

# On optimality of empirical risk minimization in linear aggregation

ADRIEN SAUMARD

CREST, ENSAI, Université Bretagne Loire, Campus de Ker-Lann, Rue Blaise Pascal – BP 37203 35712 BRUZ cedex, France.

E-mail: [adrien.saumard@ensai.fr](mailto:adrien.saumard@ensai.fr)

In the first part of this paper, we show that the small-ball condition, recently introduced by (*J. ACM* 62 (2015) Art. 21, 25), may behave poorly for important classes of localized functions such as wavelets, piecewise polynomials or for trigonometric polynomials, in particular leading to suboptimal estimates of the rate of convergence of ERM for the linear aggregation problem. In a second part, we recover optimal rates of convergence for the excess risk of ERM when the dictionary is made of trigonometric functions. Considering the bounded case, we derive the concentration of the excess risk around a single point, which is an information far more precise than the rate of convergence. In the general setting of a  $L_2$  noise, we finally refine the small ball argument by rightly selecting the directions we are looking at, in such a way that we obtain optimal rates of aggregation for the Fourier dictionary.

*Keywords:* empirical risk minimization; excess risk's concentration; linear aggregation; optimal rates; small-ball property

## 1. Introduction

Consider the following general regression framework:  $(\mathcal{X}, \mathcal{T}_{\mathcal{X}})$  is a measurable space,  $(X, Y) \in \mathcal{X} \times \mathbb{R}$  is a pair of random variables of joint distribution  $P$  – the marginal of  $X$  being denoted  $P^X$  – and it holds

$$Y = s_*(X) + \sigma(X)\varepsilon, \quad (1.1)$$

where  $s_*$  is the regression function of the response variable  $Y$  with respect to the *random design*  $X$ ,  $\sigma(X) \geq 0$  is the *heteroscedastic noise* level and  $\varepsilon$  is the conditionally standardized noise, satisfying  $\mathbb{E}[\varepsilon|X] = 0$  and  $\mathbb{E}[\varepsilon^2|X] = 1$ . Relation (1.1) is very general and is indeed satisfied as soon as  $\mathbb{E}[Y^2] < +\infty$ . In this case,  $s_* \in L_2(P^X)$  is the orthogonal projection of  $Y$  onto the space of  $X$ -measurable functions. In particular, no restriction is made on the structure of dependence between  $Y$  and  $X$ .

We thus face a typical *learning problem*, where the statistical modelling is minimal, and the goal will be, given a sample  $(X_i, Y_i)_{i=1}^n$  of law  $P^{\otimes n}$  and a new covariate  $X_{n+1}$ , to predict the value of the associated response variable  $Y_{n+1}$ . More precisely, we want to construct a function  $\hat{s}$ , depending on the data  $(X_i, Y_i)_{i=1}^n$ , such that the least-squares risk  $R(\hat{s}) = \mathbb{E}[(Y_{n+1} - \hat{s}(X_{n+1}))^2]$  is as small as possible, the pair  $(X_{n+1}, Y_{n+1})$  being independent of the sample  $(X_i, Y_i)_{i=1}^n$ .

In this paper, we focus on the technique of *linear aggregation* via Empirical Risk Minimization (ERM). This means that we are given a dictionary  $S = \{s_1, \dots, s_D\}$  and that we produce the least-

squares estimator  $\hat{s}_m$  on its linear span  $m = \text{Span}(S)$ ,

$$\hat{s}_m \in \arg \min_{s \in m} R_n(s), \quad \text{where } R_n(s) = \frac{1}{n} \sum_{i=1}^n (Y_i - s(X_i))^2. \quad (1.2)$$

The quantity  $R_n(s)$  is called the empirical risk of the function  $s$ . The accuracy of the method is tackled through an *oracle inequality*, where the risk of the estimator  $R(\hat{s}_m)$  is compared – on an event of probability close to one – to the risk of the best possible function within the linear model  $m$ . The latter function is denoted  $s_m$  and is called the oracle, or the (orthogonal) *projection* of the regression function  $s_*$  onto  $m$ ,

$$s_m \in \arg \min_{s \in m} R(s).$$

An oracle inequality then writes, on an event  $\Omega_0$  of probability close to one,

$$R(\hat{s}_m) \leq R(s_m) + r_n(D), \quad (1.3)$$

for a positive residual term  $r_n(D)$ . An easy and classical computation gives that the *excess risk* satisfies  $R(\hat{s}_m) - R(s_m) = \|\hat{s}_m - s_m\|_2^2$ , where  $\|\cdot\|_2$  is the natural quadratic norm in  $L_2(P^X)$ , associated with the scalar product  $\langle f, g \rangle = \int f(x)g(x) dP^X(x)$ . Hence, inequality (1.3) can be rewritten as  $\|\hat{s}_m - s_m\|_2^2 \leq r_n(D)$  and the quantity  $r_n(D)$  thus corresponds to the rate of estimation of the projection  $s_m$  by the least-squares estimator  $\hat{s}_m$  in terms of excess risk, corresponding here to the squared quadratic norm.

The linear aggregation problem has been well studied in various settings linked to nonparametric regression [2,7,26,31] and density estimation [27]. It has been consequently understood that the *optimal rate*  $r_n(D)$  of linear aggregation is of the order of  $D/n$ , where  $D$  is the size of the dictionary. Recently, [19] have shown that ERM is suboptimal for the linear aggregation problem in general, in the sense that there exist a dictionary  $S$  and a pair  $(X, Y)$  of random variables for which the rate of ERM (drastically) deteriorates, even in the case where the response variable  $Y$  and the dictionary are uniformly bounded.

On the positive side, [19] also made a breakthrough by showing that if a so-called *small-ball condition* is achieved with absolute constants, uniformly over the functions in the linear model  $m$ , then the optimal rate is recovered by ERM. We recall and discuss in details the small-ball condition in Section 2, but it is worth mentioning here that one of the main advantages of the small-ball method developed in a series of papers, [15,17,19,21–23] is that it enables to prove sharp bounds under very weak moment conditions and thus to derive results that were *unachievable* with more standard concentration arguments.

In Section 2, we contribute to the growing understanding of this very recent approach by looking at the behavior of the small-ball condition when the dictionary is made of elements of some classical orthonormal bases, such as histograms, piecewise polynomials, wavelets and the Fourier basis. These examples are indeed central in various methods of nonparametric statistics.

It appears that with such functions, the small-ball condition *can't be satisfied with absolute constants* and the resulting bounds obtained in [19] are far from optimal. This lack of accuracy of the small-ball approach seems rather natural for dictionaries that are made of localized functions,

such as wavelets for instance, since these functions are very “picky” and thus hardly identifiable – see Section 2 for a more thorough discussion around these terms.

However, it seems more surprising that the Fourier dictionary also leads to suboptimal rates of linear aggregation when analyzed *via* the small ball method. In fact, the behavior of the small-ball condition on the span of some trigonometric functions is essentially unknown in the literature and this type of information, in the related context of Fourier measurements in compressed sensing, has a potentially significant impact on the theory of Fourier measurements [17].

Nevertheless, we show in Section 3 that ERM achieves optimal rates of linear aggregation, both in the bounded setting and for  $L_2$ -noise. Our result in particular outperforms previously obtained bounds [2].

More precisely, when the response variable  $Y$  is bounded, we derive *concentration inequalities* for the excess risk, which is an information far more precise than the rate of convergence. Our proofs are based on empirical process theory and substantially simplify our previous approach to concentration inequalities for the excess risk on models spanned by localized bases [25,29].

When the noise is only assumed to have a second moment, we prove optimal rates of linear aggregation for the Fourier dictionary by using a refined small-ball argument. Indeed, by imposing a light and natural smoothness condition on the regression function, we localize the analysis by only looking at some directions in the model that satisfy a uniform small-ball condition. It is important to note that such approach was suggested – but not achieved – by Lecué and Mendelson [17] for the study of Fourier measurements in compressed sensing.

Finally, complete proofs are dispatched in Sections 4 and 5, at the end of the paper.

## 2. The small-ball method for classical functional bases

We recall in Section 2.1 one of the main results of [19], linking the small-ball condition to the rate of convergence of ERM in linear aggregation. Then, we show in Section 2.2 that the constants involved in the small-ball condition behave poorly for dictionaries made of localized bases and also for the Fourier basis.

### 2.1. The small-ball condition and the rate of ERM in linear aggregation

Let us first recall the definition of the small-ball condition for a linear span, as exposed in [19].

**Definition 1.** A linear span  $m \subset L_2(P^X)$  is said to satisfy the *small-ball condition* for some positive constants  $\kappa_0$  and  $\beta_0$  if for every  $s \in m$ ,

$$\mathbb{P}\left(|s(X)| \geq \kappa_0 \|s\|_2\right) \geq \beta_0. \quad (2.1)$$

The small-ball condition thus ensures that the functions of the model  $m$  do not put too much weight around zero. From a statistical perspective, it is also explained in [19] that the small-ball condition can be viewed as a quantified version of identifiability of the model  $m$ . A more general small-ball condition – that reduces to the previous definition for linear models – is also available when the model isn’t necessary linear [23].

Under the small-ball condition, [19] derive the following result, describing the rate of convergence of ERM in linear aggregation.

**Theorem 2 ([19]).** *Let  $S = \{s_1, \dots, s_D\} \subset L_2(P^X)$  be a dictionary and assume that  $m = \text{Span}(S)$  satisfies the small-ball condition with constants  $\kappa_0$  and  $\beta_0$  (see Definition 1 above). Let  $n \geq (400)^2 D / \beta_0^2$  and set  $\zeta = Y - s_m(X)$ , where  $s_m$  is the projection of the regression function  $s_*$  onto  $m$ . Assume further that one of the following two conditions holds:*

1.  $\zeta$  is independent of  $X$  and  $\mathbb{E}\zeta^2 \leq \sigma^2$ , or
2.  $|\zeta| \leq \sigma$  almost surely.

*Then the least-squares estimator  $\hat{s}_m$  on  $m$ , defined in (1.2), satisfies for every  $x > 0$ , with probability at least  $1 - \exp(-\beta_0^2 n / 4) - (1/x)$ ,*

$$\|\hat{s}_m - s_m\|_2^2 \leq \left(\frac{16}{\beta_0 \kappa_0^2}\right)^2 \frac{\sigma^2 D x}{n}. \tag{2.2}$$

Notice that Alternative 1 in Theorem 2 is equivalent to assuming that the regression function belongs to  $m$  – that is  $s_* = s_m$  – and that the noise is independent from the design – that is  $\sigma(X) \equiv \sigma$  is homoscedastic and  $\varepsilon$  is independent of  $X$  in relation (1.1).

The main feature of Theorem 2 is that if the small-ball condition is achieved with *absolute constants*  $\kappa_0$  and  $\beta_0$  not depending on the dimension  $D$  nor the sample size  $n$ , then optimal linear aggregation rates of order  $D/n$  are recovered by ERM. If moreover the regression function belongs to  $m$  (Alternative 1), then the only moment assumption required is that the noise is in  $L_2$ . Otherwise, Alternative 2 asks for a uniformly bounded noise. Some variants of Theorem 2 are also presented in [19], showing for instance, that optimal rates can be also derived for ERM when the noise as a fourth moment.

In the analysis of optimal rates in linear aggregation, it is thus worth understanding when the small ball condition stated in Definition 1 is achieved with absolute constants.

