

Lipschitzness  $\forall x, y$

$$L \|x - y\| \geq |f(x) - f(y)|$$

Smoothness  $f$  is  $\beta$ -smooth  
if  $\nabla f$  is  $\beta$ -Lipschitz:  $\forall x, y$   
$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

Ex:  $f(x) = x^2$  is 2-smooth over  $\mathbb{R}$ ,  
 $g(x) = x^3$  isn't  $c$ -smooth over  $\mathbb{R}$  for  
any constant  $c$ .

Convexity  $\forall \alpha \in (0, 1) \forall x, y \in \mathbb{C}$   
-  $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$   
(if differentiable,  $\forall x, y$ )  
$$f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle$$

Goal: Is loss min.  
a good plan even if  
we don't have finite VC?  
A: sometimes!

consider regularized  
loss minimization:

$$A(S) = \underset{w}{\operatorname{argmin}} L_S(w) + R(w)$$

where  $R$  is a regularizer  
of  $w$ . We'll focus on  
 $L^2$  reg,  $R(w) = \lambda \|w\|^2$

Example:  $L^2$ -regularized logistic regression

$$A(S) = \underset{w}{\operatorname{argmin}} \left[ \frac{1}{m} \sum_i \log(1 + \exp(-y_i \langle w, x_i \rangle)) + \lambda \|w\|^2 \right]$$

- No closed form solution, but we can solve efficiently.  
- This is strictly convex

Theorem. Fix  $\mathcal{D}$  over  $\mathcal{X} \times [-1, 1]$ ,  $\mathcal{X} = \{x \in \mathbb{R}^d \mid \|x\| \leq 1\}$ .  
Let  $\mathcal{H} = \{w \in \mathbb{R}^d \mid \|w\| \leq B\}$ . For any  $\epsilon \in (0, 1)$ ,  $m \geq \frac{8B^2}{\epsilon^2}$   
 $\lambda = \frac{\epsilon}{2B^2}$ ,  $L^2$ -regularized logistic regression satisfies  
$$\mathbb{E}_S [L_D(A_{RLM}(S))] \leq \min_{w \in \mathcal{H}} L_D(w) + \epsilon$$

[A similar statement about ridge regression...]

Why? Stability of regularized learning.

$$S = (z_1, \dots, z_m)$$

$$S^{(i)} = (z_1, \dots, z'_i, \dots, z_m)$$

$$\frac{\text{what is}}{L(A(S^{(i)}), z_i) - L(A(S), z_i)} \quad ?$$

(Probably)  $\geq 0$

If large, model is likely overfitting on  $z_i$ .

Def: On-average replace one stability An alg  $A$  is  $\epsilon(m)$

on-avg replacement stable if,  $\forall D$

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

Thm Fix  $S = (z_1, \dots, z_m) \neq z'_i \sim D$ . Let  $U[m]$  be uniform d. over  $[m]$ .  
Then  $\forall$  Algorithms  $A$ ,

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

Generalization gap

Pf

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

$$\mathbb{E}_S [L_D(A(S))] = \mathbb{E}_{S, z'} [L(A(S), z')]$$

$$= \mathbb{E}_{S, z'} [L(A(S^{(i)}), z')]$$

=

$\uparrow$   
 $z_i$

since  $S^{(i)} + z_i, z'$  are  $\perp$

So, On-avg replacement stability  $\Rightarrow$  Bd on gen. gap.

We'll show that RLM is OARS if  $L$  is convex +  $\frac{\text{smooth}}{\text{or Lipschitz}}$ .

Def (Strong convexity)  $f$  is  $\lambda$ -strongly convex if  $\forall x, y$

$$\alpha \in (0, 1), \quad f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \frac{\lambda \alpha(1-\alpha)}{2} \|x-y\|^2$$

( $\lambda$  increasing makes this more difficult to satisfy,  $\lambda=0$  all  $\text{convex}$  fns satisfy)

Note: if  $f$  is convex &  $g$  is  $\lambda$  strongly convex,  $f+g$  is  $\lambda$ -SC.

$$A(S) = \underset{w}{\operatorname{argmin}} \underbrace{L(S) + \lambda \|w\|^2}_{f(S)}$$

So, RLM for a convex  $L$  is  $2\lambda$ -SC.

