

Fast Convergence Rates for Excess Regularized Risk with Application to SVM

Karthik Sridharan, Nathan Srebro, Shai Shalev-Shwartz

1 Introduction

We consider the stochastic minimization problem

$$\operatorname{argmin}_{\mathbf{w} \in \mathbf{W}} F(\mathbf{w}) \quad (1)$$

where

$$F(\mathbf{w}) = \mathbf{E}_{\theta} [f(\mathbf{w}; \theta)] \quad (2)$$

is the expectation of a random function $f(\cdot; \theta)$, and the optimization is based on a sample $f(\cdot; \theta_1), \dots, f(\cdot; \theta_n)$ of i.i.d. instantiations of $f(\cdot; \theta)$. Instead of thinking of a “random function”, one can think of a fixed and known function $f(\mathbf{w}; \theta)$, and a distribution over the parameter θ . A special case is the familiar prediction setting where $\theta = (\mathbf{x}, y)$ is an instance-label pair and, e.g., $f(\mathbf{w}; \mathbf{x}, y) = \ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ for some loss function ℓ . However, we would like to emphasize that we do not assume any metric or other structure on the parameter θ .

In this paper, we are given a random sample of functions (i.e. parameters θ) and can find the empirical minimizer

$$\hat{\mathbf{w}} = \operatorname{arg min}_{\mathbf{w}} \hat{F}(\mathbf{w}) \quad (3)$$

where \hat{F} is the empirical average

$$\hat{F}(\mathbf{w}) = \hat{\mathbb{E}} [f(\mathbf{w}; \theta)] = \frac{1}{n} \sum_{i=1}^n f(\mathbf{w}; \theta_i). \quad (4)$$

Our guarantees are then on the expected value $F(\hat{\mathbf{w}})$ (generalization performance in the prediction setting) of the empirical minimizer. In fact, we do not need to assume that we perform the minimization (3) exactly: we provide uniform (over all $\mathbf{w} \in \mathbf{W}$) guarantees on the expected regret

$$F(\mathbf{w}) - F(\mathbf{w}^*) \quad (5)$$

in terms of the empirical regret

$$\hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*). \quad (6)$$

This is a stronger type of result than what can be obtained with online-to-batch conversions, as it applies to any possible solution \mathbf{w} , and not only some specific algorithmically defined solution. For example, it can be used to analyze the performance of approximate minimizers obtained through approximate optimization techniques.

2 Concentration of Strongly convex Linear Functions

We first introduce two conditions under which we show that the rate of convergence is of order $1/n$. The first condition assumes that the function at hand can be decomposed into a part only dependent on parameter \mathbf{w} and into a part that is Lipschitz w.r.t. $\langle \mathbf{w}, P(\theta) \rangle$ where P is a fixed function of θ

Condition 1. For all $\theta \in \Theta$ and $\mathbf{w} \in \mathbf{W}$:

$$f(w, \theta) = h(w, \theta) + N(w)$$

where N is a function of \mathbf{w} alone and $\forall \theta, \|P(\theta)\| \leq B_x$ Also for all $\theta \in \Theta$ and $\mathbf{w}_1, \mathbf{w}_2 \in \mathbf{W}$:

$$\begin{aligned} |h(\mathbf{w}_1, \theta) - h(\mathbf{w}_2, \theta)| &\leq L |\langle \mathbf{w}_1 - \mathbf{w}_2, P(\theta) \rangle| \\ |N(\mathbf{w}_1) - N(\mathbf{w}_2)| &\leq k \|\mathbf{w}_1 - \mathbf{w}_2\| \end{aligned}$$

Note that the condition implies that $|f(\mathbf{w}_1, \theta) - f(\mathbf{w}_2, \theta)| \leq (LB_x + k)\|\mathbf{w}_1 - \mathbf{w}_2\|$. The second condition we require is that the function in expectation is strongly convex.

