) + \sum_{k=1}^s \xi_k \mathcal{K}_T^{[1,1]}(\theta, \theta_k^*).$$

We show in Section B that for any function  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  there exists a unique choice of  $\alpha$  and  $\xi$  such that  $\eta_{\alpha, \xi}$  becomes an interpolating certificate, that is,  $\eta_{\alpha, \xi} = v$  and  $\eta_{\alpha, \xi}^{[1]} = 0$  on  $\mathcal{Q}^*$ , and  $p_{\alpha, \xi}$  satisfies Points (i)-(iv) of Assumption 6.1.

Moreover, for any function  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  there exists another unique choice of  $\alpha$  and  $\xi$  such that  $\eta_{\alpha, \xi}$  is an interpolating derivative certificate, that is,  $\eta_{\alpha, \xi} = 0$  and  $\eta_{\alpha, \xi}^{[1]} = v$  on  $\mathcal{Q}^*$ , and  $p_{\alpha, \xi}$  satisfies Points (i)-(iii) of Assumption 6.2.

## 8. Gaussian sparse spike deconvolution

We develop here in full details the particular example of a mixture of Gaussian features observed in a discrete regression model with regular design. In particular, we check the numerous but not very restrictive assumptions, and we illustrate that our general and more restrictive sufficient conditions for the existence of certificates can turn simpler and far less restrictive on concrete examples. The model is presented in Section 8.1, where we also check the first assumptions. The technical Section 8.2 on the existence of the certificates allows to point out the separation distance in (54) and with the simpler expression in (55). This separation distance is usually very pessimistic, but one can rely on numerical estimations to be more realistic, see Remark 8.2 in this direction. Eventually, we apply to this context our main Theorem 2.1 in Section 8.3 as Corollary 8.3 and illustrate a particular choice of the tuning parameter in Remark 8.4 in the spirit of [12, 47] established for the specific dictionary of complex exponentials.

### 8.1. Model and first assumptions of Theorem 2.1

Consider a real-valued process  $y$  observed over a regular grid  $t_1 < \dots < t_T$  of a symmetric interval  $[a_T, b_T]$ , with  $T \geq 2$ ,  $b_t = -a_T > 0$ ,  $t_j = a_T + j\Delta_T$  for  $j = 1, \dots, T$  and grid step:

$$\Delta_T = \frac{b_T - a_T}{T}.$$

Assuming that all the observations have the same weight amounts to considering  $y$  as an element of the Hilbert space  $H_T = L^2(\lambda_T)$  of real valued functions defined in  $\mathbb{R}$  and square integrable with respect to the atomic measure  $\lambda_T$  on  $\{t_1, \dots, t_T\}$ :

$$\lambda_T(dt) = \Delta_T \sum_{j=1}^T \delta_{t_j}(dt).$$

We consider a noise process  $w_T(t) = \sum_{j=1}^T G_j \mathbf{1}_{\{t_j=t\}}$  for  $t \in \mathbb{R}$ , where  $(G_1, \dots, G_T)$  is a centered Gaussian vector such that, for some noise level  $\sigma_1 > 0$ :

$$\mathbb{E}[G_j^2] = \sigma_1^2 \quad \text{and} \quad |\mathbb{E}[G_j G_i]| \leq \sigma_1^2/T \quad \text{for } j \neq i \text{ in } \{1, \dots, T\}.$$

Thus, the norm of the noise  $\|w_T\|_T$  is finite almost surely, and for any  $f \in L^2(\lambda_T)$  we have:

$$\text{Var}(\langle f, w_T \rangle_T) = \text{Var}\left(\Delta_T \sum_{j=1}^T f(t_j) G_j\right) \leq 2\sigma_1^2 \Delta_T \|f\|_T^2.$$

Hence, Assumption 1.1 on the noise is satisfied with  $\sigma^2 = 2\sigma_1^2$ . (Notice that if the random variables  $G_1, \dots, G_T$  are independent, then  $\text{Var}(\langle f, w_T \rangle_T) = \sigma^2 \Delta_T \|f\|_T^2$  with  $\sigma^2 = \sigma_1^2$ .) This gives that Point (i) of Theorem 2.1 holds.

We consider the dictionary given by the translation model of Section 3.2.1 with Gaussian features and fixed scaling parameter  $\sigma_0 > 0$ , that is the dictionary does not depend on  $T$  and is given by:

$$\left(\varphi(\theta) = k\left(\frac{\cdot - \theta}{\sigma_0}\right), \theta \in \Theta\right) \quad \text{with} \quad k(t) = e^{-t^2/2} \quad \text{and} \quad \Theta = \mathbb{R}.$$

Thus, the signal  $\beta^* \Phi(\vartheta^*)$  in model (1) can indeed be written as the convolution product of the function  $k$  and an atomic measure. It is elementary to check that Assumption 3.1 on the regularity of the features holds. Furthermore, the functions  $\varphi(\theta)$  and  $\partial_\theta \varphi(\theta)$  are linearly independent  $\lambda_T - a.e$  for all  $\theta \in \Theta$  as  $T \geq 2$ . Hence the function  $g_T$  is positive on  $\Theta$  by Lemma 3.1 and thus Assumption 3.2 holds. This gives that Point (ii) of Theorem 2.1 holds.



We now define the limit kernel  $\mathcal{K}_\infty$ . To do so, we shall assume that  $(b_T, T \geq 2)$  is a sequence of positive numbers, such that:

$$(48) \quad \lim_{T \rightarrow \infty} b_T = +\infty \quad \text{and} \quad \lim_{T \rightarrow \infty} \Delta_T = 0.$$

This in particular implies that the sequence of measures  $(\lambda_T, T \geq 2)$  converges with respect to the vague topology towards the Lebesgue measure, say  $\lambda_\infty$ , on  $\Theta_\infty = \mathbb{R}$ . We also consider the Hilbert space  $H_\infty = L^2(\lambda_\infty)$  endowed with its usual scalar product denoted  $\langle \cdot, \cdot \rangle_\infty$  and corresponding norm denoted  $\|\cdot\|_\infty$  (not to be confused with the supremum norm!). Note that the kernel  $\mathcal{K}_T$  and the associated quantities such as  $\varepsilon_T$  and  $\nu_T$  defined in (38) and (39), respectively, or the uniform bounds on  $\mathcal{K}_T^{[i,j]}$ , are difficult to calculate. However the uniform bounds on  $\Theta_\infty = \mathbb{R}$  for the kernel  $\mathcal{K}_\infty$ , defined by (29) with  $T$  replaced by  $\infty$ , are easily computed. Elementary calculations give for  $\theta, \theta' \in \Theta$ :

$$\|\varphi(\theta)\|_\infty^2 = \sqrt{\pi} \sigma_0, \quad \phi_\infty(\theta) = \frac{1}{\pi^{1/4} \sqrt{\sigma_0}} \varphi(\theta), \quad \mathcal{K}_\infty(\theta, \theta') = k\left(\frac{\theta - \theta'}{\sqrt{2} \sigma_0}\right) \quad \text{and} \quad g_\infty(\theta) = \frac{1}{2\sigma_0^2}.$$

In particular, we have  $g'_\infty(\theta) = 0$ . The Riemannian metric is equal to the Euclidean distance up to a multiplicative factor, for all  $\theta, \theta' \in \Theta_\infty = \mathbb{R}$ :

$$(49) \quad \mathfrak{d}_\infty(\theta, \theta') = \frac{|\theta - \theta'|}{\sqrt{2} \sigma_0}.$$

We see that  $\mathcal{K}_\infty$  is of class  $\mathcal{C}^{\infty, \infty}$  and that:

$$(50) \quad \mathcal{K}_\infty^{[i,j]}(\theta, \theta') = (-1)^j k^{(i+j)}\left(\frac{\theta - \theta'}{\sqrt{2} \sigma_0}\right) \quad \text{and} \quad k^{(i)}(t) = P_i(t) k(t),$$

where we give for convenience the formulae for some of the polynomials  $P_i$ :

$$P_1(t) = -t, \quad P_2(t) = -1 + t^2, \quad P_3(t) = 3t - t^3, \quad P_4(t) = 3 - 6t^2 + t^4, \quad P_6(t) = -15 + 45t^2 - 15t^4 + t^6.$$

Then, we explicitly compute the constants  $L_{i,j}$  for  $i, j \in \{0, \dots, 2\}$  and  $L_3$  defined in (34):

$$m_g = (2\sigma_0^2)^{-1}, \quad L_{0,0} = 1, \quad L_{1,0} = L_{0,1} = e^{-1/2}, \quad L_{1,1} = L_{2,0} = L_{0,2} = 1, \\ L_{2,1} = L_{1,2} = \sqrt{18 - 6\sqrt{6}} e^{\sqrt{3/2} - 3/2} \leq \sqrt{2}, \quad L_{2,2} = 3 \quad \text{and} \quad L_3 = 15.$$

Notice the constants  $L_{i,j}$  and  $L_3$  do not depend on the scaling factor  $\sigma_0$ . Thus Assumption 5.1 holds. This gives that Point (iii) of Theorem 2.1 holds.

We now check the proximity of the kernel  $\mathcal{K}_T$  to the limit kernel  $\mathcal{K}_\infty$ . The support of  $\lambda_T$  is spread over the window  $[a_T, b_T]$  where the signal is observed. Hence it is legitimate to look for the location parameters on a smaller subset of this window, and thus restrict the optimization (4) to the compact set:

$$\Theta_T = [(1 - \epsilon)a_T, (1 - \epsilon)b_T] \subset [a_T, b_T] \quad \text{with a given shrinkage } \epsilon \in (0, 1).$$

The proof of the next lemma is given in Section C.6. Recall  $\rho_T$  and  $\mathcal{V}_T$  defined in (35) and (37). Set:

$$\gamma_T = 2\Delta_T \sigma_0^{-1} + \sqrt{\pi} e^{-\epsilon^2 b_T^2 / 2\sigma_0^2}.$$

**Lemma 8.1.** *There exist finite positive universal constants  $c_0$ ,  $c_1$  and  $c_2$ , such that  $\gamma_T < c_0$  implies:*

$$(51) \quad \mathcal{V}_T \leq c_1 \gamma_T \quad \text{and} \quad |1 - \rho_T| \leq c_2 \gamma_T.$$

This implies that Assumption 5.2 holds for  $T$  such that  $\gamma_T \leq c_0$  and  $c_1 \gamma_T \leq 3$ , which holds for  $T$  large enough thanks to (48). Thus Point (iv) of Theorem 2.1 holds for  $T$  large enough.

## 8.2. Existence of certificates

We keep the model and the notations from Section 8.1. In order to get the prediction error from Theorem 2.1, we only need to check that Point (iv) therein on the existence of the certificates holds. To check the existence of the certificates, we can use Propositions 7.4 and 7.6, and check that all the hypotheses required in those two propositions hold.

We first concentrate on the hypotheses of Proposition 7.4. Assumption (i) on the regularity of the dictionary holds, see Section 8.1.

We recall that  $L_{0,2} = 1$  and thus  $1/\sqrt{2L_{0,2}} = 1/\sqrt{2} > 1/2$ . Recall  $\varepsilon_\infty(r)$  and  $\nu_\infty(r)$  defined in (38) and (39), and thanks to the explicit form of the Riemannian metric, we get for  $r \in (0, 1)$ :

$$\varepsilon_\infty(r) = 1 - e^{-r^2/2} > 0 \quad \text{and} \quad \nu_\infty(r) = (1 - r^2) e^{-r^2/2}.$$

This and the regularity of the kernel  $\mathcal{K}_\infty$  from Section 8.1 imply that Assumption (ii) holds for all  $r \in (0, 1/(\rho \vee \sqrt{2}))$ .

We obtain from (50) that  $\lim_{q \rightarrow \infty} \sup_{|\theta - \theta'| \geq q} |\mathcal{K}_\infty^{[i,j]}(\theta, \theta')| = 0$  for all  $i, j \in \{0, 1, 2\}$ . Thus, we deduce from the definition (41) of  $\delta_\infty$  that  $\delta_\infty(u, s)$  is finite for all  $s \in \mathbb{N}^*$  and  $u > 0$ . This implies that Assumption (iii) on the separation of the parameters holds.

To simplify, we set  $\rho = 2$  (but we could take any value of  $\rho > 1$ ). We deduce from Lemma 8.1, that for  $T$  large enough  $\rho_T \leq \rho = 2$ , and thus Assumption (iv) on the closeness of the metrics  $\mathfrak{d}_T$  and  $\mathfrak{d}_\infty$  holds.

Recall the definition of  $H_\infty^{(1)}$  and  $H_\infty^{(2)}$  from (45). To get the smallest separation distance, we also set:

$$(52) \quad r = \operatorname{argmax}_{0 < r' < 1/2} H_\infty^{(2)}(r', \rho) \approx 0.49.$$

Notice that the function is not *a priori* monotone in  $\rho$ . We have  $\varepsilon_\infty(r/2) \approx 2.9 \times 10^{-2}$ ,  $\nu_\infty(2r) \approx 3.7 \times 10^{-2}$ ,  $H_\infty^{(1)}(r, 2) \approx 2.9 \times 10^{-3}$  and  $H_\infty^{(2)}(r, 2) \approx 3.7 \times 10^{-3}$ . Again in order to get a “small” separation distance, we choose  $u_\infty$  close to  $H_\infty^{(2)}(r, 2)$ , say  $u_\infty = \eta_0 H_\infty^{(2)}(r, 2)$  for some  $\eta_0 < 1$  close to 1. For simplicity set  $\eta_0 = 9/10$ . Thanks to hypothesis (48), we get  $\lim_{T \rightarrow \infty} \gamma_T = 0$  and Lemma 8.1 implies that for  $T$  large enough, depending on  $\sigma_0$ ,  $\epsilon$  and the sparsity parameter  $s$ , we have:

$$(53) \quad \rho_T \leq 2, \quad \mathcal{V}_T \leq H_\infty^{(1)}(r, 2) \quad \text{and} \quad (s-1)\mathcal{V}_T \leq (1-\eta_0)H_\infty^{(2)}(r, 2),$$

and thus Assumption (v) on the proximity of the kernels  $\mathcal{K}_T$  and  $\mathcal{K}_\infty$  holds.

Thus, the assumptions of Proposition 7.4 are satisfied, and we deduce that Assumption 6.1 holds with, thanks to (46):

$$C_N \approx 2 \times 10^{-4}, \quad C'_N \approx 1.3, \quad C_B = 2 \quad \text{and} \quad C_F \approx 2.9 \times 10^{-3}.$$

We now concentrate on the hypotheses of Proposition 7.6. Assumptions (i)-(iii) clearly hold for the same reasons as Assumptions (i)-(iii) of Proposition 7.4.

Again in order to get a “small” separation distance, there is no need to choose  $u'_\infty$  larger than  $u_\infty$ , and for this reason we take  $u'_\infty = u_\infty$ . We deduce from (53) that for  $T$  large enough, depending on  $\sigma_0$ ,  $\epsilon$  and the sparsity parameter  $s$ :

$$\mathcal{V}_T \leq 1 \quad \text{and} \quad (s-1)\mathcal{V}_T + u'_\infty \leq 1/6,$$

and thus Assumption (iv) on the proximity of the kernels  $\mathcal{K}_T$  and  $\mathcal{K}_\infty$  holds.

Thus, the assumptions of Proposition 7.6 are satisfied, and we deduce, thanks to (47), that Assumption 6.2 holds with the same value of  $r$  given by (52):

$$c_N \approx 1.9, \quad c_B = 2, \quad \text{and} \quad c_F \approx 2.6.$$

In conclusion, we get that Assumptions 6.1 and 6.2 hold for  $T$  large enough, and thus Point (v) of Theorem 2.1 holds for  $T$  large enough and  $\mathcal{Q}^*$  such that for all  $\theta \neq \theta' \in \mathcal{Q}^*$  the distance  $\mathfrak{d}_T(\theta, \theta')$  is larger than the separation distance:

$$(54) \quad 2 \max(r, \rho_T \delta_\infty(u_\infty, s), \rho_T \delta_\infty(u'_\infty, s)).$$

Notice that since  $u_\infty = u'_\infty$ ,  $\rho_T \mathfrak{d}_T(\theta, \theta') \geq \mathfrak{d}_\infty(\theta, \theta')$  and  $\rho_T \leq 2$ , we deduce from (49), that a slightly stronger condition is to assume that  $|\theta - \theta'|$  is larger than:

$$(55) \quad \sqrt{2} \sigma_0 \max(1, 4\delta_\infty(u_\infty, s)).$$

*Remark 8.2* (On the separation distance (54)). The separation distance (54) is a non-decreasing function of  $s$ . We now provide an upper bound. Let  $(i, j) \in \{0, 1\} \times \{0, 1, 2\}$ . By considering the kernel  $\mathcal{K}_T$  and its derivative given by (50) and the bound  $M = \max_{0 \leq i \leq 3} \sup |P_i| \sqrt{k}$ , we deduce that  $|\mathcal{K}_\infty^{[i,j]}(\theta, \theta')| \leq M e^{-\mathfrak{d}_\infty(\theta, \theta')^2/2}$  for all  $\theta, \theta' \in \Theta$ . We easily obtain that for  $\vartheta = (\theta_1, \dots, \theta_s) \in \Theta_{\infty, \delta}^s$  with  $\delta > 0$ :

$$\max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s |\mathcal{K}_\infty^{[i,j]}(\theta_\ell, \theta_k)| \leq \psi_s(\delta) \quad \text{with} \quad \psi_s(\delta) = 2M \int_0^{s/2+1} e^{-t^2 \delta^2/4} dt.$$

The function  $\psi_s$  is decreasing and one to one from  $\mathbb{R}_+$  to  $(0, M(s+2)]$ . Setting  $\psi_s^{-1}(u) = 0$  for  $u > M(s+2)$ , we deduce from (41) that for  $u > 0$ :

$$\delta_\infty(u, s) \leq \psi_s^{-1}(u).$$

Since the map  $\delta \mapsto \psi_s(\delta)$  is decreasing and the map  $s \mapsto \psi_s(\delta)$  is increasing with limit  $\psi_\infty(\delta) = 2\sqrt{\pi} M/\delta$ , we deduce that for  $s \in \mathbb{N}^*$ :

$$\delta_\infty(u, s) \leq \frac{2\sqrt{\pi} M}{u},$$

so that the separation distance (54) (or (55)) can be bounded uniformly in  $s$  for given  $r$  and  $u_\infty = u'_\infty$ .

In fact, we shall illustrate for  $s = 2$  that the separation distance (54) is largely overestimated. We can compute  $\delta_\infty(u, s)$  thanks to its expression (43). For  $s = 2$  and with the values chosen in this section for  $u_\infty = u'_\infty$ , we obtain  $\delta_\infty(u_\infty, 2) \approx 4.5$ . We deduce that the separation distance (54) expressed with respect to the metric  $\mathfrak{d}_T$  is approximately  $9\rho_T$  (which gives  $13\sigma_0\rho_T^2$  in terms of the Euclidean metric), which is unconveniently large. However, a detailed numerical approach (using the very certificates provided in the proof of Propositions 7.4 and 7.6) with  $T$  large so that the kernel  $\mathcal{K}_T$  is indeed well approximated by  $\mathcal{K}_\infty$  (and thus  $\rho_T \approx 1$ ), gives that one can take for  $s = 2$  the separation distance with respect to the Euclidean metric equal to  $3.1 \times \sigma_0$  (that is approximately equal to 2.2 with respect to the metric  $\mathfrak{d}_\infty$ ), which is much more realistic. Therefore, the theoretical separation distance (54) is in general largely overestimated.

