

1. Project proposal due to day on grade scope

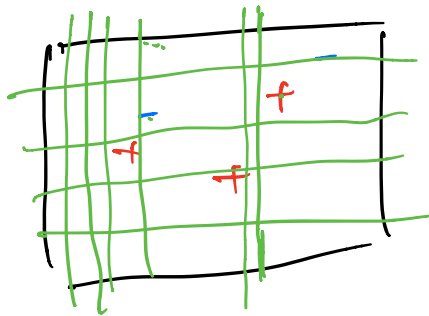
2. HW. 2 Problem 4

$n = 1, 2$ in hint

HW. 2

Problem 1

$S=1$
 $i=1$



$$B = \left\{ \begin{array}{l} s \cdot (2I(x_i \geq t) - 1) : \\ s \in \{\pm 1\} \\ i \in \{1, \dots, d\} \\ t \in \mathbb{R} \end{array} \right\}$$

$t \in \{x_i : (x_i, y) \text{ in training set}\}$

Stability: another view of generalization.

$$\hat{w} \leftarrow \underset{w}{\text{argmin}} R(w) + \mathbb{E}_S L(w)$$

intuition: alg stable if small change in input dataset does not change the output but much.

(i) Setting:

training set $S = (z_1, \dots, z_m) \stackrel{\text{iid}}{\sim} D$.

learning model: parametrized by w . $z = (x, y)$

loss fn: $l(w, z) \in \mathbb{R}$

$$l(w, (x, y)) = I(y \leq w \cdot x) \quad \text{0-1 loss}$$

$$= \max(0, 1 - y \cdot w \cdot x) \quad \text{huge loss}$$

$$L_D(w) = \mathbb{E}_{z \sim D} l(w, z) \quad : \text{generalization loss}$$

$$L_S(w) = \mathbb{E}_S l(w, z) = \frac{1}{|S|} \sum_{z \in S} l(w, z) \quad , \text{empirical loss.}$$

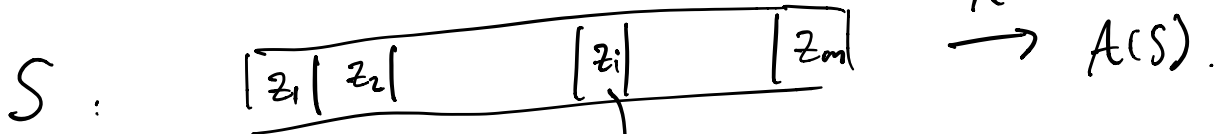
A : learning algorithm

$$A(S) = \hat{w}$$

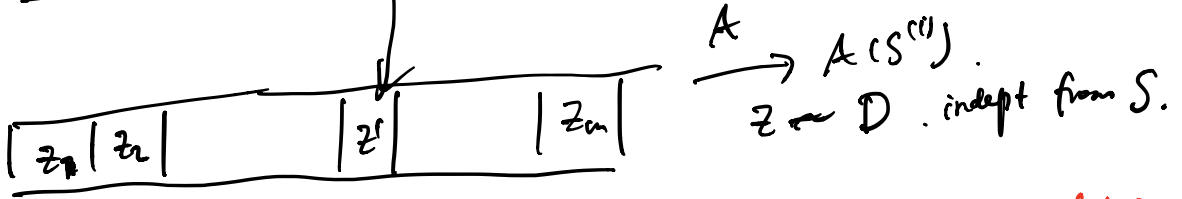
want to bound: $\mathbb{E}_{S \sim D^m} \left[\underbrace{L_D(\hat{w}) - L_S(\hat{w})}_{\text{generalization gap}} \right]$

(previously, want to give a high prob. upper bound of gen gap)

(2)



$i \in \{1, \dots, m\}$
 $S^{(i)}$:

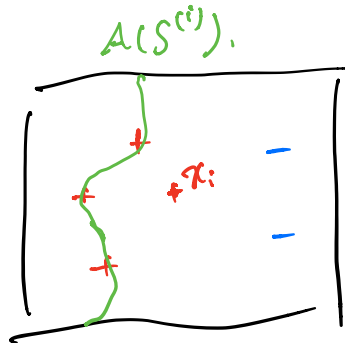
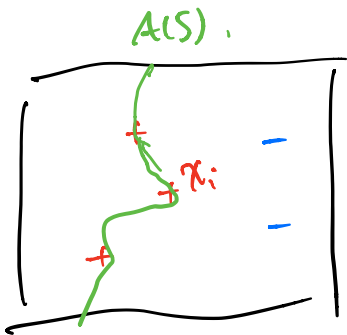


$$l(A(S), z_i)$$

vs.

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

when A is overfitting on data set S.



