

# 10. VC dimension and Rademacher Complexity

Sunday, January 23, 2022

5:44 PM

§6 in [SS] w, Rademacher Complexity taken from §3.1 in Mohri et al.

We've covered ①  $|H| < \infty$  (restrictive!)

② HW: axis-aligned rectangles,  $|H| = \infty$ .  $X = \mathbb{R}^d$ ,  $\dim(H) = 2d$

Can we generalize?

We'll cover

① Rademacher Complexity (Mohri, and essentially used in later chapters of [SS])

Simple proofs, but computing R.C. may be impossible (eg, NP-Hard)  
especially if  $ERM_H$  is difficult to compute

② Growth Function

③ VC-dimension, a way to bound the growth function, and  
easier to compute or bound

④ Result: for binary classification, finite VC-dim is

necessary and sufficient for PAC learnability

"Fundamental Thm. of ML"

Rademacher Complexity (§3.1 Mohri, notation adapted a bit)

Will depend on  $H$  and loss function

$$\ell : H \times Z \rightarrow \mathbb{R} \text{ in [SS]}$$

$$\lambda : Y \times Y \rightarrow \mathbb{R} \text{ in [Mohri et al.]}$$

$$\text{e.g. } \ell(h, (x, y)) = \lambda(h(x), y)$$

we'll apply to a family of functions

$$F = \{ f : (x, y) \mapsto \ell(h, (x, y)) \mid h \in H \}$$

$$= \lambda \circ H$$

but it'll work for any family of functions  $F$ , not just  $\lambda \circ H$

$$F \subseteq \mathbb{R}^Z \quad Z = X \times Y$$

Idea

Rademacher Complexity (RC) measures the richness / expressiveness of  $\mathcal{F}$  by measuring how well it can fit noise

### Def Empirical Rademacher Complexity

$\mathcal{F}$  a family of fcn  $f: \mathcal{Z} \rightarrow [a, b]$ . Fix  $S = (z_1, \dots, z_m)$   
then empirical R.C. of  $\mathcal{F}$  w.r.t.  $S$  is

$$\hat{R}_S(\mathcal{F}) = \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i \cdot f(z_i) \right]$$

where  $\sigma = (\sigma_1, \dots, \sigma_m)$ ,  $\sigma_i$  iid Rademacher variables

i.e.  $\sigma_i = \begin{cases} +1 & \text{w.p. } .5 \\ -1 & \text{w.p. } .5 \end{cases}$  aka symmetric Bernoulli  
or uniform on  $\{-1, 1\}$

i.e., let  $f_S := \begin{bmatrix} f(z_1) \\ \vdots \\ f(z_m) \end{bmatrix} \in \mathbb{R}^m$  then

$$\hat{R}_S(\mathcal{F}) = \frac{1}{m} \mathbb{E}_{\sigma} \underbrace{\sup_{f \in \mathcal{F}} \langle \sigma, f_S \rangle}_{\text{best correlation w/ noise}}$$

i.e., best correlation w/ noise

Extremes:  $\mathcal{F} = \{f\}$ ,  $\hat{R}_S(\{f\}) = \frac{1}{m} \mathbb{E}_{\sigma} \langle \sigma, f_S \rangle = 0$ . Best possible

vs.  $\mathcal{F}$  = all functions, say  $f: \mathcal{Z} \rightarrow \{0, 1\}$ , then possible for some  $S$

for  $\{f_S : f \in \mathcal{F}\} = \{0, 1\}^m$

Then  $\sup_{f \in \mathcal{F}} \langle \sigma, f_S \rangle = m$ , so  $\hat{R}_S(\mathcal{F}) = \frac{1}{m} \mathbb{E}_{\sigma} m = \underline{\underline{1}}$

worst-possible  
(if  $[a, b] = [0, 1]$ )

### Def Rademacher Complexity (not "empirical")

$$R_m(\mathcal{F}) := \mathbb{E}_{S \sim \mathcal{Z}^m} \hat{R}_S(\mathcal{F})$$

Careful: [SS] uses different terminology:

we'll follow  
Mohri  
et al.

concept

$\hat{R}_S(\mathcal{F})$   
"Empirical R.C."

concept

$R_m(\mathcal{F}) = \mathbb{E}_{S \sim \mathcal{Z}^m} \hat{R}_S(\mathcal{F})$   
"R.C."

Shalev-Shwartz  
+ Ben-David

$R(F \circ S)$   
and  $F = \mathcal{L} \circ \mathcal{H}$   
"R.C."

$\mathbb{E}_S R(F \circ S)$   
(no special notation)  
"Expected R.C."

How to use?

