Regularization + Stability

CPSC 532S: Modern Statistical Learning Theory 28 February 2022 <u>cs.ubc.ca/~dsuth/532S/22/</u>

Admin

- Hope your break was good!
- Office hour times have moved:
 - Monday 12-1pm (and still Thursdays 4-5pm)
 - Both are available both on Zoom or in person (ICICS X563)
- A1 grades are up on Gradescope
- A2 solutions are posted
- A3 will be posted tonight or tomorrow, due in ~2 weeks
 - If you don't have a group and want one, post on Piazza (asap)

rsdays 4-5pm) m or in person (ICICS X563)

rrow, due in ~2 weeks vant one, post on Piazza (asap)

steps to reach ε (expected) error:



 ρ -Lipschitz + λ strongly

β-smooth

β-smooth + λ strongly convex

 $B\beta$ 2ε

$$\frac{B^2\beta}{2\varepsilon^2}$$
(avg)

steps to reach $\boldsymbol{\varepsilon}$ (expected) error:



Lipschitz	
λ strongly	
convex	

β-smooth

β-smooth
+ λ strongly
convex

$$\left(\frac{1}{\varepsilon}\right)$$

$$\frac{B\beta}{2\epsilon}$$

$$12\frac{B^2\beta}{2\varepsilon^2}$$
(avg)

steps to reach $\boldsymbol{\varepsilon}$ (expected) error:



-Lipschitz	
λ strongly	β-s
convex	

β-smooth

β-smooth
+ λ strongly
convex

$$\hat{\sigma}\left(\frac{1}{\varepsilon}\right) \qquad \frac{B\beta}{2\varepsilon}$$

$$15\frac{\rho^2}{\lambda\varepsilon} \qquad 12\frac{B^2\beta}{2\varepsilon^2}$$
ail avg) (avg)

steps to reach $\boldsymbol{\varepsilon}$ (expected) error:



-Lipschitz λ strongly convex	β-smooth	β-smooth+ λ stronglyconvex
$\mathcal{O}\left(\frac{1}{\varepsilon}\right)$	$\frac{B\beta}{2\varepsilon}$	$\frac{\beta}{\lambda}\log\frac{\beta B}{2\varepsilon}$
$15 \frac{\rho^2}{\lambda \varepsilon}$ ail avg)	$\frac{B^2\beta}{2\varepsilon^2}$	

steps to reach ε (expected) error:



-Lipschitz λ strongly convex	β-smooth	β-smooth+ λ stronglyconvex	
$\mathcal{O}\left(\frac{1}{\varepsilon}\right)$	$\frac{B\beta}{2\varepsilon}$	$\frac{\beta}{\lambda}\log\frac{\beta B}{2\varepsilon}$	
$15 \frac{\rho^2}{\lambda \varepsilon}$	$\frac{B^2\beta}{2\varepsilon^2}$	$\frac{\beta \rho^2}{\lambda^2 \varepsilon} \rho^{-1}$ (last)	(also _ipsch



steps to reach ε (expected) error:



ρ-Lipschitz	<pre> ρ-Lipschitz + λ strongly convex </pre>	β-smooth	β-smooth + λ strongly convex	
$\frac{B^2 \rho^2}{\varepsilon^2}$	$\mathcal{O}\left(\frac{1}{\varepsilon}\right)$	$\frac{B\beta}{2\varepsilon}$	$\frac{\beta}{\lambda}\log\frac{\beta B}{2\varepsilon}$	
$\frac{B^2 \rho^2}{\varepsilon^2}$ (avg iterate)	$\frac{\rho^2}{15 \frac{\lambda \varepsilon}{\lambda \varepsilon}}$ (tail avg)	$\frac{B^2\beta}{2\varepsilon^2}$	$\frac{\beta \rho^2}{2 \lambda^2 \varepsilon} \rho$ (last)	(also -Lipscł
Better rates v	with strong	convexity!		



• Some problems just aren't strongly convex (or λ is very small)

- Some problems just aren't strongly convex (or λ is very small)

• If f is convex and g is λ -strongly convex, then f + g is also λ -strongly convex

- Some problems just aren't strongly convex (or λ is very small)
- $g(w) = \frac{\lambda}{2} ||w||^2$ is λ -strongly convex

• If f is convex and g is λ -strongly convex, then f + g is also λ -strongly convex

- Some problems just aren't strongly convex (or λ is very small)
- If f is convex and g is λ -strongly convex, then f + g is also λ -strongly convex
- $g(w) = \frac{\lambda}{2} ||w||^2$ is λ -strongly convex
- Regularized loss minimization (RLM): $\operatorname{argmin}_{w} L_{S}(w) + R(w)$