Condition 2. $F(\mathbf{w})$ is λ -strongly convex. That is, for all $\mathbf{w}_1, \mathbf{w}_2 \in \mathbf{W}$:

$$F\left(\frac{\mathbf{w}_1 + \mathbf{w}_2}{2}\right) + \frac{\lambda}{8} \|\mathbf{w}_1 - \mathbf{w}_2\|^2 \leq \frac{1}{2}(F(\mathbf{w}_1) + F(\mathbf{w}_2))$$

Condition 2 ensures that:

$$F(\mathbf{w}) \geq F(\mathbf{w}^*) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}^*\|^2 \quad (7)$$

for all $\mathbf{w} \in \mathbf{W}$. Recall that $\mathbf{w}^* = \arg \min_{\mathbf{w}} F(\mathbf{w})$. Our results will actually only use (7).

Condition 2 only requires that the expected function $F(\mathbf{w})$ is strongly convex. Of course, requiring that $f(\mathbf{w}; \theta)$ is λ -strongly convex for all θ (with respect to \mathbf{w}) is enough to ensure the condition.

Theorem 1. If $f(\mathbf{w}, \theta)$ and \mathbf{W} satisfy Conditions 1 and 2, then, for any $\delta > 0$ and any $a > 0$, with probability of at least $1 - \delta$ over the sample, we have that for all $\mathbf{w} \in \mathbf{W}$:

$$\begin{aligned} F(\mathbf{w}) - F(\mathbf{w}^*) &\leq (1 + a)(\hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*)) \\ &\quad + O\left(\frac{(1 + \frac{1}{a})L^2(1 + \log(1/\delta))}{\lambda n}\right). \end{aligned}$$

as well as:

$$\begin{aligned} \hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*) &\leq (1 + a)(F(\mathbf{w}) - F(\mathbf{w}^*)) \\ &\quad + O\left(\frac{(1 + \frac{1}{a})L^2(1 + \log(1/\delta))}{\lambda n}\right). \end{aligned}$$

2.1 Proof of Theorem 1

To prove the theorem we use techniques of reweighing and peeling following ? (?). First, we need the following definitions. Let $g : \mathbf{W} \times \Theta \rightarrow \mathbb{R}$ be the function

$$g(\mathbf{w}, \theta) = f(\mathbf{w}; \theta) - f(\mathbf{w}^*; \theta) . \quad (8)$$

We also define the class of functions

$$\mathcal{G} = \{g_{\mathbf{w}}(\theta) = g(\mathbf{w}, \theta) : w \in \mathbf{W}\} . \quad (9)$$

We next define a weighted class of functions. To do so, let \mathbb{Z}_+ be the set of non-negative integers and let $\gamma \geq 1$ be some fixed scalar to be set later. For any $r > 0$ we define

$$\mathcal{G}_r = \left\{ \frac{g}{\gamma^k} : g \in \mathcal{G}, k = \min\{k' \in \mathbb{Z}_+ : \mathbb{E}[g] \leq r\gamma^{k'}\} \right\} \quad (10)$$

In other words, for any function $g \in \mathcal{G}$, we have a function $g_r \in \mathcal{G}_r$ such that g_r is just a scaled version of g and the scaling factor ensures that

$$\mathbb{E}[g_r] \leq r . \quad (11)$$

The following lemma upper bounds the range of the functions in \mathcal{G}_r .

Lemma 2. *Assume that the conditions stated in Theorem 1 hold and let \mathcal{G}_r be as defined in (10). Then, for all $g_r \in \mathcal{G}_r$ and $\theta \in \Theta$ we have $|g_r(\theta)| \leq L\sqrt{2r/\lambda}$.*

Proof. Using the definition of g_r we have that there exists $w \in \mathbf{W}$ such that

$$g_r(\theta) = \frac{f(\mathbf{w}, \theta) - f(\mathbf{w}^*, \theta)}{\gamma^k} .$$

