

# Smoothness, Low-Noise and Fast Rates

Nathan Srebro

Karthik Sridharan

Ambuj Tewari

Toyota Technological Institute at Chicago  
6045 S Kenwood Ave., Chicago, IL 60637

November 18, 2010

## Abstract

We establish an excess risk bound of  $\tilde{O}\left(H\mathcal{R}_n^2 + \sqrt{HL^*}\mathcal{R}_n\right)$  for ERM with an  $H$ -smooth loss function and a hypothesis class with Rademacher complexity  $\mathcal{R}_n$ , where  $L^*$  is the best risk achievable by the hypothesis class. For typical hypothesis classes where  $\mathcal{R}_n = \sqrt{R/n}$ , this translates to a learning rate of  $\tilde{O}(RH/n)$  in the separable ( $L^* = 0$ ) case and  $\tilde{O}\left(RH/n + \sqrt{L^*RH/n}\right)$  more generally. We also provide similar guarantees for online and stochastic convex optimization of a smooth non-negative objective.

## 1 Introduction

Consider empirical risk minimization for a hypothesis class  $\mathcal{H} = \{h : \mathcal{X} \rightarrow \mathbb{R}\}$  w.r.t. some non-negative loss function  $\phi(t, y)$ . That is, we would like to learn a predictor  $h$  with small risk

$$L(h) = \mathbb{E}[\phi(h(X), Y)]$$

by minimizing the empirical risk

$$\hat{L}(h) = \frac{1}{n} \sum_{i=1}^n \phi(h(x_i), y_i)$$

of an i.i.d. sample  $(x_1, y_1), \dots, (x_n, y_n)$ .

Statistical guarantees on the excess risk are well understood for *parametric* (i.e. finite dimensional) hypothesis classes. More formally, these are hypothesis classes with finite VC-subgraph dimension [24] (aka pseudo-dimension). For such classes learning guarantees can be obtained for any bounded loss function (i.e. s.t.  $|\phi| \leq b < \infty$ ) and the relevant measure of complexity is the VC-subgraph dimension.

Alternatively, even for some non-parametric hypothesis classes (i.e. those with infinite VC-subgraph dimension), e.g. the class of low-norm linear predictors

$$\mathcal{H}_B = \{h_w : \mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle \mid \|\mathbf{w}\| \leq B\},$$

guarantees can be obtained in terms of *scale-sensitive* measures of complexity such as fat-shattering dimensions [1], covering numbers [24] or Rademacher complexity [2]. The classical statistical learning theory approach for obtaining learning guarantees for such scale-sensitive classes is to rely on the Lipschitz constant  $D$  of  $\phi(t, y)$  w.r.t.  $t$  (i.e. bound on its derivative w.r.t.  $t$ ). The excess risk can then be bounded as (in expectation over the sample):

$$L(\hat{h}) \leq L^* + 2D\mathcal{R}_n(\mathcal{H}) = L^* + 2\sqrt{D^2 \frac{R}{n}} \tag{1}$$

where  $\hat{h} = \arg \min \hat{L}(h)$  is the empirical risk minimizer (ERM),  $L^* = \inf_h L(h)$  is the approximation error, and  $\mathcal{R}_n(\mathcal{H})$  is the Rademacher complexity of the class, which typically scales as  $\mathcal{R}_n(\mathcal{H}) = \sqrt{R/n}$ . E.g. for  $\ell_2$ -bounded linear predictors,  $R = B^2 \sup \|X\|_2^2$ . The Rademacher complexity can be bounded by other scale-sensitive complexity measures, such as the fat-shattering dimensions and covering numbers, yielding similar guarantees in terms of these measures.

In this paper we address two deficiencies of the guarantee (1). First, the bound applies only to loss functions with bounded derivative, like the hinge loss and logistic loss popular for classification, or the absolute-value ( $\ell_1$ ) loss for regression. It is not directly applicable to the squared loss  $\phi(t, y) = \frac{1}{2}(t - y)^2$ , for which the second derivative is bounded, but not the first. We could try to simply bound the derivative of the squared loss in terms of a bound on the magnitude of  $h(x)$ , but e.g. for norm-bounded linear predictors  $\mathcal{H}_B$  this results in a very disappointing excess risk bound of the form  $O(\sqrt{B^4(\max \|X\|)^4/n})$ . One aim of this paper is to provide clean bounds on the excess risk for smooth loss functions such as the squared loss with a bounded second, rather than first, derivative.

The second deficiency of (1) is the dependence on  $1/\sqrt{n}$ . The dependence on  $1/\sqrt{n}$  might be unavoidable in general. But at least for finite dimensional (parametric) classes, we know it can be improved to a  $1/n$  rate when the distribution is separable, i.e. when there exists  $h \in \mathcal{H}$  with  $L(h) = 0$  and so  $L^* = 0$ . In particular, if  $\mathcal{H}$  is a class of bounded functions with VC-subgraph-dimension  $d$  (e.g.  $d$ -dimensional linear predictors), then (in expectation over the sample) [23]:

$$L(\hat{h}) \leq L^* + O\left(\frac{dD \log n}{n} + \sqrt{\frac{dDL^* \log n}{n}}\right) \quad (2)$$

The  $1/\sqrt{n}$  term disappears in the separable case, and we get a graceful degradation between the  $1/\sqrt{n}$  non-separable rate and the  $1/n$  separable rate. Could we get a  $1/n$  separable rate, and such a graceful degradation, also in the non-parametric case?

As we will show, the two deficiencies are actually related. For non-parametric classes, and non-smooth Lipschitz loss, such as the hinge-loss, the excess risk might scale as  $1/\sqrt{n}$  and not  $1/n$ , *even in the separable case*. However, for  $H$ -smooth non-negative loss functions, where the second derivative of  $\phi(t, y)$  w.r.t.  $t$  is bounded by  $H$ , a  $1/n$  separable rate *is* possible. In Section 2 we obtain the following bound on the excess risk (up to logarithmic factors):

$$\begin{aligned} L(\hat{h}) &\leq L^* + \tilde{O}\left(H\mathcal{R}_n^2(\mathcal{H}) + \sqrt{HL^*\mathcal{R}_n(\mathcal{H})}\right) \\ &= L^* + \tilde{O}\left(\frac{HR}{n} + \sqrt{\frac{HRL^*}{n}}\right) \leq 2L^* + \tilde{O}\left(\frac{HR}{n}\right). \end{aligned} \quad (3)$$

In particular, for  $\ell_2$ -norm-bounded linear predictors  $\mathcal{H}_B$  with  $\sup \|X\|_2^2 \leq 1$ , the excess risk is bounded by  $\tilde{O}(HB^2/n + \sqrt{HB^2L^*/n})$ . Another interesting distinction between parametric and non-parametric classes, is that even for the squared-loss, the bound (3) is tight and the non-separable rate of  $1/\sqrt{n}$  is unavoidable. This is in contrast to the parametric (fine dimensional) case, where a rate of  $1/n$  is always possible for the squared loss, regardless of the approximation error  $L^*$  [16]. The differences between parametric and scale-sensitive classes, and between non-smooth, smooth and strongly convex (e.g. squared) loss functions are discussed in Section 4 and summarized in Table 1.

The guarantees discussed thus far are general learning guarantees for the stochastic setting that rely only on the Rademacher complexity of the hypothesis class, and are phrased in terms of minimizing some scalar loss function. In Section 3 we consider also the online setting, in addition to the stochastic setting, and present similar guarantees for online and stochastic convex optimization [34, 25]. The guarantees of Section 3 match equation (3) for the special case of a convex loss function and norm-bounded linear predictors, but Section 3 capture a more general setting of optimizing an arbitrary non-negative convex objective, which we require to be smooth (there is no separate discussion of a “predictor” and a scalar loss function in Section 3). Results

in Section 3 are expressed in terms of properties of the norm, rather than a measure of concentration like the Radamacher complexity as in (3) and Section 2. However, the online and stochastic convex optimization setting of Section 3 is also more restrictive, as we require the objective be convex (in Section 2 and for the bound (3) we make no assumption about the convexity of the hypothesis class  $\mathcal{H}$  nor the loss function  $\phi$ ).

Specifically, for a non-negative  $H$ -smooth *convex* objective (see exact definition in Section 3), over a domain bounded by  $B$ , we prove that the average online regret (and so also the excess risk of stochastic optimization) is bounded by  $O(HB^2/n + \sqrt{HB^2L^*/n})$ . Comparing with the bound of  $O(\sqrt{D^2B^2/n})$  when the loss is  $D$ -Lipschitz rather than  $H$ -smooth [34, 22], we see the same relationship discussed above for ERM. Unlike the bound (3) for the ERM, the convex optimization bound avoids polylogarithmic factors. The results in Section 3 also generalize to smoothness and boundedness with respect to non-Euclidean norms.

Studying the online and stochastic convex optimization setting (Section 3), in addition to ERM (Section 2), has several advantages. First, it allows us to obtain a learning guarantee for an efficient single-pass learning methods, namely stochastic gradient descent (or mirror descent), as well as for the non-stochastic regret. Second, the bound we obtain in the convex optimization setting (Section 3) is actually better than the bound for the ERM (Section 2) as it avoids all polylogarithmic and large constant factors. Third, the bound is applicable to other non-negative online or stochastic optimization problems beyond classification, including problems for which ERM is not applicable (see, e.g., [25]).

## 2 Empirical Risk Minimization with Smooth Loss

Recall that the Rademacher complexity of  $\mathcal{H}$  for any  $n \in \mathbb{N}$  given by [2]:

$$\mathcal{R}_n(\mathcal{H}) = \sup_{x_1, \dots, x_n \in \mathcal{X}} \mathbb{E}_{\sigma \sim \text{Unif}(\{\pm 1\}^n)} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \left| \sum_{i=1}^n h(x_i) \sigma_i \right| \right]. \quad (4)$$

Throughout we shall consider the “worst case” Rademacher complexity.

