

12. More on Rademacher Complexity

Wednesday, January 26, 2022 9:36 AM

$$\text{Recall } \hat{\mathcal{R}}_s(\mathcal{F}) = \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sigma_i f(z_i)$$

or, to generalize (simplify): Rademacher Complexity of a set:

$$\underline{\text{def}} \quad \hat{\mathcal{R}}(A) = \frac{1}{m} \mathbb{E}_{\sigma} \sup_{a \in A} \langle a, \sigma \rangle, \quad A \subseteq \mathbb{R}^m$$

$$\left(\begin{array}{l} \text{By setting } A = \left\{ \begin{bmatrix} f(z_1) \\ f(z_m) \end{bmatrix} : f \in \mathcal{F} \right\} \\ \text{we recover } \hat{\mathcal{R}}(A) = \hat{\mathcal{R}}_s(\mathcal{F}) \\ \therefore \hat{\mathcal{R}}_s(\mathcal{F}) := \hat{\mathcal{R}}(\mathcal{F} \circ S) \end{array} \right)$$

N.B. Very related ideas have been discovered / proposed in other subfields

e.g. Gaussian Width of a set A is $G(A) = \mathbb{E}_{g \sim N(0, I_m)} \sup_{a \in A} \langle g, a \rangle$

Also "Local Rademacher Complexity" sometimes in abs. value

$$\text{and Statistical Dimension} = \mathbb{E}_{g \sim N(0, I_m)} \| \text{Proj}_A(g) \|^2$$

Calculus / tools

$$\underline{\text{Lemma 26.2 [SS]}} \quad \hat{\mathcal{R}}(c \cdot A + a_0) \leq |c| \cdot \hat{\mathcal{R}}(A)$$

"shift-invariant", scales as expected

$$\underline{\text{Lemma 26.7 [SS]}} \quad \hat{\mathcal{R}}(\text{conv}(A)) = \hat{\mathcal{R}}(A)$$

where $\text{conv}(A) = \text{convex hull of } A$
 $= \text{smallest convex set containing } A$
 $= \text{intersection of all convex sets containing } A$
 $= \left\{ b = \sum_{i=1}^p \alpha_i \cdot a_i \text{ s.t. } p \in \mathbb{N}, a_i \in A, \alpha_i \geq 0, \sum \alpha_i \leq 1 \right\}$

$$\underline{\text{Lemma 26.8 [SS] Massart Lemma}} \quad \hat{\mathcal{R}}(A) \leq \max_{a \in A} \|a - \bar{a}\| \cdot \frac{1}{m} \sqrt{2 \cdot \log(N)}$$

Pascal Massart (b. 1958)
 $\bar{a} := \frac{1}{N} \sum_{i=1}^N a_i$, a_i is center of mass
 (w.l.o.g can shift to 0)

e.g. if $A = l \circ h \circ S$,
and l is bounded, then $\max \|l \circ h\|$ is bounded
and $\hat{R}_S(l \circ h) = O(\sqrt{n} \log(1/\delta))$. Combine w/ thm to bound L_D

Lemma 26.9 [SS] If $\phi: \mathbb{R}^m \rightarrow \mathbb{R}^n$ is ρ -Lipschitz in each component $i \in [m]$
then $\hat{R}(\phi \circ A) \leq \rho \cdot \hat{R}(A)$

Using these tools you can sometimes compute / bound R for simple/nice enough sets

Ex $H_2 := \{x \mapsto \langle \omega, x \rangle : \|\omega\|_2 \leq 1\}$, an ∞ set of linear classifiers

By Lemma 26.10 [SS], if $S_x = \{x_1, \dots, x_m\} \subseteq \mathbb{R}^n$

$$\hat{R}(H_2 \circ S_x) \leq \sqrt{m} \cdot \max_i \|x_i\|_2$$

Ex $H_1 := \{x \mapsto \langle \omega, x \rangle : \|\omega\|_1 \leq 1\}$

Lemma 26.11 [SS],

$$\hat{R}(H_1 \circ S_x) \leq \sqrt{m} \max_i \sqrt{2 \cdot \log(2n)} \|x_i\|_\infty$$

(we might do in detail for HW)

Perspective

R.C. used for uniform convergence ~ 2000
(often tools have been around, but only "popular" once ML discovers)

R.C. mainly for classification

- [- For regression, if loss is bounded (less common), R.C. can work
- otherwise need specialized measures

Alternatives: (ideas like VC-dim are inherently combinatorial and won't work)

- Pseudo-dimension (like VC-dim, needs boundedness)
- Fat-shattering dim. (still needs boundedness)
see Mohri for brief discussion

Another notion of complexity:

