

# Lecture 16: end of strong convexity; weak implicit regularization

Remark: (failure of unification.)  
Literature has incomparable results across settings even within near initialization (e.g., smooth activation vs ReLU, least squares vs logistic).

} Recall Lipschitz & bounded setting, which in strongest form required  
 $\|\partial \mathcal{L}\| \leq \beta \cdot \mathcal{L}$  (holds for logistic).

Plan:

- \* Today is lecture 6/6 of near-init opt.
- \* Next week: Clarke differential & GF
- \* Outside NTK ("margins" and other stuff).

\* At least 3 gen lectures including 2 on double descent & interpolation.

\* HW 3 (lol) Tuesday, typed notes with it,  
1 "coding" problem

Theorem. Suppose  $\hat{R} \stackrel{\Delta}{=} \mathbb{2} \circ H$   $\beta$ -smooth,  $\mathcal{L}$  is  $\lambda$ -sc,  $\sigma_{\mathcal{L}} := \min_{\text{set}} \sigma_n(H(w_s)) > 0$

define  $\bar{y} := \arg \min_{\mathcal{Y}} \mathcal{L}(\bar{y})$ , and  $\eta \leq \frac{1}{\beta}$ . Then

$$\hat{R}(w_t) - \mathcal{L}(\bar{y}) \leq (\hat{R}(w_0) - \mathcal{L}(\bar{y})) \exp(-t\eta \lambda \sigma_{\mathcal{L}}^2).$$

Proof. (last time, based on  $\hat{R}(w_{s+1}) \leq \hat{R}(w_s) - \frac{\eta}{2} \|\partial \hat{R}(w_s)\|^2 \leq \hat{R}(w_s) - \frac{\eta}{2} \|\partial H(w_s)\|_2^2 \cdot \|\mathbb{2}\|_2^2$ )

Issues:  $\beta$  is time varying:  $\|\partial \hat{R}(a) - \partial \hat{R}(b)\|^2 \stackrel{?}{\leq} \beta \|a-b\|$   
 $= \|\partial H(a) \partial \mathbb{2}(H(a)) - \partial H(b) \partial \mathbb{2}(H(b))\|$

$\sigma_{\mathcal{L}}$  depends on algorithm

Lemma. Suppose  $\|\partial H(a) - \partial H(w_0)\| \leq \frac{1}{2} \sigma_n(H(w_0)) \quad \forall \|a-w_0\| \leq \beta$ .

Then  $\sigma_1(\partial H(a)) \geq \frac{3}{2} \sigma_1(\partial H(w_0))$ , and  $\sigma_n(\partial H(a)) \geq \frac{1}{2} \sigma_n(H(w_0))$ .

Proof. UB:  $\sigma_1(\partial H(a)) = \|\partial H(a)\|_2 \leq \|\partial H(w_0)\|_2 + \|\partial H(a) - \partial H(w_0)\|_2$   
 $\leq \sigma_1(\partial H(w_0)) + \frac{1}{2} \sigma_n(H(w_0)) \leq \frac{3}{2} \sigma_1(\partial H(w_0))$

Reasonable if  $\partial H$  smooth:  
 $\|\partial H(a) - \partial H(w_0)\| \leq \beta_H \|a-w_0\|$ ,  
 so we assumed  $\|a-w_0\| \leq \frac{\sigma_n(H(w_0))}{2\beta_H}$

LB:  $J_0 := \partial H(w_0)$ ,  $J := \partial H(a)$ . Then

$\sigma_n(J)^2 = \lambda_{\min}(JJ^T) = \min_{\|v\|=1} (v^T JJ^T v) = \min_{\|v\|=1} v^T (J-J_0+J_0)(J-J_0+J_0)^T v$   
 $= \min_{\|v\|=1} v^T J_0^T J_0 v + 2 v^T (J-J_0) J_0^T v + v^T (J-J_0)(J-J_0)^T v$   
 $\geq \min_{\|v\|=1} \|J_0^T v\|^2 - 2 \|(J-J_0)^T v\| \cdot \|J_0^T v\| + \|(J-J_0)^T v\|^2$   
 $= \min_{\|v\|=1} \left( \underbrace{\|J_0^T v\|}_{\geq \frac{\sigma_n(J_0)}{2}} - \underbrace{\|(J-J_0)^T v\|}_{\leq \frac{1}{2} \sigma_n(J_0)} \right)^2 \geq \left( \frac{1}{2} \sigma_n(J_0) \right)^2$

Recall  $J \in \mathbb{R}^{n \times p}$

$\min_{\|v\|=1} (\|Jv\|)^2$

$\|(J-J_0+J_0)v\| \leq \|Jv\| + \|J_0v\|$

$\|A\| = \|A-B+B\| \leq \|A-B\| + \|B\|$

$\Rightarrow \|B\| \geq \|A\| - \|A-B\|$   
 $\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$   
 $\frac{\sigma_n}{2} \quad \sigma_n \quad \frac{\sigma_n}{2} \quad \sigma_n$

Easier version: (Thanks!)

