



**HAL**  
open science

## Off-the-grid prediction and testing for linear combination of translated features

Cristina Butucea, Jean-François Delmas, Anne Dutfoy, Clément Hardy

► **To cite this version:**

Cristina Butucea, Jean-François Delmas, Anne Dutfoy, Clément Hardy. Off-the-grid prediction and testing for linear combination of translated features. 2024. hal-03880134v2

**HAL Id: hal-03880134**

**<https://hal.science/hal-03880134v2>**

Preprint submitted on 18 Jul 2024

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

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

# Off-the-grid prediction and testing for linear combination of translated features

Cristina Butucea<sup>1</sup>, Jean-François Delmas<sup>2</sup>  
Anne Duffoy<sup>3</sup> and Clément Hardy<sup>2</sup>

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

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

<sup>3</sup>*EDF R&D, Palaiseau, France, e-mail: [anne.duffoy@edf.fr](mailto:anne.duffoy@edf.fr)*

**Abstract:** We consider a model where a signal (discrete or continuous) is observed with an additive Gaussian noise process. The signal is issued from a linear combination of a finite but increasing number of translated features. The features are continuously parameterized by their location and depend on some scale parameter. First, we extend previous prediction results for off-the-grid estimators by taking into account here that the scale parameter may vary. The prediction bounds are analogous, but we improve the minimal distance between two consecutive features locations in order to achieve these bounds.

Next, we propose a goodness-of-fit test for the model and give non-asymptotic upper bounds of the testing risk and of the minimax separation rate between two distinguishable signals. In particular, our test encompasses the signal detection framework. We deduce upper bounds on the minimal energy, expressed as the  $\ell_2$ -norm of the linear coefficients, to successfully detect a signal in presence of noise. The general model considered in this paper is a non-linear extension of the classical high-dimensional regression model. It turns out that, in this framework, our upper bound on the minimax separation rate matches (up to a logarithmic factor) the lower bound on the minimax separation rate for signal detection in the high-dimensional linear model associated to a fixed dictionary of features. We also propose a procedure to test whether the features of the observed signal belong to a given finite collection under the assumption that the linear coefficients may vary, but have prescribed signs under the null hypothesis. A non-asymptotic upper bound on the testing risk is given.

We illustrate our results on the spikes deconvolution model with Gaussian features on the real line and with the Dirichlet kernel, frequently used in the compressed sensing literature, on the torus.

**MSC2020 subject classifications:** Primary 62G05, 62G10; secondary 62G08.

**Keywords and phrases:** Goodness-of-fit testing, Mixture model, Non-linear regression model, Non-parametric hypotheses testing, Off-the-grid methods, Spikes deconvolution.

## 1. Introduction

In many fields, a signal of interest can be described as a linear combination of shifted source signals having the same shape. Thus, the source signal is supposed

to belong to a parametric set of functions (for example, Gaussian, Cauchy or sinusoidal-shaped functions) parameterized by its location parameter. The signal is observed with an additive noise process in discrete or continuous time. We assume that the noise and the observation space can vary with some parameter  $T$  increasing with the quality of the observations.

For example, the chemical analysis of a material is done through spectroscopy and each chemical component is represented by a spiked Gaussian-shaped signal located at some prescribed frequency, see [5]. The final signal is a linear combination of such spikes. In multiple source detection, sound or image may present a similar structure.

More general non-linear models (not necessarily location models) for the features have been discussed in [6], and the particular case of location families has been discussed in Section 8 therein. However, we allow here the features to depend on a scale parameter which varies with  $T$ . This makes the proof technique very different from the previous one.

We are interested in estimating both the coefficients of the linear combination and the location parameters of the different features appearing in the signal. We give sufficient conditions in order to obtain upper bounds for the quadratic prediction risk of the same order as if the non-linear parameters were known. We show that these sufficient conditions are milder than those in [6] without loosing on the prediction risk bounds.

We are also interested in testing problems. First, we want to test whether the observations are issued from a given linear combination of features. We remark that it includes the case of signal detection. This test problem finds an application in spectroscopy to detect the presence of a chemical compound in a material. Finally, we are interested in testing whether the observed signal is a linear combination of features located at a prescribed list of values with linear coefficients having prescribed signs under the null hypothesis. This is of interest in spectroscopy: in a material we expect a list of chemical components. This test problem detects ageing or important damage to the material which can be detected if unexpected chemical components are present.

### 1.1. Model

Let  $T \in \mathbb{N}$ . We observe a random element  $y$  in the Hilbert space  $L^2(\lambda_T)$  of square integrable functions with respect to the measure  $\lambda_T$  on the Borel  $\sigma$ -field of some metric space. The observation is the sum of a deterministic signal and a noise process  $w_T$  in  $L^2(\lambda_T)$ . We assume that the signal is an unknown linear combination of a finite unknown number  $s$  of features belonging to a continuously parameterized subfamily  $(\varphi_T(\theta), \theta \in \Theta)$  of  $L^2(\lambda_T)$ . We call this family a continuous dictionary, the weights of the linear combination - the linear coefficients, and the parameters of the features - the non-linear parameters. Moreover, we assume that the noise is a Gaussian random process. Thus, the general model is fully specified by the choice of the Hilbert space of our observation, of the continuous dictionary of features and of the noise process.

The Hilbert space  $L^2(\lambda_T)$  is endowed with the natural scalar product noted  $\langle \cdot, \cdot \rangle_{L^2(\lambda_T)}$  and norm  $\|\cdot\|_{L^2(\lambda_T)}$ . Let us define the normalized function  $\phi_T$  defined on  $\Theta$  by:

$$\phi_T(\theta) = \varphi_T(\theta) / \|\varphi_T(\theta)\|_{L^2(\lambda_T)}. \quad (1)$$

We assume that the signal is a linear combination with unknown non-zero linear coefficients  $\beta^* = (\beta_1^*, \dots, \beta_s^*)$  in  $(\mathbb{R}^*)^s$  of an unknown number  $s \in \mathbb{N}$  of active features with unknown distinct non-linear parameters  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta^s$ . We use the notation  $\mathbb{R}^* = \mathbb{R} \setminus \{0\}$ .

Thus, we observe  $y$  in the model:

$$y = \sum_{k=1}^s \beta_k^* \Phi_T(\theta_k^*) + w_T \quad \text{in } L^2(\lambda_T). \quad (2)$$

Let us define the multivariate function  $\Phi_T$  on  $\Theta^s$  by:

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

Model (2) writes

$$y = \beta^* \Phi_T(\vartheta^*) + w_T \quad \text{in } L^2(\lambda_T).$$

When  $s = 0$ , we set by convention that  $\beta^* \Phi_T(\vartheta^*) = 0$  as well as  $A^s = \{0\}$  for any set  $A$ . We denote by  $\mathcal{Q}^* = \{\theta_\ell^*, 1 \leq \ell \leq s\}$  the set of the non-linear parameters associated to active features.

In this paper we consider a dictionary given by a one dimensional location model scaled with a given  $\sigma_T > 0$ :

$$\left( \varphi_T(\theta) = h(\theta - \cdot, \sigma_T), \theta \in \Theta \right) \quad (3)$$

where the set  $\Theta$  is the real line  $\mathbb{R}$  or the torus  $\mathbb{R}/\mathbb{Z}$ , the real-valued function  $h$  is defined on  $\Theta \times \mathfrak{S}$ , smooth with respect to its first variable and normalized so that  $\|h(\cdot, \sigma_T)\|_{L^2(\text{Leb})} = 1$ , and  $\sigma_T$  is an element of the set  $\mathfrak{S}$  of admissible positive scale parameter values. Note that  $\varphi_T$  depends on  $T$  only through the argument  $\sigma_T$ . See Section 2.1 for examples of functions  $h$  including the Gaussian scaled-spikes and the low-pass filter.

The process  $y$  is observed over the support of the measure  $\lambda_T$ . Therefore it is legitimate to consider models whose location parameters belong to the smallest interval covering the support of the measure  $\lambda_T$ . Hence, we introduce the set  $\Theta_T$ , a compact interval of  $\Theta$  (when  $\Theta$  is the torus, then we can take  $\Theta_T = \Theta$ ), and we shall assume that  $\mathcal{Q}^*$  is a subset of  $\Theta_T$ . We denote by  $|\Theta_T|$  the Euclidean diameter of the set  $\Theta_T$ .

We consider a large variety of Gaussian noise processes. Indeed, we only assume the following mild assumption on  $w_T$ , where the decay rate  $\Delta_T > 0$  controls the noise variance decay as the parameter  $T$  grows and  $\bar{\sigma} > 0$  is the intrinsic noise level. A wide range of noise processes satisfy our assumptions, see Section 2.2; they can be discrete or continuous, white or coloured under these constraints.

**Assumption 1.1** (Admissible noise). *Let  $T \in \mathbb{N}$ . The Gaussian noise process  $w_T$  satisfies  $\mathbb{E} \left[ \|w_T\|_{L^2(\lambda_T)}^4 \right] < +\infty$ , and there exist a noise level  $\bar{\sigma} > 0$  and a decay rate  $\Delta_T > 0$  such that for all  $f \in L^2(\lambda_T)$ , the random variable  $\langle f, w_T \rangle_{L^2(\lambda_T)}$  is a centered Gaussian random variable satisfying:*

$$\text{Var} \left( \langle f, w_T \rangle_{L^2(\lambda_T)} \right) \leq \bar{\sigma}^2 \Delta_T \|f\|_{L^2(\lambda_T)}^2. \quad (4)$$

We assume that the quantity  $\mathbb{E} \left[ \|w_T\|_{L^2(\lambda_T)}^2 \right]$  is known for the considered models. Using Cauchy-Schwarz inequality, we get:

$$\text{Var} \left( \langle f, w_T \rangle_{L^2(\lambda_T)} \right) \leq \mathbb{E} \left[ \|w_T\|_{L^2(\lambda_T)}^2 \right] \|f\|_{L^2(\lambda_T)}^2, \quad (5)$$

which is in some examples not as sharp as (4), see Section 2.2.2. We shall also consider the finite variance of the squared norm of the noise:

$$\Xi_T = \text{Var} \left( \|w_T\|_{L^2(\lambda_T)}^2 \right). \quad (6)$$

To sum up, the quality of the information provided by our observation  $y$  depends on the support of the measure  $\lambda_T$  and on the noise  $w_T$  through  $\Delta_T$ . It increases with the parameter  $T$ . Due to the particular form of the features, we refer to our model as a Linear combination of translation features (LCTF-model).

In this paper, we are interested both in building estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  of the parameters  $\beta^*$  and  $\vartheta^*$ , respectively, and in hypothesis testing problems concerning our model. Our goal is two-fold: on the one hand, we attain best known non asymptotic prediction bounds for the risk measure:

$$\|\hat{\beta}\Phi_T(\hat{\vartheta}) - \beta^*\Phi_T(\vartheta^*)\|_{L^2(\lambda_T)}$$

under less restrictive conditions than previous works. Moreover, we use the certificate functions designed as tools in these proofs in order to build test procedures in our model that generalize the signal detection problem in a linear regression model. On the other hand, we treat the goodness-of-fit test problem and then, the more general problem of testing whether the signal in our observation presents only features included in a prescribed list, with associated linear coefficients that may vary but cannot change signs.

## 1.2. Previous work

Estimating the linear coefficients and the parameters of model (2) from an observation  $y$  has attracted a lot of attention over the past decade. A major contribution in this field comes from the formulation of the BLasso problem in [10]. This optimization problem on a space of measures allows to estimate both linear coefficients and non-linear parameters without using a grid on the parameter space. This off-the-grid method has successfully been used in [8] and

[7] in the context of super-resolution as well as in [11] for spikes deconvolution. High probability bounds for the prediction error have been given in [20], [19] and [4] for the specific dictionary of complex exponential functions continuously parameterized by their frequencies and more recently in [6] for a wide range of dictionaries parameterized over a one-dimensional space. These results are based on certificate functions whose existence have been proven in a very general framework in [18] provided that the non-linear parameters of the mixture are well-separated with respect to a Riemannian metric.

Goodness-of-fit tests are used to check whether observations are indeed derived from a given statistical model. We refer to the monograph [14] for a comprehensive presentation of goodness-of-fit testing. When we consider a finite dictionary of features  $(\varphi_T(\theta), \theta \in \mathcal{Q})$  with  $\mathcal{Q}$  a known finite subset of  $\Theta$ , the model (2) can be rewritten as a linear regression model, possibly of high dimension depending on the size of the finite dictionary  $p := \text{Card}(\mathcal{Q})$ . In this case, testing the goodness-of-fit of the model amounts to testing whether the linear coefficients in the mixture are equal to some given linear coefficients. When the dictionary is known, the testing problem is homogeneous in the linear coefficients  $\beta$  and is therefore equivalent to testing  $\beta \equiv 0$ , which is a signal detection problem.

Signal detection has raised a lot of interest over the past decades. It is well known that the alternative hypothesis  $H_1$  (presence of signal) must be well separated from the null hypothesis  $H_0$  (only noise) in order to have tests with small risks. The separation can be seen as a minimal signal intensity allowing the detection. Then, it is a matter of interest to evaluate the minimax separation rate, i.e., the smallest separation that allows to distinguish the tested hypotheses. In [12], asymptotic rates for the minimax separation in the framework of signal detection are derived for the non-parametric Gaussian white noise model. Non-asymptotic rates were then derived in [3] and later in [16] to tackle the case of heterogeneous variances. We refer to the monograph [13] for an overview of non-parametric hypotheses testing. Regarding the high dimensional regression model where the observation is of dimension  $T$  and the dictionary is fixed, known and of size  $p$ , the work of [15] established the following asymptotic minimax separation rates under coherence assumptions on the dictionary:

$$\frac{1}{T^{\frac{1}{4}}} \wedge \sqrt{\frac{s}{T} \log(p)} \wedge \frac{p^{\frac{1}{4}}}{\sqrt{T}}.$$

The signal intensity is expressed by the  $\ell_2$ -norm of the linear coefficients. Their lower bounds on the asymptotic minimax separation stand for both fixed and random designs whereas their upper bounds stand for random designs. The work of [2] does not tackle the high dimension but provides tests achieving the minimax separation for fixed designs under coherence assumptions on the dictionary. We note that the existing results do not apply to our context.

For the non-linear extension of linear regression models that we consider here, goodness-of-fit testing does not reduce to signal detection as the mixture is not

homogeneous with respect to the non-linear parameters. Therefore, we introduce new testing procedures. We stress that one of the test statistics is not derived from estimators of the linear coefficients. In fact, depending on the sparsity of the signal, the dimension of the observation and the size of the dictionary, plug-in methods using sparse estimators might not be the best way to proceed. They do not always lead to the minimal separation. In this sense, testing is a very different statistical problem from estimation.

### 1.3. Description of the results

The aim of this paper is twofold. First, we improve on [6] in the case of linear combination of translated features by giving bounds on the prediction error under milder separation constraints between the unknown non-linear parameters in  $\mathcal{Q}^*$ . Indeed, the sufficient separation conditions between two neighboring non-linear parameters are difficult to track explicitly. In all generality, they can be rather restrictive and scale with a factor  $s$  for arbitrary dictionaries satisfying the conditions. In the particular case of Gaussian-shaped features, more explicit calculations are possible and the minimal separation reduces to some constant value.

In this paper, due to the shape of our dictionary of features, *i.e.* a location model scaled by some  $\sigma_T$ , we get more explicit sufficient separation conditions which are less restrictive. This is achieved by taking the scale parameter of the features  $\sigma_T$  into account. In particular, in the case of Gaussian-shaped features, the minimal separation is of order  $\sigma_T$ . Intuitively, this can be explained by the fact that for peaked features (with small scaling parameter  $\sigma_T$ ) we may distinguish spikes located at smaller (by a factor  $\sigma_T$ ) distance.

The second goal of this paper is to study hypotheses testing problems in these models. We give procedures for the goodness-of-fit of the mixture model in order to determine whether the unknown signal  $\beta^* \Phi_T(\vartheta^*)$  is equal to a reference signal  $\beta^0 \Phi_T(\vartheta^0)$  for some known vectors  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and  $\vartheta^0 \in \Theta_T^{s^0}$ . Under our assumptions, the model is identifiable, thus the null hypothesis is equivalent to testing that  $\beta^*, \vartheta^*$  coincide with  $\beta^0, \vartheta^0$  up to a permutation. This setup includes the case of signal detection where the null hypothesis is  $\beta^* \equiv 0$ , that is  $s = 0$ . On this aspect, our minimal intensity rates allowing signal detection are similar up to a log factor to the rates obtained in [15] for high dimensional linear models. We propose a combined procedure based on differences between the reference signal  $\beta^0 \Phi_T(\vartheta^0)$  and either the observation  $y$  or a reconstructed signal obtained from estimators of the model parameters. In order to successfully perform the test, we remove from the alternative hypothesis the signals whose proximity to the reference signal  $\beta^0 \Phi_T(\vartheta^0)$  is below some separation parameter, with respect to the norm  $\|\cdot\|_{L^2(\lambda_T)}$ . We give a non-asymptotic upper bound of the testing risk and deduce an upper bound on the minimal separation needed to distinguish two different signals. This upper bound yields two regimes according to the test procedures that we define and study. In the case of signal detection, the separation can be expressed as the  $\ell_2$ -norm of the linear coefficients of the

observed mixture. In particular, when the observation  $y$  is issued from a non-linear extension of the classical high-dimensional regression model, our upper bound matches (up to logarithmic factors) the asymptotic lower bound of the minimal separation needed to distinguish two signals that are mixture of features from a finite high-dimensional dictionary.

Moreover, we test the presence of at most  $s_0$  prescribed features in the mixture with arbitrary linear coefficients of given sign. That is, we test whether for each  $\epsilon \in \{+, -\}$  the unknown set  $\mathcal{Q}^{\star, \epsilon} = \{\theta_k^* \in \mathcal{Q}^* : \epsilon \beta_k^* > 0\}$  is a subset of  $\mathcal{Q}^{0, \epsilon}$ , with  $\mathcal{Q}^{0, +}$  and  $\mathcal{Q}^{0, -}$  being given disjoint finite subsets of  $\Theta_T$ . This setup is issued from an application to spectroscopy (see [5]), where the presence of other chemical components than the prescribed ones are indicating ageing or substantial modifications of the analyzed material. To separate the null hypothesis from the alternative hypothesis, we introduce a discrepancy that is 0 if and only if the parameters  $(\beta^*, \vartheta^*)$  belong to the null hypothesis. We give an upper bound on the minimal separation to successfully perform our test. The test statistic introduced and studied in this context makes explicit use of the construction of certificates used in compressed sensing [9, 20, 18], super resolution [8], spikes deconvolution [11], as well as in [6, 19, 4] for establishing the prediction rates of the estimators of  $(\beta^*, \vartheta^*)$ . We stress the fact that the test statistic is not an estimator of the discrepancy measure separating the null and the alternative hypotheses, as is usually the case in non-parametric tests.

