

## Lecture 12: Uniform bounds of generalization gap

December 6, 2023

Lecturer: Lei Wu

Scribe: Lei Wu

## Reading

- Section 26 and 27 of [Shalev-Shwartz and Ben-David, 2014].

## 1 Setup

Let  $z = (x, y)$ ,  $\ell_h(z) = \ell(h(x), y)$ , and

$$\begin{aligned}\hat{\mathcal{R}}(h) &= \frac{1}{n} \sum_{i=1}^n \ell_h(z_i) \\ \mathcal{R}(h) &= \mathbb{E}_z[\ell_h(z)]\end{aligned}\tag{1}$$

be the empirical risk and population risk, respectively. Let  $\mathcal{H}$  be a hypothesis class. Consider the estimator:

$$\hat{h}_n = \operatorname{argmin}_{h \in \mathcal{H}} \hat{\mathcal{R}}(h).$$

This type of estimator ensures that  $\hat{\mathcal{R}}(\hat{h}_n)$ . But our question is: How small is the true error  $\mathcal{R}(\hat{h}_n)$ ?

For any  $h \in \mathcal{H}$ , consider the decomposition:

$$\mathcal{R}(h) = \underbrace{\hat{\mathcal{R}}(h)}_{\text{training error}} + \underbrace{\mathcal{R}(h) - \hat{\mathcal{R}}(h)}_{\text{gen-gap}},$$

where the generalization gap satisfies

$$\text{gen-gap}(h) := \mathcal{R}(h) - \hat{\mathcal{R}}(h) = \mathbb{E}_z[\ell_h(z)] - \frac{1}{n} \sum \ell_h(z_i).\tag{2}$$

One may expect that  $\text{gen-gap}(h) = O(1/\sqrt{n})$ . By concentration inequality, this is true for  $h$  that is independent of training data  $(z_1, \dots, z_n)$ . However, our task is bound of gen-gap of  $\hat{h}_n$ :

$$\text{gen-gap}(\hat{h}_n) = \mathbb{E}_z[\ell_{\hat{h}_n}(z)] - \frac{1}{n} \sum \ell_{\hat{h}_n}(z_i).$$

Note that  $\hat{h}_n$  depends on  $(z_1, \dots, z_n)$  and hence  $\{\ell_{\hat{h}_n}(z_i)\}$  are not i.i.d. . Consequently, gen-gap may not be in the order of  $O(1/\sqrt{n})$ . In fact that  $\text{gen-gap}(\hat{h}_n)$  can be arbitrarily large if  $h$  is complex.

## 2 Uniform bounds

To deal with the dependence issue, we can consider the uniform bound

$$|\mathcal{R}(\hat{h}_n) - \hat{\mathcal{R}}(\hat{h}_n)| \leq \sup_{h \in \mathcal{H}} |\mathcal{R}(h) - \hat{\mathcal{R}}(h)|. \quad (3)$$

Obviously, when the hypothesis space  $\mathcal{H}$  is sufficiently “small”, e.g., the extreme case:  $\mathcal{H} = \{h\}$ , it is expected that

$$\sup_{h \in \mathcal{H}} |\mathcal{R}(h) - \hat{\mathcal{R}}(h)| \sim \frac{1}{\sqrt{n}}.$$

Some natural questions go as follows.

- What kind of  $\mathcal{H}$  can guarantee the smallness of uniform bound?
- What is the rate? Do we still have  $O(1/\sqrt{n})$ ?

Let us first look at a simple example: finite hypothesis class.

**Lemma 2.1** (Finite class). *Let  $\mathcal{H}$  be a collection of finite hypotheses and denote by  $|\mathcal{H}|$  the number of hypotheses. Assume  $\sup_{y, y'} |\ell(y, y')| \leq 1$ . For any  $\delta \in (0, 1)$ , with probability  $1 - \delta$  over the random sampling of training set  $S$ , we have*

$$\sup_{h \in \mathcal{H}} |\mathcal{R}(h) - \hat{\mathcal{R}}(h)| \leq \sqrt{\frac{2 \ln(2|\mathcal{H}|/\delta)}{n}}.$$

*Proof.* WLOG, suppose  $\mathcal{H} = \{h_1, \dots, h_m\}$ . Let  $z = (x, y)$  and  $Q_h(z) = \ell(h(x), y)$ . Taking the union bound gives us

$$\mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum_{i=1}^n Q(h, z_i) - \mathbb{E}_z[Q(h, z)] \right| \geq t \right\} \leq \sum_{j=1}^m \mathbb{P} \left\{ \left| \frac{1}{n} \sum_{i=1}^n Q(h_j, z_i) - \mathbb{E}_z[Z(h_j, z)] \right| \geq t \right\} \quad (4)$$

$$\leq m 2e^{-\frac{2nt^2}{2}} = 2me^{-\frac{nt^2}{2}}, \quad (5)$$

where the last step follows from the Hoeffding’s inequality. Let the failure probability  $2me^{-\frac{nt^2}{2}} = \delta$ , which leads to  $t = \sqrt{\frac{2 \ln(2m/\delta)}{n}}$ . □

We see that the upper bound only depends on the cardinality of hypothesis class  $|\mathcal{H}|$  logarithmically. This implies that even when the hypothesis class has exponentially many functions, the generalization gap can be still well controlled.

**Definition 2.2** (Empirical process). Let  $\mathcal{F}$  be a class of real-valued functions  $f : \Omega \mapsto \mathbb{R}$  where  $(\Omega, \Sigma, \mu)$  is a probability space. Let  $X \sim \mu$  and  $X_1, \dots, X_n$  be independent copies of  $X$ . Then, the random process  $(X_f)_{f \in \mathcal{F}}$  defined by

$$X_f := \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X)$$

is called an *empirical process* indexed by  $\mathcal{F}$ .

