

Lecture 17.

Logistics [Project Schedule Out
HW3 due Thursday

Recap. [Lipschitzness]. a function $f(\cdot)$ is β -Lipschitz if

$$\|f(w_1) - f(w_2)\| \leq \beta \|w_1 - w_2\|$$

measures how fast the function changes over the input.

Chapter 13.

We will learn a new learning paradigm: **Regularized Loss Minimization (RLM)** and show that **Convex-Lipschitz-Bounded** and **Convex-smooth-Bounded** families of learning problems are learnable.

Key insight: **regularization** makes learning algorithms more **stable**

Definition: regularized loss minimization

For a regularizer $R: \mathbb{R}^d \rightarrow \mathbb{R}$, RLM outputs

$$A(S) \leftarrow \arg \min_{w \in \mathbb{R}^d} (L_S(w) + R(w))$$

- the regularizer $R(w)$ is a **measure of complexity** of hypothesis w
- we use $R(w) = \lambda \cdot \|w\|_2^2$, also called Tikhonov Regularizer.

Example: Ridge Regression from CSE 446/546

$$A(S) \leftarrow \arg \min_{w \in \mathbb{R}^d} \left(\underbrace{\frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\langle w, x_i \rangle - y_i)^2}_{\text{Quadratic loss on linear regression}} + \underbrace{\lambda \|w\|_2^2}_{\text{regularizer}} \right)$$

- this has a closed form solution, which we get by setting the gradient to zero.

Our goal today is to show the following

[Theorem] D is a distribution over $\mathcal{X} \times [-1, 1]$, where

$$\{x \in \mathbb{R}^d : \|x\| \leq 1\}$$

$\mathcal{H} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$, loss is ρ -Lipschitz.

For any $\varepsilon \in (0, 1)$, if $m \geq \frac{8\rho^2 B^2}{\varepsilon^2}$ then

Ridge Regression with $\lambda = \sqrt{\frac{2\rho^2}{B^2 m}}$ satisfies

$$\mathbb{E}_S[\underbrace{L_D(A(S))}_{\text{RLM}}] \leq \min_{w \in \mathcal{H}_B} L_D(w) + \varepsilon$$

\uparrow Random
 \uparrow

Remarks:

- both \mathcal{X} and \mathcal{H} are **bounded**, which is important.
- relatedly, **regularization** is important.
- we are bounding the **expectation**, and not the $L_D(A(S))$ with high probability, because expectation is directly related to **stability**.

Definition:

A learning algorithm is **stable** if "a small change in the input" does not "change the output much". We make this formal

- small change in input:

$$S = (z_1, z_2, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_m)$$

$$S^{(i)} = (z_1, z_2, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m)$$

- measure of effect: loss on $z_i = (x_i, y_i)$, i.e.,

$$l(A(S^{(i)}), z_i) - l(A(S), z_i)$$

* intuitively, this is non-negative, because an algorithm that saw z_i in training will likely have less error on z_i .

* A **stable** algorithm will make the above difference small.

Definition [on-average-replace-one-stable]

We say an algorithm A is on-average-replace-one-stable

with rate $\epsilon(m)$ if $\mathbb{E}_{S, Z', i} [l(A(S^{(i)}), z_i) - l(A(S), z_i)] \leq \epsilon(m)$

$\sum_{S \in D^m}$ $\sum_{Z \in D}$ $\sum_{U \in U[m]}$ *Uniform*

for all D , for some monotonically decreasing $\epsilon(\cdot)$.

* this notion of *stability* is further justified by the following, which shows stable algorithms generalize.

[Theorem 13.2] $S = (z_1, \dots, z_m)$ i.i.d. from D , and Z' is another iid sample, let $U[m]$ be uniform distribution over $[m]$. then for any algorithm A ,

$$\underbrace{\mathbb{E}_S [L_D(A(S)) - L_S(A(S))]}_{\text{generalization error PAC.}} = \mathbb{E}_{S, Z', i} [l(A(S^{(i)}), z_i) - l(A(S), z_i)]$$

bounded by one-replace-stable $\epsilon(m)$

Proof \rangle by definition,

$$\text{2nd term} = \mathbb{E}_S [L_S(A(S))] = \mathbb{E}_S \left[\frac{1}{m} \sum_{i=1}^m l(A(S), z_i) \right]$$

(z_1, \dots, z_m)

* this is indep z_i !