### 8.3. Prediction error

We keep the model and the notations from Section 8.1 and the values chosen in Section 8.2. We deduce from Theorem 2.1 the following result.

**Corollary 8.3.** *For  $T$  large enough, depending on  $\sigma_0$ ,  $\epsilon$  and the sparsity parameter  $s$ , such that (53) holds and for all  $\theta \neq \theta' \in \mathcal{Q}^* = \{\theta_k^*, k \in S^*\}$ , with  $S^* = \text{Supp}(\beta^*)$  such that  $|\theta - \theta'|$  is larger than the separation parameter  $\sqrt{2}\sigma_0 \max(1, 4\delta_\infty(u_\infty, s))$  given by (55), then, with some universal finite constants  $\mathcal{C}_0, \dots, \mathcal{C}_3 > 0$ , for any  $\tau > 1$  and a tuning parameter:*

$$(56) \quad \kappa \geq \mathcal{C}_1 \sigma \sqrt{\Delta_T \log(\tau)},$$

*we have the prediction error bound of the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in (4) given by:*

$$(57) \quad \sqrt{\Delta_T} \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{\ell_2} \leq \mathcal{C}_0 \sqrt{s} \kappa,$$

*with probability larger than  $1 - \mathcal{C}_2 \left( \frac{\sqrt{2b_T}}{\sigma_0 \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right)$ . Moreover, with the same probability, we have that  $\left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| \leq \mathcal{C}_3 \kappa s$  as well as the inequalities (11) of Theorem 2.5.*

The values of the universal constants  $\mathcal{C}_i$ ,  $i = 0, \dots, 3$ , can be given explicitly and they are large, but they could be improved numerically.

*Remark 8.4* (A particular choice of the tuning parameter). Let  $\gamma > 0$  and  $\gamma' \geq \gamma$  such that  $1 > \gamma' - \gamma$ . Set  $\tau = T^{\gamma'}$ ,  $b_T = \sigma_0 T^{\gamma' - \gamma} \sqrt{\log(T)}$  and  $\kappa = \mathcal{C}_1 \sigma \sqrt{\Delta_T \log(\tau)}$  (which corresponds to the equality in (56)). Then, we get under the assumptions of Corollary 8.3 (and thus  $T$  large enough) that:

$$\frac{1}{\sqrt{T}} \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{\ell_2} \leq \mathcal{C}_0'' \sigma \sqrt{s \frac{\log(T)}{T}},$$

with probability larger than  $1 - \mathcal{C}_2''/T^\gamma$  where  $\mathcal{C}_0'' = \sqrt{\gamma'} \mathcal{C}_0 \mathcal{C}_1$  and  $\mathcal{C}_2'' = \sqrt{2/\gamma'} \mathcal{C}_2$ . Hence, we obtain a similar prediction error bound as the one given in Remark 2.2, see (10). Notice however that in the model and references given in Remark 2.2, the Riemannian diameter of the parameter set  $\Theta_T$  is bounded by a constant free of  $T$ , whereas in this section it grows (sublinearly) with  $T$  without degrading the prediction error bound.

## 9. Scaled exponential model

We develop in this section an example involving a dictionary that is not translation invariant and for which the associated metric differs from the Euclidean metric. We consider a continuous dictionary composed of exponential functions continuously scaled which is used in microscopy where it is often necessary to invert a Laplace transform (see for instance [41], [24]).

### 9.1. The model

Consider a real-valued process  $y$  observed continuously over  $\mathbb{R}_+$  and assume that this process is an element of the Hilbert space  $H_T = L^2(\mathbb{R}_+, \text{Leb})$  where  $\text{Leb}$  denotes here the Lebesgue measure over  $\mathbb{R}_+$ . We write  $H$  instead of  $H_T$  for the Hilbert space and we write  $\langle \cdot, \cdot \rangle$  its scalar product and  $\|\cdot\|$  its associated norm.

We consider a truncated white noise as in Section 1.2.2 such that  $w_T = \sum_{k=1}^T (1/\sqrt{T}) G_k \psi_k$ , where  $(G_k, k \in \mathbb{N})$  are independent centered Gaussian random variables with variance  $\sigma^2$  and  $(\psi_k, k \in \mathbb{N})$  denotes an orthonormal basis of  $H$ . Hence Assumption 1.1 holds as  $\|w_T\|^2 = \sum_{k=1}^T G_k^2/T$  is a.s. finite and  $\text{Var}(\langle f, w_T \rangle) \leq \sigma^2 \Delta_T \|f\|^2$  with  $\Delta_T = 1/T$ . This gives that Point (i) of Theorem 2.1 holds.

*Remark 9.1.* We stress that by the law of large numbers  $\|w_T\|$  tends almost surely to  $\sigma > 0$ . Therefore the upper bounds in previous results on super-resolution and BLasso (see [26] or [41]) which hold when  $\|w_T\|$  tends to zero do not apply here.

We consider the dictionary given by the scaling exponential model of Section 3.2.4 given by:

$$\left( \varphi(\theta) = k(\theta \cdot), \theta \in \Theta \right) \quad \text{with} \quad k(t) = e^{-t} \quad \text{and} \quad \Theta = \mathbb{R}_+^*.$$

We insist on the fact that in this example the dictionary and the observation space  $H$  do not depend on  $T$ . For simplicity we omit the index  $T$  for the quantities which shall not depend on  $T$ . As the kernels do not depend on  $T$ , we choose the limit kernel to be the same, i.e.,  $\mathcal{K} := \mathcal{K}_T = \mathcal{K}_\infty$ . In particular, Point (iv) of Theorem 2.1 holds automatically. One easily checks that Assumption 3.1 on the regularity of the features holds, and elementary calculations give for  $\theta, \theta' \in \Theta$ :

$$\|\varphi(\theta)\|^2 = 1/(2\theta), \quad \phi(\theta) = \sqrt{2\theta} e^{-\theta t}, \quad \mathcal{K}(\theta, \theta') = \frac{2\sqrt{\theta\theta'}}{\theta + \theta'} \quad \text{and} \quad g(\theta) = \frac{1}{4\theta^2}.$$

Since the function  $g$  is positive on  $\Theta$ , we get that Assumption 3.2 holds. This gives that Point (ii) of Theorem 2.1 holds. The Riemannian metric obtained from  $g$  is given by, for  $\theta, \theta' \in \Theta$ :

$$(58) \quad \mathfrak{d}(\theta, \theta') = \frac{1}{2} \left| \log \left( \frac{\theta}{\theta'} \right) \right|.$$

Notice it is not equivalent to the Euclidean distance on  $\Theta$ . We see that  $\mathcal{K}$  is of class  $\mathcal{C}^{3,3}$  and that:

$$\mathcal{K}^{[i,j]}(\theta, \theta') = (-1)^j f^{(i+j)} \left( \frac{1}{2} \log \left( \frac{\theta}{\theta'} \right) \right) \quad \text{with} \quad f(x) = \frac{1}{\cosh(x)}.$$

We shall retrieve scaling parameters over a compact set whose diameter may depend on  $T$ , for example we can take:

$$\Theta_T = [M_T^{-1}, M_T] \quad \text{with} \quad M_T > 1.$$

Assumption 5.1 holds on  $\Theta_\infty = \mathbb{R}_+^*$ . This gives that Point (iii) of Theorem 2.1 holds.

## 9.2. Existence of certificates

In order to get the prediction error from Theorem 2.1, it remains to show that Point (v) therein on the existence of the certificates holds. To check the existence of the certificates, we can use Propositions 7.4 and 7.6, and check that all the hypotheses required in those two propositions hold.

We show first that the hypotheses of Proposition 7.4 hold. Assumption (i) on the regularity of the dictionary holds, see Section above.

Elementary calculations give that  $L_{0,2} = 1$ . Recall  $\varepsilon_\infty(r)$  and  $\nu_\infty(r)$  defined in (38) and (39) (noted simply  $\varepsilon$  and  $\nu$  in this section). Let  $\theta < \theta'$  in  $\Theta$  and let us set  $r = \mathfrak{d}(\theta, \theta')$ . We have,  $\mathcal{K}(\theta, \theta') = f(r)$ . Since  $f$  is positive and decreasing on  $\mathbb{R}_+$ , we have for  $r > 0$ ,  $\varepsilon(r) = 1 - f(r) > 0$ . Similarly we have:

$$\mathcal{K}^{[0,2]}(\theta, \theta') = f^{(2)}(r) = \frac{1}{\cosh(r)^3} (\cosh(r)^2 - 2).$$

The function  $f^{(2)}$  is increasing and negative on  $(0, \log(1 + \sqrt{2}))$ . Hence, provided  $r < \log(1 + \sqrt{2})$ , we have  $\nu(r) = -f^{(2)}(r) > 0$ . This and the regularity of the kernel  $\mathcal{K}$  imply that Assumption (ii) of Proposition 7.4 holds for  $\rho = 1$  and all  $r \in (0, 1/\sqrt{2})$ .

Notice that  $f^{(i)}$  can be written as the ratio of a polynomial of degree  $i - 1$  in  $\cosh$  and  $\sinh$  and of  $\cosh^i$ . In particular, there exists a finite constant  $M$  such that for all  $i \in \{0, \dots, 3\}$  and  $x \in \mathbb{R}$ :

$$(59) \quad |f^{(i)}(x)| \leq M f(x).$$

So, we get that  $\lim_{q \rightarrow \infty} \sup_{\mathfrak{d}(\theta, \theta') \geq q} |\mathcal{K}^{[i,j]}(\theta, \theta')| = \lim_{r \rightarrow \infty} |f^{(i+j)}(r)| = 0$  for all  $i, j \in \{0, 1, 2\}$ . Thus, we deduce from the definition (41) of  $\delta_\infty$  that  $\delta_\infty(u, s)$  is finite for all  $s \in \mathbb{N}^*$  and  $u > 0$ . This implies that Assumption (iii) on the separation of the parameters holds.

As all kernels are equal in this setup, i.e  $\mathcal{K} := \mathcal{K}_T = \mathcal{K}_\infty$ , we have  $\mathcal{V}_T = 0$  and  $\rho_T = 0$ . Thus Assumption (v) on the closeness to the limit kernel and Assumption (iv) on the closeness of the metrics  $\mathfrak{d}_T$  and  $\mathfrak{d}_\infty$  come for free with  $\rho = 1$ .

Recall the definition of  $H_\infty^{(2)}$  from (45). We choose  $u_\infty = H_\infty^{(2)}(r_0, 1)$  (as  $\mathcal{K}_\infty$  is chosen equal to  $\mathcal{K}_T$ ) for some  $r_0 \in (0, 1/\sqrt{2})$ . We remark that in order to take  $u_\infty$  as large as possible and then have a separation distance as small as possible (since it is a decreasing function of  $u_\infty$ ), one could take  $r_0$  maximizing  $H_\infty^{(2)}$ .

Thus, the assumptions of Proposition 7.4 are satisfied, and we deduce that Assumption 6.1 holds.

We now concentrate on the hypotheses of Proposition 7.6. Assumptions (i)-(iii) clearly hold for the same reasons as Assumptions (i)-(iii) of Proposition 7.4. We take  $u'_\infty = u_\infty$ . Assumption (iv) comes for free in this setting.

Thus, the assumptions of Proposition 7.6 are satisfied, and we deduce, thanks to (47), that Assumption 6.2 holds.

In conclusion, we get that Assumptions 6.1 and 6.2 hold and thus Point (v) of Theorem 2.1 holds for any set of parameters  $\mathcal{Q}^*$  such that for all  $\theta \neq \theta' \in \mathcal{Q}^*$  the distance  $\mathfrak{d}(\theta, \theta')$  is larger than the separation distance:

$$(60) \quad \max(r_0, \delta_\infty(u_\infty, s)).$$

*Remark 9.2* (On the separation distance (60)). The separation distance (60) is a non-decreasing function of  $s$ . Similarly as in remark 8.2 where an upper bound on the minimal distance in the Gaussian spike deconvolution case is given, we can provide an upper bound for this distance. Let  $(i, j) \in \{0, 1\} \times \{0, 1, 2\}$ . By considering the definition of the kernel  $\mathcal{K}$  and the bound (59), we deduce that  $|\mathcal{K}^{[i,j]}(\theta, \theta')| \leq M f(\mathfrak{d}_\infty(\theta, \theta'))$  for all  $\theta, \theta' \in \Theta$ . We then obtain that for  $\vartheta = (\theta_1, \dots, \theta_s) \in \Theta_{\infty, \delta}^s$  with  $\delta > 0$ :

$$\max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s |\mathcal{K}^{[i,j]}(\theta_\ell, \theta_k)| \leq \psi_s(\delta) \quad \text{with} \quad \psi_s(\delta) = 2M \int_0^{s/2+1} f(\delta t) dt.$$

The function  $\psi_s$  is decreasing and one to one from  $\mathbb{R}_+$  to  $(0, M(s+2)]$ . Setting  $\psi_s^{-1}(u) = 0$  for  $u > M(s+2)$ , we deduce from (41) that for  $u > 0$ :

$$\delta_\infty(u, s) \leq \psi_s^{-1}(u).$$

We can bound the quantity above independently of  $s$ . Since the map  $\delta \mapsto \psi_s(\delta)$  is decreasing and the map  $s \mapsto \psi_s(\delta)$  is increasing with limit  $\psi_\infty(\delta) = 2M \int_0^{+\infty} f(\delta t) dt = M\pi/\delta$ , we deduce that for  $s \in \mathbb{N}^*$ :

$$\delta_\infty(u, s) \leq \psi_\infty^{-1}(u) = \frac{M\pi}{u}.$$

### 9.3. Prediction error

From Theorem 2.1, we deduce the subsequent following corollary. This demonstrates that by appropriately adjusting the penalization, the prediction error decreases to zero at the expected rate as the noise level tends to 0.

**Corollary 9.3.** *For all  $\theta \neq \theta'$  belonging to  $\mathcal{Q}^* = \{\theta_k^*, k \in S^*\}$ , with  $S^* = \text{Supp}(\beta^*)$  such that  $\mathfrak{d}(\theta, \theta')$  is larger than the separation given by (60), then, with some universal finite constants  $\mathcal{C}_0, \dots, \mathcal{C}_3 > 0$ , for any  $\tau > 1$  and a tuning parameter:*

$$(61) \quad \kappa \geq \mathcal{C}_1 \sigma \sqrt{\log(\tau)/T}, \quad \text{where} \quad \Delta_T = \frac{1}{T},$$

*we have the prediction error bound of the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in (4) given by:*

$$(62) \quad \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\| \leq \mathcal{C}_0 \sqrt{s} \kappa,$$

*with probability larger than  $1 - \mathcal{C}_2 \left( \frac{\log(M_T)}{\tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right)$ . Moreover, with the same probability, we have that  $\left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| \leq \mathcal{C}_3 \kappa s$  as well as the inequalities (11) of Theorem 2.5.*

*Remark 9.4.* We consider the particular case  $M_T = T^\gamma$  and  $\tau = T^{\gamma'}$ , with  $\gamma$  and  $\gamma'$  positive. We also take  $\kappa = \mathcal{C}_1 \sigma \sqrt{\gamma' \log(T)/T}$ . The prediction error is then given by:

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\| \leq \mathcal{C}_0 \mathcal{C}_1 \sqrt{s} \sigma \sqrt{\gamma' \frac{\log(T)}{T}},$$

*with probability larger than  $1 - \mathcal{C}_2 \left( \frac{\gamma}{\sqrt{\gamma'}} \frac{\sqrt{\log(T)}}{T^{\gamma'}} \vee \frac{1}{T^{\gamma'}} \right)$ .*

## References

- [1] ABRAMOWITZ, M. and STEGUN, I. A., eds. (1992). *Handbook of mathematical functions with formulas, graphs, and mathematical tables*. Dover Publications, Inc., New York Reprint of the 1972 edition. [MR1225604 42, 46](#)
- [2] ABSIL, P. A., MAHONY, R. and SEPULCHRE, R. (2008). *Optimization algorithms on matrix manifolds*. Princeton University Press, Princeton, NJ With a foreword by Paul Van Dooren. [https://doi.org/10.1515/9781400830244](#) [MR2364186 12](#)
- [3] ALIPRANTIS, C. D. and BORDER, K. C. (2006). *Infinite dimensional analysis*, Third ed. Springer, Berlin A hitchhiker's guide. [MR2378491 41](#)
- [4] ARENDT, W., BATTY, C. J. K., HIEBER, M. and NEUBRANDER, F. (2011). *Vector-valued Laplace transforms and Cauchy problems*, second ed. *Monographs in Mathematics* **96**. Birkhäuser/Springer Basel AG, Basel. [https://doi.org/10.1007/978-3-0348-0087-7](#) [MR2798103 41](#)
- [5] AZAÏS, J.-M., DE CASTRO, Y. and GAMBOA, F. (2015). Spike detection from inaccurate samplings. *Appl. Comput. Harmon. Anal.* **38** 177–195. [https://doi.org/10.1016/j.acha.2014.03.004](#) [MR3303671 4](#)