Since $f(\mathbf{w}, \theta)$ is L -Lipschitz w.r.t. w we get that

$$|f(\mathbf{w}, \theta) - f(\mathbf{w}^*, \theta)| \leq L \|w - w^*\| .$$

Combining the above two equations we obtain that,

$$|g_r(\theta)| \leq \frac{L \|w - w^*\|}{\gamma^k} . \quad (12)$$

Next, we obtain a bound on $\|w - w^*\|$. To do so, we rewrite equation (7):

$$\|\mathbf{w} - \mathbf{w}^*\| \leq \sqrt{\frac{2}{\lambda}(F(\mathbf{w}) - F(\mathbf{w}^*))} . \quad (13)$$

Using again the definition of F we know that there exists $g \in \mathcal{G}$ such that $g_r = g/\gamma^k$ and

$$F(\mathbf{w}) - F(\mathbf{w}^*) = \mathbb{E}[g] = \gamma^k \mathbb{E}[g_r] . \quad (14)$$

Using (11) we know that $\mathbb{E}[g_r] \leq r$. Combining this with (14) and (13) we obtain that

$$\|\mathbf{w} - \mathbf{w}^*\| \leq \sqrt{\frac{2}{\lambda}\gamma^k r} .$$

Combining the above with (12) yields

$$|g_r(\theta)| \leq \frac{L \sqrt{2r/\lambda}}{\gamma^{k/2}}.$$

Finally, we note that $\gamma > 1$ which concludes our proof. \square

Denote $c(r) = L\sqrt{2r/\lambda}$. Using Theorem 3.2 of ? (?) we have that with probability of at least $1 - \delta$,

$$\sup_{g_r \in \mathcal{G}_r} |\mathbb{E}[g_r] - \hat{\mathbb{E}}[g_r]| \leq 2\mathbb{E}[R_n(\mathcal{G}_r)] + c(r) \sqrt{\frac{2 \log(1/\delta)}{n}}. \quad (15)$$

The following lemma bounds the Rademacher complexity $R_n(\mathcal{G}_r)$ of the class \mathcal{G}_r .

Lemma 3. $R_n(\mathcal{G}_r) = O\left(\frac{L \log^2(n) \sqrt{r}}{\sqrt{\lambda n}}\right).$

Proof. We prove the lemma using Lemma ?? and the peeling technique (?). First, it is convenient to center g and define the function $\tilde{g} : \tilde{\mathbf{W}} \times \Theta \rightarrow \mathbb{R}$ where $\tilde{\mathbf{W}} = \{\tilde{\mathbf{w}} : (\mathbf{w}^* - \tilde{\mathbf{w}}) \in W\}$ and

$$\tilde{g}(\tilde{\mathbf{w}}, \theta) := f(\mathbf{w}^* - \tilde{\mathbf{w}}, \theta) - f(\mathbf{w}^*, \theta) = g(\mathbf{w}^* - \tilde{\mathbf{w}}, \theta).$$

We also use the notation $\tilde{G} = \{\tilde{g}_{\tilde{\mathbf{w}}}(\theta) = \tilde{g}(\tilde{\mathbf{w}}, \theta) : \tilde{\mathbf{w}} \in \tilde{\mathbf{W}}\}$ and $\tilde{G}_r = \left\{ \frac{\tilde{g}_{\tilde{\mathbf{w}}}}{\gamma^{k(\tilde{\mathbf{w}})}} : \tilde{\mathbf{w}} \in \tilde{\mathbf{W}} \right\}$ where

$$k(\tilde{\mathbf{w}}) = \min\{k' \in \mathbb{Z}_+ : \mathbb{E}[g(\mathbf{w}^* - \tilde{\mathbf{w}}, \cdot)] \leq r\gamma^{k'}\}.$$

Comparing the definitions of \tilde{G}_r and \mathcal{G}_r , we note that these two classes of functions are identical and thus $R_n(\mathcal{G}_r) = R_n(\tilde{G}_r)$. Therefore, from now on we focus on the problem of upper bounding $R_n(\tilde{G}_r)$.