Our starting point is the learning bound (1) that applies to  $D$ -Lipschitz loss functions, i.e. such that  $|\phi'(t, y)| \leq D$  (we always take derivatives w.r.t. the first argument). What type of bound can we obtain if we instead bound the second derivative  $\phi''(t, y)$ ? We will actually avoid talking about the second derivative explicitly, and instead say that a function is  $H$ -smooth iff its derivative is  $H$ -Lipschitz. For twice differentiable  $\phi$ , this just means that  $|\phi''| \leq H$ . The central observation, which allows us to obtain guarantees for smooth loss functions, is that for a smooth loss, the derivative can be bounded in terms of the function value:

**Lemma 2.1.** *For an  $H$ -smooth non-negative function  $f : \mathbb{R} \mapsto \mathbb{R}$ , we have:  $|f'(t)| \leq \sqrt{4Hf(t)}$*

*Proof.* For any  $t, r$ , we have  $t < s < r$  for which  $f(r) = f(t) + f'(s)(r - t)$ . Now:

$$\begin{aligned} 0 &\leq f(r) = f(t) + f'(t)(r - t) + (f'(s) - f'(t))(r - t) \\ &\leq f(t) + f'(t)(r - t) + H |s - t| |r - t| \leq f(t) + f'(t)(r - t) + H(r - t)^2 \end{aligned}$$

Setting  $r = t - \frac{f'(t)}{2H}$  yields the desired bounds. □

This Lemma allows us to argue that close to the optimum value, where the *value* of the loss is small, then so is its derivative. Looking at the dependence of (1) on the derivative bound  $D$ , we are guided by the following heuristic intuition: Since we should be concerned only with the behavior around the ERM, perhaps it is enough to bound  $\phi'(\hat{\mathbf{w}}, x)$  at the ERM  $\hat{\mathbf{w}}$ . Applying Lemma 2.1 to  $L(\hat{h})$ , we can bound  $|\mathbb{E}[\phi'(\hat{\mathbf{w}}, X)]| \leq \sqrt{4HL(\hat{h})}$ . What we would actually want is to bound each  $|\phi'(\hat{\mathbf{w}}, x)|$  separately, or at least have the absolute value *inside* the expectation—this is where the non-negativity of the loss plays an

important role. Ignoring this important issue for the moment and plugging this instead of  $D$  into (1) yields  $L(\hat{h}) \leq L^* + 4\sqrt{HL(\hat{h})}\mathcal{R}_n(\mathcal{H})$ . Solving for  $L(\hat{h})$  yields the desired bound (3).

This rough intuition is captured by the following Theorem:

**Theorem 1.** *For an  $H$ -smooth non-negative loss  $\phi$  s.t.  $\forall_{x,y,h} |\phi(h(x), y)| \leq b$ , for any  $\delta > 0$  we have that with probability at least  $1 - \delta$  over a random sample of size  $n$ , for any  $h \in \mathcal{H}$ ,*

$$L(h) \leq \hat{L}(h) + K \left( \sqrt{\hat{L}(h)} \left( \sqrt{H} \log^{1.5} n \mathcal{R}_n(\mathcal{H}) + \sqrt{\frac{b \log(1/\delta)}{n}} \right) + H \log^3 n \mathcal{R}_n^2(\mathcal{H}) + \frac{b \log(1/\delta)}{n} \right)$$

and so:

$$L(\hat{h}) \leq L^* + K \left( \sqrt{L^*} \left( \sqrt{H} \log^{1.5} n \mathcal{R}_n(\mathcal{H}) + \sqrt{\frac{b \log(1/\delta)}{n}} \right) + H \log^3 n \mathcal{R}_n^2(\mathcal{H}) + \frac{b \log(1/\delta)}{n} \right)$$

where  $K < 10^5$  is a numeric constant derived from [21] and [6].

Note that only the ‘‘confidence’’ terms depended on  $b = \sup |\phi|$ , and this is typically not the dominant term—we believe it is possible to also obtain a bound that holds in expectation over the sample (rather than with high probability) and that avoids a direct dependence on  $\sup |\phi|$ .

To prove Theorem 1 we use the notion of Local Rademacher Complexity [3], which allows us to focus on the behavior close to the ERM. To this end, consider the following empirically restricted loss class

$$\mathcal{L}_\phi(r) := \left\{ (x, y) \mapsto \phi(h(x), y) : h \in \mathcal{H}, \hat{L}(h) \leq r \right\}$$

Lemma 2.2, presented below, solidifies the heuristic intuition discussed above, by showing that the Rademacher complexity of  $\mathcal{L}_\phi(r)$  scales with  $\sqrt{Hr}$ . The Lemma can be seen as a higher-order version of the Lipschitz Composition Lemma [2], which states that the Rademacher complexity of the *unrestricted* loss class is bounded by  $D\mathcal{R}_n(\mathcal{H})$ . Here, we use the second, rather than first, derivative, and obtain a bound that depends on the empirical restriction:

**Lemma 2.2.** *For a non-negative  $H$ -smooth loss  $\phi$  bounded by  $b$  and any function class  $\mathcal{H}$  bounded by  $B$ :*

$$\mathcal{R}_n(\mathcal{L}_\phi(r)) \leq \sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) \left( 16 \log^{3/2} \left( \frac{nB}{\mathcal{R}_n(\mathcal{H})} \right) - 14 \log^{3/2} \left( \frac{n\sqrt{12HB}}{\sqrt{b}} \right) \right)$$

*Proof.* In order to prove Lemma 2.2, we actually move from Rademacher complexity to covering numbers, use smoothness and Lemma 2.1 to obtain an  $r$ -dependent cover of the empirically restricted class, and then return to the Rademacher complexity. More specifically:

- We use a modified version of Dudley’s integral to bound the Rademacher complexity of the empirically restricted class in terms of its  $L_2$ -covering numbers.
- We use smoothness to get an  $r$ -dependent bound on the  $L_2$ -covering numbers of the empirically restricted loss class in terms of  $L_\infty$ -covering numbers of the unrestricted hypothesis class.
- We bound the  $L_\infty$ -covering numbers of the unrestricted class in terms of its fat-shattering dimension, which in turn can be bounded in terms of its Rademacher complexity.

Before we proceed, recall the following definitions of covering numbers and fat shattering dimension. For any  $\epsilon > 0$  and function class  $\mathcal{F} \subset \mathbb{R}^{\mathcal{Z}}$ :

The  $L_2$  covering number  $\mathcal{N}_2(\mathcal{F}, \epsilon, n)$  is the supremum over samples  $z_1, \dots, z_n$  of the size of a minimal cover  $\mathcal{C}_\epsilon$  such that  $\forall f \in \mathcal{F}, \exists f_\epsilon \in \mathcal{C}_\epsilon$  s.t.  $\sqrt{\frac{1}{n} \sum_{i=1}^n (f(z_i) - f_\epsilon(z_i))^2} \leq \epsilon$ .

The  $L_\infty$  covering number  $\mathcal{N}_\infty(\mathcal{F}, \epsilon, n)$  is the supremum over samples  $z_1, \dots, z_n$  of the size of a minimal cover  $\mathcal{C}_\epsilon$  such that  $\forall f \in \mathcal{F}, \exists f_\epsilon \in \mathcal{C}_\epsilon$  s.t.  $\max_{i \in [n]} |f(z_i) - f_\epsilon(z_i)| \leq \epsilon$ .

The fat-shattering dimension  $\text{fat}_\epsilon(\mathcal{F})$  at scale  $\epsilon$  is the maximum number of points  $\epsilon$ -shattered by  $\mathcal{F}$  (see e.g. [21]).

**Bounding  $\mathcal{R}_n(\mathcal{L}_\phi(r))$  in terms of  $\mathcal{N}_2(\mathcal{L}_\phi(r))$**  Dudley's integral bound lets us bound the Rademacher complexity of a class in terms of its empirical  $L_2$  covering number. Here we use a more refined version of Dudley's integral bound due to Mendelson [21] and more explicitly stated in [27] and included for completeness as Lemma A.3 in the Appendix:

$$\mathcal{R}_n(\mathcal{L}_\phi(r)) \leq \inf_{\alpha > 0} \left\{ 4\alpha + 10 \int_\alpha^{\sqrt{br}} \sqrt{\frac{\mathcal{N}_2(\mathcal{L}_\phi(r), \epsilon, n)}{n}} d\epsilon \right\} \quad (5)$$

**Bounding  $\mathcal{N}_2(\mathcal{L}_\phi(r))$  in terms of  $\mathcal{N}_\infty(\mathcal{H})$**  In the Appendix we show that a corollary of Lemma 2.1 is that for a non-negative  $H$ -smooth  $f(\cdot)$  we have  $(f(t) - f(r))^2 \leq 6H(f(t) + f(r))(t - r)^2$  (Lemma A.1). Using this inequality, for any sample  $(x_1, y_1), \dots, (x_n, y_n)$ :

$$\begin{aligned} \sqrt{\frac{1}{n} \sum_{i=1}^n (\phi(h(z_i), z_i) - \phi(h_\epsilon(z_i), z_i))^2} &\leq \sqrt{\frac{6H}{n} \sum_{i=1}^n (\phi(h(z_i), z_i) + \phi(h_\epsilon(z_i), z_i)) (h(z_i) - h_\epsilon(z_i))^2} \\ &\leq \sqrt{\frac{6H}{n} \sum_{i=1}^n (\phi(h(z_i), z_i) + \phi(h_\epsilon(z_i), z_i)) \sqrt{\max_{i \in [n]} (h(z_i) - h_\epsilon(z_i))^2}} \\ &\leq \sqrt{12Hr} \max_{i \in [n]} |h(z_i) - h_\epsilon(z_i)| \end{aligned}$$

That is, an empirical  $L_\infty$  cover of  $\{h \in \mathcal{H} : \hat{L}(h) \leq r\}$  at radius  $\epsilon/\sqrt{12Hr}$  is also an empirical  $L_2$  cover of  $\mathcal{L}_\phi(r)$  at radius  $\epsilon$ , and we can conclude that:

$$\mathcal{N}_2(\mathcal{L}_\phi(r), \epsilon, n) \leq \mathcal{N}_\infty\left(\{h \in \mathcal{H} : \hat{L}(h) \leq r\}, \frac{\epsilon}{\sqrt{12Hr}}, n\right) \leq \mathcal{N}_\infty\left(\mathcal{H}, \frac{\epsilon}{\sqrt{12Hr}}, n\right) \quad (6)$$

**Bounding  $\mathcal{N}_\infty(\mathcal{H})$  in terms of  $\mathcal{R}_n(\mathcal{H})$**  The  $L_\infty$  covering number at scale  $\epsilon/\sqrt{12Hr}$  can be bounded in terms of the fat shattering dimension at that scale as [21]:

$$\mathcal{N}_\infty\left(\mathcal{H}, \frac{\epsilon}{\sqrt{12Hr}}, n\right) \leq \left(\frac{n \sqrt{12Hr} B}{\epsilon}\right)^{\text{fat}_{\frac{\epsilon}{\sqrt{12Hr}}}(\mathcal{H})}. \quad (7)$$

