

Chapter 12. Convex Learning Problems.

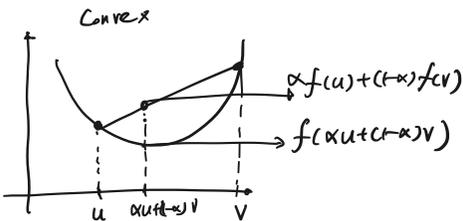
Definition [Convex Set] A set C in a vector space is convex if for any two vectors u, v in C , the line segment between u and v is contained in C . That is, for any $\alpha \in [0, 1]$, we have

$$\alpha \cdot u + (1-\alpha) \cdot v \in C, \quad \forall u, v \in C$$



Definition [Convex function] Let C be a convex set. A function $f: C \rightarrow \mathbb{R}$ is convex if for every $u, v \in C$, and $\alpha \in [0, 1]$

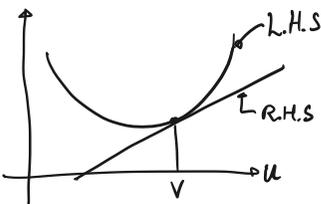
$$f(\alpha u + (1-\alpha)v) \leq \alpha f(u) + (1-\alpha)f(v).$$



• we will be using equivalent definitions of convexity:

- for convex differentiable f , $\forall u$

$$f(u) \geq f(v) + \langle \nabla f(v), u-v \rangle$$



gradient $\nabla f(v) \in \mathbb{R}^d, \forall v \in \mathbb{R}^d$
 $\nabla f(v)_i = \frac{\partial f(v)}{\partial v_i}$
 inner product $\langle a, b \rangle = \sum_i a_i \cdot b_i$

Definition [Lipschitzness] Let $C \subset \mathbb{R}^d$. A function $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ is β -Lipschitz over C if for every $w_1, w_2 \in C$ we have

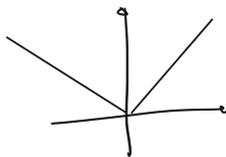
$$\|f(w_1) - f(w_2)\| \leq \beta \|w_1 - w_2\|$$

- Lipschitz constant β measures how fast the function is changing.
- it follows that if $\|\nabla f(w)\|$ is bounded, then $f(\cdot)$ is Lipschitz.

ex) $f(x) = |x|$ is 1-Lipschitz.

$$\begin{aligned} |x_1| - |x_2| &= |x_1 - x_2 + x_2| - |x_2| \\ &\leq |x_1 - x_2| + |x_2| - |x_2| \\ &= |x_1 - x_2|. \end{aligned}$$

by symmetry $||x_1| - |x_2|| \leq |x_1 - x_2|$



ex) $f(x) = x^2$ is not Lipschitz
 let $x_1 = 0, x_2 = 1+p, p > 0$
 $(1+p)^2 - 0^2 > p^2 + p = p \cdot (1+p) = p |x_2 - x_1|$
 taking p as large as we want, we can make this arbitrarily large



Definition [Smoothness] A differentiable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is β -smooth if its gradient is β -Lipschitz; $\forall v, w$ we have

$$\| \nabla f(v) - \nabla f(w) \| \leq \beta \cdot \|v - w\|$$

$$\nabla f(v) = \underbrace{\left(\frac{\partial f(v)}{\partial v_1}, \frac{\partial f(v)}{\partial v_2}, \dots, \frac{\partial f(v)}{\partial v_d} \right)}_{\mathbb{R}^d}$$

- this measures how smoothly the gradient is changing.
- we will be using the fact that, for smooth $f(\cdot)$,

$$f(v) \leq f(w) + \langle \nabla f(w), v - w \rangle + \frac{\beta}{2} \|v - w\|^2$$

ex) $f(x) = x^2$ is 2-smooth
 $f'(x) = 2x$ is 2-Lipschitz

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 is that regularizers make learning algorithm more stable.