#### 1.4. Roadmap of the paper

Section 2 gives several possible specific choices in our general model by showing examples of dictionaries of features, of observation spaces and of Gaussian processes (white or coloured under our assumptions). In Section 3, we start by presenting the assumptions needed to perform a successful estimation of the linear coefficients and location parameters of our model. After giving a prediction bound in Theorem 3.5, we show in Lemma 3.3 that the required assumptions are sufficient conditions for the identifiability of the model. In Section 4, we test whether the observation derives from a given mixture or from some other mixture sufficiently separated from the latter. We give in Theorems 4.1 and 4.3 bounds of the testing risks associated to two different test procedures. We show in Corollaries 4.2 and 4.5 that these two tests give two regimes for our upper bound on the minimal separation to distinguish two different signals from an observation contaminated by noise. We also provide a discussion on the comparison of our upper bounds with some existing lower bounds. In Section 5, we propose a procedure to test whether the active features in the observed signal belong to a given finite collection with linear coefficients of prescribed signs. Both hypotheses of this test problem are composite and a new measure of the separation between these hypotheses has been introduced. The proposed test relies on the certificates used in the proof of the prediction bounds in an original way. A bound of the testing risk is given in Theorem 5.2 and in Corollary 5.3, we provide an upper bound on the minimax separation rate. The examples



of Gaussian scaled spikes deconvolution on  $\mathbb{R}$  and low-pass filter on  $\mathbb{R}/\mathbb{Z}$  are addressed in Sections 6 and 7. Some proofs can be found in Section 8.

## 2. Specific models covered by our general model

We consider a large variety of models: discrete models where the process  $y = (y(t_1), \dots, y(t_T))$  is observed on a finite grid  $t_1 < \dots < t_T$  or continuous models where the process  $y = y(t)$  is observed on a continuous interval.

### 2.1. Examples of feature functions

Various continuous dictionaries of features can be considered under regularity conditions required later on. They include many parametric families of functions known in statistics and compressed sensing literature.

1. *Gaussian scaled-spikes deconvolution.* The noisy linear combination of translated and re-scaled Gaussian features corresponds to:

$$h(t, \sigma) \mapsto \frac{\exp(-t^2/2\sigma^2)}{\pi^{1/4}\sigma^{1/2}} \quad \text{on } \Theta \times \mathfrak{S} = \mathbb{R} \times \mathbb{R}_+^*. \quad (7)$$

The example of Gaussian spikes deconvolution is analyzed in full details in [6, Section 8] when  $\sigma_T$  does not depend on  $T$ . We shall consider here that the scale parameter  $\sigma_T$  may vary with  $T$ .

2. *Multi-resolution approximation.* We consider the normalized Shannon scaling function:

$$h(t, \sigma) \mapsto \sqrt{\sigma} \frac{\sin(\pi t/\sigma)}{\pi t} \quad \text{on } \Theta \times \mathfrak{S} = \mathbb{R} \times \mathbb{R}_+^*.$$

The associated dictionary allows to recover functions whose Fourier transform have their support in  $[-\pi/\sigma, \pi/\sigma]$  (see [17, Theorem 3.5]).

3. *Low-pass filter.* We consider the normalized Dirichlet kernel on the torus for some cut-off frequency  $f_c \in \mathbb{N}^*$  and  $T = 2f_c + 1$ :

$$h(t, \sigma) = \frac{1}{\sqrt{T}} \sum_{k=-f_c}^{f_c} e^{2i\pi kt} = \frac{\sin(T\pi t)}{\sqrt{T} \sin(\pi t)}, \quad (8)$$

with  $\sigma = 1/T$ ,  $T \in 2\mathbb{N}^* + 1$  and  $t \in \Theta = \mathbb{R}/\mathbb{Z}$ . The example of the low-pass filter is addressed in [11], where exact support recovery results are obtained for the BLasso estimators. This dictionary is also used in [7] in the context of super-resolution. Bounds on some prediction risks (different from those considered in this paper) are established therein for estimators obtained by solving the constrained formulation of the BLasso.

### 2.2. Examples of observation spaces and Gaussian noise processes

We consider both discrete-time and continuous-time processes in our general model.

### 2.2.1. Discrete-time process observed on a regular grid

Consider a real-valued process  $y$  observed over a regular grid  $t_1 < \dots < t_T$  of a symmetric interval  $[-a_T, a_T] \subset \mathbb{R}$ , with  $T \geq 1$ ,  $t_j = -a_T + j\Delta_T$  for  $j = 1, \dots, T$  and grid step:  $\Delta_T = 2a_T/T$ . We set:

$$\lambda_T = \Delta_T \sum_{j=1}^T \delta_{t_j} \quad (9)$$

Then, we see  $y$  as an element of  $L^2(\lambda_T)$ . We have for any function  $f \in L^2(\lambda_T)$  that  $\|f\|_{L^2(\lambda_T)} = \sqrt{\Delta_T} \|f\|_{\ell_2}$ , where the right-hand side is understood as the  $\ell_2$ -norm (Euclidean norm) of the vector  $(f(t_1), \dots, f(t_T))$ .

We assume that  $(a_T, T \geq 2)$  is a sequence of positive numbers, such that:  $\lim_{T \rightarrow \infty} a_T = +\infty$  and  $\lim_{T \rightarrow \infty} \Delta_T = 0$  so that the sequence of measures  $(\lambda_T, T \geq 1)$  converges with respect to the vague topology towards the Lebesgue measure, noted  $\text{Leb}$ , on  $\mathbb{R}$ . When  $\Theta = \mathbb{R}$ , it is therefore natural in this case, to consider non-linear parameters within the support of the observations and take  $\Theta_T = [-a_T, a_T]$ . When  $T$  tends to infinity, in the limit model the observation corresponds to a square integrable random process indexed on  $\Theta = \mathbb{R}$ . In the case of periodic signals, we may take the sets  $\Theta$  and  $\Theta_T$  to be the torus  $\mathbb{R}/\mathbb{Z}$ , and the limit measure is then the Haar measure identified with the Lebesgue measure.

In this formalism, the noise  $w_T \in L^2(\lambda_T)$  is given by:

$$w_T(t) = \sum_{j=1}^T G_j \mathbf{1}_{\{t_j\}}(t), \quad (10)$$

where  $\mathbf{1}_A$  denotes the indicator function of an arbitrary set  $A$  and  $(G_1, \dots, G_T)$  is a centered Gaussian random vector with independent entries of variance  $\bar{\sigma}^2$ .

In this case Assumption 1.1 holds with an equality in (4) and  $\mathbb{E}[\|w_T\|_{L^2(\lambda_T)}^4]$  is finite. Notice that  $\mathbb{E}[\|w_T\|_{L^2(\lambda_T)}^2] = \bar{\sigma}^2 \Delta_T T$ , thus the Cauchy-Schwarz inequality (5) gives an upper bound larger by a factor  $T$  than the value given by (4). We also have that  $\Xi_T = 2\bar{\sigma}^4 \Delta_T^2 T$ .

Finally, the model writes:

$$y_j := y(t_j) = \sum_{k=1}^s \beta_k^* \phi_T(\theta_k^*, t_j) + G_j, \quad j = 1, \dots, T.$$

We stress that when the noises  $(G_j)_{1 \leq j \leq T}$  are independent the model encompasses the Gaussian sequence model where the mean vector is the sampling of a linear combination of shifts of a known function.

### 2.2.2. Continuous-time processes

Assume we observe a real-valued process  $y$  on a topological state space. We note  $\lambda = \lambda_T$  for a  $\sigma$ -finite measure on the state space. In this framework,  $y$  is an

element of  $L^2(\lambda)$ . Let us assume that the noise is  $w_T = \sum_{k \in \mathbb{N}} \sqrt{\xi_k} G_k \psi_k$ , where  $(G_k, k \in \mathbb{N})$  are independent centered Gaussian random variables with variance  $\bar{\sigma}^2$ ,  $\psi = (\psi_k, k \in \mathbb{N})$  an orthonormal sequence of  $L^2(\lambda)$ , and  $\xi = (\xi_k, k \in \mathbb{N})$  a summable sequence of non-negative real numbers. The sequences  $\psi$  and  $\xi$  may depend on  $T$ . Let  $\|\xi\|_{\ell_p}$  denote the usual  $\ell_p$ -norm of the sequence  $\xi$ . We have:

$$\text{Var}(\langle f, w_T \rangle_{L^2(\lambda)}) = \bar{\sigma}^2 \sum_{k \in \mathbb{N}} \xi_k \langle f, \psi_k \rangle_{L^2(\lambda)}^2 \leq \bar{\sigma}^2 \Delta_T \|f\|_{L^2(\lambda)}^2,$$

with  $\Delta_T = \|\xi\|_{\ell_\infty} = \sup_{k \in \mathbb{N}} \xi_k$ . We also have  $\mathbb{E}[\|w_T\|_{L^2(\lambda)}^2] = \bar{\sigma}^2 \|\xi\|_{\ell_1}$  and  $\Xi_T = \text{Var}(\|w_T\|_{L^2(\lambda)}^2) = 2\bar{\sigma}^4 \|\xi\|_{\ell_2}^2$ . In particular Assumption 1.1 holds.

We may consider different choices for  $\xi$  that lead to different values for  $\Xi_T$ , the variance of the squared norm of the noise. For instance, our framework encompasses the truncated white noise by taking for all  $k \in \mathbb{N}$ ,  $\xi_k = T^{-1} \mathbf{1}_{\{1 \leq k \leq T\}}$ . In this case, we have  $\|\xi\|_{\ell_\infty} = 1/T$  and  $\|\xi\|_{\ell_1} = 1$ . In particular, we get that the inequality (5) is not as sharp as (4) since  $\Delta_T = 1/T$  whereas  $\mathbb{E}[\|w_T\|_{L^2(\lambda)}^2] = \bar{\sigma}^2$ .

### 3. Assumptions and prediction bounds

We recall in this section assumptions and definitions from Sections 3-5 of [6] in a simpler way adapted to our framework. In [6], the authors established high probability bounds for prediction and estimation errors associated to some estimators of  $\beta^*$  and  $\vartheta^*$  tackling a wider range of dictionaries.

#### 3.1. Regularity of the features

We gather in this section the hypotheses that will be required on the features defined by (3).

Recall that the parameter space  $\Theta$  is either  $\mathbb{R}$  or the torus  $\mathbb{R}/\mathbb{Z}$  endowed with the Lebesgue measure  $\text{Leb}$ . For convenience, we write  $|x - y|$  for the Euclidean distance between  $x$  and  $y$  either on  $\mathbb{R}$  or on the torus. Recall also that  $L^2(\lambda_T)$  and  $L^2(\text{Leb})$  are the sets of square integrable functions on  $\Theta$  with respect to the measures  $\lambda_T$  and  $\text{Leb}$  respectively. We denote  $\mathfrak{S}$  the set of scale parameter values.

**Assumption 3.1** (Smoothness of the features). *Let  $h$  be a function defined on  $\Theta \times \mathfrak{S}$ . Let  $T \in \mathbb{N}$  and  $\sigma_T \in \mathfrak{S}$ . We assume that the function  $\theta \mapsto h(\theta, \sigma_T)$  is of class  $\mathcal{C}^3$  on  $\Theta$ . We assume furthermore that  $\|h(\cdot, \sigma_T)\|_{L^2(\text{Leb})} = 1$ , and that for all  $\theta \in \Theta$   $\|h(\theta - \cdot, \sigma_T)\|_{L^2(\lambda_T)} > 0$  and all  $i \in \{0, \dots, 3\}$ :*

$$\|\partial_\theta^i h(\cdot, \sigma_T)\|_{L^2(\text{Leb})} < +\infty \quad \text{and} \quad \|\partial_\theta^i h(\theta - \cdot, \sigma_T)\|_{L^2(\lambda_T)} < +\infty.$$

Recall the function  $\varphi_T$  defined by (3) and notice that Assumption 3.1 implies  $\|\varphi_T(\theta)\|_{L^2(\lambda_T)} > 0$  on  $\Theta$ . We define the function:

$$g_T(\theta) = \|\partial_\theta \phi_T(\theta)\|_{L^2(\lambda_T)}^2, \quad \text{where } \phi_T(\theta) = \varphi_T(\theta) / \|\varphi_T(\theta)\|_{L^2(\lambda_T)}. \quad (11)$$

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

Let us mention that if for all  $\theta \in \Theta$ ,  $\varphi_T(\theta)$  and  $\partial_\theta \varphi_T(\theta)$  are linearly independent functions of  $L^2(\lambda_T)$  and  $\|\partial_\theta \varphi_T(\theta)\|_{L^2(\lambda_T)} > 0$ , then  $g_T > 0$  on  $\Theta$  (see [6, Lemma 3.1]).

### 3.2. Definition of the kernel and its approximation

#### 3.2.1. Measuring the colinearity of the features

We define the symmetric kernel  $\mathcal{K}_T$  on  $\Theta^2$  by:

$$\mathcal{K}_T(\theta, \theta') = \langle \phi_T(\theta), \phi_T(\theta') \rangle_{L^2(\lambda_T)}. \quad (12)$$

The kernel  $\mathcal{K}_T$  measures the colinearity of two features belonging to the continuous dictionary. It does not *a priori* have a simple form. In the following, we approximate this kernel by another kernel easier to handle.

As mentioned in the introduction, we consider in this paper a setting where the sequence of measures  $(\lambda_T, T \geq 1)$  converges in some sense towards the Lebesgue measure  $\text{Leb}$  on  $\Theta$ . In [6], the kernel  $\mathcal{K}_T$  was free of any scale parameter  $\sigma_T$  and authors have considered a pointwise limit kernel  $\mathcal{K}_\infty = \lim_{T \rightarrow \infty} \mathcal{K}_T$  which is free of  $T$  and allows to continue the proofs under some assumptions. However, due to our scale parameter  $\sigma_T$  which decreases towards zero with  $T$ , we show in the following example that the pointwise limit kernel is degenerate.

*Example 3.1* (Degenerate limit kernel). Consider the discrete-time process presented in Section 2.2.1 with the measure  $\lambda_T$  from (9) and the Gaussian features (7) from Section 2.1 scaled by the sequence  $(\sigma_T, T \geq 1)$  that tends towards zero when  $T$  grows to infinity so that  $\lim_{T \rightarrow +\infty} \Delta_T / \sigma_T = 0$ . In this case, the sequence of measures  $(\lambda_T, T \geq 1)$  converges with respect to the vague topology towards the Lebesgue measure and it is easy to check that  $\mathcal{K}_\infty$ , the pointwise limit of the kernel  $\mathcal{K}_T$ , is equal to zero almost everywhere and to 1 on the diagonal.

Thus, instead of the pointwise limit kernel  $\mathcal{K}_\infty$ , we shall approximate (for finite large enough  $T$ ) the kernel  $\mathcal{K}_T$  by a kernel  $\mathcal{K}_T^{\text{prox}}$  of the form:

$$\mathcal{K}_T^{\text{prox}} : (\theta, \theta') \mapsto F(|\theta - \theta'| / \sigma_T), \quad (13)$$

where  $F$  is a real-valued function defined on  $\mathbb{R}_+$  with  $F(0) = 1$ . (Recall that  $|\theta - \theta'|$  is the Euclidean distance between  $\theta$  and  $\theta'$  on  $\Theta$  which is either  $\mathbb{R}$  or the torus  $\mathbb{R}/\mathbb{Z}$ .) Notice that if  $F$  is of class  $\mathcal{C}^{2\ell}$  with  $F^{(2i+1)}(0) = 0$  for  $i \in \{0, \dots, \ell - 1\}$  for some integer  $\ell \geq 1$  (which is the case if  $F$  can be extended into an even function of class  $\mathcal{C}^{2\ell}$  on  $\mathbb{R}$ ), then  $\mathcal{K}_T^{\text{prox}}$  is of class  $\mathcal{C}^{\ell, \ell}$ . The choice of the function  $F$  follows from the model given by  $h$  in (3), so that  $\mathcal{K}_T$  and  $\mathcal{K}_T^{\text{prox}}$  are close (see (iii) of Assumption 3.4 below). We refer to Sections 6 and 7 for examples with  $h$  given by (7) and (8), respectively. The introduction of the kernel  $\mathcal{K}_T^{\text{prox}}$  is significantly different from the approximation developed in [6].

### 3.2.2. Covariant derivatives of the kernel

Let  $\mathcal{K}$  be a symmetric kernel of class  $\mathcal{C}^2$  such that the function  $g_{\mathcal{K}}$  defined on  $\Theta$  by:

$$g_{\mathcal{K}}(\theta) = \partial_{x,y}^2 \mathcal{K}(\theta, \theta), \quad (14)$$

is positive, where  $\partial_x$  (respectively  $\partial_y$ ) denotes the usual derivative with respect to the first (respectively second) variable. Under Assumptions 3.1 and 3.2, the definitions (11) and (14) coincide so that  $g_T = g_{\mathcal{K}_T}$  on  $\Theta$ .