One such typical situation is for *linear measurements*, that is when the functions of the dictionary are of the form  $f_i(x) = x^T t_i$ ,  $t_i \in \mathbb{R}^d$ . Indeed, very weak conditions are asked on the design  $X$  in this case to ensure the small-ball property: for instance, it suffices to assume that  $X$  has independent coordinates that are absolutely continuous with respect to the Lebesgue measure, with a density almost surely bounded (see [17] and [23], Section 6, for more details). As shown in [17] and [18], this implies that the small-ball property has important consequences in sparse recovery and analysis of regularized linear regression.

The constants  $(\kappa_0, \beta_0)$  of the small-ball condition influence the rate of convergence exposed in Theorem 2 above through the term  $V_0 := \beta_0^{-2} \kappa_0^{-4}$  and therefore, we will provide upper and lower bounds for  $V_0$  in the following section for various functional dictionaries.

## 2.2. The constants in the small-ball condition for general linear bases

Besides linear measurements discussed in Section 2.1 above, an important class of dictionaries for the linear aggregation problem consists in expansions along orthonormal bases of  $L_2(P^X)$ , which typically correspond to nonparametric estimation.

Our goal in this section is thus to investigate the behavior of the small-ball condition for some classical orthonormal bases such as piecewise polynomial functions, including histograms, wavelets or the Fourier basis.

2.2.1. *Some generic limits for the small-ball method*

Let us begin with a general proposition, describing some upper bounds for the constants  $\kappa_0$  and  $\beta_0$  appearing in the small-ball condition (1). This will enable us to deduce lower bounds for the parameter  $V_0 = \beta_0^{-2} \kappa_0^{-4}$  appearing in the rate (2.2) of Theorem 2 above and therefore, we will have some insights on the limits of the small-ball method for linear aggregation.

One can easily see from its definition that the small-ball condition is more difficult to ensure, at a heuristic level, when the model at hand contains some “picky” functions. The following proposition provides some quantifications of this fact.

**Proposition 3.** *Assume that a model  $m$  satisfies the small-ball condition (1) with constants  $(\beta_0, \kappa_0)$ . Then, it holds*

$$\beta_0 \leq \inf_{f \in m \setminus \{0\}} \mathbb{P}(f(X) \neq 0), \tag{2.3}$$

$$\kappa_0 \leq \inf_{f \in m \setminus \{0\}} \frac{\|f\|_\infty}{\|f\|_2} \tag{2.4}$$

and, for any  $q > 0$ ,

$$\beta_0 \kappa_0^q \leq \inf_{f \in m \setminus \{0\}} \left( \frac{\|f\|_q}{\|f\|_2} \right)^q. \tag{2.5}$$

In particular, we always have  $\beta_0 \kappa_0^2 \leq 1$  and if  $m$  contains the constant functions, then  $\kappa_0 \leq 1$ .

It is interesting to note that Inequalities (2.3) and (2.4) are two limiting cases of (2.5), respectively when  $q \rightarrow 0$  and when  $q \rightarrow +\infty$ . The proof of Proposition 7, which is elementary, is given in Section 4 below.

It is also worth noticing that the inequality  $\beta_0 \kappa_0^2 \leq 1$  implies that the upper bound of Theorem 2 – obtained in [19] – is always greater than  $256\sigma^2 D x/n$ .

Furthermore, consider a model of histograms on a regular partition  $\Pi$  of  $\mathcal{X} = [0, 1]^d$  made of  $D$  pieces,  $X$  being uniformly distributed on  $\mathcal{X}$ . More precisely, for any  $I \in \Pi$ , set

$$s_I = \frac{\mathbf{1}_I}{\sqrt{P^X(I)}} = \sqrt{D} \mathbf{1}_I$$

and take a dictionary  $S = \{s_I; I \in \Pi\}$ , associated to the model  $m = \text{Span}(S)$ .

Then by Inequality (2.3), one directly gets  $\beta_0 \leq D^{-1}$  and as  $m$  contains the constants, it holds  $V_0 = \beta_0^{-2} \kappa_0^{-4} \geq D^2$  and the upper bound (2.2) of Theorem 2 is greater than  $256\sigma^2 D^3 x/n$ . Hence, the rate of convergence exhibited by the small-ball method in the case of regular histograms is  $D^3/n$ , which is suboptimal since it has been proved in [1,29] that the excess risk

concentrates in this case, under Alternative 2 of Theorem 2 above, around a value exactly equal to  $\mathbb{E}[\zeta^2]D/n$ .

More generally, when considering the case of a linear model made of piecewise polynomial functions of degrees bounded by a constant  $r$  on a regular partition, we easily deduce from the previous results on histograms – that is polynomials of degree zero – that  $\beta_0 \leq rD^{-1}$  for any  $\kappa_0 \in (0, 1)$ . We thus have  $V_0 \geq r^{-2}D^2$  and the rate of convergence ensured by Theorem 2 in this case is again proportional to  $D^3/n$ . It is again suboptimal, since it is also proved in [29] that for such models of piecewise polynomial functions, the excess risk concentrates around  $\mathbb{E}[\zeta^2]D/n$ , under Alternative 2 of Theorem 2 above.

Let us now discuss the case of a dictionary made of compactly supported wavelets.

To fix ideas, let us more precisely state some notations (for more details about wavelets, see for instance [11]). We consider in this case that  $\mathcal{X} = [0, 1]$  and  $X$  is uniformly distributed on  $\mathcal{X}$ . Set  $\phi_0$  the father wavelet and  $\psi_0$  the mother wavelet. For every integers  $j \geq 0, 1 \leq k \leq 2^j$ , define

$$\psi_{j,k} : x \mapsto 2^{j/2}\psi_0(2^j x - k + 1). \tag{2.6}$$

As explained in [9], there exists several ways to consider wavelets on the interval. We apply here one of the most classical construction, that consists in using “periodized” wavelets. To this aim, we associate to a function  $\rho$  on  $\mathbb{R}$ , the 1-periodic function

$$\rho^{\text{per}}(x) = \sum_{p \in \mathbb{Z}} \rho(x + p).$$

Notice that if  $\psi$  has a compact support, then the sum at the right-hand side of the latter inequality is finite for any  $x$ .

We set for every integer  $j \geq 0, \Lambda(j) = \{(j, k); 1 \leq k \leq 2^j\}$ . Moreover, we set  $\psi_{-1,1}(x) = \phi_0(x), \Lambda(-1) = \{(-1, 1)\}$  and for any integer  $l \geq 0, \Lambda_l = \bigcup_{j=-1}^l \Lambda(j)$ . Then we consider the dictionary  $S = \{\psi_\lambda^{\text{per}}; \lambda \in \Lambda_l\}$  associated to the model  $m = \text{Span}(S)$ .

Now, it is easily seen from (2.6), that for any  $(l, k) \in \Lambda(l), \mathbb{P}(|\psi_{l,k}(X)| \neq 0) \lesssim 2^{-l} \lesssim D^{-1}$ , where  $D$  is the linear dimension of  $m$ . Consequently, as for histograms and piecewise polynomials on regular partitions, dictionaries made of compactly supported wavelets are handled through the small-ball method with a bound proportional to  $D^3/n$ . This rate is again suboptimal, as shown quite recently by Navarro and Saumard [25], who proved that for such models, the excess risk of the least-squares estimator concentrates, under Alternative 2 of Theorem 2 above, around  $\mathbb{E}[\zeta^2]D/n$ .

It is worth noting that more general multidimensional wavelets could also be considered at the price of more technicalities.

Wavelets, histograms and piecewise polynomials are models that are *formed from* “picky” functions, it is thus quite legitimate that the small-ball method implies suboptimal rates for these models. What happens when the dictionary is formed from spatially unlocalized functions such as the Fourier basis?

**Proposition 4.** *Assume that  $\mathcal{X} = [-\pi, \pi]$  and that the design  $X$  is uniformly distributed on  $\mathcal{X}$ . Then Fourier expansions (i.e., the set of trigonometric polynomials) can not satisfy the small*

ball condition (2.1) with some absolute constants  $(\beta_0, \kappa_0)$ . More precisely, let us set  $\varphi_0 \equiv 1$ ,  $\varphi_{2k}(x) = \sqrt{2} \cos(kx)$  and  $\varphi_{2k+1}(x) = \sqrt{2} \sin(kx)$  for  $k \geq 1$ , and take for some  $l \in \mathbb{N}$ ,  $l \geq 2$ , the model  $m_D = \text{Span}\{\varphi_j; j = 0, \dots, 2l\}$ , of linear dimension  $D = 2l + 1$ . If  $m_D$  satisfies the small ball condition with constants  $(\beta_0, \kappa_0)$ , then it holds  $\beta_0 \leq C\kappa_0^{-1/2} D^{-3/4}$  for some numerical constant  $C > 0$ .

**Corollary 5.** *When the design is uniform on  $\mathcal{X} = [-\pi, \pi]$ , the dictionary is made of the first  $D$  elements of the Fourier basis, the bound given in the right-hand side of Inequality (2.2) in Theorem 2 above (i.e., Theorem A of [19]) is bounded from below as follows,*

$$\left(\frac{16}{\beta_0 \kappa_0^2}\right)^2 \frac{\sigma^2 Dx}{n} \geq 256 \frac{\sigma^2 D^{5/2} x}{n}.$$

Corollary 5 shows that the rate of convergence provided by the small ball method (Theorem A of [19]) is at most  $D^{5/2}/n$  in the case of the Fourier dictionary. Therefore, we will show in Section 3 that, under Alternative 2 of Theorem 2 above, the excess risk of the least-squares estimator concentrates around  $\mathbb{E}[\zeta^2]D/n$ , just as for localized bases such as wavelets, histograms and piecewise polynomials. Hence, the small-ball method as developed in [19] gives suboptimal results for the linear aggregation of the Fourier dictionary.

The proof of Proposition 4 is based on the use of a “picky” trigonometric polynomial and can be found in Section 4. In Section 2.2.2 below, we will derive a quite general lower bound of the order  $D^{-1}$  for  $\beta_0$ . This bound is in particular valid for the Fourier dictionary, but does not match with the upper bound decaying like  $D^{-3/4}$  of Proposition 4. Therefore, an interesting open question is to determine what is the exact rate of  $\beta_0$  with respect to  $D$  in the Fourier case (at a fixed value of  $\kappa_0$ )? This question remains open.

Finally, it is important to note that Proposition 4 above is a new result, that may be of some informal interest in the related context of Fourier measurement matrices for compressed sensing, where a small-ball condition (or a slightly modified version of it) would yield optimal recovery rates, as noted by Lecué and Mendelson in [17], Remark 1.5:

One may wonder if the small-ball condition is satisfied for more structured matrices, as the argument we use here does not extend immediately to such cases. And, indeed, for structured ensembles one may encounter a different situation: a small-ball condition that is not uniform, in the sense that the constants [...] are direction-dependent.

Concerning instances of “more structured matrices”, Lecué and Mendelson add that “one notable example is a random Fourier measurement matrix”, which is designed by randomly selecting rows of a complete discrete Fourier measurement matrix.

In our setting, also dealing with the Fourier basis but in the “continuous” setting rather than discrete, we show that indeed, *the small-ball condition cannot be satisfied for constants  $(\kappa_0, \beta_0)$  that are absolute*, in the sense that they would be independent of the dimension. But, we also prove in Section 2.2.2 below that *the small-ball condition is achieved, for some constants that indeed depend on the dimension*.