- Some problems just aren't strongly convex (or λ is very small)
- If f is convex and g is λ -strongly convex, then f + g is also λ -strongly convex
- $g(w) = \frac{\lambda}{2} ||w||^2$ is λ -strongly convex
- Regularized loss minimization (RLM): $\operatorname{argmin}_{W} L_{S}(w) + R(w)$
 - Recall SRM: $\operatorname{argmin}_{h} L_{S}(h) + \varepsilon_{k_{h}}(n, \delta w_{k_{h}})$



argmin $L_{s}(w) + R(w) = argmin \pm w^{T}(X^{T}X + 7 \pm)w - y^{T}Xw$

 $\nabla = (x^{x}x + 7x)\hat{w} - ky = 0$ $\hat{w} = (x^{\tau}x + 2I)^{\prime}x^{\prime}\gamma$

Ridge regression Ridge regression: $\mathscr{H} = \{x \mapsto w^{\mathsf{T}}x\}, \ \mathscr{C}(h, (x, y)) = \frac{1}{2}(h(x) - y)^2, \ R(w) = \frac{\lambda}{2}\|w\|^2$ $L_{s}(w) = \underset{c}{\neq} \frac{1}{2} (w^{\tau} x_{c} - \gamma_{c})^{2} = \frac{1}{2} ||Xw - \gamma||^{2}$ = $\frac{1}{2} w^{\tau} X^{\tau} Xw - \gamma^{\tau} Xw + \frac{1}{2} ||ya^{2} Xw| = \frac{1}{2} w^{\tau} X^{\tau} Xw - \gamma^{\tau} Xw + \frac{1}{2} ||ya^{2} Xw| = \frac{1}{2} ||w||^{2} = w^{\tau} (\frac{2}{2} I) |w|$





Regularized loss minimization with SGD

• L2 regularizer is also called weight decay:





Regularized loss minimization with SGD

- L2 regularizer is also called weight decay:
 - $\nabla \frac{\lambda}{2} \|w\|^2 = \lambda w$, so use $w_{t+1} = w_t \eta(\lambda w_t + \hat{g}_t)$





Regularized loss minimization with SGD

- L2 regularizer is also called weight decay: • $\nabla \frac{\lambda}{2} \|w\|^2 = \lambda w$, so use $w_{t+1} =$

 $W_{t+1} = W_t - \eta_t (\lambda w_t + \hat{g})$ $= w_{t} - \frac{1}{4}w_{t} - \frac{1}{2\epsilon}$ = t-1 we - Lêge $= \frac{t-1}{t} \left(\frac{t-2}{t-1} w_{t-1} - \frac{1}{t} \hat{g}_{t-1} \right) - \frac{1}{2t} \hat{g}_{t}$ $= -\frac{1}{2} \cdot \frac{1}{t} \hat{g}_{t} \hat{g}_{t}$ $= -\frac{1}{2} \cdot \frac{1}{t} \hat{g}_{t} \hat{g}_{t}$ 42p

$$w_t - \eta(\lambda w_t + \hat{g}_t)$$

• With ρ -Lipschitz loss, SGD from 0 with $\eta_t = 1/(\lambda t)$ has step norm at most 2ρ $f(w) = f_{2}(w) = f(w) + \frac{1}{2} \lambda ||w||^{2}$







Regularized loss minimization with SGD • For ρ -Lipschitz, β -smooth, convex f, $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$: $\frac{3\rho^2}{2T} \quad \text{b/c} \quad \mathbb{E}\left[\|w_T - w_{\lambda}^*\|^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$

$$\mathbb{E}[f_{\lambda}(w_{T})] - f_{\lambda}(w_{\lambda}^{*}) \leq \frac{2\beta\rho^{2}}{\lambda^{2}T}$$

 $\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w^*_{\lambda}) + f(w^*_{\lambda}) - f(w^*)$





Regularized loss minimization with SGD • For ρ -Lipschitz, β -smooth, convex f, $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$: $\frac{2}{\tau} \quad b/c \quad \mathbb{E}\left[\|w_T - w_{\lambda}^*\|^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$

$$\mathbb{E}[f_{\lambda}(w_{T})] - f_{\lambda}(w_{\lambda}^{*}) \leq \frac{2\beta\rho^{2}}{\lambda^{2}T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) - f(w_{\lambda}) - f(w_{$$

 $\mathbb{E}f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[\|w_{\lambda}^*\|^2 - \|w_T\|^2 \right] + \frac{2\beta\rho^2}{\lambda^2 T}$