Hence by Equation (5) we have:

$$\mathcal{R}_n(\mathcal{L}_\phi(r)) \leq \inf_{\alpha \geq 0} \left\{ 4\alpha + 10 \int_\alpha^{\sqrt{br}} \sqrt{\frac{\text{fat}_{\frac{\epsilon}{\sqrt{12Hr}}}(\mathcal{H}) \log\left(\frac{n \sqrt{12Hr} B}{\epsilon}\right)}{n}} d\epsilon \right\} \quad (8)$$

choosing  $\alpha = \sqrt{12Hr} \mathcal{R}_n(\mathcal{H})$ :

$$\leq 4\sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) + 10 \int_{\sqrt{12Hr} \mathcal{R}_n(\mathcal{H})}^{\sqrt{br}} \sqrt{\frac{\text{fat}_{\frac{\epsilon}{\sqrt{12Hr}}}(\mathcal{H}) \log\left(\frac{n \sqrt{12Hr} B}{\epsilon}\right)}{n}} d\epsilon \quad (9)$$

after a change of integration variable:

$$\leq 4\sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) + 10\sqrt{12Hr} \int_{\mathcal{R}_n(\mathcal{H})}^{\sqrt{\frac{b}{12H}}} \sqrt{\frac{\text{fat}_\epsilon(\mathcal{H}) \log\left(\frac{nB}{\epsilon}\right)}{n}} d\epsilon \quad (10)$$

bounding the fat-shattering dimension in terms of the Rademacher complexity (Lemma A.2):

$$\begin{aligned} &\leq 4\sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) + 20\sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) \int_{\mathcal{R}_n(\mathcal{H})}^{\sqrt{\frac{b}{12H}}} \frac{\sqrt{\log\left(\frac{nB}{\epsilon}\right)}}{\epsilon} d\epsilon \quad (11) \\ &\leq \sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) \left( 4 + 20 \left[ -\frac{2}{3} \log^{3/2}(nB/\epsilon) \right]_{\mathcal{R}_n(\mathcal{H})}^{\sqrt{b/12H}} \right) \\ &\leq \sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) \left( 4 + 14 \left( \log^{3/2}\left(\frac{nB}{\mathcal{R}_n(\mathcal{H})}\right) - \log^{3/2}\left(\frac{nB\sqrt{12H}}{\sqrt{b}}\right) \right) \right) \\ &\leq \sqrt{12Hr} \mathcal{R}_n(\mathcal{H}) \left( 18 \log^{3/2}\left(\frac{nB}{\mathcal{R}_n(\mathcal{H})}\right) - 14 \log^{3/2}\left(\frac{n\sqrt{12HB}}{\sqrt{b}}\right) \right) \quad \square \end{aligned}$$

*Proof of Theorem 1.* Equipped with Lemma 2.2, the proof follows standard Local Rademacher arguments. Applying Theorem 6.1 of [6] to  $\psi_n(r) = 56\sqrt{Hr} \log^{1.5} n \mathcal{R}_n(\mathcal{H})$  we can show that:

$$L(h) \leq \hat{L}(h) + 106r_n^* + \frac{48b}{n} (\log \frac{1}{\delta} + \log \log n) + \sqrt{\hat{L}(h) \left( 8r_n^* + \frac{4b}{n} (\log \frac{1}{\delta} + \log \log n) \right)}$$

where  $r_n^* = 56^2 H \log^3 n \mathcal{R}_n(\mathcal{H})$  is the solution to  $\psi_n(r) = r$ . Further details can be found in the Appendix.  $\square$

## 2.1 Related Results

Rates faster than  $1/\sqrt{n}$  have been previously explored under various conditions, including when  $L^*$  is small.

**The Finite Dimensional Case** Lee et al [16] showed faster rates for squared loss, exploiting the strong convexity of this loss function, even when  $L^* > 0$ , but only with finite VC-subgraph-dimension. Panchenko [23] provides fast rate results for general Lipschitz bounded loss functions, still in the finite VC-subgraph-dimension case. Bousquet [6] provided similar guarantees for linear predictors in Hilbert spaces when the spectrum of the kernel matrix (covariance of  $X$ ) is exponentially decaying, making the situation almost finite dimensional. All these methods rely on finiteness of effective dimension to provide fast rates. In this case, smoothness is not necessary. Our method, on the other hand, establishes fast rates, when  $L^* = 0$ , for function classes that do *not* have finite VC-subgraph-dimension. We show how in this non-parametric case, smoothness is necessary and plays an important role (see also Table 1).

**Aggregation** Tsybakov [31] studied learning rates for aggregation, where a predictor is chosen from the convex hull of a finite set of base predictors. This is equivalent to an  $\ell_1$  constraint where each base predictor is viewed as a “feature”. As with  $\ell_1$ -based analysis, since the bounds depend only logarithmically on the number of base predictors (i.e. dimensionality), and rely on the scale of change of the loss function, they are of “scale sensitive” nature. For such an aggregate classifier, Tsybakov obtained a rate of  $1/n$  when zero (or small) risk is achieved by one of the base classifiers. Using Tsybakov’s result, it is not enough for zero risk to be achieved by an aggregate (i.e. bounded  $\ell_1$ ) classifier in order to obtain the faster rate. Tsybakov’s core

result is thus in a sense more similar to the finite dimensional results, since it allows for a rate of  $1/n$  when zero error is achieved by a finite cardinality (and hence finite dimension) class.

Tsybakov then used the approximation error of a small class of base predictors w.r.t. a large hypothesis class (i.e. a covering) to obtain learning rates for the large hypothesis class by considering aggregation within the small class. However these results only imply fast learning rates for hypothesis classes with very low complexity. Specifically to get learning rates better than  $1/\sqrt{n}$  using these results, the covering number of the hypothesis class at scale  $\epsilon$  needs to behave as  $1/\epsilon^p$  for some  $p < 2$ . But typical classes, including the class of linear predictors with bounded norm, have covering numbers that scale as  $1/\epsilon^2$  and so these methods do not imply fast rates for such function classes. In fact, to get rates of  $1/n$  with these techniques, even when  $L^* = 0$ , requires covering numbers that do not increase with  $\epsilon$  at all, and so actually finite VC-subgraph-dimension.

Chesneau et al [10] extend Tsybakov’s work also to general losses, deriving similar results for Lipschitz loss function. The same caveats hold: even when  $L^* = 0$ , rates faster than  $1/\sqrt{n}$  require covering numbers that grow slower than  $1/\epsilon^2$ , and rates of  $1/n$  essentially require finite VC-subgraph-dimension. Our work, on the other hand, is applicable whenever the Rademacher complexity (equivalently covering numbers) can be controlled. Although it uses some similar techniques, it is also rather different from the work of Tsybakov and Chesneau et al, in that it points out the importance of smoothness for obtaining fast rates in the non-parametric case: Chesneau et al relied only on the Lipschitz constant, which we show, in Section 4, is not enough for obtaining fast rates in the non-parametric case, even when  $L^* = 0$ .

**Local Rademacher Complexities** Bartlett et al [3] developed a general machinery for proving possible fast rates based on local Rademacher complexities. However, it is important to note that the localized complexity term typically dominates the rate and still needs to be controlled. For example, Steinwart [29] used Local Rademacher Complexity to provide fast rate on the 0/1 loss of Support Vector Machines (SVMs) ( $\ell_2$ -regularized hinge-loss minimization) based on the so called “geometric margin condition” and Tsybakov’s margin condition. Steinwart’s analysis is specific to SVMs. We also use Local Rademacher Complexities in order to obtain fast rates, but do so for general hypothesis classes, based only on the standard Rademacher complexity  $\mathcal{R}_n(\mathcal{H})$  of the hypothesis classes, as well as the smoothness of the loss function and the magnitude of  $L^*$ , but without any further assumptions on the hypothesis classes itself.

**Non-Lipschitz Loss** Beyond the strong connections between smoothness and fast rates which we highlight, we are also not aware of prior work providing an explicit and easy-to-use result for controlling a generic non-Lipschitz loss (such as the squared loss) solely in terms of the Rademacher complexity.

### 3 Online and Stochastic Optimization of Smooth Convex Objectives

We now turn to online and stochastic convex optimization. In these settings a learner chooses  $\mathbf{w} \in \mathbf{W}$ , where  $\mathbf{W}$  is a closed convex set in a normed vector space, attempting to minimize an objective (loss)  $\ell(\mathbf{w}, z)$  on instances  $z \in \mathcal{Z}$ , where  $\ell : \mathbf{W} \times \mathcal{Z} \rightarrow \mathbb{R}$  is an objective function which is convex in  $\mathbf{w}$ . This captures learning linear predictors w.r.t. a convex loss function  $\phi(t, z)$ , where  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  and  $\ell(\mathbf{w}, (x, y)) = \phi(\langle \mathbf{w}, x \rangle, y)$ , and extends well beyond supervised learning.

We consider the case where the objective  $\ell(\mathbf{w}, z)$  is  $H$ -smooth w.r.t. some norm  $\|\mathbf{w}\|$  (the reader may choose to think of  $\mathbf{W}$  as a subset of a Euclidean or Hilbert space, and  $\|\mathbf{w}\|$  as the  $\ell_2$ -norm): By this we mean that for any  $z \in \mathcal{Z}$ , and all  $\mathbf{w}, \mathbf{w}' \in \mathbf{W}$

$$\|\nabla \ell(\mathbf{w}, z) - \nabla \ell(\mathbf{w}', z)\|_* \leq H \|\mathbf{w} - \mathbf{w}'\|$$

where  $\|\cdot\|_*$  is the dual norm. The key here is to generalize Lemma 2.1 to smoothness w.r.t. a vector  $\mathbf{w}$ , rather than scalar smoothness:

**Lemma 3.1.** *For an  $H$ -smooth non-negative  $f : \mathbf{W} \rightarrow \mathbb{R}$ , for all  $\mathbf{w} \in \mathbf{W}$ :  $\|\nabla f(\mathbf{w})\|_* \leq \sqrt{4Hf(\mathbf{w})}$*

*Proof.* For any  $\mathbf{w}_0$  such that  $\|\mathbf{w} - \mathbf{w}_0\| \leq 1$ , let  $g(t) = g(\mathbf{w}_0 + t(\mathbf{w} - \mathbf{w}_0))$ . For any  $t, s \in \mathbb{R}$ ,

$$\begin{aligned} |g'(t) - g'(s)| &= |\langle \nabla f(\mathbf{w}_0 + t(\mathbf{w} - \mathbf{w}_0)) - \nabla f(\mathbf{w}_0 + s(\mathbf{w} - \mathbf{w}_0)), \mathbf{w} - \mathbf{w}_0 \rangle| \\ &\leq \|\nabla f(\mathbf{w}_0 + t(\mathbf{w} - \mathbf{w}_0)) - \nabla f(\mathbf{w}_0 + s(\mathbf{w} - \mathbf{w}_0))\|_* \|\mathbf{w} - \mathbf{w}_0\| \\ &\leq H|t - s| \|\mathbf{w} - \mathbf{w}_0\|^2 \\ &\leq H|t - s| \end{aligned}$$

Hence  $g$  is  $H$ -smooth and so by Lemma 2.1  $|g'(t)| \leq \sqrt{4Hg(t)}$ . Setting  $t = 1$  we have,  $\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}_0 \rangle \leq \sqrt{4Hf(\mathbf{w})}$ . Taking supremum over  $\mathbf{w}_0$  such that  $\|\mathbf{w}_0 - \mathbf{w}\| \leq 1$  we conclude that

$$\|\nabla f(\mathbf{w})\|_* = \sup_{\mathbf{w}_0: \|\mathbf{w} - \mathbf{w}_0\| \leq 1} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}_0 \rangle \leq \sqrt{4Hf(\mathbf{w})}$$

In order to consider general norms, we will also need to rely on a non-negative regularizer  $F : \mathbf{W} \mapsto \mathbb{R}$  that is a 1-strongly convex (see Definition in e.g. [33]) w.r.t. to the norm  $\|\mathbf{w}\|$  for all  $\mathbf{w} \in \mathbf{W}$ . For the Euclidean norm we can use the squared Euclidean norm regularizer:  $F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$ .