Prop If  $f$  is  $\lambda$ -SC,  $u^*$  minimizer of  $f$ , then  $\forall w$

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

$\forall :$

$$\begin{aligned} f(u^*) &\leq \underbrace{f(u^*)}_{(\text{minimizer})} \\ &\leq \alpha f(w) + (1-\alpha) f(u^*) - \frac{\lambda \alpha (1-\alpha)}{2} \|u^* - w\|^2 \end{aligned}$$

(By basic alg  $\Rightarrow$ )  $f(u^*) \leq f(w) - \frac{\lambda}{2} (1-\alpha) \|u^* - w\|^2$

and as  $\alpha \rightarrow 1$ ,  $\checkmark$

$\square$

Let's bound  $\|A(S) - A(S^{(i)})\|$  from above & below

or  $\|f_S(A(S)) - f_S(A(S^{(i)}))\|$  regularized loss

LB: Since  $f_S$  is  $2\lambda$ -SC,  $f_S(A(S^{(i)})) - f_S(A(S)) \geq \lambda \|A(S^{(i)}) - A(S)\|^2$   
 $u^* = A(S)$

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

$$= \underbrace{L_{S^{(i)}}(\hat{w}^{(i)}) + \lambda \|\hat{w}^{(i)}\|^2 - [L_{S^{(i)}}(\hat{w}) + \lambda \|\hat{w}\|^2]}_{\leq 0}$$

$$+ \frac{1}{n} \left[ \ell(\hat{w}^{(i)}, z_i) - \ell(\hat{w}^{(i)}, z') + \ell(\hat{w}, z') - \ell(\hat{w}, z_i) \right]$$

$$\leq \frac{1}{n} \left[ \ell(A(S^{(i)}), z_i) - \ell(A(S^{(i)}), z') + \ell(A(S), z') - \ell(A(S), z_i) \right] \quad (1)$$

When  $l(\cdot, z)$  is  $\rho$ -Lipschitz  $\forall z$ ,

$$l(A(s^{(i)}), z_i) - l(A(s), z_i) \leq \rho \|A(s^{(i)}) - A(s)\|$$

& same for  $z'$ . So, along w (1)

$$\lambda \|A(s^{(i)}) - A(s)\|^2 \leq \frac{2\rho}{m} \|A(s^{(i)}) - A(s)\|$$

$$\text{So } \|A(s^{(i)}) - A(s)\| \leq \frac{2\rho}{\lambda m}$$

$$\& l(A(s^{(i)}), z_i) - l(A(s), z_i) \leq \frac{2\rho^2}{\lambda m}$$

$\Rightarrow$  when  $l$  is convex,  $\rho$ -Lipschitz, RLM  
w.  $R(w) = \lambda \|w\|^2$  is  $\frac{2\rho^2}{m}$ -GAPOS stable.

So its exp. gen gap is  $\leq \frac{2\rho^2}{m}$ .

[Similar argument when  $l$  is  $\beta$  smooth  
but  $\frac{48\beta}{\lambda m}$  UB instead

Overall loss?

$$\mathbb{E}_S [L_D(A(s))] = \mathbb{E}_S [L_S(A(s))] + \text{Gen-GAP}$$

$$\leq \mathbb{E}_S [L_S(A(s)) + \lambda \|A(s)\|^2] + \frac{2\rho^2}{\lambda m} \quad \begin{matrix} \leq \text{stability} \\ \circ \end{matrix}$$

$$\text{for any } w^* \leq \mathbb{E}_S [L_S(w^*) + \lambda \|w^*\|^2] + \frac{2\rho^2}{\lambda m}$$

$$= \mathcal{L}_D(\omega^*) + \lambda \|\omega^*\| + \frac{2\rho^2}{\lambda m}$$

Corollary If  $\mathcal{L}(\cdot, z)$  is convex +  $\rho$ -lipschitz  
 $\forall z$  &  $\|\omega\| \leq B$ , then  $\lambda = \sqrt{\frac{2\rho^2}{B^2 m}}$