 $+f(w_{\lambda}^{*})-f(w^{*})$





R

Larized loss minimization with SC

$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:
 $\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T}$ b/c $\mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$
 $-f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$
 $-f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$
 $= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^{2}$, $\eta_{t} = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_{T})] - f_{\lambda}(w_{\lambda}^{*}) \leq \frac{2\beta\rho^{2}}{\lambda^{2}T} \quad \text{b/c} \quad \mathbb{E}\left[||w_{T} - w_{\lambda}^{*}||^{2}\right] \leq \frac{4\rho^{2}}{\lambda^{2}T}$$

$$\mathbb{E}f(w_{T}) - f(w^{*}) = \mathbb{E}f(w_{T}) - f(w_{\lambda}^{*}) + f(w_{\lambda}^{*}) - f(w^{*})$$

$$f(w_{T}) - f(w_{\lambda}^{*}) \leq \frac{\lambda}{2}\mathbb{E}\left[||w_{\lambda}^{*}||^{2} - ||w_{T}||^{2}\right] + \frac{2\beta\rho^{2}}{\lambda^{2}T}$$

$$= \frac{\lambda}{2}\mathbb{E}\left[\left(||w_{\lambda}^{*}|| + ||w_{T}||\right)\left(||w_{\lambda}^{*}|| - ||w_{T}||\right)\right] + \frac{2\beta\rho^{2}}{\lambda^{2}T}$$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$\mathbb{E}f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2}\mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2}\mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$





R

Larized loss minimization with SC

$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:
 $\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T}$ b/c $\mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$
 $-f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$
 $-f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$
 $= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$
 $\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T}$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$\mathbb{E}f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$





R

Larized loss minimization with SC

$$p$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$-f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$-f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T} \leq \frac{2B\rho}{\sqrt{T}} + \frac{2\beta\rho^2}{\lambda^2 T}$$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T} \leq \frac{2B\rho}{\sqrt{T}} + \frac{2\beta\rho^2}{\lambda^2 T}$$

• For
$$\rho$$
-Lipschitz, β -smooth, convex f , $f_{\lambda}(w) = f(w) + \frac{1}{2}\lambda ||w||^2$, $\eta_t = 1/(\lambda t)$:

$$\mathbb{E}[f_{\lambda}(w_T)] - f_{\lambda}(w_{\lambda}^*) \leq \frac{2\beta\rho^2}{\lambda^2 T} \quad \text{b/c} \quad \mathbb{E}\left[||w_T - w_{\lambda}^*||^2\right] \leq \frac{4\rho^2}{\lambda^2 T}$$

$$\mathbb{E}f(w_T) - f(w^*) = \mathbb{E}f(w_T) - f(w_{\lambda}^*) + f(w_{\lambda}^*) - f(w^*)$$

$$\mathbb{E}f(w_T) - f(w_{\lambda}^*) \leq \frac{\lambda}{2} \mathbb{E}\left[||w_{\lambda}^*||^2 - ||w_T||^2\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$= \frac{\lambda}{2} \mathbb{E}\left[\left(||w_{\lambda}^*|| + ||w_T||\right)\left(||w_{\lambda}^*|| - ||w_T||\right)\right] + \frac{2\beta\rho^2}{\lambda^2 T}$$

$$\leq \lambda B \mathbb{E}\left[||w_{\lambda}^* - w_T||\right] + \frac{2\beta\rho^2}{\lambda^2 T} \leq \frac{2B\rho}{\sqrt{T}} + \frac{2\beta\rho^2}{\lambda^2 T}$$





• So, that attempt didn't work here (though it can in some specific settings)

- So, that attempt didn't work here (though it can in some specific settings)
- Another approach which will work generically: stability

ough it can in some specific settings) enerically: **stability**

- So, that attempt didn't work here (though it can in some specific settings) Another approach which will work generically: stability
- - "If we don't change S much, the predictor A(S) doesn't change much"

- So, that attempt didn't work here (though it can in some specific settings) Another approach which will work generically: stability
- - "If we don't change S much, the predictor A(S) doesn't change much"
 - Regularizer $R(w) = \frac{1}{2} ||w||^2$ stabilizes the algorithm

- So, that attempt didn't work here (though it can in some specific settings) Another approach which will work generically: stability
- - "If we don't change S much, the predictor A(S) doesn't change much"
 - Regularizer $R(w) = \frac{1}{2} ||w||^2$ stabilizes the algorithm
- One variant is **on-average replace-one stability**: Let $S = (z_1, ..., z_n) \sim \mathcal{D}^n$, $z' \sim \mathcal{D}$, $i \sim \text{Uniform}([n])$ be independent. Let $S^{(i)} = (z_1, \ldots, z_{i-1}, z', z_{i+1}, \ldots, z_n).$

- So, that attempt didn't work here (though it can in some specific settings) Another approach which will work generically: stability
- - "If we don't change S much, the predictor A(S) doesn't change much"
 - Regularizer $R(w) = \frac{1}{2} ||w||^2$ stabilizes the algorithm
- One variant is on-average replace-one stability: Let $S = (z_1, ..., z_n) \sim \mathcal{D}^n, z' \sim \mathcal{D}, i \sim \text{Uniform}([n])$ be independent. Let $S^{(i)} = (z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_n).$
 - Theorem: $\mathbb{E}_{S}[L_{\mathcal{D}}(A(S)) L_{S}(A(S))] = \mathbb{E}_{S,z',i}[\ell(A(S^{(i)}), z_{i}) \ell(A(S), z_{i})].$