- [6] AZAÏS, J.-M. and WSCHEBOR, M. (2009). *Level sets and extrema of random processes and fields*. John Wiley & Sons, Inc., Hoboken, NJ. <https://doi.org/10.1002/9780470434642> MR2478201 42
- [7] BERNSTEIN, B. and FERNANDEZ-GRANDA, C. (2019). Deconvolution of point sources: a sampling theorem and robustness guarantees. *Comm. Pure Appl. Math.* **72** 1152–1230. <https://doi.org/10.1002/cpa.21805> MR3948555 5
- [8] BERNSTEIN, B., LIU, S., PAPADANIL, C. and FERNANDEZ-GRANDA, C. (2020). Sparse recovery beyond compressed sensing: separable nonlinear inverse problems. *IEEE Trans. Inform. Theory* **66** 5904–5926. <https://doi.org/10.1109/TIT.2020.2985015> MR4158652 5
- [9] BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.* **37** 1705–1732. <https://doi.org/10.1214/08-AOS620> MR2533469 4, 8, 15
- [10] BOYD, N., SCHIEBINGER, G. and RECHT, B. (2017). The alternating descent conditional gradient method for sparse inverse problems. *SIAM J. Optim.* **27** 616–639. <https://doi.org/10.1137/15M1035793> MR3634995 4
- [11] BOYER, C., CHAMBOLLE, A., DE CASTRO, Y., DUVAL, V., DE GOURNAY, F. and WEISS, P. (2019). On representer theorems and convex regularization. *SIAM J. Optim.* **29** 1260–1281. <https://doi.org/10.1137/18M1200750> MR3948246 2, 4
- [12] BOYER, C., DE CASTRO, Y. and SALMON, J. (2017). Adapting to unknown noise level in sparse deconvolution. *Inf. Inference* **6** 310–348. <https://doi.org/10.1093/imaia/iaw024> MR3764527 5, 6, 15, 20
- [13] BÜHLMANN, P. and VAN DE GEER, S. (2011). *Statistics for high-dimensional data*. Springer Series in Statistics. Springer, Heidelberg Methods, theory and applications. <https://doi.org/10.1007/978-3-642-20192-9> MR2807761 2
- [14] BUTUCEA, C., DELMAS, J.-F., DUTFOY, A. and HARDY, C. (2021). Modeling infra-red spectra: an algorithm for an automatic and simultaneous analysis. In *In Proceedings of the 31st European Safety and Reliability Conference* 3359–3366. 4
- [15] CANDÈS, E. and TAO, T. (2007). The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *Ann. Statist.* **35** 2313–2351. <https://doi.org/10.1214/009053606000001523> MR2382644 4
- [16] CANDÈS, E. J. and DAVENPORT, M. A. (2013). How well can we estimate a sparse vector? *Appl. Comput. Harmon. Anal.* **34** 317–323. <https://doi.org/10.1016/j.acha.2012.08.010> MR3008569 5
- [17] CANDÈS, E. J. and FERNANDEZ-GRANDA, C. (2013). Super-resolution from noisy data. *J. Fourier Anal. Appl.* **19** 1229–1254. <https://doi.org/10.1007/s00041-013-9292-3> MR3132912 4, 5, 14, 16
- [18] CANDÈS, E. J. and FERNANDEZ-GRANDA, C. (2014). Towards a mathematical theory of super-resolution. *Comm. Pure Appl. Math.* **67** 906–956. <https://doi.org/10.1002/cpa.21455> MR3193963 5, 14
- [19] CANDÈS, E. J. and PLAN, Y. (2011). A probabilistic and RIPless theory of compressed sensing. *IEEE Trans. Inform. Theory* **57** 7235–7254. <https://doi.org/10.1109/TIT.2011.2161794> MR2883653 14
- [20] CANDÈS, E. J. and TAO, T. (2005). Decoding by linear programming. *IEEE Trans. Inform. Theory* **51** 4203–4215. <https://doi.org/10.1109/TIT.2005.858979> MR2243152 14
- [21] CHIZAT, L. (2021). Sparse optimization on measures with over-parameterized gradient descent. *Mathematical Programming* 1–46. 4
- [22] DE CASTRO, Y., GADAT, S., MARTEAU, C. and MAUGIS-RABUSSEAU, C. (2021). SuperMix: sparse regularization for mixtures. *Ann. Statist.* **49** 1779–1809. <https://doi.org/10.1214/20-aos2022> MR4298881 5, 15
- [23] DE CASTRO, Y. and GAMBOA, F. (2012). Exact reconstruction using Beurling minimal extrapolation. *J. Math. Anal. Appl.* **395** 336–354. <https://doi.org/10.1016/j.jmaa.2012.05.011> MR2943626 4, 8
- [24] DENOYELLE, Q., DUVAL, V., PEYRÉ, G. and SOUBIES, E. (2020). The sliding Frank-Wolfe algorithm and its application to super-resolution microscopy. *Inverse Problems* **36** 014001, 42. <https://doi.org/10.1088/1361-6420/ab2a29> MR4040984 4, 11, 24
- [25] DONOHO, D. L., ELAD, M. and TEMLYAKOV, V. N. (2006). Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Trans. Inform. Theory* **52** 6–18. <https://doi.org/10.1109/TIT.2005.860430> MR2237332 4
- [26] DUVAL, V. and PEYRÉ, G. (2015). Exact support recovery for sparse spikes deconvolution. *Found. Comput. Math.* **15** 1315–1355. <https://doi.org/10.1007/s10208-014-9228-6> MR3394712 4, 5, 6, 14, 15, 24
- [27] DUVAL, V. and PEYRÉ, G. (2017). Sparse regularization on thin grids I: the Lasso. *Inverse Problems* **33** 055008, 29. <https://doi.org/10.1088/1361-6420/aa5e12> MR3628904 4
- [28] ELVIRA, C., GRIBONVAL, R., SOUSSEN, C. and HERZET, C. (2021). When does OMP achieve exact recovery with continuous dictionaries? *Appl. Comput. Harmon. Anal.* **51** 374–413. <https://doi.org/10.1016/j.acha.2020.12.002> MR4191904 4
- [29] EVANS, L. C. and GARIEPY, R. F. (2015). *Measure theory and fine properties of functions*, revised ed. Textbooks in Mathematics. CRC Press, Boca Raton, FL. MR3409135 42
- [30] FOUCART, S. and RAUHUT, H. (2013). *A mathematical introduction to compressive sensing*. Applied and Numerical Harmonic Analysis. Birkhäuser/Springer, New York. <https://doi.org/10.1007/978-0-8176-4948-7> MR3100033 17
- [31] GINÉ, E. and NICKL, R. (2016). *Mathematical foundations of infinite-dimensional statistical models*. Cambridge Series in Statistical and Probabilistic Mathematics, [40]. Cambridge University Press, New York. <https://doi.org/10.1017/CBO9781107337862> MR3588285 2
- [32] GOLBABAEE, M. and POON, C. (2022). An off-the-grid approach to multi-compartment magnetic resonance fingerprinting. *Inverse Problems* **38** Paper No. 085002, 31. <https://doi.org/10.1088/1361-6420/ac70da> MR4443913 4
- [33] GOLUB, G. H. and PEREYRA, V. (1973). The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate. *SIAM J. Numer. Anal.* **10** 413–432. <https://doi.org/10.1137/0710036> MR336980 4
- [34] KAUFMAN, L. (1975). A variable projection method for solving separable nonlinear least squares problems. *Nordisk Tidskr. Informationsbehandling (BIT)* **15** 49–57. <https://doi.org/10.1007/bf01932995> MR501738 4
- [35] KNEIP, A. and GASSER, T. (1988). Convergence and consistency results for self-modeling nonlinear regression. *Ann. Statist.* **16** 82–112. <https://doi.org/10.1214/aos/1176350692> MR924858 4
- [36] LANG, S. (1993). *Real and functional analysis*, third ed. Graduate Texts in Mathematics **142**. Springer-Verlag, New York. <https://doi.org/10.1007/978-1-4612-0897-6> MR1216137 41
- [37] LEE, J. M. (2018). *Introduction to Riemannian manifolds*. Graduate Texts in Mathematics **176**. Springer, Cham Second edition of [ MR1468735]. MR3887684 11
- [38] LOUNICI, K., PONTIL, M., VAN DE GEER, S. and TSYBAKOV, A. B. (2011). Oracle inequalities and optimal inference under group sparsity. *Ann. Statist.* **39** 2164–2204. <https://doi.org/10.1214/11-AOS896> MR2893865 8
- [39] MALLAT, S. (2009). *A wavelet tour of signal processing : the sparse way*, Third ed. Elsevier/Academic Press, Amsterdam With contributions from Gabriel Peyré. MR2479996 11

- [40] OLSHAUSEN, B. A. and FIELD, D. J. (1997). Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision research* **37** 3311–3325. [4](#)
- [41] POON, C., KERIVEN, N. and PEYRÉ, G. (2021). The geometry of off-the-grid compressed sensing. *Foundations of Computational Mathematics*. [5](#), [6](#), [11](#), [12](#), [15](#), [16](#), [17](#), [24](#), [34](#), [45](#)
- [42] RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2011). Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. *IEEE Trans. Inform. Theory* **57** 6976–6994. <https://doi.org/10.1109/TIT.2011.2165799> MR2882274 [5](#), [8](#)
- [43] SAKAI, T. (1996). *Riemannian geometry. Translations of Mathematical Monographs* **149**. American Mathematical Society, Providence, RI Translated from the 1992 Japanese original by the author. <https://doi.org/10.1090/mmono/149> MR1390760 [11](#)
- [44] SCHIEBINGER, G., ROBEVA, E. and RECHT, B. (2018). Superresolution without separation. *Inf. Inference* **7** 1–30. <https://doi.org/10.1093/imaiai/iax006> MR3801517 [5](#), [16](#)
- [45] TANG, G. (2015). Resolution limits for atomic decompositions via Markov-Bernstein type inequalities. In *2015 International Conference on Sampling Theory and Applications (SampTA)* 548–552. <https://doi.org/10.1109/SAMPTA.2015.7148951> [16](#)
- [46] TANG, G., BHASKAR, B. N. and RECHT, B. (2013). Sparse recovery over continuous dictionaries-just discretize. In *2013 Asilomar Conference on Signals, Systems and Computers* 1043–1047. IEEE. [4](#)
- [47] TANG, G., BHASKAR, B. N. and RECHT, B. (2015). Near minimax line spectral estimation. *IEEE Trans. Inform. Theory* **61** 499–512. <https://doi.org/10.1109/TIT.2014.2368122> MR3299978 [5](#), [8](#), [15](#), [20](#)
- [48] TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](#) [2](#), [4](#)
- [49] TROPP, J. A. (2004). Greed is good: algorithmic results for sparse approximation. *IEEE Trans. Inform. Theory* **50** 2231–2242. <https://doi.org/10.1109/TIT.2004.834793> MR2097044 [17](#)
- [50] TSYBAKOV, A. B. (2009). *Introduction to nonparametric estimation. Springer Series in Statistics*. Springer, New York Revised and extended from the 2004 French original, Translated by Vladimir Zaiats. <https://doi.org/10.1007/b13794> MR2724359 [2](#)
- [51] VAN DE GEER, S. (2016). *Estimation and testing under sparsity. Lecture Notes in Mathematics* **2159**. Springer, [Cham] Lecture notes from the 45th Probability Summer School held in Saint-Flour, 2015, École d’Été de Probabilités de Saint-Flour. [Saint-Flour Probability Summer School]. <https://doi.org/10.1007/978-3-319-32774-7> MR3526202 [4](#), [15](#)
- [52] VAN DE GEER, S. A. and BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the Lasso. *Electron. J. Stat.* **3** 1360–1392. <https://doi.org/10.1214/09-EJS506> MR2576316 [15](#)

## Appendix A: Proofs of Theorems 2.1 and 2.5

### A.1. Proof of Theorem 2.1

Let us bound the prediction error  $\hat{R}_T := \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_T$ . By definition (4) of  $\hat{\beta}$  and  $\hat{\vartheta}$  for the tuning parameter  $\kappa$ , we have:

$$\frac{1}{2} \left\| y - \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_T^2 + \kappa \|\hat{\beta}\|_{\ell_1} \leq \frac{1}{2} \left\| y - \beta^* \Phi_T(\vartheta^*) \right\|_T^2 + \kappa \|\beta^*\|_{\ell_1}.$$

We define the application  $\hat{\Upsilon}$  from  $H_T$  to  $\mathbb{R}$  by:

$$\hat{\Upsilon}(f) = \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*), f \right\rangle_T.$$

This gives, by rearranging terms and using the equation of the model  $y = \beta^* \Phi_T(\vartheta^*) + w_T$ , that:

$$(63) \quad \frac{1}{2} \hat{R}_T^2 \leq \hat{\Upsilon}(w_T) + \kappa \left( \|\beta^*\|_{\ell_1} - \|\hat{\beta}\|_{\ell_1} \right).$$

Next, we shall expand the two terms on the right hand side of (63) according to  $\hat{\beta}_\ell$  close to some  $\beta_k^*$  or not. In the rest of the proof, we fix  $r > 0$  so that Assumptions 6.1 and 6.2, are verified by  $\mathcal{Q}^*$ . In particular, for all  $k \neq k'$  in  $S^* = \{k'' \in \{1, \dots, K\}, \beta_{k''}^* \neq 0\}$  we have  $\mathfrak{d}_T(\theta_k^*, \theta_{k'}^*) > 2r$ .

Recall the definitions given in Section 2 of the sets of indices  $\hat{S}$ ,  $\tilde{S}_k(r)$  and  $\tilde{S}(r)$  for  $k \in S^*$ . Since the closed balls  $\mathcal{B}_T(\theta_k^*, r)$  with  $k \in S^*$  are pairwise disjoint, the sets  $\tilde{S}_k(r)$ , for  $k \in S^*$ , are also pairwise disjoint and one can write the following decomposition:

$$\hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) = \sum_{k=1}^K \hat{\beta}_k \phi_T(\hat{\theta}_k) - \sum_{k \in S^*} \beta_k^* \phi_T(\theta_k^*) = \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \phi_T(\hat{\theta}_\ell) + \sum_{k \in \tilde{S}(r)^c} \hat{\beta}_k \phi_T(\hat{\theta}_k) - \sum_{k \in S^*} \beta_k^* \phi_T(\theta_k^*).$$

This decomposition groups the elements of the predicted mixture according to the proximity of the estimated parameter  $\hat{\theta}_\ell$  to a true underlying parameter  $\theta_k^*$  to be estimated. We use a Taylor-type expansion with the Riemannian metric  $\mathfrak{d}_T$  for the function  $\phi_T(\theta)$  around the elements of  $\mathcal{Q}^*$ . By Assumption 3.1, the function  $\phi_T$  is twice continuously differentiable

with respect to the variable  $\theta$  and the function  $g_T$  defined in (14) is positive on  $\Theta_T$  and of class  $\mathcal{C}^1$  by Assumption 3.2. We set in this section  $\tilde{D}_{i;T}[\phi_T] = \phi_T^{[i]}$  for  $i = 0, 1, 2$ . According to Lemma 4.2, we have for any  $\theta_k^*$  and  $\hat{\theta}_\ell$  in  $\Theta_T$ :

$$\phi_T(\hat{\theta}_\ell) = \phi_T(\theta_k^*) + \text{sign}(\hat{\theta}_\ell - \theta_k^*) \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*) \phi_T^{[1]}(\theta_k^*) + \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*)^2 \int_0^1 (1-s) \phi_T^{[2]}(\gamma_s^{(k\ell)}) \text{d}s,$$

where  $\gamma^{(k\ell)}$  is a distance realizing geodesic path belonging to  $\Theta_T$  such that  $\gamma_0^{(k\ell)} = \theta_k^*$ ,  $\gamma_1^{(k\ell)} = \hat{\theta}_\ell$  and  $\mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*) = \mathcal{L}_T(\gamma^{(k\ell)})$ . Hence we obtain:

$$(64) \quad \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) = \sum_{k \in S^*} I_{0,k}(r) \phi_T(\theta_k^*) + \sum_{k \in S^*} I_{1,k}(r) \phi_T^{[1]}(\theta_k^*) + \sum_{k \in \tilde{S}(r)^c} \hat{\beta}_k \phi_T(\hat{\theta}_k) \\ + \sum_{k \in S^*} \left( \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*)^2 \int_0^1 (1-s) \phi_T^{[2]}(\gamma_s^{(k\ell)}) \text{d}s \right),$$

with

$$I_{0,k}(r) = \left( \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \right) - \beta_k^* \quad \text{and} \quad I_{1,k}(r) = \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \text{sign}(\hat{\theta}_\ell - \theta_k^*) \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*).$$

Let us introduce some notations in order to bound the different terms of the expansion above:

$$(65) \quad I_0(r) = \sum_{k \in S^*} |I_{0,k}(r)| \quad \text{and} \quad I_1(r) = \sum_{k \in S^*} |I_{1,k}(r)|,$$

$$(66) \quad I_{2,k}(r) = \sum_{\ell \in \tilde{S}_k(r)} \left| \hat{\beta}_\ell \right| \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*)^2 \quad \text{and} \quad I_2(r) = \sum_{k \in S^*} I_{2,k}(r),$$

$$(67) \quad I_3(r) = \sum_{\ell \in \tilde{S}(r)^c} \left| \hat{\beta}_\ell \right| = \left\| \hat{\beta}_{\tilde{S}(r)^c} \right\|_{\ell_1},$$

and we omit the dependence in  $r$  when there is no ambiguity.

We bound the difference  $\|\beta^*\|_{\ell_1} - \|\hat{\beta}\|_{\ell_1}$  by noticing that:

$$(68) \quad \|\beta^*\|_{\ell_1} - \|\hat{\beta}\|_{\ell_1} = \sum_{k \in S^*} \left( |\beta_k^*| - \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \right) - \sum_{k \in \tilde{S}(r)^c} |\hat{\beta}_k| \leq \sum_{k \in S^*} \left| \beta_k^* - \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \right| = I_0.$$

In the next lemma, we give an upper bound of  $I_0$ . Recall the constants  $C'_N$  and  $C_F$  from Assumption 6.1.

**Lemma A.1.** *Under the assumptions of Theorem 2.1 and with the element  $p_1 \in H_T$  from Assumption 6.1 associated to the function  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  defined by:*

$$v(\theta_k^*) = \text{sign}(I_{0,k}) \quad \text{for all } k \in S^*,$$

we get that:

$$(69) \quad I_0 \leq C'_N I_2 + (1 - C_F) I_3 + |\hat{\Upsilon}(p_1)|.$$

**Proof.** Let  $v \in \{-1, 1\}^s$  with entries  $v_k = v(\theta_k^*)$  so that:

$$I_0 = \sum_{k \in S^*} |I_{0,k}| = \sum_{k \in S^*} v_k I_{0,k} = \sum_{k \in S^*} v_k \left( \left( \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \right) - \beta_k^* \right).$$

Let  $p_1$  be an element of  $H_T$  from Assumption 6.1 associated to the application  $v$  such that properties (i)-(iv) therein hold. By adding and subtracting  $\sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \langle \phi_T(\hat{\theta}_\ell), p_1 \rangle_T$  to  $I_0$  and using the property (ii) satisfied by the element

$p_1$ , that is,  $\langle \phi_T(\theta_k^*), p_1 \rangle_T = v_k$  for all  $k \in S^*$ , we obtain:

$$I_0 = \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell \left( v_k - \left\langle \phi_T(\hat{\theta}_\ell), p_1 \right\rangle_T \right) + \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*), p_1 \right\rangle_T - \sum_{\ell \in \tilde{S}(r)^c} \hat{\beta}_\ell \left\langle \phi_T(\hat{\theta}_\ell), p_1 \right\rangle_T.$$

We deduce that:

$$I_0 \leq \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \left| v_k - \left\langle \phi_T(\hat{\theta}_\ell), p_1 \right\rangle_T \right| + |\hat{\Upsilon}(p_1)| + \sum_{\ell \in \tilde{S}(r)^c} |\hat{\beta}_\ell| \left| \left\langle \phi_T(\hat{\theta}_\ell), p_1 \right\rangle_T \right|.$$

Notice that for  $\ell \in \tilde{S}(r)^c$ ,  $\hat{\theta}_\ell \notin \bigcup_{k \in S^*} \mathcal{B}_T(\theta_k^*, r)$ . Then, by using the properties (ii) and (iii) from Assumption 6.1, we get that (69) holds with the constants  $C'_N$  and  $C_F$  from Assumption 6.1.  $\square$

In the next lemma, we give an upper bound of  $I_1$ . Recall the constants  $c_N$  and  $c_F$  from Assumption 6.2.