### 3.1 Online Optimization Setting

In the online convex optimization setting we consider an  $n$  round game played between a learner and an adversary (Nature) where at each round  $i$ , the player chooses a  $\mathbf{w}_i \in \mathbf{W}$  and then the adversary picks a  $z_i \in \mathcal{Z}$ . The player's choice  $\mathbf{w}_i$  may only depend on the adversary's choices in *previous* rounds. The goal of the player is to have low average objective value  $\frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_i, z_i)$  compared to the best single choice in hindsight [9].

A classic algorithm for this setting is Mirror Descent [4], which starts at some arbitrary  $\mathbf{w}_1 \in \mathbf{W}$  and updates  $\mathbf{w}_{i+1}$  according to  $z_i$  and a stepsize  $\eta$  (to be discussed later) as follows:

$$\mathbf{w}_{i+1} \leftarrow \arg \min_{\mathbf{w} \in \mathbf{W}} \langle \eta \nabla \ell(\mathbf{w}_i, z_i) - \nabla F(\mathbf{w}_i), \mathbf{w} \rangle + F(\mathbf{w}) \quad (12)$$

For the Euclidean norm with  $F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$ , the update (12) becomes projected online gradient descent [34]:

$$\mathbf{w}_{i+1} \leftarrow \Pi_{\mathbf{W}}(\mathbf{w}_i - \eta \nabla \ell(\mathbf{w}_i, z_i)) \quad (13)$$

where  $\Pi_{\mathbf{W}}(\mathbf{w}) = \arg \min_{\mathbf{w}' \in \mathbf{W}} \|\mathbf{w} - \mathbf{w}'\|$  is the projection onto  $\mathbf{W}$ .

**Theorem 2.** *For any  $B \in \mathbb{R}$  and  $\bar{L}^*$  if we use stepsize  $\eta = \frac{1}{HB^2 + \sqrt{H^2B^4 + HB^2nL^*}}$  for the Mirror Descent algorithm then for any instance sequence  $z_1, \dots, z_n \in \mathcal{Z}$ , the average regret w.r.t. any  $\mathbf{w}^* \in \mathbf{W}$  s.t.  $F(\mathbf{w}^*) \leq B^2$  and  $\frac{1}{n} \sum_{j=1}^n \ell(\mathbf{w}^*, z_j) \leq \bar{L}^*$  is bounded by:*

$$\frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}_i, z_i) - \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}^*, z_i) \leq \frac{4HB^2}{n} + 2\sqrt{\frac{HB^2\bar{L}^*}{n}}$$

Note that the stepsize depends on the bound  $\bar{L}^*$  on the loss in hindsight.

*Proof.* The proof follows from Lemma 3.1 and Theorem 1 of [28], using  $U_1 = B^2$  and  $U_2 = n\bar{L}^*$  in the Theorem.  $\square$



### 3.2 Stochastic Optimization I: Stochastic Mirror Descent

An online algorithm can also serve as an efficient one-pass learning algorithm in the stochastic setting. Here, we again consider an i.i.d. sample  $z_1, \dots, z_n$  from some unknown distribution (as in Section 2), and we would like to find  $\mathbf{w}$  with low risk  $L(\mathbf{w}) = \mathbb{E}[\ell(\mathbf{w}, Z)]$ . When  $z = (\mathbf{x}, y)$  and  $\ell(\mathbf{w}, z) = \phi(\langle \mathbf{w}, \mathbf{x} \rangle, y)$  this agrees with the supervised learning risk discussed in the Introduction and analyzed in Section 2. But instead of focusing on the ERM, we run Mirror Descent (or Projected Online Gradient Descent in case of a Euclidean norm) on the sample, and then take  $\tilde{\mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i$ . Standard arguments [8] allow us to convert the online regret bound of Theorem 2 to a bound on the excess risk:

**Corollary 3.** *For any  $B \in \mathbb{R}$  and  $\overline{L^*}$  if we run Mirror Descent on a random sample with stepsize  $\eta = \frac{1}{HB^2 + \sqrt{H^2B^4 + HB^2n\overline{L^*}}}$ , then for any  $\mathbf{w}^* \in \mathbf{W}$  with  $F(\mathbf{w}^*) \leq B^2$  and  $L(\mathbf{w}^*) \leq \overline{L^*}$ , with expectation over the sample:*

$$L(\tilde{\mathbf{w}}_n) - L(\mathbf{w}^*) \leq \frac{4HB^2}{n} + 2\sqrt{\frac{HB^2\overline{L^*}}{n}}.$$

Again, one must know a bound  $\overline{L^*}$  on the risk in order to choose the stepsize.

It is instructive to contrast this guarantee with similar looking guarantees derived recently in the stochastic convex optimization literature [14]. There, the model is stochastic first-order optimization, i.e. the learner gets to see an unbiased estimate  $\nabla l(\mathbf{w}, z_i)$  of the gradient of  $L(\mathbf{w})$ . The variance of the estimate is assumed to be bounded by  $\sigma^2$ . The expected accuracy after  $n$  gradient evaluations then has two terms: a ‘‘accelerated’’ term that is  $O(H/n^2)$  and a slow  $O(\sigma/\sqrt{n})$  term. While this result is applicable more generally (since it doesn’t require non-negativity of  $\ell$ ), it is not immediately clear if our guarantees can be derived using it. The main difficulty is that  $\sigma$  depends on the norm of the gradient estimates. Thus, it cannot be bounded in advance even if we know that  $L(\mathbf{w}^*)$  is small. That said, it is intuitively clear that towards the end of the optimization process, the gradient norms will typically be small if  $L(\mathbf{w}^*)$  is small because of the self bounding property (Lemma 3.1). Exploring this connection can be fruitful direction for further research.

### 3.3 Stochastic Optimization II: Regularized Batch Optimization

It is interesting to note that using stability arguments, a guarantee very similar to Corollary 3, avoiding the polylogarithmic factors of Theorem 1 as well as the dependence on the bound on the loss ( $b$  in Theorem 1), can be obtained also for a ‘‘batch’’ learning rule similar to ERM, but incorporating penalty-type regularization. For a given regularization parameter  $\lambda > 0$  define the regularized empirical loss as

$$\hat{L}_\lambda(\mathbf{w}) := \hat{L}(\mathbf{w}) + \lambda F(\mathbf{w})$$

and consider the Regularized Empirical Risk Minimizer

$$\hat{\mathbf{w}}_\lambda = \arg \min_{\mathbf{w} \in \mathbf{W}} \hat{L}_\lambda(\mathbf{w}) \tag{14}$$

The following theorem provides a bound on excess risk similar to Corollary 3:

**Theorem 4.** *For any  $B \in \mathbb{R}$  and  $\overline{L^*}$  if we set  $\lambda = \frac{128H}{n} + \sqrt{\frac{128^2H^2}{n^2} + \frac{128H\overline{L^*}}{nB^2}}$  then for all  $\mathbf{w}^* \in \mathbf{W}$  with  $F(\mathbf{w}^*) \leq B^2$  and  $L(\mathbf{w}^*) \leq \overline{L^*}$ , we have that in expectation over sample of size  $n$ :*

$$L(\hat{\mathbf{w}}_\lambda) - L(\mathbf{w}^*) \leq \frac{256HB^2}{n} + \sqrt{\frac{2048HB^2\overline{L^*}}{n}}.$$

To prove Theorem 4 we use stability arguments similar to the ones used by Shalev-Shwartz et al [25], which are in turn based on Bousquet and Elisseeff [7]. However, while Shalev-Shwartz et al [25] use the notion of uniform stability, here it is necessary to look at stability in expectation to get the faster rates (uniform stability does not hold with the desired rate).

To use stability based arguments, for each  $i \in [n]$  we consider a perturbed sample where instance  $z_i$  is replaced by instance  $z'_i$  drawn independently from same distribution as  $z_i$ . Let  $\hat{L}^{(i)}(\mathbf{w}) = \frac{1}{n}(\sum_{j \neq i} \ell(\mathbf{w}, z_j) + \ell(\mathbf{w}, z'_i))$  be the empirical risk over the perturbed sample, and consider the corresponding regularized empirical risk minimizer  $\hat{\mathbf{w}}_\lambda^{(i)} = \arg \min_{\mathbf{w}} \hat{L}_\lambda^{(i)}(\mathbf{w})$ , where  $\hat{L}_\lambda^{(i)}(\mathbf{w}) = \hat{L}^{(i)}(\mathbf{w}) + \lambda F(\mathbf{w})$ . We first prove the following Lemma on the expected stability of the regularized minimizer:

**Lemma 3.2.** *For any  $i \in [n]$  we have that*

$$\mathbb{E}_{z_1, \dots, z_n, z'_i} \left[ \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i) \right] \leq \frac{32H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)]$$

*Proof.*

$$\begin{aligned} \hat{L}_\lambda(\hat{\mathbf{w}}_\lambda^{(i)}) - \hat{L}_\lambda(\hat{\mathbf{w}}_\lambda) &= \frac{\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i)}{n} + \frac{\ell(\hat{\mathbf{w}}_\lambda, z'_i) - \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z'_i)}{n} + \hat{L}_\lambda^{(i)}(\hat{\mathbf{w}}_\lambda^{(i)}) - \hat{L}_\lambda^{(i)}(\hat{\mathbf{w}}_\lambda) \\ &\leq \frac{\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i)}{n} + \frac{\ell(\hat{\mathbf{w}}_\lambda, z'_i) - \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z'_i)}{n} \\ &\leq \frac{1}{n} \|\hat{\mathbf{w}}_\lambda^{(i)} - \hat{\mathbf{w}}_\lambda\| \left( \|\nabla \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)\|_* + \|\nabla \ell(\hat{\mathbf{w}}_\lambda, z'_i)\|_* \right) \\ &\leq \frac{2\sqrt{H}}{n} \|\hat{\mathbf{w}}_\lambda^{(i)} - \hat{\mathbf{w}}_\lambda\| \left( \sqrt{\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)} + \sqrt{\ell(\hat{\mathbf{w}}_\lambda, z'_i)} \right) \end{aligned}$$

where the last inequality follows from Lemma 3.1. By  $\lambda$ -strong convexity of  $\hat{L}_\lambda$  we have that

$$\hat{L}_\lambda(\hat{\mathbf{w}}_\lambda^{(i)}) - \hat{L}_\lambda(\hat{\mathbf{w}}_\lambda) \geq \frac{\lambda}{2} \|\hat{\mathbf{w}}_\lambda^{(i)} - \hat{\mathbf{w}}_\lambda\|^2.$$

