

# Lower bounds for the empirical minimization algorithm

Shahar Mendelson\*

February 15, 2008

## Abstract

We present a simple argument that proves that under mild geometric assumptions on the class  $F$  and the set of target functions  $\mathcal{T}$ , the empirical minimization algorithm can not yield a uniform error rate that is faster than  $1/\sqrt{k}$  in the function learning set-up. This result holds for various loss functionals and the target functions from  $\mathcal{T}$  that cause the slow uniform error rate are clearly exhibited.

## 1 Introduction

The aim of this note is to present a relatively simple proof of a lower bound on the error rate of the empirical minimization algorithm in *function learning* problems.

Let us describe the question at hand. Let  $F$  be a class of functions on the probability space  $(\Omega, \mu)$  and consider an unknown target function  $T$  that one wishes to approximate in the following sense. The learner is given a random sample  $(X_i)_{i=1}^k, (T(X_i))_{i=1}^k$ , where  $X_1, \dots, X_k$  are independent points selected according to (the unknown) probability measure  $\mu$ . The goal of the learner is to use this data to find a function  $f \in F$  that approximates  $T$  with respect to some loss-function  $\ell$ . In other words, to find  $f \in F$  such that the expected loss  $\mathbb{E}\ell(f(X), T(X))$  is close to the best possible in the

---

\*Centre for Mathematics and its Applications, Institute of Advanced Studies, The Australian National University, Canberra, ACT 0200, Australia, and Department of Mathematics, Technion, I.I.T, Haifa 32000, Israel.

email: shahar.mendelson@anu.edu.au

Supported in part by the Israel Science Foundation grant ISF 666/06 and the Australian Research Council Discovery grant DP0559465.

I would like to thank Or Neviot for many stimulating discussions.

class. A typical choice of a loss function  $\ell$  is the squared loss  $|x - y|^2$ , or more generally, the  $p$ -loss  $|x - y|^p$  for  $1 \leq p < \infty$ . There are, of course, many other choices of  $\ell$  that are used.

The function learning problem has a more general counterpart, the *agnostic learning* problem, in which the unknown target function  $T$  is replaced by a random variable  $Y$ . The data received by the learner is an i.i.d sample  $(X_i, Y_i)_{i=1}^k$  given according to the joint probability distribution of  $X$  and  $Y$ . Again, the goal is to find some  $f \in F$  for which  $\mathbb{E}\ell(f(X), Y)$  is as small as possible.

An algorithm frequently used in such prediction problems is *empirical minimization*. For every sample, the algorithm produces a function  $\hat{f} \in F$  that minimizes the empirical loss  $\sum_{i=1}^k \ell(f(X_i), T(X_i))$  ( $\sum_{i=1}^k \ell(f(X_i), Y_i)$  in the agnostic case). The hope is to obtain a high probability estimate on the way the risk of  $\hat{f}$ , defined as the conditional expectation

$$\mathbb{E} \left( \ell(\hat{f}, T) | X_1, \dots, X_k \right) - \inf_{f \in F} \mathbb{E}\ell(f, T),$$

decreases as a function of the sample size  $k$ . The expectation of this quantity is usually called the *error rate* of the problem and measures “how far” the algorithm is from choosing the best function in the class.

If  $\mathbb{E}\ell(T, \cdot)$  has a unique minimizer in  $F$ , which we denote by  $f^*$ , one can define the excess loss  $\mathcal{L}(f) = \ell(f, T) - \ell(f^*, T)$ . Then the excess loss (risk) of the empirical minimizer is

$$\mathbb{E} \left( \mathcal{L}(\hat{f}) | X_1, \dots, X_k \right).$$

There are numerous results on the performance of the empirical minimization algorithm but we shall not present any sort of survey of those results here. Roughly speaking, it turns out that the “richness” of the function class  $F$  as captured by the empirical process indexed by it, determines how well the empirical minimization behaves in the two learning problems we mentioned above. And, in order to obtain an error rate that tends to 0 asymptotically faster than  $1/\sqrt{k}$ , not only does the class have to be small, but additional information on the loss is needed; for example, a Bernstein type condition, that for every  $f \in F$ ,

$$\mathbb{E}\mathcal{L}^2(f) \leq c(\mathbb{E}\mathcal{L}(f))^\beta \tag{1.1}$$

for some constants  $c$  and  $0 < \beta \leq 1$ , would do.

For more details on richness parameters of classes and the way in which those, combined with conditions similar to (1.1) govern the error rate, we

refer the reader to [16, 2, 3] and to the surveys [7, 1, 20, 5, 14, 17, 4]. For general facts concerning empirical processes see, for example, [21, 9].

Our main interest is in a lower bound on the performance of the empirical minimization algorithm in function learning problems. There are many lower bounds that are independent of the algorithm used (see, for example [12] and references therein), but usually, these lower bounds deal with specific classes and are very different in nature from what we have in mind. The starting point of our discussion is the surprising result established in [15]: a lower bound on the performance of any learning algorithm in the agnostic case.

To formulate this result we need the following definitions. We say that  $A$  is a learning algorithm if for every integer  $k$  and any sample  $s = (X_i, T(X_i))_{i=1}^k$  (resp.  $(X_i, Y_i)_{i=1}^k$  in the agnostic case) it assigns a function  $A_s \in F$ . For a random variable  $Y$  we denote by  $\nu$  the joint probability measure endowed by  $(X, Y)$ .

**Theorem 1.1** [15] *Let  $F \subset L_2(\mu)$  be a compact class of functions bounded by 1 and set  $\ell(x, y) = (x - y)^2$ . Assume that there is a random variable  $Y$  bounded by  $\beta$  for which  $\mathbb{E}(f(X) - Y)^2$  has more than a unique minimizer in  $F$ . Then, there are constants  $c$  and  $k_0$  depending only on  $F$ ,  $\beta$  and  $\mu$  for which the following holds. If  $A$  is a learning algorithm and  $\mathcal{Y} = \{Y : \|Y\|_\infty \leq \beta\}$  then for every  $k \geq k_0$ ,*

$$\sup_Y \left( \mathbb{E} \left( \mathbb{E} \left( \ell(\hat{A}, Y) | s_k \right) \right) - \inf_{f \in F} \mathbb{E} \ell(f, Y) \right) \geq \frac{c}{\sqrt{k}},$$

where the supremum is with respect to all random variables  $Y$  taking values in  $\mathcal{Y}$ ,  $s_k = (X_i, Y_i)_{i=1}^k$  is an i.i.d sample according to the joint distribution of  $(X, Y)$  and  $\hat{A} = A_{s_k}$ .

In other words, there will be a range  $[-\beta, \beta]$  for which no matter what learning algorithm is chosen by the learner to approximate targets taking values in that range, the best uniform error rate it can guarantee with respect to all these targets can not be asymptotically better than  $1/\sqrt{k}$ .

It is important to mention that this is not the exact formulation of the result from [15]. In the original formulation, it was assumed that  $F$  is compact and nonconvex. A part of the proof was to show that under these assumptions there is some  $\beta$  and a random variable  $Y$  bounded by  $\beta$  for which  $\mathbb{E}(f(X) - Y)^2$  has multiple minimizers in  $F$ , in the following sense: there are  $f_1, f_2 \in F$  such that  $(f_1(X) - Y)^2 \neq (f_2(X) - Y)^2$  on a set of positive measure and  $f_1, f_2$  minimizer  $\mathbb{E}(f(X) - Y)^2$  in  $F$ .

