

$$\sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}_{z \sim P}[f(z)] \right]$$

이전 section에서 한 것: bound 기댓값 of \sup_{\cdot} .

=> training example z_1, \dots, z_n 에 대한 기댓값임.

=> 그러나 주로 only ONE TRAINING SET이어서, 그렇게까지 많은 training set이 없음

=> 즉, 이 bound는 우리의 ONE TRAINING SET이 bound됨을 보장해주지 못함

=> 따라서 이 값 자체 (즉 \sup_{\cdot})를 높은 확률로 bound해주는 걸 할거임!

4.5 Empirical Rademacher complexity

In the previous section, we bounded the expectation of $\sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}_{z \sim P}[f(z)] \right]$. This expectation is taken over the training examples z_1, \dots, z_n . In many instances we only have one training set, and do not have access to many training sets. Thus, the bound on the expectation does not give a guarantee for the one training set that we have. In this section, we seek to bound the quantity itself with high probability.

Definition 4.17 (Empirical Rademacher complexity). Given a dataset $S = \{z_1, \dots, z_n\}$, the empirical Rademacher complexity is defined as

$$R_S(\mathcal{F}) \triangleq \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right]. \quad \text{Function class} \quad \text{dataset} \quad \text{=< 애의 expectation이 Rademacher complexity임.} \quad (4.73)$$

$R_S(\mathcal{F})$ is a function of both the function class \mathcal{F} and the dataset S .

As the name suggests, the expectation of the empirical Rademacher complexity is the Rademacher complexity:

$$R_n(\mathcal{F}) = \mathbb{E}_{\substack{z_1, \dots, z_n \sim P \\ S = \{z_1, \dots, z_n\}}} [R_S(\mathcal{F})] = \mathbb{E}_{\substack{\sigma_1, \dots, \sigma_n \sim P \\ \sigma = \{\sigma_1, \dots, \sigma_n\}}} \left[\mathbb{E}_{\substack{\sigma_1, \dots, \sigma_n \\ \sigma = \{\sigma_1, \dots, \sigma_n\}}} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right] \right]. \quad (4.74)$$

Here is the theorem involving empirical Rademacher complexity:

Theorem 4.18. Suppose for all $f \in \mathcal{F}$, $0 \leq f(z) \leq 1$. Then, with probability at least $1 - \delta$,

원래는 이 LHS의 기댓값을 bound했었는데, * 예제 찾기: $\mathbb{E} \left[\sup \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}[f(z)] \right] \right] \leq 2R_n(\mathcal{F})$ (4.46)
(on the training examples z_1 to z_n) 이제는 LHS 자체를 bound하는거

$$\sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}[f(z)] \right] \leq 2R_S(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}. \quad \text{population ave.} \quad \text{empirical ave.} \quad (4.75)$$

Proof. For conciseness, define "위에서 찾은 "the quantity itself"

define the LHS as... \downarrow

$$g(z_1, \dots, z_n) \triangleq \sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}[f(z)] \right]. \quad \text{empirical ave.} \quad \text{pop. ave.} \quad (4.76)$$

We prove the theorem in 4 steps. \uparrow McDiarmid: Section 3.5, Theorem 3.15, page 20
함수 f '가 bounded difference condition (변수 하나가 바뀌어도 함수 값이 많이 움직이지 않음)을 만족하면,
any indep RVs X_1, \dots, X_n 에 대해, $\Pr[f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)] >= t] \leq \exp(-2t^2/\sigma^2)$

Step 1: We bound g using McDiarmid's Inequality. To use McDiarmid's Inequality, we check that the bounded difference condition holds:

(4.77)~(4.80): 함수 g 의 bounded difference condition 확인
 $g(z_1, \dots, z_n) - g(z_1, \dots, z'_i, \dots, z_n) \leq \sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{j=1}^n f(z_j) \right] - \sup_{f \in \mathcal{F}} \left[\left(\frac{1}{n} \sum_{j=1, j \neq i}^n f(z_j) \right) + \frac{f(z'_i)}{n} \right]$ (4.77)

\downarrow 대입. \downarrow $(g \text{ 찾았다})$ \downarrow $\leq \sup_{f \in \mathcal{F}} \left[\frac{1}{n} \sum_{j=1}^n f(z_j) - \left(\frac{1}{n} \sum_{j=1, j \neq i}^n f(z_j) \right) - \frac{f(z'_i)}{n} \right]$ (4.78)

$\sup_{f \in \mathcal{F}} A(f) - \sup_{f \in \mathcal{F}} B(f) \leq \sup_{f \in \mathcal{F}} [A(f) - B(f)] = \sup_{f \in \mathcal{F}} \left[\frac{1}{n} (f(z_i) - f(z'_i)) \right] \quad \frac{1}{n} f(z_i) \quad (4.79)$

$$\leq \frac{1}{n}. \quad 0 \leq f(z) \leq 1 \quad (4.80)$$

(4.78) holds because in general, $\sup_f A(f) - \sup_f B(f) \leq \sup_f [A(f) - B(f)]$, and (4.80) holds since f is bounded by $[0, 1]$. We can thus apply McDiarmid's Inequality with parameters $c_1 = \dots = c_n = 1/n$: (on g)

$$\Pr \left[g(z_1, \dots, z_n) \geq \mathbb{E}_{z_1, \dots, z_n \sim P}[g] + \epsilon \right] \leq \exp \left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2} \right) = \exp(-2n\epsilon^2). \quad (4.81)$$

$c_1 = \dots = c_n = \frac{1}{n}$

$\downarrow \sum_{i=1}^n c_i^2 = \sum_{i=1}^n \frac{1}{n^2} = \frac{1}{n}$

$$(4.81): \Pr_{z_1, \dots, z_n \sim P} [g(z_1, \dots, z_n) \geq \mathbb{E}_{z_1, \dots, z_n \sim P}[g] + \epsilon] \leq e^{-2n\epsilon^2}$$

Step 2: We apply Theorem 4.13 to get
($\Pr_{z_1, \dots, z_n \sim P} [g(z_1, \dots, z_n) \geq \mathbb{E}_{z_1, \dots, z_n \sim P}[g] + \epsilon] \leq e^{-2n\epsilon^2}$)

$$\text{Then 4.13: } \mathbb{E}_{z_1, \dots, z_n \sim P} \left[\sup_{f \in F} \left[\frac{1}{n} \sum_{j=1}^n f(z_j) - \mathbb{E}_{z_1, \dots, z_n \sim P}[f(z)] \right] \right] \leq 2R_n(F) \quad (4.82)$$

Step 3: Define

$$\tilde{g}(z_1, \dots, z_n) = R_S(F) \triangleq \mathbb{E}_{\sigma_i} \left[\sup_{f \in F} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right]. \quad (4.83)$$

Using a similar argument to that of Step 1, we show that \tilde{g} satisfies the bounded difference condition.
also

$$\begin{aligned} & \tilde{g}(z_1, \dots, z_n) - \tilde{g}(z_1, \dots, z'_i, \dots, z_n) \\ & \xrightarrow{\text{看这里}} \leq \mathbb{E}_{\sigma_i} \left[\sup_{f \in F} \left[\frac{1}{n} \sum_{j=1}^n \sigma_j f(z_j) \right] - \sup_{f \in F} \left[\left(\frac{1}{n} \sum_{j=1, j \neq i}^n \sigma_j f(z_j) \right) + \frac{1}{n} \sigma_i f(z'_i) \right] \right] \quad (4.84) \\ & \leq \mathbb{E}_{\sigma_i} \left[\sup_{f \in F} \left(\frac{1}{n} \sigma_i (f(z_i) - f(z'_i)) \right) \right] \quad (4.85) \\ & \leq \frac{1}{n}, \quad (4.86) \end{aligned}$$

since the term inside the sup is always upper bounded by 1. We can thus apply McDiarmid's Inequality with parameters $c_1 = \dots = c_n = 1/n$:

$$\Pr[\tilde{g} - \mathbb{E}[\tilde{g}] \geq \epsilon] \leq \exp(-2n\epsilon^2), \quad \text{and} \quad \Pr[\tilde{g} - \mathbb{E}[\tilde{g}] \leq -\epsilon] \leq \exp(-2n\epsilon^2). \quad (\text{similar to step 1}) \quad (4.87)$$

$$R_S(F) = \mathbb{E}[\tilde{g}] = R_n(F)$$

Step 4: We set δ such that $\exp(-2n\epsilon^2) = \delta/2$. (This implies that $\epsilon = \sqrt{\log(2/\delta)/2n}$.) Then, with probability

$$\sup_{f \in F} \left[\frac{1}{n} \sum_{i=1}^n f(z_i) - \mathbb{E}[f] \right] = g \leq \mathbb{E}[g] + \epsilon \quad (\text{Step 1}) \quad (4.88)$$

$$\leq 2R_n(F) + \epsilon \rightarrow \text{we can also use this...} \quad (\text{Step 2}) \quad (4.89)$$

$$\leq 2(R_S(F) + \epsilon) + \epsilon \quad (\text{Step 3}) \quad (4.90)$$

$$= 2R_S(F) + 3\epsilon, \quad (4.91)$$

as required. \square

Setting \mathcal{F} to be a family of loss functions (bounded by $[0, 1]$) in Theorem 4.18 gives the following corollary:

Corollary 4.19. Let \mathcal{F} be a family of loss functions $\mathcal{F} = \{(x, y) \mapsto \ell((x, y), h) : h \in \mathcal{H}\}$ with $\ell((x, y), h) \in [0, 1]$ for all $\ell, (x, y)$ and h . Then, with probability $1 - \delta$, the generalization gap is

$$\hat{L}(h) - L(h) \leq 2R_S(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \quad \text{for all } h \in \mathcal{H}. \quad (4.92)$$

Remark 4.20. If we want to bound the generalization gap by the average Rademacher complexity instead, we can replace the RHS of (4.92) with $2R_n(\mathcal{F}) + \sqrt{\frac{\log(2/\delta)}{2n}}$. negligible $R_n(\mathcal{F})$ like this!

Interpretation of Corollary 4.19. It is typically the case that $O\left(\sqrt{\frac{\log(2/\delta)}{n}}\right) \ll R_S(\mathcal{F})$ and $O\left(\sqrt{\frac{\log(2/\delta)}{n}}\right) \ll R_n(\mathcal{F})$. This is the case because $R_S(\mathcal{F})$ and $R_n(\mathcal{F})$ often take the form $\frac{c}{\sqrt{n}}$ where c is large (depends on the complexity of \mathcal{F})

is a big constant depending on the complexity of \mathcal{F} , whereas we only have a logarithmic term in the numerator of $O\left(\sqrt{\frac{\log(2/\delta)}{n}}\right)$. As a result, we can view the $3\sqrt{\frac{\log(2/\delta)}{n}}$ term in the RHS of Corollary 4.19 as negligible. Another way of seeing this is noting that a $\tilde{O}\left(\frac{1}{\sqrt{n}}\right)$ term is necessary even for the concentration bound of a single function $h \in \mathcal{H}$. Previously, we bounded $L(h) - \hat{L}(h)$ using a union bound over $h \in \mathcal{H}$, which necessarily needs to be larger than $\tilde{O}\left(\frac{1}{\sqrt{n}}\right)$. As a result, the $O\left(\sqrt{\frac{\log(2/\delta)}{n}}\right)$ term is not significant. *negligible!*

4.5.1 Rademacher complexity is translation invariant

A useful fact is that both empirical Rademacher complexity and average Rademacher complexity are translation invariant. (This is not obvious when thinking of how translation affects the picture in Figure 4.3.)

Proposition 4.5.1. Let \mathcal{F} be a family of functions mapping $Z \mapsto \mathbb{R}$ and define $\mathcal{F}' = \{f'(z) = f(z) + c_0 \mid f \in \mathcal{F}\}$ for some $c_0 \in \mathbb{R}$. Then $R_S(\mathcal{F}) = R_S(\mathcal{F}')$ and $R_n(\mathcal{F}) = R_n(\mathcal{F}')$.

Proof. We will prove here that empirical Rademacher complexity is translation invariant.

$$R_S(\mathcal{F}') = \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f' \in \mathcal{F}'} \frac{1}{n} \sum_{i=1}^n \sigma_i f'(z_i) \right] \quad (4.93)$$

$$= \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (f(z_i) + c_0) \right] \quad (4.94)$$

$$= \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\frac{1}{n} \sum_{i=1}^n \sigma_i c_0 + \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right] \quad (4.95)$$

$$= \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right] = R_S(\mathcal{F}), \quad (4.96)$$

where (4.96) follows because $\mathbb{E}_{\sigma_1, \dots, \sigma_n} \frac{1}{n} \sum_{i=1}^n \sigma_i c_0 = 0$, since the σ_i 's are Rademacher random variables. \square

4.6 Covering number upper bounds "Rademacher complexity"

In Chapter 5, we will prove Rademacher complexity bounds that hinge on elegant, ad-hoc algebraic manipulations that may not extend to more general settings. Here, we consider a more fundamental approach for proving empirical Rademacher complexity bounds based on coverings of the output space. (The trade-off is generally more tedious.) $R_S(\mathcal{F})$ Output space의 커버링을 바탕으로 $R_S(\mathcal{F})$ 의 복잡도를 줄여보자.

The first important observation is that for purposes of computing the empirical Rademacher complexity on samples z_1, \dots, z_n ,

$$R_S(\mathcal{F}) = \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right], \quad (4.97)$$

we only care about the output of function $f \in \mathcal{F}$, and not the function itself (i.e. it is sufficient for our purposes to know $f(z_1), \dots, f(z_n)$, but not f). In other words, we can characterize $f \in \mathcal{F}$ by $f(z_1), \dots, f(z_n)$. In the sequel, we will take advantage of this simplification from the (potentially large) space of all functions \mathcal{F} to the output space,

$$\mathcal{Q} \triangleq \left\{ (f(z_1), \dots, f(z_n))^\top : f \in \mathcal{F} \right\} \subseteq \mathbb{R}^n, \quad (4.98)$$

(f ∈ F) → Q: potentially large family
Q: dramatically smaller than F
Q: a vector space, Q! (we prefer to analyze vector space(Q) rather than function space(F))
elements in... output space.

which may be drastically smaller than \mathcal{F} . Correspondingly, the empirical Rademacher complexity can be rewritten as a maximization over the output space \mathcal{Q} instead of the function space \mathcal{F} :

$$\text{Def } R_S(\mathcal{F}): \quad R_S(\mathcal{F}) = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right] \quad R_S(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle \sigma, v \rangle \right].$$

complexity measure for the set \mathcal{Q} .

maximize
over output space.

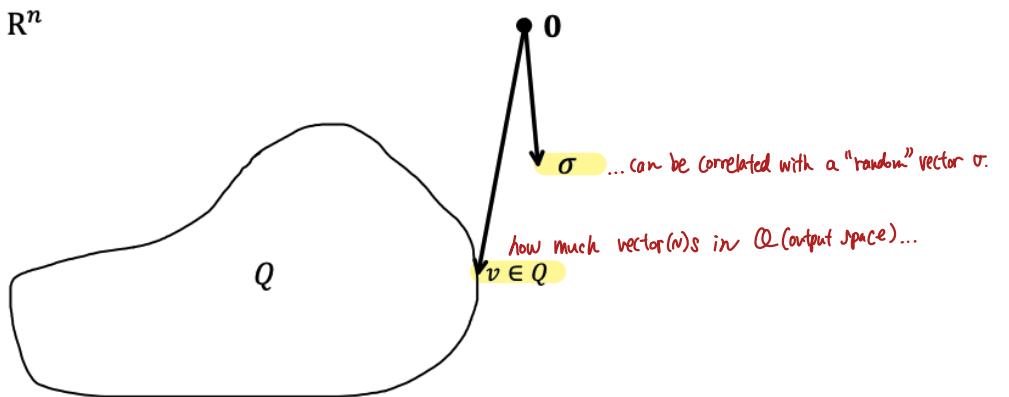
정리: 다른 벡터 v 들이, Rademacher random vector인 경우 "correlated"성이 있다!