In our case,  $f(Z) = \ell(h(X), Y)$ . Our task is to bound the supremum:

$$\sup_{f \in \mathcal{F}} |X_f|.$$

Note that the above quantity can be viewed as a “weak” distance between  $\mu$  and the empirical measure  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta(\cdot - x_i)$  with the test functions given by  $\mathcal{F}$ :

$$d_{\mathcal{F}}(\hat{\mu}_n, \mu) := \sup_{f \in \mathcal{F}} |\mathbb{E}_{\hat{\mu}_n} f - \mathbb{E}_{\mu} f|.$$

### 3 Covering number

For the finite hypothesis classes, we have shown that  $\log |\mathcal{F}|$ , i.e., the logarithm of cardinality, can be used as a good complexity measure. Then, a natural question is: can we do similar arguments for the case where  $|\mathcal{F}| = \infty$ ? One possible approach is *discretization*. This means that we choose a finite subset  $\mathcal{F}_{\varepsilon} \subset \mathcal{F}$  to “represent”  $\mathcal{F}$ .

**Definition 3.1** (Covering number). Consider a metric space  $(T, \rho)$ .

- We say  $T_{\varepsilon} \subset T$  is an  $\varepsilon$ -cover (also called  $\varepsilon$ -net) of  $T$ , if for any  $t \in T$ , there exists a  $t' \in T_{\varepsilon}$  such that  $\rho(t, t') \leq \varepsilon$ .
- The covering number  $\mathcal{N}(T, \rho, \varepsilon)$  is defined as the smallest cardinality of an  $\varepsilon$ -cover of  $T$  with respect to  $\rho$ .

**Definition 3.2** (Metric entropy). The *metric entropy* of  $T$  is defined by  $\log \mathcal{N}(T, \rho, \varepsilon)$ .

**Theorem 3.3.** Let  $\mathcal{F}$  be a function class with  $\sup_{f \in \mathcal{F}, x \in \mathcal{X}} |f(x)| \leq B$ . Let  $\|f - g\|_{\infty} = \sup_{x \in \mathcal{X}} |f(x) - g(x)|$ . Then, for any  $\delta \in (0, 1)$ , w.p. at least  $1 - \delta$  over the sampling of  $X_1, X_2, \dots, X_n$ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}[f(X)] \right| \leq 2\varepsilon + B \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon) + \log(2/\delta)}{n}}.$$

*Proof.* Let  $\mathcal{F}_{\varepsilon}$  be an  $\varepsilon$ -cover of  $\mathcal{F}$ . For any  $f \in \mathcal{F}$ , let  $f' \in \mathcal{F}_{\varepsilon}$  such that  $\|f - f'\|_{\infty} \leq \varepsilon$ . Then, we have

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}[f(X)] \right| &\leq \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \frac{1}{n} \sum_{i=1}^n f'(X_i) \right| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n f'(X_i) - \mathbb{E}[f'(X)] \right| + |\mathbb{E}[f'(X)] - \mathbb{E}[f(X)]|. \end{aligned}$$

Taking the supremum with respect to  $f \in \mathcal{F}$  gives

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}[f(X)] \right| &\leq 2\varepsilon + \sup_{f' \in \mathcal{F}_{\varepsilon}} \left| \frac{1}{n} \sum_{i=1}^n f'(X_i) - \mathbb{E}[f'(X)] \right| \\ &\leq 2\varepsilon + 2B \sqrt{\frac{\log(|\mathcal{F}_{\varepsilon}|/\delta)}{n}}, \end{aligned}$$

where the last step uses the uniform generalization bound of finite class. Notice that  $|\mathcal{F}_{\varepsilon}| \leq \mathcal{N}(\mathcal{F}, \|\cdot\|_{\infty}, \varepsilon)$ .  $\square$

**Example: Lipschitz models** Let  $f : \mathcal{X} \times \mathbb{R}^p \mapsto \mathbb{R}$  be our model. Here  $p$  denotes the number of parameters. Assume that  $f$  is  $L$ -Lipschitz in the sense that  $\sup_x |f(x; \theta_1) - f(x; \theta_2)| \leq L\rho(\theta_1, \theta_2)$ .

Let  $\mathcal{F} = \{f(\cdot; \theta) : \theta \in \Omega\}$  be the function class. Let  $\Omega_\varepsilon$  be an  $\varepsilon$ -cover of  $\Omega$  with respect to the  $\rho$  metric. Then,

$$\|f(\cdot; \theta_1) - f(\cdot; \theta_2)\|_\infty \leq L\rho(\theta_1, \theta_2)$$

implies that  $\mathcal{F}_\varepsilon = \{f(\cdot; \theta) : \theta \in \Omega_{\varepsilon/L}\}$  is an  $\varepsilon$ -cover of  $\mathcal{F}$ . Hence, we have

$$\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon) \leq \mathcal{N}(\Omega, \rho, \frac{\varepsilon}{L}). \quad (6)$$

**Linear class.** Consider the linear class:

$$\mathcal{H} = \{w^T x : \|w\|_2 \leq 1, \|x\|_2 \leq 1\}.$$

Then,

$$\sup_x |w^T x - v^T x| \leq \|w - v\| \sup_x \|x\| \leq \|w - v\|_2.$$

Let  $B_d(r) = \{x \in \mathbb{R}^d : \|x\| \leq r\}$  be the ball of radius  $r$ . Then, (6) gives

$$\mathcal{N}(\mathcal{H}, \|\cdot\|_\infty, \varepsilon) \leq \mathcal{N}(B_d(1), \|\cdot\|_2, \varepsilon).$$

The above examples demonstrate that one can reduce the estimation of covering number of a function class to the covering number of a subset in Euclidean space. The latter is often easier to estimate and given below is an example.

### 3.1 Volume argument for estimating covering number

To help the estimation of covering number, we introduce the packing number.

