

# COMS 4773: Rademacher complexity

Daniel Hsu

March 11, 2024

## 1 Uniform convergence, again

Recall the uniform convergence theorem.

**Theorem 1.** *Let  $\mathcal{F} \subset \{-1, 1\}^{\mathcal{X}}$  have VC dimension  $d < \infty$ . Let  $\mu$  be a probability distribution on  $\mathcal{X}$ , and let  $\mu_n$  be the empirical distribution based on an iid sample from  $\mu$  of size  $n$ . Then for any  $\epsilon > 0$ ,*

$$\Pr\left(\sup_{f \in \mathcal{F}} \mu_n f - \mu f \geq \epsilon\right) \leq \exp\left(-\frac{n\epsilon^2}{8}\right) + 2\binom{n}{\leq d} \exp\left(-\frac{n\epsilon^2}{32}\right).$$

*This implies that, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,*

$$\sup_{f \in \mathcal{F}} \mu_n f - \mu f \leq O\left(\sqrt{\frac{d \log n + \log(1/\delta)}{n}}\right).$$

A different way to prove Theorem 1 starts by using McDiarmid's inequality. Let  $X_1, \dots, X_n$  be an iid sample from  $\mu$ , and let  $\mathcal{F} \subset \{-1, 1\}^{\mathcal{X}}$  be a function class. For any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\sup_{f \in \mathcal{F}} \mu_n f - \mu f \leq \mathbb{E}\left[\sup_{f \in \mathcal{F}} \mu_n f - \mu f\right] + \sqrt{\frac{2 \log(1/\delta)}{n}}.$$

McDiarmid's inequality applies because the random variable on the left-hand side satisfies the  $(c_1, \dots, c_n)$ -bounded differences property with  $c_i = 2/n$  for all  $i$ . So the main task is to bound the expectation on the right-hand side.

Let  $\mu'_n$  be the empirical distribution on an independent iid sample of size  $n$ ,  $X'_1, \dots, X'_n$  (the ghost sample). Instead of using conditional expectation notations, we shall write  $\mathbb{E}$  for expectation with respect to  $X_{1:n}$ , and we write  $\mathbb{E}'$  for expectation with respect to  $X'_{1:n}$ . Then  $\mu f = \mathbb{E}'[\mu'_n f]$ , and therefore

$$\sup_{f \in \mathcal{F}} \mu_n f - \mu f = \sup_{f \in \mathcal{F}} \mathbb{E}'[\mu_n f - \mu'_n f] \leq \mathbb{E}'\left[\sup_{f \in \mathcal{F}} \mu_n f - \mu'_n f\right]$$

where the inequality follows by Jensen's inequality and the convexity of the supremum of affine functions. So we have

$$\mathbb{E}\left[\sup_{f \in \mathcal{F}} \mu_n f - \mu f\right] \leq \mathbb{E}\mathbb{E}'\left[\sup_{f \in \mathcal{F}} \mu_n f - \mu'_n f\right].$$

Letting  $\sigma = (\sigma_1, \dots, \sigma_n)$  be a Rademacher random vector, we also have

$$\mathbb{E} \mathbb{E}' \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu'_n f \right] = \mathbb{E} \mathbb{E}' \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(X_i) - f(X'_i)) \right]$$

where  $\mathbb{E}_\sigma$  denotes expectation with respect to  $\sigma$ . Since  $X_{1:n}$  and  $X'_{1:n}$  have the same distribution, and using the symmetry of  $\sigma$ ,

$$\begin{aligned} \mathbb{E} \mathbb{E}' \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(X_i) - f(X'_i)) \right] &\leq \mathbb{E} \mathbb{E}' \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) + \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (-\sigma_i) f(X'_i) \right] \\ &= 2 \mathbb{E} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) \right] \\ &= 2 \mathbb{E} \mathbb{E}_\sigma \left[ \sup_{v \in \mathcal{F}(X_{1:n})} \langle \sigma, v \rangle_n \right]. \end{aligned}$$

(Recall our notations for empirical inner product  $\langle u, v \rangle_n = \frac{1}{n} \sum_{i=1}^n u_i v_i$  and empirical norm  $\|v\|_n = \sqrt{\langle v, v \rangle_n}$ .) Since each  $v \in \mathcal{F}(X_{1:n})$  has  $\|v\|_n = 1$ , it follows by Massart's lemma that

$$\mathbb{E}_\sigma \left[ \sup_{v \in \mathcal{F}(X_{1:n})} \langle \sigma, v \rangle_n \right] \leq \sqrt{\frac{2 \log |\mathcal{F}(X_{1:n})|}{n}}.$$