Similarly to [18], we introduce the covariant derivatives which reduce to elementary expressions since the location parameters are one-dimensional. More precisely following [6, Section 4], we set for a smooth function  $f$  defined on  $\Theta$ ,  $\tilde{D}_{0;\mathcal{K}}[f] = f$ ,  $\tilde{D}_{1;\mathcal{K}}[f] = g_{\mathcal{K}}^{-1/2} f'$  and for  $i \geq 2$ :

$$\tilde{D}_{i;\mathcal{K}}[f] = \tilde{D}_{1;\mathcal{K}}[\tilde{D}_{i-1;\mathcal{K}}[f]].$$

Let us assume that the kernel  $\mathcal{K}$  has the form  $\mathcal{K}(\theta, \theta') = \langle f(\theta), f(\theta') \rangle_{L^2(\lambda)}$  for some function  $f$  of class  $\mathcal{C}^3$  and some measure  $\lambda$  on  $\Theta$ . We then define the covariant derivatives (see (27) in [6]) of  $\mathcal{K}$  for  $i, j \in \{0, \dots, 3\}$  and  $\theta, \theta' \in \Theta$  by:

$$\mathcal{K}^{[i,j]}(\theta, \theta') = \langle \tilde{D}_{i;\mathcal{K}}[f](\theta), \tilde{D}_{j;\mathcal{K}}[f](\theta') \rangle_{L^2(\lambda)}.$$

We also define the function  $h_{\mathcal{K}}$  on  $\Theta$  by:

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

The previous notation will be used both for the kernel  $\mathcal{K}_T$  in (12), which is determined by the particular choice of the features, but also for the kernel  $\mathcal{K}_T^{\text{prox}}$  in (13). The latter is determined by the function  $F$  and we derive next the particular expressions of  $g_{\mathcal{K}_T^{\text{prox}}}$  and of the covariant derivatives of  $\mathcal{K}_T^{\text{prox}}$  under additional assumptions on  $F$ .

For a real valued function  $f$  defined on a set  $A$ , we write  $\|f\|_{\infty} = \sup_{x \in A} |f(x)|$ .

**Assumption 3.3** (Properties of the function  $F$ ). *Let  $F$  be a function defined on  $\mathbb{R}_+$  of class  $\mathcal{C}^6$  with  $F(0) = 1$  and  $F^{(2i+1)}(0) = 0$  for  $i \in \{0, 1, 2\}$ . We set:*

$$g_{\infty} = -F''(0). \quad (15)$$

We assume that:

$$\begin{aligned} g_{\infty} > 0, \quad L_6 := g_{\infty}^{-3} |F^{(6)}(0)| < +\infty, \\ \text{and } L_i := g_{\infty}^{-i/2} \|F^{(i)}\|_{\infty} < +\infty \text{ for all } i \in \{0, \dots, 4\}. \end{aligned} \quad (16)$$

We give the covariant derivatives of the kernel  $\mathcal{K}_T^{\text{prox}}$  according to the definition given in [6, (27)]: for any  $\theta, \theta' \in \Theta$  and  $i, j \in \{0, \dots, 3\}$ ,

$$\mathcal{K}_T^{\text{prox}[i,j]}(\theta, \theta') = \frac{(-1)^j}{g_{\infty}^{(i+j)/2}} F^{(i+j)}(|\theta - \theta'|/\sigma_T). \quad (17)$$

We notice that we have for any  $\theta \in \Theta$ :

$$g_{\mathcal{K}_T^{\text{prox}}}(\theta) = g_{\infty}/\sigma_T^2. \quad (18)$$

### 3.2.3. Measuring the quality of the approximation

In this section, we quantify the proximity of the kernel  $\mathcal{K}_T$  and  $\mathcal{K}_T^{\text{prox}}$ .

Following [18], we define the one-dimensional Riemannian metric  $\mathfrak{d}_T(\theta, \theta')$  between  $\theta, \theta' \in \Theta$  by:

$$\mathfrak{d}_T(\theta, \theta') = |G_T(\theta) - G_T(\theta')|, \quad (19)$$

where  $G_T$  is a primitive of the function  $\sqrt{g_T}$  assumed positive on  $\Theta$  thanks to Assumption 3.2.

Recall that  $\Theta_T$ , introduced below the model (2), is a compact sub-interval of  $\Theta$ . Since  $\Theta_T$  is compact, under Assumptions 3.2 and 3.3, we deduce that the constant  $C_T$  below is positive and finite, where:

$$C_T = \max \left( \sup_{\Theta_T} \sqrt{\frac{g_{\mathcal{K}_T^{\text{prox}}}}{g_T}}, \sup_{\Theta_T} \sqrt{\frac{g_T}{g_{\mathcal{K}_T^{\text{prox}}}}} \right). \quad (20)$$

Elementary calculations show that the metric  $\mathfrak{d}_T$  defined in (19) is equivalent, up to a factor  $\sigma_T$ , to the Euclidean metric on  $\Theta_T$  as for any  $\theta, \theta' \in \Theta_T$ :

$$\frac{1}{C_T} \sqrt{g_\infty} \sigma_T^{-1} |\theta - \theta'| \leq \mathfrak{d}_T(\theta, \theta') \leq C_T \sqrt{g_\infty} \sigma_T^{-1} |\theta - \theta'|. \quad (21)$$

In order to quantify the approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_T^{\text{prox}}$ , we set:

$$\mathcal{V}_T = \max(\mathcal{V}_T^{(1)}, \mathcal{V}_T^{(2)}) \quad (22)$$

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

### 3.3. Boundedness and local concavity on the diagonal of the approximating kernel

Recall the definition of the kernel  $\mathcal{K}_T^{\text{prox}}$  given by (13) using the function  $F$ . We quantify the boundedness and local concavity on the diagonal of the kernel  $\mathcal{K}_T^{\text{prox}}$  using for  $r > 0$ :

$$\varepsilon(r) = 1 - \sup \{|F(r')|; \quad r' \geq r\}, \quad (23)$$

$$\nu(r) = - \sup \{F''(r')/g_\infty; \quad r' \in [0, r]\}. \quad (24)$$

We also quantify the colinearity between  $s \in \mathbb{N}$  features belonging to the continuous dictionary, by setting for  $u > 0$ :

$$\delta(u, s) = \inf \left\{ \delta > 0 : \max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s g_\infty^{-\frac{\delta}{2}} |F^{(i)}(x_\ell - x_k)| \leq u, \right. \\ \left. \text{for all } i \in \{0, 1, 2, 3\} \text{ and } (x_1, \dots, x_s) \in \mathbb{R}^s(\delta) \right\}, \quad (25)$$

where for any subset  $A$  of  $\mathbb{R}$  or  $\mathbb{R}/\mathbb{Z}$  and for any  $\delta \geq 0$ ,

$$A^s(\delta) = \left\{ (\theta_1, \dots, \theta_s) \in A^s : |\theta_\ell - \theta_k| > \delta \text{ for all distinct } k, \ell \in \{1, \dots, s\} \right\}. \quad (26)$$

with the conventions  $\inf \emptyset = +\infty$ , and for  $s = 0, 1$ :  $A^0(\delta) = \{0\}$  and  $A^1(\delta) = A$ .

Following [6], we define quantities which depend only on the function  $F$  and on a real parameter  $r > 0$ :

$$H_\infty^{(1)}(r) = \frac{1}{2} \wedge L_2 \wedge L_3 \wedge L_4 \wedge L_6 \wedge \frac{\nu(2r)}{10} \wedge \frac{\varepsilon(r/2)}{10},$$

$$H_\infty^{(2)}(r) = \frac{1}{6} \wedge \frac{8\varepsilon(r/2)}{10(5 + 2L_1)} \wedge \frac{8\nu(2r)}{9(2L_2 + 2L_3 + 4)},$$

where the constants  $L_i$  are defined in (16).

### 3.4. Main assumption and identifiability of the model

We summarize here all assumptions that are needed for the following results. They concern the features, the function  $F$  characterizing the proxy kernel  $\mathcal{K}_T^{\text{prox}}$ , the proximity of the kernel  $\mathcal{K}_T$  defined by the original features to the prox kernel  $\mathcal{K}_T^{\text{prox}}$  and, last but not least, the assumption that two neighbouring non-linear parameters  $\theta$  and  $\theta'$  are at least separated by some constant multiplied by  $\sigma_T$ . This is the most important improvement on the sufficient conditions in [6], as the scaling parameter  $\sigma_T$  can be chosen small in some models.

**Assumption 3.4.** *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}$ ,  $r \in (0, 1/\sqrt{2g_\infty L_2})$ ,  $\eta \in (0, 1)$  and a subset  $\mathcal{Q} \subset \Theta_T$  of cardinal  $s$ .*

- (i) **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 (11), satisfies the positivity condition of Assumption 3.2.*
- (ii) **Properties of the function  $F$ :** *Assumption 3.3 holds and we have  $\varepsilon(r/2) > 0$  and  $\nu(2r) > 0$ .*
- (iii) **Proximity to the limit setting:** *The kernel  $\mathcal{K}_T$  defined from the dictionary, see (12), is sufficiently close to the kernel  $\mathcal{K}_T^{\text{prox}}$  in the sense that we have:*

$$C_T \leq 2$$

and if  $s \geq 1$ , we have in addition:

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

- (iv) **Separation of the non-linear parameters:** *If  $s \geq 1$ , we have:*

$$\delta(\eta H_\infty^{(2)}(r), s) < +\infty \quad \text{and for any } \theta \neq \theta' \in \mathcal{Q}, \quad |\theta - \theta'| > \sigma_T \Sigma(\eta, r, s),$$

where,

$$\Sigma(\eta, r, s) = 4 \max \left( r g_\infty^{-1/2}, 2 \delta(\eta H_\infty^{(2)}(r), s) \right).$$

*Remark 3.2* (On the separation condition). The separation condition corresponds to the minimal distance between any pair of nonlinear parameters ensuring that a coherence function remains bounded from above by a specified constant dependent on the dictionary. This condition is mathematically represented in (25) and expressed with the following coherence function:

$$\max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s g_\infty^{-\frac{i}{2}} |F^{(i)}(x_\ell - x_k)|,$$

where  $\{x_1, \dots, x_s\}$  is a set of nonlinear parameters. This function is quite similar to the Babel function introduced in [21], which measures the maximum total coherence between a fixed atom and a collection of other atoms in a finite dictionary. In linear cases (when the dictionary consists of a finite number of atoms), keeping the Babel function below a certain threshold allows for the derivation of results on the recovery of sparse signals. We stress that similar separation conditions to Assumption 3.4 are common in super-resolution, compressed sensing and spikes deconvolution for recovering signals derived from continuous dictionaries, see [8, 11, 18] among many other references.

In Sections 6 and 7 we give simplified expressions of the quantities involved in the previous assumption for the particular models in hand.

Under Assumption 3.4, we shall build consistent estimators for  $\beta^*$  and  $\vartheta^*$  of the model (2) and test statistics. The following lemma gives an identifiability result for the considered model under the previous assumptions. Its proof relies on the construction of certificates from [6] and is based on ideas developed in [10] for exact reconstruction of measures, see Lemma 1.1 therein. We recall that by convention  $\beta^* \Phi_T(\vartheta^*) = 0$  when  $s = 0$ .

**Lemma 3.3** (Sufficient conditions for identifiability). *Let  $T \in \mathbb{N}$  and let  $r \in (0, 1/\sqrt{2g_\infty L_2})$ ,  $\eta \in (0, 1)$ . Suppose that Assumption 3.4 holds for the set  $\mathcal{Q}^* = \{\theta_1^*, \dots, \theta_s^*\} \subset \Theta_T$  of cardinal  $s \in \mathbb{N}$  and for the set  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_{s^0}^0\} \subset \Theta_T$  of cardinal  $s^0 \in \mathbb{N}$ . Then, for any vectors  $\beta^* \in (\mathbb{R}^*)^s, \beta^0 \in (\mathbb{R}^*)^{s^0}$ , we have that, up to the same permutation on the components of  $\beta^*$  and  $\vartheta^*$ :*

$$\beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0) \text{ in } L^2(\lambda_T), \quad \text{implies that } s = s^0, \beta^* = \beta^0, \vartheta^* = \vartheta^0. \quad (27)$$

The proof is in Section 8.1.

*Remark 3.4.* Recall that if  $s \geq 1$ , then  $\beta^*$  is a  $s$ -dimensional vector with non-zero entries. Under the assumptions of Lemma 3.3 we have that:

$$\beta^* \Phi_T(\vartheta^*) = 0 \quad \text{if and only if } s = 0.$$

### 3.5. Prediction error bound

We define the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  of  $\beta^*$  and  $\vartheta^*$  as the solution to the following regularized optimization problem with a real tuning parameter  $\kappa > 0$  and a



bound  $K$  on the unknown number  $s$  of active features in the observed mixture:

$$(\hat{\beta}, \hat{\vartheta}) \in \underset{\beta \in \mathbb{R}^K, \vartheta \in \Theta_T^K}{\operatorname{argmin}} \frac{1}{2} \|y - \beta \Phi_T(\vartheta)\|_{L^2(\lambda_T)}^2 + \kappa \|\beta\|_{\ell_1}, \quad (28)$$

where  $\|\cdot\|_{\ell_1}$  corresponds to the usual  $\ell_1$  norm. Since the interval  $\Theta_T$  on which the optimization of the non-linear parameters is performed is a compact interval and the function  $\Phi_T$  is continuous, the existence of at least a solution is guaranteed. The bound  $K$  on the number  $s$  of features in the mixture from model (2) allows to formulate an optimization problem. It can be arbitrarily large. In particular, it is not involved in the bounds on estimation and prediction risks given in [6] with high probability (see Remark 2.4 therein). We stress that the constants in [6] appearing in those bounds may *a priori* depend on  $T$  when the features are scaled by  $\sigma_T$ . We show below that, in fact, those bounds still hold with constants free of  $T$ . The results in [6] as well as the proof of Theorem 3.5 below rely on the existence of certificate functions. In [6], sufficient conditions for the certificate functions to exist are given, see Proposition 7.4 and 7.5 therein. Those conditions require the non-linear parameters in  $\mathcal{Q}^*$  to satisfy the separation condition (32). In our framework where the scaling  $\sigma_T$  decreases to zero, it turns out that this separation is in general increasing with  $s$  and decreasing with  $T$ . However, for some dictionary composed of translated spikes that vanish quickly, it converges to zero when both  $s$  and  $T$  grow to infinity. We refer to Section 6 in this direction.

Recall the definitions of  $g_\infty$  and  $L_2$  given by (15) and (16). The following theorem is a variation of [6, Theorem 2.1].

**Theorem 3.5.** *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}^*$ ,  $K \in \mathbb{N}^*$ ,  $\eta \in (0, 1)$ ,  $r \in (0, 1/\sqrt{2g_\infty L_2})$ . Assume we observe the random element  $y$  of  $L^2(\lambda_T)$  under the regression model (2) with unknown parameters  $\beta^* \in (\mathbb{R}^*)^s$  and  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*)$  a vector with distinct entries in  $\Theta_T$ , a compact interval of  $\Theta$ , such that Assumption 3.4 holds for  $\mathcal{Q}^* = \{\theta_1^*, \dots, \theta_s^*\} \subset \Theta_T$ . Assume that the unknown number of active features  $s$  is bounded by  $K$ . Suppose also that the noise process  $w_T$  satisfies Assumption 1.1 for a noise level  $\bar{\sigma} > 0$  and a decay rate for the noise variance  $\Delta_T > 0$ .*

*Then, there exist finite positive constants  $C_i$ , for  $i = 0, \dots, 3$ , depending on the function  $F$  and on  $r$  such that for any  $\tau > 1$  and a tuning parameter:*

$$\kappa \geq C_1 \bar{\sigma} \sqrt{\Delta_T \log(\tau)}, \quad (29)$$

*we have the prediction error bound of the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  defined in (28) given by:*

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{L^2(\lambda_T)} \leq C_0 \sqrt{s} \kappa, \quad (30)$$

*with probability larger than  $1 - C_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right)$  where  $|\Theta_T|$  is the Euclidean length of  $\Theta_T$ . Moreover, with the same probability, the difference of the  $\ell_1$ -norms of  $\hat{\beta}$  and  $\beta^*$  is bounded by:*

$$\left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| \leq C_3 \kappa s. \quad (31)$$

*Proof.* The proof is similar to the proof of [6, Theorem 2.1] where one replaces the limit kernel noted  $\mathcal{K}_\infty$  therein by the approximating kernel  $\mathcal{K}_T^{\text{prox}}$  defined in (13). The main difference is in checking condition (v) in Theorem 2.1 on the existence of certificate functions. This is done by using Propositions 7.4 and 7.5 therein, and by noticing that the special form of the approximating kernel  $\mathcal{K}_T^{\text{prox}}$  implies that the constants involved do not depend on the scale parameter  $\sigma_T$ . Indeed Equation (17) clearly entails that they do not depend on the scale parameter. The details of the proof are left to the interested reader.  $\square$

*Remark 3.6* (On the separation). We perform the estimation of  $\beta^*$  and  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*)$  from model (2) under the separation condition:

$$|\theta_k^* - \theta_\ell^*| \geq \sigma_T \Sigma(\eta, r, s), \quad \text{for all } 1 \leq k, \ell \leq s, k \neq \ell, \quad (32)$$

with  $\Sigma(\eta, r, s)$  given in (iv) of Assumption 3.4. Taking into account the separation condition, the number of admissible features which can be used for the prediction is at most of order  $|\Theta_T|/\sigma_T$ ; this provides a natural upper bound on  $s$ . As  $\eta$  is usually fixed, we highlight that the least separation bound tends towards zero when the scaling  $\sigma_T$  goes down to zero.

#### 4. Goodness-of-fit for the LCTF model

In this section, we build a test procedure to decide if the observation  $y$  derives from a given linear combination of translated features. We build a test  $\Psi$ , *i.e.* a measurable function of the observation  $y$  taking value in  $\{0, 1\}$ , in order to distinguish a null hypothesis  $H_0$  against an alternative  $H_1(\rho)$  depending on a nonnegative separation parameter  $\rho$ . We recall that the maximal type I and II error probabilities are  $\sup_{(\beta^*, \vartheta^*) \in H_0} \mathbb{E}_{(\beta^*, \vartheta^*)}[\Psi]$  and  $\sup_{(\beta^*, \vartheta^*) \in H_1(\rho)} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi]$ , respectively, where  $\Psi$  is a function of  $y$  which is equal to  $\beta^* \Phi_T(\vartheta^*) + w_T$  under  $\mathbb{E}_{(\beta^*, \vartheta^*)}$ . The maximal testing risk is the sum of the former quantities, that is:

$$R_\rho(\Psi) = \sup_{(\beta^*, \vartheta^*) \in H_0} \mathbb{E}_{(\beta^*, \vartheta^*)}[\Psi] + \sup_{(\beta^*, \vartheta^*) \in H_1(\rho)} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi],$$

and the minimax testing risk is:

$$R_\rho^* = \inf_{\Psi} R_\rho(\Psi), \quad (33)$$

where the infimum is taken over all the measurable functions from  $L^2(\lambda_T)$  to  $\{0, 1\}$ . The minimax separation rate of the test problem is defined for any  $\alpha \in (0, 1)$  as:

$$\rho^*(\alpha) = \inf\{\rho > 0 : R_\rho^* \leq \alpha\}. \quad (34)$$

#### 4.1. Test problem

Let  $s^0 \in \mathbb{N}$  and consider the set  $\Theta_T^{s^0}(\delta^0) \subset \Theta_T^{s^0}$  of vectors whose components are pairwise separated by a distance  $\delta^0 \geq 0$  (recall the definition (26)). Consider the vectors  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and  $\vartheta^0 = (\theta_1^0, \dots, \theta_{s^0}^0) \in \Theta_T^{s^0}(\delta^0)$ . By convention, we have for  $s^0 = 0$  that  $\beta^0 = 0$ ,  $\vartheta^0 = 0$  and  $\beta^0 \Phi_T(\vartheta^0) = 0$ .

