

On the optimality of the empirical risk minimization procedure for the Convex Aggregation problem

Guillaume Lecué^{1,3,4} Shahar Mendelson^{2,3,5}

October 14, 2011

Abstract

We study the performance of *empirical risk minimization* (ERM), with respect to the quadratic risk, in the context of *convex aggregation*, in which one wants to construct a procedure whose risk is as close as possible to the best function in the convex hull of an arbitrary finite class F . We show that ERM performed in the convex hull of F is an optimal aggregation procedure for the convex aggregation problem. We also show that if this procedure is used for the problem of model selection aggregation, in which one wants to mimic the performance of the best function in F itself, then its rate is the same as the one achieved for the convex aggregation problem, and thus is far from optimal. These results are obtained in deviation and are sharp up to logarithmic factors.

(Résumé en Français: Nous étudions les performances de la procédure de minimisation du risque empirique, par rapport au risque quadratique, pour le problème d'agrégation convexe. Dans ce problème, on souhaite construire des procédures dont le risque est aussi proche que possible du risque du meilleur élément dans l'enveloppe convexe d'une classe finie F de fonctions. Nous prouvons que la procédure obtenue par minimisation du risque empirique sur la coque convexe de F est une procédure optimale pour le problème d'agrégation convexe. Nous prouvons aussi que si cette procédure est utilisée pour le problème d'agrégation en sélection de modèle, pour lequel on souhaite imiter le meilleur dans F , alors le résiduel d'agrégation est le même que celui obtenu pour le problème d'agrégation convexe. Cette procédure est donc loin d'être optimale pour le problème d'agrégation en sélection de modèle. Ces résultats sont obtenus en déviation et sont optimaux à des facteurs logarithmiques près.)

¹CNRS, LAMA, Université Paris-Est Marne-la-vallée, 77454 France.

²Department of Mathematics, Technion, I.I.T, Haifa 32000, Israel.

³Part of this research was supported by the Centre for Mathematics and its Applications, The Australian National University, Canberra, ACT 0200, Australia, by an Australian Research Council Discovery grant DP0559465 and by the European Community's Seventh Framework Programme (FP7/2007-2013), ERC grant agreement 203134.

⁴Email: guillaume.lecue@univ-mlv.fr

⁵Email: shahar@tx.technion.ac.il

1 Introduction and main results

In this note, we study the optimality of the empirical risk minimization procedure in the aggregation framework.

Let \mathcal{X} be a probability space and let (X, Y) and $(X_1, Y_1), \dots, (X_n, Y_n)$ be $n + 1$ i.i.d. random variables with values in $\mathcal{X} \times \mathbb{R}$. From the statistical point of view, $\mathcal{D} = ((X_1, Y_1), \dots, (X_n, Y_n))$ is the family of given data.

The *quadratic risk* of a real-valued function f defined on \mathcal{X} is given by

$$R(f) = \mathbb{E}(Y - f(X))^2.$$

If \hat{f} is a function constructed using the data \mathcal{D} , the quadratic risk of \hat{f} is the random variable

$$R(\hat{f}) = \mathbb{E} \left[(Y - \hat{f}(X))^2 | \mathcal{D} \right].$$

For the sake of simplicity, throughout this article we will restrict ourselves to functions f and random variables (X, Y) for which $|Y|, |f(X)| \leq b$ almost surely, for some fixed $b \geq 1$. One should note, though, that it is possible to extend the results beyond this case, to functions with well behaved tail – though at a high technical price (cf. the chaining arguments in [21] and [20]).

In the aggregation framework, one is given a finite set F of real-valued functions defined on \mathcal{X} (usually called a *dictionary*) of cardinality M . There are three main types of aggregation problems:

1. In the *Model Selection (MS) aggregation* problem, one has to construct a procedure that produces a function whose risk is as close as possible to the risk of the best element in the given class F (cf. [2, 3, 9, 10, 11, 12, 16, 24, 25, 27]).
2. In the *Convex (C) aggregation* problem (cf. [1, 7, 8, 9, 12, 24, 28]) one wants to construct a procedure whose risk is as close as possible to the risk of the best function in the convex hull of F (later denoted by $\text{conv}(F)$).
3. In the *linear (L) aggregation* problem (cf. [9, 11, 15, 24]), one wants to construct a procedure whose risk is as close as possible to the risk of the best function in the linear span of F (later denoted by $\text{span}(F)$).

The aim in the aggregation framework is to construct a procedure \tilde{f} for which, with high probability

$$R(\tilde{f}) \leq C \min_{f \in \Delta(F)} R(f) + \psi_n^{\Delta(F)}(M) \tag{1.1}$$

with $C = 1$ and $\Delta(F)$ is either F , or $\text{conv}(F)$ or $\text{span}(F)$. It is worth mentioning that it is desirable for the constant C in (1.1) to be one in the aggregation setup for at least

two reasons. First, there are some obvious mathematical differences in the analysis leading to exact oracle inequalities ($C = 1$) and non-exact oracle inequalities ($C > 1$). In particular, the geometry of the set $\Delta(F)$ has a key role in an attempt to obtain exact oracle inequalities, whereas non-exact oracle inequalities are mainly based on complexity and concentration argument (cf. [17]). Second, an exact oracle inequality for the prediction risk $R(\cdot)$ leads to an exact oracle inequality for the estimation risk; namely, with high probability

$$\mathbb{E}[(\tilde{f}(X) - f^*(X))^2 | \mathcal{D}] \leq \min_{f \in \Delta(F)} \mathbb{E}[(f(X) - f^*(X))^2] + \psi_n^{\Delta(F)}(M),$$

where f^* denotes the regression function of Y given X . Such an estimate on the regression function cannot follow from a non-exact oracle inequality, and thus, exact oracle inequalities can provide prediction and estimation results whereas non-exact oracle inequalities only lead to prediction results.

One can define the *optimal rates of the (MS), (C) and (L) aggregation* problems, respectively denoted by $\psi_n^{(MS)}(M)$, $\psi_n^{(C)}(M)$ and $\psi_n^{(L)}(M)$ (see, for example, [24]). The optimal rates are the smallest prices in the minimax sense that one has to pay to solve the (MS), (C) or (L) aggregation problems in expectation, as a function of the cardinality M of the dictionary and of the sample size n . It has been proved in [24] (see also [12] and [28] for the (C) aggregation problem) that

$$\psi_n^{(MS)}(M) \sim \frac{\log M}{n}, \psi_n^{(C)}(M) \sim \begin{cases} \frac{M}{n} & \text{if } M \leq \sqrt{n} \\ \sqrt{\frac{1}{n} \log \left(\frac{eM}{\sqrt{n}} \right)} & \text{if } M > \sqrt{n} \end{cases} \quad \text{and } \psi_n^{(L)}(M) \sim \frac{M}{n}$$

where we denote $a \sim b$ if there are absolute positive constants c and C such that $cb \leq a \leq Cb$. Note that the rates obtained in [24] hold in expectation and in particular, the rate $\psi_n^{(C)}(M)$ was achieved in the gaussian regression model with a known variance and a known marginal distribution of the design. In [8], the authors were able to remove these assumptions at a price of an extra $\log n$ factor for $1 \leq M \leq \sqrt{n}$ (results are still in expectation). We also refer the reader to [6, 28] for non-exact oracle inequalities in the (C) aggregation context.

Lower bounds in deviation follow from the arguments of [24] for the three aggregation problems with the same rates $\psi_n^{(MS)}(M)$, $\psi_n^{(C)}(M)$ and $\psi_n^{(L)}(M)$. In other words, there exist two absolute constants $c_0, c_1 > 0$ such that for any sample cardinality $n \geq 1$, any cardinality of a dictionary $M \geq 1$ and any aggregation procedure \tilde{f}_n , there exists a dictionary F of size M such that with probability larger than c_0 ,

$$R(\tilde{f}_n) \geq \min_{f \in \Delta(F)} R(f) + c_1 \psi_n^{\Delta(F)}(M), \quad (1.2)$$

where the residual term $\psi_n^{\Delta(F)}(M)$ is $\psi_n^{(MS)}(M)$ (resp. $\psi_n^{(C)}(M)$ or $\psi_n^{(L)}(M)$) when $\Delta(F) = F$ (resp. $\Delta(F) = \text{conv}(F)$ or $\Delta(F) = \text{span}(F)$). Procedures achieving

these rates in deviation have been constructed for the (MS) aggregation problem ([2] and [16]) and the (L) aggregation problem ([15]). So far, there was no example of a procedure that achieves the rate of aggregation $\psi_n^{(C)}(M)$ with high probability for the (C) aggregation problem and the aim of this note is to prove that the most natural procedure, empirical risk minimization over the convex hull of F , achieves the rate of $\psi_n^{(C)}(M)$ in deviation (up to a $\log n$ factor for values of M close to \sqrt{n}).

Indeed, we will show that the procedure \tilde{f}^{ERM-C} minimizing the empirical risk functional

$$f \mapsto R_n(f) = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2, \quad (1.3)$$

in $\text{conv}(F)$ achieves, with high probability, the rate $\min\left(\frac{M}{n}, \sqrt{\frac{\log M}{n}}\right)$ for the (C) aggregation problem (see the exact formulation in Theorem 4.3 in the Appendix). Moreover, we will show that the rate $\psi_n^{(C)}(M)$ can be achieved by \tilde{f}^{ERM-C} for any orthogonal dictionary (formulated in Theorem B). On the other hand, it turns out that the same algorithm is far from the conjectured optimal rate $\psi_n^{(MS)}(M)$ for the (MS) aggregation problem (see Theorem A and [16] for the conjecture).

Our first main result is to prove a lower bound on the performance of \tilde{f}^{ERM-C} (ERM in the convex hull) in the context of the (MS) aggregation problem. In [16], it was proved that this procedure is suboptimal for the problem of (MS) aggregation when the size of the dictionary is of the order of \sqrt{n} . Here we complement the result by providing a lower bound for almost all values of M and n .