**Lemma A.2.** *Under the assumptions of Theorem 2.1 and with the element  $q_0 \in H_T$  from Assumption 6.2 associated to the function  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  defined by:*

$$v(\theta_k^*) = \text{sign}(I_{1,k}) \quad \text{for all } k \in S^*,$$

we get that:

$$(70) \quad I_1 \leq c_N I_2 + c_F I_3 + |\hat{\Upsilon}(q_0)|.$$

**Proof.** Let  $v \in \{-1, 1\}^s$  with entries  $v_k = v(\theta_k^*)$  so that:

$$I_1 = \sum_{k \in S^*} |I_{1,k}| = \sum_{k \in S^*} v_k I_{1,k} = \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} \hat{\beta}_\ell v_k \text{sign}(\hat{\theta}_\ell - \theta_k^*) \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*).$$

Let  $q_0 \in H_T$  from Assumption 6.2 associated to the application  $v$  such that properties (i)-(iii) therein hold. By adding and subtracting  $\sum_{\ell \in \tilde{S}(r)} \hat{\beta}_\ell \left\langle \phi_T(\hat{\theta}_\ell), q_0 \right\rangle_T = \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}), q_0 \right\rangle_T - \sum_{\ell \in \tilde{S}(r)^c} \hat{\beta}_\ell \left\langle \phi_T(\hat{\theta}_\ell), q_0 \right\rangle_T$  to  $I_1$  and using the triangle inequality, we obtain:

$$\begin{aligned} I_1 \leq \sum_{k \in S^*} \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \left| v_k \text{sign}(\hat{\theta}_\ell - \theta_k^*) \mathfrak{d}_T(\hat{\theta}_\ell, \theta_k^*) - \left\langle \phi_T(\hat{\theta}_\ell), q_0 \right\rangle_T \right| \\ + \sum_{\ell \in \tilde{S}(r)^c} |\hat{\beta}_\ell| \left| \left\langle \phi_T(\hat{\theta}_\ell), q_0 \right\rangle_T \right| + \left| \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}), q_0 \right\rangle_T \right|. \end{aligned}$$

The property (i) of Assumption 6.2 gives that  $\langle \phi_T(\theta_k^*), q_0 \rangle_T = 0$  for all  $k \in S^*$ . This implies that  $\langle \beta^* \Phi_T(\vartheta^*), q_0 \rangle_T = 0$ . Then, by using the definition of  $I_2$  and  $I_3$  from (66)-(67) and the properties (i) and (ii) of Assumption 6.2, we obtain:

$$I_1 \leq c_N I_2 + c_F I_3 + \left| \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}), q_0 \right\rangle_T \right| = c_N I_2 + c_F I_3 + |\hat{\Upsilon}(q_0)|,$$

with the constants  $c_N$  and  $c_F$  from Assumption 6.2.  $\square$

We consider the following suprema of Gaussian processes for  $i = 0, 1, 2$ :

$$M_i = \sup_{\theta \in \Theta_T} \left| \left\langle w_T, \phi_T^{[i]}(\theta) \right\rangle_T \right|.$$

By using the expansion (64) and the bounds (70) and (69) for the second inequality, we obtain:

$$(71) \quad |\hat{\Upsilon}(w_T)| \leq (I_0 + I_3) M_0 + I_1 M_1 + I_2 2^{-1} M_2$$

$$(72) \quad \leq (C'_N I_2 + (2 - C_F) I_3 + |\hat{\Upsilon}(p_1)|) M_0 + (c_N I_2 + c_F I_3 + |\hat{\Upsilon}(q_0)|) M_1 + I_2 2^{-1} M_2.$$

At this point, one needs to bound  $I_2$  and  $I_3$ . In order to do so, we will bound from above and from below the Bregman divergence  $D_B$  defined by:

$$(73) \quad D_B = \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} - \hat{\Upsilon}(p_0),$$

where  $p_0$  is the element of  $H_T$  given by the Assumption 6.1 associated to the application  $v : \mathcal{Q}^* \rightarrow \{-1, 1\}$  given by:

$$(74) \quad v(\theta_k^*) = \text{sign}(\beta_k^*) \quad \text{for all } k \in S^*.$$

The next lemma gives a lower bound of the Bregman divergence.

**Lemma A.3.** *Under the assumptions of Theorem 2.1 and with the constants  $C_N$  and  $C_F$  of Assumption 6.1, we get that:*

$$(75) \quad D_B \geq C_N I_2 + C_F I_3.$$

**Proof.** By definition (73) of  $D_B$  we have:

$$D_B = \sum_{k \in \hat{S}} |\hat{\beta}_k| - \hat{\beta}_k \left\langle \phi_T(\hat{\theta}_k), p_0 \right\rangle_T - \left( \sum_{k \in S^*} |\beta_k^*| - \beta_k^* \left\langle \phi_T(\theta_k^*), p_0 \right\rangle_T \right).$$

By using the interpolating properties of the element  $p_0$  of  $H_T$  from Assumption 6.1 associated to the function  $v$  defined in (74), we have  $\sum_{k \in S^*} |\beta_k^*| - \beta_k^* \left\langle \phi_T(\theta_k^*), p_0 \right\rangle_T = 0$ . Hence, we deduce that:

$$\begin{aligned} D_B &= \sum_{k \in \hat{S}} |\hat{\beta}_k| - \hat{\beta}_k \left\langle \phi_T(\hat{\theta}_k), p_0 \right\rangle_T \\ &\geq \sum_{k \in \hat{S}} |\hat{\beta}_k| - |\hat{\beta}_k| \left| \left\langle \phi_T(\hat{\theta}_k), p_0 \right\rangle_T \right| \\ &= \sum_{\ell \in \hat{S}(r)} |\hat{\beta}_\ell| \left( 1 - \left| \left\langle \phi_T(\hat{\theta}_\ell), p_0 \right\rangle_T \right| \right) + \sum_{k \in \hat{S}(r)^c} |\hat{\beta}_k| \left( 1 - \left| \left\langle \phi_T(\hat{\theta}_k), p_0 \right\rangle_T \right| \right). \end{aligned}$$

Thanks to properties (i) and (iii) of Assumption 6.1 and the definitions (66) and (67) of  $I_2$  and  $I_3$ , we obtain:

$$D_B \geq \sum_{k \in S^*} \sum_{\ell \in \hat{S}_k(r)} C_N |\hat{\beta}_\ell| \mathfrak{D}_T(\hat{\theta}_\ell, \theta_k^*)^2 + \sum_{k \in \hat{S}(r)^c} C_F |\hat{\beta}_k| = C_N I_2 + C_F I_3,$$

where the constants  $C_N$  and  $C_F$  are that of Assumption 6.1. □

We now give an upper bound of the Bregman divergence.

**Lemma A.4.** *Under the assumptions of Theorem 2.1, we have:*

$$(76) \quad \kappa D_B \leq I_2 (C'_N M_0 + c_N M_1 + 2^{-1} M_2) + I_3 ((2 - C_F) M_0 + c_F M_1) + |\hat{\Upsilon}(p_1)| M_0 + |\hat{\Upsilon}(q_0)| M_1 + \kappa |\hat{\Upsilon}(p_0)|.$$

**Proof.** Recall that  $\mathcal{Q}^* \subset \Theta_T$ . We deduce from (63) that:

$$(77) \quad \kappa (\|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1}) \leq \hat{\Upsilon}(w_T) - \frac{1}{2} \left\| \beta^* \Phi_T(\vartheta^*) - \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_T^2 \leq \hat{\Upsilon}(w_T).$$

Using (73), we obtain:

$$\kappa D_B \leq |\hat{\Upsilon}(w_T)| + \kappa |\hat{\Upsilon}(p_0)|.$$

Then, use (72) to get (76). □

By combining the upper and lower bounds (75) and (76), we deduce that:

$$\begin{aligned} (78) \quad I_2 \left( C_N - \frac{1}{\kappa} (C'_N M_0 + c_N M_1 + 2^{-1} M_2) \right) + I_3 \left( C_F - \frac{1}{\kappa} ((2 - C_F) M_0 + c_F M_1) \right) \\ \leq \frac{1}{\kappa} |\hat{\Upsilon}(p_1)| M_0 + \frac{1}{\kappa} |\hat{\Upsilon}(q_0)| M_1 + |\hat{\Upsilon}(p_0)|. \end{aligned}$$

We define the events:

$$(79) \quad \mathcal{A}_i = \{M_i \leq \mathcal{C} \kappa\}, \quad \text{for } i \in \{0, 1, 2\} \quad \text{and} \quad \mathcal{A} = \mathcal{A}_0 \cap \mathcal{A}_1 \cap \mathcal{A}_2,$$

where:

$$\mathcal{C} = \frac{C_F}{2(2 - C_F + c_F)} \wedge \frac{C_N}{2(C'_N + c_N + 2^{-1})}.$$

(We shall prove in (88) that the event  $\mathcal{A}$  occurs with high probability.) We get from Inequality (78), that on the event  $\mathcal{A}$ :

$$(80) \quad C_N I_2 + C_F I_3 \leq 2\mathcal{C}' \left( |\hat{\Upsilon}(p_1)| + |\hat{\Upsilon}(q_0)| + |\hat{\Upsilon}(p_0)| \right) \quad \text{with} \quad \mathcal{C}' = \mathcal{C} \vee 1.$$

By reinjecting (68), (72), (69) and (70) in (63) one gets:

$$\begin{aligned} \frac{1}{2} \hat{R}_T^2 &\leq I_2(C'_N M_0 + c_N M_1 + 2^{-1} M_2 + \kappa C'_N) + I_3((2 - C_F)M_0 + c_F M_1 + \kappa(1 - C_F)) \\ &\quad + |\hat{\Upsilon}(p_1)|(M_0 + \kappa) + |\hat{\Upsilon}(q_0)|M_1. \end{aligned}$$

Using (80), we obtain an upper bound for the prediction error on the event  $\mathcal{A}$ :

$$(81) \quad \hat{R}_T^2 \leq C \kappa (|\hat{\Upsilon}(p_0)| + |\hat{\Upsilon}(p_1)| + |\hat{\Upsilon}(q_0)|),$$

with

$$C = 4\mathcal{C}' \left( 1 + \frac{\mathcal{C}'}{C_N} (2C'_N + c_N + 1) + \frac{\mathcal{C}'}{C_F} (3 - 2C_F + c_F) \right).$$

Using the Cauchy-Schwarz inequality and the definition of  $\hat{\Upsilon}$ , we get that for  $f \in H_T$ :

$$(82) \quad |\hat{\Upsilon}(f)| \leq \hat{R}_T \|f\|_T.$$

Using Assumption 6.1 (iv) for  $p_0$  and  $p_1$ , and Assumption 6.2 (iii) for  $q_0$ , we get:

$$(83) \quad \|p_0\|_T \leq C_B \sqrt{s}, \quad \|p_1\|_T \leq C_B \sqrt{s} \quad \text{and} \quad \|q_0\|_T \leq c_B \sqrt{s}.$$

Plugging this in (81), we get that on the event  $\mathcal{A}$ :

$$(84) \quad \hat{R}_T^2 \leq \mathcal{C}_0 \kappa \hat{R}_T \sqrt{s} \quad \text{with} \quad \mathcal{C}_0 = (c_B + 2C_B)C.$$

This gives (7).

The proof of (8) is postponed to Section A.2 and will be easily deduced from the first and third inequalities in (11).

To complete the proof of Theorem 2.1 we shall give a lower bound for the probability of the event  $\mathcal{A}$  defined in (79). For  $i = 0, 1, 2$  and  $\theta \in \Theta$ , set  $X_i(\theta) = \left\langle w_T, \phi_T^{[i]}(\theta) \right\rangle_T$  a real centered Gaussian process with continuously differentiable sample paths, so that its supremum is  $M_i = \sup_{\Theta_T} |X_i|$ .

We first consider  $i = 0$ . We have, thanks to (31) and (28) for the second part:

$$\|\phi_T(\theta)\|_T^2 = 1 \quad \text{and} \quad \left\| \phi_T^{[1]}(\theta) \right\|_T^2 = \mathcal{K}_T^{[1,1]}(\theta, \theta) = 1.$$

Recall Assumption 1.1 on the noise  $w_T$  holds. We deduce from Lemma C.1 with  $C_1 = C_2 = 1$  that:

$$(85) \quad \mathbb{P}(\mathcal{A}_0^c) = \mathbb{P}\left(\sup_{\Theta_T} |X_0| > \mathcal{C} \kappa\right) \leq c_0 \left( \sigma \frac{|\Theta_T| \mathfrak{d}_T \sqrt{\Delta_T}}{\mathcal{C} \kappa} \vee 1 \right) e^{-(\mathcal{C} \kappa)^2 / (4\sigma^2 \Delta_T)},$$



where  $|\Theta_T|_{\mathfrak{d}_T}$  denotes the diameter of the set  $\Theta_T$  with respect to the metric  $\mathfrak{d}_T$  and  $c_0 = 3$ .

We consider  $i = 1$ . Thanks to (31), we get:

$$\left\| \phi_T^{[1]}(\theta) \right\|_T^2 = 1 \quad \text{and} \quad \left\| \tilde{D}_{1;T}[\phi_T^{[1]}](\theta) \right\|_T^2 = \left\| \phi_T^{[2]}(\theta) \right\|_T^2 = \mathcal{K}_T^{[2,2]}(\theta, \theta).$$

Recall  $L_{2,2}$  and  $\mathcal{V}_T$  are defined in (34) and (37). Since Assumptions 5.1 and 5.2 hold, we get that for  $\theta \in \Theta_T$ :

$$\mathcal{K}_T^{[2,2]}(\theta, \theta) \leq L_{2,2} + \mathcal{V}_T \leq 2L_{2,2}.$$

We deduce from Lemma C.1 with  $C_1 = 1$  and  $C_2 = \sqrt{2L_{2,2}}$  and taking  $c_1 = 2\sqrt{2L_{2,2}} + 1$ , that:

$$(86) \quad \mathbb{P}(\mathcal{A}_1^c) = \mathbb{P}\left(\sup_{\Theta_T} |X_1| > \mathcal{C}\kappa\right) \leq c_1 \left( \sigma \frac{|\Theta_T|_{\mathfrak{d}_T} \sqrt{\Delta_T}}{\mathcal{C}\kappa} \vee 1 \right) e^{-(\mathcal{C}\kappa)^2 / (4\sigma^2 \Delta_T)}.$$

We consider  $i = 2$ . Thanks to (31), we get:

$$\left\| \phi_T^{[2]}(\theta) \right\|_T^2 = \mathcal{K}_T^{[2,2]}(\theta, \theta) \quad \text{and} \quad \left\| \tilde{D}_{1;T}[\phi_T^{[2]}](\theta) \right\|_T^2 = \left\| \phi_T^{[3]}(\theta) \right\|_T^2 = \mathcal{K}_T^{[3,3]}(\theta, \theta).$$

Recall the definition of the function  $h_\infty$  given in (33) and the constants  $L_{2,2}$ ,  $L_3$ ,  $\mathcal{V}_T$  defined in (34) and (37). Using also Assumption 5.2 so that  $\mathcal{V}_T \leq L_{2,2} \wedge L_3$ , we get that for all  $\theta \in \Theta_T$ :

$$\mathcal{K}_T^{[2,2]}(\theta, \theta) \leq L_{2,2} + \mathcal{V}_T \leq 2L_{2,2} \quad \text{and} \quad \mathcal{K}_T^{[3,3]}(\theta, \theta) \leq L_3 + \mathcal{V}_T \leq 2L_3.$$

We deduce from Lemma C.1 with  $C_1 = \sqrt{2L_{2,2}}$  and  $C_2 = \sqrt{2L_3}$  and taking  $c_2 = 2\sqrt{2L_3} + 1$ , that:

$$(87) \quad \mathbb{P}(\mathcal{A}_2^c) = \mathbb{P}\left(\sup_{\Theta_T} |X_2| > \mathcal{C}\kappa\right) \leq c_2 \left( \sigma \frac{|\Theta_T|_{\mathfrak{d}_T} \sqrt{\Delta_T}}{\mathcal{C}\kappa} \vee 1 \right) e^{-(\mathcal{C}\kappa)^2 / (8\sigma^2 \Delta_T L_{2,2})}.$$

Since  $\mathcal{A} = \mathcal{A}_0 \cap \mathcal{A}_1 \cap \mathcal{A}_2$ , we deduce from (85), (86) and (87) that:

$$\mathbb{P}(\mathcal{A}^c) = \mathbb{P}(\mathcal{A}_0^c \cup \mathcal{A}_1^c \cup \mathcal{A}_2^c) \leq \mathcal{C}'_2 \left( \sigma \frac{|\Theta_T|_{\mathfrak{d}_T} \sqrt{\Delta_T}}{\mathcal{C}\kappa} \vee 1 \right) e^{-\kappa^2 / (\mathcal{C}'_2 \sigma^2 \Delta_T)},$$

with the finite positive constants:

$$\mathcal{C}_1 = \frac{2}{\mathcal{C}} \left( 1 \vee \sqrt{2L_{2,2}} \right) \quad \text{and} \quad \mathcal{C}'_2 = c_0 + c_1 + c_2.$$

By taking  $\kappa \geq \mathcal{C}_1 \sigma \sqrt{\Delta_T \log \tau}$ , for any positive constant  $\tau > 1$ , we get:

$$(88) \quad \mathbb{P}(\mathcal{A}_0^c \cup \mathcal{A}_1^c \cup \mathcal{A}_2^c) \leq \mathcal{C}_2 \left( \frac{|\Theta_T|_{\mathfrak{d}_T}}{\tau \sqrt{\log \tau}} \vee \frac{1}{\tau} \right) \quad \text{with} \quad \mathcal{C}_2 = \mathcal{C}'_2 \left( \frac{1}{\mathcal{C}\mathcal{C}_1} \vee 1 \right).$$

This completes the proof of the theorem.

#### A.2. Proof of Theorem 2.5 and of Equation (8)

We keep notations from Section A.1. Recall that Assumptions (i)-(v) of Theorem 2.1 are in force. We shall first provide an upper bound of  $I_i$  for  $i = 0, 1, 2, 3$ . We deduce from (82), (83) and (84), that, on the event  $\mathcal{A}$ :

$$|\hat{\Upsilon}(p_0)| \leq \mathcal{C}_0 C_B \kappa s, \quad |\hat{\Upsilon}(p_1)| \leq \mathcal{C}_0 C_B \kappa s \quad \text{and} \quad |\hat{\Upsilon}(q_0)| \leq \mathcal{C}_0 C_B \kappa s.$$

Then, we obtain from (80) that, on the event  $\mathcal{A}$ :

$$(89) \quad I_3 \leq \mathcal{C}_5 \kappa s \quad \text{and} \quad I_2 \leq \mathcal{C}_6 \kappa s \quad \text{with} \quad \mathcal{C}_5 = 2 \frac{\mathcal{C}'}{C_F} \mathcal{C}_0 (C_B + 2C_B) \quad \text{and} \quad \mathcal{C}_6 = \frac{C_F}{C_N} \mathcal{C}_5.$$

This gives the third inequality in (11), as well as Inequality (12). We also deduce from (69) that, on the event  $\mathcal{A}$ :

$$(90) \quad I_0 \leq \mathcal{C}_4 \kappa s \quad \text{with} \quad \mathcal{C}_4 = C'_N \mathcal{C}_6 + (1 - C_F) \mathcal{C}_5 + \mathcal{C}_0 C_B.$$

This gives the second inequality in (11).