In other words, the complexity of \mathcal{F} can be also interpreted as how much the vectors in \mathcal{Q} can be correlated with a random vector σ . (See Figure 4.3 for an illustration of this idea.) One can also view $\mathbb{E}_{\sigma} [\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle \sigma, v \rangle]$ as a complexity measure for the set \mathcal{Q} . If we replace σ by a Gaussian vector with spherical covariance, then the corresponding quantity (without the $\frac{1}{n}$ scaling), $\mathbb{E}_{g \sim N(0, I)} [\sup_{v \in \mathcal{Q}} \langle g, v \rangle]$, is often referred to as the Gaussian complexity of the set \mathcal{Q} . (It turns out that Gaussian complexity and Rademacher complexity are closely related.)

Another corollary of this is that the empirical Rademacher complexity only depends on the functionality of \mathcal{F} but not on the exact parameterization of \mathcal{F} . For example, suppose we have two parameterizations $\mathcal{F} = \{f(x) = \sum \theta_i x_i \mid \theta \in \mathbb{R}^d\}$ and $\mathcal{F}' = \{f(x) = \sum \theta_i^3 \cdot w_i x_i \mid \theta \in \mathbb{R}^d, w \in \mathbb{R}^d\}$. Since $Q_{\mathcal{F}}$ and $Q_{\mathcal{F}'}$ are the same, we see that $R_S(\mathcal{F}) = R_S(\mathcal{F}')$ since our earlier expression for $R_S(\mathcal{F})$ only depends on \mathcal{F} through $Q_{\mathcal{F}}$.

\mathcal{F} and \mathcal{F}' have different parameterizations
but has the same emp. Rademacher comp.

$$Q = \{(f(z_1), \dots, f(z_n))^T : f \in \mathcal{F}\} \subseteq \mathbb{R}^n$$



$$R_S(\mathcal{F}) = \mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle \sigma, v \rangle \right]$$

supremum

Figure 4.3: We can view empirical Rademacher complexity as the expectation of the maximum inner product between σ and $v \in \mathcal{Q}$.



Rademacher complexity of finite hypothesis classes. In practice, we cannot directly evaluate the Rademacher complexity, so we instead bound its value using quantities that are computable. Given finite $|\mathcal{Q}|$, we often rely on the following bound, which is also known as Massart's finite lemma:

Proposition 4.6.1. Let \mathcal{F} be a collection of functions mapping $Z \mapsto \mathbb{R}$ and let \mathcal{Q} be defined as in (4.98). Assume that $\frac{1}{\sqrt{n}} \|v\|_2 \leq M < \infty$ for all $v \in \mathcal{Q}$. Then,

finite

Massart's finite lemma:

$$R_S(\mathcal{F}) \leq \sqrt{\frac{2M^2 \log |\mathcal{Q}|}{n}}$$

$$\frac{1}{\sqrt{n}} \|v\|_2 \leq M < \infty$$

Massart's finite lemma:

$$\bullet \mathcal{Q} = h(f(z_1), \dots, f(z_n))^T : f \in \mathcal{F} \subseteq \mathbb{R}^n \quad (4.100)$$

$$\bullet \frac{1}{\sqrt{n}} \|v\|_2 \leq M < \infty$$

$$\bullet R_S(\mathcal{F}) \leq \sqrt{\frac{2M^2 \log |\mathcal{Q}|}{n}}$$

We prove a (slightly) simplified version of this result in Problem 3(c) of Homework 2, so we omit the proof of Massart's lemma here. Using Massart's lemma, we can also bound the Rademacher complexity in terms of \mathcal{F} . Restating the assumption accordingly,

Massart's lemma: However, $|\mathcal{Q}|$ is typically infinite.

$$R_S(\mathcal{F}) \leq \frac{\sqrt{2M^2 \log |\mathcal{F}|}}{n} \left(\frac{1}{\sqrt{n}} \|\mathbf{M}\|_2 \leq M < \infty \right)$$

$$\frac{1}{\sqrt{n}} \|\mathbf{M}\|_2 \leq M < \infty$$

Corollary 4.21. Let \mathcal{F} be a collection of functions mapping $Z \mapsto \mathbb{R}$. If $\sqrt{\frac{1}{n} \sum_{i=1}^n f(z_i)^2} \leq M$ for all $f \in \mathcal{F}$, then

$$R_S(\mathcal{F}) \leq \sqrt{\frac{2M^2 \log |\mathcal{F}|}{n}}. \quad (4.101)$$

을 위하여 정제하였습니다

$$\mathcal{Q} = \{(f(z_1), \dots, f(z_n))^T : f \in \mathcal{F}\} \subseteq \mathbb{R}^n, |\mathcal{Q}| \leq |\mathcal{F}|$$

Note that Corollary 4.21 yields a looser bound (than Massart's lemma) since $|\mathcal{Q}| \leq |\mathcal{F}|$.

In practice, we rarely apply Massart's lemma directly since $|\mathcal{Q}|$ is typically infinite. In the sequel, we discuss alternative approaches to bounding the Rademacher complexity that are appropriate for this setting.

Bounding "Rademacher complexity" using ϵ -covers. When $|\mathcal{Q}|$ is infinite, we can apply the same discretization trick that we used to prove the generalization bound for an infinite-hypothesis space. This time, instead of trying to cover the parameter space, we will cover the output space. To this end, we first recall a few definitions concerning ϵ -covers.

Definition 4.22. \mathcal{C} is an ϵ -cover of \mathcal{Q} (with respect to metric ρ) if for all $v' \in \mathcal{Q}$, there exists $v \in \mathcal{C}$ such that $\rho(v, v') \leq \epsilon$. (예전 방법처럼 거의 똑같음..!)

Definition 4.23. The covering number is defined as the minimum size of an ϵ -cover, or explicitly:

$$N(\epsilon, \mathcal{Q}, \rho) \triangleq (\min \text{ size of } \epsilon\text{-cover of } \mathcal{Q} \text{ w.r.t. metric } \rho). \quad (4.102)$$

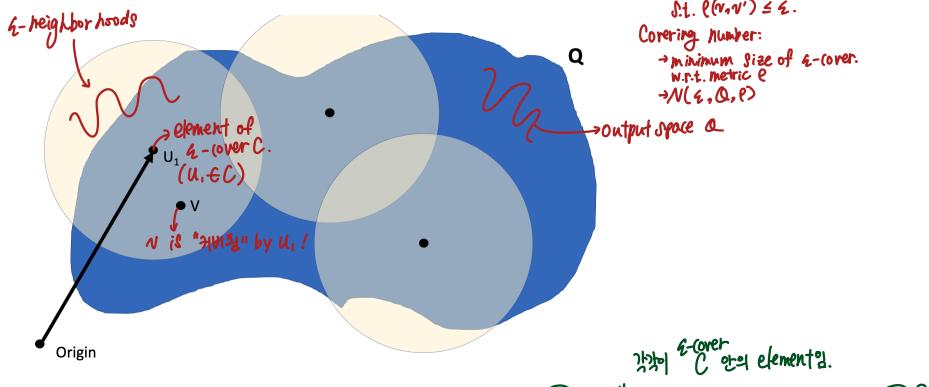


Figure 4.4: We can visualize the ϵ -cover \mathcal{C} (by depicting a set of ϵ -balls that cover the output space \mathcal{Q} .) The yellow circles denote the ϵ -neighborhoods of the covering points $u_i \in \mathcal{C}$.

In subsequent derivations, we will use the metric $\rho(v, v') = \left(\frac{1}{\sqrt{n}} \|v - v'\|_2 \right)$. (in \mathcal{Q})

Remark 4.24. We normalize the ℓ_2 norm (in ρ) by $\frac{1}{\sqrt{n}}$ to simplify comparisons to the functional analysis view of the Rademacher complexity. In the literature, the ϵ -cover of \mathcal{Q} defined above is also referred to as an ϵ -cover of the function class \mathcal{F} under the $L_2(P_n)$ metric.² In particular,

$$\rho(v, v') = \frac{1}{\sqrt{n}} \|v - v'\|_2$$

$$L_2(P_n)(f, f') = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(z_i) - f'(z_i))^2}. \quad (\text{in } \mathcal{F})$$

empirical dist. (defined)

$$\begin{aligned} \epsilon\text{-cover of } \mathcal{Q} \text{ w.r.t. metric } \rho \\ \epsilon\text{-cover of } \mathcal{F} \text{ under } L_2(P_n) \text{ metric} \\ \text{empirical distribution } \left(\frac{1}{n} \sum_{i=1}^n (f(z_i) - f'(z_i))^2 \right)^{1/2} \end{aligned}$$

² P_n denotes the empirical distribution, (i.e. the uniform distribution over the observations z_1, \dots, z_n .) More generally the $L_p(Q)$ metric is defined by $(\mathbb{E}_Q [(f(z) - f'(z))^p])^{1/p}$.

$$L_p(Q) : \left(\mathbb{E}_Q [(f(z) - f'(z))^p] \right)^{1/p}$$

oo.