We build a test procedure based on the observation  $y$  to decide, for some  $\delta^* \geq 0$ , whether:

$$\begin{cases} H_0 : & (\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*) \quad \text{s.t.} \quad \beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0), \\ H_1(\rho) : & (\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*) \quad \text{s.t.} \quad \|\beta^* \Phi_T(\vartheta^*) - \beta^0 \Phi_T(\vartheta^0)\|_{L^2(\lambda_T)} \geq \rho, \end{cases} \quad (35)$$

where  $\rho$  is a nonnegative separation parameter. When Assumption 3.4 holds for the sets  $\mathcal{Q}^* = \{\theta_1^*, \dots, \theta_s^*\}$  and  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_{s^0}^0\}$ , by Lemma 3.3, the null hypothesis implies that  $(\beta^*, \vartheta^*) = (\beta^0, \vartheta^0)$  (up to the same permutation on the components of  $\beta^*$  and  $\vartheta^*$ ). We remark that the separation condition from Point (iv) of Assumption 3.4 required between the elements of  $\mathcal{Q}^*$  (resp.  $\mathcal{Q}^0$ ) is automatically satisfied when  $\delta^* \geq \sigma_T \Sigma(\eta, r, s)$  (resp.  $\delta^0 \geq \sigma_T \Sigma(\eta, r, s^0)$ ).

We shall denote the distribution under the null hypothesis as associated to the parameters  $(\beta^0, \vartheta^0)$  and see that the maximal type I error probability writes in this case  $\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi]$  for  $\mathbb{E}_{(\beta^*, \vartheta^*)}[\Psi]$ . Furthermore, when  $s^0 = 0$ , under Assumption 3.4 for the set  $\mathcal{Q}^*$ , Lemma 3.3 implies that the null hypothesis reduces to  $H_0 : s = 0$ .

#### 4.2. Main results

We consider the test procedure  $\Psi_{\text{Test}}(t)$  associated to a real valued statistic Test (measurable function of the observation  $y$ ) and a threshold  $t > 0$  (defining a critical region) given by:

$$\Psi_{\text{Test}}(t) = \mathbf{1}_{\{|\text{Test}| > t\}}. \quad (36)$$

We recall that for a test  $\Psi$ , we accept  $H_0$  when  $\Psi = 0$  and reject it when  $\Psi = 1$ .

It is now well-known that several test statistics may be combined to cover for several regimes in the set of parameters. Our test statistics will be produced by estimating in two different ways  $\|\beta^* \Phi_T(\vartheta^*) - \beta^0 \Phi_T(\vartheta^0)\|_{L^2(\lambda_T)}^2$ , the squared  $L^2(\lambda_T)$  distance separating the null and the alternative hypothesis. On the one hand, we plug-in the estimators from the previous section into this distance and, on the other hand, we use the observed process  $y$  as a proxy for the unknown signal, in which case it is necessary to remove the known bias term  $\mathbb{E} \left[ \|w_T\|_{L^2(\lambda_T)}^2 \right]$  as follows.

Let  $s^0 \in \mathbb{N}$  and consider known linear coefficients and location parameters  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and  $\vartheta^0 = (\theta_1^0, \dots, \theta_{s^0}^0) \in \Theta_T^{s^0}$ , respectively. We define two statistics

Test<sub>1</sub> and Test<sub>2</sub> by:

$$\begin{aligned} \text{Test}_1 &= \left\| y - \beta^0 \Phi_T(\vartheta^0) \right\|_{L^2(\lambda_T)}^2 - \mathbb{E} \left[ \|w_T\|_{L^2(\lambda_T)}^2 \right], \\ \text{Test}_2 &= \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^0 \Phi_T(\vartheta^0) \right\|_{L^2(\lambda_T)}^2, \end{aligned} \quad (37)$$

where  $\hat{\beta}$  and  $\hat{\vartheta}$  denote the estimators obtained from (28) for a given value of the tuning parameter  $\kappa$  and a bound  $K$  on the unknown number  $s \in \mathbb{N}$  of active features in the observed signal.

Recall the definition (6) of  $\Xi_T$ , the variance of the squared  $L^2(\lambda_T)$ -norm of the noise  $w_T$ . The following theorem gives an upper bound of the maximal testing risk associated to the test  $\Psi_{\text{Test}_1}(t)$  for some positive threshold  $t$  and positive separation  $\rho$ . Its proof can be found in Section 8.2.

**Theorem 4.1.** *Let  $T \in \mathbb{N}$  and  $s^0 \in \mathbb{N}$ . Let:*

$$\delta^* \geq 0 \quad \text{and} \quad \delta^0 \geq 0.$$

*Assume that we observe the random element  $y$  of  $L^2(\lambda_T)$  under the regression model (2) with unknown parameters  $s \in \mathbb{N}$ ,  $\beta^* \in (\mathbb{R}^*)^s$  and  $\vartheta^* \in \Theta_T^s(\delta^*)$ . Let  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and  $\vartheta^0 \in \Theta_T^{s^0}(\delta^0)$ . Suppose that Assumption 3.1 on the smoothness of the features holds. Suppose that Assumption 1.1 holds for a noise level  $\bar{\sigma} > 0$  and a decay rate for the noise variance  $\Delta_T > 0$ .*

*Then, the test  $\Psi_{\text{Test}_1}$  in (36) using Test<sub>1</sub> in (37) satisfies:*

$$R_\rho(\Psi_{\text{Test}_1}(t)) \leq \frac{\Xi_T}{t^2} + \frac{4\Xi_T}{(\rho^2 - t)^2} + e^{-(\rho^2 - t)^2 / (32\bar{\sigma}^2 \Delta_T \rho^2)}, \quad (38)$$

*for any threshold  $t$  and any separation  $\rho$  such that  $\rho^2 > t > 0$ .*

We deduce from Theorem 4.1 upper bounds on the minimax separation  $\rho^*$  defined in (34) for the goodness-of-fit test problem (35).

**Corollary 4.2.** *Under the framework and the assumptions of Theorem 4.1, the minimax separation rate for the test problem (35) verifies for any  $\alpha \in (0, 1)$ :*

$$\rho^*(\alpha) \leq \rho^{(1)}(\alpha) \quad \text{with} \quad \rho^{(1)}(\alpha) := \max \left( \left( \frac{40\Xi_T}{\alpha} \right)^{1/4}, 8\bar{\sigma} \sqrt{2\Delta_T \log \left( \frac{2}{\alpha} \right)} \right). \quad (39)$$

*Proof of Corollary 4.2.* This result is a direct consequence of Theorem 4.1 by taking the threshold  $t$  of the test therein equal to  $\rho^2/2$ . Then, we have that for  $\rho > 0$ :

$$R_\rho^* \leq R_\rho(\Psi_{\text{Test}_1}(\rho^2/2)) \leq \frac{4\Xi_T}{\rho^4} + \frac{16\Xi_T}{\rho^4} + e^{-\rho^2/(128\bar{\sigma}^2 \Delta_T)} = \frac{20\Xi_T}{\rho^4} + e^{-\rho^2/(128\bar{\sigma}^2 \Delta_T)}.$$

We deduce that  $R_\rho^* \leq \alpha$  for any  $\alpha \in (0, 1)$  whenever the separation  $\rho$  satisfies:

$$\rho \geq \left( \frac{40\Xi_T}{\alpha} \right)^{1/4} \vee \bar{\sigma} \sqrt{128 \Delta_T \log \left( \frac{2}{\alpha} \right)}. \quad (40)$$

This implies (39).  $\square$

In the following theorem, we give a bound of the maximal testing risk associated to the test  $\Psi_{\text{Test}_2}(t)$  using  $\text{Test}_2$  in (37) for solving the test problem (35). The statistic  $\text{Test}_2$  is defined using estimators of the model parameters  $(\beta^*, \vartheta^*)$ . In view of recovering the latter, we assume that the minimal distance  $\delta^*$  (resp.  $\delta^0$ ) is large enough so that Point (iv) of Assumption 3.4 is satisfied for the components of  $\vartheta^*$  (resp.  $\vartheta^0$ ).

Recall the definitions of  $g_\infty$  and  $L_2$  given by (15) and (16), that  $|\Theta_T|$  denotes the Euclidean length of the compact set  $\Theta_T$  and  $\Sigma$  defined in (iv) of Assumption 3.4.

**Theorem 4.3.** *Let  $T \in \mathbb{N}$ ,  $s^0 \in \mathbb{N}$  and choose  $K \in \mathbb{N}$  such that  $s_0 \leq K$ . Let also  $\eta \in (0, 1)$  and  $r \in (0, 1/\sqrt{2g_\infty L_2})$ . Let*

$$\delta^* \geq \sigma_T \Sigma(\eta, r, s) \quad \text{and} \quad \delta^0 \geq \sigma_T \Sigma(\eta, r, s^0). \quad (41)$$

Assume we observe the random element  $y$  of  $L^2(\lambda_T)$  under the regression model (2) with unknown parameters  $s \in \mathbb{N}$  such that  $s \leq K$ ,  $\beta^* \in (\mathbb{R}^*)^s$  and  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta_T^s(\delta^*)$ . Let  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and  $\vartheta^0 = (\theta_1^0, \dots, \theta_{s^0}^0) \in \Theta_T^{s^0}(\delta^0)$ . Suppose that Assumption 3.4 holds for the sets  $\mathcal{Q}^* = \{\theta_1^*, \dots, \theta_s^*\} \subset \Theta_T$  of cardinal  $s$  and  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_{s^0}^0\} \subset \Theta_T$  of cardinal  $s^0$ . Suppose also that the noise process  $w_T$  satisfies Assumption 1.1 for a noise level  $\bar{\sigma} > 0$  and a decay rate for the noise variance  $\Delta_T > 0$ .

Then, there exist finite positive constants  $\mathcal{C}_0, \mathcal{C}_1, \mathcal{C}_2$ , depending on  $r$  and on the function  $F$ , such that for the tuning parameter  $\kappa$ :

$$\kappa \geq \mathcal{C}_1 \bar{\sigma} \sqrt{\Delta_T \log(\tau)}, \quad \text{for some } \tau > 1, \quad (42)$$

the test  $\Psi_{\text{Test}_2}$  using  $\text{Test}_2$  in (37) satisfies:

$$R_\rho(\Psi_{\text{Test}_2}(t)) \leq 2\mathcal{C}_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right), \quad (43)$$

for any threshold  $t$  and any separation  $\rho$  satisfying:

$$0 < t, \quad \mathcal{C}_0 \sqrt{s^0} \kappa \leq \sqrt{t} < \rho \quad \text{and} \quad \sqrt{t} + \mathcal{C}_0 \sqrt{s} \kappa \leq \rho. \quad (44)$$

The proof can be found in Section 8.3.

*Remark 4.4* (On the bound  $K$ ). The bound  $K$  on  $s$  is assumed to be known. It is needed to formulate the optimization problem (28) whose solutions are the estimators of  $\beta^*$  and  $\vartheta^*$ . However, we stress that the constants  $\mathcal{C}_0, \mathcal{C}_1, \mathcal{C}_2$  and the bound on the maximal testing risk do not depend on  $K$ . Thus,  $K$  can be taken arbitrarily large.

In the next Corollary, we obtain an additional upper bound on the minimax separation rate.

**Corollary 4.5.** *Under the framework and the assumptions of Theorem 4.3 and provided that  $|\Theta_T|/\sigma_T \geq 1$ , there exist finite positive constants  $c$  and  $C$ , depending on  $r$  and the function  $F$ , such that the minimax separation rate for the test problem (35) verifies for any  $\alpha \in (0, 1)$ :*

$$\rho^*(\alpha) \leq \rho^{(2)}(\alpha), \quad \rho^{(2)}(\alpha) := C\bar{\sigma} \sqrt{(s \vee s^0 \vee 1)\Delta_T \log\left(\frac{c|\Theta_T|}{\alpha\sigma_T}\right)}. \quad (45)$$

*Remark 4.6* (On the condition  $|\Theta_T|/\sigma_T \geq 1$ ). We recall that the set  $\Theta_T$  is a compact subset of  $\Theta$ . In the case where  $\Theta$  is the torus  $\mathbb{R}/\mathbb{Z}$ ,  $\Theta_T = \Theta$  and the scale parameter  $\sigma_T$  tends towards 0 when  $T$  grows to infinity, the condition  $|\Theta_T|/\sigma_T \geq 1$  is satisfied for  $T$  large enough. This condition also holds for  $T$  large enough in the Gaussian spikes deconvolution example, with the particular choices for  $\Theta_T$  and  $\sigma_T$  from Section 6, where  $\Theta = \mathbb{R}$ ,  $\lim_{T \rightarrow +\infty} \Theta_T = \Theta$  and  $\lim_{T \rightarrow +\infty} \sigma_T = 0$ .

*Proof of Corollary 4.5.* Notice that all the assumptions of Theorem 4.3 are in force. The result is a direct consequence of Theorem 4.3. We fix the tuning parameter  $\kappa = C_1\bar{\sigma}\sqrt{\Delta_T \log(\tau)}$  by taking the equality in (42). Then, for

$$\rho \geq C_0 \sqrt{s \vee 1} \kappa + \sqrt{t} \quad \text{and} \quad t = C_0^2 (s^0 \vee 1) \kappa^2, \quad (46)$$

we have (44) (in particular  $0 < t < \rho$ ) and by Theorem 4.3 for  $\tau > 1$ :

$$R_\rho^* \leq R_\rho(\Psi_{\text{Test}_2}(t)) \leq 2C_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right),$$

where the finite positive constants  $C_0, C_1, C_2$ , from Theorem 4.3 depend on  $r$  and  $F$ .

Then, taking  $\tau = c|\Theta_T|/(\alpha\sigma_T)$  with  $c = (2C_2) \vee e$  and using that by assumption  $|\Theta_T|/\sigma_T \geq 1$ , we get for  $\rho \geq \sqrt{2}C_0C_1\bar{\sigma}\sqrt{(s + s^0) \vee 2}\sqrt{\Delta_T \log(c|\Theta_T|/(\alpha\sigma_T))}$  and  $\alpha \in (0, 1)$  that  $R_\rho^* \leq \alpha$ . We readily deduce (45) with  $C = 2C_0C_1$ .  $\square$

*Remark 4.7* (Combining the upper bounds of Corollaries 4.2 and 4.5). Let  $\alpha \in (0, 1)$ . Suppose that the assumptions of Corollaries 4.2 and 4.5 hold. Previous results show that each procedure may perform better than the other one in convenient regimes of the parameters, involving the unknown parameter  $s$ . In order to aggregate the two procedures into an automatic one, we take the maximum of the two test procedures. This aggregated test procedure rejects as soon as at least one of the procedures rejects, and accepts otherwise.

More precisely, let  $\rho^{(1)}(\alpha/2)$  be defined by (39) with  $\alpha$  replaced by  $\alpha/2$  and set  $t^{(1)} = (\rho^{(1)}(\alpha/2))^2/2$ ; and let  $\rho^{(2)}(\alpha/2)$  be defined in (45) and  $t^{(2)}$  be given by (46) with  $\alpha$  replaced by  $\alpha/2$ . Then, Corollaries 4.2 and 4.5 imply that  $R_{\rho^{(1)}}(\Psi_{\text{Test}_1}(t^{(1)})) \leq \alpha/2$  and  $R_{\rho^{(2)}}(\Psi_{\text{Test}_2}(t^{(2)})) \leq \alpha/2$ . We define the test:

$$\Psi^{\max} = \max(\Psi_{\text{Test}_1}(t^{(1)}), \Psi_{\text{Test}_2}(t^{(2)})).$$

It is straightforward to see that the type I error probability satisfies:

$$\sup_{(\beta^*, \vartheta^*) \in H_0} \mathbb{E}_{(\beta^*, \vartheta^*)}[\Psi^{\max}] \leq \alpha.$$

Moreover, we have for  $\rho^{\min}(\alpha) = \rho^{(1)}(\alpha/2) \wedge \rho^{(2)}(\alpha/2)$  the following bound on the type II error probability:

$$\sup_{(\beta^*, \vartheta^*) \in H_1(\rho^{\min})} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi^{\max}] \leq \alpha/2.$$

Therefore, we deduce an upper bound on  $\rho^*(\alpha)$  of order  $\rho^{\min}(\alpha)$ , that is:

$$\rho^{\min}(\alpha) = \min \left( \left( \frac{80\Xi_T}{\alpha} \right)^{1/4}, C\bar{\sigma} \sqrt{(s \vee s^0 \vee 1)\Delta_T \log \left( \frac{2c|\Theta_T|}{\alpha\sigma_T} \right)} \right), \quad (47)$$

for a positive constant  $c \geq 2$ . We identify two regimes depending on whether the number of features of the observed signal is sufficiently small or not. Indeed, we notice that when  $\alpha$  is fixed and:

$$s \vee s^0 \vee 1 \ll \left( \frac{\Xi_T}{\alpha} \right)^{1/2} \cdot \left( \bar{\sigma}^2 \Delta_T \log \left( \frac{2c|\Theta_T|}{\alpha\sigma_T} \right) \right)^{-1},$$

Corollary 4.5 yields a sharper upper bound on the separation rate than Corollary 4.2.

### 4.3. Minimax separation rates for signal detection

We illustrate our results on a simple model motivated by [15] for sparse linear regression. We consider a discrete-time process  $y$  over a regular grid  $t_1 < \dots < t_T$  on  $\Theta = \mathbb{R}/\mathbb{Z}$  with grid step  $\Delta_T = 1/T$ . We set  $\lambda_T$  and  $w_T$  as in (9) and (10) from Section 2.2.1. We recall that  $\Xi_T = 2\bar{\sigma}^4 \Delta_T^2 T$  where  $\bar{\sigma} > 0$  is the noise level. In the following, we assume without any loss of generality that  $\bar{\sigma} = 1$ .

