

# Lecture 27: even looser Rademacher bounds!

Theorem (Bartlett-Mendelson '02)  $\sigma_i(0) = 0$ ,  $\sigma_i$  is  $e$ -Lipschitz  
 $URad(\sum_{k=1}^n \sigma_i(W_i \sigma_{i-1}(W_{i-1} \dots \sigma_1(W_1, x) \dots)) : \|W_i\|_{1,\infty} \leq B)_{1X}$   
 $\leq \|X\|_{2,\infty} \sqrt{2 \ln d} (2eB)^L$

Rem (looseness). Proof doesn't allow  $(W_i - W_i(0))$ , therefore  $\|W_i\|_{1,\infty} \leq \frac{1}{\sqrt{n}}$  at int; super loose.

Proof. Define  $\mathcal{F}_i :=$  "output of nodes in layer  $i$ "  
 $= \begin{cases} \mathcal{F}_0 := \{x \mapsto x_j : j \in \{1, \dots, d\}\} \\ \mathcal{F}_{i+1} := \{x \mapsto \sigma_{i+1}(\sum_{j=1}^m v_j (g_j(x))) : m \geq 0, g_j \in \mathcal{F}_i, \|v\| \leq B\} \end{cases}$

IH:  $URad(\mathcal{F}_i|_X) \leq (2eB)^i \|X\|_{2,\infty} \sqrt{2 \ln d}$

Base case:  $URad(\mathcal{F}_0|_X) \leq (\sup_{j \in \{1, \dots, d\}} \|X_{:,j}\|_2) \sqrt{2 \ln |\mathcal{F}_0|} = \|X\|_{2,\infty} \sqrt{2 \ln d} = (2eB)^0 \|X\|_{2,\infty} \sqrt{2 \ln d}$   
*Martingale finite lemma*

Ind. step:  $URad(\mathcal{F}_{i+1}|_X) = URad(\sum_{k=1}^m \sigma(\sum_{j=1}^m v_j g_j(x)) : m \geq 0, \|v\|_1 \leq B, g_j \in \mathcal{F}_i|_X)$   
 $\leq e URad(\sum_{j=1}^m v_j g_j(x)) : m \geq 0, \|v\|_1 \leq 1, g_j \in \mathcal{F}_i|_X$

$\leq eB \cdot URad(\text{conv}(\mathcal{F}_i, U - \mathcal{F}_i)|_X)$   
 $= eB \cdot URad(\mathcal{F}_i \cup -\mathcal{F}_i|_X)$   
 $\leq 2eB URad(\mathcal{F}_i|_X)$   
 $\stackrel{\text{IH}}{\leq} (2eB)^{i+1} \|X\|_{2,\infty} \sqrt{2 \ln d}$   
*using  $\sigma(0) = 0$*  *union rule*

Remarks (a) source of looseness (a) not adapted to GD (searches over too many functions) (b) worst-cases between layers.

(2) Union rule typically misstated, & source of bugs in papers. Original defn  $URad$  was  $URad_{1,1}(V) = URad(V \cup -V)$

$$= \mathbb{E}_\varepsilon \sup_{u \in V} |\langle \varepsilon, u \rangle|$$

But  $URad_{1,1}(V, V_i) \leq \sum_i URad_{1,1}(V_i)$

Union holds  $URad$  under some conditions.

(3)  $URad(\text{conv}(V)) = \mathbb{E}_\varepsilon \sup_{u \in \text{conv}(V)} \langle \varepsilon, u \rangle = \mathbb{E}_\varepsilon \sup_{\alpha \in \Delta_n} \sup_{u \in V} \langle \varepsilon, \sum_{i=1}^n \alpha_i u_i \rangle$   
 $= \mathbb{E}_\varepsilon \sup_{u \in V} \langle \varepsilon, u \rangle$   
 $= URad(V)$

Theorem (Golowich - Rokhlin - Shamir '18).  $\sigma_L(0) = 0$ ,  $\sigma_L$  coordinate-wise & 2-Lip

[check my lecture notes if need 2-hono]

$$\text{Rad}(\sum x_i \mapsto \sigma_L(W_L \dots \sigma_L(W_1 x) \dots) : \|W\|_F \leq B \beta_{1,x}) \leq B^2 \|X\|_F (1 + \sqrt{2L \ln 2})$$

Proof ideas:

know this variant of Lipschitz peeling lemma:

Lemma (Talagrand - Ledoux),  $\tilde{L} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $\tilde{L}(0) = 0$ ,  $\tilde{L}$  coordinate-wise  $e$ -Lip

$$\mathbb{E}_\varepsilon \sup_{u \in V} \exp(\sum \varepsilon_i \tilde{L}(u_i)) \leq \mathbb{E}_\varepsilon \sup_{u \in V} \exp(e \sum \varepsilon_i u_i)$$