Theorem A *There exist two absolute positive constants c_0 and c_1 for which the following holds. For any integer n and M such that $\log M \leq c_0 n^{1/3}$, there exists a dictionary F of cardinality M such that, with probability greater than $9/12$*

$$R(\tilde{f}^{ERM-C}) \geq \min_{f \in F} R(f) + c_2 \psi_n(M),$$

where $\psi_n(M) = M/n$ when $M \leq \sqrt{n}$ and $(n \log(eM/\sqrt{n}))^{-1/2}$ when $M > \sqrt{n}$. Moreover, for the same class F , if $M \geq \sqrt{n}$, then with probability larger than $7/12$,

$$R(\tilde{f}^{ERM-C}) \leq \min_{f \in F} R(f) + c_3 \psi_n(M).$$

Note that the residual term $\psi_n(M)$ of Theorem A is much larger than the optimal rate $\psi_n^{(MS)}(M) = (\log M)/n$ for the (MS) aggregation problem. It shows that ERM in the convex hull satisfies a much stronger lower bound than the one mentioned in (1.2) that holds for any algorithm. This result is of particular importance since optimal

aggregation procedures for the (MS) aggregation problem take their values in $\text{conv}(F)$, and it was thus conjectured that \tilde{f}^{ERM-C} could be an optimal aggregation procedure for the (MS) aggregation problem (cf. [16] for more details on this problem). In [16] it was proved that this not the case for $M = \sqrt{n}$; Theorem A shows that this is not the case for all the values of M and n in the significant range (when M is sub-exponential in n).

The proof of Theorem A requires two separate arguments (as in the proofs of the lower bounds in [28] and [24]). The case $M \leq \sqrt{n}$ is easier, and follows an identical path to the one used in [16] for $M = \sqrt{n}$. Its proof is presented for the sake of completeness, and to allow the reader a comparison with the situation in the other case, when $M > \sqrt{n}$. In the “large M ” range things are very different and we present a more intuitive description of the idea behind the construction in Section 2.

The performance of ERM in the convex hull has been studied for an infinite dictionary in [7], in which estimates on its performance have been obtained in terms of the metric entropy of F . The resulting upper bounds were conjectured to be suboptimal in the case of a finite dictionary, since they provide an upper bound of M/n for every n and M whereas it is possible to achieve the rate $\sqrt{(\log M)/n}$ when $M \geq \sqrt{n}$. Although this result is probably known to experts and relies on standard machinery (see for instance [15, 14]), we present its proof in the Appendix.

The residual term $\min\left(\frac{M}{n}, \sqrt{\frac{\log M}{n}}\right)$ of Theorem 4.3 behaves like $\psi_n^{(C)}(M)$ except for values of M for which $n^{1/2} < M \leq c(\epsilon)n^{1/2+\epsilon}$ for $\epsilon > 0$. And, although there is a gap in this range in the general case, under the additional assumption that the dictionary is orthogonal, this gap can be removed.

Theorem B *For every $b > 0$ there is a constant $c_1(b)$ and an absolute constant c_2 for which the following holds. Let n and M be integers which satisfy that $\log M \leq c_1(b)\sqrt{n}$. Let F be a finite dictionary F of cardinality M and (X, Y) such that $|Y|, \sup_{f \in F} |f(X)| \leq b$. If $F = \{f_1, \dots, f_M\}$ satisfies that $\mathbb{E}f_i(X)f_j(X) = 0$ for any $i \neq j \in \{1, \dots, M\}$, then \tilde{f}^{ERM-C} achieves the rate $\psi_n^{(C)}(M)$: for any $u > 0$, with probability greater than $1 - \exp(-u)$*

$$R(\tilde{f}^{ERM-C}) \leq \min_{f \in \text{conv}(F)} R(f) + c_2 b^2 \max \left[\psi_n^{(C)}(M), \frac{u}{n} \right].$$

Removing the gap in the general case is likely to be a much harder problem, although we believe that the orthogonal case should be the “worst” one.

Finally, a word about notation. Throughout, we denote absolute constants or constants that depend on other parameters by c, C, c_1, c_2 , etc., (and, of course, we will specify when a constant is absolute and when it depends on other parameters). The values of constants may change from line to line. The notation $x \sim y$ (resp.

$x \lesssim y$) means that there exist absolute constants $0 < c < C$ such that $cy \leq x \leq Cy$ (resp. $x \leq Cy$). If $b > 0$ is a parameter then $x \lesssim_b y$ means that $x \leq C(b)y$ for some constant $C(b)$ depending only on b . We denote by ℓ_p^M the space \mathbb{R}^M endowed with the ℓ_p norm. The unit ball there is denoted by B_p^M . We also denote the unit Euclidean sphere in \mathbb{R}^M by \mathcal{S}^{M-1} .

If F is a class of functions, let f^* be a minimizer in F of the true risk; in our case, f^* is the minimizer of $\mathbb{E}(f(X) - Y)^2$. For every $f \in F$ set $\mathcal{L}_f = (Y - f(X))^2 - (Y - f^*(X))^2$, and let $\mathcal{L}_F = \{\mathcal{L}_f : f \in F\}$ be the excess loss class associated with F , the target Y and the quadratic risk.

Acknowledgments

The authors would like to thank Vladimir Koltchinskii for a helpful discussion.

2 On the complexity of B_1^M with respect to ℓ_2^M

The aim of this section is to give some of the ideas needed in the proof of Theorem A in the case $M \geq \sqrt{n}$. It is also presented to explain why the seemingly unlikely fact that the rate

$$\frac{1}{\sqrt{n \log(eM/\sqrt{n})}} \quad (2.1)$$

actually improves as the size of the dictionary M increases in our construction is true.

The example used for this result is a class $F_M = \{0, \pm\phi_1, \dots, \pm\phi_M\}$ where $(\phi_i)_{i=1}^M$ is a bounded orthonormal family of $L_2(P^X)$ and $Y = \phi_{M+1}(X)$ is orthogonal to this family. We also assume that $\Phi(X) = (\phi_1(X), \dots, \phi_M(X))$ is isotropic, that is, for every $\lambda \in \mathbb{R}^n$, $\mathbb{E}\langle \Phi(X), \lambda \rangle^2 = \|\lambda\|_2^2$.

An element in $\text{conv}(F_M)$ is of the form $f_\lambda = \langle \Phi, \lambda \rangle$ for some $\lambda \in B_1^M$, its excess loss is $\mathcal{L}_{f_\lambda} = \langle \Phi, \lambda \rangle^2 - 2\langle \Phi, \lambda \rangle \phi_{M+1}$ and the process one has to minimize is indexed by B_1^M and given by

$$P_n \mathcal{L}_{f_\lambda} = \frac{1}{n} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle^2 - \frac{2}{n} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle \phi_{M+1}(X_i). \quad (2.2)$$

It follows from [21] that the oscillations of the quadratic term $\lambda \in B_1^M \rightarrow |(P_n - P)(\langle \Phi, \lambda \rangle^2)|$ are of lower order, and that the empirical process (2.2) behaves like $\lambda \in B_1^M \rightarrow \|\lambda\|_2^2 - 2n^{-1/2} \langle V, \lambda \rangle$ where $V = n^{-1/2} \sum_{i=1}^n \phi_{M+1} \Phi(X_i)$, while a gaussian approximation shows that V essentially behaves like a standard gaussian vector G in \mathbb{R}^M . Hence, the excess risk $P \mathcal{L}_{\hat{f}} = \|\hat{\lambda}\|_2^2$ of the empirical risk minimization procedure

$\widehat{f} = f_{\widehat{\lambda}}$ will be located around

$$\arg \min_{0 \leq r \leq 1} \min_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} \left(r - 2 \frac{\langle G, \lambda \rangle}{\sqrt{n}} \right) = \arg \min_{0 \leq r \leq 1} \left(r - 2n^{-1/2} \sup_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} \langle G, \lambda \rangle \right).$$

Observe that for every radius $0 < r \leq 1$, $\sup_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} \langle \cdot, \lambda \rangle$, is an interpolation norm, which will be denoted by $\|\cdot\|_{A_r^\circ}$. The problem arises because in the range $1/M \leq r \leq 1$ (which is the range we are interested in), a proportional change in the radius r only results in a logarithmic change in the value of $\mathbb{E} \|G\|_{A_r^\circ}$, which is why one has to obtain a sharp estimate on $\mathbb{E} \|G\|_{A_r^\circ}$ for every r .

It turns out that a rather accurate estimate on the complexity of $B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}$ comes from vectors of “short” support. Namely, for every $I \subset \{1, \dots, M\}$, let \mathcal{S}^I be the set of vectors in \mathcal{S}^{M-1} supported in I . Set

$$C_k = \bigcup_{|I|=k} \frac{1}{\sqrt{k}} \mathcal{S}^I \subset B_1^M \cap \frac{1}{\sqrt{k}} \mathcal{S}^{M-1}.$$

If one replaces $B_1^M \cap \frac{1}{\sqrt{k}} \mathcal{S}^{M-1}$ by C_k , it is much easier to analyze ERM over that set. Indeed, it is straightforward to verify that ERM is likely to choose a vector in C_k , where k minimizes the functional

$$k \rightarrow \frac{1}{\sqrt{k}} - \frac{2}{\sqrt{n}} \cdot \mathbb{E} \sup_{v \in C_k} \langle G, v \rangle = \frac{1}{\sqrt{k}} - \frac{2}{\sqrt{n}} \mathbb{E} \left(\sum_{i=1}^k (g_i^2)^* \right)^{1/2}, \quad (2.3)$$

where (x_i^*) is a non-increasing ordering of the vector $(|x_i|)$.

A sharp estimate on the gaussian quantity reveals that the gap between the “level” k and the “level” ℓ decrease with the dimension M . Thus, the minimum of (2.3) – which is proportional to (2.1)– decreases as M increases.

The proof of Theorem A will be a combination of two approximation arguments – first, of the measure $n^{-1/2} \sum_{i=1}^n X_i$ by a gaussian, and second, an approximation of B_1^M by the sets C_k , reducing the problem to the one described above.

One should comment that it is possible to approximate $B_1^M \cap \frac{1}{\sqrt{k}} \mathcal{S}^{M-1}$ using a completely combinatorial set $\cup_{|I|=k} k^{-1} \{-1, 1\}^I$, and the way the complexities change between the levels k and ℓ as M increases gives a more geometric explanation to why the minimizer moves closer to 0.