(OARO)

Def: learning alg A is on-average-replace-one-stable with rate ϵ for $g: \mathbb{N} \rightarrow \mathbb{R}$ if \forall distrn D. $\forall m$ (sample size).

$$\mathbb{E}_{(S, z') \sim D^{m+1}, i \in \text{unif}(\{1, \dots, m\})} \left[\underbrace{l(A(S^{(i)}), z_i')}_{\hat{w}^{(i)}} - \underbrace{l(A(S), z_i)}_{\hat{w}} \right] \leq \epsilon(m)$$

(If $A(S^{(i)}) \approx A(S)$, then this will be small)

Thm: If A is OARO-stable w/ rate g , then

$$\mathbb{E}_{S \sim D^m} [L_D(A(S)) - L_S(A(S))] \leq g(m).$$

Pf: show (4) = (Δ)

First term in (*):

$$\mathbb{E}_{(S, Z) \sim D^{m+1}} \mathbb{E}_{i \sim U(\{1, \dots, m\})} l(A(S^{(i)}), z_i) = \mathbb{E}_{S \sim D^m, Z \sim D} l(A(S), z) = \mathbb{E}_{S \sim D^m} L_D(A(S))$$

observe: fix i , $(S^{(i)}, z_i) \stackrel{d}{=} (S, z) = D^{m+1}$

Second term in (*):

$$\mathbb{E}_{(S, Z) \sim D^{m+1}} \mathbb{E}_{i \sim U(\{1, \dots, m\})} l(A(S), z_i) = \mathbb{E}_{S \sim D^m} L_S(A(S)).$$

$$\frac{1}{m} \sum_{i=1}^m l(A(S), z_i) = L_S(A(S)).$$

④ l_2 -regularization gives stability.

Assume: ① $l(w, z)$ is ρ -Lipschitz wrt w , for any z

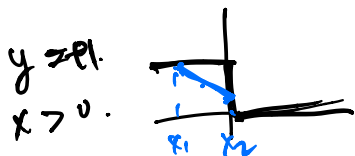
$$\forall w_1, w_2 \quad |l(w_1, z) - l(w_2, z)| \leq \rho \cdot |w_1 - w_2|.$$

(f differentiable $\max_z |f'(z)|$ $f(z_1) - f(z_2) = f'(\theta)(z_1 - z_2)$ $\theta \in [z_1, z_2]$)

② $l(w, z)$ is convex wrt w , for any z .

(hinge loss, logistic loss, exponential loss, etc).

$\mathbb{I}(y \leq w \cdot x \leq 0)$ not convex.



$$\textcircled{3} \quad \hat{w} = A(S) = \underset{w}{\operatorname{argmin}} \left(\frac{\lambda}{2} \|w\|_2^2 + L_S(w) \right).$$

we will show A is $g(w) = \frac{2\rho^2}{\lambda w}$ - OARO stable.

key tool: strong convexity

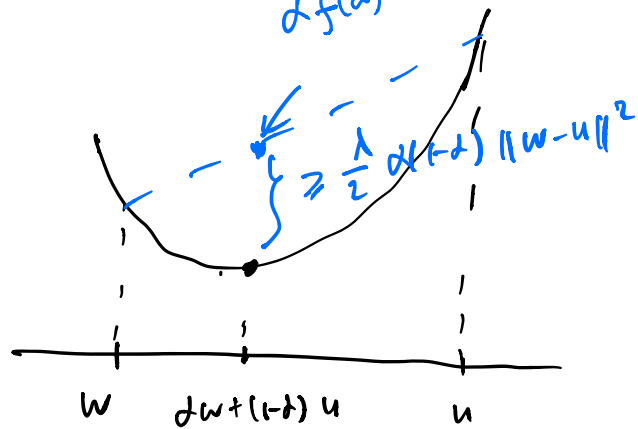
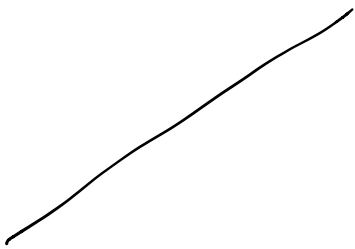
Def: fn f is λ -strongly cvx (λ -sc) if $\forall w, u$.

$\alpha \in (0, 1)$,

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

(0-sc \Leftrightarrow cvx).

$f =$



λ -sc for $\lambda > 0$?

\uparrow

key properties of sc. fns:

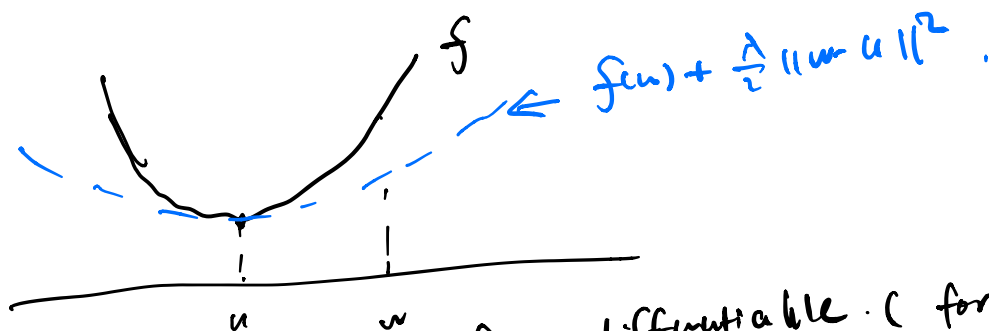
1. $f(w) = \frac{\lambda}{2} \|w\|^2$ is λ -sc. (Exercise)

2. if f is λ -sc, g is cvx, then $h = f + g$ is λ -sc.

