# 9. 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 $|\mathcal{H}| < \infty$    ( restrictive ! )

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

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 $ERM_{\mathcal{H}}$ 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 $\mathcal{H}$ and loss function

$\ell : \mathcal{H} \times Z \to \mathbb{R}$   in [SS]

$\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$   in [Mohri et al.]

e.g. $\ell ( h, (x,y) ) = \ell ( h(x), y )$

we'll apply to a family of functions

$$\mathcal{F} = \left\{ f : (x,y) \longmapsto \ell ( h, (x,y) ) \ \forall \ h \in \mathcal{H} \right\}$$
$$= \ell \circ \mathcal{H}$$

but it'll work for any family of functions $\mathcal{F}$, not just $\ell \circ \mathcal{H}$

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

Idea

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

## Def Empirical Rademacher Complexity

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

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

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

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

i.e., best correlation w/ noise

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

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

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

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

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

## Def Rademacher Complexity (not "empirical")

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

Careful: [SS] uses different terminology:

we'll follow Mohri

| Concept | concept |
|---|---|
| Mohri et al. | |
| $\hat{R}_S(F)$ | $R_m(F) = \mathbb{E}_{S \sim \mathcal{D}^m} R_S(F)$ |
| 'Empirical R.C.' | 'R.C.' |

$R(\mathcal{F} \circ S)$

and $\mathcal{F} = \ell \circ \mathcal{H}$

"R.C."

$\mathbb{E}_S \; R(\mathcal{F} \circ S)$

(no special notation)

"Expected R.C."

## How to use?

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

$$\sup_{h \in \mathcal{H}} \left| L_D(h) - \hat{L}_S(h) \right| \le \varepsilon$$

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

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

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

or more generally

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

$\Phi(S)$ in Mohri

want this small (clear)

$$= \sup \left( \mathbb{E} f - \hat{\mathbb{E}}_S f \right) \quad \text{in shorthand notation.}$$

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

Why? in $\text{Rep}_D(\mathcal{F}, S)$ 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 $S_1 = \{ i \in [m] : \sigma_i = +1 \}$, $S_2 = S \setminus S_1$, and suppose $|S_1| = m/2$ exactly

then $\hat{\mathbb{E}}_{S_1} f - \hat{\mathbb{E}}_{S_2} f = \frac{1}{m/2} \sum_{i \in S_1} f(z_i) - \frac{1}{m/2} \sum_{i \in S_2}' f(z_i)$

$$= \frac{1}{m/2} \sum_i' \sigma_i f(z_i)$$

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

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

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

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

$$= 2 \cdot \mathcal{R}_m(F).$$

proof:

$$\mathbb{E}_{S \sim D^m} \text{Rep}_D(F, S) := \mathbb{E}_{S \sim D^m} \sup_{f \in F} \mathbb{E}f - \hat{\mathbb{E}}_S f$$

$$= \mathbb{E}_{S \sim D^m} \sup_{f \in F} \underbrace{\mathbb{E}_{S' \sim D^m}}_{\mathbb{E}f} \left( \overbrace{\hat{\mathbb{E}}_{S'} f}^{\text{no effect}} - \hat{\mathbb{E}}_S f \right)$$

$$\leq \mathbb{E}_{S, S'} \sup_{f \in F} \left( \hat{\mathbb{E}}_{S'} f - \hat{\mathbb{E}}_S f \right) \qquad \text{\color{green}{sup is sub-additive [see details later]}}$$

$$:= \mathbb{E}_{S, S'} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} f(z_i') - f(z_i)$$

$$= \mathbb{E}_{S, S'} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} \sigma_i \left( f(z_i') - f(z_i) \right) \qquad \text{\color{orange}{for any } \sigma_i}$$

$$\text{\color{orange}{$\sigma_i = 1$ ok}}$$
$$\text{\color{orange}{$\sigma_i = -1$ flips $z_i', z_i$ ... but same distribution}}$$
$$\text{\color{orange}{true $\forall \sigma$ so true for $\mathbb{E}$}}$$

$$= \mathbb{E}_{S, S', \sigma} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} \sigma_i \left( f(z_i') - f(z_i) \right)$$

$$\leq \mathbb{E}_{S', \sigma} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} \sigma_i f(z_i') + \mathbb{E}_{S, \sigma} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} \sigma_i(z_i)$$

$$\text{\color{green}{sup}(x+y) \leq \text{sup}(x) + \text{sup}(y)}$$

$$= 2 \mathbb{E}_{\sigma, S} \sup_{f \in F} \frac{1}{m} \sum_{i=1}^{m} \sigma_i f(z_i) = 2 \mathbb{E}_S \hat{\mathcal{R}}_S(F)$$

$$= 2 \mathcal{R}_m(F) \quad \square$$

Sub-additive:

$$\forall a, b \qquad g(a, b) \leq \sup_{b'} g(a, b')$$

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

so

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