To recover better estimates, it seems reasonable then to look at a more refined property and searching for “direction-dependent” estimates as proposed in [17] seems a good option. Indeed, it is clear that in the directions of functions in the dictionary for instance, that is for trigonometric

functions, the constants are absolute. We follow this lead in Section 3.2 below, where we indeed prove optimal rates of convergence for aggregation on the Fourier dictionary.

2.2.2. Lower bounds for the small-ball coefficients

The following assumption, that states the equivalence between the  $L_\infty$  and  $L_2$  norms for functions in the linear model  $m$ , is satisfied by many classical functional bases:

(A1) Take  $S = \{s_1, \dots, s_D\} \subset L_2(P^X)$  a dictionary and consider its linear span  $m = \text{Span}(S)$ . Assume that there exists a positive constant  $L_0$  such that, for every  $s \in m$ ,

$$\|s\|_\infty \leq L_0 \sqrt{D} \|s\|_2. \tag{2.7}$$

**Remark 6.** As soon as we are given a finite dimensional vector space of functions  $m$ , then it holds

$$R_m := \sup_{s \in m, s \neq \{0\}} \frac{\|s\|_\infty}{\|s\|_2} < +\infty,$$

since the sup-norm and the quadratic norm are equivalent on the finite dimensional space  $m$ . In other words, Assumption (A1) is satisfied as soon as  $m$  is of finite dimension, with a parameter  $L_0$  that may depend on the dimension  $D$ . Therefore, the strength of Assumption (A1) arises when  $L_0$  can be chosen independent of the dimension.

Examples of linear models  $m$  satisfying Assumption (A1) with an absolute constant  $L_0$  are given for instance in [4] and include many classical nonparametric models for functional estimation, such as histograms and piecewise polynomials on a regular partition, compactly supported wavelets and the Fourier basis.

It appears that when a model  $m$  satisfies Assumption (A1), the small-ball condition is verified, but with constants that may depend on the dimension of the model.

**Proposition 7.** *If a linear model  $m$  is of finite linear dimension, then it achieves the small ball condition (with parameters  $(\kappa_0, \beta_0)$  that may depend on the dimension). More precisely, for any  $\kappa_0 \in (0, 1)$ ,  $m$  achieves in that case the small ball condition with parameter  $\beta_0$  achieving the following constraint,*

$$\beta_0 \geq \frac{1 - \kappa_0^2}{R_m^2} > 0, \tag{2.8}$$

where  $R_m = \sup_{s \in m, s \neq 0} \|s\|_\infty / \|s\|_2$  is defined in Remark 6 above. Consequently, if  $m$  satisfies Assumption (A1) then inequality (2.1) of the small-ball condition given in Definition 1 is verified for any  $\kappa_0 \in (0, 1)$  with  $\beta_0 = (1 - \kappa_0^2)L_0^{-2}D^{-1}$ .

The proof of Proposition 7, detailed in Section 4, is a direct application of Paley–Zygmund’s inequality (see [10]). [19] also noticed that more generally, Paley–Zygmund’s inequality could be used to prove the small-ball property when for some  $p > 2$ , the  $L_p$  and  $L_2$  norms are equivalent, or also for *subgaussian classes*, where the Orlicz  $\psi_2$  norm is controlled by the  $L_2$  norm, see [16].



These conditions are weaker than the control of the  $L_\infty$  norm by the  $L_2$  norm, however, as proved in the comments of Proposition 4 – see Section 2.2.1 –, the dependence in  $D$  for  $\beta_0$  given in Proposition 7 above is sharp for localized bases such as histograms, piecewise polynomials and wavelets. Hence, the control of the  $L_\infty$  norm by the  $L_2$  norm is in some way optimal in these cases, and weaker assumptions could not imply some improvements on the behavior of the small ball property for these models.

As for the Fourier basis, the conjunction of Propositions 4 and 7 gives that for such a basis, for any  $\kappa_0 \in (0, 1)$ ,

$$\frac{1 - \kappa_0^2}{2D} \leq \beta_0 \leq \frac{3^{1/4} \sqrt{2}}{\sqrt{\kappa_0} D^{3/4}},$$

since, for the lower bound, Assumption (A1) is satisfied with  $L_0 = \sqrt{2}$  (see, for instance, [4]). As detailed in Section 2.2.1 above, it is an open question to find the right dependence in the dimension for  $\beta_0$ . Moreover, some related questions have a potential impact on compressed sensing theory as developed in [17].

### 3. Optimal excess risks bounds for Fourier expansions

We have shown in Section 2 that the small-ball condition is satisfied for linear models such as histograms, piecewise polynomials, compactly supported wavelets or the Fourier basis, but with constants that depend on the dimension of the model in such a way that using this condition to analyze the rate of convergence of ERM on these models may lead to suboptimal bounds.

Our aim in this section is to show that optimal rates of linear aggregation can indeed be attained by ERM in the Fourier case, that is when the model  $m$  is spanned by the  $D$  first elements of the Fourier basis. We consider two different settings.

In the bounded setting, exposed in Section 3.1, we prove sharp upper and lower bounds for the excess risk that more precisely ensure its concentration around a single deterministic point.

In the general setting treated in Section 3.2, we refine the small-ball arguments developed in [19] by focusing on certain directions where the small-ball is uniform and we also obtain optimal rates of linear aggregation when the noise is only assumed to have a second moment.

#### 3.1. Excess risk’s concentration

We focus in this section on the bounded setting. Let us precisely detail our assumptions. Assume that the design  $X$  is uniformly distributed on  $\mathcal{X} = [0, 2\pi]$  and that the regression function  $s_*$  satisfies  $s_*(0) = s_*(2\pi)$ . Then the Fourier basis is orthonormal in  $L_2(P^X)$  and we consider a model  $m$  of dimension  $D$  (assumed to be odd) corresponding to the linear vector space spanned by the first  $D$  elements of the Fourier basis. More precisely, if we set  $\varphi_1 \equiv 1$ ,  $\varphi_{2k}(x) = \sqrt{2} \cos(kx)$  and  $\varphi_{2k+1}(x) = \sqrt{2} \sin(kx)$  for  $k \geq 1$ , then  $(\varphi_j)_{j=1}^D$  is an orthonormal basis of  $(m, \|\cdot\|_2)$ , for an integer  $l$  satisfying  $2l + 1 = D$ . Assume also:

- (H1) The data and the linear projection of the target onto  $m$  are bounded by a positive finite constant  $A$ :

$$|Y| \leq A \quad \text{a.s.} \tag{3.1}$$

and

$$\|s_m\|_\infty \leq A. \tag{3.2}$$

Hence, from (H1) we deduce that

$$\|s_*\|_\infty = \|\mathbb{E}[Y|X = \cdot]\|_\infty \leq A \tag{3.3}$$

and that there exists a constant  $\sigma_{\max} > 0$  such that

$$\sigma^2(X_i) \leq \sigma_{\max}^2 \leq A^2 \quad \text{a.s.} \tag{3.4}$$

- (H2) The heteroscedastic noise level  $\sigma$  is not reduced to zero:

$$\|\sigma\|_2 = \sqrt{\mathbb{E}[\sigma^2(X)]} > 0.$$

We are now in position to state our result.

**Theorem 8.** *Let  $A_+, A_-, \alpha > 0$  and let  $m$  be a linear vector space spanned by a dictionary made of the first  $D$  elements of the Fourier basis. Assume (H1)–(H2) and take  $\varphi = (\varphi_k)_{k=1}^D$  the Fourier basis of  $m$ . If it holds*

$$A_-(\ln n)^2 \leq D \leq A_+ \frac{n^{1/2}}{\ln n}, \tag{3.5}$$

*then there exists a constant  $A_0 > 0$ , only depending on  $\alpha, A_-, A_+$  and on the constants  $A, \|\sigma\|_2$  defined in assumptions (H1)–(H2), such that by setting*

$$\varepsilon_n = A_0 \max \left\{ \sqrt{\frac{\ln n}{D}}, \frac{D}{\sqrt{n}} \right\}, \tag{3.6}$$

*we have for all  $n \geq n_0(\alpha)$ ,*

$$\mathbb{P} \left[ (1 - \varepsilon_n) \frac{D}{n} C_m^2 \leq \|\hat{s}_m - s_m\|_2^2 \leq (1 + \varepsilon_n) \frac{D}{n} C_m^2 \right] \geq 1 - 3n^{-\alpha}, \tag{3.7}$$

*where  $\hat{s}_m$  is the least-squares estimator on  $m$ , defined in (1.2), and*

$$C_m^2 = \mathbb{E}[\sigma^2(X)] + \|s_* - s_m\|_2^2. \tag{3.8}$$

The rate of convergence of ERM for linear aggregation with a Fourier dictionary exhibited by Theorem 8 is thus of the order  $D/n$ , which is the optimal rate of linear aggregation. In particular, this outperforms the bounds obtained in Theorem 2.2 of [2] under same assumption as Assumption (A1), that is satisfied in the Fourier case, but also under more general moment assumptions

on the noise. Indeed, as noticed in [19], the bounds obtained by [2] are in this case of the order  $D^3/n$ , for models of dimension lower than  $n^{1/4}$ . In Theorem 8, our condition on the permitted dimension is less restrictive, since models with dimension close to  $n^{1/2}$  are allowed.

Concerning the assumptions, uniform boundedness of the projection of the target onto the model, as described in (3.2), is not so restrictive and is guaranteed as soon as the regression function belongs to a broad class of functions named the Wiener algebra, that is whenever the Fourier coefficients of the regression function are summable (in other words when the Fourier series of the regression function is absolutely convergent). For instance, functions that are Hölder continuous with index greater than  $1/2$  belong to the Wiener algebra [12]. For more on the Wiener algebra, see Section 3.2 below.

Furthermore, Theorem 8 gives an information that is far more precise than the rate of convergence of the least-squares estimator. Indeed, Inequality (3.7) of Theorem 8 actually proves the *concentration of the excess risk* of the least-squares estimator around one precise value, which is  $DC_m^2/n$ .

There are only very few and recent such concentration results for the excess risk of a M-estimator in the literature and this question constitutes an exiting new line of research in learning theory. Considering the same regression framework as ours, [29] has shown concentration bounds for the excess risk of the least-squares estimator on models of piecewise polynomial functions. Furthermore, these results have been recently extended in [25] to strongly localized bases, a class of dictionaries containing in particular compactly supported wavelets.

In a slightly different context of least-squares estimation under convex constraint, [8] also proved the concentration in  $L_2$  norm, with fixed design and Gaussian noise. Under the latter assumptions, [24] have shown the excess risk's concentration for the penalized least-squares estimator. Finally, [32] recently proved some concentration results for some regularized M-estimators. They also give an application of their results to a linearized regression context with random design and independent Gaussian noise.

The proof of Theorem 8 is developed in Section 3. We make a recurrent use along our proofs of classical Talagrand's type concentration inequalities for suprema of the empirical process with bounded arguments. We also make use of other tools from empirical process theory, such as a control of variance of the empirical process with bounded arguments – see the proof of Theorem 17 in Section 5.1.1.

### 3.2. A refined small-ball argument

As proved in Section 2 above, a direct application of results of [19] can not lead to the optimal rate of convergence for linear aggregation *via* empirical risk minimization on the Fourier dictionary.