Unfortunately, as will be explained in Section 3.2, the proof of that part of the claim from [15] is incorrect and the actual result proved there is Theorem 1.1.

The surprising point in Theorem 1.1 is that the lower bound of  $1/\sqrt{k}$  does not depend on the richness of the class. The slow rate is the best possible uniform rate regardless of how “statistically small” the class  $F$  is. Somehow, the bad geometry of  $F$  is the reason for the slow rate and understanding the geometric reasons causing the bad rates is one of our goals.

A rather delicate point in Theorem 1.1 is that it does not imply that there is a single random variable  $Y$  taking values in  $\mathcal{Y}$  for which the error rate is bounded from below by  $c/\sqrt{k}$ . What it does say is that for any fixed algorithm and any integer  $k \geq k_0$  there will be some random variable  $Y$  - depending on  $F, \mu, k$  and the chosen algorithm, on which the algorithm performs poorly after being given  $k$  data points - but even for a fixed algorithm the “bad target”  $Y$  might change with  $k$ .

Since the proof in [15] is based on a “probabilistic method” type of construction it is indirect. For each  $k$ , the fact that a “bad”  $Y_k$  exists is exhibited (based on the existence of a target  $Y$  with multiple minimizers), but what  $Y_k$  is and how it is related to the geometry of  $F$  is not revealed by the proof. Of course, this is the best that could be expected in such a general solution since the algorithm used is not specified and can be arbitrary. The proof also uses the fact that one is allowed to select an arbitrary random variable  $Y$  as a target rather than a fixed target function  $T(X)$ . Thus, the “agnostic” argument from [15] does not extend to the function learning setup.

We present a simple argument that proves the same lower bound in the function learning scenario for the empirical minimization algorithm. The argument requires minor assumptions on the loss functional, rather than assuming that  $\ell$  is the squared loss. Moreover, it enables one to pin-point a “bad” target for every sample size  $k$ . Our feeling is that this proof sheds some light on the reasons why the bad geometry of  $F$  leads to poor statistical properties.

Our starting point is similar to [15] (though the proofs take very different paths). Assume that  $E$  is a reasonable normed space of functions on  $(\Omega, \mu)$  with a norm that is naturally connected to the loss (for example, the  $p$ -loss is connected to the  $L_p$  norm). Assume further that  $F \subset E$  is “small” in an appropriate sense and that  $T$  has more than a unique best approximation in  $F$ . Fix one such best approximation  $f_* \in F$ . We will show that for every  $k \geq k_0$  and  $\lambda \sim 1/\sqrt{k}$ , the function  $(1 - \lambda)T + \lambda f_*$  is a “bad” target function

for a typical  $k$ -sample, that is, for every  $k \geq k_0$

$$\mathbb{E}_{X_1, \dots, X_k} \left( \mathbb{E} \left( \ell(\hat{f}, T_{\lambda_k}) | X_1, \dots, X_k \right) - \inf_{f \in F} \mathbb{E} \ell(f, T_{\lambda_k}) \right) \geq \frac{c}{\sqrt{k}},$$

where  $\hat{f}$  is the empirical minimizer and  $c$  is a constant that depends only on  $F$ ,  $\ell$  and properties of the space  $E$ .

A corollary of this general result is Theorem 1.2 formulated below. Recall that  $\mu$ -Donsker classes are sets  $F \subset L_2(\mu)$  that satisfy some kind of a uniform Central Limit Theorem (see [9, 21] for detailed surveys on this topic).

Let  $E$  be a normed space of functions on  $(\Omega, \mu)$  and let  $N(F, E)$  be the set of functions in  $E$  that have more than a unique best approximation in  $F$ .

**Theorem 1.2** *Let  $2 \leq p < \infty$  and set  $E = L_p(\mu)$ . Assume that  $F \subset E$  is a  $\mu$ -Donsker class of functions bounded by 1, let  $R > 0$  and assume that  $\mathcal{T} \subset E \cap B_{L_\infty}(0, R)$  is convex and contains  $F$ . If  $\ell$  is the  $p$ -loss function and  $\mathcal{T} \cap N(F, E) \neq \emptyset$ , then for  $k \geq k_0$*

$$\sup_{T \in \mathcal{T}} \left( \mathbb{E}_{X_1, \dots, X_k} \mathbb{E} \left( \ell(\hat{f}, T) | X_1, \dots, X_k \right) - \inf_{f \in F} \mathbb{E} \ell(f, T) \right) \geq \frac{c}{\sqrt{k}},$$

where  $c$  and  $k_0$  depend only on  $p$ ,  $F$  and  $R$ .

Let us note that the assumption that  $F$  and  $\mathcal{T}$  are bounded in  $L_\infty$  is only there to ensure that Lipschitz images of these functions satisfy some technical integrability properties we require and is not essential for the proof. Also, the convexity assumption on  $\mathcal{T}$  is only there to ensure that if  $T \in N(F, E)$  and  $f_* \in F$  then for any  $\lambda \in [0, 1]$  the convex combination  $(1 - \lambda)T + \lambda f_* \in \mathcal{T}$ , and thus, a “legal” target function.

It turns out that the reverse direction of Theorem 1.2 is also true [18]. Indeed, one can show that if the set  $\mathcal{T}$  is “far away” from  $N(F, E)$  then the class  $F$  satisfies a Bernstein condition. Thus, if  $F$  is “small”, the uniform error rate with respect to functions in  $\mathcal{T}$  decays faster than  $1/\sqrt{k}$ .

To formulate this reverse direction we need the following definition.

**Definition 1.3** *Let  $E$  be a Banach space, set  $F \subset E$  to be compact and assume that  $T \notin \overline{N(F, E)}$ . If  $f^*$  is the best approximation of  $T$  in  $F$ , let*

$$\lambda^*(T) = \sup \left\{ \lambda \geq 1 : \lambda f^* + (1 - \lambda)T \notin \overline{N(F, E)} \right\}.$$

In other words, if  $T \notin \overline{N(F, E)}$  and if one considers the ray originating in  $f^*$  that passes through  $T$ ,  $\lambda^*(T)$  measures how far “up” this ray one can move while still remaining at positive a distance from the set  $N(F, E)$ . Clearly, if  $d(T, N(F, E)) > 0$  then  $\lambda^*(T) > 1$ .

**Theorem 1.4** [18] *Let  $1 < p < \infty$  and set  $F \subset L_p(\mu)$  to be a compact set of functions that are bounded by 1. Let  $T$  be a function bounded by 1 for which  $d(T, N(F, E)) > 0$ . Then, for every  $f \in F$ ,*

$$\mathbb{E}\mathcal{L}_f^2 \leq B(\mathbb{E}\mathcal{L}_f)^{\alpha_p},$$

where  $\mathcal{L}_f = |f - T|^p - |f^* - T|^p$  is the  $p$ -excess loss associated with  $f$  and  $T$ ,  $\alpha_p = \min\{p/2, 2/p\}$  and

$$B = c(p) \inf_{1 < \lambda < \lambda^*(T)} \frac{\lambda}{\lambda - 1}.$$

In particular, for  $p = 2$  the combination of Theorem 1.2 and Theorem 1.4 gives an almost characterization of the uniform error rate associated with a set of targets  $\mathcal{T}$ . Indeed, if a target is “far away” from the set  $N(F, E)$ , the loss excess loss class satisfies a Bernstein condition, implying that the error rate of the empirical minimizer depends only on the statistical complexity of  $F$  and not on its geometry. In particular, if  $d(\mathcal{T}, N(F, E)) > 0$  and  $F$  is small enough, one has very fast error rates - uniformly in  $\mathcal{T}$ . On the other hand,  $\mathcal{T}$  is closed and convex and  $d(\mathcal{T}, N(F, E)) = 0$ , there is some  $T \in \mathcal{T} \cap N(F, E)$ . Hence, by Theorem 1.2 the best possible error rate is  $\sim 1/\sqrt{k}$ .