3 Proof of the lower bound for the (MS) aggregation problem (Theorem A)

The proof of Theorem A consists of two parts. The first, simpler part, is when $M \leq \sqrt{n}$. This is due to the fact that if $0 < \theta < 1$ and $\rho = \theta r \sim M/n$, the set

$B_1^M \cap \sqrt{r}S^{M-1}$ is much “larger” than the set $B_1^M \cap \sqrt{\rho}B_2^M$. This results in much larger “oscillations” of the appropriate empirical process on the former set than on the latter one, leading to very negative values of the empirical excess risk functional for functions whose excess risk larger than ρ . The case $M \geq \sqrt{n}$ is much harder because when considering the required values of r and ρ , the complexity of the two sets is very close, and comparing the two oscillations accurately involves a far more delicate analysis.

3.1 The case $M \leq \sqrt{n}$

We will follow the method used in [16]. Let $(\phi_i)_{i \in \mathbb{N}}$ be a sequence of functions defined on $[0, 1]$ and set μ to be a probability measure on $[0, 1]$ such that $(\phi_i : i \in \mathbb{N})$ is a sequence of independent Rademacher variables in $L_2([0, 1], \mu)$.

Let $M \leq \sqrt{n}$ be fixed and put (X, Y) to be a couple of random variables; X is distributed according to μ and $Y = \phi_{M+1}(X)$. Let $F = \{0, \pm\phi_1, \dots, \pm\phi_M\}$ be the dictionary, and note that any function in the convex hull of F can be written as $f_\lambda = \sum_{j=1}^M \lambda_j \phi_j$ for $\lambda \in B_1^M$. Since relative to $\text{conv}(F)$, $f^* = 0$, the excess quadratic loss function is

$$\mathcal{L}_\lambda(X, Y) = -2\phi_{M+1}(X)\langle \lambda, \Phi(X) \rangle + \langle \lambda, \Phi(X) \rangle^2$$

where we set $\Phi(\cdot) = (\phi_1(\cdot), \dots, \phi_M(\cdot))$.

The following is a reformulation of Lemma 5.4 in [16].

Lemma 3.1 *There exist absolute constants c_0, c_1 and c_2 for which the following holds. Let $(X_i, Y_i)_{i=1, \dots, n}$ be n independent copies of (X, Y) . Then, for every $r > 0$, with probability greater than $1 - 8 \exp(-c_0 M)$, for any $\lambda \in \mathbb{R}^M$,*

$$\left| \|\lambda\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle^2 \right| \leq \frac{1}{2} \|\lambda\|_2^2 \quad (3.1)$$

and

$$c_1 \sqrt{\frac{rM}{n}} \leq \sup_{\lambda \in \sqrt{r}B_2^M} \frac{1}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle \phi_{M+1}(X_i) \leq c_2 \sqrt{\frac{rM}{n}}. \quad (3.2)$$

Set $r = \beta M/n$ for some $0 < \beta \leq 1$ to be named later, and observe that $B_1^M \cap \sqrt{r}S^{M-1} = \sqrt{r}S^{M-1}$ because $r \leq 1/M$. For any $\lambda \in \sqrt{r}S^{M-1}$, $P\mathcal{L}_\lambda = \|\lambda\|_2^2 = r$, and thus applying (3.1) and (3.2), it is evident that with probability greater than

$1 - 8 \exp(-c_0 M)$,

$$\begin{aligned}
& \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} P_n \mathcal{L}_\lambda = r - \sup_{\lambda \in \sqrt{r} S^{M-1}} (P - P_n) \mathcal{L}_\lambda \\
& \leq r + \sup_{\lambda \in \sqrt{r} S^{M-1}} \left| \|\lambda\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle^2 \right| - \sup_{\lambda \in \sqrt{r} S^{M-1}} \frac{2}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle \phi_{M+1}(X_i) \\
& \leq \frac{3r}{2} - 2c_1 \sqrt{\frac{rM}{n}} = \left(\frac{3\beta}{2} - 2c_1 \sqrt{\beta} \right) \frac{M}{n} \leq -c_1 \sqrt{\beta} \frac{M}{n},
\end{aligned}$$

provided that $\beta \leq (2c_1/3)^2$.

On the other hand, let $\rho = \alpha M/n$ for some α to be chosen later. Using (3.1) and (3.2) again, it follows that with probability at least $1 - 8 \exp(-c_0 M)$, for any $\lambda \in B_1^M \cap \sqrt{\rho} B_2^M$

$$\begin{aligned}
|P_n \mathcal{L}_\lambda| & \leq P \mathcal{L}_\lambda + \left| \|\lambda\|_2^2 - \frac{1}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle^2 \right| + \left| \frac{2}{n} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle \phi_{M+1}(X_i) \right| \\
& \leq \frac{3\rho}{2} + 2c_2 \sqrt{\frac{\rho M}{n}} = \left(\frac{3\alpha}{2} + 2c_2 \sqrt{\alpha} \right) \frac{M}{n}.
\end{aligned}$$

Therefore, if $0 < \alpha < \beta$ satisfies that $3\alpha/2 + 2c_2 \sqrt{\alpha} < c_1 \sqrt{\beta}$ for some $0 < \beta \leq (2c_1/3)^2$ then with probability greater than $1 - 16 \exp(-c_0 M)$, the empirical risk function $\lambda \mapsto R_n(f_\lambda)$ achieves smaller values on $B_1^M \cap \sqrt{r} S^{M-1}$ than on $B_1^M \cap \sqrt{\rho} B_2^M$. Hence, with the same probability, $R(\tilde{f}^{ERM-C}) \geq \rho = \alpha M/n$. ■

3.2 The case $M \geq \sqrt{n}$

Let us reformulate the second part of Theorem A.

Theorem 3.2 *There exist absolute constants c_0, c_1, c_2 and n_0 for which the following holds. For every integers $n \geq n_0$ and M , if $\sqrt{n} \leq M \leq \exp(c_0 n^{1/3})$, there is a function class F_M of cardinality M consisting of functions that are bounded by 1, and a couple (X, Y) distributed according to a probability measure μ , such that with $\mu^{\otimes n}$ -probability at least $9/12$,*

$$R(\hat{f}) \geq \min_{f \in F_M} R(f) + \frac{c_1}{\sqrt{n \log(eM/\sqrt{n})}},$$

where \hat{f} is the empirical minimizer in $\text{conv}(F_M)$. Moreover, with $\mu^{\otimes n}$ -probability greater than $7/12$,

$$R(\hat{f}) \leq \min_{f \in F_M} R(f) + \frac{c_2}{\sqrt{n \log(eM/\sqrt{n})}}.$$

The proof will require accurate information on a monotone rearrangement of almost gaussian random variables.

Lemma 3.3 *There exists an absolute constant C for which the following holds. Let g be a standard gaussian random variable, set $H(x) = \mathbb{P}(|g| > x)$ and put $W(p) = H^{-1}(p)$ (the inverse function of H). Then for every $0 < p < 1$,*

$$|W^2(p) - \log(2/(\pi p^2)) + \log(\log(2/(\pi p^2)))| \leq C \frac{\log \log(2/(\pi p^2))}{\log(2/(\pi p^2))}.$$

Moreover, for every $0 < \epsilon < 1/2$ and $0 < p < 1/(1 + \epsilon)$,

$$|W^2(p) - W^2((1 + \epsilon)p)| \leq C\epsilon, \quad |W^2(p) - W^2((1 - \epsilon)p)| \leq C\epsilon.$$

Proof. The proof of the first part follows from the observation that for every $x > 0$,

$$\frac{\sqrt{2}}{x\sqrt{\pi}} \exp(-x^2/2) \left(1 - \frac{1}{x^2}\right) \leq \mathbb{P}(|g| > x) \leq \frac{\sqrt{2}}{x\sqrt{\pi}} \exp(-x^2/2), \quad (3.3)$$

where c is a suitable absolute constant (see, e.g. [22]), combined with a straightforward (yet tedious) computation. The second part of the claim follows from the first one, and is omitted. \blacksquare

The next step is a gaussian approximation of a variable $Y = n^{-1/2} \sum_{i=1}^n X_i$, where X_1, \dots, X_n are i.i.d random variables, with mean zero, variance 1, under the additional assumption that X has well behaved tails.

Definition 3.4 [18, 26] *Let $1 \leq \alpha \leq 2$. We say that a random variable X belongs to L_{ψ_α} if there exists a constant C such that*

$$\mathbb{E} \exp(|X|^\alpha / C^\alpha) \leq 2. \quad (3.4)$$

The infimum over all constants C for which (3.4) holds defines a norm called the ψ_α norm of X , and we denote it by $\|X\|_{\psi_\alpha}$.

Proposition 3.5 ([22], pg. 183) *For every L there exist constants c_1 and c_2 that depend only on L and for which the following holds. Let $(X_n)_{n \in \mathbb{N}}$ be a sequence of i.i.d., mean zero random variables with variance 1, and $\|X\|_{\psi_1} \leq L$. If $Y = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i$ then for any $0 < x \leq c_1 n^{1/6}$,*

$$\mathbb{P}[Y \geq x] = \mathbb{P}[g \geq x] \exp\left(\frac{\mathbb{E}X_1^3 x^3}{6\sqrt{n}}\right) \left[1 + c_2 \frac{x+1}{\sqrt{n}}\right]$$

and

$$\mathbb{P}[Y \leq -x] = \mathbb{P}[g \leq -x] \exp\left(\frac{-\mathbb{E}X_1^3 x^3}{6\sqrt{n}}\right) \left[1 + c_2 \frac{x+1}{\sqrt{n}}\right].$$

In particular, if $0 < x \leq c_1 n^{1/6}$ and $\mathbb{E}X_1^3 = 0$ then

$$|\mathbb{P}[|Y| \geq x] - \mathbb{P}[|g| \geq x]| = c_2 \mathbb{P}[|g| \geq x] \frac{x+1}{\sqrt{n}}.$$

Since Proposition 3.5 implies a better gaussian approximation than the standard Berry-Esséen bounds, one may consider the following family of random variables that will be used in the construction.