Let us consider the framework of signal detection when  $s^0 = 0$ . Under the assumptions of Corollary 4.5, the test problem (35) reduces to:

$$\begin{cases} H_0 : & \beta^* = 0, \\ H_1(\rho) : & (\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*) \quad \text{s.t.} \quad \|\beta^* \Phi_T(\vartheta^*)\|_{L^2(\lambda_T)} \geq \rho. \end{cases} \quad (48)$$

Moreover, under the assumptions of Corollary 4.5 (which in particular gives a lower bound on  $\delta^*$ , see (41)) and with the same arguments used to establish (69) in the proof of Lemma 3.3, we can show that:

$$5/6 \leq C_{\min} := \min_{\beta} \frac{\|\beta \Phi_T(\vartheta^*)\|_{L^2(\lambda_T)}}{\|\beta\|_{\ell_2}}, \quad C_{\max} := \max_{\beta} \frac{\|\beta \Phi_T(\vartheta^*)\|_{L^2(\lambda_T)}}{\|\beta\|_{\ell_2}} \leq 7/6. \quad (49)$$

Therefore, the separation in the alternative hypothesis  $H_1(\rho)$  can be formulated as a lower bound on  $\|\beta^*\|_{\ell_2}$  since we have:

$$C_{\min}\|\beta^*\|_{\ell_2} \leq \|\beta^*\Phi_T(\vartheta^*)\|_{L^2(\lambda_T)} \leq C_{\max}\|\beta^*\|_{\ell_2}.$$

We set  $\Theta_T = \Theta$  and thus  $|\Theta_T| = 1$ . We get from (47) the following upper bound on  $\rho^*(\alpha)$  for any  $\alpha \in (0, 1)$ :

$$\rho(\alpha) = C \min \left( \frac{1}{(\alpha T)^{\frac{1}{4}}}, \sqrt{\frac{s}{T} \log \left( \frac{c}{\alpha \sigma_T} \right)} \right), \quad (50)$$

with  $C$  a finite positive constant. Let  $(\alpha_T, T \geq 1)$  be a  $(0, 1)$ -valued sequence which converges to zero when  $T$  grows to infinity. We deduce that:

$$\lim_{s, T \rightarrow +\infty} R_{\rho(\alpha_T)}^* = 0.$$

By letting the sequence  $(\alpha_T, T \geq 1)$  converge towards 0 as slow as we want, we deduce that for a sequence of separations  $(\rho_{s,T}, T \geq 1, s \geq 1)$  such that:

$$\lim_{s, T \rightarrow +\infty} \frac{\rho_{s,T}}{\frac{1}{T^{\frac{1}{4}}} \wedge \sqrt{\frac{s}{T} \log \left( \frac{c}{\sigma_T} \right)}} = +\infty, \quad (51)$$

we have  $\rho_{s,T} \geq \rho(\alpha_T)$  and thus:

$$\lim_{s, T \rightarrow +\infty} R_{\rho_{s,T}}^* = 0.$$

Hence, we have obtained an asymptotic upper bound of the minimax separation associated to the detection of a finite linear combination of features issued from a continuous dictionary.

We now compare this upper bound to the asymptotic lower bound obtained in the case where the dictionary contains a finite number of features instead of a continuum. Assume that the dictionary is fixed, known and contains  $p$  features parametrized by the parameters in the known and fixed set  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_p^0\} \subset \Theta_T$ . We consider the high dimensional linear regression model:

$$y = \beta^* \Phi_T(\vartheta^0) + w_T \quad \text{in } L^2(\lambda_T),$$

with  $\vartheta^0 = (\theta_1^0, \dots, \theta_p^0) \in \Theta_T^p$  and where  $\beta^* \in \mathbb{R}^p$  is a  $s$ -sparse vector. Notice that in this model the entries of  $\beta^*$  can take the value 0. The high dimension comes from the fact that  $p$  can be much larger than  $T$ . Under coherence assumptions on the finite dictionary and for a sequence of separations  $(\rho_{s,T}, T \geq 1, s \geq 1)$  such that:

$$\lim_{s, T \rightarrow +\infty} \frac{\rho_{s,T}}{\frac{1}{T^{\frac{1}{4}}} \wedge \sqrt{\frac{s}{T} \log(p)} \wedge \frac{p^{\frac{1}{4}}}{\sqrt{T}}} = 0, \quad (52)$$



the authors of [15] showed for different hypotheses on the design matrix  $\Phi_T(\vartheta^0)$  that:

$$\lim_{s, T \rightarrow +\infty} R_{\rho_{s, T}}^* = 1.$$

It means that the hypotheses (48) cannot be distinguished asymptotically when the separation converges to zero faster than the rate given by (52).

*Remark 4.8* (Comparison between the rates obtained for finite and continuous dictionaries). In the high-dimensional linear case (i.e.,  $T \leq p$ ), given that  $1/T^{1/4} \leq p^{1/4}/\sqrt{T}$ , the asymptotic minimal intensity allowing signal detection given by (52) becomes:

$$\frac{1}{T^{1/4}} \wedge \sqrt{\frac{s}{T} \log(p)}.$$

This rate matches, up to a logarithmic factor, the rate given by (51) for our more general model. There are two distinct regimes: the sparse case ( $s \leq \sqrt{T}/\log(p)$ ) and the non-sparse case. Additionally, the magnitude of the size  $p$  of the finite dictionary plays an analogous role as the quantity  $1/\sigma_T$  that appears in the logarithmic terms. The term  $1/\sigma_T$  is of the order of the maximal number of shifted elements permissible in our mixture, considering a separation condition of order  $\sigma_T$  and shift parameters within a compact set possibly growing with  $T$ .

## 5. Goodness-of-fit of the dictionary

In spectroscopy, a prescribed material has known chemical components and a list of  $s_0$  corresponding location parameters of the features is provided. From a sampled material we want to decide whether its chemical components are included in the prescribed list. The linear coefficients are non-negative in this case and they are not given, which makes the null hypothesis composite, that is, fixed location parameters and varying positive linear coefficients. We generalize this setup to real valued linear coefficients. Under the null hypothesis the location parameters are still fixed, but the linear coefficients vary with fixed sign.

More precisely, let  $s^0 \in \mathbb{N}$  and let  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_{s_0}^0\} \subset \Theta_T$  be a set of known location parameters pairwise separated by a distance  $\delta^0 \geq 0$  so that the model is identifiable, see Lemma 3.3. We set the vector  $\vartheta^0 = (\theta_1^0, \dots, \theta_{s_0}^0)$ . We include in the null hypothesis all linear combinations:

$$\sum_{j=1}^{s^0} \beta_j^* \varphi_T(\theta_j^0)$$

with  $\beta_j^*$  being either 0 or with the same sign as  $\beta_j^0$ , for all  $j$  from 1 to  $s^0$ . Thus we split the set  $\mathcal{Q}^0$  into  $\mathcal{Q}^{0,+}$  and  $\mathcal{Q}^{0,-}$ , those parameters  $\theta_k^0$  associated to  $\beta_k^0 > 0$  and to  $\beta_k^0 < 0$ , respectively:

$$\mathcal{Q}^{0,\epsilon} = \{\theta_k^0 \in \mathcal{Q}^0 : \epsilon \beta_k^0 > 0\}, \quad \epsilon \in \{+, -\}.$$

Let  $s \in \mathbb{N}^*$ . Assume that we observe a random element  $y$  issued from the model (2) with linear coefficients  $\beta^* \in (\mathbb{R}^*)^s$  and non-linear parameters  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta_T^s$ . We test whether the unknown set:

$$\mathcal{Q}^{*,\epsilon} = \{\theta_k^* \in \mathcal{Q}^* : \epsilon \beta_k^* > 0\} \text{ is a subset of } \mathcal{Q}^{0,\epsilon} \text{ for each } \epsilon \in \{+, -\}.$$

If  $s^0 = 0$ , this amounts to testing that  $\mathcal{Q}^*$  is empty, which corresponds to the signal detection framework presented in Section 4 in the case  $s^0 = 0$ . Hence, we shall assume in this section that  $s_0 \geq 1$ .

For example, in spectroscopy,  $\mathcal{Q}^{0,-}$  is empty because all linear parameters are positive and this amounts to testing that the present chemical elements are in the prescribed list  $\mathcal{Q}^0$  but they may appear with various positive linear coefficients (amplitudes). Under the alternative, other chemical components are present (located at unknown frequencies not in the prescribed list).

### 5.1. A measure of discrepancy between dictionaries

We define the closed balls 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 us define for  $\epsilon \in \{+, -\}$  the set of indices  $\mathcal{I}^\epsilon = \{k \in \{1, \dots, s\}, \epsilon \beta_k^0 > 0\}$ . We introduce for  $r > 0$ ,  $\epsilon \in \{+, -\}$  and  $k \in \mathcal{I}^\epsilon$ , the set  $S_k^\epsilon(r)$  gathering the indices of the elements of  $\mathcal{Q}^{*,\epsilon}$  that are close to the element  $\theta_k^0$  of  $\mathcal{Q}^{0,\epsilon}$ :

$$S_k^\epsilon(r) = \{\ell \in \{1, \dots, s\} : \theta_\ell^* \in \mathcal{B}_T(\theta_k^0, r) \text{ and } \text{sgn}(\beta_\ell^*) = \epsilon\}. \quad (53)$$

Notice that the sets  $S_k^\epsilon(r)$  can be empty. Furthermore, we assume that  $r < \min_{\ell \neq k} \mathfrak{d}_T(\theta_\ell^0, \theta_k^0)/2$  so that the sets  $S_k^\epsilon(r)$  with  $\epsilon \in \{+, -\}$  and  $k \in \mathcal{I}^\epsilon$  are pairwise disjoint. We also set:

$$S(r) = \bigcup_{\epsilon \in \{+, -\}} S^\epsilon(r) \quad \text{with} \quad S^\epsilon(r) = \bigcup_{k \in \mathcal{I}^\epsilon} S_k^\epsilon(r).$$

We now define a discrepancy measure between the model and any approximation by a linear combination of features having their non-linear parameters in  $\mathcal{Q}^0$  and the linear parameters with the same signs, for  $r > 0$ :

$$\mathcal{D}_{T,r}(\beta^*, \vartheta^*, v^0, \vartheta^0) = \sum_{\epsilon \in \{+, -\}} \sum_{k \in \mathcal{I}^\epsilon} \sum_{\ell \in S_k^\epsilon(r)} |\beta_\ell^*| \mathfrak{d}_T(\theta_\ell^*, \theta_k^0)^2 + \sum_{k \in S(r)^c} |\beta_k^*|,$$

where  $S(r)^c$  denotes the complementary set of  $S(r)$  in  $\{1, \dots, s\}$  and  $v^0 = (v_1^0, \dots, v_{s_0}^0)$  contains the signs of all linear coefficients  $\beta^0$ ,  $v_j^0 = \text{sgn}(\beta_j^0)$ . Notice that  $\mathcal{D}_{T,r}(\beta^*, \vartheta^*, v^0, \vartheta^0) = 0$  if and only if  $\mathcal{Q}^{*,+} \subseteq \mathcal{Q}^{0,+}$  and  $\mathcal{Q}^{*,-} \subseteq \mathcal{Q}^{0,-}$ .

### 5.2. The testing hypotheses

We shall test the following hypotheses:

$$\begin{cases} H_0 : & (\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*), & \mathcal{Q}^{*,+} \subseteq \mathcal{Q}^{0,+} \text{ and } \mathcal{Q}^{*,-} \subseteq \mathcal{Q}^{0,-}, \\ H_1(\rho) : & (\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*) & \text{ and } \mathcal{D}_{T,r}(\beta^*, \vartheta^*, v^0, \vartheta^0) \geq \rho, \end{cases} \quad (54)$$

where  $\rho$  and  $\delta^*$  are separation parameters depending *a priori* on  $T$ ,  $s$  and  $s^0$  that need to be evaluated. Notice that the null hypothesis is also composite. We recall the definitions (33) and (34) of the minimax testing risk  $R_\rho^*$  and the minimax separation  $\rho^*$ . In the following, we give upper bounds on the testing risk and on the minimax separation  $\rho^*(\alpha)$  for any  $\alpha \in (0, 1)$ .

### 5.3. Main result

In this section, we build a test for (54). Under Assumptions 3.1 and 3.2, we define the element of  $L^2(\lambda_T)$ :

$$p_0 = \sum_{k=1}^{s^0} \alpha_k \phi_T(\theta_k^0) + \sum_{k=1}^{s^0} \xi_k \tilde{D}_{1,T}[\phi_T](\theta_k^0), \quad (55)$$

where  $\alpha, \xi \in \mathbb{R}^{s^0}$  solve the system:

$$\langle \phi_T(\theta_k^0), p_0 \rangle_{L^2(\lambda_T)} = \text{sgn}(\beta_k^0) \text{ and } \langle \partial_\theta \phi_T(\theta_k^0), p_0 \rangle_{L^2(\lambda_T)} = 0, \quad \forall k \in \{1, \dots, s^0\}. \quad (56)$$

*Remark 5.1.* The element  $p_0$  of  $L^2(\lambda_T)$  coincides with the vanishing derivative pre-certificate which appears in [11, Section 4] and is the solution of (56) with minimal norm  $\|p_0\|_{L^2(\lambda_T)}$ . We state in Lemma 8.1 the existence of such function and prove its further properties used in the following result.

Using the estimator  $\hat{\beta}$  from (28) for a given value of the tuning parameter  $\kappa$ , we define the test statistic:

$$\text{Test}_3 = \left\| \hat{\beta} \right\|_{\ell_1} - \langle y, p_0 \rangle_{L^2(\lambda_T)}. \quad (57)$$

and the corresponding test  $\Psi_{\text{Test}_3}(t) = \mathbf{1}_{\{|\text{Test}_3| > t\}}$ . Thus we use the certificate function as a filter of the signal and note that  $\mathbb{E} \langle y, p_0 \rangle_{L^2(\lambda_T)} = \|\beta^*\|_{\ell_1}$  under the null hypothesis.

**Theorem 5.2.** *Let  $T \in \mathbb{N}$ ,  $s^0 \in \mathbb{N}^*$  and choose  $K \in \mathbb{N}$  such that  $s_0 \leq K$ . Let also  $\eta \in (0, 1)$  and  $r \in (0, 1/\sqrt{2}g_\infty L_2)$ . Let:*

$$\delta^* \geq \sigma_T \Sigma(\eta, r, s) \quad \text{and} \quad \delta^0 \geq \sigma_T \Sigma(\eta, r, s^0).$$

*Assume we observe the random element  $y$  of  $L^2(\lambda_T)$  under the regression model (2) with unknown parameters  $s \in \mathbb{N}^*$  such that  $s \leq K$ ,  $\beta^* \in (\mathbb{R}^*)^s$  and*

$\vartheta^* = (\theta_1^*, \dots, \theta_s^*) \in \Theta_T^s(\delta^*)$ . Let  $v^0 \in \{-1, 1\}^{s^0}$  be a sign vector and let  $\vartheta^0 = (\theta_1^0, \dots, \theta_{s^0}^0) \in \Theta_T^{s^0}(\delta^0)$ . Suppose that Assumption 3.4 holds for the sets  $\mathcal{Q}^* = \{\theta_1^*, \dots, \theta_s^*\} \subset \Theta_T$  of cardinal  $s$  and  $\mathcal{Q}^0 = \{\theta_1^0, \dots, \theta_{s^0}^0\} \subset \Theta_T$  of cardinal  $s^0$ . Suppose also that the noise process  $w_T$  satisfies Assumption 1.1 for a noise level  $\bar{\sigma} > 0$  and a decay rate for the noise variance  $\Delta_T > 0$ .

Then, the test statistic  $\text{Test}_3$  is uniquely defined and there exist finite positive constants,  $a$  and  $\mathcal{C}_i$  with  $i = 1, \dots, 5$ , (depending on  $r$  and on the function  $F$ ) such that for any  $\tau > 1$  and any tuning parameter  $\kappa$ :

$$\kappa \geq \mathcal{C}_1 \bar{\sigma} \sqrt{\Delta_T \log(\tau)}, \quad (58)$$

the test  $\Psi_{\text{Test}_3}$  satisfies:

$$R_\rho(\Psi_{\text{Test}_3}(t)) \leq 2\mathcal{C}_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right) + \frac{2}{\tau^a s^0}, \quad (59)$$

for any threshold  $t > 0$  and any separation  $\rho > 0$  satisfying:

$$t \geq 2\mathcal{C}_3 s^0 \kappa \quad \text{and} \quad \rho \geq \mathcal{C}_4 s \kappa + \mathcal{C}_5 t. \quad (60)$$

The proof is given in Section 8.4.

#### 5.4. Separation rates

We give in this section an upper bound on the minimax separation  $\rho^*$  to test the goodness-of-fit of the dictionary, that is to distinguish the assumptions  $H_0$  and  $H_1(\rho)$  presented in Section 5.

**Corollary 5.3.** *Under the framework and the assumptions of Theorem 5.2, there exist finite positive constants  $c$  and  $C$  (depending on  $r$  and the function  $F$ ) such that provided that  $|\Theta_T|/\sigma_T \geq 1$ , we have for any  $\alpha \in (0, 1)$ :*

$$\rho^*(\alpha) \leq C \bar{\sigma} (s \vee s^0) \sqrt{\Delta_T \log \left( \frac{c |\Theta_T|}{\alpha \sigma_T} \right)}. \quad (61)$$

*Proof.* The result is a direct consequence of Theorem 5.2. We fix the tuning parameter  $\kappa = \mathcal{C}_1 \bar{\sigma} \sqrt{\Delta_T \log(\tau)}$  by taking the equality in (58). Then, for  $\rho \geq \mathcal{C}_4 s \kappa + \mathcal{C}_5 t$  and  $t = 2\mathcal{C}_3 s^0 \kappa$  we have by Theorem 5.2 for  $\tau > 1$  and since  $s_0 \geq 1$ :

$$R_\rho^* \leq R_\rho(\Psi_{\text{Test}_3}(t)) \leq 2\mathcal{C}_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right) + \frac{2}{\tau^a},$$

where the finite positive constants  $a$ ,  $\mathcal{C}_i$  with  $i \in \{1, \dots, 5\}$ , from Theorem 5.2 depend on  $r$  and the function  $F$ .