- So, that attempt didn't work here (though it can in some specific settings) Another approach which will work generically: stability
- - "If we don't change S much, the predictor A(S) doesn't change much"
 - Regularizer $R(w) = \frac{1}{2} ||w||^2$ stabilizes the algorithm
- One variant is on-average replace-one stability: Let $S = (z_1, ..., z_n) \sim \mathcal{D}^n, z' \sim \mathcal{D}, i \sim \text{Uniform}([n])$ be independent. Let $S^{(i)} = (z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_n).$
 - Theorem: $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S)) L_{S}(A(S))] = \mathbb{E}_{S,z',i}[\ell(A(S^{(i)}), z_{i}) \ell(A(S), z_{i})].$
 - A is on-average-replace-one-stable with rate $\varepsilon(n)$ if for all \mathcal{D} , $\mathbb{E}_{S,z',i}[\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)] \leq \varepsilon(n).$



 $f_{S}(v) - f_{S}(u) = L_{S}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S}(u) - \frac{\lambda}{2} \|u\|^{2}$

9



 $f_S(v) - f_S(u) = L_S(v) + \frac{\lambda}{2} \|v\|^2 - L_S(u) - \frac{\lambda}{2} \|v\|^2$ $= L_{S^{(i)}}(v) + \frac{\lambda}{2} \|v\|^2 - L_{S^{(i)}}(u) - \frac{\lambda}{2} \|u\|^2$

$$\frac{\frac{\lambda}{2}}{\|u\|^2} + \frac{\ell(v, z_i) - \ell(u, z_i)}{n} + \frac{\ell(u, z') - \ell(v, z_i)}{n}$$



On-average replace-one stability of RLM Let $f_S(w) = L_S(w) + \frac{1}{2}\lambda ||w||^2$ $f_S(v) - f_S(u) = L_S(v) + \frac{\lambda}{2} \|v\|^2 - L_S(u) - \frac{\lambda}{2} \|v\|^2 - L$

 $= \underbrace{L_{S^{(i)}}(v) + \frac{\lambda}{2} \|v\|_{1}^{2} - L_{S^{(i)}}(u) - \frac{\lambda}{2} \|u\|_{1}^{2} + \frac{v(v, \lambda_{i}) - v(v, \lambda_{i})}{k_{s}^{(i)}(v)}}_{f_{s}^{(i)}(v)}$ Plug in $v = A(S^{(i)}), u = A(S)$; note that v minimizes $f_{S^{(i)}}$:

$$\frac{\frac{\lambda}{2}}{\|u\|^{2}} + \frac{\ell(v, z_{i}) - \ell(u, z_{i})}{n} + \frac{\ell(u, z') - \ell(v, z_{i})}{n}$$



$$f_{S}(v) - f_{S}(u) = L_{S}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S}(u) - \frac{\lambda}{2} \|u\|^{2}$$

= $L_{S^{(i)}}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S^{(i)}}(u) - \frac{\lambda}{2} \|u\|^{2} + \frac{\ell(v, z_{i}) - \ell(u, z_{i})}{n} + \frac{\ell(u, z') - \ell(v, z_{i})}{n}$

Plug in $v = A(S^{(i)})$, u = A(S); note that v minimizes $f_{S^{(i)}}$: $f_{S}(A(S^{(i)}) - f_{S}(A(S)) \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{1 + \ell(A(S), z_{i})} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{1 + \ell(A(S), z') - \ell(A(S^{(i)}), z')}$

n



$$f_{S}(v) - f_{S}(u) = L_{S}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S}(u) - \frac{\lambda}{2} \|u\|^{2}$$
$$= L_{S^{(i)}}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S^{(i)}}(u) - \frac{\lambda}{2} \|u\|^{2} + \frac{\ell(v, z_{i}) - \ell(u, z_{i})}{n} + \frac{\ell(u, z') - \ell(v, z_{i})}{n}$$

Plug in $v = A(S^{(i)})$, u = A(S); note that v minimizes $f_{S^{(i)}}$: $f_{S}(A(S^{(i)}) - f_{S}(A(S))) \leq \frac{\ell(A(S^{(i)}), z_{i}) - f_{S}(A(S))}{2}$

 f_S is λ -strongly convex, so $f_S(v) - f_S(A(S)) \ge \frac{\lambda}{2} ||v - A(S)||^2$

$$-\ell(A(S), z_i) + \ell(A(S), z') - \ell(A(S^{(i)}), z')$$

n



$$f_{S}(v) - f_{S}(u) = L_{S}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S}(u) - \frac{\lambda}{2} \|u\|^{2}$$