So, if  $\mathcal{F}$  has VC dimension  $d < \infty$ , we conclude by Sauer's lemma that

$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right] \leq 2 \sqrt{\frac{2 \log \binom{n}{\leq d}}{n}}.$$

## 2 Rademacher complexity

Going back a few steps in this development, we have the inequality

$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right] \leq 2 \mathbb{E} [\text{Rad}_n(\mathcal{F}(X_{1:n}))]. \quad (1)$$

where, for any  $V \subseteq \mathbb{R}^n$ ,

$$\text{Rad}_n(V) = \mathbb{E}_\sigma \left[ \sup_{v \in V} \langle \sigma, v \rangle_n \right].$$

The quantity  $\mathbb{E}[\text{Rad}_n(\mathcal{F}(X_{1:n}))]$  is called the *(one-sided) Rademacher complexity*<sup>1</sup> of  $\mathcal{F}$  (which also depends on  $\mu$  and  $n$ ).<sup>2</sup> The Rademacher complexity measures how well functions from  $\mathcal{F}$  can be used to “correlate” the sample  $X_{1:n}$  with random signs.

<sup>1</sup>The one-sided Rademacher complexity is somewhat non-standard, but it is more convenient in some technical respects and sufficient for our purposes. The usual (two-sided) Rademacher complexity of  $\mathcal{F}$  is defined with  $|\langle \sigma, v \rangle_n|$  in place of  $\langle \sigma, v \rangle_n$  in the definition of  $\text{Rad}_n$ .

<sup>2</sup>Sometimes  $\text{Rad}_n(\mathcal{F}(X_{1:n}))$  itself is called the *empirical Rademacher complexity* of  $\mathcal{F}$ .

- A “complex” class is one that is able to make this correlation  $\text{Rad}_n(\mathcal{F}(X_{1:n}))$  large (in expectation with respect to  $X_{1:n}$ ). For example, the set of *all*  $\{-1, 1\}$ -valued functions on  $\mathcal{X}$  has  $\text{Rad}_n(\mathcal{F}(X_{1:n})) = 1$ .
- A class that contains only a single function (which should be considered “simple” by any measure...) has  $\text{Rad}_n(\mathcal{F}(X_{1:n})) = 0$ .

Note that Rademacher complexity is well-defined not just for  $\{-1, 1\}$ -valued functions, but for any class of functions real-valued functions (although some normalization is needed to make the quantity meaningful). Another feature of Rademacher complexity is that it is sensitive to the data distribution  $\mu$ . These two “features” of Rademacher complexity distinguish it from VC dimension.

### 3 Properties of (empirical) Rademacher complexity

**Proposition 1.** *Let  $A$  and  $B$  be subsets of  $\mathbb{R}^n$ . Then the following hold.*

1. If  $A \subseteq B$ , then  $\text{Rad}_n(A) \leq \text{Rad}_n(B)$ .
2.  $\text{Rad}_n(A + B) = \text{Rad}_n(A) + \text{Rad}_n(B)$ .
3.  $\text{Rad}_n(cA) = |c| \text{Rad}_n(A)$ .
4.  $\text{Rad}_n(\text{conv}(A)) = \text{Rad}_n(A)$ .
5. (Lipschitz contraction.) Let  $\phi_1, \dots, \phi_n$  be  $L$ -Lipschitz  $\mathbb{R}$ -valued functions on a domain  $D \subseteq \mathbb{R}$ : i.e., for each  $i \in [n]$ ,

$$\phi_i(t) - \phi_i(t') \leq L|t - t'| \quad \text{for all } t, t' \in D.$$

Define

$$\phi(A) = \{(\phi_1(a_1), \dots, \phi_n(a_n)) : (a_1, \dots, a_n) \in A\}.$$

If  $A \subseteq D^n$ , then

$$\text{Rad}_n(\phi(A)) \leq L \text{Rad}_n(A).$$

*Proof.* 1. Since  $A \subseteq B$ , we have  $\{\langle \sigma, v \rangle_n : v \in A\} \subseteq \{\langle \sigma, v \rangle_n : v \in B\}$ ; since a supremum over a set can only stay the same or increase by adding more vectors to the set, it follows that  $\sup_{v \in A} \langle \sigma, v \rangle_n \leq \sup_{v \in B} \langle \sigma, v \rangle_n$  for all  $\sigma$ , and so  $\text{Rad}_n(A) \leq \text{Rad}_n(B)$ .

