

Ch 13: Regularization and Stability

APPM 7400 Theory of Machine Learning, Spring 2020

Stephen Becker

March 16 2020

1 Intro

Intro to regularization and stability Regularized Loss Minimization (RLM) just means ERM plus a regularizer; the regularizer makes it stable to slight changes in input

$$\operatorname{argmin}_{\mathbf{w}} \widehat{L}_S(\mathbf{w}) + R(\mathbf{w}) \quad (\text{RLM})$$

Ideas behind regularization: (1) penalizes complexity (often imperfectly), (2) stabilizes problem

We'll show that if a loss function is (1) convex, (2) Lipschitz or smooth, (3) and bounded \mathcal{H} , then by adding a strongly convex regularizer, we can get PAC learning bounds

We focus on $R(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2$ (write $\|\cdot\|$ for $\|\cdot\|_2$ now)

In particular, **ridge regression** for least-squares

$$\min_{\mathbf{w}} f(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\langle \mathbf{x}_i, \mathbf{w} \rangle - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

Ridge Regression Ridge regression objective is

$$\begin{aligned} f(\mathbf{w}) &\stackrel{\text{def}}{=} \underbrace{\frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\langle \mathbf{x}_i, \mathbf{w} \rangle - y_i)^2}_{\widehat{L}_S(\mathbf{w})} + \underbrace{\lambda \|\mathbf{w}\|^2}_{R(\mathbf{w})} \\ &= \frac{1}{m} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2 \end{aligned}$$

To find solution, we can solve the normal equations (in practice, for large systems and/or ill-conditioned, there are many alternatives, such as SGD, conjugate gradient, etc.). We derive this by solving $\nabla f(\mathbf{w}) = 0$ (in this case, a necessary and sufficient condition for optimality).

$$\begin{aligned} 0 = \nabla f(\mathbf{w}) &= \frac{1}{m} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\lambda \mathbf{w} \\ \implies \underbrace{(\mathbf{X}^\top \mathbf{X})}_A + 2\lambda m I \mathbf{w} &= \underbrace{\mathbf{X}^\top \mathbf{y}}_b \end{aligned} \quad (\text{Normal Eq'n})$$

2 Analysis Setup

Analysis Framework, 1 Recall we've already talked about the traditional **bias-variance** decomposition

$$L_{\mathcal{D}}(h) = \underbrace{\left(L_{\mathcal{D}}(h) - \min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') \right)}_{\text{variance}} + \underbrace{\min_{h' \in \mathcal{H}} L_{\mathcal{D}}(h')}_{\text{bias}}$$

and most of our existing analysis has been controlling the variance, e.g., via uniform convergence to get $|L_{\mathcal{D}}(h) - \hat{L}_S(h)| < \epsilon/2$ and $\hat{L}_S(H)$ small if $h \in \text{ERM}$.

Now, instead of uniform convergence, introduce **average** or **expected risk**

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathbf{A}(S))] \quad \text{instead of} \quad (\forall h) L_{\mathcal{D}}(h) \leq \dots$$

where we are acknowledging that the classifier h (or \mathbf{w}) is chosen by an algorithm \mathbf{A} based on the data S .

Exercise 13.1 shows how expected risk can be used to get an agnostic PAC learning bound.

Analysis Framework, 2 Our goal is a bound like

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathbf{A}(S))] \leq \min_{\mathbf{w} \in \mathcal{H}} L_{\mathcal{D}}(\mathbf{w}) + \epsilon$$

and we'll get there in to parts: just like the bias-variance tradeoff, we'll do

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathbf{A}(S))] = \underbrace{\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(\mathbf{A}(S)) - \hat{L}_S(\mathbf{A}(S))]}_{\text{I}} + \underbrace{\mathbb{E}_{S \sim \mathcal{D}^m} [\hat{L}_S(\mathbf{A}(S))]}_{\text{II}}$$

3 Analysis of I (stability)

Analysis Framework, 2 Our notion of stability is that if we take

$$S = (\mathbf{z}_1, \dots, \mathbf{z}_{i-1}, \mathbf{z}_i, \mathbf{z}_{i+1}, \dots, \mathbf{z}_m) \quad \text{and replace it with} \\ S^{(i)} = (\mathbf{z}_1, \dots, \mathbf{z}_{i-1}, \mathbf{z}', \mathbf{z}_{i+1}, \dots, \mathbf{z}_m)$$

then $\mathbf{A}(S) \approx \mathbf{A}(S^{(i)})$. We'll want

$$\underbrace{0 \leq \ell(\mathbf{A}(S^{(i)}), \mathbf{z}_i) - \ell(\mathbf{A}(S), \mathbf{z}_i)}_{\text{usually}} \leq \underbrace{\epsilon}_{\text{hopefully}}$$

