

---

---

---

---

---



# Generalization & Rademacher complexities

Thm Let  $g$  be a function satisfying  $|g(z_1, \dots, z_i, \dots, z_n) - g(z_1, \dots, z_i', \dots, z_n)| \leq c_i \quad \forall 1 \leq i \leq n$  one coordinate doesn't change  $f_n$  too much

For independent RVs  $Z_i$ 's,

$$P(g(Z^n) - \mathbb{E}g(Z^n) \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right) \quad \text{Generalization of Hoeffding bound}$$

Proof

Recap We call  $\{D_i := M_i - M_{i-1}\}_{i=1}^n$  a martingale difference sequence w.r.t.  $Z_i^n$  if  $\mathbb{E}[D_i | X^{i-1}] = 0$

Lemma Let  $D_i$  be a martingale difference sequence w.r.t.  $Z_i^n$  s.t.  $\exists \sigma_i^2$

$$\mathbb{E}[e^{\lambda D_i} | Z_i^{i-1}] \leq \exp\left(\frac{\sigma_i^2 \lambda^2}{2}\right) \quad \forall i \quad \dots (*)$$

Then,  $M_n - M_0 = \sum_{i=1}^n D_i$  is  $(\sum \sigma_i^2)$ -sub-Gaussian.

Define the Doob martingale  $M_i = \mathbb{E}[g(Z^n) | Z_i]$   $\begin{pmatrix} M_0 = \mathbb{E}g(Z^n) \\ M_n = g(Z^n) \end{pmatrix}$

So we'd like to bound  $P(M_n - M_0 \geq t)$ .

Note that  $|D_i| = |\mathbb{E}[g(Z^n) | Z_i] - \mathbb{E}[g(Z^n) | Z_i^{i-1}]|$

$$\leq \sup_{z_i^{i-1}} |\mathbb{E}_{z_i^n} g(z_i^{i-1}, z_i, z_i^n) - \mathbb{E}_{z_i^{i-1}} g(z_i^{i-1}, z_i, z_i^n)| \leq c_i \quad \dots (*)$$

So  $\mathbb{E}[e^{\lambda D_i} | Z_i^{i-1}] = \mathbb{E}[e^{\lambda(D_i - \mathbb{E}(D_i | Z_i^{i-1}))} | Z_i^{i-1}] \leq \exp\left(\frac{\lambda^2}{2} \cdot \frac{c_i^2}{4}\right)$

From previous lemma, and tail inequality for sub-G RVs, we have the result  $\square$ .

Again, our goal is to show an optimality guarantee for ERM

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum \ell(\theta; z_i)$$

Now, we will use bdd diff to show the following uniform concentration result:

$$\underline{\epsilon}_n := \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\theta; z_i) - \mathbb{E} \ell(\theta; Z), \quad \bar{\epsilon}_n := \sup_{\theta \in \Theta} \mathbb{E} \ell(\theta; Z) - \frac{1}{n} \sum_{i=1}^n \ell(\theta; z_i)$$

are small w.h.p.

Why is this useful?

$$\begin{aligned} \mathbb{E}l(\hat{\theta}_n; Z) &\leq \frac{1}{n} \sum l(\hat{\theta}_n; z_i) + \bar{\epsilon}_n && \text{by def of } \epsilon_n \\ &\leq \frac{1}{n} \sum l(\theta; z_i) + \bar{\epsilon}_n && \text{by def of } \hat{\theta}_n, \text{ for any arbitrary } \theta \in \Theta \\ &\leq \mathbb{E}l(\theta; Z) + \bar{\epsilon}_n + \epsilon_n && \text{by def of } \epsilon_n \end{aligned}$$

Taking inf over  $\theta$ , we get

$$\mathbb{E}l(\hat{\theta}_n; Z) \leq \inf_{\theta \in \Theta} \mathbb{E}l(\theta; Z) + \bar{\epsilon}_n + \epsilon_n$$

So if  $\epsilon_n$  is small, then  $\hat{\theta}_n$  is near-optimal.

We will focus on finite-sample results today. Traditionally, ML guarantees are finite sample since it allows quantifying dimension dependence.

This is useful for high-dim, large-scale models.