**Definition 3.4** (Packing number). Consider a metric space  $(T, \rho)$ .  $T_\varepsilon \subset T$  is said to be  $\varepsilon$ -separated if  $\rho(x, y) > \varepsilon$  for any  $x, y \in T_\varepsilon$  and  $x \neq y$ . The packing number is defined as

$$\mathcal{P}(\mathcal{F}, \rho, \varepsilon) = \sup_{T_\varepsilon \subset T \text{ is } \varepsilon\text{-separated}} |T_\varepsilon|$$

**Lemma 3.5.**  $\mathcal{N}(T, \rho, \varepsilon) \leq \mathcal{P}(T, \rho, \varepsilon)$ .

*Proof.* Let  $T_\varepsilon$  be the maximal  $\varepsilon$ -separated subset. Then, we claim that  $T_\varepsilon$  is also an  $\varepsilon$ -cover of  $T$ , i.e.,  $T \subset \cup_{x \in T_\varepsilon} B_x(\varepsilon)$ . If not, there exists a  $y \in T$  such that  $d(y, x) > \varepsilon$  for any  $x \in T_\varepsilon$ . Hence,  $T_\varepsilon \cup \{y\}$  is also  $\varepsilon$ -separated, which is contradictory with the assumption.  $\square$

**Lemma 3.6.**  $(1/\varepsilon)^d \leq \mathcal{N}(B^d(1), \|\cdot\|_2, \varepsilon) \leq (1 + 2/\varepsilon)^d$ .

The proof follows from a volume argument.

*Proof. Lower bound.* Let  $N_\varepsilon$  be an  $\varepsilon$ -cover of  $B^d(1)$ . Then,  $B^d(1) \subset \cup_{x \in N_\varepsilon} B_x^d(\varepsilon)$ . Therefore,

$$\text{Vol}(B^d(1)) \leq \sum_{x \in N_\varepsilon} \text{Vol}(B_x^d(\varepsilon)) = |N_\varepsilon| \text{Vol}(B_x^d(\varepsilon)).$$

Hence,

$$\mathcal{N}(B_1^d, \|\cdot\|_2, \varepsilon) = |N_\varepsilon| \geq \frac{\text{Vol}(B^d(1))}{\text{Vol}(B_x^d(\varepsilon))} = \left(\frac{1}{\varepsilon}\right)^d$$

**Upper bound.** Let  $P_\varepsilon \subset B^d(1)$  be  $\varepsilon$ -separated. Then, by definition of packing, we have

$$\bigcup_{x \in P_\varepsilon} B_x^d(\varepsilon/2) \subset B^d(1 + \varepsilon/2) \Rightarrow \sum_{x \in P_\varepsilon} \text{Vol}(B_x^d(\varepsilon/2)) \leq \text{Vol}(B^d(1 + \varepsilon/2)).$$

Let  $C_d r^d$  be the volume of a  $\ell_2$  ball of radius  $r$ . Then,

$$|P_\varepsilon| C_d (\varepsilon/2)^d \leq C_d (1 + \varepsilon/2)^d \Rightarrow |P_\varepsilon| \leq (1 + 2/\varepsilon)^d.$$

Then, the upper bound follows from Lemma 3.5.  $\square$

*Remark 3.7.* It should be remarked that the above volume argument can be applied to estimate the covering number of other classes and different metrics.

## 4 Rademacher complexity

The following inequality

**Lemma 4.1** (Symmetrization of empirical processes).

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right] \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n \xi_i f(X_i) \right],$$

where  $\xi_1, \dots, \xi_n$  are i.i.d. Rademacher random variable:  $\mathbb{P}(\xi = 1) = \mathbb{P}(\xi = -1) = \frac{1}{2}$

*Proof.* Let  $X'_i$  be an independent copy of  $X_i$ . To simplify the notation, we use  $\mathbb{E}_{X_i}$  and  $\mathbb{E}_{X'_i}$  to denote the expectation with respect to  $\{X_i\}_{i=1}^n$  and  $\{X'_i\}_{i=1}^n$ , respectively. Then,

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right] = \mathbb{E}_{X_i} \sup_{f \in \mathcal{F}} \mathbb{E}_{X'_i} \left[ \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right] \quad (7)$$

$$\leq \mathbb{E}_{X_i, X'_i} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n (f(X_i) - f(X'_i)) \right] \quad (8)$$

Due to that  $f(X_i) - f(X'_i)$  is symmetric, for any  $\{\xi_i\} \in \{\pm 1\}^n$ , we have

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

$$= 2 \mathbb{E}_{X_i, \xi} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \xi_i f(X_i)$$

□

**Definition 4.2** (Rademacher complexity). The empirical Rademacher complexity of a function class  $\mathcal{F}$  on finite samples is defined as

$$\widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E}_{\xi} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \xi_i f(X_i) \right].$$

The population Rademacher complexity is given by

$$\text{Rad}_n(\mathcal{F}) = \mathbb{E}_S [\widehat{\text{Rad}}_n(\mathcal{F})].$$

The symmetrization lemma 4.1 implies that

$$\mathbb{E} \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right] \leq 2 \text{Rad}_n(\mathcal{F}). \quad (9)$$

**Theorem 4.3.** Assume that  $0 \leq f \leq B$  for all  $f \in \mathcal{F}$ . For any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over the choice of the training set  $S = \{X_1, \dots, X_n\}$ , we have

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right| \leq 2 \text{Rad}_n(\mathcal{F}) + B \sqrt{\frac{2 \log(2/\delta)}{n}}, \quad (10)$$

and the sample-dependent version:

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right| \leq 2 \widehat{\text{Rad}}_n(\mathcal{F}) + 4B \sqrt{\frac{2 \log(4/\delta)}{n}}. \quad (11)$$