Regularized Loss Minimization

for a regularizer is $R: \mathbb{R}^d \rightarrow \mathbb{R}$, the RLM outputs

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

the value of the regularization measures the complexity of the hypothesis, w

we use $R(w) = \lambda \cdot \|w\|^2$, also called Tikhonov regularization

example) Ridge Regression from CSE 446/546

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

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

our goal today is to show the following

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

$\mathcal{X} = \{x \in \mathbb{R}^d : \|x\| \leq 1\}$ and $\mathcal{H} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$, β -Lipschitz loss,

For any $\epsilon \in (0, 1)$, if $m \geq \frac{8\beta^2 B^2}{\epsilon^2}$ then ridge regression

with $\lambda = \sqrt{\frac{2\beta^2}{\delta m}}$ satisfies

$$\mathbb{E}_S [L_{\mathcal{D}}(A(S))] \leq \underbrace{\min_{w \in \mathcal{H}} L_{\mathcal{D}}(w)}_{\text{RLM}} + \epsilon.$$

Remarks:

- Both \mathcal{X} and \mathcal{H} are bounded, and we will learn that this is important
- relatedly, regularization is important
- we are proving a bound on Expected test error (or Risk) as opposed to the usual high probability bound on the test error. Note that test error is random because hypothesis $A(S)$ is random.
- Bounded expected risk implies agnostic PAC learnability, which we will not cover in this lecture.
- we use expected Risk because it is related to stability.

A learning algorithm is stable if "a small change in the input" does not

"change the output much!" we make this formal.

- $S = (z_1, z_2, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n)$
- $S^{(i)} = (z_1, z_2, \dots, z_{i-1}, z', z_{i+1}, \dots, z_n)$ "small change in the input"
- measure of the effect: loss on $z_i = C(x_i, y_i)$

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

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

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

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

this is further justified by the following, which shows stable algorithms generalize.

[Theorem 13.2] $S = (z_1, \dots, z_n)$ iid from D and z' is another iid sample.

let $U(m)$ be the uniform distribution over $[m]$. Then for any Algorithm, A ,

$$\mathbb{E}_S [L_D(A(S)) - L_S(A(S))] = \mathbb{E}_{S, z', i \sim U(m)} [l(A(S^{(i)}), z_i) - l(A(S), z_i)]$$

proof by definition

$$\mathbb{E}_S [L_S(A(S))] = \mathbb{E}_{S, i} [l(A(S), z_i)] \text{ , and}$$

$$\begin{aligned} \mathbb{E}_S [L_D(A(S))] &= \mathbb{E}_{S, z'} [l(A(S), z')] \\ &= \mathbb{E}_{S, z'} [l(A(S^{(i)}), z_i)] \text{ for all } i. \end{aligned}$$

independent

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

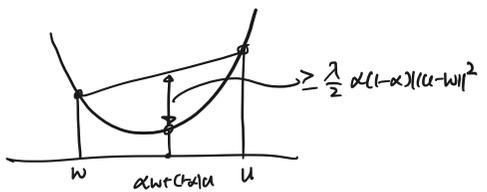
Definition [strong convexity]

A function is λ -strongly convex if for all w, u and $\alpha \in (0, 1)$ we have

$$f(\alpha w + (1-\alpha)u) \leq \alpha f(w) + (1-\alpha)f(u) - \frac{\lambda}{2} \alpha(1-\alpha) \|w-u\|^2$$

how large can this λ be.

example $f(w) = \lambda \|w\|^2$ is 2λ -strongly convex



in particular, we use the fact that

- f is convex and g is λ -strongly convex $\rightarrow f+g$ is λ strongly convex
- if f is λ -strongly convex, then

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

Recall, for RLM

$$A(S) \leftarrow \arg \min_w (\underbrace{L_S(w)}_{f_S(w)} + \lambda \|w\|^2)$$

$f_S(w)$ - 2λ -strongly convex.

lower bound: $f_S(A(S^{(i)})) - f_S(A(S)) \geq \lambda \|A(S^{(i)}) - A(S)\|^2 \leftarrow (**)$

upper bound:

$$f_S(A(S^{(t)})) - f_S(A(S)) = L_S(\hat{w}^{(t)}) + \lambda \|\hat{w}^{(t)}\|^2 - L_S(\hat{w}) - \lambda \|\hat{w}\|^2$$

$$= \underbrace{L_S(\hat{w}^{(t)}) + \lambda \|\hat{w}^{(t)}\|^2 - (L_S(\hat{w}) + \lambda \|\hat{w}\|^2)}_{\leq 0 \text{ because } \hat{w}^{(t)} \text{ is optimal solution for } f_S(\cdot)} + \frac{\ell(\hat{w}^{(t)}, \mathcal{Z}_1) - \ell(\hat{w}, \mathcal{Z}_1)}{m} + \frac{\ell(\hat{w}, \mathcal{Z}') - \ell(\hat{w}^{(t)}, \mathcal{Z}')}{m}$$