Finally, a word about notation. Throughout, all constant will be denoted by  $c, c_1$ , etc. Their vales may change from line to line. We will also emphasize when a constant is absolute (that is, a fixed positive number) and when it depends on other parameters of the problem (for example, the diameter of the set  $F$  with respect to the norm). Constants that will remain fixed throughout this article will be denoted by  $C_1, C_2$ , etc.

Let  $P_k g = k^{-1} \sum_{i=1}^k g(X_i)$  where  $X_1, \dots, X_k$  are independent, distributed according to  $\mu$  and set  $\|P_k - P\|_G$  to be the supremum of the empirical process indexed by  $G$ , that is,  $\sup_{g \in G} |k^{-1} \sum_{i=1}^k g(X_i) - \mathbb{E}g|$ . If  $E$  is a normed space, let  $B(x, r)$  be the closed ball centered at  $x$  and of radius  $r$  and set  $B_E$  to be the closed unit ball.

## 2 Preliminaries

Let  $(E, \|\cdot\|)$  be a normed space of functions on the probability space  $(\Omega, \mu)$ . We need to assume that  $E$  has a nice structure, namely that it is smooth and uniformly convex (defined below). For more information on these geometric notions we refer the reader to [8, 6, 10].

We say that a normed space is smooth if the norm is Gâteaux differentiable in any  $x \neq 0$ . There is an equivalent geometric formulation of this notion, that for every  $x$  on the unit sphere of  $E$  there is a unique, norm one linear functional that supports the unit ball in  $x$  (that is, a unique functional  $x^*$ , such that  $\|x^*\| = 1$  and  $x^*(x) = 1$ ).

Uniform convexity measures how far “inside” the unit ball  $(x + y)/2$  is, where  $x$  and  $y$  are distant points on the unit sphere.

**Definition 2.1**  *$E$  is called uniformly convex if there is a positive function  $\delta_E(\varepsilon)$  satisfying that for every  $0 < \varepsilon < 2$  and every  $x, y \in B_E$  for which  $\|x - y\| \geq \varepsilon$ ,  $\|x + y\| \leq 2 - 2\delta(\varepsilon)$ . In other words,*

$$\delta_E(\varepsilon) = \inf \left\{ 1 - \frac{1}{2}\|x + y\| : \|x\|, \|y\| \leq 1, \|x - y\| \geq \varepsilon \right\}.$$

*The function  $\delta(\varepsilon)$  is called the modulus of convexity of  $E$ .*

A fact we shall use later is that  $\delta_E(\varepsilon)$  is an increasing function of  $\varepsilon$  and that if  $0 < \varepsilon_1 \leq \varepsilon_2 \leq 2$  then  $\delta_E(\varepsilon_1)/\varepsilon_1 \leq \delta_E(\varepsilon_2)/\varepsilon_2$  (see. e.g., [8], Chapter 8).

Next, let us turn to some of the properties of the learning problem we study. Consider the sets  $F$  and  $\mathcal{T}$  that are closed subsets of  $E$ . The aim of the learner is to approximate an unknown target function  $T \in \mathcal{T}$  using functions from  $F$ , and the notion of approximation is via the loss functional  $\ell(x, y)$ . The assumptions on the loss  $\ell$  are that it is a Lipschitz function from  $\mathbb{R}^2$  to  $\mathbb{R}$  and that the expected loss is compatible with the norm.

**Assumption 2.1** *Assume that  $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$  is a nonnegative Lipschitz function with constant  $\|\ell\|_{\text{lip}}$ . Assume further that there is some function  $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  that is differentiable, strictly increasing and convex, such that for every  $f, g \in E$ ,  $\mathbb{E}\ell(f, g) = \phi(\|f - g\|)$ .*

This assumption is natural in the context of function learning. For example, the  $p$  loss function  $\ell_p(x, y) = |x - y|^p$  satisfies Assumption 2.1 when  $E = L_p(\mu)$  and  $\phi(t) = t^p$ .

In our construction we consider the excess loss functional that measures how far  $f$  is from being the best in the class. To that end, set  $F_{*,T} = \{f \in$

$F : \mathbb{E}\ell(f, T) = \inf_{g \in F} \mathbb{E}\ell(g, T)\}$ . It is well known (see, e.g. [15, 16, 17]) that under various assumptions on the class and on the loss functional,  $F_{*,T}$  is a singleton, which we denote by  $f^*$ . In such a case one can define the excess loss functional

$$\mathcal{L}_T(f) = \ell(f, T) - \ell(f^*, T),$$

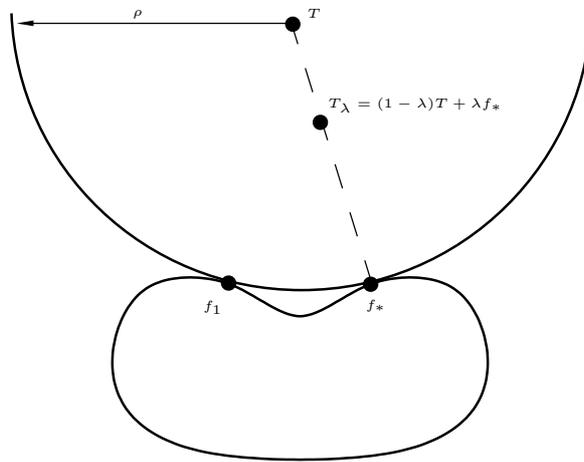
and the excess loss class  $\mathcal{L}_T(F) = \{\mathcal{L}_T(f) : f \in F\}$ . If  $F_{*,T}$  contains more than one element the excess loss class is not well defined, but we will not have to tackle this issue here.

## 2.1 An outline of the proof

Let  $\mathcal{T}$  be a convex set containing  $F$  and put  $N(F, \ell)$  to be the set of functions  $T \in E$  for which  $\mathbb{E}\ell(T, \cdot)$  has multiple minimizers in  $\{\ell(f, T) : f \in F\}$ . Assume that  $\mathcal{T} \cap N(F, \ell) \neq \emptyset$  and let  $T \in \mathcal{T} \cap N(F, \ell)$ . Thus, there is a set  $F_{*,T} \subset F$  of cardinality strictly larger than 1 with the following properties:

1. If  $f_1, f_2 \in F_{*,T}$  and  $f_1 \neq f_2$  then  $\ell(f_1, T) \neq \ell(f_2, T)$  on a set of positive  $\mu$ -probability.
2. For every  $g \in F_{*,T}$ ,  $\mathbb{E}\ell(g, T) = \min_{f \in F} \mathbb{E}\ell(f, T) \equiv R$ .