This relates to our error **I** via this theorem:

Theorem 3.1 (Thm. 13.2 in Shalev-Shwartz and Ben-David). *If $S \stackrel{iid}{\sim} \mathcal{D}^m$, $\mathbf{z}' \sim \mathcal{D}$ (independent of S), $i \sim \text{Uniform}([m])$, then \forall algorithms \mathbf{A}*

$$\text{I} \stackrel{\text{def}}{=} \mathbb{E}_S [L_{\mathcal{D}}(\mathbf{A}(S)) - \hat{L}_S(\mathbf{A}(S))] = \mathbb{E}_{\substack{S \\ \mathbf{z}', i}} [\ell(\mathbf{A}(S^{(i)}), \mathbf{z}_i) - \ell(\mathbf{A}(S), \mathbf{z}_i)]$$

Proof. We'll show the left terms on both sides equal, then the right terms. For the left terms,

$$\mathbb{E}_{\substack{S \\ \mathbf{z}', i}} [\ell(\mathbf{A}(S^{(i)}), \mathbf{z}_i)] = \mathbb{E}_{\substack{S \\ \mathbf{z}'}} [\ell(\mathbf{A}(S), \mathbf{z}')] = \mathbb{E}_S [L_{\mathcal{D}}(\mathbf{A}(S))]$$

and similarly for the right terms

$$\mathbb{E}_i [\ell(\mathbf{A}(S), \mathbf{z}_i)] = \mathbb{E}_S \left[\frac{1}{m} \sum_{i=1}^m \ell(\mathbf{A}(S), \mathbf{z}_i) \right] = \mathbb{E}_S [\hat{L}_S(\mathbf{A}(S))]$$

□

Analysis Framework, 4 Informally, if ① is small, the algorithm **A** is **stable**

Formally, say **A** is **(on-average-replacement) stable** with rate $\epsilon(m)$ (non-increasing in m) if $\forall \mathcal{D}, \textcircled{1} \leq \epsilon(m)$.

We'll investigate how to prove an algorithm is stable, using our theorem to characterize ①. We'll assume regularizer is 2λ -strongly convex, e.g., $R(\mathbf{w}) = \lambda \|\mathbf{w}\|^2$.

Recall if f is μ -strongly convex and non-negative, then if $\mathbf{u} \in \operatorname{argmin} f$,

$$(\forall \mathbf{v} \in \mathbb{R}^d) \frac{\mu}{2} \|\mathbf{v} - \mathbf{u}\|^2 \leq f(\mathbf{v}) - f(\mathbf{u}) \leq f(\mathbf{v}) \quad (\text{self-boundedness})$$

Showing stability Write $f_S(\mathbf{w}) = \widehat{L}_S(\mathbf{w}) + \lambda \|\mathbf{w}\|^2$, or just $f(\mathbf{w})$ when S is clear from context. This is 2λ strongly convex. Our algorithm **A** is RLM, so $\boxed{\mathbf{u} = \mathbf{A}(S)} \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$

Similarly, define $\boxed{\mathbf{v} = \mathbf{A}(S^{(i)})}$.

$$\begin{aligned} f(\mathbf{v}) - f(\mathbf{u}) &\stackrel{\text{def}}{=} \widehat{L}_S(\mathbf{v}) + \lambda \|\mathbf{v}\|^2 - (\widehat{L}_S(\mathbf{u}) + \lambda \|\mathbf{u}\|^2) \\ &= \underbrace{\widehat{L}_{S^{(i)}}(\mathbf{v}) + \lambda \|\mathbf{v}\|^2}_{\textcircled{a}} - \underbrace{\widehat{L}_{S^{(i)}}(\mathbf{u}) + \lambda \|\mathbf{u}\|^2}_{\textcircled{b}} + \frac{1}{m} (\ell(\mathbf{v}, \mathbf{z}_i) - \ell(\mathbf{u}, \mathbf{z}_i)) + \frac{1}{m} (\ell(\mathbf{u}, \mathbf{z}') - \ell(\mathbf{v}, \mathbf{z}')) \end{aligned}$$