*Proof.* Let

$$G(x_1, \dots, x_n) = \sup_{f \in \mathcal{F}} \left[ \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E} f(X) \right].$$

Let  $\tilde{X}_1, \dots, \tilde{X}_n$  be a copy of  $X_1, \dots, X_n$  with only  $\tilde{X}_i \neq X_i$  for  $i \in [n]$ . Then, we have

$$\begin{aligned} & G(X_1, \dots, X_n) - G(\tilde{X}_1, \dots, \tilde{X}_n) \\ &= \sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right) - \sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n f(\tilde{X}_i) - \mathbb{E} f(X) \right) \\ &\leq \sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) - \left( \frac{1}{n} \sum_{i=1}^n f(\tilde{X}_i) - \mathbb{E} f(X) \right) \right) \\ &\leq \sup_{f \in \mathcal{F}} \frac{1}{n} \left( f(X_i) - f(\tilde{X}_i) \right) \leq \frac{2B}{n}. \end{aligned}$$

Similarly, we have

$$G(\tilde{X}_1, \dots, \tilde{X}_n) - G(X_1, \dots, X_n) \geq -\frac{2B}{n}.$$

Therefore, the variation satisfies

$$\|D_i G\|_\infty := \sup_{X, \tilde{X}} |G(X_1, \dots, X_n) - G(\tilde{X}_1, \dots, \tilde{X}_n)| \leq 2B/n,$$

where  $X = (\tilde{X}_1, \dots, \tilde{X}_n)$  and  $\tilde{X} = (X_1, \dots, X_n)$  are different for only the  $i$ -th component.

Therefore,  $\sigma^2 = \frac{1}{4} \sum_{i=1}^n \|D_i G\|_\infty^2 \leq \frac{B^2}{n}$ . By McDiarmid's inequality,

$$\mathbb{P}\{|G(X_1, \dots, X_n) - \mathbb{E} G| \geq t\} \leq 2e^{-\frac{nt^2}{2B^2}}.$$

Let the failure probability  $2e^{-\frac{nt^2}{2B^2}} = \delta$ , which leads to  $t = \sqrt{\frac{2B^2 \log(2/\delta)}{n}}$ . Restating the above inequality gives the bound (10).

Analogously, we can applying McDiarmid's inequality to the Rademacher complexity  $Q(x_1, \dots, x_n) = \mathbb{E}_\xi \sup_{f \in \mathcal{F}} [\frac{1}{n} \sum_{i=1}^n \xi_i f(x_i)]$ , which leads to the sample-dependent bound (11).  $\square$

### Examples.

- Let  $\mathcal{F} = \{f\}$ . Then,

$$\widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E}_\xi [\frac{1}{n} \sum_{i=1}^n \xi_i f(x_i)] = 0.$$

- Two functions. Let  $\mathcal{F} = \{f_{-1}, f_1\}$  where  $f_{-1} \equiv -1$  and  $f_1 \equiv 1$ .

$$\sqrt{n} \widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E}_\xi \sup_{f \in \{-1, +1\}} f \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i = \mathbb{E}_\xi \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \right| \rightarrow \mathbb{E}_{Z \sim \mathcal{N}(0,1)} |Z| = \sqrt{\frac{2}{\pi}}.$$

Hence, when  $n$  is sufficiently large,

$$\text{Rad}_n(\mathcal{F}) \sim \sqrt{\frac{2}{n\pi}}.$$

**Remark:** This implies that it is impossible to obtain a rate faster than  $O(1/\sqrt{n})$  using Rademacher complexity since it saturates even for learning/distinguishing two constant functions. This is a bad news!

**Lemma 4.4** (Massart's lemma). *Assume that  $\sup_{x \in \mathcal{X}, f \in \mathcal{F}} |f(x)| \leq B$  and  $\mathcal{F}$  is finite. Then,*

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq B \sqrt{\frac{2 \log |\mathcal{F}|}{n}}.$$

*Proof.* Let  $Z_f = \sum_{i=1}^n \xi_i f(x_i)$ . Then,

$$\log \mathbb{E}[e^{\lambda Z_f}] = \log \left( \prod_{i=1}^n \mathbb{E}[e^{\lambda \xi_i f(x_i)}] \right) \leq \sum_{i=1}^n \log \mathbb{E} e^{\lambda \xi_i f(x_i)} \stackrel{(i)}{\leq} \sum_{i=1}^n \lambda^2 \frac{(B - (-B))^2}{8} = \frac{nB^2}{2} \lambda^2,$$

where (i) follows from the Hoeffding's lemma, which provides an upper bound of the log-moment generating functions of a bounded random variable. Hence,  $Z_f$  is sub-Gaussian with the variance proxy  $\sigma^2 = nB^2$ . Using the maximal inequality, we have

$$\widehat{\text{Rad}}_n(\mathcal{F}) = \frac{1}{n} \mathbb{E}_\xi [\sup_{f \in \mathcal{F}} Z_f] \leq \frac{1}{n} \cdot \sqrt{n} B \sqrt{2 \log |\mathcal{F}|} = B \sqrt{\frac{2 \log |\mathcal{F}|}{n}}. \quad (12)$$

$\square$

Applying Massart's lemma to bound the generalization gap recovers Lemma 2.1.