= $L_{S^{(i)}}(v) + \frac{\lambda}{2} \|v\|^{2} - L_{S^{(i)}}(u) - \frac{\lambda}{2} \|u\|^{2} + \frac{\ell(v, z_{i}) - \ell(u, z_{i})}{n} + \frac{\ell(u, z') - \ell(v, z_{i})}{n}$

Plug in $v = A(S^{(i)})$, u = A(S); note that v minimizes $f_{S^{(i)}}$: $f_{S}(A(S^{(i)}) - f_{S}(A(S))) \leq \frac{\ell(A(S^{(i)}), z_{i}) - f_{S}(A(S))}{2}$

$$f_{S} \text{ is } \lambda \text{-strongly convex, so } f_{S}(v) - f_{S}(A(S)) \ge \frac{\lambda}{2} \|v - A(S)\|^{2}$$
$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \le \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{n} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{n}$$

$$-\ell(A(S), z_i) + \ell(A(S), z') - \ell(A(S^{(i)}), z')$$

n


$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^2 \le \frac{\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)}{2} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{2}$

On-average replace-one stability for Lipschitz RLM



 $\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^2 \le \frac{\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)}{2} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{2}$

If $\ell(\cdot, z_i)$ is ρ -Lipschitz:

On-average replace-one stability for Lipschitz RLM



On-average replace-one

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S^{(i)}))}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

 $\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^2 \le \frac{\rho \|A(S^{(i)}) - A(S)\|}{1 + \frac{\rho \|A(S) - A(S^{(i)})\|}{1 + \frac{\rho \|A(S^{(i)})\|}{1 + \frac{\rho \|A(S) - A(S^{(i)})\|}{1 + \frac{\rho \|A(S^{(i)})\|}{1 + \frac{\rho \|A(S^{(i)}$ n

e stability for Lipschitz RLM $-\ell(A(S), z_i) + \ell(A(S), z') - \ell(A(S^{(i)}), z')$ n

n



On-average replace-one

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S^{(i)}))}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

n

e stability for Lipschitz RLM $-\ell(A(S), z_i) + \ell(A(S), z') - \ell(A(S^{(i)}), z')$ n

 $\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^2 \le \frac{\rho \|A(S^{(i)}) - A(S)\|}{2} + \frac{\rho \|A(S) - A(S^{(i)})\|}{2} = \frac{2\rho \|A(S^{(i)}) - A(S)\|}{2}$ n n



On-average replace-one stability for Lipschitz RL

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{n} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

 $\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^2 \le \frac{\rho \|A(S^{(i)}) - A(S)\|}{2} + \frac{\rho \|A(S) - A(S^{(i)})\|}{2} = \frac{2\rho \|A(S^{(i)}) - A(S)\|}{2}$ n

N n $\|A(S^{(i)}) - A(S)\| \le \frac{4\rho}{\lambda n}$



On-average replace-one stability for Lipschitz RL

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{n} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\rho \|A(S^{(i)}) - A(S)\|}{n} + \frac{\rho \|A(S) - A(S^{(i)})\|}{n} = \frac{2\rho \|A(S^{(i)}) - A(S)\|}{n}$$

$$\|A(S^{(i)}) - A(S)\| \leq \frac{4\rho}{\lambda n}$$
So, because $\ell(\cdot, z_{i})$ is ρ -Lipschitz:



On-average replace-one stability for Lipschitz RI

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{n} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\rho \|A(S^{(i)}) - A(S)\|}{n} + \frac{\rho \|A(S) - A(S^{(i)})\|}{n} = \frac{2\rho \|A(S^{(i)}) - A(S)\|}{n}$$

$$\|A(S^{(i)}) - A(S)\| \leq \frac{4\rho}{\lambda n}$$
So, because $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

$$\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i}) \leq \frac{4\rho^{2}}{\lambda n}$$

 λn



On-average replace-one stability for Lipschitz RL

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i})}{n} + \frac{\ell(A(S), z') - \ell(A(S^{(i)}), z')}{n}$$
If $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

$$\frac{\lambda}{2} \|A(S^{(i)}) - A(S)\|^{2} \leq \frac{\rho \|A(S^{(i)}) - A(S)\|}{n} + \frac{\rho \|A(S) - A(S^{(i)})\|}{n} = \frac{2\rho \|A(S^{(i)}) - A(S^{(i)})\|}{n}$$
So, because $\ell(\cdot, z_{i})$ is ρ -Lipschitz:

$$\ell(A(S^{(i)}), z_{i}) - \ell(A(S), z_{i}) \leq \frac{4\rho^{2}}{\lambda n}$$
and $\mathbb{E}_{S}[L_{\mathfrak{D}}(A(S)) - L_{S}(A(S))] \leq \frac{4\rho^{2}}{\lambda n}$