Definition 3.6 *We say that a random variable Y is (L, n) -almost gaussian for $L > 0$ and $n \in \mathbb{N}$, if $Y = n^{-1/2} \sum_{i=1}^n X_i$, where X_1, \dots, X_n are independent copies of X , which is a non-atomic random variable with mean 0, variance 1, and satisfies that $\mathbb{E}X^3 = 0$ and $\|X\|_{\psi_1} \leq L$.*

Let X_1, \dots, X_n and Y be such that $Y = n^{-1/2} \sum_{i=1}^n X_i$ is (L, n) -almost gaussian. For $0 < p < 1$ set

$$U(p) = \{x > 0 : \mathbb{P}(|Y| > x) = p\}.$$

Since X is non-atomic then $U(p)$ is non-empty and let

$$u^+(p) = \sup U(p) \quad \text{and} \quad u^-(p) = \inf U(p).$$

We shall apply Lemma 3.3 and Proposition 3.5 in the following case to bound $u^+(i/M)$ and $u^-(i/M)$ for every i , as long as M is not too large (i.e. $\log M \leq c_1 n^{1/3}$). To that end, set $\epsilon_{M,n} = [(\log M)/n]^{1/2}$, and for fixed values of M and n , and $1 \leq i \leq M$ let

$$u_i^+ = u^+(i/M) \quad \text{and} \quad u_i^- = u^-(i/M).$$

Corollary 3.7 *For every $L > 0$ there exist a constant C_0 that depends on L and an absolute constant C_1 for which the following holds. Assume that Y is (L, n) -almost gaussian and that $\log M \leq C_0 n^{1/3}$. Then, for every $1 \leq i \leq M/2$,*

$$(u_i^+)^2 \leq \log\left(\frac{2M^2}{\pi i^2}\right) - \log\left(\log\left(\frac{2M^2}{\pi i^2}\right)\right) + C_1 \max\left\{\frac{\log(\log(2M^2/(\pi i^2)))}{\log(2M^2/(\pi i^2))}, \epsilon_{M,n}\right\},$$

and

$$(u_i^-)^2 \geq \log\left(\frac{2M^2}{\pi i^2}\right) - \log\left(\log\left(\frac{2M^2}{\pi i^2}\right)\right) - C_1 \max\left\{\frac{\log(\log(2M^2/(\pi i^2)))}{\log(2M^2/(\pi i^2))}, \epsilon_{M,n}\right\}.$$

Proof. Since $\sqrt{\log M} \leq C_0 n^{1/6}$, one may use the gaussian approximation from Proposition 3.5 to obtain

$$\begin{aligned} \mathbb{P}[|Y| \geq \sqrt{4 \log M}] &\leq \mathbb{P}[|g| \geq \sqrt{4 \log M}] \left(1 + c_1 \left(\frac{\sqrt{4 \log M} + 1}{\sqrt{n}}\right)\right) \\ &\leq \sqrt{\frac{2}{4\pi \log M}} \exp(-2 \log M) \left(1 + c_1 \left(\frac{\sqrt{4 \log M} + 1}{\sqrt{n}}\right)\right) \leq \frac{1}{M^2}. \end{aligned}$$

Thus, for every $1 \leq i \leq M$, if $x \in U(i/M)$ then $x \leq \sqrt{4 \log M}$.

Let $1 \leq i \leq M/2$ and $x \in U(i/M)$. Since $x \leq 2C_0 n^{1/6}$ (because $x \leq \sqrt{4 \log M} \leq 2C_0 n^{1/6}$), it follows from Proposition 3.5 that

$$|i/M - H(x)| \leq c_3 H(x) \frac{x+1}{\sqrt{n}} \leq c_4 H(x) \epsilon_{M,n}, \quad (3.5)$$

where $H(x) = \mathbb{P}[|g| \geq x]$. Observe that if $W(p) = H^{-1}(p)$, then

$$|W^2(i/M) - x^2| \leq c_5 \epsilon_{M,n}.$$

Indeed, since $H(x)(1 - c_4 \epsilon_{M,n}) \leq i/M \leq H(x)(1 + c_4 \epsilon_{M,n})$, then by the monotonicity of W and the second part of Lemma 3.3, setting $p = H(x)$,

$$W^2(i/M) \leq W^2((1 + c_4 \epsilon_{M,n})H(x)) \leq W^2(H(x)) + c_6 \epsilon_{M,n} = x^2 + c_6 \epsilon_{M,n}.$$

One obtains the lower bound in a similar way. The claim follows by using the approximate value of $W^2(i/M)$ provided in the first part of Lemma 3.3. ■

The parameters u_i^+ and u_i^- can be used to estimate the distribution of a non-increasing rearrangement $(Y_i^*)_{i=1}^M$ of the absolute values of M independent copies of Y .

Lemma 3.8 *There exists constants $c > 0$ and $j_0 \in \mathbb{N}$ for which the following holds. Let Y_1, \dots, Y_M be i.i.d. non-atomic random variables. For every $1 \leq s \leq M$, with probability at least $1 - 2 \exp(-cs)$,*

$$|\{i : |Y_i| \geq u_s^-\}| \geq s/2 \quad \text{and} \quad |\{i : |Y_i| \geq u_s^+\}| \leq 3s/2.$$

In particular, with probability at least $11/12$, for every $j_0 \leq j \leq M/2$,

$$u_{2j}^- \leq Y_j^* \leq u_{\lceil 2(j-1)/3 \rceil}^+,$$

where $\lceil x \rceil = \min\{n \in \mathbb{N} : x \leq n\}$.

Proof. Fix $0 < p < 1$ to be named later and let $(\delta_i)_{i=1}^M$ be independent $\{0, 1\}$ -valued random variables with $\mathbb{E}\delta_i = p$. A straightforward application of Bernstein's inequality [26] shows that

$$\mathbb{P}\left(\left|\frac{1}{M}\sum_{i=1}^M\delta_i - p\right|\geq t\right)\leq 2\exp(-cM\min\{t^2/p, t\}).$$

In particular, with probability at least $1 - 2\exp(-c_1Mp)$,

$$(1/2)Mp\leq\sum_{i=1}^M\delta_i\leq(3/2)Mp.$$

We will apply this observation to the independent random variables $\delta_i = \mathbb{1}_{\{|Y_i|>a\}}$, $1 \leq i \leq M$ for an appropriate choice of a . Indeed, if we take a for which $\mathbb{P}(|Y_1| > a) = s/M$ (such an a exists because Y_1 is non-atomic), then with probability at least $1 - 2\exp(-c_1s)$, at least $s/2$ of the $|Y_i|$ will be larger than a , and at most $3s/2$ will be larger than a . Since this result holds for any $a \in U(s/M)$ the first part of the claim follows.

Now take s_0 to be the smallest integer such that $1 - 2\sum_{s=s_0}^M\exp(-cs) \geq 11/12$ (in particular $c^{-1}\log 24 \leq s_0 \leq c^{-1}(\log 48 + 1)$). Applying the union bound and a change of variables, it is evident that with probability at least $5/6$, for every $\lfloor(3s_0)/2\rfloor + 1 \leq j \leq M/2$,

$$|\{i : |Y_i| \geq u_{2j}^-\}| \geq j \quad \text{and} \quad |\{i : |Y_i| \geq u_{\lfloor(2(j-1))/3\rfloor}^+\}| \leq j - 1,$$

and thus $u_{2j}^- \leq Y_j^* \leq u_{\lfloor(2(j-1))/3\rfloor}^+$. ■

With Lemma 3.8 and Corollary 3.7 in hand, one can bound the following functional of the random variables $(Y_i^*)_{i=1}^M$.

Lemma 3.9 *For every $L > 0$ there exist constants c_1, \dots, c_4, j_0 and $\alpha < 1$ that depend only on L for which the following holds. Let Y be (L, n) -almost gaussian and let Y_1, \dots, Y_M be independent copies of Y . Then, with probability at least $11/12$, for every $j_0 \leq \ell \leq k \leq \alpha M$,*

$$c_1\left(\frac{\log(ek/\ell) - \epsilon_{M,n}}{\sqrt{\log(eM/\ell)}}\right)\leq Y_\ell^* - Y_k^* \leq c_2\left(\frac{\log(ek/\ell) + \epsilon_{M,n}}{\sqrt{\log(eM/\ell)}}\right).$$

Moreover, with probability at least $10/12$, for every $j_0 \leq \ell \leq k \leq \alpha M$

$$Y_\ell^* - Y_k^* - \left(\frac{1}{k}\sum_{i=1}^k(Y_i^* - Y_k^*)^2\right)^{1/2} \geq c_3\frac{\log(ek/\ell)}{\sqrt{\log(eM/\ell)}} - \frac{c_4}{\sqrt{\log(eM/k)}},$$

and if $j_0 \leq k \leq \alpha M$, then $u_{2k}^- \leq Y_k^* \leq u_{\lceil 2(k-1)/3 \rceil}^+$ and

$$\left(\frac{1}{k} \sum_{i=1}^k (Y_i^* - Y_k^*)^2 \right)^{1/2} \leq \frac{c_4}{\sqrt{\log(eM/k)}},$$

provided that $\log^2 M \lesssim_L k$ and that $\epsilon_{M,n} = \sqrt{(\log M)/n} \leq 1$.

Proof. The first part of the claim follows from Lemma 3.8 and Corollary 3.7, combined with a straightforward computation. For the second part, observe that, for some well chosen constant $c_1(L)$ depending only on L , with probability at least $11/12$, $Y_1^* \leq c_1(L)\sqrt{\log M}$. Hence, applying the first part of the claim, with probability at least $10/12$,

$$\begin{aligned} \frac{1}{k} \sum_{i=1}^k (Y_i^* - Y_k^*)^2 &\leq c_1(L) \frac{j_0 \log M}{k} + \frac{1}{k} \sum_{i=j_0}^k (Y_i^* - Y_k^*)^2 \\ &\leq c_1(L) \frac{j_0 \log M}{k} + \frac{c_2}{k} \sum_{i=j_0}^k \left(\frac{\log^2(ek/i)}{\log(eM/i)} + \frac{\epsilon_{M,n}^2}{\log(eM/i)} \right) \\ &\leq c_1(L) \frac{j_0 \log M}{k} + c_3 \frac{1 + \epsilon_{M,n}^2}{\log(eM/k)} \leq \frac{c_4}{\log(eM/k)}, \end{aligned}$$

provided that $\log^2 M \lesssim_L k$ and that $\epsilon_{M,n} \leq 1$. Note that to estimate the sum we have used that

$$\frac{1}{k} \sum_{i=j_0}^k \frac{\log^2(ek/i)}{\log(eM/i)} \leq \frac{1}{\log(eM/k)} \frac{1}{k} \sum_{i=j_0}^k \log^2(ek/i) \leq \frac{c_3}{\log(eM/k)}.$$

Now the second and the third parts follow from the first one. \blacksquare

The next preliminary step we need is a simple bound on the dual norm to the one whose unit ball is $A_r = B_1^M \cap \sqrt{r}B_2^M$. Recall that for a convex body $C \subset \mathbb{R}^M$, the polar body of C is $C^\circ = \{x \in \mathbb{R}^M : \sup_{y \in C} \langle x, y \rangle \leq 1\}$, and in our case, $A_r^\circ = \text{conv}(B_\infty^M \cup r^{-1/2}B_2^M)$ (see, for example, [23]). From here on, given $v \in \mathbb{R}^M$, set

$$\|v\|_{A_r^\circ} = \sup_{w \in A_r} \langle v, w \rangle,$$

and, as always, $(v_i^*)_{i=1}^M$ is the monotone rearrangement of $(|v_i|)_{i=1}^M$.