**Linear functions.** Let  $\mathcal{F} = \{w^T x : \|w\|_p \leq 1\}$ . Let  $q$  be the conjugate of  $p$ , i.e.,  $1/q + 1/p = 1$ . Then,

$$\widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E}_\xi \sup_{\|w\|_p \leq 1} \frac{1}{n} \sum_{i=1}^n \xi_i w^T X_i = \mathbb{E}_\xi \sup_{\|w\|_p \leq 1} w^T \left( \frac{1}{n} \sum_{i=1}^n \xi_i X_i \right) = \mathbb{E}_\xi \left\| \frac{1}{n} \sum_{i=1}^n \xi_i X_i \right\|_q. \quad (13)$$

**Lemma 4.5.** Assume that  $\|x_i\|_q \leq 1$  for all  $x_i \in S$ . Then,

- If  $p = 2$ , then

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq \sqrt{\frac{1}{n}}.$$

- If  $p = 1$ , then,

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq \sqrt{\frac{2 \log(2d)}{n}}.$$

*Proof.* For the case where  $p = 2$ ,

$$\begin{aligned} \widehat{\text{Rad}}_n(\mathcal{F}) &\leq \mathbb{E}_\xi \left\| \frac{1}{n} \sum_{i=1}^n \xi_i x_i \right\|_2 \leq \sqrt{\mathbb{E}_\xi \left\| \frac{1}{n} \sum_{i=1}^n \xi_i x_i \right\|_2^2} \\ &= \sqrt{\frac{1}{n^2} \sum_{i,j=1}^n x_i x_j \mathbb{E}[\xi_i \xi_j]} = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2} \leq \sqrt{\frac{1}{n}}. \end{aligned}$$

The case of  $p = 1$  leaves to homework. □

We have shown the Rademacher complexity of linear functions. To obtain the estimates of more general classes, we need follow results.

**Lemma 4.6** (Rademacher calculus). *The Rademacher complexity has the following properties.*

- $\text{Rad}_n(\lambda \mathcal{F}) = |\lambda| \text{Rad}_n(\mathcal{F})$ .
- $\text{Rad}_n(\mathcal{F} + f_0) = \text{Rad}_n(\mathcal{F})$ .
- Let  $\text{Conv}(\mathcal{F})$  denote the convex hull of  $\mathcal{F}$  defined by

$$\text{Conv}(\mathcal{F}) = \left\{ \sum_{j=1}^m a_j f_j : a_j \geq 0, \sum_{j=1}^m a_j = 1, f_1, \dots, f_m \in \mathcal{F}, m \in \mathbb{N}_+ \right\}.$$

Then, we have  $\text{Rad}_n(\text{Conv}(\mathcal{F})) = \text{Rad}_n(\mathcal{F})$ .

*Proof.* Here, we only prove the third result. By definition,

$$n \widehat{\text{Rad}}_n(\text{Conv}(\mathcal{F})) = \mathbb{E} \sup_{f_j \in \mathcal{F}, \|\alpha\|_1 = 1} \sum_{i=1}^n \xi_i \sum_{j=1}^m a_j f_j(X_i)$$



$$\begin{aligned}
&= \mathbb{E} \sup_{f_j \in \mathcal{F}, \|\alpha\|_1=1} \sum_{j=1}^m a_j \sum_{i=1}^n \xi_i f_j(X_i) \\
&= \mathbb{E} \sup_{f_j \in \mathcal{F}} \max_j \sum_{i=1}^n \xi_i f_j(X_i) \\
&= \mathbb{E} \sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i f(X_i) = n \widehat{\text{Rad}}_n(\mathcal{F})
\end{aligned}$$

□

The third property suggests that convex combinations does not change the Rademacher complexity.

**Lemma 4.7** (Ledoux & Talagrand 2011, Contraction lemma). *Let  $\varphi_i : \mathbb{R} \mapsto \mathbb{R}$  with  $i = 1, \dots, n$  be  $\beta$ -Lispchitz continuous. Then,*

$$\frac{1}{n} \mathbb{E}_\xi \sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i \varphi_i \circ f(x_i) \leq \beta \widehat{\text{Rad}}_n(\mathcal{F}).$$

*Proof.* WLOG, assume  $\beta = 1$ . Let  $\hat{\xi} = (\xi_1, \dots, \xi_n)$  and  $Z_k(f) = \sum_{i=1}^k \xi_i \varphi_i \circ f(x_i)$ . Then,

$$\begin{aligned}
\mathbb{E}_{\xi_n} \sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i \varphi_i \circ f(x_i) &= \frac{1}{2} \left[ \sup_{f \in \mathcal{F}} (Z_{n-1}(f) + \varphi_n \circ f(x_n)) + \sup_{f \in \mathcal{F}} (Z_{n-1}(f) - \varphi_n \circ f(x_n)) \right] \\
&= \frac{1}{2} \sup_{f, f' \in \mathcal{F}} \left( Z_{n-1}(f) + Z_{n-1}(f') + \varphi_n \circ f(x_n) - \varphi_n \circ f'(x_n) \right) \\
&\leq \frac{1}{2} \sup_{f, f' \in \mathcal{F}} \left( Z_{n-1}(f) + Z_{n-1}(f') + |f(x_n) - f'(x_n)| \right) \\
&= \frac{1}{2} \sup_{f, f' \in \mathcal{F}} \left( Z_{n-1}(f) + Z_{n-1}(f') + (f(x_n) - f'(x_n)) \right) \quad (\text{Use the symmetry}) \\
&= \frac{1}{2} \left[ \sup_{f \in \mathcal{F}} (Z_{n-1}(f) + f(x_n)) + \sup_{f \in \mathcal{F}} (Z_{n-1}(f) - f(x_n)) \right] \\
&= \mathbb{E}_{\xi_n} \sup_{f \in \mathcal{F}} (Z_{n-1}(f) + \xi_n f(x_n)).
\end{aligned}$$

Hence, by induction, we have

$$\begin{aligned}
\mathbb{E}_{\hat{\xi}} [\sup_{f \in \mathcal{F}} Z_n(f)] &\leq \mathbb{E}_{\hat{\xi}} \sup_{f \in \mathcal{F}} (Z_{n-1}(f) + \xi_n f(x_n)) \\
&\leq \mathbb{E}_{\hat{\xi}} \sup_{f \in \mathcal{F}} (Z_{n-2}(f) + \xi_{n-1} f(x_{n-1}) + \xi_n f(x_n)) \\
&\leq \mathbb{E}_{\hat{\xi}} \sup_{f \in \mathcal{F}} (\xi_1 f(x_1) + \dots + \xi_n f(x_n)) \\
&= n \widehat{\text{Rad}}_n(\mathcal{F}).
\end{aligned} \tag{14}$$