On-average replace-one stability for Smooth RLM

If $\ell(\cdot, z_i)$ is β -smooth and nonnegative, can show (SSBD section 13.3.2):

$$\begin{split} \ell(A(S^{(i)}), z_i) - \ell(A(S), z_i) &\leq \frac{48\beta}{\lambda n} \left(\ell(A(S), z_i) + \ell(A(S^{(i)}), z') \right) \\ & \text{if } \lambda \geq \beta/n \\ lies \mathbb{E}_S[L_{\mathscr{D}}(A(S)) - L_S(A(S))] &\leq \frac{96\beta}{\lambda n} \mathbb{E}_S[L_S(A(S))] \end{split}$$

which impl



On-average replace-one stability for Smooth RLM

If $\ell(\cdot, z_i)$ is β -smooth and nonnegative, can show (SSBD section 13.3.2):

$$\begin{aligned} \ell(A(S^{(i)}), z_i) - \ell(A(S), z_i) &\leq \frac{48\beta}{\lambda n} \left(\ell(A(S), z_i) + \ell(A(S^{(i)}), z') \right) \\ & \text{if } \lambda \geq \beta/n \\ \text{lies } \mathbb{E}_S[L_{\mathscr{D}}(A(S)) - L_S(A(S))] &\leq \frac{96\beta}{\lambda n} \mathbb{E}_S[L_S(A(S))] \end{aligned}$$

which impl

(and note, e.g., $\mathbb{E}_{S}[L_{S}(A(S))] \leq \mathbb{E}_{S}[L_{S}(0)] = L_{S}(0)$, often bounded by a constant)

OOP5





$\mathbb{E}_{S}[L_{\mathcal{D}}(A(S))] = \mathbb{E}_{S}[L_{S}(A(S))] + \mathbb{E}_{S}[L_{\mathcal{D}}(A(S)) - L_{S}(A(S))]$

$\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] = \mathbb{E}_{S}[L_{S}(A(S))] + \mathbb{E}_{S}[L_{\mathscr{D}}(A(S)) - L_{S}(A(S))]$ • Second term is the on-average replace-one stability

$\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] = \mathbb{E}_{S}[L_{S}(A(S))] + \mathbb{E}_{S}[L_{\mathscr{D}}(A(S)) - L_{S}(A(S))]$ Second term is the on-average replace-one stability

- - Bigger λ means more stable

$\mathbb{E}_{S}[L_{\mathcal{D}}(A(S))] = \mathbb{E}_{S}[L_{S}(A(S))] + \mathbb{E}_{S}[L_{\mathcal{D}}(A(S)) - L_{S}(A(S))]$ Second term is the on-average replace-one stability

- - Bigger λ means more stable
- First term is how well the algorithm fits the training data

$\mathbb{E}_{S}[L_{\mathcal{D}}(A(S))] = \mathbb{E}_{S}[L_{S}(A(S))] + \mathbb{E}_{S}[L_{\mathcal{D}}(A(S)) - L_{S}(A(S))]$ Second term is the on-average replace-one stability

- - Bigger λ means more stable
- First term is how well the algorithm fits the training data
 - Bigger λ means worse fit

$L_{S}(A(S)) \leq L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2}$

For any fixed vector w^*

 $L_{S}(A(S)) \le L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2} \le L_{S}(w^{*}) + \frac{\lambda}{2} ||w^{*}||^{2}$

For any fixed vector w^*

 $L_{S}(A(S)) \le L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2} \le L_{S}(w^{*}) + \frac{\lambda}{2} ||w^{*}||^{2}$ $\leq \mathbb{E}_{S} L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$

 $\mathbb{E}_{S}L_{S}(A(S))$

For any fixed vector w^*

 $L_{S}(A(S)) \le L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2} \le L_{S}(w^{*}) + \frac{\lambda}{2} ||w^{*}||^{2}$

 $\mathbb{E}_{S}L_{S}(A(S))$

 $\leq \mathbb{E}_{S} L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2} = L_{\mathcal{D}}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$



For any fixed vector w^*

 $L_{S}(A(S)) \le L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2} \le$

 $\mathbb{E}_{S}L_{S}(A(S))$

 $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^{*}) + \frac{\lambda}{2} \| v$

$$\leq L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$$

$$\leq \mathbb{E}_{S}L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2} = L_{\mathcal{D}}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$$

$$w^* \|^2 + \mathbb{E}_S[L_{\mathcal{D}}(A(S)) - L_S(A(S))]$$



For any fixed vector w^*

 $L_{S}(A(S)) \leq L_{S}(A(S)) + \frac{\lambda}{2} ||A(S)||^{2} \leq$

 $\mathbb{E}_{S}L_{S}(A(S))$

 $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^{*}) + \frac{\lambda}{2} \| v$