$\sigma_n(J)^2 = \min_{\|v\|=1} v^T JJ^T v = \left( \min_{\|v\|=1} \|Jv\| \right)^2 \geq \left( \min_{\|v\|=1} \|J_0v\| - \|(J-J_0)v\| \right)^2$   
 $\geq \left( \min_{\|v\|=1} \|J_0v\| - \|v\| \cdot \|J-J_0\| \right)^2$   
 $\geq \left( \min_{\|v\|=1} \|J_0v\| - \frac{\sigma_n(J_0)}{2} \right)^2 = \left( \frac{\sigma_n(J_0)}{2} \right)^2$

(In the typed notes: show how to control  $\|w_t - z\|$  for this strongly convex setting.)

Weak implicit regularization: GD automatically finds good predictors & stays close to initialization.

Theorem (from before): Assume  $w_{st+1} := \Pi_U(w_s - \eta \partial \hat{R}(w_s))$ , let  $z \in U$  be given &

$$\varepsilon := \max_{S \subseteq U} \max_{z \in S} |H(z); - (H(w_0); - \langle \partial H(w_0); \cdot, w_s - z \rangle)|,$$

$$\|\partial \mathcal{L}(w_s)\| \leq A \mathcal{L}(w_s) \quad \text{vs} \quad (\text{true for Logistic with } A=1)$$

$$B := \max_{S \subseteq U} \|\partial H(w_s)\|_{2, \infty} \quad (\leq 1 \text{ for shallow ReLU}).$$

$$\Rightarrow \frac{1}{2\eta t} \|w_t - z\|^2 + \frac{1}{t} \left[ \left(1 - \varepsilon - \frac{\eta}{2} AB^2\right) \sum_{s=0}^{t-1} \hat{R}(w_s) \right] \leq \frac{1}{2\eta t} \|w_0 - z\|^2 + \hat{R}(z).$$

Goal: show don't need to project.

Why: (a) no projection in practice (b) some other settings where projection isn't an option.

Pick  $t$ , & pick  $z$  &  $\varepsilon$  so that near optimal  $\hat{R}(z) \leq \frac{R^2}{2t}$

and near initialization  $\|w_0 - z\| \leq \frac{R}{2}$

Recall for ReLU:  $\varepsilon \leq \frac{1}{m^{\frac{1}{2}}} (4R^{\frac{4}{3}} + 2R \ln(\frac{1}{\delta})^{\frac{1}{2}})$ ,

$$\Rightarrow m \geq \left(\frac{1}{\varepsilon} (4R^{\frac{4}{3}} + 2R \ln(\frac{1}{\delta})^{\frac{1}{2}})\right)^2, \text{ then } \varepsilon \leq \frac{1}{4}$$

$$\Rightarrow \left(1 - \varepsilon - \frac{\eta}{2} AB^2\right) \geq \frac{1}{2}.$$

Apply theorem with  $U := \{ \|w - w_0\|^2 \leq (4\eta)R^2 : w \in \mathbb{R}^p \}$ .

Suppose  $\exists s$  with  $w_s - \eta \partial \hat{R}(w_s) \notin U$ , then

$$\frac{1}{2\eta t} \|w_{st+1} - z\|^2 \leq \frac{1}{2\eta t} \|z - w_0\|^2 + \hat{R}(z) \leq \frac{1}{2\eta t} \left(\frac{R}{2}\right)^2 + \frac{R^2}{2t} \leq \frac{R^2}{\varepsilon}$$

$$\begin{aligned} \text{Note also } \|w_{st+1} - w_0\|^2 &= \|w_{st+1} - z\|^2 + 2\langle w_{st+1} - z, z - w_0 \rangle + \|z - w_0\|^2 \\ &\leq 2\|w_{st+1} - z\|^2 + 2\|z - w_0\|^2 \end{aligned}$$

$$\Rightarrow \|w_{st+1} - w_0\|^2 \leq 2\|w_{st+1} - z\|^2 + 2\|z - w_0\|^2 \leq 4\eta R^2 + \frac{R^2}{2}$$

$\Rightarrow$  did not exit  $U$

$S$  was arbitrary  $\Rightarrow$  projection never encountered

$$\max_{S \subseteq U} \|w_s - z\|^2 \leq (4\eta)R^2$$