To recover better estimates, it seems reasonable then to look at a more refined property and searching for “direction-dependent” estimates as proposed in [17] – see the quotation in Section 2.2.1 above – seems a good option. Indeed, it is clear that in the directions of functions in the dictionary for instance, that is for trigonometric functions, the constants are absolute. We follow here this lead and this enables us to prove optimal rates of convergence for linear aggregation on the Fourier dictionary.

As explained in the comments following Theorem 8 above, the assumptions needed for Theorem 8 and especially Assumption (3.2) of uniform boundedness of the projection of the regression

function, are ensured if the target belongs to the Wiener algebra, that is if Fourier coefficients are summable. In this case of course, the projection of the target on a Fourier dictionary (with any cardinality) is again in the Wiener algebra. We now denote  $A(\mathbb{T})$  the Wiener algebra. It holds, by definition,

$$A(\mathbb{T}) = \left\{ f = \sum_{k \geq 1} \beta_k \varphi_k; \sum_{k \geq 1} |\beta_k| < +\infty \right\}.$$

We look here at some subsets of the Wiener algebra.

**Definition 9.** Let us take  $\nu > 0$  and denote, for a function  $f$   $2\pi$ -periodic,  $\beta_k(f) = \langle f, \varphi_k \rangle$ . We define the set

$$\Lambda_\nu(L_1, L_2) = \left\{ f \in L_\infty(\mathbb{T}); \sum_{k \geq 1} k^\nu |\beta_k(f)| \leq L_1 \ \& \ \|f\|_\infty \geq L_2 \right\}.$$

In the perspective of the small-ball approach, the interest of the set  $\Lambda_\nu(L_1, L_2)$  lies in the following proposition, ensuring that the small-ball condition (1) is fulfilled uniformly on  $\Lambda(L_1, L_2)$  whenever  $\nu > 1/2$  and  $L_2 > 0$ , for some constants  $(\kappa_0, \beta_0)$  that only depend on  $\nu, L_1$  and  $L_2$ .

**Proposition 10.** Fix  $\nu > 1/2$  and  $L_1, L_2 > 0$ . Take some function  $f \in \Lambda_\nu(L_1, L_2)$ . Then for any  $\kappa_0 \in (0, 1)$ , it holds

$$\mathbb{P}(|f(X)| \geq \kappa_0 \|f\|_2) \geq \frac{(1 - \kappa_0^2) L_2^2}{4C_\nu^2 L_1^2} > 0,$$

with  $C_\nu = \sum_{k \geq 1} k^{-2\nu} < +\infty$ . In other words, the small-ball condition (1) is satisfied uniformly over  $\Lambda_\nu(L_1, L_2)$  with constants  $(\kappa_0, \beta_0)$ , for  $\kappa_0 \in (0, 1)$  and  $\beta_0 = C_\nu^{-2} L_2^2 L_1^{-2} (1 - \kappa_0^2)/4$ .

It is clear from Definition 9 that for any  $\nu > 0$ ,  $\Lambda_\nu(L_1, L_2) \subset A(\mathbb{T})$ . Furthermore, any function of sup-norm greater than the constant  $L_2$  and belonging to a (periodic) Sobolev space  $W_\gamma$  of parameter  $\gamma$  belongs to  $\Lambda_\nu(L_1, L_2)$ , for some constant  $L_1$  and  $\nu < \gamma - 1/2$ .

Recall that periodic Sobolev spaces  $W_\gamma := \bigcup_{L>0} W(\gamma, L)$  are defined as follows (see for instance [30], Section 1.10), for any  $\gamma \in \mathbb{N}_*$ ,

$$W(\gamma, L) := \left\{ f \in L_2(\mathbb{T}); f^{(\gamma-1)} \text{ is absolutely continuous,} \right. \\ \left. \frac{1}{2\pi} \int_0^{2\pi} (f^{(\gamma)}(x))^2 dx \leq L \ \& \ f^{(j)}(0) = f^{(j)}(1), j = 0, 1, \dots, \gamma - 1 \right\}.$$

In addition, the regularity of periodic functions in Sobolev spaces can be directly read on the order of magnitude of their Fourier coefficients. More precisely, for any  $\gamma \in \mathbb{N}_*$ ,  $W_\gamma =$

$\bigcup_{Q>0} \tilde{W}(\gamma, Q)$ , where

$$\tilde{W}(\gamma, Q) := \left\{ f \in L_2(\mathbb{T}); f = \sum_{k \geq 1} \beta_k \varphi_k \ \& \ \sum_{k \geq 1} k^{2\gamma} \beta_k^2 \leq Q \right\}.$$

This second characterization of Sobolev spaces  $W_\gamma$  allow to extend their definition to any  $\gamma > 0$  and not only to integer valued  $\gamma$ . Thus, this is the definition we use in the following proposition.

**Proposition 11.** *With the previous notations, it holds for any  $\nu > 0$ ,*

$$\left[ \bigcup_{\{\gamma: 1/2 + \nu < \gamma\}} \tilde{W}(\gamma, Q) \cap \{f \in L_\infty(\mathbb{T}); \|f\|_\infty \geq L_2\} \right] \subset \Lambda_\nu(L_1, L_2),$$

whenever  $Q \leq L_1^2 (\sum_{k \geq 1} k^{2(\nu-\gamma)})^{-1} < +\infty$ .

Proposition 11 is appealing since the Fourier dictionary is known to achieve minimax rates of convergences for the estimation of a regression function, whenever it lies in a Sobolev space  $W_\gamma$  of parameter  $\gamma > 1$  ([30]). Indeed, by Proposition 10, we are interested by the sets  $\Lambda_\nu(L_1, L_2)$  for  $\nu > 1/2$  and Proposition 11 implies that such sets contain function of Sobolev regularity  $\gamma > \nu + 1/2 > 1$ . This latter fact thus legitimate the focus on the sets  $\Lambda_\nu(L_1, L_2)$ ,  $\nu > 1/2$ , to deal with the performance of linear aggregation from the Fourier dictionary.

Let us turn now to the main result of this section.

**Theorem 12.** *Fix  $\nu > 1/2$ ,  $L_1, L_2 > 0$  and assume that  $s_* \in \Lambda_\nu(L_1, L_2)$ . Let  $S = \{\varphi_1, \dots, \varphi_D\}$  be a dictionary made of the  $D$  first elements of the Fourier basis. Set  $\zeta = Y - s_m(X)$ , where  $s_m$  is the projection of the regression function  $s_*$  onto  $m$ . Assume that  $\zeta$  is independent of  $X$  and  $\mathbb{E}\zeta^2 \leq \sigma^2$ . Then there exists three constants  $L_\nu, L_{L_1, L_2, \sigma, \nu}, C_{\nu, L_1, L_2} > 0$  and an integer  $n_0(\nu, L_1, L_2)$  such that, if*

$$0 < (2\sqrt{2}L_1L_2^{-1})^{1/\nu} \leq D \leq L_\nu(n/\ln n)^{\frac{1}{2(\nu+1)}} \tag{3.9}$$

and  $x \in (0, L_{L_1, L_2, \sigma, \nu} n / D^{2(\nu+1)})$ , the least-squares estimator  $\hat{s}_m$  on  $m$ , defined in (1.2), satisfies for any  $n \geq n_0(\nu, L_1, L_2)$ , on an event of probability at least  $1 - \exp(-\beta_0^2 n/4) - n^{-2} - (2/x)$ ,

$$\|\hat{s}_m - s_m\|_2^2 \leq C_{\nu, L_1, L_2} \frac{\sigma^2 D x}{n}. \tag{3.10}$$

The bound (3.10) obtained in Theorem 12 is optimal in the sense that it achieves the optimal rate  $D/n$  of linear aggregation. Moreover, the only moment needed on the noise term  $\zeta$  is a second moment, which is a minimal assumption. The proof, exposed in Section 5.2, is based on a localization of the least-squares estimator on directions of uniform small-ball property.

Compared to Theorem 8, where we also derived optimal rates of aggregation, but in the bounded setting, we have a stronger assumption on the regularity of the target. Indeed, in Theorem 12  $s_*$  is assumed to belong to some  $\Lambda_\nu(L_1, L_2)$ ,  $\nu > 1/2$ , whereas in Theorem 8, we only

assume that the projection  $s_m$  of the target onto the model  $m$  is uniformly bounded by a constant independent of the dimension, which is achieved as soon as  $s_*$  belongs to the Wiener algebra  $A(\mathbb{T})$ . It appears to be the price to pay to deal with a general noise term, but as explained earlier in this section, the sets  $\Lambda_\nu(L_1, L_2)$ ,  $\nu > 1/2$ , are natural when dealing with the performance of the Fourier dictionary.

Finally, the range of considered dimensions in (3.9) is fairly reasonable, the upper bound being polynomial in  $n$ . In addition, the lower bound, of the order of a constant, is very mild and ensures that the projection of the regression function onto the model does not vanish in sup-norm.

## 4. Proofs related to Section 2

**Proof of Proposition 3.** Take  $f \in m \setminus \{0\}$ . Then, it holds

$$\beta_0 \leq \mathbb{P}(|f(X)| \geq \kappa_0 \|f\|_2) \leq \mathbb{P}(f(X) \neq 0)$$

which readily gives (2.3). Furthermore, as  $\mathbb{P}(|f(X)| \geq \kappa \|f\|_2) = 0$  for  $\kappa > \|f\|_\infty / \|f\|_2$ , it holds  $\kappa_0 \leq \|f\|_\infty / \|f\|_2$ , which implies (2.4) by minimizing the latter bound over all  $f \in m \setminus \{0\}$ . Now,

$$\begin{aligned} \beta_0 &\leq \mathbb{P}(|f(X)| \geq \kappa_0 \|f\|_2) = \mathbb{P}\left(\frac{|f(X)|}{\kappa_0 \|f\|_2} \geq 1\right) \\ &\leq \int_{\mathcal{X}} \left(\frac{|f(x)|}{\kappa_0 \|f\|_2}\right)^q dP^X(x) = \left(\frac{\|f\|_q}{\kappa_0 \|f\|_2}\right)^q, \end{aligned}$$

and by taking the infimum over  $f \in m \setminus \{0\}$ , this thus proves (2.5) and imply  $\beta_0 \kappa_0^2 \leq 1$  for  $q = 2$ . Finally, when  $m$  contains the function identically equal to one, then  $\inf_{f \in m \setminus \{0\}} \|f\|_\infty / \|f\|_2 = 1$ , which implies  $\kappa_0 \leq 1$ .  $\square$

**Proof of Proposition 4.** Recall that  $D = 2l + 1$  is the linear dimension of  $m_D$  and take  $(\beta_0, \kappa_0)$  satisfying the small ball condition on  $m_D$ . Define the  $l$ th Fejér kernel  $F_l$  as follows,

$$F_l(t) = \begin{cases} \frac{\sin^2((l+1)t/2)}{(l+1)\sin^2(t/2)}, & t \in [-\pi, \pi] \setminus \{0\}, \\ l+1, & t = 0. \end{cases}$$