Lemma 3.10 For every $v \in \mathbb{R}^M$ and any $0 < \rho < r \leq 1$ such that $1/r$ and $1/\rho$ are integers,

$$\|v\|_{A_r^\circ} - \|v\|_{A_\rho^\circ} \geq v_{1/r}^* - v_{1/\rho}^* - \left(\rho \sum_{i=1}^{1/\rho} (v_i^* - v_{1/\rho}^*)^2 \right)^{1/2}$$

and in general for any $0 < r \leq 1$,

$$v_{[1/r]}^* \leq \|v\|_{A_r^\circ} \leq v_{[1/r]}^* + \sqrt{[1/r]} \left(\sum_{i=1}^{[1/r]} (v_i^* - v_{[1/r]}^*)^2 \right)^{1/2}.$$

Proof. First, observe that for every $v \in \mathbb{R}^M$,

$$\|v\|_{A_r^\circ} = \min_{1 \leq j \leq M} \left(\sqrt{r} \left(\sum_{i=1}^j (v_i^* - v_j^*)^2 \right)^{1/2} + v_j^* \right). \quad (3.6)$$

Indeed, since $A_r^\circ = \text{conv}(B_\infty^M \cup (1/\sqrt{r})B_2^M)$, it is evident that $\|v\|_{A_r^\circ} = \inf\{\|u\|_\infty + \sqrt{r}\|w\|_2, v = u + w\}$. One may verify that if $v = u + w$ is an optimal decomposition then $\text{supp}(w) \subset \{i : |u_i| = \|u\|_\infty\}$. Hence, if $\|u\|_\infty = K$ then for every $1 \leq i \leq M$, $u_i = K \text{sgn}(v_i) \mathbb{1}_{\{|v_i| \geq K\}} + v_i \mathbb{1}_{\{|v_i| < K\}}$, and thus, $w_i = \mathbb{1}_{\{|v_i| \geq K\}}(v_i - \text{sgn}(v_i)K)$. Therefore,

$$\|v\|_{A_r^\circ} = \inf_{K > 0} \left\{ K + \sqrt{r} \left(\sum_{\{i: |v_i| \geq K\}} (|v_i| - K)^2 \right)^{1/2} \right\}.$$

Moreover, since it is enough to consider only values of K in $\{v_j^* : 1 \leq j \leq M\}$, (3.6) is verified. In particular, if $1/r$ is an integer then

$$\|v\|_{A_r^\circ} \leq \sqrt{r} \left(\sum_{i=1}^{1/r} (v_i^* - v_{1/r}^*)^2 \right)^{1/2} + v_{1/r}^*.$$

On the other hand, if $T_r = \{u \in \mathbb{R}^M : \|u\|_2 \leq \sqrt{r}, |\text{supp}(u)| \leq 1/r\}$ then $T_r \subset B_1^M \cap \sqrt{r}B_2^M$. Hence,

$$\|v\|_{A_r^\circ} \geq \sup_{w \in T_r} \langle v, w \rangle = \sqrt{r} \left(\sum_{i=1}^{1/r} (v_i^*)^2 \right)^{1/2}.$$

Therefore, if $1/r$ and $1/\rho$ are integers, it follows that

$$\begin{aligned} \|v\|_{A_r^\circ} - \|v\|_{A_\rho^\circ} &\geq \sqrt{r} \left(\sum_{i=1}^{1/r} (v_i^*)^2 \right)^{1/2} - \left(\sqrt{\rho} \left(\sum_{i=1}^{1/\rho} (v_i^* - v_{1/\rho}^*)^2 \right)^{1/2} + v_{1/\rho}^* \right) \\ &\geq v_{1/r}^* - v_{1/\rho}^* - \sqrt{\rho} \left(\sum_{i=1}^{1/\rho} (v_i^* - v_{1/\rho}^*)^2 \right)^{1/2}, \end{aligned}$$

because $(r \sum_{i=1}^{1/r} (v_i^*)^2)^{1/2} \geq v_{1/r}^*$.

The second part follows in a similar fashion and it omitted. \blacksquare

Proof of the lower bound of Theorem 3.2. Let ϕ_1, \dots, ϕ_M, X and $a > 0$ be such that $\phi_1(X), \dots, \phi_M(X)$ are uniformly distributed on $[-a, a]$ and have variance 1 (in particular $a = \sqrt{3}$). Set $T(X) = \phi_{M+1}(X) = Y$ to be a Rademacher variable. Assume further that $(\phi_i)_{i=1}^{M+1}$ are independent in $L_2(P^X)$ and let $F_M = \{0, \pm\phi_1, \dots, \pm\phi_M\}$. Note that the functions in $\text{conv}(F_M)$ are given by $f_\lambda = \langle \Phi, \lambda \rangle$ where $\Phi = (\phi_1, \dots, \phi_M)$ and $\lambda \in B_1^M$.

It is straightforward to verify that the excess loss function of f_λ relative to $\text{conv}(F_M)$ is

$$\mathcal{L}_{f_\lambda} = (f_\lambda - \phi_{M+1})^2 - (0 - \phi_{M+1})^2 = \langle \Phi, \lambda \rangle^2 - 2\langle \Phi, \lambda \rangle \phi_{M+1}$$

(since $f^* = 0$), implying that $\mathbb{E}\mathcal{L}_{f_\lambda} = \|\lambda\|_2^2$.

Let us consider the problem of empirical minimization in $\text{conv}(F_M) = \{\langle \lambda, \Phi \rangle : \lambda \in B_1^M\}$. Recall that $A_r = B_1^M \cap \sqrt{r}B_2^M$ and, for an independent sample $(\Phi(X_i), \phi_{M+1}(X_i))_{i=1}^n$, define the functional

$$\psi(r, \rho) = n \left(\inf_{\lambda \in A_r} R_n(f_\lambda) - \inf_{\mu \in A_\rho} R_n(f_\mu) \right).$$

If we show that for some $r \geq \rho$, $\psi(r, \rho) < 0$, then for that sample, $\mathbb{E}\mathcal{L}_{\hat{f}} \geq \rho$.

Note that, for any $r, \rho > 0$,

$$\begin{aligned} \psi(r, \rho) &\leq \sup_{\lambda \in A_r} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle^2 - 2 \sup_{\lambda \in A_r} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle \phi_{M+1}(X_i) \\ &\quad + 2 \sup_{\mu \in A_\rho} \sum_{i=1}^n \langle \Phi(X_i), \mu \rangle \phi_{M+1}(X_i), \end{aligned}$$

and let us estimate the supremum of the process

$$\lambda \in A_r \rightarrow \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle^2 = n \left((P_n - P)(\langle \Phi, \lambda \rangle^2) + \|\lambda\|_2^2 \right).$$

Observe that $\Phi(X)$ is isotropic (that is, for every $\lambda \in \mathbb{R}^M$, $\mathbb{E}\langle \lambda, \Phi(X) \rangle^2 = \|\lambda\|_2^2$), and subgaussian – since $\|\langle \lambda, \Phi(X) \rangle\|_{\psi_2} \leq 4a \|\lambda\|_2$. Hence, applying the results from [21], it is evident that with probability at least $11/12$,

$$\sup_{\lambda \in A_r} \left| (P_n - P)(\langle \lambda, \Phi \rangle^2) \right| \leq c(a) \max \left\{ \text{diam}(A_r, \|\cdot\|_2) \frac{\gamma_2(A_r, \|\cdot\|_2)}{\sqrt{n}}, \frac{\gamma_2^2(A_r, \|\cdot\|_2)}{n} \right\}. \quad (3.7)$$

Recall that for $r \geq 1/M$, $\gamma_2(A_r, \|\cdot\|_2) \sim \sqrt{\log(eMr)}$ (see, for instance, [21]), and thus, if $r \geq \max(1/M, 1/n)$, then with probability at least $11/12$,

$$n \sup_{\lambda \in A_r} \left((P_n - P)(\langle \Phi, \lambda \rangle^2) + \|\lambda\|_2^2 \right) \leq nr + c_1 \sqrt{nr \log(eMr)},$$

where c_1 is a constant that depends only on a .

Next, to estimate the first two terms, let

$$Y_j = n^{-1/2} \sum_{i=1}^n \phi_{M+1}(X_i) \langle \Phi(X_i), e_j \rangle$$

and observe that $(Y_j)_{j=1}^M$ are independent copies of a $(2, n)$ -almost gaussian variable. If we set $V = (Y_i)_{i=1}^M$ then

$$\begin{aligned} & \sup_{\lambda \in A_r} \sum_{i=1}^n \langle \lambda, \Phi(X_i) \rangle \phi_{M+1}(X_i) - \sup_{\theta \in A_\rho} \sum_{i=1}^n \langle \theta, \Phi(X_i) \rangle \phi_{M+1}(X_i) \\ &= \sqrt{n} \left(\sup_{\lambda \in A_r} \langle \lambda, V \rangle - \sup_{\theta \in A_\rho} \langle \theta, V \rangle \right) = \sqrt{n} (\|V\|_{A_r^\circ} - \|V\|_{A_\rho^\circ}) = (*). \end{aligned}$$

By Lemma 3.10, if $1/r = \ell$ and $1/\rho = k$ are integers, then

$$(*) \geq \sqrt{n} \left[Y_\ell^* - Y_k^* - \left(\frac{1}{k} \sum_{i=1}^k (Y_i^* - Y_k^*)^2 \right)^{1/2} \right]$$

and thus, if ℓ, k, M and n are as in Lemma 3.9, then with probability at least $9/12$,

$$\begin{aligned} (*) & \geq \sqrt{n} \left(\frac{c_2 \log(ek/\ell)}{\sqrt{\log(eM/\ell)}} - \frac{c_3}{\sqrt{\log(eM/k)}} \right) \\ & \geq c_4 \sqrt{n} \frac{\log(ek/\ell)}{\sqrt{\log(eM/\ell)}}, \end{aligned}$$

provided that $k \geq c_5 \ell$ for c_5 large enough.