It is convenient to divide the class \tilde{G} into rings according to the expectation. Formally, for any $0 < a < b$ we define

$$\tilde{G}(a, b) = \{\tilde{g}_{\tilde{\mathbf{w}}} : \mathbb{E}[\tilde{g}_{\tilde{\mathbf{w}}}] \in [a, b]\}.$$

We also let

$$\tilde{\mathbf{W}}(a, b) = \{\tilde{\mathbf{w}} \in \tilde{\mathbf{W}} : \mathbb{E}[\tilde{g}_{\tilde{\mathbf{w}}}] \in [a, b]\}.$$

Consider the set $\tilde{\mathbf{W}}(0, r')$ for some arbitrary $r' > 0$. If $\tilde{\mathbf{w}} \in \tilde{\mathbf{W}}(0, r')$ then we have

$$r' \geq \mathbb{E}[\tilde{g}_{\tilde{\mathbf{w}}}] = \mathbb{E}[g_{\mathbf{w}^* - \tilde{\mathbf{w}}}] = F(\mathbf{w}^* - \tilde{\mathbf{w}}) - F(\mathbf{w}^*).$$

Using (7) we know that

$$F(\mathbf{w}^* - \tilde{\mathbf{w}}) - F(\mathbf{w}^*) \geq \frac{\lambda}{2} \|\tilde{\mathbf{w}}\|^2.$$

Therefore, overall we have shown that

$$\forall \tilde{\mathbf{w}} \in \tilde{\mathbf{W}}(0, r'), \quad \|\tilde{\mathbf{w}}\| \leq \sqrt{\frac{2r'}{\lambda}}.$$

This establishes Condition ??, with $B = \sqrt{2r'/\lambda}$, for $\tilde{\mathbf{W}}$. It is straightforward to verify that since the Lipschitz and convexity Condition 1 holds for f it also holds for \tilde{g} (strictly speaking, we should be discussing its convex hull $\text{conv } \tilde{\mathbf{W}}$, on which the norm is of course also bounded, so that convexity inside it can be properly defined). We can therefore apply Lemma ?? and get that there exists a constant C such that

$$R_n(\tilde{G}(0, r')) \leq \frac{C L \sqrt{r'} \log^2(n)}{\sqrt{\lambda n}}. \quad (16)$$

We now get back to the problem of bounding $R_n(\tilde{G}_r)$. To do so, we note that

$$\tilde{G}_r = \cup_{j=0}^{\infty} \left\{ \frac{\tilde{g}_{\tilde{\mathbf{w}}}}{\gamma^j} : \tilde{\mathbf{w}} \in \tilde{\mathbf{W}}, k(\tilde{\mathbf{w}}) = j \right\}.$$

Therefore

$$\begin{aligned} R_n(\tilde{G}_r) &\leq \sum_{j=0}^{\infty} R_n \left(\left\{ \frac{\tilde{g}_{\tilde{\mathbf{w}}}}{\gamma^j} : \tilde{\mathbf{w}} \in \tilde{\mathbf{W}}, k(\tilde{\mathbf{w}}) = j \right\} \right) \\ &\leq \sum_{j=0}^{\infty} \gamma^{-j} R_n(\tilde{G}(0, r \gamma^j)) \\ &\leq \frac{C L \log^2(n) \sqrt{r}}{\sqrt{\lambda n}} \sum_{j=0}^{\infty} \gamma^{-j/2} \end{aligned}$$

Finally, setting $\gamma = 4$ concludes our proof. \square

Equipped with the above We are now ready to prove Theorem 1.