* note that this holds

even if A treats (z_1, \dots, z_m) differently, for example train on every even index data

$$\text{1st term} = \mathbb{E}_S [L_D(A(S))] = \mathbb{E}_{S, Z'} [l(A(S), z'_i)]$$

$$= \mathbb{E}_{S, Z'} [l(A(\{z_1, \dots, z_{i-1}, \underbrace{z_i}_{\text{indep.}}, z_{i+1}, \dots, z_m\}), z'_i)]$$

$$= \mathbb{E}_{z_1, \dots, z_m, z'_i} [l(A(S^{(i)}), z'_i)]$$

* Assuming Convex loss function with Lipschitzness or Smoothness,
we show that RLM is stable, because it is strongly convex.

Definition [Strong Convexity]

A function is λ -strongly convex if $\forall u, v$ & $\alpha \in [0, 1]$,

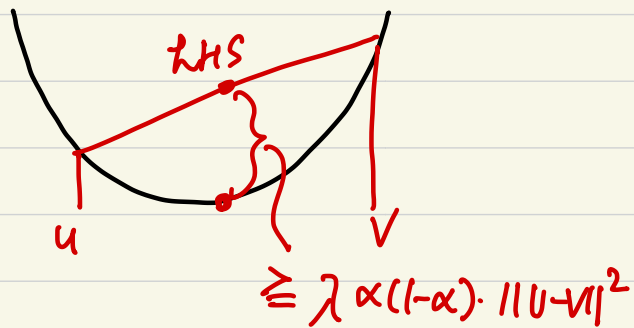
$$\alpha f(u) + (1-\alpha) f(v) \geq f(\alpha u + (1-\alpha)v) + \frac{\lambda}{2} \alpha(1-\alpha) \cdot \|u-v\|^2$$

example \succ

piecewise linear,
0-strongly convex



$f(u) = \lambda \cdot \|u\|^2$ is 2λ -strongly convex,



• Smallest eigenvalue of Hessian.

* in particular, we use the fact that

① f is convex & g is λ -strongly convex $\rightarrow f+g$ is λ -strongly convex.

[Proposition]

② if f is λ -strongly convex, and w^* is minimizer of f .

$$f(w) - f(w^*) \geq \frac{\lambda}{2} \|w - w^*\|^2$$

We assume $\begin{cases} l(w, z) \text{ is Convex in } w \in \mathbb{R}^d \\ R(w) = \lambda \|w\|^2, \text{ which is } 2\lambda\text{-strongly Convex.} \end{cases}$

Recall that RLM outputs,

$$A(S) \leftarrow \arg \min_w [f_S(w)] \triangleq L_S(w) + \lambda \|w\|^2$$

by ①, $f_S(w)$ is 2λ -strongly Convex

let $\hat{w}^{(i)} = A(S^{(i)})$, and $\hat{w} = A(S)$, then

Lower Bound: by ② and 2λ -strong Convexity, and the fact that

\hat{w} is minimizer of $f_S(w)$,

$$f_S(\hat{w}^{(i)}) - f_S(\hat{w}) \geq \lambda \cdot \|\hat{w}^{(i)} - \hat{w}\|_2^2 \quad \begin{matrix} \geq 0 \\ \uparrow \\ \text{Convexity.} \end{matrix}$$

Upper Bound:

$$f_S(\hat{w}^{(i)}) - f_S(\hat{w}) = \underbrace{L_S(\hat{w}^{(i)}) + \lambda \|\hat{w}^{(i)}\|^2}_{L_{S^{(i)}}(\hat{w}^{(i)}) + \frac{1}{m} (l(\hat{w}^{(i)}, z_i) - l(\hat{w}^{(i)}, z_i'))} - \underbrace{L_S(\hat{w}) + \lambda \|\hat{w}\|^2}_{L_{S^{(i)}}(\hat{w}) + \frac{1}{m} (l(\hat{w}, z_i) - l(\hat{w}, z_i'))}$$

$$= \underbrace{\left(L_{S^{(i)}}(\hat{w}^{(i)}) + \lambda \|\hat{w}^{(i)}\|^2 \right)}_{f_{S^{(i)}}(\hat{w}^{(i)})} - \underbrace{\left(L_{S^{(i)}}(\hat{w}) + \lambda \|\hat{w}\|^2 \right)}_{f_{S^{(i)}}(\hat{w})} + \frac{l(\hat{w}^{(i)}, z_i) - l(\hat{w}, z_i)}{m} + \frac{l(\hat{w}, z_i') - l(\hat{w}^{(i)}, z_i')}{m}$$

$\underbrace{\hspace{10em}}_{\substack{\uparrow \\ \text{minimizer}}}$

≤ 0

$$\leq \frac{l(\hat{w}^{(i)}, z_i) - l(\hat{w}, z_i)}{m} + \frac{l(\hat{w}, z_i') - l(\hat{w}^{(i)}, z_i')}{m}$$

How fast l changes as we change $\hat{w}^{(i)}$, related to Lipschitzness of l .