We now establish the first inequality in (11). We deduce from (63) that:

$$(91) \quad \kappa(\|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1}) \leq \hat{\Upsilon}(w_T).$$

Then, using the bounds (90) and (89) on  $I_0$ ,  $I_2$  and  $I_3$ , we deduce from (71) and (70) that, on the event  $\mathcal{A}$ :

$$(92) \quad |\hat{\Upsilon}(w_T)| \leq \mathcal{C}_7 s \kappa^2 \quad \text{with} \quad \mathcal{C}_7 = \mathcal{C}(\mathcal{C}_4 + \mathcal{C}_5(1 + c_F) + \mathcal{C}_6(1 + c_N) + \mathcal{C}_0 c_B).$$

Thus, (91) and (92) imply that, on the event  $\mathcal{A}$ :

$$(93) \quad \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \leq \mathcal{C}_7 s \kappa.$$

Then, use (68) and (90) to deduce that, on the event  $\mathcal{A}$ :

$$|\|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1}| \leq (\mathcal{C}_4 \vee \mathcal{C}_7) s \kappa.$$

This proves (8) (we shall take  $\mathcal{C}_3 = \mathcal{C}_7 + 2\mathcal{C}_4$ , see below). Let  $\mathcal{I}^+$  (resp.  $\mathcal{I}^-$ ) be the set of indices  $k \in S^*$  such that the quantity  $\left(\sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| - |\beta_k^*|\right)$  is non negative (resp. negative). We have the following decomposition:

$$(94) \quad \begin{aligned} \sum_{k \in S^*} \left| \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| - |\beta_k^*| \right| &= \sum_{k \in \mathcal{I}^+} \left( \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| - |\beta_k^*| \right) + \sum_{k \in \mathcal{I}^-} \left( |\beta_k^*| - \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \right) \\ &\leq \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} + 2 \sum_{k \in \mathcal{I}^-} \left( |\beta_k^*| - \sum_{\ell \in \tilde{S}_k(r)} |\hat{\beta}_\ell| \right) \\ &\leq \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} + 2I_0. \end{aligned}$$

Then, use (90) and (93) to obtain the first inequality (11) with  $\mathcal{C}_3 = \mathcal{C}_7 + 2\mathcal{C}_4$ . This ends the proof of Theorem 2.5.

## Appendix B: Construction of certificate functions

### B.1. Proof of Proposition 7.4 (Construction of an interpolating certificate)

This section is devoted to the proof of Proposition 7.4. We closely follow the proof of [41] taking into account the approximation of the kernel  $\mathcal{K}_T$  by the kernel  $\mathcal{K}_\infty$ , which is measured through the quantity  $\mathcal{V}_T$  defined in (37).

Let  $T \in \mathbb{N}$  and  $s \in \mathbb{N}^*$ . Recall Assumptions 3.2 (and thus 3.1 on the regularity of  $\varphi_T$ ) and 5.1 on the regularity of the asymptotic kernel  $\mathcal{K}_\infty$  are in force. Let  $\rho \geq 1$ , let  $r \in (0, 1/\sqrt{2L_{0,2}})$  and  $u_\infty \in (0, H_\infty^{(2)}(r, \rho))$  such that (ii), (iii), (iv) and (v) of Proposition 7.4 hold. We denote by  $\|\cdot\|_{\text{op}}$  the operator norm associated to the  $\ell_\infty$  norm on  $\mathbb{R}^s$ .

By assumption  $\delta_\infty(u_\infty, s)$  is finite. Let  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta_{T, 2\rho_T \delta_\infty(u_\infty, s)}^s$ . We note  $\mathcal{Q}^* = \{\theta_i^*, 1 \leq i \leq s\}$  the set of parameters of cardinal  $s$ . By Lemma 7.3, we have:

$$\Theta_{T, \rho_T \delta_\infty(u_\infty, s)}^s \subseteq \Theta_{T, \delta_T(u_T(s), s)}^s \quad \text{where} \quad u_T(s) = u_\infty + (s-1)\mathcal{V}_T.$$

Hence we have:

$$(95) \quad \vartheta^* \in \Theta_{T, \delta_T(u_T(s), s)}^s.$$

Set

$$(96) \quad \Gamma^{[i,j]} = \mathcal{K}_T^{[i,j]}(\vartheta^*) \quad \text{and} \quad \Gamma = \begin{pmatrix} \Gamma^{[0,0]} & \Gamma^{[1,0]\top} \\ \Gamma^{[1,0]} & \Gamma^{[1,1]} \end{pmatrix}.$$

We deduce from (43) and (95) that:

$$(97) \quad \left\| I - \Gamma^{[0,0]} \right\|_{\text{op}} \leq u_T(s), \quad \left\| I - \Gamma^{[1,1]} \right\|_{\text{op}} \leq u_T(s), \quad \left\| \Gamma^{[1,0]} \right\|_{\text{op}} \leq u_T(s) \quad \text{and} \quad \left\| \Gamma^{[1,0]^\top} \right\|_{\text{op}} \leq u_T(s).$$

For simplicity, for an expression  $A$  we write  $A_T$  for  $A_{\mathcal{K}_T}$ . Using this convention, recall the definition of the derivative operator  $\tilde{D}_{i;T}$  and write  $\phi_T^{[1]}$  for  $\tilde{D}_{1;T}[\phi_T]$ .

Let  $\alpha = (\alpha_1, \dots, \alpha_s)^\top$  and  $\xi = (\xi_1, \dots, \xi_s)^\top$  be elements of  $\mathbb{R}^s$ . Let  $p_{\alpha,\xi}$  be an element of  $H_T$  defined by:

$$(98) \quad p_{\alpha,\xi} = \sum_{k=1}^s \alpha_k \phi_T(\theta_k^*) + \sum_{k=1}^s \xi_k \phi_T^{[1]}(\theta_k^*),$$

and, using (31) in Lemma 4.3, set the interpolating real-valued function  $\eta_{\alpha,\xi}$  defined on  $\Theta$  by:

$$(99) \quad \eta_{\alpha,\xi}(\theta) = \langle \phi_T(\theta), p_{\alpha,\xi} \rangle_T = \sum_{k=1}^s \alpha_k \mathcal{K}_T(\theta, \theta_k^*) + \sum_{k=1}^s \xi_k \mathcal{K}_T^{[0,1]}(\theta, \theta_k^*).$$

By Assumption 3.2 on the regularity of  $\varphi_T$  and the positivity of  $g_T$  and Lemma 4.3, we get that the function  $\eta_{\alpha,\xi}$  is of class  $\mathcal{C}^3$  on  $\Theta$ , and using (23), we get that:

$$(100) \quad \eta_{\alpha,\xi}^{[1]} := \tilde{D}_{1;T}[\eta_{\alpha,\xi}](\theta) = \sum_{k=1}^s \alpha_k \mathcal{K}_T^{[1,0]}(\theta, \theta_k^*) + \sum_{k=1}^s \xi_k \mathcal{K}_T^{[1,1]}(\theta, \theta_k^*).$$

We give a preliminary technical lemma.

**Lemma B.1.** *Let  $v = (v_1, \dots, v_s)^\top \in \{-1, 1\}^s$  be a sign vector. Assume that (97) holds with  $u_T(s) < 1/2$ . Under Assumption 3.2, there exist unique  $\alpha, \xi \in \mathbb{R}^s$  such that:*

$$(101) \quad \eta_{\alpha,\xi}(\theta_k^*) = v_k \in \{-1, 1\} \quad \text{and} \quad \eta_{\alpha,\xi}^{[1]}(\theta_k^*) = 0 \quad \text{for} \quad 1 \leq k \leq s.$$

Furthermore, we have:

$$(102) \quad \|\alpha\|_{\ell_\infty} \leq \frac{1 - u_T(s)}{1 - 2u_T(s)}, \quad \|\alpha - v\|_{\ell_\infty} \leq \frac{u_T(s)}{1 - 2u_T(s)} \quad \text{and} \quad \|\xi\|_{\ell_\infty} \leq \frac{u_T(s)}{1 - 2u_T(s)}.$$

**Proof of Lemma B.1.** Thanks to (31), (28) and (100), we have:

$$\left( \eta_{\alpha,\xi}(\theta_1^*), \dots, \eta_{\alpha,\xi}(\theta_s^*), \eta_{\alpha,\xi}^{[1]}(\theta_1^*), \dots, \eta_{\alpha,\xi}^{[1]}(\theta_s^*) \right)^\top = \Gamma \begin{pmatrix} \alpha \\ \xi \end{pmatrix}.$$

Thus, solving (101) is equivalent to solving,

$$(103) \quad \Gamma \begin{pmatrix} \alpha \\ \xi \end{pmatrix} = \begin{pmatrix} v \\ 0_s \end{pmatrix},$$

with  $0_s$  the vector of size  $s$  with all its components equal to zero.

We first show that  $\Gamma$  is non singular so that  $\alpha$  and  $\xi$  exist and are uniquely defined. Using Lemma C.3 based on the Schur complement,  $\Gamma$  has an inverse provided that  $\Gamma^{[1,1]}$  and  $\Gamma_{SC} := \Gamma^{[0,0]} - \Gamma^{[1,0]^\top} [\Gamma^{[1,1]}]^{-1} \Gamma^{[1,0]}$  are non singular. We recall that if  $M$  is a matrix such that,  $\|I - M\|_{\text{op}} < 1$ , then  $M$  is non singular,  $M^{-1} = \sum_{i \geq 0} (I - M)^i$  and  $\|M^{-1}\|_{\text{op}} \leq$

$$\left( 1 - \|I - M\|_{\text{op}} \right)^{-1}.$$

Recall that by assumption  $u_T(s) \leq 1/2$ . Then, the second inequality in (97) imply that  $\|I - \Gamma^{[1,1]}\|_{\text{op}} < 1$  and thus  $\Gamma^{[1,1]}$  is non singular. We now prove that  $\Gamma_{SC}$  is also non singular. Using the triangle inequality we have:

$$\begin{aligned} \|I - \Gamma_{SC}\|_{\text{op}} &= \left\| I - \Gamma^{[0,0]} + \Gamma^{[1,0]^\top} [\Gamma^{[1,1]}]^{-1} \Gamma^{[1,0]} \right\|_{\text{op}} \\ &\leq \left\| I - \Gamma^{[0,0]} \right\|_{\text{op}} + \left\| \Gamma^{[1,0]^\top} [\Gamma^{[1,1]}]^{-1} \Gamma^{[1,0]} \right\|_{\text{op}}. \end{aligned}$$

Let us bound the terms on the right hand side of the inequality above. To bound  $\|\Gamma^{[1,0]\top}[\Gamma^{[1,1]}]^{-1}\Gamma^{[1,0]}\|_{\text{op}}$  notice that:

$$\left\|\Gamma^{[1,0]\top}[\Gamma^{[1,1]}]^{-1}\Gamma^{[1,0]}\right\|_{\text{op}} \leq \|\Gamma^{[1,0]}\|_{\text{op}} \left\|\Gamma^{[1,0]\top}\right\|_{\text{op}} \left\|[\Gamma^{[1,1]}]^{-1}\right\|_{\text{op}}.$$

We have, thanks to (97) for the second inequality:

$$(104) \quad \left\|[\Gamma^{[1,1]}]^{-1}\right\|_{\text{op}} \leq \frac{1}{1 - \|I - \Gamma^{[1,1]}\|_{\text{op}}} \leq \frac{1}{1 - u_T(s)}.$$

Using (97), we get:

$$\|I - \Gamma_{SC}\|_{\text{op}} \leq u_T(s) + \frac{u_T(s)^2}{1 - u_T(s)} = \frac{u_T(s)}{1 - u_T(s)}.$$

By assumption, we have  $u_T(s) \leq H_\infty^{(2)}(r, \rho) < 1/2$ . Hence, we have  $\frac{u_T(s)}{1 - u_T(s)} < 1$  and thus,  $\Gamma_{SC}$  is non singular. Furthermore, we get:

$$(105) \quad \|\Gamma_{SC}^{-1}\|_{\text{op}} \leq \frac{1}{1 - \|I - \Gamma_{SC}\|_{\text{op}}} \leq \frac{1 - u_T(s)}{1 - 2u_T(s)}.$$

As the matrices  $\Gamma^{[1,1]}$  and  $\Gamma_{SC}$  are non singular, we deduce that the matrix  $\Gamma$  is non singular.

We now give bounds related to  $\alpha$  and  $\xi$ . The Lemma C.3 on the Schur complement gives also that:

$$\alpha = \Gamma_{SC}^{-1}v \quad \text{and} \quad \xi = -[\Gamma^{[1,1]}]^{-1}\Gamma^{[1,0]}\Gamma_{SC}^{-1}v.$$

Hence, we deduce that:

$$\begin{aligned} \|\alpha\|_{\ell_\infty} &\leq \|\Gamma_{SC}^{-1}\|_{\text{op}} \|v\|_{\ell_\infty} \leq \frac{1 - u_T(s)}{1 - 2u_T(s)}, \\ \|\xi\|_{\ell_\infty} &\leq \left\|[\Gamma^{[1,1]}]^{-1}\Gamma^{[1,0]}\Gamma_{SC}^{-1}\right\|_{\text{op}} \|v\|_{\ell_\infty} \leq \left\|[\Gamma^{[1,1]}]^{-1}\right\|_{\text{op}} \left\|\Gamma^{[1,0]}\right\|_{\text{op}} \|\Gamma_{SC}^{-1}\|_{\text{op}} \leq \frac{u_T(s)}{1 - 2u_T(s)}, \\ \|\alpha - v\|_{\ell_\infty} &\leq \|(\Gamma_{SC}^{-1} - I)\|_{\text{op}} \|v\|_{\ell_\infty} \leq \|\Gamma_{SC} - I\|_{\text{op}} \|\Gamma_{SC}^{-1}\|_{\text{op}} \leq \frac{u_T(s)}{1 - 2u_T(s)}. \end{aligned}$$

This finishes the proof.  $\square$

We now fix a sign vector  $v = (v_1, \dots, v_s)^\top \in \{-1, 1\}^s$  and consider  $p_{\alpha, \xi}$  and  $\eta_{\alpha, \xi}$  with  $\alpha$  and  $\xi$  characterized by (101) from Lemma B.1. Let  $e_\ell \in \mathbb{R}^s$  be the vector with all the entries equal to zero but the  $\ell$ -th which is equal to 1.

**Proof of (iii) from Assumption 6.1** with  $C_F = \varepsilon_\infty(r/\rho)/10$ . Let  $\theta \in \Theta_T$  such that  $\mathfrak{d}_T(\theta, \mathcal{Q}^*) > r$  (far region). It is enough to prove that  $|\eta_{\alpha, \xi}(\theta)| \leq 1 - C_F$ . Let  $\theta_\ell^*$  be one of the elements of  $\mathcal{Q}^*$  closest to  $\theta$  in terms of the metric  $\mathfrak{d}_T$ . Since  $\vartheta^* \in \Theta_{T, 2\rho_T\delta_\infty(u_\infty, s)}^s$ , we have, by the triangle inequality that for any  $k \neq \ell$ :

$$2\rho_T\delta_\infty(u_\infty, s) < \mathfrak{d}_T(\theta_\ell^*, \theta_k^*) \leq \mathfrak{d}_T(\theta_\ell^*, \theta) + \mathfrak{d}_T(\theta, \theta_k^*) \leq 2\mathfrak{d}_T(\theta, \theta_k^*).$$

Hence, we have  $\vartheta_{\ell, \theta}^* \in \Theta_{T, \rho_T\delta_\infty(u_\infty, s)}^s$ , where  $\vartheta_{\ell, \theta}^*$  denotes the vector  $\vartheta^*$  whose  $\ell$ -th coordinate has been replaced by  $\theta$ . Then, we obtain from Lemma 7.3 that  $\Theta_{T, \rho_T\delta_\infty(u_\infty, s)}^s \subseteq \Theta_{T, \delta_T(u_T(s), s)}^s$  and thus:

$$(106) \quad \vartheta_{\ell, \theta}^* \in \Theta_{T, \delta_T(u_T(s), s)}^s.$$

We denote by  $\Gamma_{\ell, \theta}$  (resp.  $\Gamma_{\ell, \theta}^{[i, j]}$ ) the matrix  $\Gamma$  (resp.  $\Gamma^{[i, j]}$ ) in (96) where  $\vartheta^*$  has been replaced by  $\vartheta_{\ell, \theta}^*$ . Notice the upper bounds (97) also hold for  $\Gamma_{\ell, \theta}$  because of (106). Recall we have Equalities (32) on the diagonal of the kernel  $\mathcal{K}_T$  and its derivatives. Elementary calculations give with  $\eta_{\alpha, \xi}$  from Lemma B.1 that:

$$(107) \quad \eta_{\alpha, \xi}(\theta) = e_\ell^\top \left( \Gamma_{\ell, \theta}^{[0, 0]} - I \right) \alpha + \mathcal{K}_T(\theta, \theta_\ell^*) \alpha_\ell + e_\ell^\top \Gamma_{\ell, \theta}^{[1, 0]\top} \xi + \mathcal{K}_T^{[0, 1]}(\theta, \theta_\ell^*) \xi_\ell.$$

We deduce that:

$$(108) \quad |\eta_{\alpha,\xi}(\theta)| \leq \left\| \Gamma_{\ell,\theta}^{[0,0]} - I \right\|_{\text{op}} \|\alpha\|_{\ell_\infty} + \|\alpha\|_{\ell_\infty} |\mathcal{K}_T(\theta, \theta_\ell^*)| + \left\| \Gamma_{\ell,\theta}^{[1,0]\top} \right\|_{\text{op}} \|\xi\|_{\ell_\infty} + |\mathcal{K}_T^{[0,1]}(\theta, \theta_\ell^*)| \|\xi\|_{\ell_\infty}.$$

Since  $\theta$  belongs to the “far region”, we have by definition of  $\varepsilon_T(r)$  given in (38) that:

$$(109) \quad |\mathcal{K}_T(\theta, \theta_\ell^*)| \leq 1 - \varepsilon_T(r).$$

The triangle inequality, the definitions (37) of  $\mathcal{V}_T$  and (34) of  $L_{1,0}$ , give:

$$(110) \quad |\mathcal{K}_T^{[0,1]}(\theta, \theta_\ell^*)| \leq L_{0,1} + \mathcal{V}_T.$$

Then, using (97) (which holds for  $\Gamma_{\ell,\theta}$  thanks to (106)), we get that:

$$|\eta_{\alpha,\xi}(\theta)| \leq 1 - \varepsilon_T(r) + \frac{u_T(s)}{1 - 2u_T(s)} (2 + L_{1,0} + \mathcal{V}_T).$$

Notice that the function  $r \mapsto \varepsilon_\infty(r)$  is increasing. Since  $\rho_T \leq \rho$ , we get by Lemma 7.1 that:

$$(111) \quad \varepsilon_T(r) \geq \varepsilon_\infty(r/\rho_T) - \mathcal{V}_T \geq \varepsilon_\infty(r/\rho) - \mathcal{V}_T.$$

By assumption, we have  $u_T(s) \leq H_\infty^{(2)}(r, \rho) \leq 1/4$ . Hence, we have  $\frac{1}{1-2u_T(s)} \leq 2$ . We also have  $\mathcal{V}_T \leq 1/2$ . Therefore, we get:

$$|\eta_{\alpha,\xi}(\theta)| \leq 1 - \varepsilon_\infty(r/\rho) + \mathcal{V}_T + u_T(s) (5 + 2L_{1,0}).$$

The assumption  $u_T(s) \leq H_\infty^{(2)}(r, \rho)$  gives:

$$(112) \quad u_T(s) \leq \frac{8}{10(5 + 2L_{1,0})} \varepsilon_\infty(r/\rho).$$

The assumption  $\mathcal{V}_T \leq H_\infty^{(1)}(r, \rho)$  gives  $\mathcal{V}_T \leq \varepsilon_\infty(r/\rho)/10$ . Hence, we have  $|\eta_{\alpha,\xi}(\theta)| \leq 1 - \frac{\varepsilon_\infty(r/\rho)}{10}$ . Thus, Property (iii) from Assumption 6.1 holds with  $C_F = \varepsilon_\infty(r/\rho)/10$ .