Hence, with probability at least 9/12,

$$\psi(r, \rho) \leq -2c_4\sqrt{n} \frac{\log(ek/\ell)}{\sqrt{\log(eM/\ell)}} + nr + c_1\sqrt{rn \log(eMr)}.$$

It follows that if we select $r \sim 1/\sqrt{n \log(eM/\sqrt{n})}$ and $\rho \sim r$ with $\rho < r$ so that the conditions of Lemma 3.9 are satisfied, then with probability at least 9/12, $\psi(r, \rho) < 0$. Hence, with the same probability,

$$R(\hat{f}) - \min_{f \in F_M} R(f) = \mathbb{E} \mathcal{L}_{\hat{f}} \geq \frac{c_6}{\sqrt{n \log(eM/\sqrt{n})}}.$$

■

Proof of the upper bound in Theorem 3.2. We will show that with constant probability,

$$\inf_{0 \leq r \leq r_0} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} R_n(f_\lambda) < \inf_{r_0 \leq r \leq 1} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} R_n(f_\lambda) \quad (3.8)$$

for $r_0 \sim 1/\sqrt{n \log(eM/\sqrt{n})}$, and thus, on that event, $R(\hat{f}) \leq r_0$. To that end, one has to show that

$$\inf_{0 \leq r \leq r_0} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} P_n \mathcal{L}_{f_\lambda} < \inf_{r_0 \leq r \leq 1} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} P_n \mathcal{L}_{f_\lambda}.$$

Let

$$Q(r) = \sup_{\lambda \in B_1^M \cap \sqrt{r} B_2^M} |(P_n - P)(\langle \Phi, \lambda \rangle^2)|$$

and set $r^* = \inf \{r > 0 : \mathbb{E} Q(r) \leq r/2\}$. Applying (3.7) and since $\gamma_2(A_r, \|\cdot\|_2) \sim \sqrt{\log(eMr)}$, then $r^* \leq c_0 \sqrt{\log(eM/\sqrt{n})}/n$. Hence, by a standard fixed point argument (see for instance, [4]), it follows that with probability greater than 11/12, if $\lambda \in B_1^M$ and $\|\lambda\|_2^2 \geq r^*$, then

$$\frac{\|\lambda\|_2^2}{2} \leq \frac{1}{n} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle^2 \leq \frac{3\|\lambda\|_2^2}{2}.$$

In particular, by Lemma 3.9, Lemma 3.10 and Corollary 3.7, with probability larger than 9/12, for every $r \geq r^*$,

$$\begin{aligned}
\inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda} &= \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} \left(\frac{1}{n} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle^2 - \frac{2}{n} \sum_{i=1}^n \langle \Phi(X_i), \lambda \rangle \phi_{M+1}(X_i) \right) \\
&\geq \frac{r}{2} - \frac{2}{\sqrt{n}} \sup_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} \langle \lambda, V \rangle \geq \frac{r}{2} - \frac{2}{\sqrt{n}} \sup_{\lambda \in B_1^M \cap \sqrt{r} B_2^M} \langle \lambda, V \rangle = \frac{r}{2} - \frac{2 \|V\|_{A_r^\circ}}{\sqrt{n}} \\
&\geq \frac{r}{2} - \frac{2}{\sqrt{n}} \left(Y_{\lceil 1/r \rceil}^* + \sqrt{\lceil 1/r \rceil} \left(\sum_{i=1}^{\lceil 1/r \rceil} (Y_i^* - Y_{\lceil 1/r \rceil}^*)^2 \right)^{1/2} \right) \\
&\geq \frac{r}{2} - \frac{2}{\sqrt{n}} \left(u_{\lceil 2(\lceil 1/r \rceil - 1)/3 \rceil}^+ + \frac{c_1}{\sqrt{\log(c_2 M r)}} \right) \\
&\geq \frac{r}{2} - \frac{2}{\sqrt{n}} \left(c_3 \sqrt{\log(c_4 M r)} + \frac{c_1}{\sqrt{\log(c_5 M r)}} \right) > 0
\end{aligned} \tag{3.9}$$

provided that $r \geq c_6 \sqrt{\log(eM/\sqrt{n})/n}$ for some constant c_6 large enough. Therefore, on that event, if $\|\lambda\|_2^2 \geq c_7 \sqrt{\log(eM/\sqrt{n})/n}$ then $P_n \mathcal{L}_{f_\lambda} > 0$. On the other hand, $P_n \mathcal{L}_{f_0} = 0$, and thus $\|\hat{\lambda}\|_2^2 \leq c_7 \sqrt{\log(eM/\sqrt{n})/n}$ (where $\hat{f} = f_{\hat{\lambda}}$) with probability at most 9/12.

It remains to show that with sufficiently high constant probability

$$\inf_{0 \leq r \leq r_0} \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda} < \inf_{r_0 \leq r \leq r_1} \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda}$$

for $r_0 \sim 1/\sqrt{n \log(eM/\sqrt{n})}$ and $r_1 = c_7 \sqrt{\log(eM/\sqrt{n})/n}$.

Using the same argument as in (3.9) and applying Lemma 3.9, Lemma 3.10 and Corollary 3.7, it is evident that with probability at least 10/12,

$$\begin{aligned}
&\inf_{r_0 \leq r \leq r_1} \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda} \\
&\geq \inf_{r_0 \leq r \leq r_1} \left(r - \frac{2}{\sqrt{n}} \left(\log(C_0 M r) - c_8 \log \log(C_0 M r_0) \right)^{1/2} \right) - \frac{c_1}{\sqrt{n \log(c_5 M r_0)}} - Q(r_1)
\end{aligned} \tag{3.10}$$

and for some $r_2 \leq r_0$ to be named later,

$$\begin{aligned}
&\inf_{0 \leq r \leq r_0} \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda} \leq \inf_{0 \leq r \leq r_2} \inf_{\lambda \in B_1^M \cap \sqrt{r} \mathcal{S}^{M-1}} P_n \mathcal{L}_{f_\lambda} \\
&\leq \inf_{0 \leq r \leq r_2} \left(r - \frac{2}{\sqrt{n}} \left(\log(C_1 M r) - c_9 \log \log(C_1 M r_2) \right)^{1/2} \right) + Q(r_2).
\end{aligned} \tag{3.11}$$

Moreover, thanks to (3.7), with probability greater than 10/12,

$$Q(r_1) + Q(r_2) \leq c_{10} \sqrt{r_1 \log(eMr_1)/n}.$$

Fix $0 < \beta_2 < \beta_0$ to be named later and set

$$r_0 = \frac{\beta_0}{\sqrt{n \log(eM/\sqrt{n})}} \quad \text{and} \quad r_2 = \frac{\beta_2}{\sqrt{n \log(eM/\sqrt{n})}}.$$

For β_0 large enough (resp. β_2 small enough), the infimum in (3.10) (resp. (3.11)) is achieved in r_0 (resp. r_2). Therefore, with probability greater than 8/12

$$\begin{aligned} & \inf_{0 \leq r \leq r_0} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} P_n \mathcal{L}_{f_\lambda} - \inf_{r_0 \leq r \leq r_1} \inf_{\lambda \in B_1^M \cap \sqrt{r} S^{M-1}} P_n \mathcal{L}_{f_\lambda} \\ & \leq \frac{c_{11} \log((C_0 r_0)/(C_1 r_2))}{\sqrt{\log(C_1 M r_2)}} + r_2 - r_0 + \frac{c_1}{\sqrt{n \log(c_5 M r_0)}} + c_{12} \sqrt{r_1 \log\left(\frac{eMr_1}{\sqrt{n}}\right)} \\ & \leq \frac{c_{13} \log((C_0 \beta_0)/(C_1 \beta_2))}{\sqrt{n \log(eM/\sqrt{n})}} + \frac{\beta_2 - \beta_0}{\sqrt{n \log(eM/\sqrt{n})}} + \frac{c_{14}}{\sqrt{n \log(eM/\sqrt{n})}} + \frac{c_{14} \log(eM/\sqrt{n})}{n}. \end{aligned}$$

Therefore, there exists some β_0 for which the latter quantity is negative and thus (3.8) holds for $r_0 = \beta_0/\sqrt{n \log(C_0 \beta_0 M/\sqrt{n})}$.

4 Proof of Theorem B

Our starting point is to describe the machinery developed in [4], leading to the desired estimates on the performance of ERM in a general class of functions. Let G be a class of functions and denote by $\mathcal{L}_G = \{(x, y) \mapsto (y - g(x))^2 - (y - g_G^*(x))^2 : g \in G\}$ the associated class of quadratic excess loss functions, where g_G^* is the minimizer of the quadratic risk in G . Let $V = \text{star}(\mathcal{L}_G, 0) = \{\theta \mathcal{L} : 0 \leq \theta \leq 1, \mathcal{L} \in \mathcal{L}_G\}$ and for every $\lambda > 0$ set $V_\lambda = \{h \in V : \mathbb{E}h \leq \lambda\}$.

Theorem 4.1 ([4]) *For every positive B and b there exists a constant $c_0 = c_0(B, b)$ for which the following holds. Let G be a class of functions for which \mathcal{L}_G consists of functions that are bounded by b almost surely. Assume further that for any $\mathcal{L} \in \mathcal{L}_G$, $\mathbb{E}\mathcal{L}^2 \leq B\mathbb{E}\mathcal{L}$. If $x > 0$, $\lambda^* > 0$ satisfies that $\mathbb{E}\|P - P_n\|_{V_{\lambda^*}} \leq \lambda^*/8$ and*

$$\lambda^*(x) = c_0 \max\left(\lambda^*, \frac{x}{n}\right),$$

then with probability greater than $1 - \exp(-x)$, the empirical risk minimization procedure \hat{g} in G satisfies

$$R(\hat{g}) \leq \inf_{g \in G} R(g) + \lambda^*(x).$$

Let F be the given dictionary and set $G = \text{conv}(F)$. Using the notation of Theorem 4.1, put $\mathcal{L}_{\text{conv}(F)} = \{\mathcal{L}_f : f \in \text{conv}(F)\}$, consider the star-shaped hull $V = \text{star}(\mathcal{L}_{\text{conv}(F)}, 0)$ and its localizations $V_\lambda = \{g \in V : \mathbb{E}g \leq \lambda\}$ for any $\lambda > 0$. Thanks to convexity, the following observation holds in our case (see [19] for the proof).