2. Since  $A + B = \{a + b : a \in A, b \in B\}$ , it follows that  $\sup_{v \in A+B} \langle \sigma, v \rangle_n = \sup_{a \in A} \sup_{b \in B} \langle \sigma, a + b \rangle_n = \sup_{a \in A} \langle \sigma, a \rangle_n + \sup_{b \in B} \langle \sigma, b \rangle_n$  for all  $\sigma$ , so  $\text{Rad}_n(A + B) = \text{Rad}_n(A) + \text{Rad}_n(B)$ .

3. If  $c \geq 0$ , then  $\sup_{v \in cA} \langle \sigma, v \rangle_n = \sup_{a \in A} \langle \sigma, ca \rangle_n = \sup_{a \in A} |c| \langle \sigma, a \rangle_n = |c| \sup_{a \in A} \langle \sigma, a \rangle_n$  for all  $\sigma$ . If  $c < 0$ , then  $\sup_{v \in cA} \langle \sigma, v \rangle_n = \sup_{a \in A} \langle \sigma, ca \rangle_n = \sup_{a \in A} |c| \langle -\sigma, a \rangle_n = |c| \sup_{a \in A} \langle -\sigma, a \rangle_n$  for all  $\sigma$ . In either case,  $\sigma$  and  $-\sigma$  have the same distribution, so we conclude that  $\text{Rad}_n(cA) = |c| \text{Rad}_n(A)$ .

4. By definition, for any  $v \in \text{conv}(A)$ , we can write  $v = c_1 v_1 + \dots + c_k v_k$  for some  $k \in \mathbb{N}$ , some  $c = (c_1, \dots, c_k) \in \Delta^{k-1}$ , and some  $v_1, \dots, v_k \in A$ . For such  $v$ , we have  $\langle \sigma, v \rangle_n = \sum_{i=1}^k c_i \langle \sigma, v_i \rangle_n \leq \max_{i \in [k]} \langle \sigma, v_i \rangle_n$ . So  $\sup_{v \in \text{conv}(A)} \langle \sigma, v \rangle_n \leq \sup_{a \in A} \langle \sigma, a \rangle_n$  for all  $\sigma$ , and therefore  $\text{Rad}_n(\text{conv}(A)) \leq \text{Rad}_n(A)$ . Since  $A \subseteq \text{conv}(A)$ , it also follows that  $\text{Rad}_n(A) \leq \text{Rad}_n(\text{conv}(A))$  by the first property above.

5. We show that by replacing  $\phi_1(\cdot)$  with  $L\cdot$  can only increase the Rademacher average. Write  $\mathbb{E}_{\sigma_1}$  for expectation over  $\sigma_1$  only (conditioning on  $\sigma_2, \dots, \sigma_n$ ). Then

$$\begin{aligned}
& \mathbb{E}_{\sigma_1} \left[ \sup_{v \in \phi(A)} \langle \sigma, v \rangle_n \right] \\
&= \mathbb{E}_{\sigma_1} \left[ \sup_{a \in A} \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_i(a_i) \right] \\
&= \frac{1}{2} \left( \sup_{a \in A} \frac{1}{n} \phi_1(a_1) + \underbrace{\frac{1}{n} \sum_{i=2}^n \sigma_i \phi_i(a_i)}_B \right) + \frac{1}{2} \left( \sup_{a' \in A} -\frac{1}{n} \phi_1(a'_1) + \underbrace{\frac{1}{n} \sum_{i=2}^n \sigma_i \phi_i(a'_i)}_{B'} \right) \\
&= \sup_{a, a' \in A} \frac{1}{2n} \phi_1(a_1) - \frac{1}{2n} \phi_1(a'_1) + \frac{1}{2} (B + B') \\
&\leq \sup_{a, a' \in A} \frac{L}{2n} |a_1 - a'_1| + \frac{1}{2} (B + B') \\
&= \sup_{a, a' \in A} \frac{L}{2n} (a_1 - a'_1) + \frac{1}{2} (B + B') \\
&= \frac{1}{2} \left( \sup_{a \in A} \frac{1}{n} L a_1 + \frac{1}{n} \sum_{i=2}^n \sigma_i \phi_i(a_i) \right) + \frac{1}{2} \left( \sup_{a' \in A} -\frac{1}{n} L a'_1 + \frac{1}{n} \sum_{i=2}^n \sigma_i \phi_i(a'_i) \right) \\
&= \mathbb{E}_{\sigma_1} \left[ \sup_{a \in A} \frac{1}{n} \sigma_1 (L a_1) + \frac{1}{n} \sum_{i=2}^n \sigma_i \phi_i(a_i) \right].
\end{aligned}$$

Repeat this for  $\phi_2, \dots, \phi_n$  to get

$$\mathbb{E}_{\sigma} \left[ \sup_{v \in \phi(A)} \langle \sigma, v \rangle_n \right] \leq \mathbb{E}_{\sigma} \left[ \sup_{v \in LA} \langle \sigma, v \rangle_n \right]$$

(where  $LA = \{La : a \in A\}$ ). Then apply the third property to prove the claim.  $\square$

## 4 Lower bound

The following is a complement to (1).