RLM for  $R(\omega) = \lambda \|\omega\|^2$  satisfies

$$\mathbb{E}_S[\mathcal{L}_D(A(S))] \leq \min_{\omega \in \mathcal{H}} \mathcal{L}_D(\omega) + \rho B \sqrt{\frac{8}{m}}$$

or, if  $m \geq \frac{8\rho^2 B^2}{\epsilon^2}$

$$\leq \epsilon$$

We'll now switch to starting to talking more about knowing how good our predictions are (in a more precise way)

Let's start by recalling squared error of a predictor  $f$  (also called the Brier Score)

$$B(f, D) = \mathbb{E}_{(x, y) \sim D} [(f(x) - y)^2]$$

Let's consider regression, where  $f: \mathcal{X} \rightarrow [0, 1]$

We'd like to learn  $f^*$ :  $f^*(x) = \mathbb{E}_{y|x}[y]$

which we can't do, (Bayes-opt)

---

but maybe we can have marginal  
mean consistency w. error  $\alpha$ :

$$\left| \mathbb{E}_{x \sim D_x} [f(x)] - \mathbb{E}_{(x, y)} [y] \right| \leq \alpha$$

Since  
it's only  
a statement  
abt global f-avg

If  $f$  isn't  $\alpha$ -MMC, it's easy to

fix: let  $\Delta: \mathbb{E}_{(x,y) \sim \mathcal{D}}[y] - \mathbb{E}_{x \sim \mathcal{D}_x}[f(x)]$

Then,  $\hat{f}(x) = f(x) + \Delta$  is MMC.

But did we make  $f \rightarrow \hat{f}$  worse in some way?

Lemma: Fix any  $\mathcal{D}$ ,  $f: X \rightarrow [0,1]$ , and  $\Delta, \hat{f}$  as above. Then,

$$l_2(\hat{f}) = B(\hat{f}, \mathcal{D}) = B(f, \mathcal{D}) - \Delta^2$$

Eg, the squared loss is lower than for  $f$ !

Proof

$$B(\hat{f}, \mathcal{D}) - B(f, \mathcal{D}) = \mathbb{E}[(f(x) - y)^2 - (\hat{f}(x) - y)^2]$$

$$= \mathbb{E}\left[\begin{matrix} f^2(x) - 2f(x)y + y^2 \\ \hat{f}^2(x) - 2\hat{f}(x)y + y^2 \end{matrix}\right]$$

$$= \mathbb{E}\left[\begin{matrix} f^2(x) - 2f(x)y \\ -[f(x) + \Delta]^2 + 2[f(x) + \Delta]y \end{matrix}\right]$$

$$= \mathbb{E}[-2f(x)y - 2\Delta f(x) - \Delta^2 + 2f(x)y + 2\Delta y]$$

$$= \mathbb{E}[-2\Delta f(x) - \Delta^2 + 2\Delta y]$$

$$= 2\Delta \underbrace{\mathbb{E}[y - f(x)]}_{\Delta} - \Delta^2 = \Delta^2 \quad \checkmark$$

Ok, so marginal mean consistency of  $f$  isn't hard (or costly in terms of MSE).

What other things about  $(x, y) \sim D$  can we find models to satisfy/learn?