Proposition 4.2 *If $f \in \text{conv}(F)$ then $\mathbb{E}\mathcal{L}_f \geq \|f - f^*\|_{L_2(P^X)}^2$ where f^* is the minimizer of the quadratic risk in $\text{conv}(F)$. In particular,*

1. $\mathbb{E}\mathcal{L}^2 \leq 4b^2\mathbb{E}\mathcal{L}$ for any $\mathcal{L} \in \mathcal{L}_{\text{conv}(F)}$;
2. For $\mu > 0$, if $f \in \text{conv}(F)$ satisfies that $\mathbb{E}\mathcal{L}_f \leq \mu$, then $f \in f^* + K_\mu$, where

$$K_\mu = 2[\text{conv}\{\pm f_1, \dots, \pm f_M\} \cap \sqrt{\mu}\mathcal{B}(L_2(P^X))].$$

The first part of Proposition 4.2 shows that $\mathcal{L}_{\text{conv}(F)}$ satisfies the assumptions of Theorem 4.1 with $B = 4b^2$. To apply Theorem 4.1 one has to find $\lambda^* > 0$ for which $\mathbb{E}\|P - P_n\|_{V_{\lambda^*}} \leq \lambda^*/8$, and to that end we will use the second part of Proposition 4.2. First, observe that it was shown in [5] that

$$\mathbb{E}\|P - P_n\|_{V_\lambda} \leq \sum_{i \geq 0} 2^{-i} \mathbb{E}\|P - P_n\|_{\mathcal{L}_{2^{i+1}\lambda}}, \quad (4.1)$$

where from here on we set $\mathcal{L}_\mu = \{\mathcal{L} \in \mathcal{L}_{\text{conv}(F)} : \mathbb{E}\mathcal{L} \leq \mu\}$. Applying the second part of Proposition 4.2 it is evident that $\{f \in \text{conv}(F) : \mathbb{E}\mathcal{L}_f \leq \mu\} \subset f^* + K_\mu$.

Proof of Theorem B. By the Giné-Zinn symmetrization Theorem [26],

$$\mathbb{E}\|P - P_n\|_{\mathcal{L}_\mu} \leq 2\mathbb{E}\mathbb{E}_\epsilon \sup_{\mathcal{L} \in \mathcal{L}_\mu} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathcal{L}(X_i, Y_i) \right|. \quad (4.2)$$

Note that if $\mathcal{L} \in \mathcal{L}_\mu$ and $f \in \text{conv}(F)$ satisfies that $\mathcal{L} = \mathcal{L}_f$, then for any (x, y) ,

$$\begin{aligned} |\mathcal{L}(x, y)| &= |(y - f(x))^2 - (y - f^*(x))^2| \\ &= |(f^*(x) - f(x))(2y - f(x) - f^*(x))| \leq 4b|f(x) - f^*(x)|. \end{aligned}$$

Thus, by the contraction principle (see, e.g. [18]) and Proposition 4.2,

$$\mathbb{E}\|P - P_n\|_{\mathcal{L}_\mu} \leq \frac{8b}{\sqrt{n}} \mathbb{E}\mathbb{E}_\epsilon \sup_{f \in K_\mu} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right|.$$

Observe that since the dictionary consists of an orthogonal family, if (e_1, \dots, e_M) is the standard basis in ℓ_2^M and $F(\cdot) = (f_1(\cdot), \dots, f_M(\cdot))$, then

$$K_\mu = \{2\langle \lambda, F \rangle : \lambda \in B_1^M \cap \sqrt{\mu}\mathcal{E}\},$$

where \mathcal{E} is an ellipsoid with principal axes $(\|f_i\|_{L_2} e_i)_{i=1}^M$. From here on we will assume that $(\|f_i\|_{L_2})_{i=1}^M$ is a non-increasing sequence.

Now, we want to bound

$$\begin{aligned} \mathbb{E} \sup_{f \in K_\mu} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| &= \mathbb{E} \sup_{\lambda \in B_1^M \cap \sqrt{\mu} \mathcal{E}} \left| \frac{2}{n} \sum_{i=1}^n \epsilon_i \langle \lambda, F(X_i) \rangle \right| \\ &= \frac{2}{\sqrt{n}} \mathbb{E} \left\| \sum_{i=1}^n \frac{1}{\sqrt{n}} \epsilon_i F(X_i) \right\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ}, \end{aligned}$$

where $\|\cdot\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ}$ denotes the dual norm to the one whose unit ball is $B_1^M \cap \sqrt{\mu} \mathcal{E}$. We will use two different strategies to bound this process depending on $M \leq \sqrt{n}$ or $M > \sqrt{n}$. First start with the case $M \geq \sqrt{n}$. Since both B_1^M and \mathcal{E} are unconditional with respect to the coordinate structure given by $(e_i)_{i=1}^M$, it follows that

$$\|v\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ} \sim \inf_{I \subset \{1, \dots, M\}} \left[\sqrt{\mu} \left(\sum_{i \in I} \left(\frac{v_i}{\|f_i\|_{L_2}} \right)^2 \right)^{1/2} + \max_{i \in I^c} |v_i| \right], \quad (4.3)$$

and in our case, $v = (v_j)_{j=1}^M = ((1/\sqrt{n}) \cdot \sum_{i=1}^n \epsilon_i f_j(X_i))_{j=1}^M$.

Let

$$J_0 = \{j : \|f_j\|_{L_2} \geq c_0 b \sqrt{\log M} / \sqrt{n}\},$$

where c_0 is a constant to be named later. A straightforward application of Bernstein inequality [26] shows that, for $t \geq c_1$,

$$\begin{aligned} \mathbb{P}(\exists j \in J_0 : P_n f_j^2 \geq (t+1) \|f_j\|_{L_2}^2) &\leq \sum_{j \in J_0} \exp(-c_2 n (\|f_j\|_{L_2}^2 / b^2) \min(t^2, t)) \\ &\leq M \exp(-c_3 t \log M) \leq \exp(-c_4 t \log M), \end{aligned}$$

and

$$\mathbb{P}(\exists j \in J_0^c : P_n f_j^2 \geq (t+1) b^2 n^{-1} \log M) \leq \exp(-c_4 t \log M).$$

For every integer $\ell \geq c_1$, let

$$\mathcal{A}_\ell = \{\forall j \in J_0 : P_n f_j^2 \leq (\ell+1) \|f_j\|_{L_2}^2\} \cap \{\forall j \in J_0^c : P_n f_j^2 \leq (\ell+1) b^2 n^{-1} \log M\}.$$

Set $\mathcal{B}_\ell = \mathcal{A}_{\ell+1} \cap \mathcal{A}_\ell^c$ and note that $\mathbb{P}(\mathcal{B}_\ell) \leq \mathbb{P}(\mathcal{A}_\ell^c) \leq 2 \exp(-c_4 \ell \log M)$ for any $\ell \geq c_1$.

For every $\ell \geq c_1$, consider the random variables conditioned on \mathcal{B}_ℓ ,

$$U_{j,\ell} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f_j(X_i) / \|f_j\|_{L_2} \Big| \mathcal{B}_\ell \quad \forall j \in J_0$$

and

$$U_{j,\ell} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f_j(X_i) | \mathcal{B}_\ell \quad \forall j \in J_0^c.$$

Hence, by Hoeffding's inequality (cf. [26]), there exists an absolute constant c_5 such that, for any $j \in J_0$,

$$\|U_{j,\ell}\|_{\psi_2(\epsilon)}^2 \leq c_5 \frac{n^{-1} \sum_{i=1}^n f_j^2(X_i)}{\|f_j\|_{L_2}^2} \leq c_5(\ell + 1),$$

and for any $j \in J_0^c$,

$$\|U_{j,\ell}\|_{\psi_2(\epsilon)}^2 \leq c_5(\ell + 1)b^2(\log M)/n.$$

By a result due to Klartag [13], it follows that for every such ℓ and any $1 \leq j \leq |J_0|$,

$$\mathbb{E}_\epsilon \left(\sum_{i=1}^j (U_{i,\ell}^*)^2 \right)^{1/2} \leq c_6 \sqrt{\ell} \sqrt{j \log(e|J_0|/j)},$$

where $(U_{j,\ell}^*)_{j=1}^{|J_0|}$ is a decreasing rearrangement of $(|U_{j,\ell}|)_{j \in J_0}$. Moreover, by a standard maximal inequality (see, e.g. [26])

$$\mathbb{E}_\epsilon \max_{j \in J_0^c} U_{j,\ell} \leq c_7 \sqrt{\log |J_0^c|} \max_{j \in J_0^c} \|U_{j,\ell}\|_{\psi_2(\epsilon)} \leq c_8 \sqrt{\ell} b \frac{\log M}{\sqrt{n}}.$$

For every $1 \leq j \leq |J_0|$, let I be the set of the j largest coordinates of $(|U_{j,\ell}|)_{j \in J_0}$. Hence, by (4.3) and since $\|f_j\|_{L_2} \leq b$,

$$\begin{aligned} & \mathbb{E}_\epsilon \left(\left\| \sum_{i=1}^n \frac{1}{\sqrt{n}} \epsilon_i F(X_i) \right\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ} | (X_i)_{i=1}^n \in \mathcal{B}_\ell \right) \\ & \lesssim \mathbb{E}_\epsilon \sqrt{\mu} \left(\sum_{i=1}^j (U_{i,\ell}^*)^2 \right)^{1/2} + \mathbb{E}_\epsilon \max \left(\|f_j\|_{L_2} U_{j,\ell}^*, \max_{j \in J_0^c} |U_{j,\ell}| \right) \\ & \lesssim \sqrt{\ell} \left(\sqrt{\mu} \sqrt{j \log(e|J_0|/j)} + b \sqrt{\log(e|J_0|/j)} \right) + \mathbb{E}_\epsilon \max_{j \in J_0^c} |U_{j,\ell}| \\ & \leq \sqrt{\ell} \left(\sqrt{\mu} \sqrt{j \log(e|J_0|/j)} + b \sqrt{\log(e|J_0|/j)} \right) + \sqrt{\ell} b \frac{\log M}{\sqrt{n}}. \end{aligned}$$