Proof: tricky case analysis.

$$X_0 = \begin{bmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad X_L^w = \sigma_L(X_{l-1}^w W_{l-1}^T)$$

$$\begin{aligned} \text{Rad}(F_L) &= \mathbb{E}_\varepsilon \sup_w \varepsilon^T X_L^w \\ &= \mathbb{E}_\varepsilon \frac{1}{t} \ln \sup_w \exp(t \varepsilon^T X_L^w) \\ &= \frac{1}{t} \ln \mathbb{E}_\varepsilon \sup_w \exp(t \varepsilon^T X_L^w) \end{aligned}$$

① induction:  $\mathbb{E}_\varepsilon \sup_w (t \|\varepsilon^T X_L^w\|) \leq \mathbb{E}_\varepsilon \sup_w (t B^L \|\varepsilon^T X_0\|)$ .  
many steps

②  $\|\varepsilon^T X_0\|$  is  $\|X\|_F^2$ -subGaussian &  $\mathbb{E} \|\varepsilon^T X_0\|_2 \leq \|X\|_F$ .

③ pick  $t := \sqrt{\frac{2 \ln 2}{B^{2L} \|X\|_F}}$ .

Theorem (\*) (Buckdett - Foster - Telgarsky '17)  $\sigma: \mathcal{D} = \sigma$ ,  $e_i$  - Lip, lot coordinate-wise

$$\text{URad} \left( \left\{ x \mapsto \sigma_x(w_1 \dots \sigma_x(w_k, x) \dots) : \begin{array}{l} \| (w_i - w_i(\mathcal{D}))^T \|_{2,1} \leq b_i \\ \| w_i \|_2 \leq s_i \end{array} \right\}_{1 \times k} \right) \\
 \leq \tilde{O} \left( \|K\|_F \left( \prod_i e_i s_i \right) \left( \sum_i \left( \frac{b_i}{s_i} \right)^{2/3} \right)^{3/2} \right).$$

Proof remarks.

① uses "covering numbers";

$C(u) \approx u$ ,  $|C|$  small

$$\mathbb{E}_\varepsilon \sup_{u \in \mathcal{U}} \langle \varepsilon, u \rangle = \mathbb{E}_\varepsilon \sup_{u \in \mathcal{U}} \left( \langle \varepsilon, u - C(u) \rangle + \langle \varepsilon, C(u) \rangle \right) \\
 + \int_{\mathcal{U}} \sup_{u \in \mathcal{U}} \|u - C(u)\| + \underbrace{\text{URad}(C)}_{\int \ln(C)}.$$

② (\*) has not been proved with Rademacher complexity.

Theorem (VC bound) or Rad

$$URad \left( \left\{ x \mapsto \text{sgn}(\sigma_L(W_L \dots \sigma_1(W_1 x + b_1) + \dots + b_L)) : \begin{matrix} p \text{ parameters} \\ L \text{ layers} \end{matrix} \right\} \Big| X \right) \\ \leq \sqrt{2n(1+VC(\text{---}))} \ln(n+1) \\ \text{where } VC(\text{---}) \leq 6pl \ln(pl).$$

Proof remarks

$$* URad(\text{sgn}(f)|_X) \leq \left( \sup_{u \in \text{sgn}(f)} \|u\|_2 \right) \sqrt{2n \ln |\text{sgn}(f)|_X} \\ \leq \sqrt{2n} \sqrt{\ln 2^n} = \sqrt{2n} \sqrt{n} = \sqrt{2} n$$


Sh(f; X)

$$\text{Sh}(f; n) = \sup_{|X| \leq n} |\text{sgn}(f)|_X$$

"shatter coefficient"

$$VC(f) := \sup \{ i \geq 0 : \text{sh}(f; i) = 2^i \}$$

$$\text{Sh}(f; n) \leq 1 + n^{1+VC(f)}$$

\* Activation matters (f convex-concave monotone & bounded activation) s.t.  $VC(\{x \mapsto \phi(w_1 \dots \phi(w_L x) \dots\}) = \infty$

\* Specific Rad comments:

$$x \mapsto \sigma_L(W_L \dots \sigma_1(W_1 x) \dots) \stackrel{w_L}{\circ} (w_{L-1} \dots (\stackrel{w_2}{\circ} (w_1 (\stackrel{w_1}{\circ} (W_1 x))))))$$

for a fixed  $S_L \dots S_1$ , then this is linear in x  
L-degree polynomial in w } hard to prove (Warren '62)

⇒ need to count  $(S_L^w \dots S_1^w)$ .

recursively refine partitions  $P_1, \dots, P_L$

where  $\forall S \in P_i \exists s' \in P_i$  s.t.  $S \subseteq S'$