Therefore, as  $\mathbb{E}\ell(f, T) = \phi(\|f - T\|)$  and since  $\phi$  is strictly increasing, then  $F_{*,T} \subset F \cap B(T, \rho)$  where  $\rho = \phi^{-1}(R)$ . Fix any  $f_* \in F_{*,T}$  and for every  $\lambda \in (0, 1]$  define  $T_\lambda = (1 - \lambda)T + \lambda f_*$ .



Our construction has three components. First of all, one has to show that  $T_\lambda \notin N(F, \ell)$ , and that the unique minimizer of  $\mathbb{E}\ell(T_\lambda, \cdot)$  is  $f_*$ . This

follows from the fact that  $T_\lambda$  has a unique nearest point in  $F$  with respect to the norm. In particular, the excess loss functional with respect to the target  $T_\lambda$  (denoted by  $\mathcal{L}_\lambda$ ) is well defined and is given by

$$\mathcal{L}_\lambda(f) = \ell(f, T_\lambda) - \ell(f_*, T_\lambda).$$

Note that  $T_\lambda$  is a convex combination of points in  $\mathcal{T}$  and thus  $T_\lambda \in \mathcal{T}$ . Therefore, if the empirical minimization algorithm is to give fast rates it must produce with high probability functions for which the conditional expectation satisfies

$$\lim_{k \rightarrow \infty} \sqrt{k} \mathbb{E} \left( \mathcal{L}_\lambda(\hat{f}) | X_1, \dots, X_k \right) = 0.$$

On the other hand, take any other (fixed)  $f_1 \in F_{*,T}$ . It is clear that  $\mathbb{E}|\ell(f_1, T) - \ell(f_*, T)| \equiv \Delta > 0$ , otherwise  $\ell(f, T) = \ell(f_*, T)$  almost surely, contradicting our assumption that those are distinct points in  $F_{*,T}$ .

Observe that  $\text{var}(\mathcal{L}_\lambda(f_1))$  tends to  $\text{var}(\ell(f_1, T) - \ell(f_*, T))$  as  $\lambda \rightarrow 0$ . Thus, there is a constant  $\lambda_0$  such that for  $\lambda \leq \lambda_0$ ,  $\mathcal{L}_\lambda(f_1)$  has a “large” variance  $\sigma^2$  satisfying  $\sigma \geq \Delta/2$ . On the other hand, a rather simple calculation shows that the expectation of  $\mathcal{L}_\lambda(f_1)$  is at most  $c\lambda$ .

since typical values of  $P_k \mathcal{L}_\lambda(f_1)$  are

$$\frac{1}{k} \sum_{i=1}^k (\mathcal{L}_\lambda(f_1))(X_i) \sim \mathbb{E}(\mathcal{L}_\lambda(f_1)) \pm \sigma/\sqrt{k}$$

then by a quantitative version of the Central Limit Theorem, there is an integer  $k_0$  (that depends only on  $T$ ,  $f_1$  and  $\ell$ ) such that for every  $k \geq k_0$  there is a set of samples  $(X_i)_{i=1}^k$ , denoted by  $S_k$ , with the following properties:

1. The measure  $\mu^k(S_k) \geq 1/4$ , and
2. for every sample in  $S_k$ ,

$$P_k \mathcal{L}_\lambda(f_1) \leq -c_1 \frac{\Delta}{\sqrt{k}} + c_2 \lambda \leq -c_3 \frac{\Delta}{\sqrt{k}},$$

if one takes  $\lambda \leq \min\{c_4(\Delta)/\sqrt{k}, \lambda_0\}$ .

Hence, for such a choice of  $\lambda$ , if  $f$  is a potential empirical minimizer for  $\mathcal{L}_\lambda$  with respect to a sample  $(X_i)_{i=1}^k \in S_k$ , then its empirical error must be smaller than  $-c_3 \Delta/\sqrt{k}$ .

On the other hand, by the uniform law of large numbers, for every  $r > 0$  the empirical process indexed by  $\mathcal{L}_{\lambda,r}(F) = \{\mathcal{L}_\lambda(f) : \mathbb{E}\mathcal{L}_\lambda(f) \leq r\}$  satisfies that with high  $\mu^k$ -probability (say, at least  $5/6$ ),

$$\|P_k - P\|_{\mathcal{L}_{\lambda,r}(F)} \leq c\mathbb{E}\|P_k - P\|_{\mathcal{L}_{\lambda,r}(F)} \equiv W_\lambda(r),$$

where  $c$  is an absolute constant. In other words, if  $g \in F$  satisfies that  $\mathcal{L}_\lambda(g) \leq r$  then

$$|P_k\mathcal{L}_\lambda(g)| \leq W_\lambda(r) + r.$$

It remain to show that if we select  $\lambda \sim 1/\sqrt{k}$  and  $r \sim 1/\sqrt{k}$ , then  $W_\lambda(r) + r \leq c_3\Delta/(2\sqrt{k})$ . Thus, if  $\mathcal{L}_\lambda(f)$  has small expectation, its empirical expectation is not negative enough to defeat the empirical error generated by  $f_1$  and the empirical minimizer  $\hat{f}$  must have a relatively large risk.

We will show that one can find such choices of  $\lambda$  and  $r$  if the Gaussian processes indexed by random coordinate projections of  $F$  are continuous in some sense with respect to the norm on  $E$ . And, this continuity assumption is satisfied, for example, if  $F$  is a  $\mu$ -Donsker class and  $\ell$  is the squared loss.

This analysis shows that for  $\lambda_k \sim 1/\sqrt{k}$ , the target  $T_{\lambda_k}$  causes the risk of the empirical minimizer to be larger than  $c/\sqrt{k}$  for typical samples of size  $k$ . However, as we mentioned in the introduction, it was proved in [18] that for each  $0 < \lambda < 1$ , the error rate associated with the fixed target  $T_\lambda$ , that is, the rate at which

$$\mathbb{E}_{X_1, \dots, X_k} \left( \mathbb{E} \left( \mathcal{L}_\lambda(\hat{f}) | X_1, \dots, X_k \right) \right)$$

tends to 0 as a function of  $k$  will be significantly better than  $1/\sqrt{k}$  if  $F$  is small enough. Thus, the slow uniform error rate is not caused by a single function, and the guilty party truly changes with the sample size  $k$ .

### 3 A detailed proof

In this section we will present the details of the construction leading to the promised lower bound.

Consider  $\mathcal{L}_\lambda(F)$  and  $W_\lambda(r)$  as above and set  $F_{\lambda,r} = \{f : \mathbb{E}\mathcal{L}_\lambda(f) \leq r\}$ . The first part of the proof will be to show that with high probability, if  $\lambda \sim 1/\sqrt{k}$  then functions  $\mathcal{L}_\lambda(f)$  with expectation smaller than  $c/\sqrt{k}$  have an empirical error that is close to 0, and in particular, not “very negative”. This, of course, requires some assumption on the richness of  $F$ , which will be captured by a variant of the notion of asymptotic equicontinuity (see, e.g. [13, 9]).