Covering Numbers

§ 27 [SS], § 3.5 Mohri

Def The $N_p(\varepsilon, A)$ covering number of a set A is the minimum number of points needed to cover A , i.e., s.t. A is contained in the union of ε -balls centered at these points

$$B_\varepsilon(y) = \{x : \|x-y\|_p \leq \varepsilon\}$$

(If p not specified, assume $p=2$)

related to... packing numbers $M_p(\varepsilon, A)$, the max # balls you can fit in A w/o overlapping. of radius ε wrt ℓ^p norm

$$\text{Fact: } M_p(2\varepsilon, A) \leq N_p(\varepsilon, A) \leq M_p(\varepsilon, A)$$



We'll prove shortly that you can bound $\hat{R}(A)$ by $N_2(\varepsilon, A)$
 "Dudley's Thm" ('67, '87) ... and can bound N_2 by VC-dim.
 uses "Chaining technique" which we'll cover

Facts

$\log(\text{covering number})$ is sometimes called "metric entropy"

$$\cdot N_2(\varepsilon, cA + a_0) = N_2(c\varepsilon, A) \quad \text{translation / scaling}$$

$$\cdot \phi \text{ componentwise Lipschitz, then } N_2(\varepsilon, \phi \circ A) \leq N_2(\varepsilon/\rho, A) \quad \text{w/ constant } \rho$$

$$\cdot \underline{\text{Lemma 2.1}} \quad [\text{SS}] \quad \text{If } A \subseteq \mathbb{R}^m \text{ is in a d-dimensional subspace and } A \subseteq B_c(o) \quad (\text{i.e. radius } \leq c) \quad \text{then } N_2(\varepsilon, A) \leq \left(2c \frac{\sqrt{d}}{\varepsilon}\right)^d$$

$$\cdot \underline{\text{Lemma 2.2}} \quad [\text{Woodruff}] \quad A \subseteq \mathbb{R}^m \text{ is in a d-dim. subspace and } \forall a \in A, \|a\|_2 = 1 \quad (\text{i.e. subspace } \cap \text{sphere}) \\ \text{then } N_2(\varepsilon, A) \leq (1 + 2/\varepsilon)^d$$

Proof (fin)

wlog work in unit sphere inside \mathbb{R}^d .

We'll find a covering (maybe not optimal) in a greedy fashion that will at least be maximal! nothing redundant

Let C be the centers, so covering = $\bigcup_{c \in C} B_\varepsilon(c)$

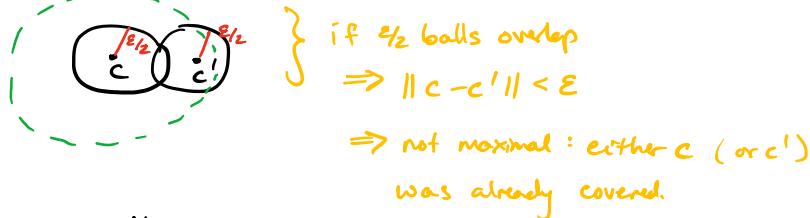
Pick $c \in A$ arbitrarily. Check if A is covered.

- If so, stop
- If not, $\exists c' \in A$ that hasn't been covered
(ie. not with ε of any $c \in C$). Add c' to C

Repeat. This is maximal.

(This is the intuitive way) How big is $|C|$?

Balls of radius ε may overlap, but balls of radius $\varepsilon/2$ (around centers) cannot overlap, otherwise non maximal.



Furthermore, all these $\varepsilon/2$ balls are within $1 + \varepsilon/2$ of origin
(centers on unit sphere)

Recall in dimension d , volume (ball radius r) $\propto r^d$

By disjointness of $\varepsilon/2$ balls,

$$\begin{aligned} 1D: & 2r \\ 2D: & \pi r^2 \\ 3D: & 4/3\pi r^3 \end{aligned}$$

$$\begin{aligned} \text{Volume (all } \varepsilon/2 \text{ balls)} &= \sum_{\text{all } \varepsilon/2 \text{ balls}} \text{Vol.}(\varepsilon/2 \text{ ball}) \\ &= |C| \cdot \text{const.} \cdot (\varepsilon/2)^d \end{aligned}$$

but also

$$\text{Vol. (all } \varepsilon/2 \text{ balls)} \leq \text{Vol} (1 + \varepsilon/2 \text{ ball}) = \text{const.} (1 + \varepsilon/2)^d$$

$$\begin{aligned} \Rightarrow |C| \cdot (\varepsilon/2)^d &\leq (1 + \varepsilon/2)^d, \quad |C| \leq \left(\frac{1 + \varepsilon/2}{\varepsilon/2}\right)^d = (1 + 2/\varepsilon)^d \\ &\Rightarrow N_2(\varepsilon, A) \leq (1 + 2/\varepsilon)^d. \quad \square \end{aligned}$$

(seems loose since we didn't exploit $\|a\|=1$ only $\|a\| \leq 1$,
but in high dim. the surface is where all the volume is, so not a big deal)