Together,

$$\lambda \|\hat{w}^{(t)} - \hat{w}\|^2 \leq \frac{\ell(\hat{w}^{(t)}, \mathcal{Z}_1) - \ell(\hat{w}, \mathcal{Z}_1)}{m} - \frac{\ell(\hat{w}, \mathcal{Z}') - \ell(\hat{w}^{(t)}, \mathcal{Z}')}{m} \quad (**)$$

we use (**) to show that RLM is stable.

Case 1. Lipschitz loss

If $\ell(\cdot, \mathcal{Z}_i)$ is β -Lipschitz, then

$$\ell(\hat{w}^{(t)}, \mathcal{Z}_1) - \ell(\hat{w}, \mathcal{Z}_1) \leq \beta \cdot \|\hat{w}^{(t)} - \hat{w}\|, \text{ and } (***)$$

$$\ell(\hat{w}, \mathcal{Z}') - \ell(\hat{w}^{(t)}, \mathcal{Z}') \leq \beta \cdot \|\hat{w}^{(t)} - \hat{w}\|, \text{ plugging into (**),}$$

$$\lambda \|\hat{w}^{(t)} - \hat{w}\|^2 \leq \frac{2\beta}{m} \|\hat{w}^{(t)} - \hat{w}\|$$

$$\rightarrow \|\hat{w}^{(t)} - \hat{w}\| \leq \frac{2\beta}{\lambda m}, \text{ plugging into (***)}$$

$$\ell(\hat{w}^{(t)}, \mathcal{Z}_1) - \ell(\hat{w}, \mathcal{Z}_1) \leq \frac{2\beta^2}{\lambda m}. \quad \leftarrow \text{which is monotonically decreasing in } m.$$

This implies

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

$$\mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \leq \frac{2\beta^2}{\lambda m}.$$

$\lambda \uparrow \rightarrow$ more strongly convex \rightarrow better generalization

$\beta \downarrow \rightarrow$ more smooth loss \rightarrow better generalization.

we skip the proof for smooth losses, but for β -smooth losses,

$$\text{for } \lambda \geq \frac{2\beta}{m}, \quad \mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \leq \frac{48\beta}{\lambda m} \cdot \mathbb{E}[L_S(A(S))]$$

* Fitting - Stability Tradeoff

$$\mathbb{E}_S [L_D(A(S))] = \underbrace{\mathbb{E}_S [L_S(A(S))]}_{\text{fit to train data}} + \underbrace{\left\{ \mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \right\}}_{\approx \text{stability: Corollary 13.6}}$$

$$\leq \mathbb{E}_S [L_S(A(S)) + \lambda \|\hat{w}\|^2] \leq \frac{2\beta^2}{\lambda m}$$

$$\begin{aligned} \text{A(S) is optimal} &\rightarrow \leq \mathbb{E}_S [L_S(w^*) + \lambda \|w^*\|^2] \\ \text{Solution to RLM} &= L_D(w^*) + \lambda \|w^*\|^2, \quad \forall w^* \end{aligned}$$

$$\text{This gives } \mathbb{E}_S [L_D(A(S))] \leq L_D(w^*) + \lambda \|w^*\|^2 + \frac{2\beta^2}{\lambda m}$$

Trade off by λ .

[Corollary 13.9] if ℓ is β -Lipschitz, and $\|w^*\| \leq B$, then

$$\lambda = \sqrt{\frac{2\beta^2}{B^2 m}} \text{ achieves}$$

$$\mathbb{E}_D (L_D(A(S))) \leq \min_{w \in \mathcal{H}} L_D(w) + \beta B \sqrt{\frac{8}{m}}$$

RLM

This implies the [Theorem].