Hence, by taking  $\tau = c' / (\sigma_T \alpha / (2|\Theta_T|))^{c''}$  with  $c'' = 1 \vee (1/a)$  and  $c' = (2\mathcal{C}_2) \vee e \vee 2^{1/a}$ , we get for  $\rho \geq 2\mathcal{C}_1 ((2\mathcal{C}_3 \mathcal{C}_5) \vee \mathcal{C}_4) \bar{\sigma} (s \vee s^0) \sqrt{\Delta_T \log(c' / (\sigma_T \alpha / (2|\Theta_T|))^{c''})}$  and  $\alpha \in (0, 1)$  that  $R_\rho^* \leq \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha$ . We then deduce (61) with  $c = 2c'(1/c'')$ .  $\square$

## 6. Gaussian scaled-spikes deconvolution

In this section, we consider the discrete time process observed on a regular grid of  $\mathbb{R}$  given in Section 2.2.1. We recall that Assumption 1.1 holds with:

$$\lambda_T = \Delta_T \sum_{j=1}^T \delta_{t_j} \quad \text{with} \quad t_j = -a_T + j\Delta_T \quad \text{and} \quad \Delta_T = \frac{2a_T}{T},$$

and  $w_T$  given by (10), where  $T \in \mathbb{N}^*$ . We consider the scaled Gaussian features associated to the function:

$$h(t, \sigma) \mapsto \frac{\exp(-t^2/2\sigma^2)}{\pi^{1/4}\sigma^{1/2}} \quad \text{defined on} \quad \Theta \times \mathfrak{S} = \mathbb{R} \times \mathbb{R}_+^*.$$

We shall see below that the natural choice for the function  $F$  appearing in (13) is given by:

$$F = h^0 * h^0 = \pi^{1/4} h^0(\cdot/\sqrt{2}) \quad \text{with} \quad h^0(\cdot) = h(\cdot, 1).$$

In the following, we check that Assumption 3.4 holds. Then, using Theorem 3.5 on a particular example, we provide a prediction bound for the estimator of  $(\beta^*, \vartheta^*)$  solution of the optimization problem (28).

### 6.1. Choice of the approximating kernel

We denote the unscaled feature  $\varphi^0$  on  $\theta \in \Theta$  by:

$$\varphi^0(\theta) = h(\theta - \cdot, 1) = h^0(\theta - \cdot).$$

We define the mapping  $f_T : \Theta \rightarrow \Theta$  by  $f_T(\theta) = \theta/\sigma_T$  for any  $\theta \in \Theta$  and the (pushforward) measure  $\lambda_T^0 = \lambda_T \circ f_T^{-1}$  so that for any  $g \in L^1(\lambda_T^0)$ :

$$\int g(\theta/\sigma_T) \lambda_T(d\theta) = \int g(\theta) \lambda_T^0(d\theta).$$

The Hilbert space  $L^2(\lambda_T^0)$  is endowed with its natural scalar product  $\langle \cdot, \cdot \rangle_{L^2(\lambda_T^0)}$  and norm  $\|\cdot\|_{L^2(\lambda_T^0)}$ . We define on  $\Theta^2$  the kernel:

$$\mathcal{K}_T^0(\theta, \theta') = \langle \phi_T^0(\theta), \phi_T^0(\theta') \rangle_{L^2(\lambda_T^0)} \quad \text{with} \quad \phi_T^0(\theta) = \varphi^0(\theta) / \|\varphi^0(\theta)\|_{L^2(\lambda_T^0)}.$$

The kernel  $\mathcal{K}_T$  can be seen as a scaled kernel derived from  $\mathcal{K}_T^0$  as for  $\theta, \theta' \in \Theta$ :

$$\mathcal{K}_T(\theta, \theta') = \mathcal{K}_T^0(\theta/\sigma_T, \theta'/\sigma_T).$$

When the measure  $\lambda_T^0$  converges in some sense, as  $T$  goes to infinity, towards the Lebesgue measure  $\text{Leb}$  on  $\mathbb{R}$ , it is natural to consider the approximation  $\mathcal{K}_\infty^0$  of  $\mathcal{K}_T^0$  on  $\Theta^2$  by:

$$\mathcal{K}_\infty^0(\theta, \theta') = \langle \phi_\infty^0(\theta), \phi_\infty^0(\theta') \rangle_{L^2(\text{Leb})} \quad \text{with} \quad \phi_\infty^0(\theta) = \varphi^0(\theta) / \|\varphi^0(\theta)\|_{L^2(\text{Leb})}.$$

Thanks to the definition of  $F$ , we also have on  $\Theta^2$  that:

$$F(\theta - \theta') = \mathcal{K}_\infty^0(\theta, \theta').$$

The approximating kernel  $\mathcal{K}_T^{\text{prox}}$  is then given by (13) on  $\Theta^2$ , that is,  $\mathcal{K}_T^{\text{prox}}(\cdot, \cdot) = \mathcal{K}_\infty^0(\cdot/\sigma_T, \cdot/\sigma_T)$ .

## 6.2. Checking Assumption 3.4

### 6.2.1. Regularity of the dictionary

We refer to [6, Section 8] to check that Assumption 3.4 (i) holds for the feature  $\varphi_T$  defined by (3) and any scale parameter  $\sigma_T \in \mathfrak{S} = \mathbb{R}_+^*$ .

### 6.2.2. Boundedness and local concavity on the diagonal

Elementary calculations show that  $g_\infty = -F''(0) = 1/2$ . By definition of  $F$ , we directly deduce that Assumption 3.3 holds. We also get that for  $r \in (0, \sqrt{2})$ :

$$\varepsilon(r) = 1 - e^{-r^2/4} > 0 \quad \text{and} \quad \nu(r) = \left(1 - \frac{r^2}{2}\right) e^{-r^2/4}.$$

We fix  $r \in (0, 1/2)$ . We readily check that Assumption 3.4 (ii) is verified.

### 6.2.3. Proximity to the approximating kernel

In order for the kernel  $\mathcal{K}_T^{\text{prox}}$  to be a good approximation of  $\mathcal{K}_T$  in the sense of Assumption 3.4 (iii), we shall consider the set  $\Theta_T$  over which the optimization is performed:

$$\Theta_T = [-(1 - \xi)a_T, (1 - \xi)a_T] \subset [-a_T, a_T],$$

with a given shrinkage parameter  $\xi \in (0, 1)$ . Intuitively, one does not expect the estimation of the location parameter to perform well near the lower and upper bounds of the observation grid (given by the support of  $\lambda_T$ ). Following [6, Section 8], we set:

$$\gamma_T = 2\Delta_T \sigma_T^{-1} + \sqrt{\pi} e^{-\xi^2 a_T^2 / 2\sigma_T^2}. \quad (62)$$

Recall  $\mathcal{V}_T$  and  $C_T$  defined by (20) and (22). Using Lemma [6, Lemma 8.1], there exist finite positive universal constants  $c_0$ ,  $c_1$  and  $c_2$ , such that  $\gamma_T < c_0$  implies:

$$\mathcal{V}_T \leq c_1 \gamma_T \quad \text{and} \quad |1 - C_T| \leq c_2 \gamma_T. \quad (63)$$

Assume that  $(a_T, T \geq 2)$  and  $(\sigma_T, T \geq 2)$  are sequences of positive numbers, such that:

$$\lim_{T \rightarrow \infty} a_T = +\infty, \quad \lim_{T \rightarrow \infty} \sigma_T = 0 \quad \text{and} \quad \lim_{T \rightarrow \infty} \Delta_T \sigma_T^{-1} = 0. \quad (64)$$

Therefore, we have  $\lim_{T \rightarrow +\infty} \mathcal{V}_T = 0$  and  $\lim_{T \rightarrow +\infty} C_T = 1$ .

Let  $\eta \in (0, 1)$  be fixed. We deduce that under (64), Assumption 3.4 (iii) is satisfied provided that  $T$  is larger than some constant depending on  $\eta$ ,  $r$ , the sparsity  $s$  and the sequences  $(a_T, T \geq 2)$  and  $(\sigma_T, T \geq 2)$ .

#### 6.2.4. Separation of the non-linear parameters

We remark that  $\lim_{r'' \rightarrow \infty} \sup_{|r'| \geq r''} |F^{(i)}(r')| = 0$  for all  $i \in \{0, \dots, 3\}$ . Thus, we deduce from the definition (25) of  $\delta$  that  $\delta(u, s)$  is finite for all  $s \in \mathbb{N}^*$  and  $u > 0$ . Let us stress that  $\sup_{s \in \mathbb{N}^*} \delta(u, s) \leq M/u$  for some universal finite constant  $M$ , see [6, Remark 8.2]. Therefore, the quantity  $\Sigma(\eta, r, s)$  is bounded by a constant depending only on  $\eta$  and  $r$ .

So Assumption 3.4 (iv) is verified as soon as  $|\theta - \theta'| > \sigma_T \Sigma(\eta, r, s)$  for all for all  $\theta \neq \theta' \in \mathcal{Q}^*$ . (Notice this happens for the scaling parameter  $\sigma_T$  small enough depending on  $\mathcal{Q}^*$ .)

### 6.3. Prediction error bound in a particular case

Recall the shrinkage parameter  $\xi \in (0, 1)$  in (62). Let us assume that:

$$a_T = \log(T) \quad \text{and} \quad \sigma_T = 1/\sqrt{\xi \log(T)}.$$

In particular, condition (64) holds. In this case, there exists a finite positive constant  $c$  depending on  $r$ ,  $\eta$  and  $\xi$  such that for  $T \geq c \log(T)^{3/2} s$ , Assumption 3.4 holds (notice that the separation condition (32) of the location parameters in  $\mathcal{Q}^*$  is also verified for  $T$  large enough, depending on  $\mathcal{Q}^*$ , as  $\lim_{T \rightarrow +\infty} \sigma_T = 0$ ). By Theorem 3.5 with  $\tau = T$  and  $\kappa$  given by the equality in (29), we get that:

$$\frac{1}{\sqrt{T}} \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{\ell_2} \leq \mathcal{C}_0 \mathcal{C}_1 \bar{\sigma} \sqrt{\frac{s \log(T)}{T}},$$

with probability larger than  $1 - \mathcal{C}_2 \left( \frac{2\sqrt{\xi} \log(T)}{T} \vee \frac{1}{T} \right)$ , where the constants  $\mathcal{C}_0$ ,  $\mathcal{C}_1$  and  $\mathcal{C}_2$  do not depend on  $T$ .

## 7. Low-pass filter

In this section, we consider the continuous-time process described in Section 2.2.2 on the torus  $\Theta = \mathbb{R}/\mathbb{Z}$  with  $\lambda_T$  the Haar measure on  $\Theta$ , which is identified with the Lebesgue measure  $\text{Leb}$ , and the noise:

$$w_T = \sum_{k \in \mathbb{N}} \sqrt{\xi_k} G_k \psi_k,$$

where  $(G_k, k \in \mathbb{N})$  are independent centered Gaussian random variables with variance  $\bar{\sigma}^2$ ,  $\psi = (\psi_k, k \in \mathbb{N})$  is an o.n.b. of  $L^2(\text{Leb})$  on  $\Theta$  and  $\xi = (\xi_k, k \in \mathbb{N})$  is a summable sequence of non-negative real numbers. The sequences  $\psi$  and  $\xi$  may depend on  $T$ . Recall from Section 2.2.2 that the noise satisfies Assumption 1.1 for a positive noise level  $\bar{\sigma}$  and a decay on the noise variance  $\Delta_T = \sup_{k \in \mathbb{N}} \xi_k$ .

We consider the normalized Dirichlet kernel, see (8), on  $\Theta$ :

$$h(t, \sigma) = \frac{\sin(T\pi t)}{\sqrt{T} \sin(\pi t)} \quad \text{for } t \in \Theta = \mathbb{R}/\mathbb{Z} \text{ and } \sigma = \frac{1}{T}, \quad T \in 2\mathbb{N}^* + 1. \quad (65)$$

The parameter  $T$  is related to the so-called cut-off frequency  $f_c \in \mathbb{N}^*$  by  $T = 2f_c + 1$ . We shall see below that the natural choice for the function  $F$  appearing in (13) is given by:

$$F(t) = \frac{\sin(\pi t)}{\pi t} \quad \text{for } t \in \mathbb{R}. \quad (66)$$

We get from the definition (15) that  $g_\infty = -F''(0) = \pi^2/3$ .

*Remark 7.1.* Note that, if we consider the Shannon scaling function from multi-resolution approximation in Section 2.1 with  $\sigma_T = 1/T$ , then its kernel  $\mathcal{K}_T$  (see (12)) is exactly equal to  $\mathcal{K}_T^{\text{prox}}$  (see (13)) with  $F$  from (66). Therefore there is no approximation in this case. This example can be treated similarly to the low-pass filter.

In the following, we check that Assumption 3.4 hold. Then, using Theorem 3.5, we provide a prediction bound for the estimator of  $(\beta^*, \vartheta^*)$  solution of the optimization problem (28).

### 7.1. The approximating kernel

We define the features  $\varphi_T$  using (3) with  $\sigma_T = 1/T$ . Elementary calculations give that for  $\theta, \theta' \in \Theta$ :

$$\mathcal{K}_T(\theta, \theta') = \frac{\sin(T\pi(\theta - \theta'))}{T \sin(\pi(\theta - \theta'))}.$$

Recall that by convention  $|\theta - \theta'|$  is the Euclidean distance between  $\theta$  and  $\theta'$  in  $\Theta$ , and in particular it belongs to  $[0, 1/2]$ . We define the approximating kernel  $\mathcal{K}_T^{\text{prox}}$  on  $\Theta$  by:

$$\mathcal{K}_T^{\text{prox}}(\theta, \theta') = F(T|\theta - \theta'|) \quad \text{with } |\theta - \theta'| \in [0, 1/2].$$

Since  $F$  is even, we get also that  $F(T|\theta - \theta'|) = F(T(\theta - \theta'))$  where, for  $\theta, \theta' \in \Theta$ , their representers in  $\mathbb{R}$  are chosen so that  $\theta - \theta'$  belongs to  $[-1/2, 1/2]$ .

### 7.2. Checking Assumption 3.4

#### 7.2.1. Regularity of the dictionary

It is elementary to check that  $g_T$  is a constant function on  $\Theta$  equal to  $(T^2 - 1)g_\infty$  and that Assumption 3.4 (i) on the regularity of the dictionary holds.



### 7.2.2. Boundedness and local concavity on the diagonal

There exists  $R > 0$  such that for any  $r \in (0, R)$ :

$$\varepsilon(r) = 1 - \frac{\sin(\pi r)}{\pi r} > 0 \quad \text{and} \quad \nu(r) = - \left( \frac{6}{\pi^3 r^3} - \frac{3}{\pi r} \right) \sin(\pi r) + \frac{6 \cos(\pi r)}{\pi^2 r^2} > 0.$$

We fix  $r \in (0, (1/\sqrt{2g_\infty L_2}) \wedge (R/2))$ . This and the fact that  $F$  is  $\mathcal{C}^\infty$  with bounded derivatives implies that Assumption 3.4 (ii) on the boundedness and the local concavity of the approximating kernel holds.

### 7.2.3. Proximity to the approximating kernel

We set  $\Theta_T = \Theta$ . The proof of the next lemma on the uniform approximation of  $\mathcal{K}_T$  by  $\mathcal{K}_T^{\text{prox}}$  on the torus is postponed to Section 8.5.

**Lemma 7.2.** *There exists a universal positive finite constant  $c_3$  such that for any  $T \in 2\mathbb{N}^* + 1$ :*

$$\mathcal{V}_T \leq \frac{c_3}{T} \quad \text{and} \quad |1 - C_T| \leq \frac{1}{2(T^2 - 1)}. \quad (67)$$

Let  $\eta \in (0, 1)$  be fixed. We deduce from (67) that Assumption 3.4 (iii) is satisfied provided that  $T$  is larger than some constant depending on  $\eta$ ,  $r$ , and the sparsity  $s$ .

### 7.2.4. Separation of the non-linear parameters

Notice that  $\lim_{r'' \rightarrow \infty} \sup_{|r'| \geq r''} |F^{(i)}(r')| = 0$  for all  $i \in \{0, \dots, 3\}$ . Thus, we deduce from the definition (25) of  $\delta(u, s)$  is finite for all  $s \in \mathbb{N}^*$  and  $u > 0$ .

So Assumption 3.4 (iv) is verified as soon as  $|\theta - \theta'| > \sigma_T \Sigma(\eta, r, s)$  for all  $\theta \neq \theta' \in \mathcal{Q}^*$ . (Notice this happens for  $T$  large enough depending on  $\mathcal{Q}^*$  as  $\sigma_T = 1/T$ .)

## 7.3. Prediction error bound

There exists a constant  $c$  depending on  $\eta$  and  $r$  such that for any  $T \in 2\mathbb{N}^* + 1$  such that  $T \geq cs$ , and provided that (32) is satisfied, Assumption 3.4 holds. Using Theorem 3.5 with  $\kappa$  given by an equality in (29) with  $\tau > 1$ , we obtain the prediction bound:

$$\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^* \Phi_T(\vartheta^*) \right\|_{L^2(\text{Leb})} \leq \mathcal{C}_0 \mathcal{C}_1 \bar{\sigma} \sqrt{s \Delta_T \log(\tau)},$$

with probability larger than  $1 - \mathcal{C}_2 \left( \frac{T}{\tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right)$ , where the constants  $\mathcal{C}_0$ ,  $\mathcal{C}_1$  and  $\mathcal{C}_2$  do not depend on  $T$ .

*Remark 7.3.* Exact support recovery results were obtained in [11]. The authors considered a small noise regime, that is:

$$\|w_T\|_{L^2(\text{Leb})} \leq C\kappa, \quad (68)$$

for some finite constant  $C$ . They assumed that the location parameters satisfy for any distinct  $k, \ell \in \{1, \dots, s\}$ , the separation condition  $|\theta_k^* - \theta_\ell^*| \geq C/f_c$  for  $T = 2f_c + 1$ , for some positive constant  $C$  and with  $f_c \geq s$  ( $s$  being the number of active features in the mixture). They showed that there exist finite constants  $C'$  and  $C''$  such that for all  $k \in \{1, \dots, s\}$ :

$$|\tilde{\theta}_k - \theta_k^*| \leq C'\|w_T\|_{L^2(\text{Leb})} \quad \text{and} \quad |\tilde{\beta}_k - \beta_k^*| \leq C''\|w_T\|_{L^2(\text{Leb})},$$

for some estimators  $(\tilde{\beta}, \tilde{\vartheta} = (\tilde{\theta}_1, \dots, \tilde{\theta}_s))$  obtained by solving the BLasso problem.