□

**Corollary 4.8.** *Given a function class  $\mathcal{F}$  and  $\varphi : \mathbb{R} \mapsto \mathbb{R}$ , let  $\varphi \circ \mathcal{F} = \{\varphi \circ f : f \in \mathcal{F}\}$ . Then,*

$$\text{Rad}_n(\varphi \circ \mathcal{F}) \leq \text{Lip}(\varphi) \text{Rad}_n(\mathcal{F}).$$

**Rademacher complexity of neural networks.** In the following, we provide an example showing the power of combining the contraction lemma with Rademacher calculus. They together can bound the Rademacher complexity of many complex models.

Consider two-layer neural networks. Suppose the activation function  $\sigma : \mathbb{R} \mapsto \mathbb{R}$  is  $\sigma_{\text{Lip}}$ -Lipschitz continuous. Let

$$\mathcal{F}_m = \left\{ f_m(x; \theta) = \sum_{j=1}^m a_j \sigma(w_j^T x) : \sum_j |a_j| \leq A, \|w_j\|_2 \leq B \right\}.$$

be the collection of two-layer neural networks  $f_m(\cdot; \theta)$ .

**Lemma 4.9.** Suppose  $\|x_i\|_2 \leq 1$  for  $i = 1, \dots, n$ . Then, we have

$$\widehat{\text{Rad}}_n(\mathcal{F}_m) \leq \frac{2\sigma_{\text{Lip}}AB}{\sqrt{n}}.$$

The above lemma implies that Rademacher complexity only depends on the parameter norm, independent of the network width. This implies that the capacity of over-parameterized networks can be well-controlled by enforcing a constraint on an appropriate parameter norm. It is worth noting that for different networks, we may need to identify the appropriate norm of parameters.

*Proof.*

$$\begin{aligned} \widehat{\text{Rad}}_n(\mathcal{F}_m) &= \frac{1}{n} \mathbb{E}_{\xi} \sup_{f \in \mathcal{F}_m} \sum_{i=1}^n f(x_i) \xi_i \\ &= \frac{1}{n} \mathbb{E}_{\xi} \sup_{\theta \in \Theta} \sum_{i=1}^n \xi_i \sum_{j=1}^m a_j \sigma(w_j^T x_i) \\ &= \frac{1}{n} \mathbb{E}_{\xi} \sup_{\theta \in \Theta} \sum_{j=1}^m a_j \sum_{i=1}^n \xi_i a_j \sigma(w_j^T x_i) \\ &\leq \frac{1}{n} \mathbb{E}_{\xi} \sup_{\theta \in \Theta} \sum_{j=1}^m |a_j| \left| \sup_{\|w\| \leq B} \sum_{i=1}^n \xi_i \sigma(w^T x_i) \right| \\ &\stackrel{(i)}{\leq} A \frac{1}{n} \mathbb{E}_{\xi} \sup_{\|w\| \leq B} \left| \sum_{i=1}^n \xi_i \sigma(w^T x_i) \right| \\ &= A \frac{1}{n} \mathbb{E}_{\xi} \left( \sup_{\|w\| \leq B} \sum_{i=1}^n \sigma(\xi_i w^T x_i) \right) + A \frac{1}{n} \mathbb{E}_{\xi} \left( - \sup_{\|w\| \leq B} \sum_{i=1}^n \sigma(\xi_i w^T x_i) \right) \\ &\stackrel{(ii)}{\leq} 2A \frac{1}{n} \mathbb{E}_{\xi} \left( \sup_{\|w\| \leq B} \sum_{i=1}^n \sigma(\xi_i w^T x_i) \right) \\ &\stackrel{iii}{\leq} 2A \sigma_{\text{Lip}} \frac{1}{n} \mathbb{E}_{\xi} \left( \sup_{\|w\| \leq B} \sum_{i=1}^n \xi_i w^T x_i \right) \\ &\stackrel{(iii)}{\leq} \frac{\sigma_{\text{Lip}}AB}{\sqrt{n}}, \end{aligned}$$

where (i) is due to  $\sum_{j=1}^m |a_j| \leq A$ ; (ii) use the symmetry of  $\xi_i$ ; (iii) follows from the contraction property (Lemma 4.7); (iii) follows from Lemma 4.5.  $\square$

## 5 Bounding Rademacher complexity using covering number

Consider the function space  $(\mathcal{F}, L^2(\mathbb{P}_n))$ , where  $\mathcal{F}$  is the hypothesis class and  $L^2(\mathbb{P}_n)$  is defined by

$$\|f - f'\|_{L^2(\mathbb{P}_n)} = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - f'(x_i))^2},$$

where  $x_1, \dots, x_n$  denote the finite training samples. Since only the  $n$  samples are available, we can really think of these functions as a  $n$ -dimensional vector:

$$\hat{f} = (f(x_1), f(x_2), \dots, f(x_n))^T \in \mathbb{R}^n,$$

Obviously, we cannot distinguish functions using information beyond these  $n$ -dimensional vectors.

**Example 1.** Let  $\mathcal{F} = \{f : \mathbb{R} \mapsto [0, 1] : f \text{ is non-decreasing}\}$ . Then,  $\mathcal{N}(\mathcal{F}, L_2(\mathbb{P}_n), \varepsilon) = n^{1/\varepsilon}$ .