We can conclude that

$$\|\hat{\mathbf{w}}_\lambda^{(i)} - \hat{\mathbf{w}}_\lambda\| \leq \frac{4\sqrt{H}}{\lambda n} \left( \sqrt{\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)} + \sqrt{\ell(\hat{\mathbf{w}}_\lambda, z'_i)} \right)$$

This gives us:

$$\begin{aligned} \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i) &\leq \|\nabla \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)\|_* \|\hat{\mathbf{w}}_\lambda^{(i)} - \hat{\mathbf{w}}_\lambda\| \\ &\leq \sqrt{4H\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)} \left( \frac{4\sqrt{H}}{\lambda} \left( \sqrt{\ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i)} + \sqrt{\ell(\hat{\mathbf{w}}_\lambda, z'_i)} \right) \right) \\ &\leq \frac{16H}{\lambda n} \left( \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) + \ell(\hat{\mathbf{w}}_\lambda, z'_i) \right) \end{aligned}$$

Taking expectation:

$$\begin{aligned} \mathbb{E}_{z_1, \dots, z_n, z'_i} \left[ \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i) \right] &\leq \frac{16H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n, z'_i} \left[ \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) + \ell(\hat{\mathbf{w}}_\lambda, z'_i) \right] \\ &= \frac{16H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n, z'_i} \left[ L(\hat{\mathbf{w}}_\lambda^{(i)}) + L(\hat{\mathbf{w}}_\lambda) \right] = \frac{32H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)] \quad \square \end{aligned}$$

*Proof of Theorem 4.* By Lemma 3.2 we have :

$$\begin{aligned} \mathbb{E}_{z_1, \dots, z_n} [L_\lambda(\hat{\mathbf{w}}_\lambda) - L_\lambda(\mathbf{w}_\lambda^*)] &\leq \mathbb{E}_{z_1, \dots, z_n} [L_\lambda(\hat{\mathbf{w}}_\lambda) - \hat{L}_\lambda(\hat{\mathbf{w}}_\lambda)] = \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda) - \hat{L}(\hat{\mathbf{w}}_\lambda)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z_1, \dots, z_n, z'_i} \left[ \ell(\hat{\mathbf{w}}_\lambda^{(i)}, z_i) - \ell(\hat{\mathbf{w}}_\lambda, z_i) \right] \leq \frac{32H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)] \end{aligned}$$

Noting the definition of  $\hat{L}_\lambda(\mathbf{w})$  and rearranging we get

$$\mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda) - L(\mathbf{w}^*)] \leq \frac{32H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)] + \lambda F(\mathbf{w}^*) - \lambda F(\hat{\mathbf{w}}_\lambda) \leq \frac{32H}{\lambda n} \mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)] + \lambda F(\mathbf{w}^*)$$

Rearranging further we get

$$\mathbb{E}_{z_1, \dots, z_n} [L(\hat{\mathbf{w}}_\lambda)] - L(\mathbf{w}^*) \leq \left( \frac{1}{1 - \frac{32H}{\lambda n}} - 1 \right) L(\mathbf{w}^*) + \frac{\lambda}{1 - \frac{32H}{\lambda n}} F(\mathbf{w}^*)$$

plugging in the value of  $\lambda$  gives the result.  $\square$

## 4 Tightness

In this Section we return to the learning rates for the ERM for parametric and for scale-sensitive hypothesis classes (i.e. in terms of the dimensionality and in terms of scale sensitive complexity measures), discussed in the Introduction and analyzed in Section 2. We compare the guarantees on the learning rates in different situations, identify differences between the parametric and scale-sensitive cases and between the smooth and non-smooth cases, and argue that these differences are real by showing that the corresponding guarantees are tight. Although we discuss the tightness of the learning guarantees for ERM in the stochastic setting, similar arguments can also be made for online learning.

Table 1 summarizes the bounds on the excess risk of the ERM implied by Theorem 1 as well previous bounds for Lipschitz loss on finite-dimensional [23] and scale-sensitive [2] classes, and a bound for squared-loss on finite-dimensional classes [9, Theorem 11.7] that can be generalized to any smooth strongly convex loss.

Loss function is:	Parametric $\dim(\mathcal{H}) \leq d$ , $ h  \leq 1$	Scale-Sensitive $\mathcal{R}_n(\mathcal{H}) \leq \sqrt{R/n}$
$D$ -Lipschitz	$\frac{dD}{n} + \sqrt{\frac{dDL^*}{n}}$	$\sqrt{\frac{D^2 R}{n}}$
$H$ -smooth	$\frac{dH}{n} + \sqrt{\frac{dHL^*}{n}}$	$\frac{HR}{n} + \sqrt{\frac{HRL^*}{n}}$
$H$ -smooth and $\lambda$ -strongly Convex	$\frac{H}{\lambda} \frac{dH}{n}$	$\frac{HR}{n} + \sqrt{\frac{HRL^*}{n}}$

Table 1: Bounds on the excess risk, up to polylogarithmic factors.

We shall now show that the  $1/\sqrt{n}$  dependencies in Table 1 are unavoidable. To do so, we will consider the class  $\mathcal{H} = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\| \leq 1\}$  of  $\ell_2$ -bounded linear predictors (all norms in this Section are Euclidean), with different loss functions, and various specific distributions over  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq 1\}$  and  $Y = [0, 1]$ . For the non-parametric lower-bounds, we will allow the dimensionality  $d$  to grow with the sample size  $n$ .

### Infinite dimensional, Lipschitz (non-smooth), separable

Consider the absolute difference loss  $\phi(h(\mathbf{x}), y) = |h(\mathbf{x}) - y|$ , take  $d = 2n$  and consider the following distribution:  $X$  is uniformly distributed over the  $d$  standard basis vectors  $\mathbf{e}_i$  and if  $X = \mathbf{e}_i$ , then  $Y = \frac{1}{\sqrt{n}} r_i$ , where  $r_1, \dots, r_d \in \{\pm 1\}$  is an arbitrary sequence of signs unknown to the learner (say drawn randomly beforehand). Taking  $\mathbf{w}^* = \frac{1}{\sqrt{n}} \sum_{i=1}^n r_i \mathbf{e}_i$ ,  $\|\mathbf{w}^*\| = 1$  and  $L^* = L(\mathbf{w}^*) = 0$ . However any sample  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  reveals at most  $n$  of  $2n$  signs  $r_i$ , and no information on the remaining  $\geq n$  signs. This means that for any algorithm used by the learner, there exists a choice of  $r_i$ 's such that on at least  $n$  of the remaining points not seen by the learner the learner has to suffer a loss of at least  $1/\sqrt{n}$ , yielding an overall risk of at least  $1/(2\sqrt{n})$ .

### Infinite dimensional, smooth, non-separable, even if strongly convex

Consider the squared loss  $\phi(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2$  which is 2-smooth and 2-strongly convex. For any  $\sigma \geq 0$

let  $d = \sqrt{n}/\sigma$  and consider the following distribution:  $X$  is uniform over  $\mathbf{e}_i$  as before, but this time  $Y|X$  is random, with  $Y|(X = \mathbf{e}_i) \sim \mathcal{N}(\frac{r_i}{2\sqrt{d}}, \sigma)$ , where again  $r_i$  are pre-determined, unknown to the learner, random signs. The minimizer of the expected risk is  $\mathbf{w}^* = \sum_{i=1}^d \frac{r_i}{2\sqrt{d}} \mathbf{e}_i$ , with  $\|\mathbf{w}^*\| = \frac{1}{2}$  and  $L^* = L(\mathbf{w}^*) = \sigma^2$ . Furthermore, for any  $\mathbf{w} \in \mathbf{W}$ ,

$$L(\mathbf{w}) - L(\mathbf{w}^*) = \mathbb{E}[\langle \mathbf{w} - \mathbf{w}^*, \mathbf{x} \rangle]^2 = \frac{1}{d} \sum_{i=1}^d (\mathbf{w}[i] - \mathbf{w}^*[i])^2 = \frac{1}{d} \|\mathbf{w} - \mathbf{w}^*\|^2$$

If the norm constraint becomes tight, i.e.  $\|\hat{\mathbf{w}}\| = 1$ , then  $L(\hat{\mathbf{w}}) - L(\mathbf{w}^*) \geq 1/(4d) = \sigma/(4\sqrt{n}) = \sqrt{L^*}/(4\sqrt{n})$ . Otherwise, each coordinate is a separate mean estimation problem, with  $n_i$  samples, where  $n_i$  is the number of appearances of  $\mathbf{e}_i$  in the sample. We have  $\mathbb{E}[(\hat{\mathbf{w}}[i] - \mathbf{w}^*[i])^2] = \sigma^2/n_i$  and so

$$L(\hat{\mathbf{w}}) - L^* = \frac{1}{d} \|\hat{\mathbf{w}} - \mathbf{w}^*\|^2 = \frac{1}{d} \sum_{i=1}^d \frac{\sigma^2}{n_i} \geq \frac{\sigma^2}{d} \frac{d^2}{\sum_i n_i} = \frac{\sigma^2 d}{n} = \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{L^*}{n}}$$

**Finite dimensional, smooth, not strongly convex, non-separable:**

Take  $d = 1$ , with  $X = 1$  with probability  $q$  and  $X = 0$  with probability  $1 - q$ . Conditioned  $X = 0$  let  $Y = 0$  deterministically, while conditioned on  $X = 1$  let  $Y = +1$  with probability  $p = \frac{1}{2} + \frac{0.2}{\sqrt{qn}}$  and  $Y = -1$  with probability  $1 - p$ . Consider the following 1-smooth loss function, which is quadratic around the correct prediction, but linear away from it:

$$\phi(h(\mathbf{x}), y) = \begin{cases} (h(\mathbf{x}) - y)^2 & \text{if } |h(\mathbf{x}) - y| \leq 1/2 \\ |h(\mathbf{x}) - y| - 1/4 & \text{if } |h(\mathbf{x}) - y| \geq 1/2 \end{cases}$$

First note that irrespective of choice of  $\mathbf{w}$ , when  $\mathbf{x} = 0$  and so  $y = 0$  we always have  $h(\mathbf{x}) = 0$  and so suffer no loss. This happens with probability  $1 - q$ . Next observe that for  $p > 1/2$ , the optimal predictor is  $\mathbf{w}^* \geq 1/2$ . However, for  $n > 20$ , with probability at least 0.25,  $\sum_{i=1}^n y_i < 0$ , and so the empirical minimizer is  $\hat{\mathbf{w}} \leq -1/2$ . We can now calculate

$$L(\hat{\mathbf{w}}) - L^* > L(-1/2) - L(1/2) = q(2p - 1) + (1 - q)0 = \frac{0.4 q}{\sqrt{qn}} = \frac{0.4 \sqrt{q}}{\sqrt{n}}$$

However note that for  $p > 1/2$ ,  $\mathbf{w}^* = \frac{3}{2} - \frac{1}{2p}$  and so for  $n > 20$ :

$$L^* > \frac{q}{2}$$

Hence we conclude that with probability 0.25 over the sample,

$$L(\hat{\mathbf{w}}) - L^* > \sqrt{\frac{0.32L^*}{n}}$$

## 5 Implications

We demonstrate the implications of our results in several settings.