**Proposition 2.** For any  $\mathcal{F} \subset [-1, 1]^{\mathcal{X}}$  and any probability distribution  $\mu$  on  $\mathcal{X}$ ,

$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right] + \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu f - \mu_n f \right] \geq \mathbb{E} \text{Rad}_n(\mathcal{F}(X_{1:n})) - \sup_{f \in \mathcal{F}} |\mu f| \frac{1}{\sqrt{n}}.$$

Together, (1) and Proposition 2 show that Rademacher complexity essentially characterizes distribution-specific uniform convergence.

*Proof of Proposition 2.*

$$\begin{aligned}
\text{Rad}_n(\mathcal{F}(X_{1:n})) &= \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) + \frac{1}{n} \sum_{i=1}^n \sigma_i \mu f \right] \\
&\leq \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) \right] + \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mu f \right] \\
&\leq \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) \right] + \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} |\mu f| \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \right| \right] \\
&= \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) \right] + \sup_{f \in \mathcal{F}} |\mu f| \mathbb{E}_\sigma \left[ \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \right| \right] \\
&\leq \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) \right] + \sup_{f \in \mathcal{F}} |\mu f| \frac{1}{\sqrt{n}}.
\end{aligned}$$

Now let  $X'_1, \dots, X'_n$  be an independent iid sample from  $\mu$ , and write  $\mathbb{E}'$  for expectation with respect to  $X'_{1:n}$ . Then we have

$$\begin{aligned}
\mathbb{E} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mu f) \right] &= \mathbb{E} \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mathbb{E}' f(X'_i)) \right] \\
&\leq \mathbb{E} \mathbb{E}' \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - f(X'_i)) \right] \\
&= \mathbb{E} \mathbb{E}' \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i(f(X_i) - \mathbb{E} f(X_i)) - \sigma_i(f(X'_i)) - \mathbb{E}' f(X'_i) \right] \\
&= \mathbb{E} \mathbb{E}' \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E} f(X_i)) - (f(X'_i)) - \mathbb{E}' f(X'_i) \right] \\
&\leq \mathbb{E} \mathbb{E}' \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E} f(X_i)) + \sup_{f \in \mathcal{F}} -(f(X'_i)) - \mathbb{E}' f(X'_i) \right] \\
&= \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right] + \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu f - \mu_n f \right].
\end{aligned}$$

We conclude that

$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right] + \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu f - \mu_n f \right] \geq \mathbb{E} \text{Rad}_n(\mathcal{F}(X_{1:n})) - \sup_{f \in \mathcal{F}} |\mu f| \frac{1}{\sqrt{n}}. \quad \square$$

A simple corollary of Proposition 2 is that, for any  $\mathcal{F} \subset [-1, 1]^\mathcal{X}$  and any probability distribution  $\mu$  on  $\mathcal{X}$ ,

$$\max \left\{ \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu_n f - \mu f \right], \mathbb{E} \left[ \sup_{f \in \mathcal{F}} \mu f - \mu_n f \right] \right\} \geq \frac{1}{2} \mathbb{E} \text{Rad}_n(\mathcal{F}(X_{1:n})) - \sup_{f \in \mathcal{F}} |\mu f| \frac{1}{2\sqrt{n}}.$$

Notice that if  $\mathcal{F}$  is closed under negation (i.e.,  $\mathcal{F} = \mathcal{F} \cup (-\mathcal{F})$ ), then both terms in the max are the same. But it is not generally possible to replace the max with min.<sup>3</sup>

---

<sup>3</sup>To see this, consider the class  $\mathcal{F}$  of all characteristic functions  $f_S(x) = \mathbf{1}\{x \in S\}$  of finite subsets of  $[0, 1]$ , and let  $\mu$  be the uniform distribution on  $[0, 1]$ . The Rademacher complexity of  $\mathcal{F}$  is  $1/2$ , but  $\mathbb{E}[\sup_{f \in \mathcal{F}} \mu f - \mu_n f] = 0$ .