Assumption  $l(\theta; z) \in [0, M]$

We proceed in two parts to bound  $\epsilon_n$  &  $\bar{\epsilon}_n$ . As we'll see, the case for  $\bar{\epsilon}_n$  is symmetric, so we focus on  $\epsilon_n$  below.

Part I We show  $\epsilon_n$  is concentrated around its mean w.h.p.

Define  $g(z_1, \dots, z_n) := \sup_{\theta \in \Theta} \frac{1}{n} \sum l(\theta; z_i) - \mathbb{E}l(\theta; Z)$  so that  $g(Z_1^n) = \epsilon_n$ . We'll apply Hoeffding's lemma.

As a notational shorthand, we use  $\hat{P}_n(\cdot) = \frac{1}{n} \sum \mathbb{1}\{z_i \in \cdot\}$ , and write  $Ql(\theta; z) = \mathbb{E}z \cdot l(\theta; z)$ .

$$\begin{aligned} &|g(z_1, \dots, z_n) - g(z_1', \dots, z_n')| \\ &= \left| \sup_{\theta \in \Theta} \left\{ \frac{1}{n} \sum l(\theta; z_i) - \mathbb{E}l(\theta; Z) \right\} - \sup_{\theta \in \Theta} \left\{ \frac{1}{n} \sum l(\theta; z_i') - \mathbb{E}l(\theta; Z) - \frac{1}{n} l(\theta; z_i) + \frac{1}{n} l(\theta; z_i') \right\} \right| \leq \frac{2M}{n} \end{aligned}$$

From Hoeffding's lemma,  $\mathbb{P}(\epsilon_n - \mathbb{E}\epsilon_n \geq t) \leq \exp\left(-\frac{nt^2}{2M^2}\right)$ . Equivalently,  $\epsilon_n \leq \mathbb{E}\epsilon_n + M\sqrt{\frac{2t}{n}}$  w.p.  $\geq 1 - e^{-t}$

Part II We bound  $\mathbb{E}\epsilon_n$  via symmetrization. cf. see VWV Ch. 2.2-3, 2.14 for more.

↳ This is tricky to bound. Think about how you would approach this.

Let  $z_1', \dots, z_n'$  be indep copies of  $z_1, \dots, z_n$ .

$$\mathbb{E}\epsilon_n = \mathbb{E}\left[\sup_{\theta \in \Theta} \frac{1}{n} \sum l(\theta; z_i) - \mathbb{E}\left[\frac{1}{n} \sum l(\theta; z_i') \mid z_1^n\right]\right] \leq \mathbb{E}\sup_{\theta \in \Theta} \frac{1}{n} \sum (l(\theta; z_i) - l(\theta; z_i'))$$

Let  $\sigma_i$  be i.i.d. random signs (Rademacher RVs), indep of everything else.

From  $\sigma_i (l(\theta_i; z_i) - l(\theta_i; z_i')) \stackrel{D}{=} l(\theta_i; z_i) - l(\theta_i; z_i')$ ,

$$\begin{aligned} \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n} \sum (\sigma_i l(\theta_i; z_i) - l(\theta_i; z_i')) &= \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n} \sum \sigma_i (l(\theta_i; z_i) - l(\theta_i; z_i')) \\ &\leq \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n} \sum \sigma_i l(\theta_i; z_i) + \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n} \sum (-\sigma_i) l(\theta_i; z_i) \\ &= 2 \mathbb{E} \sup_{\theta \in \Theta} \frac{1}{n} \sum \sigma_i l(\theta_i; z_i) \end{aligned}$$

Def The (empirical) Rademacher complexity of a class  $\mathcal{H}$  of functions  $h: \mathcal{Z} \rightarrow \mathbb{R}$  is

$$\mathcal{R}_n \mathcal{H} := \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(z_i) \mid \mathcal{Z}^n \right]$$

↳ Interpretation: How well can  $\mathcal{H}$  fit random noise  $\sigma_i$ 's? (where  $\sigma_i h(z_i)$  is the margin)

Note that  $\mathcal{R}_n \mathcal{H} = \mathcal{R}_n (-\mathcal{H})$ . So the case for  $\bar{\mathcal{E}}_n$  is symmetric.