Properties of  $F_l$  are well-known, see, for instance, [3], Section 4.15. In particular,  $F_l \in m_D$ ,  $F_l \geq 0$ ,  $\|F_l\|_\infty \leq l+1 = (D+1)/2$ . Furthermore,  $\int_{-\pi}^{\pi} F_l(t) dt = 2\pi$  which by positivity of  $F_l$  gives  $\|F_l\|_1 = 1$ . We also have, for all  $t \in [-\pi, \pi]$ ,

$$F_l(t) = \sum_{k=-l+1}^{l-1} \left(1 - \frac{|k|}{l}\right) e^{ikt}.$$

Using this formula, one easily computes the quadratic norm of the Fejér kernel, for any  $l \geq 2$ ,

$$\|F_l\|_2^2 = 1 + 2 \sum_{k=1}^{l-1} \left(1 - \frac{k}{l}\right)^2 = 1 + \frac{2}{l^2} \sum_{j=1}^{l-1} j^2 \geq \frac{l}{6}.$$

Now, since for any  $\varepsilon \in (0, \pi]$ ,

$$\sup_{\varepsilon \leq |t| \leq \pi} F_l(t) \leq \frac{1}{l+1} \frac{1}{\sin^2(\varepsilon/2)} \leq \frac{1}{l+1} \left(\frac{\pi}{\varepsilon}\right)^2,$$

it holds

$$\mathbb{P}(|F_l(X)| \geq \kappa_0 \|F_l\|_2) \leq (\kappa_0 \|F_l\|_2 (l+1))^{-1/2}.$$

Consequently,  $\beta_0 \leq (\kappa_0 \|F_l\|_2 (l+1))^{-1/2} \leq C \kappa_0^{-1/2} D^{-3/4}$ , which gives the result. □

**Proof of Corollary 5.** From Proposition 4, it holds

$$\beta_0 \kappa_0^{1/2} \leq C D^{-3/4}.$$

Furthermore, as the model contains the constants, we have  $\kappa_0 \leq 1$  and by combining the two inequalities,  $\beta_0 \kappa_0^2 \leq C D^{-3/4}$ , which gives the result. □

**Proof of Proposition 7.** Take  $s \in m \setminus \{0\}$  and  $\kappa_0 \in (0, 1)$ . Set  $\Omega_{\kappa_0}(s) = \{|s(X)| \geq \kappa_0 \|s\|_2\}$ . By Paley–Zygmund’s inequality (Corollary 3.3.2 in [10]), it holds

$$\mathbb{P}(\Omega_{\kappa_0}(s)) \geq (1 - \kappa_0^2) \frac{\|s\|_2^2}{\|s\|_\infty^2} \geq \frac{1 - \kappa_0^2}{R_m^2},$$

which gives (2.8). The rest of Proposition 7 follows from the latter bound via a simple application of assumption (A1) to bound from above the term  $R_m$ . □

## 5. Proofs related to Section 3

### 5.1. Proof of Theorem 8

Aiming at clarifying the proofs, we generalize a little bit the Fourier framework by invoking the following assumption, that is satisfied for Fourier expansions. From now on,  $m \subset L_2(P^X)$  is considered to be a linear model of dimension  $D$ , not necessarily built from the Fourier basis.

- (H3) Uniformly bounded basis: there exists an orthonormal basis  $\varphi = (\varphi_k)_{k=1}^D$  in  $(m, \|\cdot\|_2)$  that satisfies, for a positive constant  $u_m$ ,

$$\|\varphi_k\|_\infty \leq u_m.$$

Notice that in the Fourier case, (H3) is valid by taking  $u_m \leq \sqrt{2}$ .

**Remark 13.** By Cauchy–Schwarz inequality, we also see that when (H3) is valid, it holds

$$\sup_{s \in m, \|s\|_2 \leq 1} \|s\|_\infty \leq u_m \sqrt{D}. \tag{5.1}$$

Let us denote  $\psi_m(x, y) = y - s_m(x)$ . Then, if  $(\varphi_k)_{k=1}^D$  is formed by the first  $D$  elements of the Fourier basis, the quantity  $\mathcal{C}_m$  defined in (3.8) satisfies

$$\mathcal{C}_m^2 = \frac{1}{D} \sum_{k=1}^D \text{Var}(\psi_m \cdot \varphi_k). \tag{5.2}$$

We will thus prove a slightly more general version than Theorem 8, assuming that (H3) holds and proving Inequality (3.7) with the term  $\mathcal{C}_m$  given by (5.2).

We are now in position to prove Theorem 8.

**Proof of Theorem 8.** Take  $s = \sum_{k=1}^D \beta_k \varphi_k \in m$ . The empirical risk on  $s$  writes

$$\begin{aligned} P_n(\gamma(s)) &= P_n \left[ \left( y - \left( \sum_{k=1}^D \beta_k \varphi_k(x) \right) \right)^2 \right] \\ &= P_n y^2 - 2 \sum_{k=1}^D \beta_k P_n(y \varphi_k(x)) + \sum_{k,l=1}^D \beta_k \beta_l P_n(\varphi_k \varphi_l). \end{aligned}$$

By taking the derivative with respect to  $\beta_l$  in the last quantity, we get

$$\begin{aligned} \frac{1}{2} \frac{\partial}{\partial \beta_l} P_n \left[ \left( y - \left( \sum_{k=1}^D \beta_k \varphi_k(x) \right) \right)^2 \right] \\ = -P_n(y \varphi_l(x)) + \sum_{k=1}^D \beta_k P_n(\varphi_k \varphi_l). \end{aligned} \tag{5.3}$$

Hence, we see that if  $\hat{\beta}_m = (\hat{\beta}_k)_{k=1}^D \in \mathbb{R}^D$  is a critical point of the empirical risk (seen as a function on  $\mathbb{R}^D$ ), then it satisfies the following random linear system,

$$(I_D + A_{n,D}) \hat{\beta}_m = E_{y,n}, \tag{5.4}$$

where  $E_{y,n} = (P_n(y \varphi_k(x)))_{k=1}^D \in \mathbb{R}^D$ ,  $I_D$  is the identity matrix of dimension  $D$  and  $A_{n,D} = ((P_n - P)(\varphi_k \varphi_l))_{k,l=1,\dots,D}$  is a  $D \times D$  matrix. Now, by Inequality (5.11) in Lemma 14 below, a positive integer  $n_0(u_m, \alpha)$  can be found such that for all  $n \geq n_0$ , we have on an event  $\Omega_n$  of probability at least  $1 - n^{-\alpha}$ ,

$$\|A_{n,D}\| \leq L_{A-,u_m,\alpha} \frac{D}{\sqrt{n}} \leq \frac{1}{2}, \tag{5.5}$$



where for a  $D \times D$  matrix  $A$ , the operator norm  $\| \cdot \|$  associated to the quadratic norm  $| \cdot |_2$  on vectors is

$$\|A\| = \sup_{x \neq 0} \frac{|Ax|_2}{|x|_2}.$$

We restrict now on analysis on the event  $\Omega_n$ . Then we deduce from (5.5) that  $(I_D + A_{n,D})$  is a non-singular  $D \times D$  matrix and, as a consequence, that the linear system (5.4) admits a unique solution  $\hat{\beta}_m$  for any  $n \geq n_0(u_m, \alpha)$ . Moreover, since  $P_n(y - (\sum_{k=1}^D \beta_k \varphi_k(x)))^2$  is a nonnegative quadratic functional with respect to  $(\beta_k)_{k=1}^D \in \mathbb{R}^D$  we deduce that for any  $n \geq n_0(u_m, \alpha)$ ,  $\hat{\beta}_m$  achieves on  $\Omega_n$  the unique minimum of  $P_n(y - (\sum_{k=1}^D \beta_k \varphi_k(x)))^2$  on  $\mathbb{R}^D$ , thus  $\hat{s}_m = \sum_{k=1}^D \hat{\beta}_k \varphi_k$ .

Now, if we denote  $\beta_m = (\beta_{*,k})_{k=1}^D$  the vector such that  $s_m = \sum_{k=1}^D \beta_{*,k} \varphi_k$ , then from (5.4) we obtain

$$(I_D + A_{n,D})(\hat{\beta}_m - \beta_m) = F_{y,n},$$

where  $F_{y,n} := E_{y,n} - (I_D + A_{n,D})\beta_m \in \mathbb{R}^D$ . Furthermore, straightforward computations give,

$$F_{y,n} = ((P_n - P)(\psi_{1,m} \varphi_k))_{k=1}^D, \tag{5.6}$$

where  $\psi_{1,m}(x, y) = y - s_m(x)$ ,  $(x, y) \in \mathcal{X} \times \mathbb{R}$ . Finally, for any  $n \geq n_0(u_m, \alpha)$  we get that,

$$\hat{\beta}_m - \beta_m = (I_D + A_{n,D})^{-1} F_{y,n} \tag{5.7}$$

and

$$\|\hat{s}_m - s_m\|_2^2 = |\hat{\beta}_m - \beta_m|_2^2 = |(I_D + A_{n,D})^{-1} F_{y,n}|_2^2. \tag{5.8}$$

By setting  $B_{n,D} = (I_D + A_{n,D})^{-1} - I_D$ , it thus holds,

$$\begin{aligned} \|\hat{s}_m - s_m\|_2^2 - |F_{y,n}|_2^2 &= |(I_D + B_{n,D})F_{y,n}|_2^2 - |F_{y,n}|_2^2 \\ &= |B_{n,D}F_{y,n}|_2^2 + 2\langle F_{y,n}, B_{n,D}F_{y,n} \rangle \\ &\leq (\|B_{n,D}\|^2 + 2\|B_{n,D}\|)|F_{y,n}|_2^2 \end{aligned} \tag{5.9}$$

and for any  $n \geq n_0(u_m, \alpha)$ ,

$$\|B_{n,D}\| \leq \frac{\|A_{n,D}\|}{1 - \|A_{n,D}\|} \leq 2\|A_{n,D}\| \leq L_{A-,u_m,\alpha} \frac{D}{\sqrt{n}}. \tag{5.10}$$

Combining (5.8) and (5.10) implies that, for any  $n \geq n_0(u_m, \alpha)$ ,

$$\|\hat{s}_m - s_m\|_2^2 - |F_{y,n}|_2^2 \leq L_{A-,u_m,\alpha} \frac{D}{\sqrt{n}} |F_{y,n}|_2^2,$$

and the proof simply follows by using Lemma 15 together with the latter inequality. □

**Lemma 14.** Recall that  $A_{n,D} = ((P_n - P)(\varphi_k \varphi_l))_{k,l=1,\dots,D}$  is a  $D \times D$  matrix and that for a  $D \times D$  matrix  $A$ , the operator norm  $\|\cdot\|$  associated to the quadratic norm on the vectors is

$$\|A\| = \sup_{x \neq 0} \frac{|Ax|_2}{|x|_2}.$$

Then, under Assumption (H3), the following inequalities hold on an event of probability at least  $1 - n^{-\alpha}$ ,

$$\|A_{n,D}\| \leq L_{u_m,\alpha} \frac{D}{\sqrt{n}} \left(1 + \sqrt{\frac{\ln n}{D}}\right) \leq \frac{1}{2}. \tag{5.11}$$

**Proof.** Let us denote  $B_2$  the unit ball of  $(m, \|\cdot\|_2)$ . It holds,

$$\begin{aligned} \|A_{n,D}\|^2 &= \sup_{|x|_2=1} |A_{n,D}x|_2^2 \\ &= \sup_{|x|_2=1} \sum_{k=1}^D \left( \sum_{l=1}^D x_l (P_n - P)(\varphi_k \varphi_l) \right)^2 \\ &= \sup_{s \in B_2} \sum_{k=1}^D ((P_n - P)(\varphi_k s))^2 \\ &= \sup_{s,t \in B_2} ((P_n - P)(s \cdot t))^2. \end{aligned}$$