However the small noise regime assumption is restrictive as it does not encompass the example of Section 2.2.2 where for all  $k \in \mathbb{N}$ ,  $\xi_k = T^{-1}\mathbf{1}_{\{1 \leq k \leq T\}}$  and thus  $\Delta_T = 1/T$  and  $\mathbb{E}[\|w_T\|_{L^2(\text{Leb})}]$  is of order 1. So taking  $\kappa$  given by (29) with an equality and  $\tau = T$ , we deduce that (68) does not hold for  $T$  large. Recall that in (31) we obtain that our estimators satisfy:

$$\left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| \leq C \frac{s \sqrt{\log(T)}}{\sqrt{T}}$$

for some constant  $C > 0$  with high probability. Thus our prediction and estimation rates are smaller by a factor  $\sqrt{\log(T)}/\sqrt{T}$  due to the probabilistic bounds on linear functionals of the noise process that we used in the proof, and this holds under an analogous separation condition on any  $\theta_k^*$  and  $\theta_\ell^*$ , for  $k \neq \ell$  in  $\{1, \dots, s\}$ .

## 8. Technical proofs

### 8.1. Proof of Lemma 3.3

First, for  $s \geq 1$  and  $\vartheta^* = (\theta_1^*, \dots, \theta_s^*)$  such that Assumption 3.4 stands for the set  $\mathcal{Q}^*$ , we show that the application  $\beta \mapsto \beta\Phi_T(\vartheta^*)$  defined from  $\mathbb{R}^s$  to  $L^2(\lambda_T)$  is injective.

We have that  $\|\beta\Phi_T(\vartheta^*)\|_{L^2(\lambda_T)} = \beta\Gamma\beta^\top$ , where  $\Gamma \in \mathbb{R}^{s \times s}$  is the symmetric matrix defined by  $\Gamma_{k,\ell} = \mathcal{K}_T(\theta_k^*, \theta_\ell^*)$ . Let  $\lambda_{\min}$  be the smallest eigenvalue of  $\Gamma$ . Using Gershgorin's theorem and the definition of  $\mathcal{V}_T$  given by (22), we have that:

$$\begin{aligned} \lambda_{\min} &\geq 1 - \max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s |\mathcal{K}_T(\theta_\ell^*, \theta_k^*)| \\ &\geq 1 - \max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s \left| F \left( \frac{|\theta_\ell^* - \theta_k^*|}{\sigma_T} \right) \right| - (s-1)\mathcal{V}_T. \end{aligned}$$

The separation condition from Point (iv) of Assumption 3.4 implies that for all  $k, \ell \in \{1, \dots, s\}$  such that  $k \neq \ell$  we have  $|\theta_k^* - \theta_\ell^*| \geq \sigma_T \Sigma(\eta, r, s) \geq 8\sigma_T \delta(\eta H_\infty^{(2)}(r), s)$ . Recall the definition of  $\delta(u, s)$  given by (25). We deduce that:

$$\max_{1 \leq \ell \leq s} \sum_{k=1, k \neq \ell}^s \left| F \left( \frac{|\theta_\ell^* - \theta_k^*|}{\sigma_T} \right) \right| \leq \eta H_\infty^{(2)}(r).$$

By Point (iii) of Assumption 3.4, we have  $(s-1)\mathcal{V}_T \leq (1-\eta)H_\infty^{(2)}(r)$  and  $H_\infty^{(2)}(r) \leq 1/6$ . Thus, we get:

$$\lambda_{\min} \geq 5/6. \quad (69)$$

Hence, the symmetric matrix  $\Gamma$  is positive-definite. This proves that the application  $\beta \mapsto \beta \Phi_T(\vartheta^*)$  is injective from  $\mathbb{R}^s$  to  $L^2(\lambda_T)$ . By symmetry, we obtain for  $s^0 \geq 1$  that the application  $\beta \mapsto \beta \Phi_T(\vartheta^0)$  is injective from  $\mathbb{R}^{s^0}$  to  $L^2(\lambda_T)$ .

If  $s = 0$ , we have  $\beta^* \Phi_T(\vartheta^*) = 0$ . For  $s^0 \geq 1$ , we have  $\beta^0 \in (\mathbb{R}^*)^{s^0}$  and since  $\beta \mapsto \beta \Phi_T(\vartheta^0)$  is injective, we deduce that  $\beta^0 \Phi_T(\vartheta^0) \neq 0$ . Thus,  $s = 0$  and  $\beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0)$  implies that  $s^0 = 0$ . By symmetry,  $s^0 = 0$  and  $\beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0)$  implies also that  $s = 0$ .

Assume from now on that  $s, s^0 \in \mathbb{N}^*$  and that  $\beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0)$ . Let us consider the application  $v : \mathcal{Q}^* \mapsto \{-1, 1\}$  defined by:  $v(\theta_k^*) = \text{sgn}(\beta_k^*)$  for any  $k \in \{1, \dots, s\}$ . According to Lemma 8.1, there exists  $p^* \in L^2(\lambda_T)$  such that:

$$\|\beta^*\|_{\ell_1} = \sum_{k=1}^s \beta_k^* \langle \phi_T(\theta_k^*), p^* \rangle_{L^2(\lambda_T)} = \langle \beta^* \Phi_T(\vartheta^*), p^* \rangle_{L^2(\lambda_T)}.$$

Using the fact that  $\beta^* \Phi_T(\vartheta^*) = \beta^0 \Phi_T(\vartheta^0)$  and Properties (i) and (ii) of  $p^*$  in Lemma 8.1, we get:

$$\|\beta^*\|_{\ell_1} = \sum_{k=1}^{s^0} \beta_k^0 \langle \phi_T(\theta_k^0), p^* \rangle_{L^2(\lambda_T)} \leq \|\beta^0\|_{\ell_1}. \quad (70)$$

The role of  $(\beta^*, \vartheta^*)$  and  $(\beta^0, \vartheta^0)$  being symmetric, we also get  $\|\beta^0\|_{\ell_1} \leq \|\beta^*\|_{\ell_1}$ . Hence, we have  $\|\beta^0\|_{\ell_1} = \|\beta^*\|_{\ell_1}$  and  $\text{sgn}(\beta_k^0) = \langle \phi_T(\theta_k^0), p^* \rangle_{L^2(\lambda_T)}$  for  $k \in \{1, \dots, s^0\}$ . Using Properties (i) and (ii) of  $p^*$  in Lemma 8.1, we remark that for any  $\theta \notin \mathcal{Q}^*$

$$\left| \langle \phi_T(\theta), p^* \rangle_{L^2(\lambda_T)} \right| < 1.$$

Thus, we deduce from (70) that  $\mathcal{Q}^0 \subseteq \mathcal{Q}^*$  and by symmetry  $\mathcal{Q}^0 = \mathcal{Q}^*$ . Hence, we obtain  $\vartheta^* = \vartheta^0$  (up to a permutation on the components of  $\vartheta^*$ ) and  $s = s^0$ . Then use the injectivity of the function  $\beta \mapsto \beta \Phi_T(\vartheta^*)$  to get that  $\beta^* = \beta^0$  (up to the same permutation). This finishes the proof of the Lemma.

### 8.2. Proof of Theorem 4.1

We give a bound of the type I error probability. Using that under  $H_0$  we have  $y = \beta^0 \Phi_T(\vartheta^0) + w_T$ , we get:

$$\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_1}(t)] = \mathbb{P}\left(\left|\|w_T\|_{L^2(\lambda_T)}^2 - \mathbb{E}\left[\|w_T\|_{L^2(\lambda_T)}^2\right]\right| > t\right).$$

Using Chebyshev's inequality, we obtain:

$$\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_1}(t)] \leq \frac{\Xi_T}{t^2}. \quad (71)$$

We now give a bound of the type II error probability. We set:

$$R = \|\beta^0 \Phi_T(\vartheta^0) - \beta^* \Phi_T(\vartheta^*)\|_{L^2(\lambda_T)},$$

where  $(\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*)$ . Using the decomposition of  $y$  from the model (2) and the triangle inequality, we have:

$$\begin{aligned} |\text{Test}_1| \geq R^2 - \left| \|w_T\|_{L^2(\lambda_T)}^2 - \mathbb{E}[\|w_T\|_{L^2(\lambda_T)}^2] \right| \\ - 2 \left| \langle \beta^0 \Phi_T(\vartheta^0) - \beta^* \Phi_T(\vartheta^*), w_T \rangle_{L^2(\lambda_T)} \right|. \end{aligned}$$

Notice that by Assumption 1.1, the random variable

$$\langle \beta^0 \Phi_T(\vartheta^0) - \beta^* \Phi_T(\vartheta^*), w_T \rangle_{L^2(\lambda_T)},$$

is Gaussian with zero mean and variance bounded by  $\bar{\sigma}^2 \Delta_T R^2$ . Hence, using that under  $H_1(\rho)$  we have  $R \geq \rho$ , we obtain:

$$\begin{aligned} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi_{\text{Test}_1}(t)] \leq \mathbb{P}\left((\rho^2 - t)/2 \leq \left| \|w_T\|_{L^2(\lambda_T)}^2 - \mathbb{E}[\|w_T\|_{L^2(\lambda_T)}^2] \right|\right) \\ + \mathbb{P}\left((R^2 - t)/2 \leq 2\bar{\sigma}\sqrt{\Delta_T} R |G|\right), \end{aligned} \quad (72)$$

where  $G$  is a standard Gaussian random variable. On the one hand, for  $t < \rho^2$ , using Chebyshev's inequality we get:

$$\mathbb{P}\left((\rho^2 - t)/2 \leq \left| \|w_T\|_{L^2(\lambda_T)}^2 - \mathbb{E}[\|w_T\|_{L^2(\lambda_T)}^2] \right|\right) \leq \frac{4\Xi_T}{(\rho^2 - t)^2}. \quad (73)$$

On the other hand, we have:

$$\mathbb{P}\left((R^2 - t)/2 \leq 2\bar{\sigma}\sqrt{\Delta_T} R |G|\right) \leq \mathbb{P}\left(\frac{\rho^2 - t}{4\bar{\sigma}\sqrt{\Delta_T}\rho} \leq |G|\right) \leq e^{-(\rho^2 - t)^2 / (32\bar{\sigma}^2 \Delta_T \rho^2)}. \quad (74)$$

where we used that  $\rho \leq R$  and the tail bound (see [1, Formula 7.1.13]):

$$\frac{1}{\sqrt{2\pi}} \int_u^{+\infty} e^{-t^2/2} dt \leq \frac{1}{2} e^{-u^2/2}, \quad \text{for } u > 0. \quad (75)$$

By combining (72) with (73) and (74), we get the following bound on the type II error probability:

$$\mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi_{\text{Test}_1}(t)] \leq \frac{4\Xi_T}{(\rho^2 - t)^2} + e^{-(\rho^2 - t)^2/(32\bar{\sigma}^2 \Delta_T \rho^2)}. \quad (76)$$

Then, by putting together (71) and (76), we obtain (38).

### 8.3. Proof of Theorem 4.3

**Case  $s > 0$ .** Let  $(\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*)$ . We consider the estimators  $(\hat{\beta}, \hat{\vartheta})$  defined in (28). Notice that the hypotheses of Theorem 3.5 are in force. We use the constants  $\mathcal{C}_0, \mathcal{C}_1, \mathcal{C}_2$  defined therein. Under  $H_0$ , we have  $s = s^0$ . Thus, for  $\sqrt{t} \geq \mathcal{C}_0 \sqrt{s} \kappa$ , we get the following bound on the type I error probability:

$$\begin{aligned} \mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_2}(t)] &\leq \mathbb{P}\left(\left\|\hat{\beta}\Phi_T(\hat{\vartheta}) - \beta^*\Phi_T(\vartheta^*)\right\|_{L^2(\lambda_T)} > \mathcal{C}_0 \sqrt{s} \kappa\right) \\ &\leq \mathcal{C}_2 \left(\frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau}\right), \end{aligned} \quad (77)$$

where we used that  $\beta^0\Phi_T(\vartheta^0) = \beta^*\Phi_T(\vartheta^*)$  and that  $\sqrt{t} \geq \mathcal{C}_0 \sqrt{s} \kappa$  for the first inequality and Theorem 3.5 for the second.

We now bound the type II error probability. Under  $H_1(\rho)$ , since

$$\left\|\beta^*\Phi_T(\vartheta^*) - \beta^0\Phi_T(\vartheta^0)\right\|_{L^2(\lambda_T)} \geq \rho,$$

we obtain that:

$$\begin{aligned} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi_{\text{Test}_2}(t)] &\leq \mathbb{P}\left(\rho - \sqrt{t} \leq \left\|\hat{\beta}\Phi_T(\hat{\vartheta}) - \beta^*\Phi_T(\vartheta^*)\right\|_{L^2(\lambda_T)}\right) \\ &\leq \mathcal{C}_2 \left(\frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau}\right), \end{aligned} \quad (78)$$

where we used the triangle inequality for the first inequality and Theorem 3.5 as well as  $\rho - \sqrt{t} \geq \mathcal{C}_0 \sqrt{s} \kappa$  for the second.

**Case  $s = 0$ .** Since  $s = 0$ , we have  $y = w_T$  according to (2). Let us first bound the type I error probability  $\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_2}(t)]$ . Assume that the hypothesis  $H_0$  holds so that  $s = s^0 = 0$ . By definition we have:

$$\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_2}(t)] = \mathbb{P}\left(\left\|\hat{\beta}\Phi_T(\hat{\vartheta})\right\|_{L^2(\lambda_T)}^2 > t\right).$$

We get from the definition of the estimators  $\hat{\beta}$  and  $\hat{\vartheta}$  from (28) that:

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

By rearranging some terms in the equation above, we get:

$$\begin{aligned} \frac{1}{2} \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_{L^2(\lambda_T)}^2 &\leq \left\langle \hat{\beta} \Phi_T(\hat{\vartheta}), w_T \right\rangle_{L^2(\lambda_T)} - \kappa \left\| \hat{\beta} \right\|_{\ell_1} \\ &\leq \left\| \hat{\beta} \right\|_{\ell_1} \left( \sup_{\theta \in \Theta_T} |\langle \phi_T(\theta), w_T \rangle_{L^2(\lambda_T)}| - \kappa \right). \end{aligned} \quad (79)$$

Let us define the event:

$$\mathcal{A} = \left\{ \sup_{\theta \in \Theta_T} |\langle \phi_T(\theta), w_T \rangle_{L^2(\lambda_T)}| < \kappa \right\}. \quad (80)$$

We deduce from (79) that on the event  $\mathcal{A}$  we have  $\left\| \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_{L^2(\lambda_T)} = 0$ . Therefore we get:

$$\mathbb{E}_{(\beta^0, \vartheta^0)}[\Psi_{\text{Test}_2}(t)] \leq \mathbb{P} \left( \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_{L^2(\lambda_T)} > 0 \right) \leq \mathbb{P}(\mathcal{A}^c). \quad (81)$$

We shall bound later  $\mathbb{P}(\mathcal{A}^c)$ , see (83).

We now consider the type II error probability. We assume  $H_1$ , that is

$$\left\| \beta^0 \Phi_T(\vartheta^0) \right\|_{L^2(\lambda_T)} \geq \rho.$$

We obtain:

$$\begin{aligned} \mathbb{E}_{(\beta^*, \vartheta^*)}[1 - \Psi_{\text{Test}_2}(t)] &= \mathbb{P} \left( \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) - \beta^0 \Phi_T(\vartheta^0) \right\|_{L^2(\lambda_T)} \leq \sqrt{t} \right) \\ &\leq \mathbb{P} \left( \rho - \sqrt{t} \leq \left\| \hat{\beta} \Phi_T(\hat{\vartheta}) \right\|_{L^2(\lambda_T)} \right) \leq \mathbb{P}(\mathcal{A}^c). \end{aligned} \quad (82)$$

where we used the definition of  $\text{Test}_2$  and the triangle inequality for the first inequality, the second inequality of (81) as well as  $\rho - \sqrt{t} > 0$  for the second.

We shall apply [6, Lemma A.1] to bound  $\mathbb{P}(\mathcal{A}^c)$ . It amounts to controlling the supremum of the Gaussian process  $\theta \mapsto \langle \phi_T(\theta), w_T \rangle_{L^2(\lambda_T)}$ . Recall that Assumptions 3.1 and 3.2 hold. The function  $\phi_T$  is of class  $\mathcal{C}^1$  from the interval  $\Theta_T$  to  $L^2(\lambda_T)$ , with  $\Theta_T$  a sub-interval of  $\Theta$ . We have also, with  $\phi_T^{[1]} = \tilde{D}_{1, \kappa_T}[\phi_T]$ , that:

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

Since Assumption 1.1 on the noise  $w_T$  holds, the hypotheses of [6, Lemma A.1] hold and we deduce from [6, Lemma A.1] (with  $C_1 = C_2 = 1$  therein) that:

$$\begin{aligned} \mathbb{P}(\mathcal{A}^c) &= \mathbb{P} \left( \sup_{\theta \in \Theta_T} |\langle \phi_T(\theta), w_T \rangle_{L^2(\lambda_T)}| \geq \kappa \right) \\ &\leq 3 \cdot \left( \frac{2\bar{\sigma} \sqrt{g_\infty} |\Theta_T| \sqrt{\Delta_T}}{\sigma_T \kappa} \vee 1 \right) e^{-\kappa^2 / (4\bar{\sigma}^2 \Delta_T)}, \end{aligned}$$

where the diameter  $|\Theta_T|_{\mathfrak{d}_T}$  of the set  $\Theta_T$  with respect to the metric  $\mathfrak{d}_T$  is bounded by  $2\sqrt{g_\infty}|\Theta_T|/\sigma_T$  using (21) and the fact that  $C_T \leq 2$ . By taking  $\kappa \geq 2\bar{\sigma}\sqrt{\Delta_T \log(\tau)}$ , we get:

$$\mathbb{P}(\mathcal{A}^c) = \mathbb{P}\left(\sup_{\theta \in \Theta_T} |\langle \phi_T(\theta), w_T \rangle_{L^2(\lambda_T)}| \geq \kappa\right) \leq 3 \cdot \left(\frac{\sqrt{g_\infty}|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau}\right). \quad (83)$$

Notice that the constant  $C_2$  from Theorem 3.5 is equal to  $2\sqrt{g_\infty}C'_2$  where  $C'_2$  is given by [6,  $C_2$  from Eq. (84) therein] and is greater than 3. The constant  $C_2$  depends only on  $r$  and the function  $F$ . Finally, by putting together (77), (78), (81) and (82), we obtain for  $\kappa \geq C_1\bar{\sigma}\sqrt{\Delta_T \log(\tau)}$  (where the constant  $C_1$  is defined in [6, Proof of Theorem 2.1 (p.32)] and is superior to 4) the bound on the maximal testing risk from Theorem 4.3. This finishes the proof.