Proof of Theorem 1 contd. Given any $r > 0$, note that for any $g \in \mathcal{G}$ we have for the corresponding $g_r = g/\gamma^k \in \mathcal{G}_r$ such that,

$$\mathbb{E}[g_r] - \hat{\mathbb{E}}[g_r] \leq \sup_{g'_r \in \mathcal{G}_r} (\mathbb{E}[g'_r] - \hat{\mathbb{E}}[g'_r]).$$

and so

$$\mathbb{E}[g] - \hat{\mathbb{E}}[g] \leq \gamma^k \sup_{g'_r \in \mathcal{G}_r} (\mathbb{E}[g'_r] - \hat{\mathbb{E}}[g'_r]).$$

Now we consider two possible cases, the first when $k = 0$ and the second when $k > 0$. If $k = 0$ then

$$\mathbb{E}[g] - \hat{\mathbb{E}}[g] \leq \sup_{g'_r \in \mathcal{G}_r} (\mathbb{E}[g'_r] - \hat{\mathbb{E}}[g'_r]). \quad (17)$$

When $k > 0$ then, since $k = \min\{k' \in \mathbb{Z}_+ : \mathbb{E}[g] \leq r\gamma^{k'}\}$, we have that $\gamma^{k-1}r \leq \mathbb{E}[g]$. Hence,

$$\mathbb{E}[g] - \hat{\mathbb{E}}[g] \leq \frac{\gamma \mathbb{E}[g]}{r} \sup_{g'_r \in \mathcal{G}_r} (\mathbb{E}[g'_r] - \hat{\mathbb{E}}[g'_r]). \quad (18)$$

For a given $K > 1$, we will select an $r > 0$ such that

$$\sup_{g'_r \in \mathcal{G}_r} (\mathbb{E}[g'_r] - \hat{\mathbb{E}}[g'_r]) \leq \frac{r}{\gamma K}. \quad (19)$$

In that case when $k = 0$, from (17) we have that

$$\mathbb{E}[g] - \hat{\mathbb{E}}[g] \leq \frac{r}{\gamma K}.$$

and in the case when $k > 0$ we have using (18) that

$$\mathbb{E}[g] - \hat{\mathbb{E}}[g] \leq \frac{\mathbb{E}[g]}{K}$$

or equivalently

$$\mathbb{E}[g] \leq \frac{K}{K-1} \hat{\mathbb{E}}[g].$$

Combining both the case when $k = 0$ and $k > 0$ we see that

$$\mathbb{E}[g] \leq \frac{K}{K-1} \hat{\mathbb{E}}[g] + \frac{r}{\gamma K}. \quad (20)$$

All we need to do now is to find an r that satisfies the condition given in (19). To do so, we use (15), the upper bound on $c(r)$ given in Lemma 2, and the upper bound on the Rademacher complexity given in Lemma 3. We get that there exists some constant C such that with probability of at least $1 - \delta$ we have

$$\sup_{g_r \in \mathcal{G}_r} |\mathbb{E}[g_r] - \hat{\mathbb{E}}[g_r]| \leq C L \sqrt{\frac{r}{\lambda n}} \left(\log^2(n) + \sqrt{\log(\frac{1}{\delta})} \right)$$

Hence, the condition given in (19) is satisfied if the following holds

$$\frac{r}{\gamma k} \geq C L \sqrt{\frac{r}{\lambda n}} \left(\log^2(n) + \sqrt{\log(\frac{1}{\delta})} \right)$$

or equivalently if

$$r = \frac{(C\gamma KL)^2}{\lambda n} \left(\log^2(n) + \sqrt{\log(\frac{1}{\delta})} \right)^2.$$

Substituting the above value of r into (20) yields

$$\mathbb{E}[g] \leq \frac{K}{K-1} \hat{\mathbb{E}}[g] + \frac{C^2 L^2 \gamma K}{\lambda n} \left(\log^2(n) + \sqrt{\log(\frac{1}{\delta})} \right)^2.$$