**Proof of (i) from Assumption 6.1** with  $C_N = \nu_\infty(\rho r)/180$ . Let  $\theta \in \Theta_T$  such that  $\mathfrak{d}_T(\theta, \mathcal{Q}^*) \leq r$ . Let  $\ell \in \{1, \dots, s\}$  such that  $\theta \in \mathcal{B}_T(\theta_\ell^*, r)$  (“near region”). Thus, it is enough to prove that  $|\eta_{\alpha,\xi}(\theta)| \leq 1 - C_N \mathfrak{d}_T(\theta_\ell^*, \theta)^2$ . This will be done by using Lemma C.4 to obtain a quadratic decay on  $\eta_{\alpha,\xi}$  from a bound on its second Riemannian derivative.

Recall that the function  $\eta_{\alpha,\xi}$  is twice continuously differentiable. Set  $\eta_{\alpha,\xi}^{[2]} = \tilde{D}_{2,T}[\eta_{\alpha,\xi}]$ . Differentiating (100) and using that  $\mathcal{K}_T^{[2,0]}(\theta, \theta) = -1$  and  $\mathcal{K}_T^{[2,1]}(\theta, \theta) = 0$ , see (32), we deduce that:

$$(113) \quad \eta_{\alpha,\xi}^{[2]}(\theta) = e_\ell^\top (I + \Gamma_{\ell,\theta}^{[2,0]}) \alpha + \mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*) e_\ell^\top \alpha + e_\ell^\top \Gamma_{\ell,\theta}^{[2,1]} \xi + \mathcal{K}_T^{[2,1]}(\theta, \theta_\ell^*) e_\ell^\top \xi.$$

Since  $v = (v_1, \dots, v_s)^\top \in \{-1, 1\}^s$  is a sign vector, we get:

$$(114) \quad \eta_{\alpha,\xi}^{[2]}(\theta) - v_\ell \mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*) = e_\ell^\top (I + \Gamma_{\ell,\theta}^{[2,0]}) \alpha + \mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*) e_\ell^\top (\alpha - v) + e_\ell^\top \Gamma_{\ell,\theta}^{[2,1]} \xi + \mathcal{K}_T^{[2,1]}(\theta, \theta_\ell^*) e_\ell^\top \xi.$$

The triangle inequality and the definition of  $\mathcal{V}_T$  give:

$$(115) \quad |\mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*)| \leq L_{2,0} + \mathcal{V}_T \quad \text{and} \quad |\mathcal{K}_T^{[2,1]}(\theta, \theta_\ell^*)| \leq L_{2,1} + \mathcal{V}_T,$$

where  $L_{2,0}$  and  $L_{1,2}$  are defined in (34). We deduce from (106), the definition of  $\delta_T$  in (43) and (44) that:

$$(116) \quad \left\| I + \Gamma_{\ell,\theta}^{[2,0]} \right\|_{\text{op}} \leq u_T(s) \quad \text{and} \quad \left\| \Gamma_{\ell,\theta}^{[2,1]} \right\|_{\text{op}} \leq u_T(s).$$

We deduce from (114) that:

$$\begin{aligned} |\eta_{\alpha,\xi}^{[2]}(\theta) - v_\ell \mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*)| &\leq \|\alpha\|_{\ell_\infty} \left\| I + \Gamma_{\ell,\theta}^{[2,0]} \right\|_{\text{op}} + \|\alpha - v\|_{\ell_\infty} (L_{2,0} + \mathcal{V}_T) + \|\xi\|_{\ell_\infty} \left( \left\| \Gamma_{\ell,\theta}^{[2,1]} \right\|_{\text{op}} + L_{2,1} + \mathcal{V}_T \right) \\ &\leq \frac{u_T(s)}{1 - 2u_T(s)} (1 + L_{2,0} + L_{2,1} + 2\mathcal{V}_T). \end{aligned}$$

By assumption, we have  $u_T(s) \leq H_\infty^{(2)}(r, \rho) \leq 1/4$ . Hence, we have  $\frac{1}{1-2u_T(s)} \leq 2$ . Furthermore, we have by assumption  $\mathcal{V}_T \leq H_\infty^{(1)}(r, \rho) \leq 1/2$  and  $u_T(s) \leq H_\infty^{(2)}(r, \rho)$ . In particular, we have:

$$u_T(s) \leq \frac{8}{9(2L_{2,0} + 2L_{2,1} + 4)} \nu_\infty(\rho r).$$

Therefore, we obtain:

$$(117) \quad |\eta_{\alpha,\xi}^{[2]}(\theta) - v_\ell \mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*)| \leq \frac{8}{9} \nu_\infty(\rho r).$$

We now check that the hypotheses of Lemma C.4-(ii) hold in order to obtain a quadratic decay on  $\eta_{\alpha,\xi}$  from the bound (117). First recall that  $\eta_{\alpha,\xi}$  is twice continuously differentiable and have the interpolation properties (101). By the triangle inequality and since by assumption  $\mathcal{V}_T \leq L_{2,0}$  we have:

$$\sup_{\Theta_T^*} |\mathcal{K}_T^{[2,0]}| \leq L_{2,0} + \mathcal{V}_T \leq 2L_{2,0}.$$

Then, Lemma 7.1 ensures that for any  $\theta, \theta'$  in  $\Theta_T$  such that  $\mathfrak{d}_T(\theta, \theta') \leq r$  we have:

$$-\mathcal{K}_T^{[2,0]}(\theta, \theta') \geq \nu_\infty(r\rho_T) - \mathcal{V}_T \geq \nu_\infty(\rho r) - \mathcal{V}_T \geq \frac{9}{10} \nu_\infty(\rho r),$$

where we used that the function  $r \mapsto \nu_\infty(r)$  is decreasing and  $\rho_T \leq \rho$  for the second inequality and that  $\mathcal{V}_T \leq H_\infty^{(1)}(r, \rho) \leq \nu_\infty(\rho r)/10$  for the last inequality.

Set  $\delta = \frac{8}{9} \nu_\infty(\rho r)$ ,  $\varepsilon = \frac{9}{10} \nu_\infty(\rho r)$ ,  $L = 2L_{2,0}$ . As  $r < L^{-\frac{1}{2}}$  and  $\delta < \varepsilon$ , we apply Lemma C.4-(ii) and get for  $\theta \in \mathcal{B}_T(\theta_\ell^*, r)$ :

$$|\eta_{\alpha,\xi}(\theta)| \leq 1 - \frac{\nu_\infty(\rho r)}{180} \mathfrak{d}_T(\theta, \theta_\ell^*)^2.$$

**Proof of (ii) from Assumption 6.1** with  $C'_N = (5L_{2,0} + L_{2,1} + 4)/8$ . Let  $\theta \in \Theta_T$  such that  $\mathfrak{d}_T(\theta, \mathcal{Q}^*) \leq r$ . Let  $\ell \in \{1, \dots, s\}$  such that  $\theta \in \mathcal{B}_T(\theta_\ell^*, r)$  ("near region"). We shall prove that  $|\eta_{\alpha,\xi}(\theta) - v_\ell| \leq C'_N \mathfrak{d}_T(\theta, \theta_\ell^*)^2$ .

Let us consider the function  $f : \theta \rightarrow \eta_{\alpha,\xi}(\theta) - v_\ell$ . We will bound the second covariant derivative  $f^{[2]} = \tilde{D}_{2,T}[f]$  of  $f$  and apply Lemma C.4-(i) on  $f$  to prove the property (ii) for  $\eta_{\alpha,\xi}$ . Notice that  $f$  is twice continuously differentiable. By construction, see (101), we have  $f(\theta_\ell^*) = 0$  and  $f^{[1]}(\theta_\ell^*) = 0$ . Since  $f^{[2]} = \eta_{\alpha,\xi}^{[2]}$ , we deduce from (113), the bounds (115) that:

$$|f^{[2]}(\theta)| \leq \|\alpha\|_{\ell_\infty} \left\| I + \Gamma_{\ell,\theta}^{[2,0]} \right\|_{\text{op}} + \|\alpha\|_{\ell_\infty} (L_{2,0} + \mathcal{V}_T) + \|\xi\|_{\ell_\infty} \left\| \Gamma_{\ell,\theta}^{[2,1]} \right\|_{\text{op}} + \|\xi\|_{\ell_\infty} (L_{2,1} + \mathcal{V}_T).$$

Using (116), and the bounds on  $\alpha$  and  $\xi$  from Lemma B.1, we get:

$$|f^{[2]}(\theta)| \leq \frac{1 - u_T(s)}{1 - 2u_T(s)} (L_{2,0} + \mathcal{V}_T + u_T(s)) + \frac{u_T(s)}{1 - 2u_T(s)} (L_{2,1} + \mathcal{V}_T + u_T(s)).$$

Since  $u_T(s) \leq H_\infty^{(2)}(r, \rho) \leq 1/6$  and  $\mathcal{V}_T \leq H_\infty^{(1)}(r, \rho) \leq 1/2$ , we get:

$$|f^{[2]}(\theta)| \leq \frac{5}{4} L_{2,0} + \frac{1}{4} L_{2,1} + 1.$$



We get thanks to Lemma C.4-(i) on the function  $f$  that for any  $\theta \in \mathcal{B}_T(\theta_\ell^*, r)$ :

$$|\eta_{\alpha,\xi}(\theta) - v_\ell| \leq \frac{1}{8} (5L_{2,0} + L_{1,2} + 4) \mathfrak{d}_T(\theta, \theta_\ell^*)^2.$$

**Proof of (iv) from Assumption 6.1** with  $C_B = 2$ . Recall the definition of  $p_{\alpha,\xi}$  in (98). Elementary calculations give using the definitions of  $\Gamma^{[0,0]}$ ,  $\Gamma^{[1,1]}$  and  $\Gamma^{[1,1]}$  in (96):

$$\begin{aligned} \|p_{\alpha,\xi}\|_T^2 &\leq 2 \left\| \sum_{k=1}^s \alpha_k \phi_T(\theta_k^*) \right\|_T^2 + 2 \left\| \sum_{k=1}^s \xi_k \phi_T^{[1]}(\theta_k^*) \right\|_T^2 \\ &= 2\alpha^\top \Gamma^{[0,0]} \alpha + 2\xi^\top \Gamma^{[1,1]} \xi \\ &\leq 2\|\alpha\|_{\ell_1} \|\alpha\|_{\ell_\infty} \left\| \Gamma^{[0,0]} \right\|_{\text{op}} + 2\|\xi\|_{\ell_1} \|\xi\|_{\ell_\infty} \left\| \Gamma^{[1,1]} \right\|_{\text{op}}. \end{aligned}$$

Using that  $\|I\|_{\text{op}} = 1$  and (97), we get that:

$$\left\| \Gamma^{[0,0]} \right\|_{\text{op}} \leq (1 + u_T(s)) \quad \text{and} \quad \left\| \Gamma^{[1,1]} \right\|_{\text{op}} \leq (1 + u_T(s)).$$

By assumption we have  $u_T(s) \leq H_\infty^{(2)}(r, \rho) \leq \frac{1}{6}$ . We deduce that:

$$\|p_{\alpha,\xi}\|_T^2 \leq 2(1 + u_T(s)) \frac{(1 - u_T(s))^2 + u_T(s)^2}{(1 - 2u_T(s))^2} s \leq 4s.$$

This gives:

$$(118) \quad \|p_{\alpha,\xi}\|_T \leq 2\sqrt{s}.$$

We proved that (i)-(iv) from Assumption 6.1 stand. By assumption we also have that for all  $\theta \neq \theta' \in \mathcal{Q}^* : \mathfrak{d}_T(\theta, \theta') > 2r$ , therefore Assumption 6.1 holds.

This finishes the proof of Proposition 7.4.

## B.2. Proof of Proposition 7.6 (Construction of an interpolating derivative certificate)

This section is devoted to the proof of Proposition 7.6 and is close to Section B.1. Let  $T \in \mathbb{N}$  and  $s \in \mathbb{N}^*$ . Recall Assumptions 3.2 (and thus 3.1 on the regularity of  $\varphi_T$ ) and 5.1 on the regularity of the limit kernel  $\mathcal{K}_\infty$  are in force. Set  $u'_\infty \in (0, 1/6)$ . We denote by  $\|\cdot\|_{\text{op}}$  the operator norm associated to the  $\ell_\infty$  norm on  $\mathbb{R}^s$ . By assumption  $\delta_\infty(u'_\infty, s)$  is finite. Let  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta_{T, 2\rho_T}^s \delta_\infty(u'_\infty, s)$ . We note  $\mathcal{Q}^* = \{\theta_i^*, 1 \leq i \leq s\}$  the set of parameters of cardinal  $s$ . Let  $\alpha = (\alpha_1, \dots, \alpha_s)^\top$  and  $\xi = (\xi_1, \dots, \xi_s)^\top$  be elements of  $\mathbb{R}^s$ . Recall  $p_{\alpha,\xi}$ ,  $\eta_{\alpha,\xi}$  and  $\eta_{\alpha,\xi}^{[1]} = \tilde{D}_{1,T}[\eta_{\alpha,\xi}]$  given by (98), (99) and (100).

The next lemma is similar to Lemma B.1, but notice that in Lemma B.2 the function  $\eta_{\alpha,\xi}$  vanished on  $\mathcal{Q}^*$  and has a derivative that interpolates a sign vector, whereas in Lemma B.1 it is the opposite.

Recall the definition of  $\mathcal{V}_T$  from (37) and define  $u'_T(s) = u'_\infty + (s-1)\mathcal{V}_T$ . We remark that (97) holds with  $u_T(s)$  replaced by  $u'_T(s)$  because of (95).

**Lemma B.2.** *Let  $v = (v_1, \dots, v_s)^\top \in \{-1, 1\}^s$  be a sign vector. Assume that (97) holds with  $u_T(s)$  replaced by  $u'_T(s) < 1/2$ . Under Assumption 3.2, there exist unique  $\alpha, \xi \in \mathbb{R}^s$  such that:*

$$(119) \quad \eta_{\alpha,\xi}(\theta_k^*) = 0 \quad \text{and} \quad \eta_{\alpha,\xi}^{[1]}(\theta_k^*) = v_k \quad \text{for} \quad 1 \leq k \leq s.$$

Furthermore, we have:

$$(120) \quad \|\alpha\|_{\ell_\infty} \leq \frac{u'_T(s)}{1 - 2u'_T(s)} \quad \text{and} \quad \|\xi\|_{\ell_\infty} \leq \frac{1 - u'_T(s)}{1 - 2u'_T(s)}.$$

**Proof.** Thus, with  $0_s$  the vector of size  $s$  with all its components equal to zero and  $\Gamma$  defined by (96), Equation (119) is equivalent to:

$$(121) \quad \Gamma \begin{pmatrix} \alpha \\ \xi \end{pmatrix} = \begin{pmatrix} 0_s \\ v \end{pmatrix}.$$

According to the proof of Lemma B.1, the matrices  $\Gamma_{SC} = \Gamma^{[0,0]} - \Gamma^{[1,0]\top} [\Gamma^{[1,1]}]^{-1} \Gamma^{[1,0]}$ ,  $\Gamma^{[1,1]}$  and  $\Gamma$  are non singular. Thus the vectors  $\alpha$  and  $\xi$  exist and are uniquely determined by (121). From Lemma C.3, we deduce that:

$$\alpha = -\Gamma_{SC}^{-1} \Gamma^{[1,0]\top} [\Gamma^{[1,1]}]^{-1} v \quad \text{and} \quad \xi = \left( I + [\Gamma^{[1,1]}]^{-1} \Gamma^{[1,0]} \Gamma_{SC}^{-1} \Gamma^{[1,0]\top} \right) [\Gamma^{[1,1]}]^{-1} v.$$

Using (105), (97) and (104) and replacing  $u_T(s)$  by  $u'_T(s)$ , we easily obtain the inequalities (120).  $\square$

We fix the sign vector  $v = (v_1, \dots, v_s)^\top \in \{-1, 1\}^s$  and consider  $p_{\alpha, \xi}$  and  $\eta_{\alpha, \xi}$  given by (98) and (99), with  $\alpha$  and  $\xi$  given by Lemma B.2.

**Proof of (i) from Assumption 6.2** with  $c_N = (L_{0,2} + L_{2,1} + 7)/8$ . We define the function  $f : \theta \mapsto \eta_{\alpha, \xi}(\theta) - v_\ell \text{sign}(\theta - \theta_\ell^*) \mathfrak{d}_T(\theta, \theta_\ell^*)$  on  $\Theta$ . To prove the Property (i), we will bound the second covariant derivative of  $f$ , that is  $f^{[2]} := \tilde{D}_{2;T}[f]$ , and apply Lemma C.4-(i). Recall  $\mathfrak{d}_T(\theta, \theta_\ell^*) = |G_T(\theta) - G_T(\theta_\ell^*)|$  with  $G_T$  a primitive of  $\sqrt{g_T}$ , and thus  $f(\theta) = \eta_{\alpha, \xi}(\theta) - v_\ell (G_T(\theta) - G_T(\theta_\ell^*))$ . We deduce that  $f$  is twice continuously differentiable on  $\Theta$ ; and elementary calculations give  $f^{[2]} = \eta_{\alpha, \xi}^{[2]}$ .