### 5.1 Improved Margin Bounds

“Margin bounds” provide a bound on the expected zero-one loss of a classifiers based on the margin zero-one error on the training sample. Koltchinskii and Panchenko [13] provides margin bounds for a generic class

$\mathcal{H}$  based on the Rademacher complexity of the class. This is done by using a non-smooth Lipschitz “ramp” loss that upper bounds the zero-one loss and is upper-bounded by the margin zero-one loss. However, such an analysis unavoidably leads to a  $1/\sqrt{n}$  rate even in the separable case, since as we discuss in Section 4, it is not possible to get a faster rate for a non-smooth loss. Following the same idea we use the following smooth “ramp”:

$$\phi(t) = \begin{cases} 1 & t \leq 0 \\ \frac{1+\cos(\pi t/\gamma)}{2} & 0 < t < \gamma \\ 0 & t \geq \gamma \end{cases}$$

This loss function is  $\frac{\pi^2}{4\gamma^2}$ -smooth and is lower bounded by the zero-one loss and upper bounded by the  $\gamma$  margin loss. Using Theorem 1 we can now provide improved margin bounds for the zero-one loss of any classifier based on empirical margin error. Denote  $\text{err}(h) = \mathbb{E} [\mathbb{1}_{\{h(x) \neq y\}}]$  the zero-one risk and for any  $\gamma > 0$  and sample  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathcal{X} \times \{\pm 1\}$  define the  $\gamma$ -margin empirical zero one loss as

$$\widehat{\text{err}}_\gamma(h) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i h(\mathbf{x}_i) < \gamma\}}$$

**Theorem 5.** *For any hypothesis class  $\mathcal{H}$ , with  $|h| \leq b$ , and any  $\delta > 0$ , with probability at least  $1 - \delta$ , simultaneously for all margins  $\gamma > 0$  and all  $h \in \mathcal{H}$ :*

$$\text{err}(h) \leq \widehat{\text{err}}_\gamma(h) + K \left( \sqrt{\widehat{\text{err}}_\gamma(h)} \left( \frac{\log^{1.5} n}{\gamma} \mathcal{R}_n(\mathcal{H}) + \sqrt{\frac{\log(\log(\frac{4b}{\gamma})/\delta)}{n}} \right) + \frac{\log^3 n}{\gamma^2} \mathcal{R}_n^2(\mathcal{H}) + \frac{\log(\log(\frac{4b}{\gamma})/\delta)}{n} \right)$$

where  $K$  is a numeric constant from Theorem 1

In particular, the above bound implies:

$$\text{err}(h) \leq 1.01 \widehat{\text{err}}_\gamma(h) + K \left( \frac{2 \log^3 n}{\gamma^2} \mathcal{R}_n^2(\mathcal{H}) + \frac{2 \log(\log(\frac{4b}{\gamma})/\delta)}{n} \right)$$

where  $K$  is an appropriate numeric constant.

Improved margin bounds of the above form have been previously shown specifically for linear prediction in a Hilbert space (as in Support Vector Machines) based on the PAC Bayes theorem [20, 15]. However these PAC-Bayes based results are specific to the linear function class. Theorem 5 is, in contrast, a generic concentration-based result that can be applied to any function class with and yields rates dominated by  $\mathcal{R}^2(\mathcal{H})$ .

## 5.2 Interaction of Norm and Dimension

Consider the problem of learning a low-norm linear predictor with respect to the squared loss  $\phi(t, z) = (t-z)^2$ , where  $\mathcal{X} \in \mathbb{R}^d$ , for finite but very large  $d$ , and where the expected norm of  $X$  is low. Specifically, let  $X$  be Gaussian with  $\mathbb{E} [\|X\|^2] = B$ ,  $Y = \langle \mathbf{w}^*, X \rangle + \mathcal{N}(0, \sigma^2)$  with  $\|\mathbf{w}^*\| = 1$ , and consider learning a linear predictor using  $\ell_2$  regularization. What determines the sample complexity? How does the error decrease as the sample size increases?

From a scale-sensitive statistical learning perspective, we expect that the sample complexity, and the decrease of the error, should depend on the norm  $B$ , especially if  $d \gg B^2$ . However, for any fixed  $d$  and  $B$ , even if  $d \gg B^2$ , asymptotically as the number of samples increase, the excess risk of norm-constrained or norm-regularized regression actually behaves as  $L(\hat{\mathbf{w}}) - L^* \approx \frac{d}{n} \sigma^2$ , and depends (to first order) only on the dimensionality  $d$  and not at all on  $B$  [17]. How does the scale sensitive complexity come into play?

The asymptotic dependence on the dimensionality alone can be understood through Table 1. In this non-separable situation, parametric complexity controls can lead to a  $1/n$  rate, ultimately dominating the  $1/\sqrt{n}$

rate resulting from  $L^* > 0$  when considering the scale-sensitive, non-parametric complexity control  $B$ . (The dimension-dependent behavior here is actually a bit better than in the generic situation—the well-posed Gaussian model allows the bound to depend on  $\sigma^2 = L^*$  rather than on  $\sup(w'x - y)^2 \approx B^2 + \sigma^2$ ).

Combining Theorem 4 with the asymptotic  $\frac{d}{n}\sigma^2$  behavior, and noting that at the worst case we can predict using a zero vector, yields the following overall picture on the expected excess risk of ridge regression with an optimally chosen  $\lambda$ :

$$L(\hat{\mathbf{w}}_\lambda) - L^* \leq O\left(\min\left(B^2, \frac{B^2}{n} + \frac{B\sigma}{\sqrt{n}}, \frac{d\sigma^2}{n}\right)\right)$$

Roughly speaking, each term above describes the behavior in a different regime of the sample size:

- The first (“random”) regime until  $n = \Theta(B^2)$  where the excess is  $B^2$ .
- The second (“low-noise”) regime, where the excess risk is dominated by the norm and behaves as  $B^2/n$ , until  $n = \Theta(B^2/\sigma^2)$  and  $L(\hat{\mathbf{w}}) = \Theta(L^*)$ .
- The third (“slow”) regime, where the excess risk is controlled by the norm and the approximation error and behaves as  $B\sigma/\sqrt{n}$ , until  $n = \Theta(d^2\sigma^2/B^2)$  and  $L(\hat{\mathbf{w}}) = L^* + \Theta(B^2/d)$ .
- the fourth (“asymptotic”) regimes, where the excess risk is dominated by the dimensionality and behaves as  $d/n$ .

This sheds further light on recent work on this phenomena by Liang and Srebro based on exact asymptotics of simplified situations [18].

### 5.3 Sparse Prediction

The use of the  $\ell_1$  norm has become very popular for learning sparse predictors in high dimensions, as in the LASSO. The LASSO estimator [30]  $\hat{\mathbf{w}}$  is obtained by considering the squared loss  $\phi(z, y) = (z - y)^2$  and minimizing  $\hat{L}(\mathbf{w})$  subject to  $\|\mathbf{w}\|_1 \leq B$ . Let us assume there is some (unknown) sparse reference predictor  $\mathbf{w}^0$  that has low expected loss and sparsity (number of non-zeros)  $\|\mathbf{w}^0\|_0 = k$ , and that  $\|\mathbf{x}\|_\infty \leq 1, y \leq 1$ . In order to choose  $B$  and apply Theorem 1 in this setting, we need to bound  $\|\mathbf{w}^0\|_1$ . This can be done by, e.g., assuming that the features  $\mathbf{x}[i]$  *in the support of  $\mathbf{w}^0$*  are mutually uncorrelated. Under such an assumption, we have:  $\|\mathbf{w}^0\|_1^2 \leq k\mathbb{E}[\langle \mathbf{w}^0, x \rangle^2] \leq 2k(L(\mathbf{w}^0) + \mathbb{E}[y^2]) \leq 4k$ . Thus, Theorem 1 along with Rademacher complexity bounds from [11] gives us,

$$L(\hat{\mathbf{w}}) \leq L(\mathbf{w}^0) + \tilde{O}\left(\frac{k \log(d)}{n} + \sqrt{\frac{k L(\mathbf{w}^0) \log(d)}{n}}\right). \quad (15)$$

It is possible to relax the no-correlation assumption to a bound on the correlations, as in mutual incoherence, or to other weaker conditions [26]. But in any case, unlike typical analysis for compressed sensing, where the goal is recovering  $\mathbf{w}^0$  itself, here we are only concerned with correlations *inside the support of  $\mathbf{w}^0$* . Furthermore, we do not need to require that the optimal predictor is sparse or close to being sparse, or that the model is well specified: only that there exists a good (low risk) predictor using a small number of fairly uncorrelated features.

Bounds similar to (15) have been derived using specialized arguments [12, 32, 5]—here we demonstrate that a simple form of these bounds can be obtained under very simple conditions, using the generic framework we suggest.

It is also interesting to note that the methods and results of Section 3 can also be applied to this setting. But since  $\|\mathbf{w}\|_1^2$  is *not* strongly convex with respect to  $\|\mathbf{w}\|_1$ , we must instead use the entropy regularizer

$$F(\mathbf{w}) = B \sum_i \mathbf{x}[i] \log \left( \frac{\mathbf{x}[i]}{1/d} \right) + \frac{B^2}{e} \quad (16)$$

which is non-negative and 1-strongly convex with respect to  $\|\mathbf{w}\|_1$  on  $\mathbf{W} = \{\mathbf{w} \in \mathbb{R}^d \mid \mathbf{w}[i] \geq 0, \|\mathbf{w}\|_1 \leq B\}$ , with  $F(\mathbf{w}) \leq B^2(1 + \log d)$  (we consider here only non-negative weights—in order to allow  $\mathbf{w}[i] < 0$  we can include also each features negation, doubling the dimensionality). Recalling that  $\|\mathbf{w}^0\|_1 \leq 2\sqrt{k}$  and using  $B = 2\sqrt{k}$  in (16), we have from Theorem 4 we that:

$$L(\hat{\mathbf{w}}_\lambda) \leq L(\mathbf{w}^0) + O \left( \frac{k \log(d)}{n} + \sqrt{\frac{k L(\mathbf{w}^0) \log(d)}{n}} \right). \quad (17)$$

where  $\hat{\mathbf{w}}_\lambda$  is the regularized empirical minimizer (14) using the entropy regularizer (16) with  $\lambda$  as in Theorem 4. The advantage here is that using Theorem 4 instead of Theorem 1 avoids the extra logarithmic factors (yielding a clean big- $O$  dependence in (17) as opposed to big- $\tilde{O}$  in (15)).