Finally, writing $a = 1/(K - 1)$ concludes the proof of the first inequality. For the second inequality we start by noting that for all $g \in \mathcal{G}$,

$$\hat{\mathbb{E}}[g] - \mathbb{E}[g] \leq \gamma^k \sup_{g'_r \in \mathcal{G}_r} (|\hat{\mathbb{E}}[g'_r] - \mathbb{E}[g'_r]|).$$

and then essentially repeat the same steps as above with appropriate changes, □

3 A $1/n$ Rate for the SVM Objective

We consider the problem of L_2 -regularized linear prediction of target labels $y \in \mathcal{Y}$, using vectors \mathbf{x} in a subset \mathcal{X} of a Hilbert space, with $\|\mathbf{x}\| \leq B_x$. Predictions are given by $\langle \mathbf{w}, \mathbf{x} \rangle$, and we shall consider a convex loss function $\ell : \mathbb{R} \times \mathcal{Y} \rightarrow \mathbb{R}$ that is L_ℓ -Lipschitz in its first argument. An example of such a loss function is the hinge loss used in SVMs:

$$\ell(z, y) = [1 - yz]_+ = \max(1 - yz, 0) \quad (21)$$

In this section we aim at analyzing the rate of convergence of the regularized risk, e.g. the SVM objective, rather than the risk directly. To this end, we define the regularized objective function,

$$f(\mathbf{w}; \mathbf{x}, y) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \ell(\langle \mathbf{w}, \mathbf{x} \rangle, y) \quad (22)$$

As before, we use $F(\mathbf{w})$ and $\hat{F}(\mathbf{w})$ to denote the expectation and empirical averages of the objective function (22). When $\ell(\cdot, \cdot)$ is the hinge-loss, the empirical average $\hat{F}(\mathbf{w})$ is precisely the SVM training objective.

Since we take $\ell(z, y)$ to be convex in z , then $\mathbf{w} \mapsto \ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ is also convex in \mathbf{w} for all \mathbf{x}, y . Adding the λ -strongly convex function $\mathbf{w} \mapsto \frac{\lambda}{2} \|\mathbf{w}\|^2$, we see that f is λ -strongly convex w.r.t. \mathbf{w} , establishing Condition 2.

We will also need to assume that the loss is non-negative and that $\ell(0, y)$ is bounded for all y . This later condition is trivially true when the label space \mathcal{Y} is finite, and holds for any reasonable loss function and a bounded label space. To reduce excess notation, we will just assume $\ell(0, y) \leq 1$, as is the case for the hinge-loss. This implies that

$$\mathbb{E} [\ell(\langle \mathbf{w}^*, \mathbf{x} \rangle, y)] + \frac{\lambda}{2} \|\mathbf{w}^*\|^2 = F(\mathbf{w}^*) \leq F(0) \leq 1 \quad (23)$$

and hence we conclude that $\|\mathbf{w}^*\| \leq \sqrt{\frac{2}{\lambda}}$. We will therefore restrict our analysis only to weight vectors whose norm is bounded by $\sqrt{\frac{2}{\lambda}}$. That is, we take:

$$\mathbf{W} = \left\{ \mathbf{w} \mid \|\mathbf{w}\| \leq \sqrt{\frac{2}{\lambda}} \right\} \quad (24)$$

We already saw that $f(\mathbf{w}; \mathbf{x}, y)$ is convex w.r.t. \mathbf{w} . To bound the Lipschitz constant of f , note that $\mathbf{w} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle$ is $\|\mathbf{x}\|$ -Lipschitz, and so $\mathbf{w} \mapsto \ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)$ is $L_\ell \cdot B_x$ -Lipschitz for all $\|\mathbf{x}\| \leq B_x$. To this we must add the Lipschitz constant of $\frac{\lambda}{2} \|\mathbf{w}\|^2$, which is $\frac{\lambda}{2} (2 \sup \|\mathbf{w}\|) \leq \frac{\lambda}{2} (2 \sqrt{\frac{2}{\lambda}}) = \sqrt{2\lambda}$. We therefore establish Condition 1 with

$$L \leq (L_\ell B_x + \sqrt{2\lambda})$$

Having established Conditions 1 and 2, we can apply Theorem 1. We further note that the additional factor of $\log^2(n)$ in Theorem 1 results from bounding the Rademacher complexity using the fat shattering dimension. However, for objective functions of the form (22), we can bound the Rademacher complexity directly using

standard results on the Rademacher complexity of linear predictors, and avoid the spurious log-factors: In equation (16) of Lemma 3, instead of using Lemma ?? we can use standard results on linear predictors to bound $R_n(\tilde{G}(0, r')) \leq 2 \frac{L\sqrt{r'}}{\sqrt{\lambda n}}$. We can therefore conclude:

Theorem 4. *Let f be of the form (22), with ℓ convex, L_ℓ -Lipschitz, non-negative, and $\ell(0, y) \leq 1$, and $\|\mathbf{x}\| \leq B_x$. Then for any $a, \delta > 0$, with probability of at least $1 - \delta$ over the sample, we have that for all \mathbf{w} with $\|\mathbf{w}\| \leq \sqrt{2/\lambda}$:*

$$F(\mathbf{w}) - F(\mathbf{w}^*) \leq (1 + a)(\hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*)) + O\left(\frac{(1 + \frac{1}{a})(L_\ell^2 B_x^2 + \lambda)(\log(1/\delta))}{\lambda n}\right).$$

as well as:

$$\hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*) \leq (1 + a)(F(\mathbf{w}) - F(\mathbf{w}^*)) + O\left(\frac{(1 + \frac{1}{a})(L_\ell^2 B_x^2 + \lambda)(\log(1/\delta))}{\lambda n}\right).$$

Using the above we see that for SVMs (i.e. ℓ is the hinge loss $L_\ell = 1$ and so for $\lambda = O(B_x^2)$ fixed, we can conclude that the generalization regret (in terms of the SVM objective) of the empirical optimizer $\hat{\mathbf{w}}$ converges as $O(\frac{B_x^2 \log(1/\delta)}{\lambda n})$. Studying the convergence rate of the SVM objective allows us to better understand and appreciate analysis of computational optimization approaches for this objective, as well as obtain oracle inequalities on the generalization loss of $\hat{\mathbf{w}}$.

4 Oracle Inequalities for SVMs

In this Section we apply the concentration results of Section 3 to obtain an oracle inequality on the generalization error of an approximate minimizer of the SVM objective. The basic setup is the same as in Section 3. We are now mostly concerned with the generalization error (i.e. expected loss—as we are dealing here only with convex losses, we slightly abuse standard terminology, and refer to the “loss” also as “error”):

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}[\ell(\langle \mathbf{w}, \mathbf{x} \rangle, y)] \quad (25)$$

We assume, as an oracle assumption, that there exists a good predictor \mathbf{w}_o with low norm $\|\mathbf{w}_o\|$ and which attains low generalization error $\mathcal{L}(\mathbf{w}_o)$. We now want to analyze the generalization error $\mathcal{L}(\hat{\mathbf{w}})$ of the empirical minimizer of the SVM objective $\hat{F}(\mathbf{w})$, or perhaps an approximate minimizer of this objective.

For any $\lambda > 0$ and $\mathbf{w} \in W$, we first decompose:

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \mathcal{L}(\mathbf{w}_o) \\ &+ (F(\mathbf{w}) - F(\mathbf{w}^*)) \\ &+ (F(\mathbf{w}^*) - F(\mathbf{w}_o)) \\ &+ \frac{\lambda}{2} \|\mathbf{w}_o\|^2 - \frac{\lambda}{2} \|\mathbf{w}\|^2 \end{aligned} \quad (26)$$