**Definition 3.1** We say that  $F$  is compatible with the norm  $\|\cdot\|$  if the following holds. For every  $u > 0$  there is some integer  $k_0$  and  $q > 0$  for which

$$\sup_{k \geq k_0} \frac{1}{\sqrt{k}} \mathbb{E} \sup_{\{f, h \in F: \|f-h\| \leq q\}} \left| \sum_{i=1}^k g_i(f-h)(X_i) \right| \leq u.$$

In other words, the oscillation of the Gaussian process indexed by a “typical” coordinate projection of  $F$  is well behaved with respect to the norm on  $E$ .

A fundamental fact due to Giné and Zinn (see, e.g. [13], Theorem 14.6 and [9]), is that a  $\mu$ -Donsker class is compatible with the  $L_2$  norm, and therefore it is compatible with any  $L_p$  norm for  $2 \leq p < \infty$ .

**Theorem 3.2** Let  $F$  be a class of functions that is compatible with the norm on  $E$  and set  $F_{\lambda, r} = \{f : \mathbb{E} \mathcal{L}_\lambda(f) \leq r\}$ . For every  $\alpha, u > 0$  there is an integer  $k_0$  and  $0 < \beta \leq u$  such that for  $k \geq k_0$ , with probability  $5/6$ ,

$$\sup_{f \in F_{\alpha/\sqrt{k}, \beta/\sqrt{k}}} |P_k \mathcal{L}_{\alpha/\sqrt{k}}(f)| \leq \frac{2u}{\sqrt{k}}.$$

The first step in the proof of Theorem 3.2 is the following standard lemma.

**Lemma 3.3** There exists an absolute constant  $c$  such that for every  $r, \lambda > 0$ ,

$$\mathbb{E} \sup_{f \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k (\mathcal{L}_\lambda(f))(X_i) - \mathbb{E} \mathcal{L}_\lambda(f) \right| \leq \frac{c \|\ell\|_{\text{lip}}}{k} \mathbb{E} \sup_{f, h \in F_{\lambda, r}} \left| \sum_{i=1}^k g_i(f-h)(X_i) \right|, \quad (3.1)$$

where  $(g_i)_{i=1}^k$  are independent, standard Gaussian variables.

**Proof.** Fix  $r, \lambda > 0$ . By the Giné-Zinn symmetrization argument [11] and the fact that the expectation of the supremum of a Rademacher processes is dominated by the expectation of the supremum of the Gaussian process indexed by the same set (see, for example, [13], Chapter 4),

$$\begin{aligned} \mathbb{E} \sup_{f \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k (\mathcal{L}_\lambda(f))(X_i) - \mathbb{E} \mathcal{L}_\lambda(f) \right| &\leq 2 \mathbb{E}_X \mathbb{E}_\varepsilon \sup_{f \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k \varepsilon_i (\mathcal{L}_\lambda(f))(X_i) \right| \\ &\leq c \mathbb{E}_X \mathbb{E}_g \sup_{f \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k g_i (\mathcal{L}_\lambda(f))(X_i) \right|, \end{aligned}$$

where  $c$  is an absolute constant.

Clearly, for every  $k$ -sample  $(X_i)_{i=1}^k$  and every  $h, f \in F_{\lambda, r}$ ,

$$\begin{aligned} & \mathbb{E}_g \left( \sum_{i=1}^k g_i(\mathcal{L}_\lambda(f))(X_i) - \sum_{i=1}^k g_i(\mathcal{L}_\lambda(h))(X_i) \right)^2 \\ & \leq \|\ell\|_{\text{lip}}^2 \sum_{i=1}^k ((f - f_*) - (f_* - h))^2(X_i) \\ & = \|\ell\|_{\text{lip}}^2 \mathbb{E}_g \left( \sum_{i=1}^k g_i(f - f_*) - \sum_{i=1}^k g_i(f_* - h)(X_i) \right)^2. \end{aligned}$$

Therefore, by Slepian's Lemma [9], for every  $k$ -sample,

$$\mathbb{E}_g \sup_{f \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k g_i(\mathcal{L}_\lambda(f))(X_i) \right| \leq \mathbb{E}_g \sup_{f, h \in F_{\lambda, r}} \left| \frac{1}{k} \sum_{i=1}^k g_i(f - h)(X_i) \right|,$$

and the claim follows by taking the expectation with respect to  $X_1, \dots, X_k$ .  $\blacksquare$

The next step is to show that the set  $F_{\lambda, r}$  is contained in a relatively small ball (with respect to the norm) around  $f_*$ . The main point in the proof of this fact is the following geometric lemma.

**Lemma 3.4** *Let  $E$  be a uniformly convex, smooth normed space and consider  $x, y \in E$  and  $\rho \in \mathbb{R}_+$  such that  $\|y - x\| = \rho$ . Let  $0 < \lambda < 1$  and set  $x_\lambda = (1 - \lambda)x + \lambda y$ . If  $z$  satisfies that  $\|z - x\| \geq \rho$  then*

$$\|x_\lambda - z\| - (1 - \lambda)\rho \geq 2\lambda\|z - x\|\delta_E \left( \frac{\|z - y\|}{\|z - x\|} \right).$$

**Proof.** Without loss of generality we can assume that  $x = 0$ . Fix  $z \neq y$  and by our assumption,  $\|z\| \geq \|y\|$ . Define the function

$$H(\lambda) = \frac{\|x_\lambda - z\|}{\|z\|} = \frac{\|\lambda y - z\|}{\|z\|}$$

and observe that  $H$  is a convex function and  $H(0) = 1$ . Also, since  $E$  is smooth,  $H$  is differentiable in  $\lambda = 0$ . Thus,  $H(\lambda) - H(0) = H(\lambda) - 1 \geq H'(0)\lambda$ . Therefore,

$$H(\lambda) - (1 - \lambda) \geq (H'(0) + 1)\lambda,$$

and to complete the proof one has to bound  $H'(0)$  from below.



**Corollary 3.5** Consider the excess loss  $\mathcal{L}_\lambda(f) = \ell(f, T_\lambda) - \ell(f_*, T_\lambda)$ , set  $D = \sup_{f \in F} \|T - f\|$  and  $\rho = \|T - f_*\|$ . Then, for every  $0 < \lambda \leq 1$  and every  $f \in F$ ,

$$\mathbb{E}\mathcal{L}_\lambda(f) \geq \phi'((1 - \lambda)\rho) \cdot 2\lambda D \delta_E \left( \frac{\|f - f_*\|}{D} \right).$$

In particular, for every  $0 < \lambda \leq 1$ ,  $f_*$  is the unique minimizer of  $\mathbb{E}\ell(\cdot, T_\lambda)$  in  $F$ .