Hence,

$$\|A_{n,D}\| = \sup_{s,t \in B_2} (P_n - P)(st). \tag{5.12}$$

We will now apply Bousquet’s concentration inequality (5.19) to control the deviations of the supremum of the empirical process (5.12). We have,

$$\begin{aligned} \mathbb{E}[\|A_{n,D}\|] &\leq \mathbb{E}^{1/2}[\|A_{n,D}\|^2] \leq \mathbb{E}^{1/2} \left[ \sum_{k,l=1}^D (P_n - P)^2(\varphi_k \varphi_l) \right] \\ &\leq \sqrt{\frac{\sum_{k,l=1}^D \mathbb{E}[\varphi_k^2 \varphi_l^2]}{n}} \leq \frac{u_m D}{\sqrt{n}}, \end{aligned}$$

where we used Assumption (H3) in the last inequality. Furthermore, using (H3) and Remark 13,

$$\sup_{s,t \in B_2} \mathbb{V}(st) \leq \sup_{s \in B_2} \|s\|_\infty^2 \leq u_m^2 D \quad \text{and} \quad \sup_{s,t \in B_2} \|st\|_\infty \leq \sup_{s \in B_2} \|s\|_\infty^2 \leq u_m^2 D.$$

Hence, Bousquet’s concentration inequality (5.19) gives (by taking  $\mathcal{F} = \{st; s, t \in B_2\}$  and  $\varepsilon = 1$ ), for any  $x \geq 0$ ,

$$\mathbb{P}\left[\|A_{n,D}\| \geq \frac{u_m D}{\sqrt{n}} + u_m \sqrt{\frac{2Dx}{n}} + \frac{u_m^2 Dx}{3n}\right] \leq \exp(-x).$$

Now, we get (5.11) by taking  $x = \alpha \ln n$  in the latter inequality. □

**Lemma 15.** *Let us denote  $\psi_m(x, y) = y - s_m(x)$ . Assume that (H1)–(H3) and recall that  $F_{y,n} = ((P_n - P)(\psi_m \varphi_k))_{k=1}^D \in \mathbb{R}^D$ . Then*

$$\mathbb{P}\left(\left(1 - L_{A,A+,A-,u_m,\|\sigma\|_2,\alpha} \sqrt{\frac{\ln n}{D}}\right) \frac{D}{n} C_m^2 \leq \|F_{y,n}\|_2^2\right) \geq 1 - n^{-\alpha} \tag{5.13}$$

and

$$\mathbb{P}\left(\|F_{y,n}\|_2^2 \leq \left(1 + L_{A,A+,A-,u_m,\|\sigma\|_2,\alpha} \sqrt{\frac{\ln n}{D}}\right) \frac{D}{n} C_m^2\right) \geq 1 - n^{-\alpha}, \tag{5.14}$$

where

$$C_m^2 = \frac{1}{D} \sum_{k=1}^D \text{Var}(\psi_m \cdot \varphi_k).$$

**Proof.** It holds

$$\|F_{y,n}\|_2 = \sqrt{\sum_{k=1}^D ((P_n - P)(\psi_m \varphi_k))^2} = \sup_{s \in B_2} (P_n - P)(\psi_m s).$$

We are thus reduced to the study of the supremum of an empirical process. We have, by the hypotheses (H1), (H3) and Remark 13,

$$\sigma^2 := \sup_{s \in B_2} \text{Var}(\psi_m s) \leq \|\psi_m\|_\infty^2 \leq 4A^2 \quad \text{and} \quad b := \sup_{s \in B_2} \|\psi_m s\|_\infty \leq 2A u_m \sqrt{D}. \tag{5.15}$$

Furthermore, it holds

$$\mathbb{E}[\|F_{y,n}\|_2^2] = \frac{D}{n} C_m^2,$$

which gives that for  $\varkappa_n = 2AC_m^{-1}D^{-1/2} \max\{1; \sqrt{A+}u_m\}$ , the two following inequalities are satisfied,

$$\varkappa_n^2 \mathbb{E}[\|F_{y,n}\|_2^2] \geq \frac{\sigma^2}{n}$$

and

$$\varkappa_n^2 \sqrt{\mathbb{E}[\|F_{y,n}\|_2^2]} \geq \frac{b}{n}.$$

Hence, by Theorem 17 applied with  $\mathcal{F} = B_2$ , we have

$$\left(1 - \frac{L_{A, A+, u_m, \|\sigma\|_2}}{\sqrt{D}}\right) C_m \sqrt{\frac{D}{n}} \leq \mathbb{E}[\|F_{y,n}\|_2]. \tag{5.16}$$

We also have

$$\mathbb{E}[\|F_{y,n}\|_2] \leq \sqrt{\mathbb{E}[\|F_{y,n}\|_2^2]} = C_m \sqrt{\frac{D}{n}}. \tag{5.17}$$

Now, by combining the bounds obtained in (5.17) and (5.16) with Inequality (5.21) applied with  $\mathcal{F} = B_2$ ,  $\varepsilon = n^{-1/4} \sqrt{\ln n}$  and  $x = \alpha \ln n$ , we get that on an event of probability at least  $1 - n^{-\alpha}$ ,

$$\begin{aligned} \|F_{y,n}\|_2 &\geq -\sqrt{\frac{2\sigma^2 \alpha \ln n}{n}} + (1 - \varepsilon) \mathbb{E}[\|F_{y,n}\|_2] - \left(\frac{1}{\varepsilon} + 1\right) \frac{b\alpha \ln n}{n} \\ &\geq \left(1 - L_{A, A+, u_m, \|\sigma\|_2, \alpha} \sqrt{\frac{\ln n}{D}}\right) \sqrt{\frac{D}{n}} C_m. \end{aligned}$$

Then easy calculations allow to derive Inequality (5.14) from the latter lower bound.

Finally, combining the bounds obtained in (5.17) and (5.16) with Inequality (5.19) applied with  $\mathcal{F} = B_2$ ,  $\varepsilon = n^{-1/4} \sqrt{\ln n}$  and  $x = \alpha \ln n$ , we also get that on an event of probability at least  $1 - n^{-\alpha}$ ,

$$\begin{aligned} \|F_{y,n}\|_2 &\leq \sqrt{\frac{2\sigma^2 \alpha \ln n}{n}} + (1 + \varepsilon) \mathbb{E}[\|F_{y,n}\|_2] + \left(\frac{1}{\varepsilon} + \frac{1}{3}\right) \frac{b\alpha \ln n}{n} \\ &\leq \left(1 + L_{A, A+, u_m, \|\sigma\|_2, \alpha} \sqrt{\frac{\ln n}{D}}\right) \sqrt{\frac{D}{n}} C_m, \end{aligned}$$

which readily gives (5.13). □

### 5.1.1. Probabilistic tools

We recall here the main probabilistic results that are instrumental in the proof of Theorem 8 above.

Denote by

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i}$$

the empirical measure associated to the sample  $(\xi_1, \dots, \xi_n)$  and by

$$\|P_n - P\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |(P_n - P)(f)|$$

the supremum of the empirical process over  $\mathcal{F}$ .

We turn now to concentration inequalities for the empirical process around its mean. Bousquet’s inequality [6] provides optimal constants for the deviations at the right. Klein–Rio’s inequality [14] gives sharp constants for the deviations at the left, that slightly improves Klein’s inequality [13].

**Theorem 16.** *Let  $(\xi_1, \dots, \xi_n)$  be  $n$  i.i.d. random variables having common law  $P$  and taking values in a measurable space  $\mathcal{Z}$ . If  $\mathcal{F}$  is a class of measurable functions from  $\mathcal{Z}$  to  $\mathbb{R}$  satisfying*

$$|f(\xi_i) - Pf| \leq b \quad \text{a.s., for all } f \in \mathcal{F}, i \leq n,$$

then, by setting

$$\sigma_{\mathcal{F}}^2 = \sup_{f \in \mathcal{F}} \{P(f^2) - (Pf)^2\},$$

we have, for all  $x \geq 0$ ,

Bousquet’s inequality:

$$\begin{aligned} \mathbb{P} \left[ \|P_n - P\|_{\mathcal{F}} - \mathbb{E}[\|P_n - P\|_{\mathcal{F}}] \geq \sqrt{2(\sigma_{\mathcal{F}}^2 + 2b\mathbb{E}[\|P_n - P\|_{\mathcal{F}}])\frac{x}{n}} + \frac{bx}{3n} \right] \\ \leq \exp(-x) \end{aligned} \tag{5.18}$$

and we can deduce that, for all  $\varepsilon, x > 0$ , it holds

$$\begin{aligned} \mathbb{P} \left[ \|P_n - P\|_{\mathcal{F}} - \mathbb{E}[\|P_n - P\|_{\mathcal{F}}] \geq \sqrt{2\sigma_{\mathcal{F}}^2\frac{x}{n}} + \varepsilon\mathbb{E}[\|P_n - P\|_{\mathcal{F}}] + \left(\frac{1}{\varepsilon} + \frac{1}{3}\right)\frac{bx}{n} \right] \\ \leq \exp(-x). \end{aligned} \tag{5.19}$$

Klein–Rio’s inequality:

$$\begin{aligned} \mathbb{P} \left[ \mathbb{E}[\|P_n - P\|_{\mathcal{F}}] - \|P_n - P\|_{\mathcal{F}} \geq \sqrt{2(\sigma_{\mathcal{F}}^2 + 2b\mathbb{E}[\|P_n - P\|_{\mathcal{F}}])\frac{x}{n}} + \frac{bx}{n} \right] \\ \leq \exp(-x) \end{aligned} \tag{5.20}$$

and again, we can deduce that, for all  $\varepsilon, x > 0$ , it holds

$$\begin{aligned} \mathbb{P} \left[ \mathbb{E}[\|P_n - P\|_{\mathcal{F}}] - \|P_n - P\|_{\mathcal{F}} \geq \sqrt{2\sigma_{\mathcal{F}}^2\frac{x}{n}} + \varepsilon\mathbb{E}[\|P_n - P\|_{\mathcal{F}}] + \left(\frac{1}{\varepsilon} + 1\right)\frac{bx}{n} \right] \\ \leq \exp(-x). \end{aligned} \tag{5.21}$$

The following theorem is proved in [29], Corollary 25. It can be derived from a Theorem by Rio [28], improving on previous results by Ledoux, and controlling the variance of the supremum of an empirical process with bounded arguments (see also Theorem 11.10 in [5]).