Recall for uniform convergence,  $S$  was " $\epsilon$ -representative" if

$$\sup_{h \in \mathcal{H}} |L_D(h) - \hat{L}_S(h)| \leq \epsilon$$

(which implied ERM worked:  $L_D(\text{ERM}_{\mathcal{H}}(S)) \leq \epsilon + \min_{h \in \mathcal{H}} L_D(h)$ )

Something very similar is the "representativeness" of  $S$   
(w.r.t.  $\mathcal{H}, \mathcal{L}$ ) as

$$\text{Rep}_D((\mathcal{H}, \mathcal{L}), S) := \sup_{h \in \mathcal{H}} \underbrace{\mathbb{E}_{z \sim D} \mathcal{L}(h, z)}_{L_D(h)} - \underbrace{\frac{1}{m} \sum_{i=1}^m \mathcal{L}(h, z_i)}_{\hat{L}_S(h)}$$

or more generally

$$\underbrace{\Phi(S)}_{\text{in Mohri}} = \text{Rep}_D(F, S) = \sup_{f \in F} \underbrace{\mathbb{E}_{z \sim D} f(z)}_{\mathbb{E} f} - \underbrace{\frac{1}{m} \sum_{i=1}^m f(z_i)}_{\hat{\mathbb{E}}_S f}$$

want this small.  
clear  $= \sup (\mathbb{E} f - \hat{\mathbb{E}}_S f)$  in shorthand notation.

Intuitively,  $\hat{R}_S(F)$  is a reasonable estimate for  $\text{Rep}_D(F, S)$

Why? in  $\text{Rep}_D(F, S)$  we have  $\mathbb{E} f - \hat{\mathbb{E}}_S f$ . Split  $S = S_1 \cup S_2$  at random

estimate  $\mathbb{E} f - \hat{\mathbb{E}}_S f$  by  $\underbrace{\hat{\mathbb{E}}_{S_1} f - \hat{\mathbb{E}}_{S_2} f}_{\text{rewrite}}$

Let  $S_1 = \{z_i \in [m] : \sigma_i = +1\}$ ,  $S_2 = S \setminus S_1$ , and suppose  $|S_1| = m/2$  exactly

$$\begin{aligned} \text{then } \hat{\mathbb{E}}_{S_1} f - \hat{\mathbb{E}}_{S_2} f &= \frac{1}{m/2} \sum_{i \in S_1} f(z_i) - \frac{1}{m/2} \sum_{i \in S_2} f(z_i) \\ &= \frac{1}{m/2} \sum_i \sigma_i f(z_i) \end{aligned}$$

Thus taking a  $\sup_{f \in F} (\dots)$ , as we do in  $\text{Rep}_D$  and  $\hat{R}_S$ ,

we get  $\text{Rep}_D(F, S) \approx 2 \cdot \hat{R}_S(F)$

Now, let's be slightly more careful and formalize the above:

Lemma 26.2 (Mohri)  $\mathbb{E}_{S \sim \mathcal{D}^m} \text{Rep}_{\mathcal{D}}(\mathcal{F}, S) \leq 2 \cdot \mathbb{E}_{S \sim \mathcal{D}^m} \hat{\mathcal{R}}_S(\mathcal{F})$   
 $= 2 \cdot \mathcal{R}_m(\mathcal{F})$ .

proof:

$$\begin{aligned}
 \mathbb{E}_{S \sim \mathcal{D}^m} \text{Rep}_{\mathcal{D}}(\mathcal{F}, S) &:= \mathbb{E}_{S \sim \mathcal{D}^m} \sup_{f \in \mathcal{F}} \mathbb{E} f - \mathbb{E}_S^1 f \\
 &= \mathbb{E}_{S \sim \mathcal{D}^m} \sup_{f \in \mathcal{F}} \underbrace{\mathbb{E}_{S' \sim \mathcal{D}^m}}_{\mathbb{E} f} \left( \underbrace{\mathbb{E}_{S'} f}_{\text{no effect}} - \mathbb{E}_S^1 f \right) \\
 &\leq \mathbb{E}_{S, S'} \sup_{f \in \mathcal{F}} \left( \mathbb{E}_{S'}^1 f - \mathbb{E}_S^1 f \right) \quad \text{sup is sub-additive [see details later]} \\
 &:= \mathbb{E}_{S, S'} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m f(z_i') - f(z_i) \\
 &= \mathbb{E}_{S, S'} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i (f(z_i') - f(z_i)) \\
 &\quad \sigma_i = 1 \text{ ok} \quad \text{for any } \sigma_i \\
 &\quad \sigma_i = -1 \text{ flips } z_i', z_i \dots \text{ but same distribution} \\
 &= \mathbb{E}_{S, S', \sigma} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i (f(z_i') - f(z_i)) \quad \text{true } \forall \sigma \text{ so true for } \mathbb{E} \\
 &\leq \mathbb{E}_{S', \sigma} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i') + \mathbb{E}_{S, \sigma} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i (-f(z_i)) \\
 &\quad \text{Sup}(x+y) \leq \text{Sup}(x) + \text{Sup}(y) \\
 &= 2 \mathbb{E}_{S, \sigma} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i) = 2 \mathbb{E}_S \hat{\mathcal{R}}_S(\mathcal{F}) \\
 &= 2 \mathcal{R}_m(\mathcal{F}) \quad \square
 \end{aligned}$$

sub-additive:

$$\forall a, b \quad g(a, b) \leq \sup_b g(a, b')$$

$$\text{so } \forall b \quad \mathbb{E}_a g(a, b) \leq \mathbb{E}_a \sup_b g(a, b')$$

so

$$\sup_b \mathbb{E}_a g(a, b) \leq \mathbb{E}_a \sup_b g(a, b')$$