**Proof.** Using Assumption 2.1 and applying Lemma 3.4 for  $x = T$ ,  $y = f_*$  and  $z = f$ ,

$$\begin{aligned} \mathbb{E}\mathcal{L}_\lambda(f) &= \phi(\|T_\lambda - f\|) - \phi(\|T_\lambda - f_*\|) \\ &\geq \phi'(\|T_\lambda - f_*\|) \cdot (\|T_\lambda - f\| - \|T_\lambda - f_*\|) \\ &\geq 2\lambda \|f - T\| \delta_E \left( \frac{\|f - f_*\|}{\|f - T\|} \right) \cdot \phi'(\|T_\lambda - f_*\|) \\ &\geq 2\lambda D \delta_E \left( \frac{\|f - f_*\|}{D} \right) \cdot \phi'(\|T_\lambda - f_*\|), \end{aligned}$$

where the first inequality follows from the convexity of  $\phi$ , the second one is Lemma 3.4 and the last one is the monotonicity property of  $\delta_E$ . Indeed, let  $\varepsilon_2 = \|f - f_*\|/\|f - T\| \geq \varepsilon_1 = \|f - f_*\|/D$ . Thus

$$\frac{D}{\|f - f_*\|} \cdot \delta_E \left( \frac{\|f - f_*\|}{D} \right) \leq \frac{\|f - T\|}{\|f - f_*\|} \cdot \delta_E \left( \frac{\|f - f_*\|}{\|f - T\|} \right),$$

as claimed.  $\blacksquare$

Let us clarify the meaning of Corollary 3.5. First of all, we will be interested in “small” values of  $\lambda$ . With this in mind,  $\phi'(\|T_\lambda - f_*\|) = \phi'((1 - \lambda)\rho)$  is a positive constant (when  $\lambda \rightarrow 0$ , it tends to the derivative of  $\phi$  at  $\rho$ , which is a fixed, positive number) and the term  $D \delta_E \left( \frac{\|f - f_*\|}{D} \right)$  also does not depend on  $\lambda$ , but rather on properties of  $E$ ,  $F$  and  $T$ , and the distance between  $f$  and  $f_*$ . In particular, the minimizer of  $\mathbb{E}\mathcal{L}_\lambda$  in  $F$  is unique, and moreover, for  $\lambda$  sufficiently small,

$$\mathbb{E}\mathcal{L}_\lambda(f_1) \geq c\lambda,$$

where  $c$  depends only on properties of  $E$ ,  $F$ ,  $\phi$  and  $\|f_1 - f_*\|$ .

The second outcome is that functions with risk  $\mathbb{E}\mathcal{L}_\lambda(f) \leq r$  are contained in a small ball around  $f_*$ . Since this is a straightforward application of Corollary 3.5 we omit its proof.

**Corollary 3.6** *Assume that  $\delta_E(\varepsilon) \geq \eta\varepsilon^p$  for some fixed  $\eta$ . Then, for every  $0 < \lambda < 1$  and  $r > 0$*

$$F_{\lambda,r} \subset B\left(f_*, C_0 \left(\frac{r}{\lambda}\right)^{1/p}\right),$$

where  $C_0$  depends only on  $\rho, \phi, D$  and  $\eta$ .

Set  $B = B\left(f_*, C_0 \left(\frac{r}{\lambda}\right)^{1/p}\right)$ . Combining Corollary 3.6 with (3.1) shows that there is a constant  $C_1$  such that for every  $k$ , with  $\mu^k$ -probability of at least  $5/6$ ,

$$\sup_{\{f: \mathbb{E}\mathcal{L}_\lambda(f) \leq r\}} |P_k(\mathcal{L}_\lambda(f))(X_i)| \leq r + \frac{C_1 \|\ell\|_{\text{lip}}}{k} \mathbb{E} \sup_{f,h \in F \cap B} \left| \sum_{i=1}^k g_i(f-h)(X_i) \right|, \quad (3.2)$$

where the expectation is with respect to both  $(g_i)_{i=1}^k$  and  $(X_i)_{i=1}^k$ .

**Proof of Theorem 3.2.** Fix  $\alpha, u > 0$ , let  $\beta > 0$  to be named later and set  $C_1$  as in (3.2). Applying the compatibility assumption for  $u' = u/C_1 \|\ell\|_{\text{lip}}$ , there is some  $k_0$  and  $q > 0$  such that for  $k \geq k_0$ ,

$$\frac{C_1 \|\ell\|_{\text{lip}}}{k} \mathbb{E} \sup_{\{f,h \in F: \|f-h\| \leq q\}} \left| \sum_{i=1}^k g_i(f-h)(X_i) \right| \leq \frac{u}{\sqrt{k}}.$$

Set  $\beta$  to satisfy that  $q = 2C_0 \left(\frac{\beta}{\alpha}\right)^{1/p}$  and since  $q$  can be made smaller if we wish, we can assume that  $\beta \leq u$ .

By Corollary 3.6, taking  $\lambda = \alpha/\sqrt{k}$  and  $r = \beta/\sqrt{k}$  it is evident that

$$F_{\alpha/\sqrt{k}, \beta/\sqrt{k}} \subset B\left(f_*, C_0 \left(\frac{\beta}{\alpha}\right)^{1/p}\right) \cap F.$$

Thus, with probability at least  $5/6$ ,

$$\begin{aligned} & \sup_{f \in \mathcal{L}_{\alpha/\sqrt{k}, \beta/\sqrt{k}}} |P_k F_{\alpha/\sqrt{k}}(f)| \\ & \leq \frac{\beta}{\sqrt{k}} + \frac{C_1 \|\ell\|_{\text{lip}}}{k} \mathbb{E} \sup_{\{f,h \in F: \|f-h\| \leq q\}} \left| \sum_{i=1}^k g_i(f-h)(X_i) \right| \\ & \leq \frac{2u}{\sqrt{k}}. \end{aligned}$$

■

### 3.1 Constructing “bad” functions

So far, we showed that with high  $\mu^k$ -probability, if  $\lambda \sim 1/\sqrt{k}$  then functions  $\mathcal{L}_\lambda(f)$  with an expectation smaller than  $\sim 1/\sqrt{k}$  have an empirical expectation that is not very negative. To complete the proof of the lower bound one has to construct functions with a very negative empirical excess loss  $P_k \mathcal{L}_\lambda(f)$ .

To make things even more difficult, recall that for every  $f \in F$ ,  $\mathbb{E} \mathcal{L}_\lambda(f) \geq 0$ , and by the central limit Theorem  $k^{-1/2} \sum_{i=1}^k (\mathcal{L}_\lambda(f) - \mathbb{E} \mathcal{L}_\lambda(f))$  converges to a centered Gaussian variable. Thus, the hope of generating some excess loss function  $\mathcal{L}_\lambda(g)$  with a relatively large expectation but with a negative empirical loss on a typical sample of cardinality  $k$  is realistic only if  $\mathbb{E} \mathcal{L}_\lambda(g) \leq c_1 \sigma / \sqrt{k}$ , where  $\sigma^2$  is the variance of  $\mathcal{L}_\lambda(g)$ , because typical values of  $P_k \mathcal{L}_\lambda(g)$  are  $\mathbb{E} \mathcal{L}_\lambda(g) \pm c_1 \sigma / \sqrt{k}$ . Thus, it seems natural to expect that the bad behavior, if indeed exists, is generated by a family of functions, each one for a different value of  $k$ .

We will show that  $\mathcal{L}_{\lambda_k}(f_1)$  is such a family of functions for the choice of  $\lambda_k \sim 1/\sqrt{k}$ .

Our starting point is the Berry-Esséen Theorem (see, e.g., [19])

**Theorem 3.7** *There exists an absolute constant  $c$  for which the following holds. Let  $(\xi)_{i=1}^k$  be independent, identically distributed random variables, with mean 0 and variance  $\sigma^2$  and denote by  $F_k(x)$  the distribution function of  $\frac{1}{\sigma\sqrt{k}} \sum_{i=1}^k \xi_i$ . Then,*

$$\sup_{x \in \mathbb{R}} |F_k(x) - \Phi(x)| \leq c \frac{\mathbb{E}|\xi_1|^3}{\sigma^3 \sqrt{k}},$$

where  $\Phi$  is the distribution function of a standard Gaussian variable.