Consider an optimization algorithm for $\hat{F}(\mathbf{w})$ that is guaranteed to find $\tilde{\mathbf{w}}$ such that $\hat{F}(\tilde{\mathbf{w}}) \leq \min \hat{F}(\mathbf{w}) + \epsilon_{\text{opt}}$. Now, using Theorem 4 (ensuring $\lambda = O(B_x^2)$), and setting $a = 1$, we have, with probability at least $1 - \delta$:

$$F(\tilde{\mathbf{w}}) - F(\mathbf{w}^*) \leq 2\epsilon_{\text{opt}} + O\left(\frac{B_x^2 \log(1/\delta)}{\lambda n}\right). \quad (27)$$

Optimizing to within $\epsilon_{\text{opt}} = O(\frac{B_x^2}{\lambda n})$ is then enough to ensure

$$F(\tilde{\mathbf{w}}) - F(\mathbf{w}^*) = O\left(\frac{B_x^2 \log(1/\delta)}{\lambda n}\right). \quad (28)$$

We can use (28) to bound the second term of (26). The optimality of \mathbf{w}^* guarantees us that $F(\mathbf{w}^*) - F(\mathbf{w}_o) \leq 0$ and certainly $\frac{\lambda}{2} \|\tilde{\mathbf{w}}\|^2 \geq 0$. We therefore obtain:

$$\mathcal{L}(\tilde{\mathbf{w}}) \leq \mathcal{L}(\mathbf{w}_o) + \frac{\lambda}{2} \|\mathbf{w}_o\|^2 + O\left(\frac{B_x^2 \log(1/\delta)}{\lambda n}\right) \quad (29)$$

This might seem like a rate of $1/n$ on the generalization error, but we need to choose λ so as to balance the second and third terms. The optimal choice for λ is

$$\lambda_n = O\left(\frac{B_x \sqrt{\log(1/\delta)}}{\|\mathbf{w}_o\| \sqrt{n}}\right). \quad (30)$$

which yields:

$$\mathcal{L}(\tilde{\mathbf{w}}) \leq \mathcal{L}(\mathbf{w}_o) + O\left(\sqrt{\frac{B_x^2 \|\mathbf{w}_o\|^2 \log(1/\delta)}{n}}\right) \quad (31)$$

Even though we rely only on a single fixed \mathbf{w}_o , the choice of λ should change as the sample size increases.

It is interesting to repeat the analysis of this Section using the more standard result:

$$F(\mathbf{w}) - F(\mathbf{w}^*) \leq \hat{F}(\mathbf{w}) - \hat{F}(\mathbf{w}^*) + O\left(\sqrt{\frac{B_x^2}{\lambda n}}\right) \quad (32)$$

where we ignore the dependence on δ , and use $\|\mathbf{w}\| \leq \sqrt{2/\lambda}$. The bound (32) is immediate from concentration results on the unregularized loss, as the regularization terms cancel out. However, using (32) instead of Theorem 4 in our analysis yields the oracle inequality:

$$\mathcal{L}(\tilde{\mathbf{w}}) \leq \mathcal{L}(\mathbf{w}_o) + O\left(\left(\frac{B_x^2 \|\mathbf{w}_o\|^2 \log^2(1/\delta)}{n}\right)^{1/3}\right) \quad (33)$$

The oracle analysis studied here is very simple—our oracle assumption involves only a single predictor \mathbf{w}_o , and we make no assumptions about the kernel or the noise. We note that a more sophisticated analysis has been carried out by Steinwart et al.

(2006). Steinwart et al. shows that rates faster than $1/\sqrt{n}$ are possible under certain conditions on noise and complexity of kernel class. In Steinwart et al.'s analyses the estimation rates (i.e. rates for expected regularized risk) are given in terms of the approximation error quantity $\frac{\lambda}{2}\|\mathbf{w}^*\|^2 + \mathcal{L}(\mathbf{w}^*) - \mathcal{L}^*$ where \mathcal{L}^* is the Bayes risk. In our result we consider the estimation rate for regularized objective independent of the approximation error.

References

Steinwart, I., Hush, D., & Scovel, C. (2006). A new concentration result for regularized risk minimizers. *High-dimensional Probability IV, in IMS Lecture Notes, 51*, 260–275.