So, for convex ρ -Lipschitz loss, RLM with regularizer $\frac{\lambda}{2} ||w||^2$ has $\mathbb{E}_{S}[L_{\mathcal{D}}(A(S))] \leq L_{\mathcal{D}}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2} + \frac{4\rho^{2}}{\lambda n}$

$$\leq L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$$

$$\leq \mathbb{E}_{S}L_{S}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2} = L_{\mathcal{D}}(w^{*}) + \frac{\lambda}{2} \|w^{*}\|^{2}$$

$$v^* \|^2 + \mathbb{E}_S[L_{\mathcal{D}}(A(S)) - L_S(A(S))]$$



• RLM with regularizer $\frac{\lambda}{2} \|w\|^2$ has $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \le L_{\mathscr{D}}(w^*) + \frac{\lambda}{2} \|w^*\|^2 + \frac{4\rho^2}{\lambda n}$



• Taking
$$\lambda = \frac{\rho}{B} \sqrt{\frac{8}{n}}$$
 gives $\mathbb{E}_{S}[L_{\mathbb{S}}]$

• RLM with regularizer $\frac{\lambda}{2} \|w\|^2$ has $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^*) + \frac{\lambda}{2} \|w^*\|^2 + \frac{4\rho^2}{\lambda n}$

 $\mathcal{D}(A(S))] \le L_{\mathcal{D}}(w^*) + \rho B \sqrt{\frac{8}{n}}$



• Taking
$$\lambda = \frac{\rho}{B} \sqrt{\frac{8}{n}}$$
 gives $\mathbb{E}_{S}[L_{\mathbb{S}}]$

• So, for $n \ge 8\rho^2 B^2 / \varepsilon^2$, $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \le \inf_{w \in \mathscr{H}} L_{\mathscr{D}}(w) + \varepsilon$

• RLM with regularizer $\frac{\lambda}{2} \|w\|^2$ has $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \le L_{\mathscr{D}}(w^*) + \frac{\lambda}{2} \|w^*\|^2 + \frac{4\rho^2}{\lambda n}$

 $\mathcal{D}(A(S))] \le L_{\mathcal{D}}(w^*) + \rho B \sqrt{\frac{8}{n}}$



• Taking
$$\lambda = \frac{\rho}{B} \sqrt{\frac{8}{n}}$$
 gives $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^{*}) + \rho B \sqrt{\frac{8}{n}}$

- So, for $n \ge 8\rho^2 B^2 / \varepsilon^2$, $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \le \inf_{w \in \mathscr{H}} L_{\mathscr{D}}(w) + \varepsilon$

• RLM with regularizer $\frac{\lambda}{2} \|w\|^2$ has $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^*) + \frac{\lambda}{2} \|w^*\|^2 + \frac{4\rho^2}{\lambda n}$

• Similar result for convex-smooth-bounded (SSBD Corollaries 13.10, 13.11)



• RLM with regularizer $\frac{\lambda}{2} ||w||^2$ has \mathbb{E}_S

• Taking
$$\lambda = \frac{\rho}{B} \sqrt{\frac{8}{n}}$$
 gives $\mathbb{E}_{S}[L_{\mathscr{D}}(A(S))] \leq L_{\mathscr{D}}(w^{*}) + \rho B \sqrt{\frac{8}{n}}$

- So, for $n \ge 8\rho^2 B^2 / \varepsilon^2$, $\mathbb{E}_S[L_{\mathscr{D}}(A(S))] \le \inf_{w \in \mathscr{H}} L_{\mathscr{D}}(w) + \varepsilon$

$$S[L_{\mathscr{D}}(A(S))] \le L_{\mathscr{D}}(w^*) + \frac{\lambda}{2} \|w^*\|^2 + \frac{4\rho^2}{\lambda n}$$

• Similar result for convex-smooth-bounded (SSBD Corollaries 13.10, 13.11)

Can convert these expectation bounds into high-probability: SSBD exercise 13.1





• Uniform stability instead says $|\ell(A(S), z) - \ell(A(S^{(i)}, z))| \leq \gamma$ for all possible training sets S, slightly different sets $S^{(i)}$, test points z

- Uniform stability instead says $|\ell(A(S), z) \ell(A(S^{(i)}, z))| \leq \gamma$ for all possible training sets S, slightly different sets $S^{(i)}$, test points z
- Stronger condition: implies on-average replace-one stability

- Uniform stability instead says $\ell(A(S))$ for all possible training sets S, slightly different sets $S^{(l)}$, test points z
- Stronger condition: implies on-average replace-one stability
 - We (*almost*) showed $\gamma = \frac{4\rho^2}{\lambda n}$ in our proof for ρ -Lipschitz f with $\frac{\lambda}{2} ||w||^2$

$$S(z), z) - \ell(A(S^{(i)}, z)) \le \gamma$$

- Uniform stability instead says $\ell(A(S))$ for all possible training sets S, slightly different sets $S^{(i)}$, test points z
- Stronger condition: implies on-average replace-one stability
 - We (almost) showed $\gamma = \frac{4\rho^2}{\lambda n}$ in our proof for ρ -Lipschitz f with $\frac{\lambda}{2} ||w||^2$