More interestingly, following Corollary 3, one can use stochastic mirror descent, taking steps of the form (12) with the entropy regularizer (16), to obtain the same performance guarantee as in (17). This provides an efficient, single-pass optimization approach to sparse prediction as an alternative to batch optimization with an  $\ell_1$ -norm constraint, and yielding the same (if not somewhat better) guarantees.

## References

- [1] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. *FOCS*, 0:292–301, 1993.
- [2] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *JMLR*, 3:463–482, 2002.
- [3] P.L. Bartlett, O. Bousquet, and S. Mendelson. Local rademacher complexities. *Annals of Statistics*, 33(4):1497–1537, 2005.
- [4] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31:167–175, 2003.
- [5] P.J. Bickel, Y. Ritov, and A.B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- [6] O. Bousquet. *Concentration Inequalities and Empirical Processes Theory Applied to the Analysis of Learning Algorithms*. PhD thesis, Ecole Polytechnique, 2002.
- [7] Olivier Bousquet and André Elisseeff. Stability and generalization. *J. Mach. Learn. Res.*, 2:499–526, 2002.
- [8] N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of on-line learning algorithms. In *NIPS*, pages 359–366, 2002.
- [9] N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- [10] Christophe Chesneau and Guillaume Lecu. Adapting to unknown smoothness by aggregation of thresholded wavelet estimators. 2006.

- [11] S.M. Kakade, K. Sridharan, and A. Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *NIPS*, 2008.
- [12] V. Koltchinskii. Sparsity in penalized empirical risk minimization. *Ann. Inst. H. Poincaré Probab. Statist.*, 45(1):7–57, 2009.
- [13] V. Koltchinskii and D. Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *Ann. of Stats.*, 30(1):1–50, 2002.
- [14] G. Lan. *Convex Optimization Under Inexact First-order Information*. PhD thesis, Georgia Institute of Technology, 2009.
- [15] J. Langford and J. Shawe-Taylor. PAC-Bayes & margins. In *Advances in Neural Information Processing Systems 15*, pages 423–430, 2003.
- [16] Wee Sun Lee, Peter L. Bartlett, and Robert C. Williamson. The importance of convexity in learning with squared loss. *IEEE Trans. on Information Theory*, 1998.
- [17] P. Liang, F. Bach, G. Bouchard, and M. I. Jordan. Asymptotically optimal regularization in smooth parametric models. In *NIPS*, 2010.
- [18] P. Liang and N. Srebro. On the interaction between norm and dimensionality: Multiple regimes in learning. In *ICML*, 2010.
- [19] P. Massart. Some applications of concentration inequalities to statistics. *Annales de la Faculté des Sciences de Toulouse*, IX(2):245–303, 2000.
- [20] D. A. McAllester. Simplified PAC-Bayesian margin bounds. In *COLT*, pages 203–215, 2003.
- [21] Shahar Mendelson. Rademacher averages and phase transitions in glivenko-cantelli classes. *IEEE Trans. On Information Theory*, 48(1):251–263, 2002.
- [22] A. Nemirovski and D. Yudin. *Problem complexity and method efficiency in optimization*. Nauka Publishers, Moscow, 1978.
- [23] D. Panchenko. Some extensions of an inequality of vapnik and chervonenkis. *Electronic Communications in Probability*, 7:55–65, 2002.
- [24] David Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, 1984.
- [25] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In *COLT*, 2009.
- [26] S. Shalev-Shwartz, N. Srebro, and T. Zhang. Trading accuracy for sparsity. Technical report, TTI-C, 2009. Available at [ttic.uchicago.edu/~shai](http://ttic.uchicago.edu/~shai).
- [27] N. Srebro and K. Sridharan. Note on refined dudley integral covering number bound. In <http://ttic.uchicago.edu/~karthik/dudley.pdf>, 2010.
- [28] S. Shalev-Shwartz. *Online Learning: Theory, Algorithms, and Applications*. PhD thesis, Hebrew University of Jerusalem, 2007.
- [29] I. Steinwart and C. Scovel. Fast rates for support vector machines using gaussian kernels. *ANNALS OF STATISTICS*, 35:575, 2007.
- [30] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B.*, 58(1):267–288, 1996.



- [31] A. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Annals of Statistics*, 32:135–166, 2004.
- [32] S. A. van de Geer. High-dimensional generalized linear models and the lasso. *Annals of Statistics*, 36(2):614–645, 2008.
- [33] C. Zalinescu. *Convex analysis in general vector spaces*. World Scientific Publishing Co. Inc., River Edge, NJ, 2002.
- [34] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *ICML*, 2003.

## A Technical proofs

**Lemma A.1.** For any  $H$ -smooth non-negative function  $f : \mathbb{R} \mapsto \mathbb{R}$  and any  $t, r \in \mathbb{R}$  we have that

$$(f(t) - f(r))^2 \leq 6H(f(t) + f(r))(t - r)^2$$

*Proof.* We start by noting that by the mean value theorem for any  $t, r \in \mathbb{R}$  there exists  $s$  between  $t$  and  $r$  such that

$$f(t) - f(r) = f'(s)(t - r) \tag{18}$$

By smoothness we have that

$$|f'(s) - f'(t)| \leq H|t - s| \leq H|t - r|.$$

Hence we see that

$$|f'(s)| \leq |f'(t)| + H|t - r| \tag{19}$$

We now consider two cases:

**Case I:** If  $|t - r| \leq \frac{|f'(t)|}{5H}$  then by Equation (19),  $|f'(s)| \leq 6/5|f'(t)|$ , and combining this with Equation (18) we have:

$$(f(t) - f(r))^2 \leq f'(s)^2(t - r)^2 \leq \frac{36}{25}f'(t)^2(t - r)^2$$

But Lemma 2.1 ensures  $f'(t)^2 \leq 4Hf(t)$  yielding:

$$\leq \frac{144}{25}Hf(t)(t - r)^2 < 6Hf(t)(t - r)^2 \tag{20}$$

**Case II:** On the other hand, when  $|t - r| > \frac{|f'(t)|}{5H}$ , we have from Equation (19) that  $|f'(s)| \leq 6H|t - r|$ . Plugging this into Equation (18) yields:

$$\begin{aligned} (f(t) - f(r))^2 &= |f(t) - f(r)| \cdot |f(t) - f(r)| \leq |f(t) - f(r)| (|f'(s)| |t - r|) \\ &\leq |f(t) - f(r)| (6H|t - r| \cdot |t - r|) = 6H|f(t) - f(r)| (t - r)^2 \\ &\leq 6H \max\{f(t), f(r)\} (t - r)^2 \end{aligned} \tag{21}$$

Combining the two cases, we have from Equations (20) and (21) and the non-negativity of  $f(\cdot)$ , that in either case:

$$(f(t) - f(r))^2 \leq 6H(f(t) + f(r))(t - r)^2 \quad \square$$

### Relating Fat-shattering Dimension and Rademacher complexity :

The following lemma upper bounds the fat-shattering dimension at scale  $\epsilon \geq \mathcal{R}_n(\mathcal{H})$  in terms of the Rademacher Complexity of the function class. The proof closely follows the arguments of Mendelson [21, discussion after Definition 4.2].

**Lemma A.2.** For any hypothesis class  $\mathcal{H}$ , any sample size  $n$  and any  $\epsilon > \mathcal{R}_n(\mathcal{H})$  we have that

$$\text{fat}_\epsilon(\mathcal{H}) \leq \frac{4n\mathcal{R}_n(\mathcal{H})^2}{\epsilon^2}$$

In particular, if  $\mathcal{R}_n(\mathcal{H}) = \sqrt{R/n}$  (the typical case), then  $\text{fat}_\epsilon(\mathcal{H}) \leq R/\epsilon^2$ .

*Proof.* Consider any  $\epsilon \geq \mathcal{R}_n(\mathcal{H})$ . Let  $x_1^*, \dots, x_{\text{fat}_\epsilon}^*$  be the set of  $\text{fat}_\epsilon$  shattered points. This means that there exists  $s_1, \dots, s_{\text{fat}_\epsilon}$  such that for any  $J \subset [\text{fat}_\epsilon]$  there exists  $h_J \in \mathcal{H}$  such that  $\forall i \in J, h_J(x_i) \geq s_i + \epsilon$  and  $\forall i \notin J, h_J(x_i) \leq s_i - \epsilon$ . Now consider a sample  $x_1, \dots, x_{n'}$  of size  $n' = \lceil \frac{n}{\text{fat}_\epsilon} \rceil \text{fat}_\epsilon$ , obtained by taking each  $x_i^*$  and repeating it  $\lceil \frac{n}{\text{fat}_\epsilon} \rceil$  times, i.e.  $x_i = x_{\lfloor \frac{i}{\text{fat}_\epsilon} \rfloor}^*$ . Now, following Mendelson's arguments:

$$\begin{aligned}
\mathcal{R}_{n'}(\mathcal{H}) &\geq \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \frac{1}{n'} \sup_{h \in \mathcal{H}} \left| \sum_{i=1}^{n'} \sigma_i h(x_i) \right| \right] \\
&\geq \frac{1}{2} \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \frac{1}{n'} \sup_{h, h' \in \mathcal{H}} \left| \sum_{i=1}^{n'} \sigma_i (h(x_i) - h'(x_i)) \right| \right] \quad (\text{triangle inequality}) \\
&= \frac{1}{2} \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \frac{1}{n'} \sup_{h, h' \in \mathcal{H}} \left| \sum_{i=1}^{\text{fat}_\epsilon} \left( \sum_{j=1}^{\lceil n/\text{fat}_\epsilon \rceil} \sigma_{(i-1)\text{fat}_\epsilon + j} \right) (h(x_i^*) - h'(x_i^*)) \right| \right] \\
&\geq \frac{1}{2} \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \frac{1}{n'} \sum_{i=1}^{\text{fat}_\epsilon} \left( \sum_{j=1}^{\lceil n/\text{fat}_\epsilon \rceil} \sigma_{(i-1)\text{fat}_\epsilon + j} \right) (h_R(x_i^*) - h_{\overline{R}}(x_i^*)) \right]
\end{aligned}$$

where for each  $\sigma_1, \dots, \sigma_{n'}$ ,  $R \subseteq [\text{fat}_\epsilon]$  is given by  $R = \left\{ i \in [\text{fat}_\epsilon] \mid \text{sign} \left( \sum_{j=1}^{\lceil n/\text{fat}_\epsilon \rceil} \sigma_{(i-1)\text{fat}_\epsilon + j} \right) \geq 0 \right\}$ ,  $h_R$  is the function in  $\mathcal{H}$  that  $\epsilon$ -shatters the set  $R$  and  $h_{\overline{R}}$  be the function that shatters the complement of set  $R$ .