Because we chose \mathbf{v} as above, it minimizes $\widehat{L}_{S^{(i)}}(\mathbf{v}) + \lambda \|\mathbf{v}\|^2$, so $\textcircled{a} \leq \textcircled{b}$, hence

$$f(\mathbf{v}) - f(\mathbf{u}) \leq \frac{1}{m} (\ell(\mathbf{v}, \mathbf{z}_i) - \ell(\mathbf{u}, \mathbf{z}_i)) + \frac{1}{m} (\ell(\mathbf{u}, \mathbf{z}') - \ell(\mathbf{v}, \mathbf{z}'))$$

By self-boundedness, $f(\mathbf{v}) - f(\mathbf{u}) \geq \lambda \|\mathbf{v} - \mathbf{u}\|^2$, so combining this with above,

$$\lambda \|\mathbf{v} - \mathbf{u}\|^2 \leq \frac{1}{m} \left(\ell(\mathbf{v}, \mathbf{z}_i) - \ell(\mathbf{u}, \mathbf{z}_i) \right) + \frac{1}{m} \left(\ell(\mathbf{u}, \mathbf{z}') - \ell(\mathbf{v}, \mathbf{z}') \right) \quad (1)$$

From here, there are two ways to proceed:

Case 1: assuming $(\forall \mathbf{z}), \mathbf{w} \mapsto \ell(\mathbf{w}, \mathbf{z})$ is ρ -Lipschitz

So, directly from Lipschitz property,

$$\begin{aligned} \ell(\mathbf{v}, \mathbf{z}_i) - \ell(\mathbf{u}, \mathbf{z}_i) &\leq \rho \|\mathbf{v} - \mathbf{u}\| \\ \ell(\mathbf{u}, \mathbf{z}') - \ell(\mathbf{v}, \mathbf{z}') &\leq \rho \|\mathbf{v} - \mathbf{u}\| \end{aligned} \quad (2)$$

so substitute this into Eq. (1) gives

$$\lambda \|\mathbf{v} - \mathbf{u}\|^2 \leq \frac{1}{m} \rho \|\mathbf{v} - \mathbf{u}\| + \frac{1}{m} \rho \|\mathbf{v} - \mathbf{u}\|$$

and either $\mathbf{v} = \mathbf{u}$ or $\|\mathbf{v} - \mathbf{u}\| > 0$ and then we can divide by it; either way,

$$\|\mathbf{v} - \mathbf{u}\| \leq \frac{2\rho}{\lambda m}$$

and put this back into Eq. (2) to get

$$\ell(\underbrace{\mathbf{v}}_{\mathbf{A}(S^{(i)})}, \mathbf{z}_i) - \ell(\underbrace{\mathbf{u}}_{\mathbf{A}(S)}, \mathbf{z}_i) \leq \rho \frac{2\rho}{\lambda m}$$

and thus by Thm. 3.1

$$\textcircled{1} = \mathbb{E}_{\substack{S \\ \mathbf{z}', i}} \left[\ell(\mathbf{A}(S^{(i)}), \mathbf{z}_i) - \ell(\mathbf{A}(S), \mathbf{z}_i) \right] \leq \frac{2\rho^2}{\lambda m} \stackrel{\text{def}}{=} \epsilon(m)$$

leading to **Corollary 13.6** which states that if ℓ is uniformly ρ -Lipschitz and strongly convex with parameter $\mu > 0$ then it is (on-average-replace-one) **stable** with rate $\epsilon(m) = \frac{4\rho^2}{\mu m}$.

Note that we did not need to assume \mathbf{x} or \mathbf{w} was bounded.

Case 2: assuming $(\forall \mathbf{z}), \mathbf{w} \mapsto \ell(\mathbf{w}, \mathbf{z})$ is β -smooth

(And assume ℓ is non-negative, but we already made that assumption; of course, all we really need is that it is bounded below, with a known bound, since there is nothing special about 0).