Collecting bounds in Parts I & II, we arrive at

$$\underline{\mathcal{E}}_n \leq 2 \mathbb{E} \mathcal{R}_n \mathcal{H} + M \sqrt{\frac{\Phi}{2n}}, \quad \bar{\mathcal{E}}_n \leq 2 \mathbb{E} \mathcal{R}_n \mathcal{H} + M \sqrt{\frac{2\Phi}{n}} \quad \text{w.p. } \geq 1 - 2e^{-t}.$$

So we conclude

$$\mathbb{E} l(\hat{\theta}_n; \mathcal{Z}) \leq \inf_{\theta \in \Theta} \mathbb{E} l(\theta; \mathcal{Z}) + 4 \mathbb{E} \mathcal{R}_n \mathcal{H} + 2M \sqrt{\frac{2\Phi}{n}} \quad \text{w.p. } \geq 1 - 2e^{-t} //$$

Basic properties of Rademacher complexity:

1) Contraction Principle: Let  $\phi$  be a  $C_\phi$ -Lipschitz function with  $\phi(0) = 0$ ,

$$\mathcal{R}_n \phi \circ \mathcal{H} \leq 2C_\phi \mathcal{R}_n \mathcal{H}$$

↳ This will be useful for HW2.

↳ think LP, sup obtained at vertices.

2)  $\mathcal{R}_n(\text{convex-hull}(\mathcal{H})) = \mathcal{R}_n(\mathcal{H})$  for finite  $\mathcal{H}$

3) Consider any finite  $\mathcal{H}$ .

$$\text{Then, } \mathcal{R}_n \mathcal{H} \leq \sqrt{\frac{2 \log |\mathcal{H}|}{n}} \cdot \sqrt{\sup_{h \in \mathcal{H}} \frac{1}{n} \sum h(z_i)^2}$$

↳ You'll show this in HW2.



Now, we analyze the Rademacher complexity of regularized linear models.

Example

$$l(\theta; X, Y) = (1 - Y\theta^T X)_+ = \phi(Y\theta^T X), \quad \Theta = \{\theta \in \mathbb{R}^d : \|\theta\|_p \leq r\}$$

$$\mathbb{E} \mathcal{R}_n \{(X, Y) \mapsto l(\theta; X, Y) : \theta \in \Theta\} = \mathbb{E} \mathcal{R}_n \{(X, Y) \mapsto \phi(Y\theta^T X) - \phi'(0) : \theta \in \Theta\} \\ \leq \mathbb{E} \mathcal{R}_n \{(X, Y) \mapsto Y \cdot \theta^T X : \theta \in \Theta\} \quad \text{by contraction principle}$$

Define  $Z := Y \cdot X$ . Then,  $\nearrow = \mathbb{E} \mathcal{R}_n \{Z \mapsto \theta^T Z : \theta \in \Theta\}$ .

We now derive scale-sensitive bounds on this quantity.

Theorem  $\mathcal{H} := \{Z \mapsto \theta^T Z : \|\theta\|_2 \leq r\}$  If  $\mathbb{E} \|Z\|_2^2 \leq C_2^2$ , then  $\mathbb{E} \mathcal{R}_n \mathcal{H} \leq \frac{C_2}{\sqrt{n}} r$

Pf)

$$\mathbb{E} \mathcal{R}_n \mathcal{H} = \frac{1}{n} \mathbb{E} \sup_{\|\theta\|_2 \leq r} \theta^T \left( \sum_i \sigma_i z_i \right) \leq \frac{r}{n} \mathbb{E} \left\| \sum \sigma_i z_i \right\|_2 \quad \text{by Cauchy-Schwarz} \\ \leq \frac{r}{n} \sqrt{\mathbb{E} \left\| \sum \sigma_i z_i \right\|_2^2} \quad \text{by Jensen's inequality}$$

Write out  $\left\| \sum \sigma_i z_i \right\|_2^2$  and note that cross terms have mean zero.

$$= \frac{r}{n} \sqrt{\mathbb{E} \sum \|\sigma_i z_i\|_2^2} = \frac{r}{n} \sqrt{\mathbb{E} \sum \|z_i\|_2^2} \leq \frac{r}{\sqrt{n}} \cdot C_2. \quad \square$$

What if you are interested in high-dimensional features, but think the model is sparse?