$$\begin{aligned}
&\geq \frac{1}{2} \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \frac{1}{n'} \sum_{i=1}^{\text{fat}_\epsilon} \left| \sum_{j=1}^{\lceil n/\text{fat}_\epsilon \rceil} \sigma_{(i-1)\text{fat}_\epsilon + j} \right| 2\epsilon \right] \\
&\geq \frac{\epsilon}{n'} \sum_{i=1}^{\text{fat}_\epsilon} \mathbb{E}_{\sigma \sim \text{Unif}\{\pm 1\}^{n'}} \left[ \left| \sum_{j=1}^{\lceil n/\text{fat}_\epsilon \rceil} \sigma_{(i-1)\text{fat}_\epsilon + j} \right| \right] \\
&\geq \frac{\epsilon \text{fat}_\epsilon}{n'} \sqrt{\frac{\lceil n/\text{fat}_\epsilon \rceil}{2}} \quad (\text{Kinchine's inequality}) \\
&= \sqrt{\frac{\epsilon^2 \text{fat}_\epsilon}{2 n'}}.
\end{aligned}$$

We can now conclude that:

$$\text{fat}_\epsilon \leq \frac{2n' \mathcal{R}_{n'}^2(\mathcal{H})}{\epsilon^2} \leq \frac{4n \mathcal{R}_n^2(\mathcal{H})}{\epsilon^2}$$

where last inequality is because Rademacher complexity decreases with increase in number of samples and  $n \leq n' \leq 2n$  (because  $\epsilon \geq \mathcal{R}_n(\mathcal{H})$  which implies that  $\text{fat}_\epsilon < n$ ).  $\square$

### Dudley Upper Bound :

In order to state and prove the following Lemma, we shall find it simpler to use the empirical Rademacher complexity for a given sample  $x_1, \dots, x_n$  [2]:

$$\hat{\mathcal{R}}_n(\mathcal{H}) = \mathbb{E}_{\sigma \sim \text{Unif}(\{\pm 1\}^n)} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \left| \sum_{i=1}^n h(x_i) \sigma_i \right| \right]. \quad (22)$$

and the  $L_2$  covering number at scale  $\epsilon > 0$  specific to a sample  $x_1, \dots, x_n$ , denoted by  $N_2(\epsilon, \mathcal{F}, (x_1, \dots, x_n))$  as the size of a minimal cover  $\mathcal{C}_\epsilon$  such that

$$\forall f \in \mathcal{F}, \exists f_\epsilon \in \mathcal{C}_\epsilon \text{ s.t. } \sqrt{\frac{1}{n} \sum_{i=1}^n (f(z_i) - f_\epsilon(z_i))^2} \leq \epsilon.$$

We will also denote  $\hat{\mathbb{E}}[f^2] = \frac{1}{n} \sum_{i=1}^n f^2(x_i)$ .

The following Lemma is stated in terms of the empirical Rademacher complexity and covering numbers. Taking a supremum over samples of size  $n$ , we get the same relationship between the worst-case Rademacher complexity and covering numbers, as is used in Section 2.

**Lemma A.3** ([27] following [21]). *For any function class  $\mathcal{F}$  containing functions  $f : \mathcal{X} \mapsto \mathbb{R}$ , we have that*

$$\hat{R}_n(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ 4\alpha + 10 \int_{\alpha}^{\sup_{f \in \mathcal{F}} \sqrt{\hat{\mathbb{E}}[f^2]}} \sqrt{\frac{\log \mathcal{N}(\epsilon, \mathcal{F}, (x_1, \dots, x_n))}{n}} d\epsilon \right\}$$

*Proof.* Let  $\beta_0 = \sup_{f \in \mathcal{F}} \sqrt{\hat{\mathbb{E}}[f^2]}$  and for any  $j \in \mathbb{Z}_+$  let  $\beta_j = 2^{-j} \sup_{f \in \mathcal{F}} \sqrt{\hat{\mathbb{E}}[f^2]}$ . The basic trick here is the idea of chaining. For each  $j$  let  $T_j$  be a (proper)  $L_2$ -cover at scale  $\beta_j$  of  $\mathcal{F}$  for the given sample. For each  $f \in \mathcal{F}$  and  $j$ , pick an  $\hat{f}_j \in T_j$  such that  $\hat{f}_j$  is a  $\beta_j$  approximation of  $f$ . Now for any  $N$ , we express  $f$  by chaining as

$$f = f - \hat{f}_N + \sum_{i=1}^N (\hat{f}_i - \hat{f}_{i-1})$$

where  $\hat{f}_0 = 0$ . Hence for any  $N$  we have that

$$\begin{aligned} \hat{R}_n(\mathcal{F}) &= \frac{1}{n} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i \left( f(\mathbf{x}_i) - \hat{f}_N(\mathbf{x}_i) + \sum_{j=1}^N (\hat{f}_j(\mathbf{x}_i) - \hat{f}_{j-1}(\mathbf{x}_i)) \right) \right] \\ &\leq \frac{1}{n} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i (f(\mathbf{x}_i) - \hat{f}_N(\mathbf{x}_i)) \right] + \sum_{j=1}^N \frac{1}{n} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i (\hat{f}_j(\mathbf{x}_i) - \hat{f}_{j-1}(\mathbf{x}_i)) \right] \\ &\leq \frac{1}{n} \sqrt{\sum_{i=1}^n \sigma_i^2} \sup_{f \in \mathcal{F}} \sqrt{\sum_{i=1}^n (f(x_i) - \hat{f}_N(x_i))^2} + \sum_{j=1}^N \frac{1}{n} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i (\hat{f}_j(\mathbf{x}_i) - \hat{f}_{j-1}(\mathbf{x}_i)) \right] \\ &\leq \beta_N + \sum_{j=1}^N \frac{1}{n} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i (\hat{f}_j(\mathbf{x}_i) - \hat{f}_{j-1}(\mathbf{x}_i)) \right] \end{aligned} \quad (23)$$

where the step before last is due to Cauchy-Schwarz inequality and  $\sigma = [\sigma_1, \dots, \sigma_n]^{\top}$ . Now note that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\hat{f}_j(x_i) - \hat{f}_{j-1}(x_i))^2 &= \frac{1}{n} \sum_{i=1}^n \left( (\hat{f}_j(x_i) - f(x_i)) + (f(x_i) - \hat{f}_{j-1}(x_i)) \right)^2 \\ &\leq \frac{2}{n} \sum_{i=1}^n (\hat{f}_j(x_i) - f(x_i))^2 + \frac{2}{n} \sum_{i=1}^n (f(x_i) - \hat{f}_{j-1}(x_i))^2 \\ &\leq 2\beta_j^2 + 2\beta_{j-1}^2 = 6\beta_j^2 \end{aligned}$$

Now Massart's finite class lemma [19] states that if for any function class  $\mathcal{G}$ ,  $\sup_{g \in \mathcal{G}} \sqrt{\frac{1}{n} \sum_{i=1}^n g(x_i)^2} \leq R$ , then  $\hat{R}_n(\mathcal{G}) \leq \sqrt{\frac{2R^2 \log(|\mathcal{G}|)}{n}}$ . Applying this to function classes  $\{f - f' : f \in T_j, f' \in T_{j-1}\}$  (for each  $j$ ) we

get from Equation (23) that for any  $N$ ,

$$\begin{aligned}
\hat{R}_n(\mathcal{F}) &\leq \beta_N + \sum_{j=1}^N \beta_j \sqrt{\frac{12 \log(|T_j| |T_{j-1}|)}{n}} \\
&\leq \beta_N + \sum_{j=1}^N \beta_j \sqrt{\frac{24 \log |T_j|}{n}} \\
&\leq \beta_N + 10 \sum_{j=1}^N (\beta_j - \beta_{j+1}) \sqrt{\frac{\log |T_j|}{n}} \\
&\leq \beta_N + 10 \sum_{j=1}^N (\beta_j - \beta_{j+1}) \sqrt{\frac{\log \mathcal{N}(\beta_j, \mathcal{F}, (x_1, \dots, x_n))}{n}} \\
&\leq \beta_N + 10 \int_{\beta_{N+1}}^{\beta_0} \sqrt{\frac{\log \mathcal{N}(\epsilon, \mathcal{F}, (x_1, \dots, x_n))}{n}} d\epsilon
\end{aligned}$$

where the third step is because  $2(\beta_j - \beta_{j+1}) = \beta_j$  and we bounded  $\sqrt{24}$  by 5. Now for any  $\alpha > 0$ , pick  $N = \sup\{j : \beta_j > 2\alpha\}$ . In this case we see that by our choice of  $N$ ,  $\beta_{N+1} \leq 2\alpha$  and so  $\beta_N = 2\beta_{N+1} \leq 4\alpha$ . Also note that since  $\beta_N > 2\alpha$ ,  $\beta_{N+1} = \frac{\beta_N}{2} > \alpha$ . Hence we conclude that

$$\hat{R}_n(\mathcal{F}) \leq 4\alpha + 10 \int_{\alpha}^{\sup_{f \in F} \sqrt{\mathbb{E}[f^2]}} \sqrt{\frac{\log \mathcal{N}(\epsilon, \mathcal{F}, (x_1, \dots, x_n))}{n}} d\epsilon$$

Since the choice of  $\alpha$  was arbitrary we take an infimum over  $\alpha$ . □

### Detailed Proof of Main Result :

**Detailed proof of Theorem 1.** By Theorem 6.1 of [6] (specifically the displayed equation prior to the last one in the proof of the theorem) we have that if  $\psi_n$  is any sub-root function that satisfies for all  $r > 0$ ,  $\mathcal{R}_n(\mathcal{L}_\phi(r)) \leq \psi_n(r)$  then,

$$L(h) \leq \hat{L}(h) + 45r_n^* + \sqrt{L(h)} \left( \sqrt{8r_n^*} + \sqrt{\frac{4b(\log(1/\epsilon) + 6 \log \log n)}{n}} \right) + \frac{20b(\log(1/\epsilon) + 6 \log \log n)}{n} \quad (24)$$

where  $r_n^*$  is the largest solution to equation  $\psi_n(r) = r$ . Now by Lemma 2.2 we have that  $\psi_n(r) = 56\sqrt{Hr} \log^{1.5} n \hat{\mathcal{R}}_n(\mathcal{H})$  satisfies the property that for all  $r > 0$ ,  $\mathcal{R}_n(\mathcal{L}_\phi(r)) \leq \psi_n(r)$  and so using this we see that

$$r_n^* = 56^2 H \log^3 n \mathcal{R}_n(\mathcal{H})$$

and for this  $r_n^*$  Equation (24) holds. Now using the simple fact that for any non-negative  $A, B, C$ ,

$$A \leq B + C\sqrt{A} \Rightarrow A \leq B + C^2 + \sqrt{BC}$$

we conclude,

$$L(h) \leq \hat{L}(h) + 106 r_n^* + \frac{48b}{n} (\log \frac{1}{\epsilon} + \log \log n) + \sqrt{\hat{L}(h) \left( 8r_n^* + \frac{4b}{n} (\log \frac{1}{\epsilon} + \log \log n) \right)}$$

plugging in  $r_n^*$  we get the required statement. □