Let  $\theta \in \Theta_T$  and let  $\theta_\ell^*$  be one of the elements of  $\mathcal{Q}^*$  closest to  $\theta$  in terms of the metric  $\mathfrak{d}_T$ . Recall the notations  $\Gamma_{\ell, \theta}$  (resp.  $\Gamma_{\ell, \theta}^{[i,j]}$ ) and  $\vartheta_{\ell, \theta}^*$  from the proof of Proposition 7.4. Since  $f^{[2]} = \eta_{\alpha, \xi}^{[2]}$ , we deduce from (113) that:

$$(122) \quad |f^{[2]}(\theta)| \leq \left\| I + \Gamma_{\ell, \theta}^{[2,0]} \right\|_{\text{op}} \|\alpha\|_{\ell_\infty} + \|\alpha\|_{\ell_\infty} |\mathcal{K}_T^{[2,0]}(\theta, \theta_\ell^*)| + \|\xi\|_{\ell_\infty} \left\| \Gamma_{\ell, \theta}^{[2,1]} \right\|_{\text{op}} + \|\xi\|_{\ell_\infty} |\mathcal{K}_T^{[2,1]}(\theta, \theta_\ell^*)|.$$

Notice that (106) holds with  $u_T(s)$  replaced by  $u'_T(s)$ . Using (115) and (116) and the bounds (120) on  $\alpha$  and  $\xi$  from Lemma B.2, we get:

$$|f^{[2]}(\theta)| \leq \frac{u'_T(s)}{1 - 2u'_T(s)} (L_{2,0} + \mathcal{V}_T + u'_T(s)) + \frac{1 - u'_T(s)}{1 - 2u'_T(s)} (L_{2,1} + \mathcal{V}_T + u'_T(s)).$$

By assumption, we have  $u'_T(s) \leq 1/6$  and  $\mathcal{V}_T \leq 1$ . Hence, we obtain:

$$|f^{[2]}(\theta)| \leq \frac{1}{4} L_{2,0} + \frac{5}{4} L_{2,1} + \frac{7}{4}.$$

Since  $f(\theta_\ell^*) = 0$  and  $f^{[1]}(\theta_\ell^*) = 0$  as well, using Lemma C.4 (i), we get, with  $c_N = (L_{2,0} + 5L_{2,1} + 7)/8$ :

$$|\eta_{\alpha, \xi}(\theta) - v_\ell \text{sign}(\theta - \theta_\ell^*) \mathfrak{d}_T(\theta, \theta_\ell^*)| = |f(\theta)| \leq c_N \mathfrak{d}_T(\theta, \theta_\ell^*)^2.$$

**Proof of (ii) from Assumption 6.2** with  $c_F = (5L_{1,0} + 7)/4$ . Let  $\theta \in \Theta_T$ , we shall prove that  $|\eta_{\alpha, \xi}(\theta)| \leq c_F$ . Let  $\theta_\ell^*$  be one of the elements of  $\mathcal{Q}^*$  closest to  $\theta$  in terms of the metric  $\mathfrak{d}_T$ . We deduce from (107) that:

$$|\eta_{\alpha, \xi}(\theta)| \leq \|\alpha\|_{\ell_\infty} \left\| \Gamma_{\ell, \theta}^{[0,0]} - I \right\|_{\text{op}} + \|\alpha\|_{\ell_\infty} |\mathcal{K}_T(\theta, \theta_\ell^*)| + \|\xi\|_{\ell_\infty} \left\| \Gamma_{\ell, \theta}^{[1,0]\top} \right\|_{\text{op}} + \|\xi\|_{\ell_\infty} |\mathcal{K}_T^{[0,1]}(\theta, \theta_\ell^*)|.$$

Using (97), (32), (110) and the bounds (120) on  $\alpha$  and  $\xi$  from Lemma B.2, we get:

$$|\eta_{\alpha, \xi}(\theta)| \leq \frac{u'_T(s)}{1 - 2u'_T(s)} (1 + u'_T(s)) + \frac{1 - u'_T(s)}{1 - 2u'_T(s)} (L_{1,0} + \mathcal{V}_T + u'_T(s)).$$

By assumption, we have  $u'_T(s) \leq 1/6$ , and thus  $\frac{1}{1 - 2u'_T(s)} \leq 3/2$ . Since  $\mathcal{V}_T \leq 1$ , we obtain:

$$|\eta_{\alpha, \xi}(\theta)| \leq \frac{5}{4} L_{1,0} + \frac{7}{4}.$$

**Proof of (iii) from Assumption 6.2** with  $c_B = 2$ . Using very similar arguments as in the proof of (118) (taking care that the upper bound of the  $\ell_\infty$  norm of  $\alpha$  and  $\xi$  are given by (120)) we also get  $\|p_{\alpha,\xi}\|_T \leq 2\sqrt{s}$ .

We proved that (i)-(ii) from Assumption 6.2 stand for any  $\theta \in \Theta_T$ . Hence Assumption 6.2 holds for any positive  $r$  such that for all  $\theta \neq \theta' \in \mathcal{Q}^* : \mathfrak{d}_T(\theta, \theta') > 2r$ .

This finishes the proof of Proposition 7.6.

## Appendix C: Auxiliary Lemmas

We recall in the next section some basic results on the Fréchet derivative and the Bochner integral. Then, we provide the proofs of the intermediate results.

### C.1. The Fréchet derivative and the Bochner integral

The Fréchet derivative and Bochner integrals are defined for Banach space valued functions, but we shall only consider the case of Hilbert space valued functions.

Let  $(H, \langle \cdot, \cdot \rangle)$  be an Hilbert space and let  $\Theta$  be an interval of  $\mathbb{R}$ . We note  $\|\cdot\|$  the norm associated to the scalar product. A function  $f$  from  $\Theta$  to  $H$  is Fréchet differentiable at  $\theta \in \Theta$  if it is continuous at  $\theta$  and there exists an element  $\partial_\theta f \in H$  such that:

$$\lim_{h \rightarrow 0; \theta+h \in \Theta} \left\| \frac{f(\theta+h) - f(\theta)}{h} - \partial_\theta f(\theta) \right\| = 0.$$

The derivative of  $f$  is the function  $\partial_\theta f : \theta \mapsto \partial_\theta f(\theta)$  defined on  $\Theta$  when it exists. We also define by recurrence the derivative  $\partial_\theta^i f$  of order  $i \in \mathbb{N}^*$  of  $f$  as the derivative of  $\partial_\theta^{i-1} f$ , with the convention that  $\partial_\theta^0 f = f$ , and say that  $f$  is of class  $\mathcal{C}^i$  if the derivatives  $\partial_\theta^j f$  exist and are continuous on  $\Theta$  for  $j \in \{0, \dots, i\}$ . The standard differentiating rules for composition, addition and multiplication apply to the Fréchet derivative. We refer to [36] for a complete presentation of the subject. By definition, if  $f$  is differentiable at  $\theta \in \Theta$ , then we have for all  $g \in H$  that:

$$(123) \quad \partial_\theta \langle f(\theta), g \rangle = \langle \partial_\theta f(\theta), g \rangle.$$

The Bochner integral extends the Lebesgue integral. We refer to [4, Chapter 1] and [3, Section 11.8] for further details on the Bochner integral. We endow the interval  $\Theta \subset \mathbb{R}$  with its usual Borel sigma field inherited from the Borel sigma field on  $\mathbb{R}$  and a measure  $\mu$ . A function  $f$  from  $\Theta$  to  $H$  is strongly measurable if it is the limit of simple functions or equivalently, see [3, Lemma 11.37], if the map  $\theta \mapsto \langle f(\theta), g \rangle$  is measurable for all  $g \in H$  and  $f(\theta)$  lies for  $\mu$ -almost every  $\theta \in \Theta$  in a closed separable subspace of  $H$ . In particular if the function  $f$  is continuous, then it is strongly measurable, see [4, Corollary 1.1.2]. If  $f$  is strongly measurable, then the norm  $\|f\|$  is a measurable function from  $\Theta$  to  $\mathbb{R}$ , see [3, Lemma 11.39]. Then a function  $f$  defined on  $\Theta$  (endowed with the Lebesgue measure) is Bochner integrable if and only if it is strongly measurable and if  $\|f\|$  is integrable; in which case, we have  $\|\int f(\theta) d\theta\| \leq \int \|f(\theta)\| d\theta$ , see [3, Theorem 11.44] (which is easily extended from finite measure to  $\sigma$ -finite measure, see also [4, Theorem 1.1.4] in this direction). We remark that the fundamental theorem of calculus is still valid in this framework, see [4, Proposition 1.2.2]. In particular, if  $f$  is continuous and Bochner integrable on  $\Theta$  and  $\theta_0 \in \Theta$ , then, we have:

$$(124) \quad F'(\theta) = f(\theta) \quad \text{where} \quad F(\theta) = \int_{\theta_0}^{\theta} f(q) dq.$$

As a particular case of [3, Lemma 11.45], if  $f$  is Bochner integrable on  $\Theta$ , then for all  $g \in H$ , we have that:

$$(125) \quad \int_{\Theta} \langle f(\theta), g \rangle d\theta = \langle \int_{\Theta} f(\theta) d\theta, g \rangle.$$

### C.2. Tail bounds for suprema of Gaussian processes

In order to prove Theorems 2.1 and 2.5, we provide in Lemma C.1 a bound with high probability of the supremum of a Gaussian process given by  $\theta \mapsto \langle w_T, h(\theta) \rangle_T$ , where  $w_T$  is a noise process and  $h$  is a function from  $\Theta$ , an interval of  $\mathbb{R}$ , to the Hilbert space  $(H_T, \langle \cdot, \cdot \rangle_T)$ . The next lemma is in the spirit of [6, Proposition 4.1] (where one assumes that the Gaussian process has unitary variance); its proof is given at the end of this section and relies on Lemma C.2.

We denote by  $\mathfrak{d}_T$  the Riemannian metric associated to the kernel  $\mathcal{K}_T$ , see also Section 4.2. Recall definitions (20) and (22) and set  $f^{[1]}(\theta) = \tilde{D}_{1,T}[f](\theta) = \partial_\theta f(\theta) / \sqrt{g_T(\theta)}$  with  $g_T$  defined in (30).

**Lemma C.1.** *Let  $T \in \mathbb{N}$  be fixed. Suppose that Assumptions 3.1 and 3.2 hold. Let  $h$  be a function of class  $\mathcal{C}^1$  from  $\Theta_T$  to  $H_T$ , with  $\Theta_T$  a sub-interval of  $\Theta$ . Assume there exist finite constants  $C_1$  and  $C_2$  such that for all  $\theta \in \Theta_T$ :*

$$(126) \quad \|h(\theta)\|_T \leq C_1 \quad \text{and} \quad \|h^{[1]}(\theta)\|_T \leq C_2.$$

Let  $w_T$  be an  $H_T$ -valued Gaussian noise such that Assumption 1.1 holds, and consider the Gaussian process  $X = (X(\theta) = \langle h(\theta), w_T \rangle_T, \theta \in \Theta)$ . Then, we have for  $u > 0$ :

$$(127) \quad \mathbb{P} \left( \sup_{\theta \in \Theta_T} |X(\theta)| \geq u \right) \leq c \cdot \left( \sigma \frac{|\Theta_T| \sqrt{\Delta_T}}{u} \vee 1 \right) e^{-u^2 / (4\sigma^2 \Delta_T C_1^2)},$$

where  $|\Theta_T|$  denotes the Riemannian length of the interval  $\Theta_T$  and  $c = 2C_2 + 1$ .

We first state a technical lemma.

**Lemma C.2.** *Let  $I \subset \mathbb{R}$  be an interval. Assume that  $X = (X(\theta), \theta \in I)$  is a real centered Gaussian process with Lipschitz sample paths. Then, for all  $u > 0$  and an arbitrary  $\theta_0 \in I$ , we have:*

$$(128) \quad \mathbb{P} \left( \sup_I X \geq u \right) \leq \frac{1}{u} \int_I \sqrt{\text{Var}(X'(\theta))} e^{-u^2 / (4\text{Var}(X(\theta)))} d\theta + \frac{1}{2} e^{-u^2 / (2\text{Var}(X(\theta_0)))}.$$

**Proof.** We first start with a general remark on Lipschitz functions on  $\mathbb{R}$ . Let  $f$  be a real-valued Lipschitz function defined on an interval  $I \subset \mathbb{R}$ . Let  $b > a$  and set  $f_{a,b} = \min(\max(f, a), b)$ . The function  $f_{a,b}$  is also Lipschitz and, thanks to [29, Theorem 3.3 p107], we get that  $f'_{a,b} = f' = 0$  a.e. on  $\{x \in I : f(x) = a \text{ or } b\}$  and thus  $f'_{a,b} = f' \mathbf{1}_{\{f(x) \in (a,b)\}}$  a.e. on  $I$ . We deduce that:

$$\sup f_{a,b} - \inf f_{a,b} \leq \int_I |f'_{a,b}(x)| dx = \int_I |f'(x)| \mathbf{1}_{\{f(x) \in (a,b)\}} dx.$$

Using this inequality, we obtain that for any  $x_0 \in I$ :

$$(129) \quad \int_a^b \mathbf{1}_{\{\sup_I f > t\}} dt = \int_a^b \mathbf{1}_{\{\sup_I f_{a,b} > t\}} dt = \sup f_{a,b} - a \leq (b - a) \mathbf{1}_{\{f(x_0) \geq a\}} + \int_I |f'(x)| \mathbf{1}_{\{f(x) \in (a,b)\}} dx.$$

Then, applying Inequality (129) to the function  $X$  and taking the expectation, we get, with  $M = \sup_I X$ ,  $a = u > 0$ ,  $b = u + \varepsilon$ ,  $\varepsilon > 0$  and  $x_0 = \theta_0$ :

$$(130) \quad \int_u^{u+\varepsilon} \mathbb{P}(M \geq t) dt \leq \varepsilon \mathbb{P}(X(\theta_0) \geq u) + \int_I \mathbb{E} [|X'(\theta)| \mathbf{1}_{\{u < X(\theta) < u+\varepsilon\}}] d\theta.$$

The random variable  $X(\theta_0)$  is a centered Gaussian variable and therefore we have:

$$(131) \quad \mathbb{P}(X(\theta_0) \geq u) = \int_u^{+\infty} \frac{e^{-x^2 / (2\text{Var}(X(\theta_0)))}}{\sqrt{2\pi \text{Var}(X(\theta_0))}} dx \leq \frac{1}{2} e^{-u^2 / 2\text{Var}(X(\theta_0))},$$

where we used for the inequality that  $\int_u^{+\infty} e^{-t^2} dt \leq \frac{\sqrt{\pi}}{2} e^{-u^2}$  holds for  $u > 0$ , see [1, Formula 7.1.13]. Notice that (131) trivially holds if  $\text{Var}(X(\theta_0)) = 0$  as  $u > 0$ .

We now give a bound of the second term in the right hand-side of (130). Since  $(X', X)$  is also a Gaussian process, we can write:

$$X'(\theta) = \alpha_\theta X(\theta) + \beta_\theta G,$$

where  $G$  is a standard Gaussian random variable independent of  $X(\theta)$  and:

$$\alpha_\theta = \frac{\mathbb{E}[X'(\theta)X(\theta)]}{\text{Var}(X(\theta))} \quad \text{and} \quad \beta_\theta^2 = \text{Var}(X'(\theta)) - \alpha_\theta^2 \text{Var}(X(\theta)),$$

with the convention that  $\alpha_\theta = 0$  if  $\text{Var}(X(\theta)) = 0$ . We get  $|X'(\theta)| \leq |\alpha_\theta X(\theta)| + |\beta_\theta| |G|$ . Since  $G$  is independent of  $X(\theta)$  and  $u > 0$ , we deduce that:

$$\mathbb{E} [|X'(\theta)| \mathbf{1}_{\{u < X(\theta) < u + \varepsilon\}}] \leq \left( |\alpha_\theta|(u + \varepsilon) + \sqrt{\frac{2}{\pi}} |\beta_\theta| \right) \mathbb{P}(u < X(\theta) < u + \varepsilon).$$

Letting  $\varepsilon$  goes to 0 in (130), using (131) the right continuity of the cdf of  $M$  and the monotonicity of the density  $p_{X(\theta)}(u)$  of the law of  $X(\theta)$ , we deduce that:

$$(132) \quad \mathbb{P}(M \geq u) \leq \frac{1}{2} e^{-u^2/2\text{Var}(X(\theta_0))} + \int_I \left( |\alpha_\theta|u + \sqrt{\frac{2}{\pi}} |\beta_\theta| \right) p_{X(\theta)}(u) d\theta,$$

where by convention  $p_{X(\theta)}(u)$  is taken equal to 0 if  $\text{Var}(X(\theta)) = 0$ . We now bound the second term of the right-hand side of (132) in two steps. Using that  $\beta_\theta^2 \leq \text{Var}(X'(\theta))$  and the inequality  $e^{-x^2} \leq e^{-x^2/2} / \sqrt{2} x$  for  $x > 0$ , we get that:

$$(133) \quad \sqrt{\frac{2}{\pi}} |\beta_\theta| p_{X(\theta)}(u) \leq \frac{1}{\pi} \frac{\sqrt{\text{Var}(X'(\theta))}}{u} e^{-u^2/4\text{Var}(X(\theta))}.$$

Thanks to the Cauchy-Schwarz inequality, we get  $|\alpha_\theta| \leq \sqrt{\text{Var}(X'(\theta))} / \sqrt{\text{Var}(X(\theta))}$ . Using also the inequality  $e^{-x^2} \leq 3e^{-x^2/2} / 4x^2$  for  $x > 0$ , we get that:

$$(134) \quad |\alpha_\theta|u p_{X(\theta)}(u) \leq \frac{3}{4} \sqrt{\frac{2}{\pi}} \frac{\sqrt{\text{Var}(X'(\theta))}}{u} e^{-u^2/4\text{Var}(X(\theta))}.$$

Notice that (133) and (134) hold also if  $\text{Var}(X(\theta)) = 0$ . Using that  $\frac{3}{4} \sqrt{\frac{2}{\pi}} + \frac{1}{\pi} \simeq 0.92 \leq 1$ , we deduce (128) from (132), (133) and (134).  $\square$

**Proof of Lemma C.1.** We first consider the case  $\Theta_T = [\theta_0, \theta_1]$  and let  $\gamma : [0, 1] \rightarrow [\theta_0, \theta_1]$  be a minimizing path with respect to the Riemannian metric  $\mathfrak{d}_T$  (see Remark 4.1); in particular we have  $|\gamma'(s)| \sqrt{g_T(\gamma(s))} = \mathfrak{d}_T(\theta_0, \theta_1)$ . Thanks to (123), the Gaussian process  $\tilde{X} = (\tilde{X}(s) = X(\gamma(s)), s \in [0, 1])$  is of class  $\mathcal{C}^1$  on  $s \in [0, 1]$ , with derivative  $\tilde{X}'(s) = \gamma'(s) X'(\gamma(s)) = \gamma'(s) \langle \partial_\theta h(\gamma(s)), w_T \rangle_T$ . Then, according to Lemma C.2, Inequality (128) holds. By Assumption 1.1, we have for all  $\theta \in \Theta_T$ :

$$\text{Var}(X(\theta)) \leq \sigma^2 \Delta_T \|h(\theta)\|_T^2 \leq \sigma^2 \Delta_T C_1^2 \quad \text{and} \quad \frac{\text{Var}(X'(\theta))}{g_T(\theta)} \leq \sigma^2 \Delta_T \|h^{[1]}(\theta)\|_T^2 \leq \sigma^2 \Delta_T C_2^2.$$

Plugging those bounds in Inequality (128) with  $|\gamma'(s)| \sqrt{g_T(\gamma(s))} = \mathfrak{d}_T(\theta_0, \theta_1)$ , we obtain:

$$\begin{aligned} \mathbb{P} \left( \sup_{[\theta_0, \theta_1]} X \geq u \right) &\leq \frac{1}{u} \sqrt{\sigma^2 \Delta_T} C_2 e^{-u^2/(4\sigma^2 \Delta_T C_1^2)} \int_0^1 |\gamma'(s)| \sqrt{g_T(\gamma(s))} ds + \frac{1}{2} e^{-u^2/(2\sigma^2 \Delta_T C_1^2)} \\ &\leq \left( C_2 + \frac{1}{2} \right) \left( \sigma \frac{\mathfrak{d}_T(\theta_0, \theta_1) \sqrt{\Delta_T}}{u} \vee 1 \right) e^{-u^2/(4\sigma^2 \Delta_T C_1^2)}. \end{aligned}$$

Since  $\mathbb{P} \left( \sup_{[\theta_0, \theta_1]} |X| \geq u \right) \leq 2 \mathbb{P} \left( \sup_{[\theta_0, \theta_1]} X \geq u \right)$ , we obtain that (127) holds for  $\Theta_T$  a bounded closed interval. Then, use monotone convergence and the continuity of  $X$  to get (127) for any interval  $\Theta_T$ .  $\square$

### C.3. Schur complement

The following Lemma is a classical result on the Schur complement.