Applying the Berry-Esséen Theorem to the random variable  $\xi = g(X) - \mathbb{E}g$  we obtain the following

**Corollary 3.8** *Let  $g$  be a function with variance  $\sigma$  and a finite third moment. Then, there is some constant  $k_1$  that depends only on  $\sigma$  and  $\mathbb{E}|g|^3$  and an absolute constant  $C_2$ , such that for every  $k \geq k_1$ , with probability at least  $1/4$ ,*

$$\frac{1}{k} \sum_{i=1}^k g(X_i) \leq \mathbb{E}g - C_2 \frac{\sigma}{\sqrt{k}}.$$

In what follows we abuse notation and write  $\mathcal{L}(f) = \ell(T, f) - \ell(T, f_*)$ .

**Lemma 3.9** *There exist constants  $C_3, C_4$  and  $\lambda_0, k_1$  that depend only on  $\rho = \|T - f_*\|$ ,  $\|\phi\|_{\text{lip}}$ ,  $\|\ell\|_{\text{lip}}$ ,  $D_3 = \sup_{f \in F} \|T - f\|_{L_3(\mu)}$  and  $\Delta = \mathbb{E}|\mathcal{L}(f_1)|$  for which the following hold. For every  $0 \leq \lambda \leq \lambda_0$ ,*

1.  $\mathbb{E}\mathcal{L}_\lambda(f_1) \leq 2\lambda\rho\|\phi\|_{\text{lip}}$ .
2.  $\text{var}(\mathcal{L}_\lambda(f_1)) \geq \Delta^2/4$ .
3. *For every  $k \geq k_1$  and  $\lambda_k \leq \min\{\lambda_0, C_3/\sqrt{k}\}$ , with probability at least  $1/4$ ,  $k^{-1} \sum_{i=1}^k (\mathcal{L}_\lambda(f_1))(X_i) \leq -C_4/\sqrt{k}$ .*

**Proof.** For the first part, since  $\|T - f_1\| = \|T - f_*\|$  and  $\phi$  is Lipschitz then

$$\begin{aligned} \mathbb{E}\mathcal{L}_\lambda(f_1) &= \phi(\|T_\lambda - f_1\|) - \phi(\|T_\lambda - f_*\|) \\ &= (\phi(\|T_\lambda - f_1\|) - \phi(\|T - f_1\|)) + (\phi(\|T - f_*\|) - \phi(\|T_\lambda - f_*\|)) \\ &\leq 2\|\phi\|_{\text{lip}}\|T - T_\lambda\| = 2\lambda\|\phi\|_{\text{lip}}\|T - f_*\| = 2\lambda\rho\|\phi\|_{\text{lip}}. \end{aligned}$$

Turning to the second part, recall that  $\mathbb{E}|\mathcal{L}(f_1)| = \Delta > 0$ . Since  $\mathcal{L}_\lambda(f_1)$  tends to  $\mathcal{L}(f_1)$  in  $L_2$  as  $\lambda \rightarrow 0$ , then by the first part and standard calculations, there is some  $\lambda_0$  such that for any  $0 < \lambda \leq \lambda_0$ ,

$$\text{var}(\mathcal{L}_\lambda(f_1)) \geq \Delta^2 - \mathbb{E}|\mathcal{L}^2(f_1) - \mathcal{L}_\lambda^2(f_1)| - (\mathbb{E}\mathcal{L}_\lambda(f_1))^2 \geq \frac{\Delta^2}{4}.$$

Finally, let  $k_1$  as in Corollary 3.8 and set  $0 < \lambda < \lambda_0$ . Since  $\text{var}(\mathcal{L}_\lambda(f_1)) \geq \Delta^2/4$  and

$$\mathbb{E}|\ell(T_\lambda, f_1) - \ell(T_\lambda, f_*)|^3 \leq \|\ell\|_{\text{lip}}^3 \mathbb{E}|f_1 - f_*|^3,$$

this choice of  $k_1$  does not depend on  $\lambda$  as long as  $\lambda \leq \lambda_0$ . Therefore, by Corollary 3.8, with probability at least  $1/4$

$$\frac{1}{k} \sum_{i=1}^k (\mathcal{L}_\lambda(f_1))(X_i) \leq \mathbb{E}\mathcal{L}_\lambda(f_1) - C_2 \frac{\Delta}{2\sqrt{k}} \leq 2\lambda\rho\|\phi\|_{\text{lip}} - C_2 \frac{\Delta}{2\sqrt{k}}.$$

Setting  $C_3 = C_2\Delta/(8\rho\|\phi\|_{\text{lip}})$  it is evident that if one takes  $\lambda \leq \min\{\lambda_0, C_3/\sqrt{k}\}$ , the claim follows.  $\blacksquare$

Now, let us formulate and prove our main result. To make the formulation simpler, we shall refer to the space  $E$ ,  $\rho = \|T - f_*\|$ ,  $\Delta = \mathbb{E}|\mathcal{L}(f_1)|$ ,  $D = \sup_{f \in F} \|T - f\|$ ,  $D_3 = \sup_{f \in F} \|T - f\|_{L_3(\mu)}$ , the functions  $\phi$  and  $\ell$  and the asymptotic equicontinuity properties of  $F$  (see Definition 3.1) as “the parameters of the problem”. The constants in the formulation and the proof of the main result depend on these parameters.

**Theorem 3.10** *Assume that  $E$  is a smooth, uniformly convex normed space with a modulus of convexity  $\delta_E(\varepsilon) \geq \eta\varepsilon^p$  for some  $2 \leq p < \infty$  and  $\eta > 0$ . Let  $\phi, \ell, T, T_\lambda$  and  $\mathcal{L}_\lambda$  as above and consider  $F \subset E$  that is compatible with the norm. There are constants  $\lambda'$  and  $k'$  and  $c_1, c_2$  depending on the parameters of the problem for which the following holds. Let  $k \geq k'$  and set  $\lambda_k = \min\{\lambda', c_1/\sqrt{k}\}$ . Then with probability at least  $1/12$ , the empirical minimizer  $\hat{f}$  of  $\sum_{i=1}^k \ell(f, T_{\lambda_k})(X_i)$  satisfies*

$$\mathbb{E}(\ell(f, T_{\lambda_k}) | X_1, \dots, X_k) \geq \inf_{f \in F} \mathbb{E}\ell(f, T_{\lambda_k}) + \frac{c_2}{\sqrt{k}}.$$

**Proof.** Let  $f_1$  be the function constructed in Lemma 3.9 and set  $k_1, \lambda_0, C_3, C_4$  and  $\lambda_k$  as in the formulation of the Lemma. Let  $\lambda' = \lambda_0$  and  $k > k_1$ . Clearly, by increasing  $k_1$  if needed, one can assume  $\lambda_k = C_3/\sqrt{k}$ . Therefore, by Lemma 3.9, with  $\mu^k$ -probability  $1/4$ ,  $P_k \mathcal{L}_\lambda(f_1) \leq -C_4/\sqrt{k}$ . Hence, for each  $k > k_1$ , with that probability, the empirical minimizer with respect to the target  $T_{C_3/\sqrt{k}}$  must have an empirical error smaller than  $-C_4/\sqrt{k}$ .