Theorem  $\mathcal{H} := \{Z \mapsto \theta^T Z : \|\theta\|_1 \leq s\}$  If  $\|Z\|_\infty \leq C_\infty$  a.s., then  $\mathbb{E} \mathcal{R}_n \mathcal{H} \leq \frac{C_\infty}{\sqrt{n}} s \cdot \sqrt{2 \log(2d)}$ .

↳ You'll show this in HW1.

\*  $\log d$  vs.  $d$  when  $s \ll d$  then  $L_1$ -regularization is nice.

These theorems say "so long as you regularize properly, your model complexity doesn't grow with problem dimension  $d$ "

Of course, all of these results compare performance against best-in-model-class. They don't say anything for whether that model class is good.

# Chaining & Dudley's entropy integral

We now give more sophisticated bounds on the Rademacher complexity. The bounds we develop play a key role in empirical process theory

e.g. uniform CLT

$$\sqrt{n} \left( \frac{1}{n} \sum h(z_i) - \mathbb{E}Z \right) \Rightarrow G(h) \quad \text{where } G \text{ is a Gaussian process indexed by } h \in \mathcal{H}$$

## Covering

We begin with notions of packing & covering numbers.

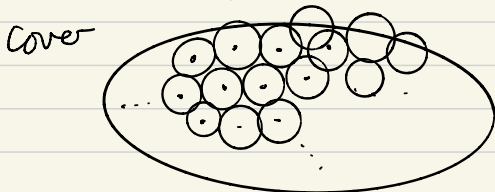
Consider a metric space  $(T, d)$  <sup>any nonempty set</sup> <sub>metric on  $T$</sub>

Def For any  $\varepsilon > 0$ ,  $\{h_i\}_{i=1}^N$  is a  $\varepsilon$ -cover of  $T$  if  $\forall h \in T \exists 1 \leq i \leq N$  s.t.  $d(h, h_i) \leq \varepsilon$ .

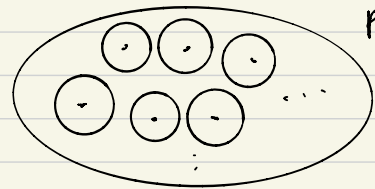
Def The  $\varepsilon$ -covering number of  $T$  is the size of the smallest  $\varepsilon$ -cover of  $T$

$$N(T, d, \varepsilon) := \inf \{ N \geq 0 : \exists \varepsilon\text{-cover } \{h_i\}_{i=1}^N \text{ of } T \}$$

We call  $\log N(T, d, \varepsilon)$  the metric entropy.



Cover



Packing

Def For any  $\delta > 0$ ,  $\{h_i\}_{i=1}^M \subseteq T$  is a  $\delta$ -packing of  $T$  if  $d(h_i, h_j) > \delta \quad \forall i \neq j$ .

The  $\delta$ -packing number of  $T$  is the size of the largest  $\delta$ -packing

$$M(T, d, \delta) := \sup \{ M \geq 0 : \exists \delta\text{-packing } \{h_i\}_{i=1}^M \text{ of } T \}$$

$$M(T, d, 2\delta) \stackrel{\textcircled{1}}{\leq} N(T, d, \delta) \stackrel{\textcircled{2}}{\leq} M(T, d, \delta)$$

Let  $\{h_i\}_{i=1}^M$  be the maximal  $\delta$ -packing. Then,  $\forall h \in T, d(h, h_i) \leq \delta \quad \forall i=1, \dots, M$ . So this is a  $\delta$ -cover of  $T$ .

Def Suppose there exists  $2\delta$ -packing  $\{h_1, \dots, h_M\}$  and  $\delta$ -cover  $\{h_1, \dots, h_N\}$ , with  $M \geq N+1$ . Then,  $\exists 1 \leq i < j \leq M$ , and  $1 \leq k \leq N$  s.t.  $d(h_i, h_j) \leq 2\delta$ ,  $d(h_j, h_k) \leq \delta$ . So  $d(h_i, h_k) \leq 2\delta \times \square$ .