Therefore, if we take $j = \min\{\lceil 1/\mu \rceil, |J_0|\}$ it is evident that

$$\mathbb{E}_\epsilon \left(\left\| \sum_{i=1}^n \frac{1}{\sqrt{n}} \epsilon_i F(X_i) \right\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ} | (X_i)_{i=1}^n \in \mathcal{B}_\ell \right) \lesssim b \sqrt{\ell} \left(\sqrt{\log(eM\mu)} + \frac{\log M}{\sqrt{n}} \right).$$

Thus, integration with respect to X_1, \dots, X_n and applying the estimates on the measure of \mathcal{B}_ℓ ,

$$\frac{2}{\sqrt{n}} \mathbb{E} \left\| \sum_{i=1}^n \frac{1}{\sqrt{n}} \epsilon_i F(X_i) \right\|_{(B_1^M \cap \sqrt{\mu} \mathcal{E})^\circ} \lesssim b \sqrt{\frac{\log(eM\mu)}{n}}.$$

Finally, by (4.1), for any $\lambda > 1/M$,

$$\begin{aligned} \mathbb{E} \|P - P_n\|_{V_\lambda} &\leq \sum_{i \geq 0} 2^{-i} \mathbb{E} \|P - P_n\|_{\mathcal{L}_{2^{i+1}\lambda}} \\ &\lesssim b^2 \sum_{i \geq 0} 2^{-i} \sqrt{\frac{\log(eM2^{i+1}\lambda)}{n}} \lesssim b^2 \sqrt{\frac{\log(eM\lambda)}{n}}, \end{aligned}$$

and, if

$$\lambda^* \sim b^2 \sqrt{\frac{1}{n} \log\left(\frac{eMb^2}{\sqrt{n}}\right)},$$

then $\mathbb{E} \|P - P_n\|_{V_{\lambda^*}} \leq \lambda^*/8$, as required.

When $M \leq \sqrt{n}$, we use the strategy developed in [15]. Let S be the linear subspace of $L^2(P^X)$ spanned by F and take $(e_1, \dots, e_{M'})$ to be an orthonormal basis of S (where $M' = \dim(S) \leq M$). Since $K_\mu \subset S \cap 2\sqrt{\mu}\mathcal{B}(L_2(P^X))$, then

$$\begin{aligned} \mathbb{E} \sup_{f \in K_\mu} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right| &\leq 2 \mathbb{E} \sup_{\|\lambda\|_{\ell_2^{M'}} \leq 2\sqrt{\lambda}} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \left(\sum_{j=1}^{M'} \lambda_j e_j(X_i) \right) \right| \\ &\lesssim \sqrt{\mu} \mathbb{E} \left(\sum_{j=1}^{M'} \left(\frac{1}{n} \sum_{i=1}^n \epsilon_i e_j(X_i) \right)^2 \right)^{1/2} \lesssim \sqrt{\frac{M'\mu}{n}}. \end{aligned}$$

The rate obtained in the case $M \leq \sqrt{n}$ follows now from (4.1). ■

Appendix

We establish the following upper bound on the risk of \tilde{f}^{ERM-C} as a (C)-aggregation procedure in the general case. Its proof follows the same path as in Section 4. But, rather than studying the empirical process indexed by the interpolation body $B_1^M \cap \sqrt{\mu} \mathcal{E}$, in the case $M \geq \sqrt{n}$, one simply uses the approximation $B_1^M \cap \sqrt{\mu} \mathcal{E} \subset B_1^M$ to get, conditionally on X_1, \dots, X_n ,

$$\mathbb{E}_\epsilon \sup_{f \in K_\mu} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i f(X_i) \right| \leq \mathbb{E}_\epsilon \sup_{\lambda \in B_1^M} \left| \left\langle \lambda, \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i F(X_i) \right\rangle \right| = \mathbb{E}_\epsilon \max_{1 \leq j \leq M} |\gamma_j|,$$

where, for all $j = 1, \dots, M$, γ_j is the subgaussian random variable $n^{-1/2} \sum_{i=1}^n \epsilon_i f_j(X_i)$ with ψ_2 -norm bounded by $n^{-1} \sum_{i=1}^n f_j(X_i)^2 \leq c_0 b^2$ and thus by a maximal inequality [18],

$$\mathbb{E}_g \max_{1 \leq j \leq M} |\gamma_j| \leq c_1 b \sqrt{\log M}.$$

The result below follows from this upper bound and (4.1) for the case $M > \sqrt{n}$, and the case $M \leq \sqrt{n}$ follows the same path as the proof of Theorem B, and thus its proof is omitted.

Theorem 4.3 *For every $b > 0$ there is a constant $c_1(b)$ and an absolute constant c_2 for which the following holds. Let n and M be integers which satisfy that $\log M \leq c_1(b)\sqrt{n}$. For any couple (X, Y) and any finite dictionary F of cardinality M such that $|Y|, \sup_{f \in F} |f(X)| \leq b$, and for any $u > 0$, with probability greater than $1 - \exp(-u)$,*

$$R(\tilde{f}^{ERM-C}) \leq \min_{f \in \text{conv}(F)} R(f) + c_2 b^2 \max \left[\min \left(\frac{M}{n}, \sqrt{\frac{\log M}{n}} \right), \frac{u}{n} \right].$$

References

- [1] Jean-Yves Audibert. Aggregated estimators and empirical complexity for least square regression. *Ann. Inst. H. Poincaré Probab. Statist.*, 40(6):685–736, 2004.
- [2] Jean-Yves Audibert. Progressive mixture rules are deviation suboptimal. *Advances in Neural Information Processing Systems (NIPS)*, 2007.
- [3] Jean-Yves Audibert. Fast learning rates in statistical inference through aggregation. *Ann. Statist.*, 37(4):1591–1646, 2009.
- [4] Peter L. Bartlett and Shahar Mendelson. Empirical minimization. *Probab. Theory Related Fields*, 135(3):311–334, 2006.
- [5] Peter L. Bartlett, Shahar Mendelson, and Joseph Neeman. ℓ_1 -regularized linear regression: Persistence and oracle inequalities. *To appear in Probab. Theory Related Fields*, 2011.
- [6] Lucien Birgé. Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Ann. Inst. H. Poincaré Probab. Statist.*, 42(3):273–325, 2006.
- [7] Olivier Bousquet, Vladimir Koltchinskii, and Dmitriy Panchenko. Some local measures of complexity of convex hulls and generalization bounds. In *Computational learning theory (Sydney, 2002)*, volume 2375 of *Lecture Notes in Comput. Sci.*, pages 59–73. Springer, Berlin, 2002.

- [8] Florentina Bunea and Andrew Nobel. Sequential procedures for aggregating arbitrary estimators of a conditional mean. *IEEE Trans. Inform. Theory*, 54(4):1725–1735, 2008.
- [9] Florentina Bunea, Alexandre B. Tsybakov, and Marten H. Wegkamp. Aggregation for Gaussian regression. *Ann. Statist.*, 35(4):1674–1697, 2007.
- [10] Arnak S. Dalalyan and Alexandre B. Tsybakov. Aggregation by exponential weighting and sharp oracle inequalities. In *Learning theory*, volume 4539 of *Lecture Notes in Comput. Sci.*, pages 97–111. Springer, Berlin, 2007.
- [11] M. Emery, A. Nemirovski, and D. Voiculescu. *Lectures on probability theory and statistics*, volume 1738 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2000. Lectures from the 28th Summer School on Probability Theory held in Saint-Flour, August 17–September 3, 1998, Edited by Pierre Bernard.
- [12] Anatoli Juditsky and Arkadii Nemirovski. Functional aggregation for nonparametric regression. *Ann. Statist.*, 28(3):681–712, 2000.
- [13] Bo’az Klartag. $5n$ Minkowski symmetrizations suffice to arrive at an approximate Euclidean ball. *Ann. of Math. (2)*, 156(3):947–960, 2002.
- [14] V. Koltchinskii. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: Ecole d’été de Probabilités de Saint-FlourXXXVIII-2008*. Lecture Notes in Mathematics / école d’été de Probabilités de Saint-Flour Series. Springer, 2011.
- [15] Vladimir Koltchinskii. Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.*, 34(6):2593–2656, 2006.
- [16] Guillaume Lecué and Shahar Mendelson. Aggregation via empirical risk minimization. *Probab. Theory Related Fields*, 145(3-4):591–613, 2009.
- [17] Guillaume Lecué and Shahar Mendelson. General non-exact oracle inequalities in the unbounded case. *Submitted*, 2010.
- [18] Michel Ledoux and Michel Talagrand. *Probability in Banach spaces*, volume 23 of *Ergebnisse der Mathematik und ihrer Grenzgebiete (3) [Results in Mathematics and Related Areas (3)]*. Springer-Verlag, Berlin, 1991. Isoperimetry and processes.
- [19] Shahar Mendelson. Obtaining fast error rates in nonconvex situations. *J. Complexity*, 24(3):380–397, 2008.

- [20] Shahar Mendelson. Empirical processes with a bounded ψ_1 diameter. *Geom. Funct. Anal.*, 20(4):988–1027, 2010.
- [21] Shahar Mendelson, Alain Pajor, and Nicole Tomczak-Jaegermann. Reconstruction and subgaussian operators in asymptotic geometric analysis. *Geom. Funct. Anal.*, 17(4):1248–1282, 2007.
- [22] Valentin V. Petrov. *Limit theorems of probability theory*, volume 4 of *Oxford Studies in Probability*. The Clarendon Press Oxford University Press, New York, 1995.
- [23] Gilles Pisier. *The volume of convex bodies and Banach space geometry*, volume 94 of *Cambridge Tracts in Mathematics*. Cambridge University Press, Cambridge, 1989.
- [24] Alexandre B. Tsybakov. Optimal rate of aggregation. In *Computational Learning Theory and Kernel Machines (COLT-2003)*, volume 2777 of *Lecture Notes in Artificial Intelligence*, pages 303–313. Springer, Heidelberg, 2003.
- [25] Alexandre B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Ann. Statist.*, 32(1):135–166, 2004.
- [26] Aad W. van der Vaart and Jon A. Wellner. *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. With applications to statistics.
- [27] Yuhong Yang. Mixing strategies for density estimation. *Ann. Statist.*, 28(1):75–87, 2000.
- [28] Yuhong Yang. Aggregating regression procedures to improve performance. *Bernoulli*, 10(1):25–47, 2004.