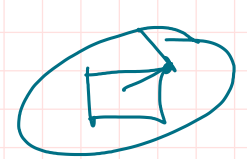


# Lecture 26: more Rademacher complexity

Definition.  $URad(V) = \mathbb{E}_{\epsilon \subseteq \mathbb{R}^n} \sup_{u \in V} \langle \epsilon, u \rangle$

Typical usage  
 $G_{1S} = \{g(z_1), \dots, g(z_n) : g \in G\} \subseteq \mathbb{R}^n$

$\epsilon_i \stackrel{i.i.d.}{\sim} \text{Unif}(\{-1, +1\})$   
 alternatively  $v_i \sim \mathcal{N}(0, 1)$  "Gaussian complexity"



Theorem (core Rademacher bound). Given  $G$  with  $g(z) \in [a, b]$  a.s.,  
 w/  $p_i \geq 1 - \delta \quad \forall g \in G$   
 $\mathbb{E} g(z_i) \leq \frac{1}{n} \sum_i g(z_i) + \frac{2}{n} URad(G_{1S}) + 3(b-a) \sqrt{\frac{\ln 2/\delta}{2n}}$

Lemma.  $URad(\{\sum_i \epsilon_i w_i \mapsto w^T x : \|w - w_0\| \leq R\}) \leq \|x\|_F \cdot R$

Rem. Regularization  $\Rightarrow$  good generalization. (Sufficient but not necessary.)

Proof.  
 $\mathbb{E}_{\epsilon} \sup_{\|w - w_0\| \leq R} \sum_i \epsilon_i w_i^T x_i$   
 $= \mathbb{E}_{\epsilon} \sup_{\|w - w_0\| \leq R} \sum_i \epsilon_i (w - w_0 + w_0)^T x_i$   
 $= \underbrace{\mathbb{E}_{\epsilon} \sum_i w_0^T x_i}_0 + \mathbb{E}_{\epsilon} \sup_{\|w - w_0\| \leq R} \sum_i \epsilon_i (w - w_0)^T x_i$   
 $= \mathbb{E}_{\epsilon} \sup_{\|v\| \leq R} \langle v, \sum_i \epsilon_i x_i \rangle = \mathbb{E}_{\epsilon} R \cdot \|\sum_i \epsilon_i x_i\|$   
 $\leq R \sqrt{\mathbb{E}_{\epsilon} \|\sum_i \epsilon_i x_i\|^2} = R \sqrt{\mathbb{E}_{\epsilon} (\sum_i \epsilon_i^2 \|x_i\|^2 + \sum_{i \neq j} \epsilon_i \epsilon_j \langle x_i, x_j \rangle)}$   
 $= R \|x\|_F$

Rem. essentially tight,  $URad(\{\sum_i \epsilon_i w_i \mapsto w^T x\}) \geq \frac{1}{\sqrt{2}} \|x\|_F \cdot R$   
 "Khintchine's inequality"

Lemma.  $\vec{l} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , where  $l_i$  is  $e$ -Lipschitz  
 then  $URad(\vec{l} \circ V) \leq e \cdot URad(V)$

Proof idea Core idea main Rademacher is proved via "symmetrization"  
 modern proof of this lemma de-symmetrizes then re-symmetrizes. //

Lemma (Massart finite lemma).  $URad(V) \leq \left( \sup_{u \in V} \|u\|_2 \right) \sqrt{2 \ln |V|}$

Proof idea. Subgaussian tricks on  $(\epsilon_1, \dots, \epsilon_n)$ . //

Remark. ① Agrees with union bound:

$$\sup_{g \in G} \left( \mathbb{E} g(z) - \frac{1}{n} \sum_i g(z_i) \right) \leq 3(b-a) \sqrt{\frac{\ln 2/\delta}{2n}} + \frac{2}{n} URad(G_{1S}) \approx \sqrt{\ln |G|}$$

② Many generalization bounds are proved via discretization + Massart lemma. //

Example (logistic regression).

logistic loss  $l(z) = \ln(1 + \exp(-z))$   
 Suppose  $\|w\| \leq R, \|x\| \leq 1$  a.s.  
 $\Rightarrow$  ①  $l$  is 1-Lip  
 ②  $l(y w^T x) \in [0, \ln 2]$  a.s.  
 ③  $\|x\|_F \leq \sqrt{n}$

$$\mathbb{E} l(w^T x_i) \leq \frac{1}{n} \sum_i l(w^T x_i) + \frac{2}{n} URad(\{l \circ \sum_i \epsilon_i x_i : \|w\| \leq R\}) + 3(b-a) \sqrt{\frac{\ln 2/\delta}{2n}}$$

$$l_i(f(x_i)) = l(y_i f(x_i)) \Rightarrow URad(\vec{l} \circ \{\sum_i \epsilon_i w_i \mapsto \dots\}) \leq 1 \cdot URad(\{\sum_i \epsilon_i w_i \mapsto \dots\}) \leq \|x\|_F \cdot R = R \sqrt{n}$$

$$\leq \frac{1}{n} \sum_i l(w^T x_i) + \frac{2R}{\sqrt{n}} + \frac{3(\ln 2)}{\sqrt{2n}} \sqrt{\ln \frac{2}{\delta}}$$

# Four deep network Rademacher bounds.

① [Bartlett - Mendelson '02.]  $\sigma_i$ : coordinate-wise  $e$ -Lipschitz,  $\sigma_i(0) = 0$

$$\text{URad} \left( \left\{ x \mapsto \sigma_L(W_L \sigma_{L-1} \dots \sigma_1(W_1 x) \dots) : \|W^T\|_{1,\infty} \leq B \right\} | X \right) \leq \|X\|_{2,\infty} (2eB)^2 \sqrt{2 \ln d}.$$

② [Golowich - Rakhtin - Shamir '18.]  $\sigma_i$ : coordinate-wise 1-Lipschitz, homogeneous

$$\text{URad} \left( \left\{ x \mapsto \sigma_L(W_L \dots \sigma_1(W_1 x)) : \|W_i\|_F \leq B \right\} | X \right) \leq \|X\|_F B^L (1 + \sqrt{2L \ln 2}).$$

③ [Bartlett - Foster - Telyarshy '17.]  $\sigma_i$ :  $e_i$ -Lipschitz

$$\text{URad} \left( \left\{ x \mapsto \sigma_L(W_L \dots \sigma_1(W_1 x) \dots) : \begin{array}{l} \|W_i - W_i(0)\|_{2,1} \leq b_i \\ \|W_i\|_2 \leq s_i \end{array} \right\} | X \right) \leq \tilde{O} \left( \|X\|_F \left( \prod_{i=1}^L s_i e_i \right) \left( \sum_{i=1}^L \left( \frac{b_i}{s_i} \right)^{1/3} \right)^{3/2} \right).$$

④ [Bartlett - '18.] ("VC dim bound")  $\text{ReLU} \subseteq \sum_{1, H}^n$

$$\text{URad} \left( \left\{ x \mapsto \text{sgn}(\sigma_L(W_L \dots \sigma_1(W_1 x) \dots)) : \begin{array}{l} p \text{ parameters} \\ L \text{ layers} \end{array} \right\} | X \right) \leq 30 \sqrt{np \cdot L \cdot (\ln pL) (\ln n)}$$

Remarks. \* All have strawman weaknesses (i.e., add garbage to weight matrices)

\* All are "vacuous".

\* Community efforts:

\* Identify "phenomena"

\* 2-layer case

\* Apply restrictions (e.g., Wei - Ma "Lipschitz Argumentation" no ReLU interesting)

\* Data & algorithmic assumptions