Bousquet and Elisseef (2002) (also in MRT chap 14): if $\ell \in [0,C]$, $L_{\mathcal{D}}(A(S)) - L_{S}(A(S)) \le \operatorname{const}\left(\sqrt{n\gamma} + C/\sqrt{n}\right)\sqrt{\log\frac{1}{\delta}}$

$$S(x), z) - \ell(A(S^{(i)}, z)) \le \gamma$$

- Uniform stability instead says $\ell(A(S))$ for all possible training sets S, slightly different sets $S^{(l)}$, test points z
- Stronger condition: implies on-average replace-one stability
 - We (almost) showed $\gamma = \frac{4\rho^2}{\lambda n}$ in our proof for ρ -Lipschitz f with $\frac{\lambda}{2} ||w||^2$
- Bousquet and Elisseef (2002) (also in MRT chap 14): if $\ell \in [0,C]$,
- Bousquet et al. (2019), building off Feldman and Vondrák (2018, 19), show

$$S(x), z) - \ell(A(S^{(i)}, z)) \le \gamma$$

 $L_{\mathcal{D}}(A(S)) - L_{S}(A(S)) \le \operatorname{const}\left(\sqrt{n\gamma} + C/\sqrt{n}\right)\sqrt{\log\frac{1}{\delta}}$

 $L_{\mathcal{D}}(A(S)) - L_{S}(A(S)) \le \operatorname{const}\left(\gamma \log(n) \log\left(\frac{1}{\delta}\right) + \frac{C}{\sqrt{n}}\sqrt{\log\frac{1}{\delta}}\right)$

- <u>Hardt, Recht, and Singer (ICML-16)</u>:

• For convex, β -smooth, ρ -Lipschitz losses, T-step **multi-pass** SGD with $\eta_t \leq \frac{2}{\beta}$ gives uniform stability of $\frac{2\rho^2}{n} \sum_{t=1}^T \eta_t$

- Hardt, Recht, and Singer (ICML-16):
 - For convex, β -smooth, ρ -Lipschitz losses, T-step **multi-pass** SGD with $\eta_t \leq \frac{2}{\beta}$ gives uniform stability of $\frac{2\rho^2}{n} \sum_{t=1}^T \eta_t$

• For λ -strongly convex, β -smooth losses, T-step **multi-pass** projected $2L^2$ SGD with constant $\eta \leq 1/\beta$ has uniform stability $\frac{--}{\lambda n}$, independent of T

- Hardt, Recht, and Singer (ICML-16):
 - For convex, β -smooth, ρ -Lipschitz losses, T-step **multi-pass** SGD with $\eta_t \leq \frac{2}{\beta}$ gives uniform stability of $\frac{2\rho^2}{n} \sum_{t=1}^T \eta_t$
 - For λ -strongly convex, β -smooth losses, T-step **multi-pass** projected $2L^2$ SGD with constant $\eta \leq 1/\beta$ has uniform stability $\frac{--}{\lambda n}$, independent of T

- Some results even for nonconvex case

- Hardt, Recht, and Singer (ICML-16):
 - For convex, β -smooth, ρ -Lipschitz losses, T-step **multi-pass** SGD with $\eta_t \leq \frac{2}{\beta}$ gives uniform stability of $\frac{2\rho^2}{n} \sum_{t=1}^T \eta_t$
 - For λ -strongly convex, β -smooth losses, T-step **multi-pass** projected $2L^2$ SGD with constant $\eta \leq 1/\beta$ has uniform stability $\frac{--}{\lambda n}$, independent of T

- Some results even for nonconvex case
- <u>Chen, Jin, Yu (2018)</u> and <u>Feldman and Vondrák (2019)</u> extend to full-batch gradient descent, other variants

Summary

- Adding an L2 regularizer makes things strongly convex, which is nicer
 - Ridge regression / weight decay
 - Regularized loss minimization (RLM), the regularized version of ERM
 - Analogous to SRM
Summary

- Adding an L2 regularizer makes things strongly convex, which is nicer
 - Ridge regression / weight decay
 - Regularized loss minimization (RLM), the regularized version of ERM
 - Analogous to SRM
- Can analyze via stability
 - Holds for convex-Lipschitz-bounded / convex-smooth-bounded
- - On-average replace-one stability characterizes learnability Amount of regularization trades off between fitting and stability

Summary

- Adding an L2 regularizer makes things strongly convex, which is nicer
 - Ridge regression / weight decay
 - Regularized loss minimization (RLM), the regularized version of ERM
 - Analogous to SRM
- Can analyze via stability
 - Holds for convex-Lipschitz-bounded / convex-smooth-bounded
- - On-average replace-one stability characterizes learnability Amount of regularization trades off between fitting and stability
- Uniform stability: stronger notion that can give better bounds • Multi-pass SGD / GD / ... are uniformly stable