**Lemma C.3** (Schur complement). *Let  $M \in \mathbb{R}^{n \times n}$  be a matrix composed of blocks  $A \in \mathbb{R}^{(n-k) \times (n-k)}$ ,  $B \in \mathbb{R}^{(n-k) \times k}$ ,  $C \in \mathbb{R}^{k \times (n-k)}$ ,  $D \in \mathbb{R}^{k \times k}$ .*

$$M = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$$

Assume that  $D$  and  $S_1 = A - BD^{-1}C$  are non singular. Then, the system:

$$(135) \quad M \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} a \\ b \end{pmatrix}.$$

with  $x \in \mathbb{R}^{n-k}$ ,  $y \in \mathbb{R}^k$ ,  $a \in \mathbb{R}^{n-k}$  and  $b \in \mathbb{R}^k$ , has a unique solution given by:

$$x = S_1^{-1}a - S_1^{-1}BD^{-1}b \quad \text{and} \quad y = D^{-1}b - D^{-1}CS_1^{-1}a + D^{-1}CS_1^{-1}BD^{-1}b.$$

#### C.4. Proofs of Lemmas in Section 4

**Proof of Lemma 4.2.** For simplicity, we remove the subscript  $\mathcal{K}$  and for example write  $f^{[1]} = \tilde{D}_1[f] = D_1[f]/\sqrt{g}$ . Recall that  $G$ , a primitive of  $\sqrt{g}$ , is continuous increasing and thus induces a one-to-one map from  $\Theta$  to its image. Following Remark 4.1, we consider the minimizing path  $\gamma : [0, 1] \rightarrow \Theta$  from  $\theta_0$  to  $\theta$  defined by  $\gamma_s = G^{-1}(as + b)$ , with  $b = G(\theta_0)$  and  $a = G(\theta) - G(\theta_0)$ . Thus, we have  $\mathcal{L}(\gamma) = \mathfrak{d}(\theta, \theta_0)$ . The minimizing path from  $\theta_0$  to  $\theta$  has constant speed thus equal to  $\mathfrak{d}(\theta_0, \theta)$ . From the explicit expression of  $\gamma$ , we get in fact that  $\dot{\gamma}_t \sqrt{g(\gamma_t)} = A$  for  $t \in [0, 1]$ , where  $A = \text{sign}(\theta - \theta_0) \mathfrak{d}(\theta, \theta_0)$ . Thus, we have:

$$(136) \quad f(\theta) - f(\theta_0) = f(\gamma_1) - f(\gamma_0) = \int_0^1 \dot{\gamma}_t f'(\gamma_t) dt = A \int_0^1 \tilde{D}_1[f](\gamma_t) dt = A \int_0^1 f^{[1]}(\gamma_t) dt,$$

where we used (124) and that the derivative of  $f \circ \gamma_t$  is  $\dot{\gamma}_t f' \circ \gamma_t$  for the second equality and the definition of  $\tilde{D}_1[f]$  as well as the equality  $\dot{\gamma}_t \sqrt{g(\gamma_t)} = A$  for the last.

Using (136) for  $f$  and  $\theta$  replaced by  $f^{[1]}$  and  $\gamma(t)$  for some  $t \in [0, 1]$ , we get thanks to (23) that:

$$f^{[1]}(\gamma_t) = f^{[1]}(\theta_0) + \tilde{A} \int_0^1 f^{[2]}(\tilde{\gamma}_s) ds,$$

where  $\tilde{\gamma}$  is a geodesic from  $\theta_0$  to  $\gamma_t$  and  $\tilde{A} = \dot{\gamma}_s \sqrt{g(\tilde{\gamma}_s)}$ . Since  $\gamma$  is itself a geodesic, we deduce that  $\tilde{\gamma}_s = \gamma_{st}$ , and thus  $\tilde{A} = tA$ . Plugging this in (136), we get:

$$f(\theta) - f(\theta_0) = A f^{[1]}(\theta_0) + A^2 \int_{[0,1]^2} f^{[2]}(\gamma_{st}) t dt ds = A f^{[1]}(\theta_0) + A^2 \int_0^1 (1-r) f^{[2]}(\gamma_r) dr.$$

This gives (24). □

**Proof of Lemma 4.3.** Recall that by Assumption 3.2 the function  $\phi_T$  is  $\mathcal{C}^3$ . According to (123), we have that for any  $i, j \in \{0, \dots, 3\}$  and any  $\theta, \theta' \in \Theta$ :

$$(137) \quad \partial_{\theta, \theta'}^{i,j} \langle \phi_T(\theta), \phi_T(\theta') \rangle_T = \left\langle \partial_{\theta}^i \phi_T(\theta), \partial_{\theta'}^j \phi_T(\theta') \right\rangle_T.$$

This and (20), (22), (25) and (26) readily imply (31). The first equality of (32) comes from Cauchy-Schwarz's inequality. The second is clear. We also have:

$$(138) \quad \langle \partial_{\theta} \phi_T(\theta), \phi_T(\theta) \rangle_T = \frac{1}{2} \partial_{\theta} \|\phi_T(\theta)\|^2 = 0$$

Since the right hand-side is also equal to  $\sqrt{g_T(\theta)} \mathcal{K}_T^{[1,0]}(\theta, \theta)$  thanks to (31), we get the third equality of (32). Taking the derivative with respect to  $\theta$  in (138) yields  $g_T(\theta) = \langle \partial_{\theta} \phi_T(\theta), \partial_{\theta} \phi_T(\theta) \rangle = -\langle \partial_{\theta}^2 \phi_T(\theta), \phi_T(\theta) \rangle$ . Thanks to (21),



we get  $\partial_\theta^2 \phi_T = g_T \tilde{D}_{2,T}[\phi_T] + (1/2g_T)g'_T \partial_\theta \phi_T$ . Using (31) and (138) again, we deduce that  $\langle \partial_\theta^2 \phi_T(\theta), \phi_T(\theta) \rangle = g_T(\theta) \mathcal{K}_T^{[2,0]}(\theta, \theta)$ . This gives the fourth equality of (32). Eventually, we deduce from (31), (21) and (22) that:

$$g_T(\theta)^{3/2} \mathcal{K}_T^{[2,1]}(\theta, \theta) = \langle \partial_\theta^2 \phi_T(\theta), \partial_\theta \phi_T(\theta) \rangle - \frac{1}{2} \frac{g'_T(\theta)}{g_T(\theta)} \langle \partial_\theta \phi_T(\theta), \partial_\theta \phi_T(\theta) \rangle.$$

Then, use that  $g'_T(\theta) = 2\langle \partial_\theta^2 \phi_T(\theta), \partial_\theta \phi_T(\theta) \rangle$  to deduce that  $\mathcal{K}_T^{[2,1]}(\theta, \theta) = 0$ .  $\square$

### C.5. Control on $f$ from its derivatives $f^{[2]}$

The proof of the next lemma is similar to the proof of [41, Lemma 2] and is left to the reader. Recall from (32) that  $\mathcal{K}_T^{[2,0]}(\theta, \theta) = -1$  on  $\Theta$ .

**Lemma C.4.** *Suppose Assumptions 3.1 and 3.2 on the dictionary hold. Let  $f$  be a real valued function defined on an interval  $\Theta$  of class  $\mathcal{C}^2$ . Let  $\theta_0 \in \Theta$ . Set for  $i = 1, 2$ ,  $f^{[i]} = \tilde{D}_{i,T}[f]$  (see (22)).*

(i) *Assume  $f(\theta_0) = 0$ ,  $f^{[1]}(\theta_0) = 0$  and that there exist  $\delta > 0$  and  $r > 0$  such that for any  $\theta \in \mathcal{B}_T(\theta_0, r)$ :*

$$(139) \quad |f^{[2]}(\theta)| \leq 2\delta.$$

*Then, we have  $|f(\theta)| \leq \delta \mathfrak{d}_T(\theta, \theta_0)^2$ , for any  $\theta \in \mathcal{B}_T(\theta_0, r)$ .*

(ii) *Let  $\Theta_T \subset \Theta$  be an interval and suppose that  $L \geq \sup_{\Theta_T^2} |\mathcal{K}_T^{[2,0]}|$  is finite and there exist  $\varepsilon > 0$  and  $r \in (0, L^{-\frac{1}{2}})$  such that for any  $\theta \in \mathcal{B}_T(\theta_0, r)$ ,  $-\mathcal{K}_T^{[2,0]}(\theta, \theta_0) \geq \varepsilon$ . Assume that  $\mathcal{B}_T(\theta_0, r) \subset \Theta_T$ ,  $f(\theta_0) = v \in \{-1; 1\}$ ,  $f^{[1]}(\theta_0) = 0$  and that there exists  $\delta \in (0, \varepsilon)$  such that for any  $\theta \in \mathcal{B}_T(\theta_0, r)$ :*

$$(140) \quad |f^{[2]}(\theta) - v \mathcal{K}_T^{[2,0]}(\theta, \theta_0)| \leq \delta.$$

*Then, we have  $|f(\theta)| \leq 1 - \frac{(\varepsilon - \delta)}{2} \mathfrak{d}_T(\theta, \theta_0)^2$ , for any  $\theta \in \mathcal{B}_T(\theta_0, r)$ .*

### C.6. Proof of Lemma 8.1

We keep the notations from Section 8.1. In order to prove that the constants  $c_0$ ,  $c_1$  and  $c_2$  do not depend on the scaling factor  $\sigma_0$ , we shall rewrite  $\rho_T$  and  $\mathcal{V}_T$  defined in (35) and (37) using a change of scale. To do so, we define  $\varphi^0(\theta) = k(\cdot - \theta)$  for  $\theta \in \Theta$ ; the grid  $t_1^0, \dots, t_T^0$  where  $t_j^0 = t_j/\sigma_0$ ; the Hilbert space  $L^2(\lambda_T^0)$  with  $\lambda_T^0 = \Delta_T \sigma_0^{-1} \sum_{j=1}^T \delta_{t_j^0}$ , endowed with its natural scalar product noted  $\langle \cdot, \cdot \rangle_{\lambda_T^0}$  and norm  $\|\cdot\|_{\lambda_T^0}$ ; the parameter space  $\Theta_T^0 = [a_T(1 - \epsilon)\sigma_0^{-1}, b_T(1 - \epsilon)\sigma_0^{-1}]$ . Since the scaling factor  $\sigma_0$  is fixed, the measures  $(\lambda_T^0, T \geq 2)$  converge vaguely towards the Lebesgue measure  $\lambda_\infty$  on  $\mathbb{R}$ . We shall also consider another kernel:

$$\mathcal{K}_T^0(\theta, \theta') = \langle \phi_T^0(\theta), \phi_T^0(\theta') \rangle_{\lambda_T^0} \quad \text{with} \quad \phi_T^0 = \varphi^0 / \|\varphi^0\|_{\lambda_T^0},$$

and the limit kernel  $\mathcal{K}_\infty^0(\theta, \theta') = \langle \phi_\infty^0(\theta), \phi_\infty^0(\theta') \rangle_\infty$  with  $\phi_\infty^0 = \varphi^0 / \|\varphi^0\|_\infty$ . For any  $T \in \mathbb{N} \cup \{+\infty\}$ , the kernel  $\mathcal{K}_T^0$  is of class  $\mathcal{C}^{3,3}$  on  $\Theta^2$  and for  $i, j \in \{0, \dots, 3\}$  and  $\theta, \theta' \in \Theta$ , we have:

$$\mathcal{K}_T^{[i,j]}(\theta, \theta') = \mathcal{K}_T^{0[i,j]} \left( \frac{\theta}{\sigma_0}, \frac{\theta'}{\sigma_0} \right) \quad \text{and} \quad \frac{1}{\sigma_0^2} g_{\mathcal{K}_T^0} \left( \frac{\theta}{\sigma_0} \right) = g_{\mathcal{K}_T}(\theta).$$

We can now rewrite  $\rho_T$  and  $\mathcal{V}_T$  by using a change of scale and we get:

$$\rho_T = \max \left( \sup_{\Theta_T^0} \sqrt{\frac{g_{\mathcal{K}_T^0}}{g_{\mathcal{K}_\infty^0}}}, \sup_{\Theta_T^0} \sqrt{\frac{g_{\mathcal{K}_\infty^0}}{g_{\mathcal{K}_T^0}}} \right),$$

and

$$\mathcal{V}_T = \max(\mathcal{V}_T^{(1)}, \mathcal{V}_T^{(2)}) \quad \text{with} \quad \mathcal{V}_T^{(1)} = \max_{i,j \in \{0,1,2\}} \sup_{(\Theta_T^0)^2} |\mathcal{K}_T^{0[i,j]} - \mathcal{K}_\infty^{0[i,j]}| \quad \text{and} \quad \mathcal{V}_T^{(2)} = \sup_{\Theta_T^0} |h_{\mathcal{K}_T^0} - h_{\mathcal{K}_\infty^0}|.$$

Thus, bounding  $\rho_T$  and  $\mathcal{V}_T$  amounts to controlling the proximity between the kernels  $\mathcal{K}_T^0$  and  $\mathcal{K}_\infty^0$ .

First, we provide an upper bound for any  $i, j \in \{0, \dots, 3\}$  of:

$$(141) \quad B_{i,j}(T) = \sup_{\theta, \theta' \in \Theta_T^0} \left| \left\langle \partial_\theta^i \varphi^0(\theta), \partial_{\theta'}^j \varphi^0(\theta') \right\rangle_{\lambda_T^0} - \left\langle \partial_\theta^i \varphi^0(\theta), \partial_{\theta'}^j \varphi^0(\theta') \right\rangle_\infty \right|.$$

Notice that:

$$\partial_\theta^i \partial_t^j \varphi^0(\theta, t) = (-1)^j k^{(i+j)}(\theta - t).$$

In what follows, we shall use at most three derivatives in  $\theta$  and one derivative in  $t$ , so that  $i + j \leq 4$  in the above formula. Recall the polynomials  $P_i$  are defined as  $k^{(i)} = P_i k$  and set  $M = \max_{0 \leq i \leq 4} \sup |P_i| \sqrt{k}$ . It is elementary to get that for  $\theta, \theta' \in \mathbb{R}$ :

$$\left| (\Delta_T / \sigma_0) \sum_{k=1}^T \partial_\theta^i \varphi^0(\theta, t_k^0) \partial_{\theta'}^j \varphi^0(\theta', t_k^0) - \int_{a_T / \sigma_0}^{b_T / \sigma_0} \partial_\theta^i \varphi^0(\theta, t) \partial_{\theta'}^j \varphi^0(\theta', t) dt \right| \leq 4\sqrt{\pi} \Delta_T M^2 \sigma_0^{-1}.$$

We have for  $\theta, \theta' \in \Theta_T^0$  that:

$$\begin{aligned} \left| \int_{\mathbb{R} \setminus [a_T / \sigma_0, b_T / \sigma_0]} \partial_\theta^i \varphi^0(\theta, t) \partial_{\theta'}^j \varphi^0(\theta', t) dt \right| &\leq \left| \int_{b_T / \sigma_0}^{+\infty} \partial_\theta^i \varphi^0(\theta, t) \partial_{\theta'}^j \varphi^0(\theta', t) dt \right| + \left| \int_{-\infty}^{a_T / \sigma_0} \partial_\theta^i \varphi^0(\theta, t) \partial_{\theta'}^j \varphi^0(\theta', t) dt \right| \\ &\leq 2M^2 \int_{\epsilon b_T / \sigma_0}^{+\infty} k(t) dt \\ &\leq 2\sqrt{\pi} M^2 e^{-\epsilon^2 b_T^2 / 2\sigma_0^2}, \end{aligned}$$

where we used that  $2 \int_u^{+\infty} e^{-t^2} dt \leq \sqrt{\pi} e^{-u^2}$  for  $u > 0$ , see formula 7.1.13 in [1]. We deduce that:

$$B_{i,j}(T) \leq 4\sqrt{\pi} \Delta_T M^2 \sigma_0^{-1} + 2\sqrt{\pi} M^2 e^{-\epsilon^2 b_T^2 / 2\sigma_0^2} \leq 2\sqrt{\pi} M^2 \gamma_T,$$

with  $\gamma_T = 2\Delta_T \sigma_0^{-1} + \sqrt{\pi} e^{-\epsilon^2 b_T^2 / 2\sigma_0^2}$ .

Similar arguments as above yield that:

$$\sup_{\theta \in \Theta_T^0} \left| \|\varphi^0(\theta)\|_{\lambda_T^0}^2 - \|\varphi^0(\theta)\|_\infty^2 \right| \leq \gamma_T.$$

so that  $\|\varphi^0(\theta)\|_{\lambda_T^0}^2 \geq \sqrt{\pi} - \gamma_T$  for all  $\theta \in \Theta_T^0$ . It is then easy to deduce that  $\sup_{\theta \in \Theta_T^0} |g_{\mathcal{K}_T^0} - g_{\mathcal{K}_\infty^0}|$  is bounded by a constant times  $\gamma_T$  when  $\gamma_T$  is smaller than a universal finite constant. Up to taking  $\gamma_T$  smaller than some universal finite constant, this and the fact that  $g_{\mathcal{K}_\infty^0} = 1/2$  give the second part of (51). Then use formulae for the derivatives of the kernels, see (29) and (22), to get the first part of (51).