**Theorem 17.** Under notations of Theorem 16, if some  $\varkappa_n \in (0, 1)$  exists such that

$$\varkappa_n^2 \mathbb{E}[\|P_n - P\|_{\mathcal{F}}^2] \geq \frac{\sigma^2}{n}$$

and

$$\varkappa_n^2 \sqrt{\mathbb{E}[\|P_n - P\|_{\mathcal{F}}^2]} \geq \frac{b}{n}$$

then we have, for a numerical constant  $A_{1,-}$ ,

$$(1 - \varkappa_n A_{1,-}) \sqrt{\mathbb{E}[\|P_n - P\|_{\mathcal{F}}^2]} \leq \mathbb{E}[\|P_n - P\|_{\mathcal{F}}].$$

### 5.2. Proofs related to Section 3.2

**Proof of Proposition 10.** Take  $f \in \Lambda_\nu(L_1, L_2)$  and  $\kappa_0 \in (0, 1)$ . Then,

$$\begin{aligned} \|f\|_\infty &\leq \sqrt{2} \sum_{k \in \mathbb{N}_*} |\beta_k(f)| \\ &\leq \sqrt{2} \sqrt{L_1} \sqrt{\sum_{k \in \mathbb{N}_*} \frac{|\beta_k(f)|}{k^\nu}}, \end{aligned} \tag{5.22}$$

where the second inequality follows from Cauchy–Schwarz inequality. Furthermore, by Cauchy–Schwarz inequality again,

$$\sum_{k \in \mathbb{N}_*} \frac{|\beta_k(f)|}{k^\nu} \leq \left( \sum_{k \in \mathbb{N}_*} \frac{1}{k^{2\nu}} \right)^{1/2} \left( \sum_{k \in \mathbb{N}_*} \beta_k^2(f) \right)^{1/2} = \sqrt{C_\nu} \|f\|_2, \tag{5.23}$$

with  $C_\nu := \sum_{k \in \mathbb{N}_*} k^{-2\nu} < +\infty$  since  $\nu > 1/2$ . Combining (5.22), (5.23) and the fact that  $\|f\|_\infty \geq L_2 > 0$ , we get

$$\|f\|_\infty \leq \frac{\|f\|_\infty^2}{L_2} \leq 2C_\nu \frac{L_1}{L_2} \|f\|_2.$$

The conclusion then follows from Paley–Zygmund’s inequality (Corollary 3.3.2 in [10]), since it holds

$$\mathbb{P}(|f(X)| \geq \kappa_0 \|f\|_2) \geq (1 - \kappa_0^2) \frac{\|f\|_2^2}{\|f\|_\infty^2} \geq \frac{(1 - \kappa_0^2) L_2^2}{4C_\nu^2 L_1^2} > 0. \quad \square$$

We turn now to the proof of Theorem 12. The idea is to localize the calculations on a subset of the model  $m$ , containing the estimator  $\hat{s}_m$  w.h.p. and achieving the small-ball condition with some absolute constants. Therefore, we first need the following result, which is a direct extension of Theorem A in [19].

**Theorem 18.** Let  $S = \{s_1, \dots, s_D\} \subset L_2(P^X)$  be a dictionary. Assume that a set  $m_0 \subset m := \text{Span}(S)$  satisfies the small-ball condition with constants  $\kappa_0$  and  $\beta_0$  (see Definition 1 above) and contains, on an event  $\Omega_0$ , the least-squares estimator  $\hat{s}_m$  on  $m$ , defined in (1.2). Let  $n \geq (400)^2 D/\beta_0^2$  and set  $\zeta = Y - s_m(X)$ , where  $s_m$  is the projection of the regression function  $s_*$  onto  $m$ . Assume further that one of the following two conditions holds:

1.  $\zeta$  is independent of  $X$  and  $\mathbb{E}\zeta^2 \leq \sigma^2$ , or
2.  $|\zeta| \leq \sigma$  almost surely.

Then the estimator  $\hat{s}_m$  satisfies for every  $x > 0$ , with probability at least  $1 - \mathbb{P}(\Omega_0^c) - \exp(-\beta_0^2 n/4) - (1/x)$ ,

$$\|\hat{s}_m - s_m\|_2^2 \leq \left(\frac{16}{\beta_0 \kappa_0^2}\right)^2 \frac{\sigma^2 D x}{n}.$$

Theorem 18 ensures that if an information is available w.h.p. on the location of the estimator on the model  $m$ , then it may be used to derive better rates by taking advantage of better small-ball constants achieved on the restricted set containing the estimator.

The proof is omitted, since a careful reading of the proof of Theorem A in [19] allows to conclude that using a localization of the estimator  $\hat{s}_m$  does not change the reasoning, neither the validity of the arguments.

The following proposition states that indeed, when the regression function  $s_*$  is sufficiently regular, then so is the least-squares estimator on the first elements of the Fourier basis.

**Proposition 19.** Take  $v, L_1, L_2, z > 0$  and assume that  $s_* \in \Lambda_v(L_1, L_2)$ . For a dimension  $D$  satisfying

$$0 < (2\sqrt{2}L_1L_2^{-1})^{1/v} \leq D \leq L_v(n/\ln n)^{\frac{1}{2(v+1)}}$$

and for  $z \leq L_{L_1, L_2, \sigma, v} n/D^{2(v+1)}$ , it holds

$$\mathbb{P}\left(\hat{s}_m \in \Lambda_v\left(2L_1, \frac{L_2}{4}\right)\right) \geq 1 - n^{-2} - 1/z.$$

We are now in a position to prove Theorem 12. The proof of Proposition 19 is thus postponed after the proof of Theorem 12.

**Proof of Theorem 12.** Apply Theorem 18 with

$$\Omega_0 = \left\{ \hat{s}_m \in \Lambda_v\left(2L_1, \frac{L_2}{4}\right) \right\}$$

and  $x = z$ . Then Proposition 10 ensures that on  $\Omega_0$  the small-ball is achieved with parameters  $\kappa_0 = 2^{-1/2}$  and  $\beta_0 = C_v^{-2} L_2^2 L_1^{-2}/8$ . Hence, the condition  $n \geq (400)^2 D/\beta_0^2$  is satisfied for  $D \leq L_v(n/\ln n)^{\frac{1}{2(v+1)}}$  whenever  $n \geq n_0(v, L_1, L_2)$ . Theorem 12 then follows from Proposition 19 and straightforward computations. □

Before proving Proposition 19, let us denote, for any  $\nu > 0$ ,

$$\Lambda_\nu = \bigcup_{L_1, L_2 > 0} \Lambda_\nu(L_1, L_2) = \left\{ f = \sum_{k \geq 1} \beta_k \varphi_k; \sum_{k \geq 1} k^\nu |\beta_k| < +\infty \right\}.$$

For any  $f \in \Lambda_\nu$ , let us write  $\|f\|_{\Lambda, \nu} = \sum_{k \geq 1} k^\nu |\langle f, \varphi_k \rangle|$ . It is easily seen that  $\|\cdot\|_{\Lambda, \nu}$  is a norm on the space  $\Lambda_\nu$ .

For a sequence  $\beta = (\beta_k)_{k \geq 1} \in \mathbb{R}^{\mathbb{N}}$ , we denote  $|\beta|_{\Lambda, \nu} = \sum_{k \geq 1} k^\nu |\beta_k| \in \mathbb{R}_+ \cup \{+\infty\}$  and  $\tilde{\Lambda}_\nu := \{\beta = (\beta_k)_{k \geq 1} \in \mathbb{R}^{\mathbb{N}}; \sum_{k \geq 1} k^\nu |\beta_k| < +\infty\}$ . Furthermore, for a  $D \times D$  matrix  $A$ , the operator norm  $\|A\|_{\Lambda, \nu}$  associated to the norm  $|\cdot|_{\Lambda, \nu}$  on the vectors (seen as sequences with finite support) is

$$\|A\|_{\Lambda, \nu} := \sup_{x \in \mathbb{R}^D, x \neq 0} \frac{|Ax|_{\Lambda, \nu}}{|x|_{\Lambda, \nu}}.$$

By simple computations it holds, for any matrix  $A = (A_{k,l})_{1 \leq k, l \leq D}$ ,

$$\begin{aligned} \|A\|_{\Lambda, \nu} &= \sup \left\{ \sum_{k=1}^D k^\nu \left| \sum_{l=1}^D A_{k,l} x_l \right|; x \in \mathbb{R}^D \text{ \& } \sum_{k=1}^D k^\nu |x_k| = 1 \right\} \\ &= \sum_{k=1}^D k^\nu \max_{l=1, \dots, D} \left| \frac{A_{k,l}}{l^\nu} \right|. \end{aligned} \tag{5.24}$$

**Proof of Proposition 19.** Let us write  $s_* = \sum_{k \geq 1} \beta_k \varphi_k$ . Thus  $s_m = \sum_{k=1}^D \beta_k \varphi_k$  and since  $s_* \in \Lambda_\nu(L_1, L_2)$ , it holds  $\sum_{k \geq 1} k^\nu |\beta_k| \leq L_1$ . Hence, it holds in particular  $\sum_{k=1}^D k^\nu |\beta_k| \leq L_1$  and

$$\|s_* - s_m\|_\infty \leq \sqrt{2} \sum_{k \geq D+1} |\beta_k| \leq \frac{\sqrt{2}}{D^\nu} \sum_{k \geq D+1} k^\nu |\beta_k| \leq \frac{\sqrt{2} L_1}{D^\nu}.$$

Consequently, we have  $\|s_* - s_m\|_\infty \leq L_2/2$  and so  $\|s_m\|_\infty \geq \|s_*\|_\infty - \|s_* - s_m\|_\infty \geq L_2/2$  whenever  $D \geq (2\sqrt{2}L_1L_2^{-1})^{1/\nu}$ . Therefore, for such dimension  $D$ , we get  $s_m \in \Lambda_\nu(L_1, L_2/2)$ .

Now, we write  $\hat{s}_m = \sum_{k=1}^D \hat{\beta}_k \varphi_k$ ,  $\hat{\beta}_m = (\hat{\beta}_k)_{k=1}^D$  and we define the following set,

$$\begin{aligned} \Omega_\Lambda &= \left\{ \|A_{n,D}\|_{\Lambda, \nu} \leq \frac{4(D+1)^{\nu+1}}{\nu+1} \sqrt{\frac{3 \ln n}{n}} \leq \frac{1}{2} \right\} \\ &\cap \left\{ |F_{y,n}|_{\Lambda, \nu} \leq \frac{(D+1)^{\nu+1}}{\nu+1} \sqrt{\frac{2\sigma^2 z}{n}} \right\}, \end{aligned}$$

where the matrix  $A_{n,D}$  and the vector  $F_{y,n}$  are defined respectively in (5.4) and (5.6), where  $(\varphi_k)_{k=1}^D$  should stand for the first  $D$  elements of the Fourier basis this time. On  $\Omega_\Lambda$ , the matrix  $I_D + A_{n,D}$  is invertible and it holds  $\hat{\beta}_m - \beta_m = (I_D + A_{n,D})^{-1} F_{y,n}$ .