What about quantiles?  $\tau$  is a  $q$ -quantile of

① if

$$\Pr_y [y \leq \tau] = q$$

" $q$  fraction of  $D$  is below  $\tau$ "

It'd be rad to find  $f$  that was conditionally (on  $x$ ) a  $q$ -quantile, generally hard/imp. But can ask for it marginally:

Def  $f$  has marginal quantile consistency error  $\alpha$  for target quantile  $q$  if

$$\left| \Pr_{(x,y) \sim D} [y \leq f(x)] - q \right| \leq \alpha$$

Squared error: <sup>marginal</sup> mean consistency  
~~Pinball loss~~: quantile

Def: The pinball loss function for quantile  $q$ :

For  $D$  &  $f$ ,

$$L_q(\tau, y) = \begin{cases} (y - \tau)q & y > \tau \\ (\tau - y)(1-q) & y \leq \tau \end{cases} \quad \text{PB}_q(f, D) = \mathbb{E}_{(x,y) \sim D} [L_q(f(x), y)]$$



Lemma: For any continuous distribution over  $y$   
 $0 \leq q \leq 1$

$$\tau_q = \underset{\tau \in [0,1]}{\operatorname{argmin}} \mathbb{E}_y [L_q(\tau, y)]$$

is a  $q$ -quantile.

$$\int f(y) \left[ (y-\tau)^+ q \right] dy - \int f(y) \left[ (\tau-y)^+ (1-q) \right] dy$$

Pf Since  $y$  is continuous  $\Phi$  is continuous  
 $\Phi$  is actually convex in  $\tau$  so is min. at a pt  
 where its subderivative is 0

$$\frac{d}{d\tau} \mathbb{E}_y [L_q(\tau, y)] = \mathbb{E}_y [(1-q) \mathbb{1}[y \leq \tau] - q \mathbb{1}[y > \tau]]$$

$$= \mathbb{E}_y [\mathbb{1}[y \leq \tau] - q]$$

$$= \Pr[y \leq \tau] - q$$

$$= 0 \text{ where } \tau \text{ is a } q\text{-quantile.}$$

Can we "patch"/fix  $f$  to be  $q$ -quantile  
 like we did w. means?

Say  $f$  is  $\alpha$ -violating  $q$ -quantile consistent

Defined:  $\Pr(y \leq f(x) + \Delta) = q$  [ $\exists$  since  $y$  is continuous]

$$\hat{f}(x) = f(x) + \Delta$$

$\hookrightarrow$  now MQC

[And only has lower pinball loss.]

But its improvement in pinball loss depends on density b/w  $f$  &  $\hat{f}$

Recall  $\rho$ -Lipschitzness applied to cond label distribution  $D(x)$  says  $\forall 0 \leq \tau \leq \tau' \leq 1$

$$\Pr_{y \sim D(x)}[y \leq \tau'] - \Pr_{y \sim D(x)}[y \leq \tau] \leq \rho[\tau' - \tau]$$

The whole of the label dist  $D$  is  $\rho$ -Lipschitz if  $\forall x, D(x)$  is  $\rho$ -Lipschitz.

(we just need marginal Lipschitzness, not 1 on  $x$ )

Lemma: Fix any  $D$  which is continuous &  $\rho$ -Lipschitz. If  $f$  has marginal consistency err  $\alpha$

wrt  $q$ ,  $\Delta$  as above,  $\hat{f} = f + \Delta$

then 
$$PB_q(\hat{f}) \leq PB_q(f) - \frac{\alpha^2}{2\rho}$$

and

$$PB(f) \leq PB(\hat{f})$$

$$+ |\Delta| \alpha - \frac{\alpha^2}{2\rho}$$

This argument is messier (more calculus /  
but morally similar (see p14-15 geometry)  
in uncertainty notes).

Some overview of argument

$$\begin{aligned}\frac{\partial \text{PB}_q(f(x) + \tau)}{\partial \tau} &= \mathbb{E}_x \left[ \frac{\partial \mathbb{E}_{y \sim D_x} L_q(f(x) + \tau, y)}{\partial \tau} \right] \\ &= \mathbb{E}_x [\Pr(y \leq f(x) + \tau) - q] \\ &= \Pr_{x,y} [y \leq f(x) + \tau] - q\end{aligned}$$

$$\begin{aligned}\text{PB}_q(f) - \text{PB}_q(\tilde{f}) &= \text{PB}_q(f(x) + \Delta) - \text{PB}_q(f(x)) \\ &= \int_0^\Delta \frac{\partial \text{PB}_q(f(x) + \tau)}{\partial \tau} d\tau \\ &= \int_0^\Delta \underbrace{\Pr_{x,y} (y \leq f(x) + \tau) - q}_{\substack{\leftarrow \frac{1}{|\Delta|} q \quad \Delta \geq 0 \\ \leftarrow \frac{1}{|\Delta|} q \quad \Delta < 0}} d\tau\end{aligned}$$