Recall we have established the following correspondences between the set of functions \mathcal{F} and the output space \mathcal{Q} :

$$f \in \mathcal{F} \iff \begin{pmatrix} f(z_1) \\ \vdots \\ f(z_n) \end{pmatrix} \in \mathcal{Q} \quad (4.104)$$

outputs!

oo.
o/p not = sufficient for our purposes!

We can write a trivial correspondence between both the output and function class points of view as follows:

$$N(\epsilon, \mathcal{F}, L_2(P_n)) = N\left(\epsilon, \mathcal{Q}, \frac{1}{\sqrt{n}} \|\cdot\|_2\right) \quad (4.105)$$

covering number

(Remark 4.24 in 8.)
Recall that...
h-cover of \mathcal{Q} w.r.t. metric ρ
 ϵ -cover of \mathcal{F} under $L_2(P_n)$ metric
empirical distribution $(\frac{1}{n} \sum_{i=1}^n (f(z_i) - f(z_i))^2)^{1/2}$

The results below will be stated in the function-space notation, but in the proofs we will shift to the \mathcal{Q} -formulation for the sake of clarity. In general, we prefer to reason about covering numbers on \mathcal{Q} as it is more natural to analyze vector spaces (compared to function spaces). $\exists \mathcal{F}$ (function space) yet \mathcal{Q} (vector space) \cong analysis.

Equipped with the definition of minimal ϵ -covers, we can prove the following Rademacher complexity bound:

• Bound $R_S(\mathcal{F})$ when $|\mathcal{Q}|$ is infinite! (Using discretization trick)

Theorem 4.25. Let \mathcal{F} be a family of functions $Z \mapsto [-1, 1]$. Then

$$R_S(\mathcal{F}) \leq \inf_{\epsilon > 0} \left(\epsilon + \sqrt{\frac{2 \log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} \right). \quad (4.106)$$

↑ is bounded.

Proposition (4.6.1) \Rightarrow Rademacher finite lemma:
 $\mathcal{Q} = h(f(z_1), \dots, f(z_n))^T: \mathcal{F} \subset \mathbb{R}^n$
 $\frac{1}{n} \|\mathcal{Q}\|_2 \leq M < \infty$
 $\hookrightarrow R_S(\mathcal{F}) \leq \sqrt{\frac{2M^2 \log |\mathcal{Q}|}{n}}$

$= N(\epsilon, \mathcal{Q}, \rho)$ (where we are using $\rho = \frac{1}{\sqrt{n}} \|\cdot\|_2$ here)

The ϵ term can be thought of as the discretization error, while the second term is the Rademacher complexity of the finite ϵ -cover. The precise form of this complexity bound follows from Proposition 4.6.1.

Proof. Fix any $\epsilon > 0$. Let \mathcal{C} be the minimal ϵ -cover of \mathcal{Q} (with respect to the metric $\rho(v, v') = \frac{1}{\sqrt{n}} \|v - v'\|_2$)

Note that $|\mathcal{C}| = N(\epsilon, \mathcal{Q}, \frac{1}{\sqrt{n}} \|\cdot\|_2) = N(\epsilon, \mathcal{F}, L_2(P_n))$.

We aim to bound $R_S(\mathcal{F}) = \mathbb{E}_\sigma [\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle]$ by approximating v with $v' \in \mathcal{C}$. In particular, for every point $v \in \mathcal{Q}$, choose $v' \in \mathcal{C}$ such that $\rho(v, v') \leq \epsilon$ (and z is small (specifically, $\frac{1}{\sqrt{n}} \|z\|_2 \leq \epsilon$). This gives

$$\frac{1}{n} \langle v, \sigma \rangle = \frac{1}{n} \langle v', \sigma \rangle + \frac{1}{n} \langle z, \sigma \rangle \quad (z \text{ is a vector b/w } v \text{ and } v') \quad (4.107)$$

v' is the closest point to v in the minimal ϵ -cover of \mathcal{Q} .
 z is a vector b/w v and v'
 $(z \triangleq v - v')$

$$\text{approximate } v \text{ with } v' \in \mathcal{C} \leq \frac{1}{n} \langle v', \sigma \rangle + \frac{1}{n} \|z\|_2 \|\sigma\|_2 \quad (z \triangleq v - v', \text{ Cauchy-Schwarz}) \quad (4.108)$$

$$\rightarrow \text{for every } v \in \mathcal{Q}, \quad \leq \frac{1}{n} \langle v', \sigma \rangle + \epsilon. \quad (\text{since } \|z\|_2 \leq \sqrt{n}\epsilon \text{ and } \|\sigma\|_2 \leq \sqrt{n}) \quad (4.109)$$

\checkmark

\downarrow s.t. $\rho(v, v') \leq \epsilon$ and $z \triangleq v - v'$ is small.

Taking the expectation of the supremum on both sides of this inequality gives $\frac{1}{n} \|\sigma\|_2 \leq 1 < \infty$ (needed when using prop. 4.6.1)

$$R_S(\mathcal{F}) = \mathbb{E}_\sigma \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle \right] \quad (4.110)$$

$$\frac{1}{n} \langle v, \sigma \rangle \leq \frac{1}{n} \langle v', \sigma \rangle + \epsilon \quad (4.109) \quad \downarrow \quad \leq \mathbb{E}_\sigma \left[\sup_{v' \in \mathcal{C}} \left(\frac{1}{n} \langle v', \sigma \rangle + \epsilon \right) \right] \quad (4.111)$$

$$= \epsilon + \mathbb{E}_\sigma \left[\sup_{v' \in \mathcal{C}} \left(\frac{1}{n} \langle v', \sigma \rangle \right) \right] \quad (4.112)$$

$$R_S(\mathcal{F}) \leq \sqrt{\frac{2 \log |\mathcal{C}|}{n}} \quad \text{with } \frac{1}{n} \|\sigma\|_2 \leq 1 < \infty \quad \leq \epsilon + \sqrt{\frac{2 \log |\mathcal{C}|}{n}} \quad (\text{Proposition 4.6.1}) \quad (4.113)$$

$$= \epsilon + \sqrt{\frac{2 \log N(\epsilon, \mathcal{Q}, \rho)}{n}} \quad (4.114)$$

$$= \epsilon + \sqrt{\frac{2 \log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} \quad (\text{Remark 4.24}) \quad (4.115)$$

holds for any $\epsilon > 0$
 \rightarrow fake infimum
over all ϵ
 \Rightarrow arrive (4.106)!

Since the argument above holds for any $\epsilon > 0$, we can take the infimum over all ϵ to arrive at Equation (4.106).

$$R_S(\mathcal{F}) \leq \inf_{\epsilon > 0} \left(\frac{1}{n} \langle z, \sigma \rangle + \sqrt{\frac{2 \log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} \right)$$

□ ↗

10/28
11/06

4.6.1 Chaining and Dudley's theorem

While Theorem 4.25 is useful, the bound in (4.108) is rarely tight as z might not be perfectly correlated with σ . It is possible to obtain a stronger theorem by constructing a "chained ϵ -covering scheme." Specifically, when we decompose $v = v' + z$, we can construct a finer-grained covering of the ball $B(v', \epsilon)$, and then we can decompose z into smaller components and so on (see Figure 4.5 for an illustration). $R_S(\mathcal{F}) \leq \epsilon$

Using this method of chaining, we can obtain the following (stronger) result:

Theorem 4.26 (Dudley's Theorem). *If \mathcal{F} is a function class from $Z \mapsto \mathbb{R}$, then*

$$R_S(\mathcal{F}) \leq 12 \int_0^\infty \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon. \quad (4.116)$$

$\hookrightarrow f$ is not bounded. (no need to be)
useful when this integral is finite)

Note that unlike in Theorem 4.25, we do not require $f \in \mathcal{F}$ to be bounded.