* Putting the Upper Bound and Lower Bound together, we get

$$\lambda \|\hat{w}^{(n)} - \bar{w}\|^2 \leq \frac{l(\hat{w}^{(n)}, \mathcal{Z}_n) - l(\bar{w}, \mathcal{Z}_n)}{m} + \frac{l(\bar{w}, \mathcal{Z}^n) - l(\hat{w}^{(n)}, \mathcal{Z}^n)}{m} \quad (*)$$

This holds for any convex RLM.

We use this to show that if the loss $l(\cdot, \mathcal{Z}_i)$ is well-behaved $\left\{ \begin{array}{l} \text{Lipschitz} \\ \text{smooth} \end{array} \right.$ then RLM is stable.

Case 1: β -Lipschitz loss $l(\cdot, \mathcal{Z}_i) \rightarrow l(w_1, \mathcal{Z}_i) - l(w_2, \mathcal{Z}_i) \leq \beta \|w_1 - w_2\|$

$$l(\hat{w}^{(n)}, \mathcal{Z}_i) - l(\bar{w}, \mathcal{Z}_i) \leq \beta \|\hat{w}^{(n)} - \bar{w}\|$$

$$l(\bar{w}, \mathcal{Z}_i) - l(\hat{w}^{(n)}, \mathcal{Z}_i) \leq \beta \|\bar{w} - \hat{w}^{(n)}\|$$

by (*) $\lambda \cdot \|\hat{w}^{(n)} - \bar{w}\|^2 \leq \frac{2\beta}{m} \|\hat{w}^{(n)} - \bar{w}\|$

$\rightarrow \|\hat{w}^{(n)} - \bar{w}\| \leq \frac{2\beta}{\lambda \cdot m}$

Recall Stability is a bound on

Stability: $\mathbb{E} [l(\hat{w}^{(n)}, \mathcal{Z}_n) - l(\bar{w}, \mathcal{Z}_n)] \leq \frac{2\beta^2}{\lambda m}$

$\underbrace{\hspace{10em}}_{\leq \frac{2\beta^2}{\lambda m}}$

\nearrow more Lipschitz, stochaste loss.
 \searrow more sample, less influence by 1
 \nearrow $\lambda \uparrow$, more strongly convex, more stable

Thm 13.2.

[Corollary 13.6] For β -Lipschitz loss function, RLM with $\lambda \|w\|^2$ regularization achieves

Generalization: $\mathbb{E} [L_{\mathcal{D}}(\text{AveS}) - L_S(\text{AveS})] \leq \frac{2\beta^2}{\lambda m}$

Case 2. Similarly, for β -smooth loss and $\lambda \geq \frac{2\beta}{m}$,

$$\mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \leq \frac{48\beta}{\lambda m} \mathbb{E}[L_S(A(S))]$$

* Fitting - Stability Tradeoff

$$\underbrace{\mathbb{E}_S [L_D(A(S))]}_{\substack{\text{what we care about,} \\ \text{Test loss}}} = \underbrace{\mathbb{E}_S [L_S(A(S))]}_{\substack{\text{train loss} \\ \text{how well you fit} \\ \text{Train data}}} + \underbrace{\left\{ \mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \right\}}_{\substack{\approx \text{Stability: Corollary 13.6} \\ \leq \frac{2\beta^2}{\lambda m}}}$$

Let's expand:

$$\begin{aligned} &\leq \mathbb{E}_S [L_S(A(S)) + \lambda \|A(S)\|^2] \\ &\quad \underbrace{\hspace{10em}}_{\text{minimizer RLM}} \\ &\leq \mathbb{E}_S [L_S(w^*) + \lambda \|w^*\|^2] \\ &\leq \underbrace{L_D(w^*)}_{\text{irreducible error}} + \underbrace{\lambda \|w^*\|^2}_{\substack{\text{misfit} \\ \text{due to reg.}}} + \underbrace{\frac{2\beta^2}{\lambda m}}_{\substack{\text{stability} \\ \text{due to reg.}}} \\ &\quad \underbrace{\hspace{15em}}_{\text{covered by } \lambda.} \end{aligned}$$

[Corollary 13.9]. (Convex-Lipschitz-Bounded Risk bound)

If l is β -Lipschitz, $\|w^*\| \leq B$, then $\lambda = \sqrt{\frac{2\beta^2}{B^2 m}}$ achieves

$$\mathbb{E}_D (L_D(A(S))) \leq \underbrace{\min_{w \in H} L_D(w)}_{\substack{\uparrow \\ \text{RLM}}} + \beta \cdot B \cdot \sqrt{\frac{8}{m}}$$

this proves the [Theorem].