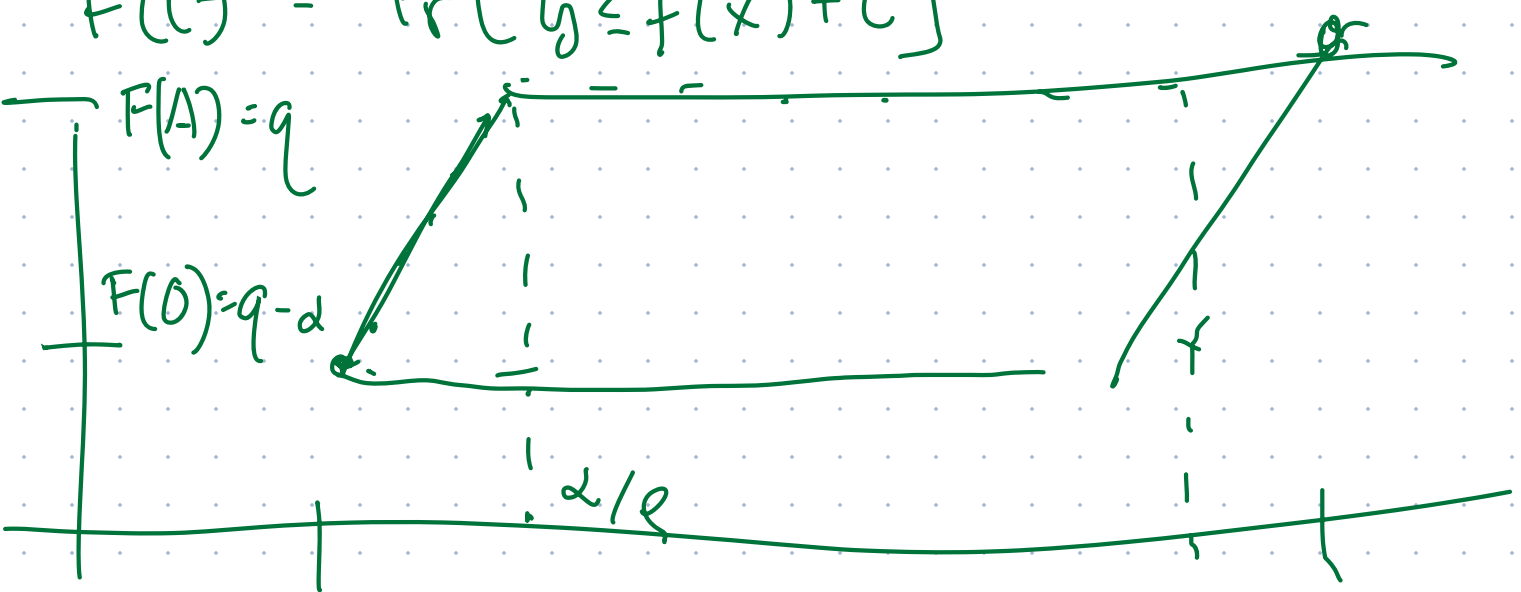
If we want to bound  $\tau$  It'll be useful to understand that integral.

lets consider  $\Delta \geq 0$  (area under curve between  $f(x)$ ,  $f(x) + \tau$ )

$f$  had quantile  $q - \Delta$

$\hat{f}$  has quantile  $q$

$$F(\tau) = \Pr[y \leq f(x) + \tau]$$



○ This area is biggest  $\Delta$  if the slopes are as steep as possible

$$[\text{slope} \leq \rho]$$