It is not obvious how (4.116) improves upon the one-step discretization bound given by (4.106). At a high level, we can interpret this bound as removing the discretization error term by averaging over different scales of ϵ . But before we can explicitly prove this claim, we motivate our approach. In the proof of Theorem 4.25, we approximated v with $v' + z$ where v' is the closest point to v in the minimal ϵ -cover of \mathcal{Q} , and z is the vector between v' and v . In particular,

$$\frac{1}{n} \langle v, \sigma \rangle = \frac{1}{n} \langle v', \sigma \rangle + \frac{1}{n} \langle z, \sigma \rangle \quad (4.117)$$

$\frac{1}{n} \|z\|_2 \leq \epsilon$
 \hookrightarrow
 $\frac{1}{n} \|z\|_2 \leq \epsilon$

Then, to obtain a bound, we take a sup of both sides, but apply the sup separately to each term on the right hand side. Namely, we show that:

$$\mathbb{E} \left[\sup_v \frac{1}{n} \langle v, \sigma \rangle \right] \leq \mathbb{E} \left[\sup_{v' \in \mathcal{C}} \frac{1}{n} \langle v', \sigma \rangle \right] + \mathbb{E} \left[\sup_{z \in B_{v'}} \frac{1}{n} \langle z, \sigma \rangle \right] \quad (4.118)$$

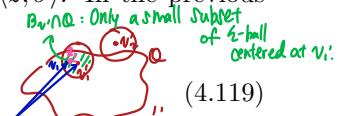
\hookrightarrow
 $\frac{1}{n} \|z\|_2 \leq \epsilon$
 \hookrightarrow
 $\frac{1}{n} \|z\|_2 \leq \epsilon$

(This bound follows by observing that $\mathbb{E}[\sup(A+B)] \leq \mathbb{E}[\sup A] + \mathbb{E}[\sup B]$ since the sup on the RHS is taken separately over both terms.) The difficult term to tightly bound is the last one, $\frac{1}{n} \langle z, \sigma \rangle$. In the previous derivation, we naively upper bounded $\langle z, \sigma \rangle$ using Cauchy-Schwarz,

$$\frac{1}{n} \langle z, \sigma \rangle \leq \frac{\|z\|_2 \cdot \|\sigma\|_2}{n}, \quad \langle z, \sigma \rangle \leq \|z\|_2 \cdot \|\sigma\|_2 \quad (4.119)$$

$\frac{1}{n} \|z\|_2 \leq \epsilon$
 \hookrightarrow
 $\frac{1}{n} \|z\|_2 \leq \epsilon$

but this bound is only tight if there exists $z \in B_{v'}$ that is perfectly correlated with σ . We claim that such perfect correlation is unlikely. (Recall that the output space is defined by possible outputs of $f \in \mathcal{F}$ given n inputs. Unless our function class is extremely expressive, the set of radius ϵ (around v' contained in \mathcal{Q}) will only be a small subset of the ϵ -ball (centered at v'), thus, $\sup_z \frac{1}{n} \langle z, \sigma \rangle \ll \frac{\|z\|_2 \cdot \|\sigma\|_2}{n}$.) Empirical



To precisely set up our approach, we observe that $\mathbb{E}[\sup_{z \in B_{v'}} \frac{1}{n} \langle z, \sigma \rangle]$ is itself a Rademacher complexity:

$R_S(B_{v'} \cap \mathcal{Q})$. To more tightly bound $\mathbb{E}[\sup_{z \in B_{v'}} \frac{1}{n} \langle z, \sigma \rangle]$, we then repeat the ϵ -covering argument again with a smaller choice of ϵ . Intuitively, this procedure amounts to decomposing $\langle z, \sigma \rangle$ from (4.117) into another pair of terms (corresponding to the new ϵ -cover and the discretization error). "Chaining" then repeats this decomposition countably many times. This procedure is illustrated visually by Figure 4.5, and we formalize this argument in the sequel.

Proof. Let $\epsilon_0 = \sup_{f \in \mathcal{F}} \max_i |f(z_i)|$, so that for all $v \in \mathcal{Q}$,

$$\begin{aligned} \epsilon_0 &= \sup_{f \in \mathcal{F}} \max_i |f(z_i)| \\ &\geq \sqrt{\frac{1}{n} \sum_{i=1}^n f(z_i)^2} = \sqrt{\frac{1}{n} \|v\|_2^2}. \end{aligned} \quad (4.120)$$

$\frac{1}{n} \|z\|_2^2 \leq \epsilon_0^2$
 \hookrightarrow
 $\frac{1}{n} \|z\|_2^2 \leq \epsilon_0^2$

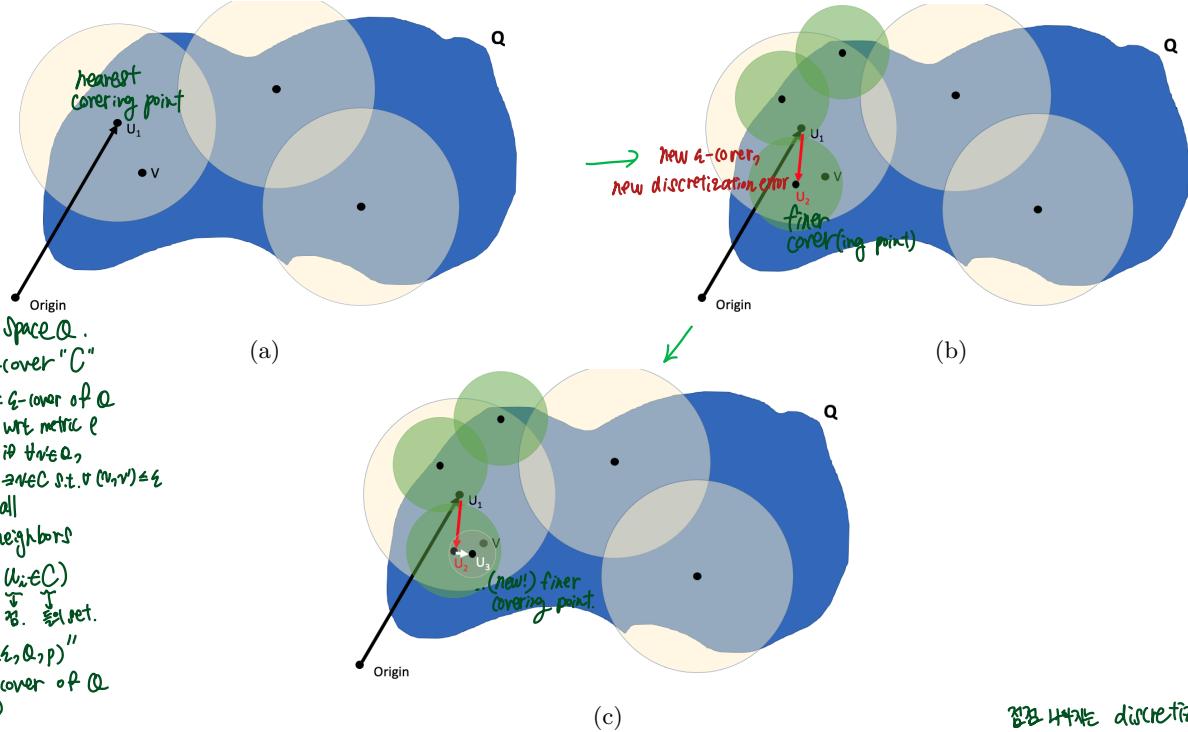


Figure 4.5: We depict how the chaining procedure approximates v using a sequence of progressively finer discretizations. Figure 4.5a illustrates how we first approximate v using the nearest covering point u_1 , while Figures 4.5b and 4.5c describe how we refine this approximation using two finer covers, whose nearest points are denoted by u_2 and u_3 , respectively.

Define $\epsilon_j = 2^{-j}\epsilon_0$ and let C_j be an ϵ_j -cover of Q . Then, C_0 is the coarsest cover of Q , and as j increases, we obtain progressively more fine-grained covers C_j . We can intuitively think of these covers as nested, but this is not necessary for the proof to hold. We next use this sequence of covers to define a telescoping series that equals v ; the terms in this series can then be analyzed using the tools that we have developed in the prequel. 이제 정리하는 것에는 텔레스코프 합이 필요!