*Proof.* WLOG, assume  $-\infty = x_0 < x_1 \leq x_2 \leq \dots \leq x_n \leq x_{n+1} = 1$ . For any  $y = (y_1, y_2, \dots, y_n) \in \mathbb{R}^n$ , define a piecewise constant function

$$f_y(x) = y_i \quad \text{for } x \in [x_i, x_{i+1}), \quad i = 1, 2, \dots, n.$$

For any  $\varepsilon \in (0, 1)$ , let  $Y_\varepsilon = (0, \varepsilon, 2\varepsilon, 3\varepsilon, \dots, 1 - \varepsilon)$ . Then,  $|Y_\varepsilon| \leq 1/\varepsilon$ . Define the following non-decreasing set:

$$S_\varepsilon := \{y \in \mathbb{R}^n : y_i \in Y_\varepsilon \text{ and } y_1 \leq y_2 \leq \dots \leq y_n\}.$$

Let  $\mathcal{F}_\varepsilon = \{f_y : y \in S_\varepsilon\}$ . Obviously,  $\mathcal{F}_\varepsilon \subset \mathcal{F}$ . Moreover, for any  $f \in \mathcal{F}$ , there exists  $y \in S_\varepsilon$  such that

$$\|f - f_y\|_{L_2(\mathbb{P}_n)}^2 = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \leq \varepsilon^2.$$

Hence,  $\mathcal{F}_\varepsilon$  is an  $\varepsilon$ -cover of  $\mathcal{F}$  and  $|\mathcal{F}_\varepsilon| = |S_\varepsilon|$ . What remains is to count the cardinality of  $|S_\varepsilon|$ . Let  $y_0 = 0, y_{n+1} = 1$  and  $\Delta_i = (y_i - y_{i-1})/\varepsilon$ . Then,  $\{\Delta_i\}_{i=1}^{n+1}$  must be non-negative integers and satisfy

$$\Delta_1 + \Delta_2 + \dots + \Delta_{n+1} = \frac{1}{\varepsilon}.$$

Hence,  $|S_\varepsilon|$  is equal to the number of solutions of the above equation:

$$|S_\varepsilon| = \binom{n + \frac{1}{\varepsilon}}{n} = \frac{(n + \frac{1}{\varepsilon})(n + \frac{1}{\varepsilon} - 1) \cdots (n + 1)}{(\frac{1}{\varepsilon})(\frac{1}{\varepsilon} - 1) \cdots 1} \leq n^{\frac{1}{\varepsilon}}.$$

$\square$

In the following, we show that the Rademacher complexity can be bounded using the metric entropy. To simplify notation, we use  $\|\cdot\|$  and  $\langle \cdot, \cdot \rangle$  to denote  $L^2(\mathbb{P}_n)$  norm and the induced inner product:  $\langle f, g \rangle = \frac{1}{n} \sum_{i=1}^n f(x_i)g(x_i)$ . Then,

$$\widehat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, f \rangle.$$

**Proposition 5.1** (One-resolution discretization). *Suppose  $\sup_{x \in \mathcal{X}, f \in \mathcal{F}} |f(x)| \leq B$ . Then,*

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq \inf_{\varepsilon} \left( \varepsilon + B \sqrt{\frac{2 \log \mathcal{N}(\mathcal{F}, L_2(\mathbb{P}_n), \varepsilon)}{n}} \right).$$

The above bound is similar to Theorem 3.3. The difference is that the above bound is determined by the  $L^2(\mathbb{P}_n)$  covering number, while Theorem 3.3 relies on the  $L^\infty$  covering number. Technically speaking, this improvement is obtained by removing the  $\mathbb{E} f(X)$  term with symmetrization.

*Proof.* Let  $\mathcal{F}_\varepsilon$  be an  $\varepsilon$ -cover of  $\mathcal{F}$  with respect to the metric  $L^2(\mathbb{P}_n)$ . For any  $f \in \mathcal{F}$ , let  $\pi(f) \in \mathcal{F}_\varepsilon$  such that  $\|f - \pi(f)\| \leq \varepsilon$ . Then,

$$\begin{aligned} \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, f \rangle &= \mathbb{E} \sup_{f \in \mathcal{F}} \left[ \langle \xi, f - \pi(f) \rangle + \langle \xi, \pi(f) \rangle \right] \\ &\leq \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, f - \pi(f) \rangle + \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, \pi(f) \rangle \\ &\leq \mathbb{E} \|\xi\| \|f - \pi(f)\| + \mathbb{E} \sup_{f \in \mathcal{F}_\varepsilon} \langle \xi, f \rangle \\ &\leq \varepsilon \sqrt{\frac{\mathbb{E} \|\xi\|_2^2}{n}} + \widehat{\text{Rad}}_n(\mathcal{F}_\varepsilon) \quad (\text{Jessen's inequality}) \\ &\leq \varepsilon + B \sqrt{\frac{2 \log |\mathcal{F}_\varepsilon|}{n}}, \quad (\text{Massart's lemma}). \end{aligned}$$

Using the definition of covering number and optimizing over  $\varepsilon$ , we complete the proof.  $\square$

For the non-decreasing functions considered previously, we have

$$\text{Rad}_n(\mathcal{F}) \leq \inf \left( \varepsilon + \sqrt{\frac{2 \log n}{\varepsilon n}} \right) = C \left( \frac{\log n}{n} \right)^{1/3}. \quad (15)$$

This rate is slower than the expected  $O(1/\sqrt{n})$ . Is it because non-decreasing functions are complex? No! It is actually just an artifact caused by the proof technique.

In many cases, the one-resolution discretization may give us sub-optimal bounds of generalization gap. To fix this problem, we need a sophisticated analysis of all the resolutions. This is typically done by using a *chaining* approach introduced by Dudley.