We now apply Theorem 3.2 for  $\alpha = C_3$  and  $u = C_4/4$ . Let  $\beta > 0$  and  $k_0$  be as in the assertion of that Theorem for those values of  $\alpha$  and  $u$ , and put  $k' = \max\{k_0, k_1\}$ . Hence, if  $k > k'$  and  $r_k = \beta/\sqrt{k}$ , then with  $\mu^k$ -probability  $5/6$ ,

$$\sup_{f \in F_{\lambda_k, r_k}} |P_k \mathcal{L}_{\lambda_k}(f)| = \sup_{f \in F_{C_3/\sqrt{k}, \beta/\sqrt{k}}} |P_k \mathcal{L}_{C_3/\sqrt{k}}(f)| \leq \frac{2u}{\sqrt{k}} = \frac{C_4}{2\sqrt{k}}.$$

Therefore, with  $\mu^k$ -probability of at least  $1/12$ , the empirical minimizer with respect to the target  $T_{C_3/\sqrt{k}}$  is outside  $F_{\lambda_k, r_k}$  and thus its risk satisfies

$$\mathbb{E}\left(\mathcal{L}_{\lambda_k}(\hat{f}) | X_1, \dots, X_k\right) \geq \frac{\beta}{\sqrt{k}}. \quad \blacksquare$$

As we mentioned before, it is well known that  $\mu$ -Donsker classes are compatible with the  $L_2(\mu)$  norm and thus are compatible with any  $L_p(\mu)$  norm for  $2 \leq p < \infty$ . Since the  $L_p$  norms are uniformly convex (with a polynomial modulus of convexity) and smooth for  $1 < p < \infty$ , Theorem 1.2 follows from Theorem 3.10.

### 3.2 Concluding remarks

Finally, let us turn to the result from [15], in which a lower bound for the agnostic learning problem was given in a nonconvex situation, with respect

to the squared loss. It is claimed there that if  $F$  is a compact, nonconvex subset of  $L_2$  the following holds: If the functions in  $F$  take values in a bounded interval  $\mathcal{Y}$  then there is another bounded interval  $\mathcal{Y}'$  that depends on  $F$  and  $\mathcal{Y}$ , such that the best error rate that holds *uniformly* for every  $Y$  that takes values in  $\mathcal{Y}'$  is  $c/\sqrt{k}$ . In the proof of Lemma 5 in [15], the authors use the fact that  $F$  is not convex to construct a “well bounded” function that has more than a unique best approximation in  $F$ . Their argument is as follows: consider a function  $f_c = \alpha f + (1 - \alpha)g \notin F$  where  $0 < \alpha < 1$  and  $f, g \in F$ . If  $f_c$  has more than a unique best approximation, the proof is complete. If not, let  $f_1$  be the unique best approximation and consider  $f(t) = tf_c + (1 - t)f_1$  for  $t \in R_+$ . The aim is to find some  $t$  for which

$$\{f \in F : \|f - f_1\| > 0 : \|f - f(t)\| = \|f - f_1\|, P_F f(t) = f_1\} \neq \emptyset,$$

where  $P_F$  is the nearest point projection onto  $F$ , since such a function  $f(t)$  will have the desired properties. To that end they look for the smallest  $t > 1$  for which

$$A_t = \{f \in F : \|f - f_1\| > 0 : \|f - f(t)\| = \|f - f_1\|\} \neq \emptyset.$$

Unfortunately, if  $t' = \inf\{t : A_t \neq \emptyset\}$  is not attained then it does not follow that there is some  $t$  for which  $f(t)$  has more than a unique best approximation in  $F$ . Moreover, it was not proved that the only way in which this infimum would not be attained is if  $f_c$  had more than a unique best approximation in  $F$ . To prove such a claim one would have to show that a “smooth dynamic” of  $P_F$  as one approaches  $f_c$  from “above” on ray  $\overrightarrow{f_1, f_c}$  is impossible.

Hence, there is a gap in the claim from [15] - that one can find a function that is “well bounded” with more than a unique best approximation in  $F$ . Of course, once such a function is found or is assumed to exist, the rest of the argument remain valid and Theorem 1.1 follows.

## References

- [1] M. Anthony, P.L. Bartlett, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, 1999.
- [2] P.L. Bartlett, O. Bousquet, S. Mendelson, Local Rademacher Complexities, *The Annals of Statistics*, 33(4), 1497-1537, 2005.
- [3] P.L. Bartlett, S. Mendelson, Empirical minimization, *Probability Theory and Related Fields*, 135, 311-334, 2006.

- [4] S. Boucheron, O. Bousquet, G. Lugosi, Theory of Classification: a Survey of Recent Advances, *ESAIM: Probability and Statistics*, 9, 323-375, 2005.
- [5] O. Bousquet, *Concentration Inequalities and Empirical Processes Theory Applied to the Analysis of Learning Algorithms*. PhD thesis, 2002.
- [6] R. Deville, G. Godefroy, V. Zizler, *Smoothness and renorming in Banach spaces* Wiley, 1993.
- [7] L. Devroye, L. Györfi, G. Lugosi, *A probabilistic theory of pattern recognition*, Springer, 1996.
- [8] J. Diestel, *Sequences and series in Banach spaces*, Graduate text in Mathematics, vol. 92, Springer, Berlin, 1984.
- [9] R.M. Dudley, *Uniform Central Limit Theorems*, Cambridge University Press, 1999.
- [10] P. Habala, P. Hajek, V. Zizler: *Introduction to Banach spaces* vol I and II, matfyzpress, Univ. Karlovy, Prague, 1996.
- [11] E. Giné, J. Zinn, Some limit theorems for empirical processes, *Annals of Probability*, 12(4), 929–989, 1984.
- [12] L. Györfi, M. Kohler, A. Krzyżak, H. Walk, *A distribution-free theory of nonparametric regression*, Springer series in Statistics, 2002.
- [13] M. Ledoux, M. Talagrand, *Probability in Banach Spaces: isoperimetry and processes*, Springer, 1991.
- [14] G. Lugosi, Pattern classification and learning theory, in L. Györfi (editor), *Principles of Nonparametric Learning* Springer, 1–56, 2002.
- [15] W.S. Lee, P.L. Bartlett, R.C. Williamson, The Importance of Convexity in Learning with Squared Loss, *IEEE Transactions on Information Theory*, 44(5), 1974-1980, 1998.
- [16] S. Mendelson, Improving the sample complexity using global data, *IEEE Transactions on Information Theory* 48(7), 1977-1991, 2002.
- [17] S. Mendelson, Geometric parameters in Learning Theory, *GAFa lecture notes*, LNM 1850, 193-236, 2004.
- [18] S. Mendelson, Obtaining fast error rates in nonconvex situations, *Journal of Complexity*, in press, doi:10.1016/j.jco.2007.09.001

- [19] D. Stroock, *Probability Theory, An analytic view*, Cambridge University Press, 1993.
- [20] S. van de Geer, *Empirical Processes in M-Estimation*, Cambridge University Press, 2000.
- [21] A.W. Van der Vaart, J.A. Wellner, *Weak convergence and Empirical Processes*, Springer-Verlag, 1996.