When $\nabla \ell$ is β -Lipschitz, we have

$$(\forall \mathbf{w})(\forall \mathbf{z}) \|\nabla \ell(\mathbf{w}, \mathbf{z})\|^2 \leq 2\beta (\ell(\mathbf{w}, \mathbf{z}) - \ell(\mathbf{w}^*, \mathbf{z})) \leq 2\beta \ell(\mathbf{w}, \mathbf{z}) \quad (3)$$

where $\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{z})$ and the 2nd inequality follows by assuming non-negativity. Also, using the quadratic upper bound property of strongly smooth functions,

$$\begin{aligned} \ell(\underbrace{\mathbf{v}}_{\mathbf{A}(S^{(i)})}, \mathbf{z}_i) - \ell(\underbrace{\mathbf{u}}_{\mathbf{A}(S)}, \mathbf{z}_i) &\leq \langle \nabla \ell(\mathbf{u}, \mathbf{z}_i), \mathbf{v} - \mathbf{u} \rangle + \frac{\beta}{2} \|\mathbf{v} - \mathbf{u}\|^2 \\ &\leq \|\nabla \ell(\mathbf{u}, \mathbf{z}_i)\| \cdot \|\mathbf{v} - \mathbf{u}\| + \frac{\beta}{2} \|\mathbf{v} - \mathbf{u}\|^2 \quad (\text{Cauchy-Schwarz}) \\ &\leq \sqrt{2\beta \ell(\mathbf{u}, \mathbf{z})} \cdot \|\mathbf{v} - \mathbf{u}\| + \frac{\beta}{2} \|\mathbf{v} - \mathbf{u}\|^2 \quad \text{via Eq. (3)} \end{aligned} \quad (4)$$

and an analogous result holds for $\ell(\mathbf{v}, \mathbf{z}') - \ell(\mathbf{u}, \mathbf{z}')$. Plug these results into Eq. (1) and divide by $\lambda \|\mathbf{v} - \mathbf{u}\|$ and re-arrange to get

$$\|\mathbf{v} - \mathbf{u}\| \leq \frac{\sqrt{2\beta}}{\lambda m - \beta} \left(\sqrt{\ell(\mathbf{u}, \mathbf{z}_i)} + \sqrt{\ell(\mathbf{v}, \mathbf{z}')} \right)$$

We can choose the value of λ , so pick it such that $\beta \leq \frac{\lambda m}{2}$ so then

$$\|\mathbf{v} - \mathbf{u}\| \leq \frac{\sqrt{8\beta}}{\lambda m} \left(\sqrt{\ell(\mathbf{u}, \mathbf{z}_i)} + \sqrt{\ell(\mathbf{v}, \mathbf{z}')} \right)$$

Now, as before, we go back to an earlier bound: plug the above eq into Eq. (4) to get (skipping a few steps)

$$\begin{aligned} \ell(\underbrace{\mathbf{v}}_{\mathbf{A}(S^{(i)})}, \mathbf{z}_i) - \ell(\underbrace{\mathbf{u}}_{\mathbf{A}(S)}, \mathbf{z}_i) &\leq \left(\frac{4\beta}{\lambda m} + \frac{8\beta^2}{(\lambda m)^2} \right) \left(\sqrt{\ell(\mathbf{u}, \mathbf{z}_i)} + \sqrt{\ell(\mathbf{v}, \mathbf{z}')} \right)^2 \quad \text{and bound } \frac{8\beta^2}{(\lambda m)^2} \geq 0 \\ &\leq \frac{24\beta}{\lambda m} \left(\ell(\mathbf{u}, \mathbf{z}_i) + \ell(\mathbf{v}, \mathbf{z}') \right)^2 \quad \text{since } (a+b)^2 \leq 3(a^2 + b^2) \end{aligned}$$

thus via Thm. 3.1

$$\begin{aligned} \textcircled{1} &= \mathbb{E}_{\substack{S \\ \mathbf{z}', i}} \left[\ell(\mathbf{A}(S^{(i)}), \mathbf{z}_i) - \ell(\mathbf{A}(S), \mathbf{z}_i) \right] \leq \frac{24\beta}{\lambda m} \mathbb{E}_{\substack{S \\ \mathbf{z}', i}} [\ell(\mathbf{u}, \mathbf{z}_i) + \ell(\mathbf{v}, \mathbf{z}')] \\ &= \frac{48\beta}{\lambda m} \mathbb{E}_S \hat{L}_S(\mathbf{A}(S)) \end{aligned}$$

and typically the loss function is bounded for all \mathbf{z} , e.g., $\ell(0, \mathbf{z}) \leq c$, hence $\mathbb{E}_S \hat{L}_S(\mathbf{A}(S))$, so in this case, we have **Corollary 13.7** which is that if ℓ is uniformly β -smooth and strongly convex with parameter $\mu > 0$ then it is (on-average-replace-one)stable with rate $\epsilon(m) = \frac{96\beta c}{\mu m}$.

4 Analysis of II (bias/underfitting)