From Lemmas 20 and 21, there exists an integer  $n_0(v)$  such that for any  $n \geq n_0(v)$ ,  $\mathbb{P}(\Omega_\Lambda) \geq 1 - n^{-2} - 1/z$ . Furthermore, on  $\Omega_\Lambda$ , we have,

$$\begin{aligned} \|\hat{s}_m - s_m\|_{\Lambda, v} &= |\hat{\beta}_m - \beta_m|_{\Lambda, v} \leq \|(I_D + A_{n, D})^{-1}\| \|F_{y, n}|_{\Lambda, v} \\ &\leq (1 + 2\|A_{n, D}\|_{\Lambda, v}) |F_{y, n}|_{\Lambda, v} \leq \frac{2(D+1)^{v+1}}{v+1} \sqrt{\frac{\sigma^2 z}{n}} \end{aligned} \tag{5.25}$$

and

$$\|\hat{s}_m - s_m\|_\infty \leq \sqrt{2} \sum_{k=1}^D |\hat{\beta}_k - \beta_k| \leq \sqrt{2} |\hat{\beta}_m - \beta_m|_{\Lambda, v} \leq \frac{2\sqrt{2}(D+1)^{v+1}}{v+1} \sqrt{\frac{\sigma^2 z}{n}}. \tag{5.26}$$

Finally, it is easily seen from (5.25) and (5.26) that there exists a constant  $L_{L_1, L_2, \sigma, v}$  such that if  $z \leq L_{L_1, L_2, \sigma, v} n / D^{2(v+1)}$ , then  $\|\hat{s}_m - s_m\|_{\Lambda, v} \leq L_1$  and  $\|\hat{s}_m - s_m\|_\infty \leq L_2/4$ .  $\square$

**Lemma 20.** *Recall that  $A_{n, D} = ((P_n - P)(\varphi_k \varphi_l))_{k, l=1, \dots, D}$  is a  $D \times D$  matrix. Then the following inequalities hold on an event of probability at least  $1 - D^2 n^{-\alpha}$ ,*

$$\|A_{n, D}\|_{\Lambda, v} \leq \frac{2(D+1)^{v+1}}{v+1} \sqrt{\frac{\alpha \ln n}{n}} \left(1 + \sqrt{\frac{\alpha \ln n}{n}}\right). \tag{5.27}$$

Consequently, there exists a constant  $L_v > 0$  such that for  $D \leq L_v (n / \ln n)^{\frac{1}{2(v+1)}}$ , it holds for any  $n \geq n_0(v)$ , with probability at least  $1 - n^{-2}$ ,

$$\|A_{n, D}\|_{\Lambda, v} \leq \frac{4(D+1)^{v+1}}{v+1} \sqrt{\frac{3 \ln n}{n}} \leq \frac{1}{2}. \tag{5.28}$$

**Proof.** By (5.24) we have,

$$\|A_{n, D}\|_{\Lambda, v} = \sum_{k=1}^D k^v \max_{l=1, \dots, D} \left| \frac{(P_n - P)(\varphi_k \varphi_l)}{l^v} \right|. \tag{5.29}$$

Furthermore, for any  $k, l = 1, \dots, D$ , it holds

$$\mathbb{V}(\varphi_k \varphi_l) \leq \|\varphi_k\|_\infty^2 \mathbb{E}[\varphi_l^2] \leq 2 \quad \text{and} \quad \|\varphi_k \varphi_l\|_\infty \leq 2.$$

Hence, for any  $x > 0$ , we get by Bernstein's inequality (see, for instance, [20]), that on an event  $\Omega_{k, l}(x)$  of probability at least  $1 - 2 \exp(-x)$ ,

$$\left| (P_n - P)(\varphi_k \varphi_l) \right| \leq 2\sqrt{\frac{x}{n}} + \frac{2x}{n}.$$

Then Identity (5.29) implies that, for any  $\alpha > 0$ , on the event  $\Omega_D = \bigcap_{1 \leq l \leq k \leq D} \Omega_{k,l}(\alpha \ln n)$  of probability greater than  $1 - D^2/n^\alpha$ ,

$$\begin{aligned} \|A_{n,D}\|_{\Lambda,v} &\leq 2 \left( \sqrt{\frac{\alpha \ln n}{n}} + \frac{\alpha \ln n}{n} \right) \sum_{k=1}^D k^\nu \\ &\leq \frac{2(D+1)^{\nu+1}}{\nu+1} \sqrt{\frac{\alpha \ln n}{n}} \left( 1 + \sqrt{\frac{\alpha \ln n}{n}} \right). \end{aligned} \tag{5.30}$$

Thus (5.27) is proved and Inequality (5.28) can be deduced from it by simply taking  $\alpha = 3$ .  $\square$

**Lemma 21.** *Let us denote  $\psi_m(x, y) = y - s_m(x)$ . Recall that  $F_{y,n} = ((P_n - P)(\psi_m \varphi_k))_{k=1}^D \in \mathbb{R}^D$ . Then, for any  $z > 0$ ,*

$$\mathbb{P} \left( \|F_{y,n}\|_{\Lambda,v} \leq \frac{(D+1)^{\nu+1}}{\nu+1} \sqrt{\frac{2\sigma^2 z}{n}} \right) \geq 1 - \frac{1}{z}. \tag{5.31}$$

**Proof.** It holds

$$\begin{aligned} \sqrt{\mathbb{E}[\|F_{y,n}\|_{\Lambda,v}^2]} &= \sqrt{\mathbb{E} \left[ \left( \sum_{k=1}^D k^\nu |(P_n - P)(\psi_m \varphi_k)| \right)^2 \right]} \\ &\leq \sum_{k=1}^D k^\nu \sqrt{\mathbb{E}[(P_n - P)^2(\psi_m \varphi_k)]} \\ &\leq \max_{k=1, \dots, D} \sqrt{\frac{\mathbb{E}[\psi_m^2 \varphi_k^2]}{n}} \sum_{k=1}^D k^\nu \\ &\leq \frac{(D+1)^{\nu+1}}{\nu+1} \sqrt{\frac{2\sigma^2}{n}}. \end{aligned}$$

Then Lemma 21 follows from Markov's inequality.  $\square$

## Acknowledgements

I would like to express my gratitude to the Associate Editor and two anonymous referees for their remarks and questions that helped to greatly improve on the content of the first submission. I also warmly thank Guillaume Lecué for fruitful discussions about the small-ball approach, for his great scientific open mind and for his enthusiasm.

## References

- [1] Arlot, S. and Massart, P. (2009). Data-driven calibration of penalties for least-squares regression. *J. Mach. Learn. Res.* **10** 245–279 (electronic).
- [2] Audibert, J.-Y. and Catoni, O. (2011). Robust linear least squares regression. *Ann. Statist.* **39** 2766–2794.
- [3] Bachman, G., Narici, L. and Beckenstein, E. (2000). *Fourier and Wavelet Analysis. Universitext*. New York: Springer. [MR1729490](#)
- [4] Barron, A., Birgé, L. and Massart, P. (1999). Risk bounds for model selection via penalization. *Probab. Theory Related Fields* **113** 301–413. [MR1679028](#)
- [5] Boucheron, S., Lugosi, G. and Massart, P. (2013). *Concentration Inequalities. A Nonasymptotic Theory of Independence*. Oxford: Oxford Univ. Press. [MR3185193](#)
- [6] Bousquet, O. (2002). A Bennett concentration inequality and its application to suprema of empirical processes. *C. R. Math. Acad. Sci. Paris* **334** 495–500. [MR1890640](#)
- [7] Bunea, F., Tsybakov, A.B. and Wegkamp, M.H. (2007). Aggregation for Gaussian regression. *Ann. Statist.* **35** 1674–1697.
- [8] Chatterjee, S. (2014). A new perspective on least squares under convex constraint. *Ann. Statist.* **42** 2340–2381. [MR3269982](#)
- [9] Cohen, A., Daubechies, I. and Vial, P. (1993). Wavelets on the interval and fast wavelet transforms. *Appl. Comput. Harmon. Anal.* **1** 54–81.
- [10] de la Peña, V.H. and Giné, E. (1999). From dependence to independence, randomly stopped processes. *U-statistics and processes*. In *Decoupling. Probability and Its Applications (New York)*. New York: Springer.
- [11] Härdle, W., Kerkycharian, G., Picard, D. and Tsybakov, A. (1998). *Wavelets, Approximation, and Statistical Applications. Lecture Notes in Statistics* **129**. New York: Springer. [MR1618204](#)
- [12] Katznelson, Y. (2004). *An Introduction to Harmonic Analysis*, 3rd ed. *Cambridge Mathematical Library*. Cambridge: Cambridge Univ. Press. [MR2039503](#)
- [13] Klein, T. (2002). Une inégalité de concentration à gauche pour les processus empiriques. *C. R. Math. Acad. Sci. Paris* **334** 501–504. [MR1890641](#)
- [14] Klein, T. and Rio, E. (2005). Concentration around the mean for maxima of empirical processes. *Ann. Probab.* **33** 1060–1077.
- [15] Koltchinskii, V. and Mendelson, S. (2015). Bounding the smallest singular value of a random matrix without concentration. *Int. Math. Res. Not. IMRN* **23** 12991–13008. [MR3431642](#)
- [16] Lecué, G. and Mendelson, S. (2013). Learning subgaussian classes: Upper and minimax bounds. In *Topics in Learning Theory – Societe Mathematique de France* (S. Boucheron and N. Vayatis, eds.). To appear.
- [17] Lecué, G. and Mendelson, S. (2014). Sparse recovery under weak moment assumptions. *J. Eur. Math. Soc. (JEMS)*. Technical report. To appear.
- [18] Lecué, G. and Mendelson, S. (2016). Regularization and the small-ball method i: Sparse recovery. [arXiv:1601.05584](#).
- [19] Lecué, G. and Mendelson, S. (2016). Performance of empirical risk minimization in linear aggregation. *Bernoulli* **22** 1520–1534. [MR3474824](#)
- [20] Massart, P. (2007). *Concentration Inequalities and Model Selection. Lecture Notes in Math.* **1896**. Berlin: Springer. [MR2319879](#)
- [21] Mendelson, S. (2014). Learning without concentration for general loss functions. Technical report. Australia: Technion, Israel and ANU.
- [22] Mendelson, S. (2014). A remark on the diameter of random sections of convex bodies. In *Geometric Aspects of Functional Analysis. Lecture Notes in Math.* **2116** 395–404. Springer, Cham. [MR3364699](#)

- [23] Mendelson, S. (2015). Learning without concentration. *J. ACM* **62** Art. 21, 25. [MR3367000](#)
- [24] Muro, A. and van de Geer, S. (2015). Concentration behavior of the penalized least squares estimator. Preprint. Available at [arXiv:1511.08698](#).
- [25] Navarro, F. and Saumard, A. (2017). Slope heuristics and V-Fold model selection in heteroscedastic regression using strongly localized bases. *ESAIM: Probability and Statistics*, to appear.
- [26] Nemirovski, A. (2000). Topics in non-parametric statistics. In *Lectures on Probability Theory and Statistics. Lecture Notes in Math.* **1738** 85–277. Berlin: Springer.
- [27] Rigollet, Ph. and Tsybakov, A.B. (2007). Linear and convex aggregation of density estimators. *Math. Methods Statist.* **16** 260–280. [MR2356821](#)
- [28] Rio, E. (2001). Inégalités de concentration pour les processus empiriques de classes de parties. *Probab. Theory Related Fields* **119** 163–175. [MR1818244](#)
- [29] Saumard, A. (2012). Optimal upper and lower bounds for the true and empirical excess risks in heteroscedastic least-squares regression. *Electron. J. Stat.* **6** 579–655. [MR2988421](#)
- [30] Tsybakov, A.B. (1996). *Introduction à L'estimation Non-paramétrique*. Berlin: Springer.
- [31] Tsybakov, A.B. (2003). Optimal rates of aggregation. In *Learning Theory and Kernel Machines* 303–313. Springer.
- [32] van de Geer, S. and Wainwright, M. (2016). On concentration for (regularized) empirical risk minimization. Preprint. Available at [arXiv:1512.00677](#).

Received October 2016