**Theorem 5.2** (Dudley's integral inequality). *Let  $D = \sup_{f, f' \in \mathcal{F}} \|f - f'\|_{L^2(\mathbb{P}_n)}$  be the diameter of  $\mathcal{F}$ . Then,*

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq 12 \inf_{\alpha < D} \left( \alpha + \int_{\alpha}^D \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, L^2(\mathbb{P}_n), \varepsilon)}{n}} d\varepsilon \right).$$

Then, for the for non-decreasing functions, we have

$$\text{Rad}_n(\mathcal{F}) \lesssim \int_0^2 \sqrt{\frac{\log n}{n\varepsilon}} d\varepsilon \lesssim \sqrt{\frac{\log n}{n}}.$$

Figure 1 visualizes the difference between the upper bound given in Proposition 5.1 and the one in Theorem 5.2. Clearly, the latter is smaller.

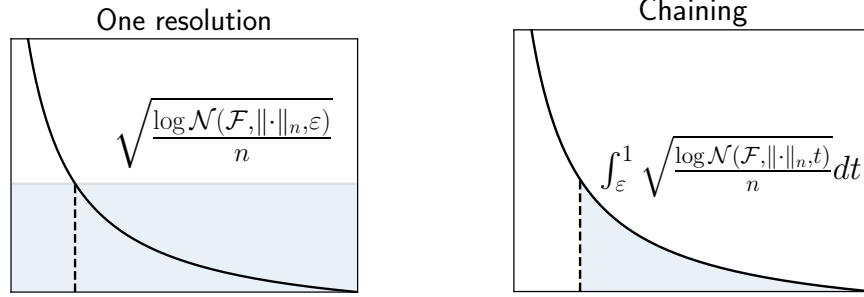


Figure 1: (Left) The result of one-resolution analysis; (Right) The result of chaining with all resolutions. In this case, the diameter  $D = 1$ . The comparison of two figures provides a visual illustration of how the chaining bound is tighter than the one-resolution bound.

*Proof.* Let  $\varepsilon_j = 2^{-j}D$  be the dyadic scale and  $\mathcal{F}_j$  be an  $\varepsilon_j$ -cover of  $\mathcal{F}$ . Given  $f \in \mathcal{F}$ , let  $f_j \in \mathcal{F}_j$  such that  $\|f_j - f\| \leq \varepsilon_j$ . Consider the decomposition

$$f = f - f_m + \sum_{j=1}^m (f_j - f_{j-1}), \quad (16)$$

where  $f_0 = 0$ . Notice that

- $\|f - f_m\| \leq \varepsilon_m$ .
- $\|f_j - f_{j-1}\| \leq \|f_j - f\| + \|f - f_{j-1}\| \leq \varepsilon_j + \varepsilon_{j-1} \leq 3\varepsilon_j$ .

Then,

$$\begin{aligned} \widehat{\text{Rad}}_n(\mathcal{F}) &= \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, f \rangle \\ &= \mathbb{E} \sup_{f \in \mathcal{F}} \left( \langle \xi, f - f_m \rangle + \sum_{j=1}^m \langle \xi, f_j - f_{j-1} \rangle \right) \\ &\leq \varepsilon_m + \mathbb{E} \sup_{f \in \mathcal{F}} \sum_{j=1}^m \langle \xi, f_j - f_{j-1} \rangle \\ &\leq \varepsilon_m + \sum_{j=1}^m \mathbb{E} \sup_{f \in \mathcal{F}} \langle \xi, f_j - f_{j-1} \rangle \end{aligned}$$

$$\begin{aligned}
&= \varepsilon_m + \sum_{j=1}^m \mathbb{E} \sup_{f_j \in \mathcal{F}_j, f_{j-1} \in \mathcal{F}_{j-1}} \langle \xi, f_j - f_{j-1} \rangle \\
&= \varepsilon_m + \sum_{j=1}^m \widehat{\text{Rad}}_n(\mathcal{F}_j \cup \mathcal{F}_{j-1}).
\end{aligned}$$

Using the Massart lemma and the fact that  $\sup_{f \in \mathcal{F}_j, f' \in \mathcal{F}_{j-1}} \|f_j - f_{j-1}\| \leq 3\varepsilon_j$ ,

$$\begin{aligned}
\widehat{\text{Rad}}_n(\mathcal{F}) &\leq \varepsilon_m + \sum_{j=1}^m 3\varepsilon_j \sqrt{\frac{2 \log(|\mathcal{F}_j| |\mathcal{F}_{j-1}|)}{n}} \\
&\leq \varepsilon_m + \sum_{j=1}^m 6\varepsilon_j \sqrt{\frac{\log |\mathcal{F}_j|}{n}} \\
&= \varepsilon_m + \sum_{j=1}^m 12(\varepsilon_j - \varepsilon_{j+1}) \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, L^2(\mathbb{P}_n), \varepsilon_j)}{n}}.
\end{aligned}$$

Taking  $m \rightarrow \infty$ , we obtain

$$\widehat{\text{Rad}}_n(\mathcal{F}) \leq 12 \int_0^D \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, L^2(\mathbb{P}_n), t)}{n}} dt.$$

□

Similarly, we can obtain that

$$\widehat{\text{Rad}}_n(\mathcal{F}) \lesssim \inf_{\alpha > 0} \left( \alpha + \int_{\alpha}^D \sqrt{\frac{\log \mathcal{N}(\mathcal{F}, L^2(\mathbb{P}_n), t)}{n}} dt \right).$$

The key ingredient of proceeding analysis is the multi-resolution decomposition (16). The technical reason why chaining provides a better estimate is as follows. In the one-resolution discretization, we apply Massart's lemma to functions whose range in  $[-1, 1]$ , whereas in chaining, we apply Massart's lemma to functions whose range has size  $O(\varepsilon_j)$ .

*Remark 5.3.* Metric entropy is actually a more intuitive complexity measure than Rademacher complexity. The essence is discretization and applying Massart's lemma. Moreover, metric entropy is sometimes more convenient to estimate.

## References

[Shalev-Shwartz and Ben-David, 2014] Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.