For $v \in Q$, let u_i denote the nearest neighbor of v in C_i . (Note that by definition $\rho(u_i, v) \leq \epsilon_i$.) Taking $u_0 = 0$, it follows from our definition of C_i that as $j \rightarrow \infty$, $\epsilon_j \rightarrow 0$ and $u_j \rightarrow v$. Leveraging these observations, we can express v using the following series:

$$v = u_1 + (u_2 - u_1) + (u_3 - u_2) + \dots \quad (4.121)$$

$$= (u_1 - u_0) + (u_2 - u_1) + (u_3 - u_2) + \dots \quad (4.122)$$

$$= \sum_{i=1}^{\infty} (u_i - u_{i-1}). \quad (4.123)$$

$$N = \sum_{i=1}^{\infty} (u_i - u_{i-1})$$

모든 디스크리트이션은, $(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$
 1/2 approximate set!

$$N = \sum_{x=1}^{\infty} N(\epsilon_x, \mathcal{F}_{\epsilon_x})$$

Substituting (4.123) in the Rademacher complexity we aim to bound, we obtain

$$\mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle \right] = \mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \sum_{i=1}^{\infty} \langle u_i - u_{i-1}, \sigma \rangle \right] \quad (4.124)$$

$$\text{LHS} \quad \sup_{v \in \mathcal{Q}} \leq \sup_{u_i \in \mathcal{C}_i, u_{i-1} \in \mathcal{C}_{i-1}} \frac{1}{n} \langle u_i - u_{i-1}, \sigma \rangle \quad (4.125)$$

$$= \sum_{i=1}^{\infty} \mathbb{E} \left[\sup_{(u_i \in \mathcal{C}_i, u_{i-1} \in \mathcal{C}_{i-1})} \frac{1}{n} \langle u_i - u_{i-1}, \sigma \rangle \right]. \quad (4.126)$$

Observe that

$$\mathbb{E} \left[\sup_{u_i \in \mathcal{C}_i, u_{i-1} \in \mathcal{C}_{i-1}} \frac{1}{n} \langle u_i - u_{i-1}, \sigma \rangle \right] \quad (4.127)$$

is a Rademacher complexity defined over the *finite* space $\mathcal{C}_i \times \mathcal{C}_{i-1}$, so we can use Proposition 4.6.1 (Massart's lemma) to obtain a tractable upper bound. To do so, we must first compute an upper bound on $\frac{1}{\sqrt{n}} \|u_i - u_{i-1}\|_2$:

$$\frac{1}{\sqrt{n}} \|u_i - u_{i-1}\|_2 = \frac{1}{\sqrt{n}} \| (u_i - v) - (u_{i-1} - v) \|_2 \quad (4.128)$$

$$\leq \frac{1}{\sqrt{n}} (\|u_i - v\|_2 + \|u_{i-1} - v\|_2) \quad (4.129)$$

$$\leq \epsilon_i + \epsilon_{i-1} \quad \rho(u_i, v) \leq \epsilon_i \quad \text{distance metric} \quad (4.130)$$

$$= 3\epsilon_i \quad 2^{-(i-1)} \epsilon_0 = 2\epsilon_i \quad (\epsilon_{i-1} \triangleq 2\epsilon_i) \quad (4.131)$$

Now we apply Proposition 4.6.1 with $M = 3\epsilon_i$ and $|\mathcal{Q}| = |\mathcal{C}_i \times \mathcal{C}_{i-1}| \leq |\mathcal{C}_i| \cdot |\mathcal{C}_{i-1}|$

$$\mathbb{E} \left[\sup_{u_i \in \mathcal{C}_i, u_{i-1} \in \mathcal{C}_{i-1}} \frac{1}{n} \langle u_i - u_{i-1}, \sigma \rangle \right] \leq \sqrt{\frac{2(3\epsilon_i)^2 \log(|\mathcal{C}_i| \cdot |\mathcal{C}_{i-1}|)}{n}} \quad (4.132)$$

- Massart's finite lemma:
- $\mathcal{Q} = \{f(f(z_1), \dots, f(z_n))^\top : f \in \mathcal{F}\} \subseteq \mathbb{R}^n$
- $\frac{1}{n} \|M\|_2 \leq M < \infty$ true
- $R_s(F) \leq \sqrt{\frac{2n^2 \log |P|}{n}}$

$$= \frac{3\epsilon_i}{\sqrt{n}} \sqrt{2(\log |\mathcal{C}_i| + \log |\mathcal{C}_{i-1}|)} \quad (4.133)$$

$$\leq \frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log |\mathcal{C}_i|} \quad (|\mathcal{C}_i| \geq |\mathcal{C}_{i-1}|) \quad (4.134)$$

Substitute w. N($\epsilon_i, \mathcal{F}, L_2(P_n)$) (Cause they are the same)

Applying (4.134) to each term in (4.126) and substituting the covering number $N(\epsilon_i, \mathcal{F}, L_2(P_n))$ for $|\mathcal{C}_i|$, we obtain the following upper bound on the Rademacher complexity:

$$\mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle \right] \leq \sum_{i=1}^{\infty} \left(\frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log N(\epsilon_i, \mathcal{F}, L_2(P_n))} \right) = \sum_{i=1}^{\infty} \frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log |\mathcal{C}_i|} \quad (4.135)$$

X Finally, we must relate (4.135) to the target upper bound of $12 \int \frac{1}{\sqrt{n}} \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon$. Examining Figure 4.6, we can make two crucial observations. First, for sufficiently large ϵ , $\log N(\epsilon, \mathcal{F}, L_2(P_n)) = 0$ since one point is sufficient to construct a cover. Second, we observe that (for $\epsilon > \epsilon_0$)

$$(\epsilon_i - \epsilon_{i+1}) \sqrt{\log |\mathcal{C}_i|} \leq \int_{\epsilon_{i+1}}^{\epsilon_i} \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon \quad (4.136)$$

from the graph (for $\epsilon > \epsilon_0$) ① for ϵ large, $\log N(\epsilon, \mathcal{F}, L_2(P_n)) = 0$ (only point is enough to cover the set)

then, size of ϵ -cover is 1
if only need one dot, $|\mathcal{C}_i| = 1$
then $\log |\mathcal{C}_i| = \log 1 = 0$
thus $\log N(\epsilon, \mathcal{F}, L_2(P_n)) = 0$

since the LHS of (4.136) is the area of the dotted rectangle illustrated in Figure 4.6 while the RHS is the area under the curve for that interval. (Formally, this result is equivalent to observing that the right Riemann sum underestimates the integral for monotone decreasing functions f .

The RHS $\rightarrow N(\epsilon_i, \mathcal{F}, L_2(P_n))$ is monotone dec. func. (in ϵ)

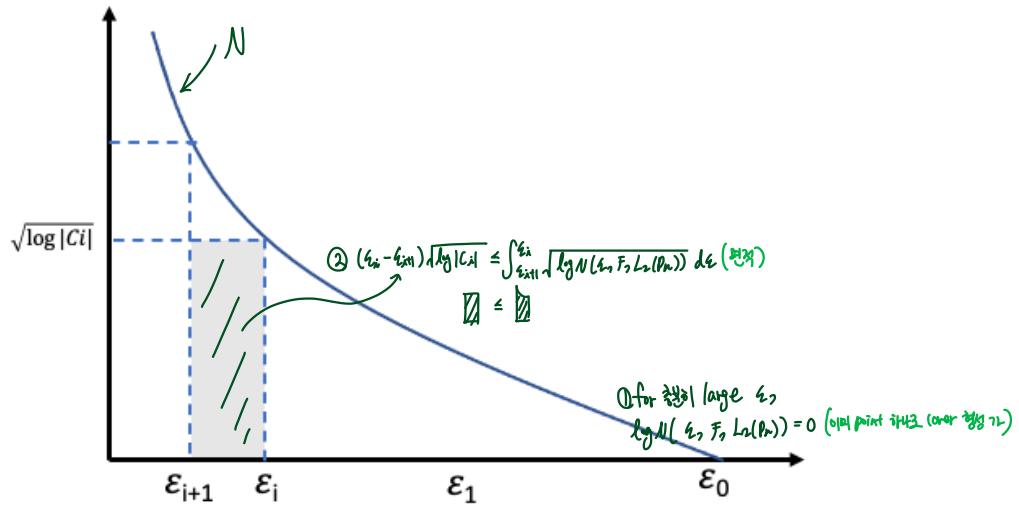


Figure 4.6: We observe that $\log N(\epsilon, \mathcal{F}, L_2(P_n))$ is monotone decreasing in ϵ . The area of the dotted rectangle formed by the vertical lines at ϵ_{i+1} and ϵ_i equals (up to a constant factor) the i -th term of the infinite sum derived in our proof of Dudley's theorem (4.135). The figure shows that the area of this rectangle is no larger than the integral of $\log N(\epsilon, \mathcal{F}, L_2(P_n))$ over this same interval.