(write down the λ -sc. cvx defs of f & g & and sum).

3. If f is λ -sc & $u = \underset{w}{\operatorname{argmin}} f(w)$, then.

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



df: focus on the case when f is differentiable. (for general f , we need the defn of subgradient).

$$\nabla f(u) = 0.$$

$$\frac{f(u + \alpha(w-u)) - f(u)}{\alpha \rightarrow 0} \leq \frac{f(w) - f(u) - \frac{\lambda}{2} (1-\alpha) \|w-u\|^2}{\alpha \rightarrow 0}$$

taking limit w/ $\alpha \rightarrow 0^+$

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

$$g(\alpha) = f(u + \alpha(w-u))$$

$$\text{LHS} \rightarrow g'(0) = 0$$

$$g'(\alpha) = \langle \nabla f(u + \alpha(w-u)), w-u \rangle$$

$$\stackrel{\alpha=0}{=} 0$$

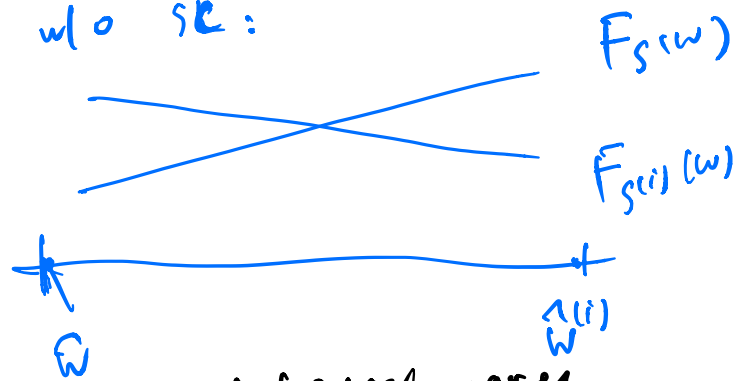
Back to stability:

$$\hat{w} = \arg \min_w \left(L_S(w) + \frac{\lambda}{2} \|w\|^2 \right) = F_S(w) \stackrel{\lambda-SL}{\approx}$$

$$\hat{w}^{(i)} = \arg \min_w \left(L_{S^{(i)}}(w) + \frac{\lambda}{2} \|w\|^2 \right) \stackrel{\lambda-SL}{\approx}$$

want to show \hat{w} is close to $\hat{w}^{(i)}$.

caveat: w/o SL:



w/ SL: can rule out such pathological cases.

property 3: $F_S(\hat{w}^{(i)}) - F_S(\hat{w}) \geq \frac{\lambda}{2} \|\hat{w}^{(i)} - \hat{w}\|^2$

$F_{S^{(i)}}(\hat{w}) - F_{S^{(i)}}(\hat{w}^{(i)}) \geq \frac{\lambda}{2} \|\hat{w}^{(i)} - \hat{w}\|^2$

summing up:

$(F_S(\hat{w}^{(i)}) - F_{S^{(i)}}(\hat{w}^{(i)})) - (F_S(\hat{w}) - F_{S^{(i)}}(\hat{w})) \geq \lambda \|\hat{w}^{(i)} - \hat{w}\|^2$

$(\frac{1}{m} l(\hat{w}^{(i)}, z_i) - \frac{1}{m} l(\hat{w}^{(i)}, z'_i)) - (\frac{1}{m} l(\hat{w}, z_i) - \frac{1}{m} l(\hat{w}, z'_i))$

$\leq \rho \cdot \|\hat{w}^{(i)} - \hat{w}\| \cdot \frac{l(\hat{w}^{(i)}, z_i) - l(\hat{w}, z_i)}{m} - \frac{l(\hat{w}^{(i)}, z'_i) - l(\hat{w}, z'_i)}{m}$

regroup w.r.t z_i or z'_i

$\frac{2\rho}{m} \|\hat{w} - \hat{w}^{(i)}\|$

$\Rightarrow \|\hat{w}^{(i)} - \hat{w}\| \leq \frac{2\rho}{m \cdot \lambda}$

$\Rightarrow l(\hat{w}^{(i)}, z_i) - l(\hat{w}, z_i) \leq \rho \cdot \|\hat{w}^{(i)} - \hat{w}\|$

$$\leq \frac{2\rho^2}{m\lambda}.$$

$\Rightarrow A$ is $g(m) = \frac{2\rho^2}{m\lambda}$ - OADR - stable.

Next class:

- Application of stability argument to Regularized loss minimization w/ the goal of minimizing final generalization loss
- online learning