#### 8.4. Proof of Theorem 5.2

This proof is based on the certificate function. Following [6], we give the existence and properties of the interpolating certificate function.

**Lemma 8.1** (Interpolating certificate). *Let  $T \in \mathbb{N}$ ,  $s \in \mathbb{N}^*$ ,  $\eta \in (0, 1)$ ,  $r \in (0, 1/\sqrt{2g_\infty L_2})$  and  $\mathcal{Q} = \{\theta_1, \dots, \theta_s\} \subset \Theta_T$ . Suppose that Assumption 3.4 holds.*

*Then, there exist finite positive constants  $C_N, C_F, C_B$  with  $C_F < 1$ , depending on  $r$  and the function  $F$ , such that for any application  $v : \mathcal{Q} \mapsto \{-1, 1\}$ , there exist unique  $\alpha, \xi \in \mathbb{R}^s$  such that  $p \in L^2(\lambda_T)$  uniquely defined by:*

$$\begin{cases} p = \sum_{k=1}^s \alpha_k \phi_T(\theta_k) + \sum_{k=1}^s \xi_k \tilde{D}_{1,T}[\phi_T](\theta_k), \\ \langle \phi_T(\theta), p \rangle_{L^2(\lambda_T)} = v(\theta) \quad \text{and} \quad \langle \partial_\theta \phi_T(\theta), p \rangle_{L^2(\lambda_T)} = 0, \quad \text{for all } \theta \in \mathcal{Q}, \end{cases} \quad (84)$$

satisfies:

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

$$|\langle \phi_T(\theta'), p \rangle_{L^2(\lambda_T)}| \leq 1 - C_N \mathfrak{d}_T(\theta, \theta')^2.$$

(ii) For all  $\theta$  in  $\Theta_T$ ,  $\theta \notin \bigcup_{\theta' \in \mathcal{Q}} \mathcal{B}_T(\theta', r)$  (far region), we have:

$$|\langle \phi_T(\theta), p \rangle_{L^2(\lambda_T)}| \leq 1 - C_F.$$

(iii) We have  $\|p\|_{L^2(\lambda_T)} \leq \sqrt{s} C_B$ .

*Proof.* Using similar arguments as those developed in the proof of Theorem 3.5, we get that all the hypotheses of [6, Proposition, 7.4] are satisfied. The existence and uniqueness of  $p$  is then guaranteed by [6, Lemma, 10.1]. The properties satisfied by  $p$  are direct consequences of [6, Proposition, 7.4].  $\square$

Recall the test problem given by (54). Assumption 3.4 holds for the set  $\mathcal{Q}^0$ . Thanks to Lemma 8.1, the element  $p_0$  of  $L^2(\lambda_T)$  is uniquely defined by  $v^0$ , (55) and (56). Hence, the test statistic  $\text{Test}_3$  from (57) is well-defined.

We first bound the type I error probability. Let us fix  $(\beta^*, \vartheta^*) \in (\mathbb{R}^*)^s \times \Theta_T^s(\delta^*)$  such that  $H_0$  holds. Using that  $y = \beta^* \Phi_T(\vartheta^*) + w_T$  and the triangle inequality, we obtain:

$$\begin{aligned} |\text{Test}_3| &= \left| \left\| \hat{\beta} \right\|_{\ell_1} - \|\beta^*\|_{\ell_1} + \|\beta^*\|_{\ell_1} - \langle \beta^* \Phi_T(\vartheta^*), p_0 \rangle_{L^2(\lambda_T)} - \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| \\ &\leq \left| \left\| \hat{\beta} \right\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| + |B| + \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right|, \end{aligned} \quad (85)$$

where:

$$B = \|\beta^*\|_{\ell_1} - \langle \beta^* \Phi_T(\vartheta^*), p_0 \rangle_{L^2(\lambda_T)}. \quad (86)$$

Since  $\mathcal{Q}^{*,+} \subseteq \mathcal{Q}^{0,+}$ ,  $\mathcal{Q}^{*,-} \subseteq \mathcal{Q}^{0,-}$ , we have for all  $k \in \{1, \dots, s\}$ :

$$|\beta_k^*| - \langle \beta_k^* \phi_T(\theta_k^*), p_0 \rangle_{L^2(\lambda_T)} = 0,$$

we deduce that  $B = 0$  under  $H_0$ . Hence, we have that:

$$\mathbb{E}_{(\beta^*, \vartheta^*)}[\Psi_{\text{Test}_3}(t)] \leq \mathbb{P} \left( \left| \left\| \hat{\beta} \right\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| > t/2 \right) + \mathbb{P} \left( \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| > t/2 \right). \quad (87)$$

Recall that under  $H_0$ , we have  $s \leq s^0$ . Therefore, since  $\mathcal{C}_3 \kappa s^0 \leq t/2$ , we have  $\mathcal{C}_3 \kappa s \leq t/2$ . We get from Theorem 3.5 that:

$$\mathbb{P} \left( \left| \left\| \hat{\beta} \right\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right| > t/2 \right) \leq \mathcal{C}_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right). \quad (88)$$

Then, thanks to Assumptions 1.1 and Lemma 8.1, the quantity  $\langle w_T, p_0 \rangle_{L^2(\lambda_T)}$  is a centered Gaussian random variable of variance bounded by  $\bar{\sigma}^2 C_B^2 \Delta_T s_0$  where  $C_B$  is the finite positive constant from Lemma 8.1. Hence we have, provided that  $t \geq 2\mathcal{C}_3 \kappa s^0$  with  $\kappa \geq \mathcal{C}_1 \bar{\sigma} \sqrt{\Delta_T \log(\tau)}$ , that is,  $t^2 \geq (2\mathcal{C}_1 \mathcal{C}_3 \bar{\sigma} s_0)^2 \Delta_T \log(\tau)$ :

$$\begin{aligned} \mathbb{P} \left( \langle w_T, p_0 \rangle_{L^2(\lambda_T)} > t/2 \right) &\leq \int_{t/2}^{+\infty} \frac{e^{-x^2/(2\bar{\sigma}^2 \Delta_T C_B^2 s_0)}}{\sqrt{2\pi \bar{\sigma}^2 \Delta_T C_B^2 s_0}} dx \\ &\leq \frac{1}{2} e^{-\frac{t^2}{8(\bar{\sigma}^2 \Delta_T C_B^2 s_0)}} \leq \frac{1}{2\tau a s_0}, \end{aligned}$$

with  $a = (\mathcal{C}_1 \mathcal{C}_3 / C_B)^2 / 2$  and where we used the tail bound (75). It gives by symmetry that:

$$\mathbb{P} \left( \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| > t/2 \right) \leq \frac{1}{\tau a s_0}. \quad (89)$$



Plugging (88) and (89) in (87), we get:

$$\sup_{(\beta^*, \vartheta^*) \in H_0} \mathbb{E}_{(\beta^*, \vartheta^*)} [\Psi_{\text{Test}_3}(t)] \leq C_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right) + \frac{1}{\tau^{a_{s_0}}}. \quad (90)$$

We now bound the type II error probability. Assume that  $H_1$  holds, that is  $\mathcal{D}_{T,r}(\beta^*, \vartheta^*, v^0, \vartheta^0) \geq \rho$ . We have, using the first equality of (85) and the triangle inequality, that:

$$|\text{Test}_3| \geq |B| - \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| - \left| \|\hat{\beta}\|_{\ell_1} - \|\beta^*\|_{\ell_1} \right|,$$

with  $B$  defined in (86). Using the definitions (53) of  $S(r)$  and  $S_k^\epsilon(r)$  with  $\epsilon \in \{+, -\}$  and  $k \in \mathcal{I}^\epsilon$ , we get:

$$\begin{aligned} B = & \sum_{\substack{\epsilon \in \{+, -\} \\ k \in \mathcal{I}^\epsilon, \ell \in S_k^\epsilon(r)}} |\beta_\ell^*| \left( 1 - \text{sgn}(\beta_\ell^*) \langle \phi_T(\theta_\ell^*), p_0 \rangle_{L^2(\lambda_T)} \right) \\ & + \sum_{k \in S(r)^c} |\beta_k^*| \left( 1 - \text{sgn}(\beta_k^*) \langle \phi_T(\theta_k^*), p_0 \rangle_{L^2(\lambda_T)} \right). \end{aligned}$$

Thanks to Lemma 8.1 (i)-(ii) of , we obtain:

$$\begin{aligned} B & \geq \sum_{\substack{\epsilon \in \{+, -\} \\ k \in \mathcal{I}^\epsilon, \ell \in S_k^\epsilon(r)}} C_N |\beta_\ell^*| \mathfrak{d}_T(\theta_\ell^*, \theta_k^0)^2 + \sum_{k \in S(r)^c} C_F |\beta_k^*| \\ & \geq (C_N \wedge C_F) \mathcal{D}_{T,r}(\beta^*, \vartheta^*, v^0, \vartheta^0) \geq (C_N \wedge C_F) \rho, \end{aligned}$$

where the constants  $C_N$  and  $C_F$  are defined in Lemma 8.1 and depend on  $r$  and on the function  $F$ . Therefore, we have with  $a_t = (C_N \wedge C_F) \rho - t$ :

$$\begin{aligned} \mathbb{E}_{(\beta^*, \vartheta^*)} [1 - \Psi_{\text{Test}_3}(t)] & \leq \mathbb{P} \left( \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| + \left| \|\beta^*\|_{\ell_1} - \|\hat{\beta}\|_{\ell_1} \right| \geq a_t \right) \\ & \leq \mathbb{P} \left( \left| \langle w_T, p_0 \rangle_{L^2(\lambda_T)} \right| \geq a_t/2 \right) \\ & \quad + \mathbb{P} \left( \left| \|\beta^*\|_{\ell_1} - \|\hat{\beta}\|_{\ell_1} \right| \geq a_t/2 \right). \end{aligned}$$

Provided that  $\rho \geq \mathcal{C}_4 s \kappa + \mathcal{C}_5 t$  with  $\mathcal{C}_4 = 2\mathcal{C}_3/(C_N \wedge C_F)$  and  $\mathcal{C}_5 = 2/(C_N \wedge C_F)$  we have  $a_t/2 \geq (\mathcal{C}_3 \kappa s) \vee (t/2)$ . By using (88) and (89), we obtain:

$$\sup_{(\beta^*, \vartheta^*) \in H_1(\rho)} \mathbb{E}_{(\beta^*, \vartheta^*)} [1 - \Psi_{\text{Test}_3}(t)] \leq C_2 \left( \frac{|\Theta_T|}{\sigma_T \tau \sqrt{\log(\tau)}} \vee \frac{1}{\tau} \right) + \frac{1}{\tau^{a_{s_0}}}. \quad (91)$$

Finally, by adding both sides of (90) and (91), we get (59). This concludes the proof.

### 8.5. Proof of Lemma 7.2

It is easy to check that the functions  $g_T$  and  $g_{\mathcal{K}_T^{\text{prox}}}$  are constant functions with:

$$g_T = g_\infty (T^2 - 1) \quad \text{and} \quad g_{\mathcal{K}_T^{\text{prox}}} = g_\infty T^2. \quad (92)$$

Thus, we easily deduce the second inequality of (67) from the definition (20) of  $C_T$ .

We now consider the bound on  $\mathcal{V}_T$ . For  $i, j \in \{0, \dots, 3\}$  and  $\ell = i + j$ , we have with  $\alpha_T = 1 - 1/T^2$ :

$$\sup_{\Theta^2} |\mathcal{K}_T^{[i,j]} - \mathcal{K}_T^{\text{prox}[i,j]}| = g_\infty^{-\ell/2} (T^2 \alpha_T)^{-\ell/2} A_{\ell,T}, \quad (93)$$

where

$$A_{\ell,T} = \sup_{t \in [-\frac{1}{2}, \frac{1}{2}]} \left| \partial_t^\ell \left[ D_T(t) + \left(1 - \alpha_T^{\ell/2}\right) \frac{\sin(T\pi t)}{T\pi t} \right] \right|,$$

and, for  $t \in [-1/2, 1/2]$  and the convention  $J(0) = 0$ :

$$D_T(t) = \frac{\sin(T\pi t)}{T} J(t) \quad \text{and} \quad J(t) = \frac{1}{\sin(\pi t)} - \frac{1}{\pi t}.$$

It is easy to check that the function  $J$  can be expanded as a power series at 0 with positive convergence radius, and thus is of class  $\mathcal{C}^\infty$  on  $[-1/2, 1/2]$ . Thus the following constant is finite:

$$M = \sup_{0 \leq \ell \leq 6} \sup_{t \in [-1/2, 1/2]} |J^{(\ell)}| < +\infty.$$

Using the Leibniz rule, we have that for  $\ell \in \{1, \dots, 6\}$  and  $t \in [-1/2, 1/2]$ :

$$|\partial_t^\ell D_T(t)| = \frac{1}{T} \left| \sum_{j=0}^{\ell} \binom{\ell}{j} (T\pi)^j \sin^{(j)}(T\pi t) J^{(\ell-j)}(t) \right| \leq M \frac{(T\pi + 1)^\ell}{T}.$$

We deduce from (93) that for  $i, j \in \{0, \dots, 3\}$  and  $\ell = i + j$ :

$$\begin{aligned} \sup_{\Theta^2} |\mathcal{K}_T^{[i,j]} - \mathcal{K}_T^{\text{prox}[i,j]}| &\leq g_\infty^{-\ell/2} (T^2 \alpha_T)^{-\ell/2} \left( M \frac{(T\pi + 1)^\ell}{T} + (1 - \alpha_T^{\ell/2}) \right) \\ &\leq M 3^\ell T^{-1}, \end{aligned}$$

where we used that  $T \geq 3$  and  $g_\infty \alpha_T \geq 1$ , and that  $1 - \alpha_T^{\ell/2} = 0$  for  $\ell = 0$ . Recall the definition (22) of  $\mathcal{V}_T$  to get  $\mathcal{V}_T \leq M 3^\ell T^{-1}$ . This finishes the proof.

## 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](#)
- [2] ARIAS-CASTRO, E., CANDÈS, E. J. and PLAN, Y. (2011). Global testing under sparse alternatives: ANOVA, multiple comparisons and the higher criticism. *Ann. Statist.* **39** 2533–2556. [MR2906877](#)
- [3] BARAUD, Y. (2002). Non-asymptotic minimax rates of testing in signal detection. *Bernoulli* **8** 577–606. [MR1935648](#)
- [4] BOYER, C., DE CASTRO, Y. and SALMON, J. (2017). Adapting to unknown noise level in sparse deconvolution. *Inf. Inference* **6** 310–348. [MR3764527](#)
- [5] 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.
- [6] BUTUCEA, C., DELMAS, J.-F., DUTFOY, A. and HARDY, C. (2022). Off-the-grid learning of sparse mixtures from a continuous dictionary. *arXiv preprint arXiv:2207.00171*.
- [7] CANDÈS, E. J. and FERNANDEZ-GRANDA, C. (2013). Super-resolution from noisy data. *J. Fourier Anal. Appl.* **19** 1229–1254. [MR3132912](#)
- [8] CANDÈS, E. J. and FERNANDEZ-GRANDA, C. (2014). Towards a mathematical theory of super-resolution. *Comm. Pure Appl. Math.* **67** 906–956. [MR3193963](#)
- [9] CANDÈS, E. J. and PLAN, Y. (2011). A probabilistic and RIPless theory of compressed sensing. *IEEE Trans. Inform. Theory* **57** 7235–7254. [MR2883653](#)
- [10] DE CASTRO, Y. and GAMBOA, F. (2012). Exact reconstruction using Beurling minimal extrapolation. *J. Math. Anal. Appl.* **395** 336–354. [MR2943626](#)
- [11] DUVAL, V. and PEYRÉ, G. (2015). Exact support recovery for sparse spikes deconvolution. *Found. Comput. Math.* **15** 1315–1355. [MR3394712](#)
- [12] ERMAKOV, M. S. (1990). Minimax detection of a signal in Gaussian white noise. *Teor. Veroyatnost. i Primenen.* **35** 704–715. [MR1090496](#)
- [13] 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. [MR3588285](#)
- [14] INGSTER, Y. I. and SUSLINA, I. A. (2003). *Nonparametric goodness-of-fit testing under Gaussian models*. *Lecture Notes in Statistics* **169**. Springer-Verlag, New York. [MR1991446](#)
- [15] INGSTER, Y. I., TSYBAKOV, A. B. and VERZELEN, N. (2010). Detection boundary in sparse regression. *Electron. J. Stat.* **4** 1476–1526. [MR2747131](#)
- [16] LAURENT, B., LOUBES, J.-M. and MARTEAU, C. (2012). Non asymptotic minimax rates of testing in signal detection with heterogeneous variances. *Electron. J. Stat.* **6** 91–122. [MR2879673](#)

- [17] MALLAT, S. (2009). *A wavelet tour of signal processing : the sparse way*, Third ed. Elsevier/Academic Press, Amsterdam With contributions from Gabriel Peyré. [MR2479996](#)
- [18] POON, C., KERIVEN, N. and PEYRÉ, G. (2021). The geometry of off-the-grid compressed sensing. *Foundations of Computational Mathematics*.
- [19] TANG, G., BHASKAR, B. N. and RECHT, B. (2015). Near minimax line spectral estimation. *IEEE Trans. Inform. Theory* **61** 499–512. [MR3299978](#)
- [20] TANG, G., BHASKAR, B. N., SHAH, P. and RECHT, B. (2013). Compressed sensing off the grid. *IEEE Trans. Inform. Theory* **59** 7465–7490. [MR3124655](#)
- [21] TROPP, J. A. (2004). Greed is good: algorithmic results for sparse approximation. *IEEE Trans. Inform. Theory* **50** 2231–2242. [MR2097044](#)