Recognizing that $\epsilon_i - \epsilon_{i+1} = \frac{\epsilon_i}{2}$, we note that the LHS of (4.136) is equal (up to a constant factor) to the i -th term of (4.135). Thus,

$$(4.135): \mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle \right] \leq \sum_{i=1}^{\infty} \frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log N(\epsilon_i, \mathcal{F}, L_2(P_n))} = \sum_{i=1}^{\infty} \frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log |C_i|} \quad \text{with term of } \frac{1}{2} \leq \frac{1}{2}$$

$$\sum_{i=1}^{\infty} \frac{6\epsilon_i}{\sqrt{n}} \sqrt{\log N(\epsilon_i, \mathcal{F}, L_2(P_n))} = \frac{12}{\sqrt{n}} \sum_{i=1}^{\infty} (\epsilon_i - \epsilon_{i+1}) \sqrt{\log N(\epsilon_i, \mathcal{F}, L_2(P_n))} \quad (4.137)$$

$$6\epsilon_i = 12\epsilon_i - 6\epsilon_i = 12\epsilon_i - 12\epsilon_{i+1} \leq \frac{12}{\sqrt{n}} \int_{\epsilon_{i+1}}^{\epsilon_i} \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon \quad (4.138)$$

$$= \frac{12}{\sqrt{n}} \int_0^{\epsilon_0} \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon. \quad (4.139)$$

To complete the proof, observe that $\log N(\epsilon, \mathcal{F}, L_2(P_n)) = 0$ for all $\epsilon > \epsilon_0$. This allows us to extend the upper limit of the integral given by (4.139) to ∞ and yields the desired result:

$$\mathbb{E} \left[\sup_{v \in \mathcal{Q}} \frac{1}{n} \langle v, \sigma \rangle \right] \leq \frac{12}{\sqrt{n}} \int_0^{\infty} \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon. \quad (4.140) \quad \text{(From 4.116 to 4.140)}$$

□

Remark 4.27. If \mathcal{F} consists of functions bounded in $[-1, 1]$, then we have that for all $\epsilon > 1$, $N(\epsilon, \mathcal{F}, L_2(P_n)) = 1$. To see this, choose $\{f \equiv 0\}$, which is a complete cover for $\epsilon > 1$. Hence, the limits of integration in (4.116) can be truncated to $[0, 1]$:

$$R_S(\mathcal{F}) \leq 12 \int_0^1 \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon, \quad (4.141)$$

Dudley's Theorem

since $\log N(\epsilon, \mathcal{F}, L_2(P_n)) = 0$ for $\epsilon > 1$.

✓

$$\text{Dudley's} \quad (4.116) \quad \mathbb{E} \left[\sup_{\eta \in \mathcal{N}} \frac{1}{n} \langle \eta, \tau \rangle \right] \leq \frac{R}{n} \int_0^\infty \sqrt{\log N(\epsilon, \mathcal{F}, L_2(P_n))} d\epsilon.$$

$$R_S(\mathcal{F}) \leq n \int_0^\infty \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon$$

4.6.2 Translating Covering Number Bounds to Rademacher Complexity

Of course, the bound in (4.116) is only useful if the integral on the RHS is finite. Here are some setups where this is the case (we continue to assume that the functions in \mathcal{F} are bounded in $[-1, 1]$):

1. If after ignoring multiplicative and additive constants,

$$N(\epsilon, \mathcal{F}, L_2(P_n)) \approx (1/\epsilon)^R, \quad \text{then we have } \log N(\epsilon, \mathcal{F}, L_2(P_n)) \approx R \log(1/\epsilon). \quad (4.142)$$

then we have $\log N(\epsilon, \mathcal{F}, L_2(P_n)) \approx R \log(1/\epsilon)$. We can plug this into the RHS of (4.116) to get

$$\int_0^1 \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon = \int_0^1 \sqrt{\frac{R \log(1/\epsilon)}{n}} d\epsilon \approx \sqrt{\frac{R}{n}}. \quad (4.143)$$

2. If after ignoring multiplicative and additive constants, for some a ,

$$N(\epsilon, \mathcal{F}, L_2(P_n)) \approx a^{R/\epsilon}, \quad \text{then we have } \log N(\epsilon, \mathcal{F}, L_2(P_n)) \approx \frac{R}{\epsilon} \log a. \quad (4.144)$$

then we have $\log N(\epsilon, \mathcal{F}, L_2(P_n)) \approx \frac{R}{\epsilon} \log a$. The bound in (4.116) becomes

$$\int_0^1 \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon \approx \int_0^1 \sqrt{\frac{R}{n\epsilon} \log a} d\epsilon \quad (4.145)$$

$$= \sqrt{\frac{R}{n} \log a} \int_0^1 \sqrt{\frac{1}{\epsilon}} d\epsilon \quad (4.146)$$

$$= \tilde{O}\left(\sqrt{\frac{R}{n}}\right). \quad (4.147)$$

3. If the covering number has the form $N(\epsilon, \mathcal{F}, L_2(P_n)) \approx a^{R/\epsilon^2}$, then $\log N(\epsilon, \mathcal{F}, L_2(P_n)) \approx \frac{R}{\epsilon^2} \log a$. In this case we have:

$$\int_0^1 \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon \approx \sqrt{\frac{R}{n} \log a} \underbrace{\int_0^1 \frac{1}{\epsilon} d\epsilon}_{\text{cannot be bounded}} = \infty, \quad (4.148)$$

i.e. the bound in (4.116) is vacuous. This is because of the behavior of $\epsilon \mapsto 1/\epsilon^2$ near 0: the function goes to infinity too quickly for us to upper bound its integral. Fortunately, there is an “improved” version of Dudley’s theorem that is applicable here:

Theorem 4.28 (Localized Dudley’s Theorem). *If \mathcal{F} is a function class from $Z \mapsto \mathbb{R}$, then for any fixed cutoff $\alpha \geq 0$ we have the bound*

$$R_S(\mathcal{F}) \leq 4\alpha + 12 \int_\alpha^\infty \sqrt{\frac{\log N(\epsilon, \mathcal{F}, L_2(P_n))}{n}} d\epsilon. \quad (4.149)$$

The proof of this theorem is similar to the proof of the original Dudley’s theorem, except that the iterative covering procedure is stopped at the threshold $\epsilon = \alpha$, at the cost of the extra 4α term above. Stopped before $0 \rightarrow \infty$??

Theorem 4.28 allows us to avoid the problematic region around $\epsilon = 0$ in the integral in (4.116). If we let $\alpha = 1/\text{poly}(n)$, where $\text{poly}(n)$ denotes some polynomial function of n , the bound in (4.149) becomes

*with the covering number
(with the form of)*
 $N(\epsilon, \mathcal{F}, L_2(P_n)) = \alpha^{\frac{R}{\epsilon^2}}$

$$\text{LHS (4.149)} \leq \frac{1}{\text{poly}(n)} + \frac{\sqrt{R \log a}}{\sqrt{n}} \int_{\alpha}^1 \frac{1}{\epsilon} d\epsilon \quad (4.150)$$

$$= \frac{1}{\text{poly}(n)} + \frac{\sqrt{R \log a}}{\sqrt{n}} \log(1/\alpha) \quad (4.151)$$

$$= \tilde{O}\left(\sqrt{\frac{R}{n}}\right). \quad (4.152)$$

$y_a = \text{poly}(n)$

The last line follows by observing that $\log(1/\alpha) = \log \text{poly}(n)$.

In summary, we have that $R_S(\mathcal{F}) \leq \tilde{O}\left(\sqrt{\frac{R}{n}}\right)$ for covering numbers of the form $R \log(1/\epsilon)$, $\frac{R}{\epsilon} \log a$, or $\frac{R}{\epsilon^2} \log a$ for some a . Note that if the dependence on ϵ is $1/\epsilon^c$ for $c > 2$, then even the improved Dudley's theorem does not help us. This is because the $\log(1/\alpha)$ term above becomes $\alpha^{1-c/2}$; then, for $\alpha = 1/\text{poly}(n)$, the second term in Dudley's integral is no longer $\tilde{O}\left(\sqrt{\frac{R}{n}}\right)$.

✓

4.6.3 Lipschitz composition

Covering numbers also interact nicely with composition by Lipschitz functions. The following result is the analog of Talagrand's lemma for Rademacher complexity (Lemma 5.3), but its proof is much more elementary as given below. We will use this Lemma in Section 5.5 when bounding the covering number of deep nets.

Lemma 4.29. Suppose ϕ is κ -Lipschitz, and $\rho = L_2(P_n)$. Then,

$$\log N(\epsilon, \phi \circ \mathcal{F}, \rho) \leq \log N(\epsilon/\kappa, \mathcal{F}, \rho) \quad (4.153)$$

Proof. Let \mathcal{C} denote an ϵ/κ -cover for \mathcal{F} . Then $\phi \circ \mathcal{C}$ is an ϵ -cover of $\phi \circ \mathcal{F}$.

$$\rho(\phi \circ f', \phi \circ f) = \sqrt{\frac{1}{n} \sum (\phi(f'(z_i)) - \phi(f(z_i)))^2} \quad (4.154)$$

$$\leq \sqrt{\frac{1}{n} \cdot \kappa^2 \sum (f'(z_i) - f(z_i))^2} \quad (4.155)$$

$$\leq \kappa \cdot \frac{\epsilon}{\kappa} = \epsilon \quad (4.156)$$

□ ✓

4.7 VC dimension and its limitations

Upper bound of Rademacher complexity

In this section, we briefly discuss a classical notion of complexity measure of function class, VC dimension. We will show that VC dimension is an upper bound on the Rademacher complexity. We will focus on classification and will be working within the framework of supervised learning stated in Chapter 1. The labels belong to the output space $\mathcal{Y} = \{-1, 1\}$, each classifier is a function $h : \mathcal{X} \rightarrow \mathbb{R}$ for all $h \in \mathcal{H}$, and the prediction is the sign of the output, i.e. $\hat{y} = \text{sgn}(h(x))$. We will look at zero-one loss, i.e. $\ell_{0-1}((x, y), h) = \mathbb{1}(\text{sgn}(h(x)) \neq y)$. (Note that we can re-express the loss function as

$$\ell_{0-1}((x, y), h) = \frac{1 - \text{sgn}(h(x))y}{2}. \quad (4.157)$$

*loss function (zero-one)
labels: $y = \{-1, 1\}$
classifier: $h: \mathcal{X} \rightarrow \mathbb{R}$ for all $h \in \mathcal{H}$
prediction: sign of the output
 $\text{sgn}(h(x))$*

*VC dimension.
 $\ell_{0-1}(x, y, h) = \mathbb{1}(\text{sgn}(h(x)) \neq y)$,
 $= \frac{1 - \text{sgn}(h(x))y}{2}$*

*↑
Rademacher Complexity*

The first approach is to reason directly about the Rademacher complexity of ℓ_{0-1} loss, i.e. considering the family of functions $\mathcal{F} = \{z = (x, y) \mapsto \ell_{0-1}((x, y), h) : h \in \mathcal{H}\}$. Define Q to be the set of all possible outputs

$$\mathcal{F} = \{(x, y) \mapsto \ell_{0-1}(x, y), h : h \in \mathcal{H}\}$$

\mathbb{Q} : set of all possible outcomes on our dataset

$\exists \mathbb{Q} = \{\text{sgn}(h(x^{(1)})), \dots, \text{sgn}(h(x^{(n)})) \mid h \in \mathcal{H}\}$

possible outcome

on our dataset: $Q = \{\text{sgn}(h(x^{(1)})), \dots, \text{sgn}(h(x^{(n)})) \mid h \in \mathcal{H}\}$. Then, using our earlier remark about viewing the empirical Rademacher complexity as an inner product between $v \in Q$ and σ , we have

$$R_S(\mathcal{F}) = \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{\ell_{0-1}(x_i, h_i)}{2} \right] \quad (4.158)$$

$$= \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{\text{sgn}(h(x^{(i)}))}{2} \right] \quad y = f_{1,1} \rightarrow \text{sgn}(y)$$

$$= \frac{1}{2} \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{v \in Q} \frac{1}{n} \langle \sigma, v \rangle \right]. \quad \begin{array}{l} R_S(\mathcal{F}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right] \\ = \mathbb{E}_{\sigma} \left[\sup_{v \in Q} \frac{1}{n} \langle \sigma, v \rangle \right] \\ (\text{how much the vector in } Q \text{ is } v \text{ correlated with a random vec } v) \end{array} \quad (4.160)$$

Notice that the supremum is now over Q (instead of \mathcal{F}). If n is sufficiently large, then it is typically the case that $|Q| > |\mathcal{F}|$. To see why this is the case, note that each function f corresponds to a single element in Q . However, as n increases, $|Q|$ increases as well. For any particular $v \in Q$, notice that $\langle v, \sigma \rangle$ is a sum of bounded random variables, so we can use Hoeffding's inequality to obtain

$$\Pr \left[\frac{1}{n} \langle \sigma, v \rangle \geq t \right] \leq \exp(-nt^2/2).$$

$$\begin{aligned} &\stackrel{(4.17)(24p)}{=} \frac{2|\mathcal{H}| \exp(-2ne^2)}{\delta}, \text{ then it follows that} \\ &\frac{e^{-2ne^2}}{2|\mathcal{H}|} \leq \frac{\delta}{2n} \Rightarrow e^{-2ne^2} \leq \frac{\delta^2}{4n^2} \Rightarrow n^2 \leq \frac{4\ln(2/\delta)}{e^2} \end{aligned} \quad (4.161)$$

Taking the union bound over $v \in Q$, we see that

$$\Pr \left[\exists v \in Q \text{ such that } \frac{1}{n} \langle \sigma, v \rangle \geq t \right] \leq |Q| \exp(-nt^2/2). \quad \begin{array}{l} \text{at most} \\ \hookrightarrow \Pr \left[\sup_{v \in Q} \frac{1}{n} \langle \sigma, v \rangle \geq t \right] \leq |Q| e^{-\frac{nt^2}{2}} \Rightarrow \delta. \end{array}$$

Thus, with probability at least $1 - \delta$, it is true that $\sup_{v \in Q} \frac{1}{n} \langle v, \sigma \rangle \leq \sqrt{\frac{2(\log |Q| + \log(2/\delta))}{n}}$. (Similarly, we can show that $\mathbb{E} [\sup_{v \in Q} \frac{1}{n} \langle v, \sigma \rangle] \leq O \left(\sqrt{\frac{\log |Q| + \log(2/\delta)}{n}} \right)$ holds.)

The key point to notice here is that the upper bound on $R_S(\mathcal{F})$ depends on $\log |Q|$. VC dimension is one way that we deal with bounding the size of Q . (We will not delve into the details of this approach (for those interested, see Section 3.11 of [Liang, 2016]).) VC dimension, however, has a number of limitations. For one, we will always end up with a bound that depends somehow on the dimension. For linear models, we obtain a bound $\log |Q| \lesssim d \log n$, corresponding to a bound on Rademacher complexity that looks like

$$R_S(\mathcal{F}) \leq \tilde{O} \left(\sqrt{\frac{d}{n}} \right), \quad \begin{array}{l} \text{still somehow depends on the dimension. (of the model)} \\ \text{d} \end{array} \quad (4.163)$$

so we still have a \sqrt{d} term. This will not be a good bound for high-dimensional models. For general models, we will arrive at a bound of the form

$$R_S(\mathcal{F}) \leq \tilde{O} \left(\sqrt{\frac{\# \text{ of parameters}}{n}} \right). \quad (4.164)$$

This upper bound only depends on the number of parameters in our model, and does not take into the account the scale and norm of the parameters. Additionally, this doesn't work with kernel methods since the explicit parameterization is possibly infinite-dimensional, and therefore this upper bound becomes useless.

These limitations motivate the use of margin theory, which does take into account the norm of parameters and provides a theoretical basis for regularization techniques such as L_1 and L_2 regularization.