## Lemma

Consider  $\|\cdot\|, \|\cdot\|'$  on  $\mathbb{R}^d$ . Let  $B, B'$  be corresponding unit balls. Then,

$$\left( \frac{1}{\delta} \right)^d \frac{\text{Vol}(B)}{\text{Vol}(B')} \stackrel{\textcircled{1}}{\leq} N(B, \|\cdot\|, \delta) \stackrel{\textcircled{2}}{\leq} \frac{\text{Vol}(\frac{2}{\delta}B + B')}{\text{Vol}(B')}$$

Pf) 1: Let  $\{h_j\}_{j=1}^N$  be a  $\delta$ -cover (in  $\|\cdot\|$ ) of  $B$ , so  $B \subseteq \bigcup_{j=1}^N \{h_j + \delta B\}$ . This implies  $\text{Vol}(B) \leq N \text{Vol}(\delta B) = N \delta^d \text{Vol}(B')$ .

2: Let  $\{h_i\}_{i=1}^M$  be a maximal  $\frac{\delta}{2}$ -packing of  $B$  (in  $\|\cdot\|$ ). By def of packing,  $\{h_j + \frac{\delta}{2}B'\}_{j=1}^M$  are disjoint and contained in  $B + \frac{\delta}{2}B'$ .

$$\text{Vol} \left( \bigcup_{j=1}^M \left\{ h_j + \frac{\delta}{2}B' \right\} \right) = M \text{Vol} \left( \frac{\delta}{2}B' \right) = M \cdot \left( \frac{\delta}{2} \right)^d \text{Vol}(B') \leq \text{Vol} \left( B + \frac{\delta}{2}B' \right) = \left( \frac{\delta}{2} \right)^d \text{Vol} \left( \frac{2}{\delta}B + B' \right) \quad \square$$

Example Consider  $\mathcal{H} = \{l(\theta; \cdot) : \theta \in \Theta\}$ . Let  $\|h\|_{L^2(\mathbb{R}^d)} := \sqrt{\frac{1}{n} \sum h(z_i)^2}$ .  
 Assume  $|l(\theta; z) - l(\theta'; z)| \leq L(z) \|\theta - \theta'\|$  for some norm  $\|\cdot\|$  on  $\mathbb{R}^d$ .

Then, any  $\varepsilon$ -cover of  $\Theta$  induces a  $\|L\|_\infty \varepsilon$ -cover on  $\mathcal{H}$  in  $\|L\|_{L^2(\mathbb{R}^d)}$   
 (Let  $\{\theta_j\}_{j=1}^M$  be a  $\varepsilon$ -cover. Then, consider  $\{l(\theta_j; \cdot)\}_{j=1}^M$  is a  $\|L\|_{L^2(\mathbb{R}^d)} \varepsilon$ -cover of  $\mathcal{H}$ .  
 $\forall \theta \in \Theta$ , let  $j$  be st.  $\|\theta - \theta_j\| \leq \varepsilon$ . Take  $\|l(\theta; \cdot) - l(\theta_j; \cdot)\|_{L^2(\mathbb{R}^d)} \leq \|L\|_{L^2(\mathbb{R}^d)} \|\theta - \theta_j\| \leq \|L\|_{L^2(\mathbb{R}^d)} \varepsilon$

So we conclude  $N(\mathcal{H}, \|L\|_{L^2(\mathbb{R}^d)}, \varepsilon \|L\|_{L^2(\mathbb{R}^d)}) \leq N(\Theta, \|\cdot\|, \varepsilon)$ .

SubG processes Instead of the (empirical) Rademacher complexity, we consider more general processes.

Def A collection of zero mean RVs  $\{V_h : h \in T\}$  is a sub-Gaussian process w.r.t.  $d$  if  
 $\mathbb{E} e^{\lambda(V_h - V_{h'})} \leq \exp\left(\frac{\lambda^2}{2} d(h, h')^2\right) \quad \forall h, h' \in T, \forall \lambda \in \mathbb{R}$ .

↳ tail of  $V_h - V_{h'}$  is  $d(h, h')^2$ -subG.

Example (Rademacher process) Consider  $R_{n,h} := \frac{1}{\sqrt{n}} \sum \sigma_i h(z_i)$  where  $\sigma_i$  : i.i.d. random signs,  $h \in \mathcal{H}$ .  
Conditional on  $Z^n$ ,  $h \mapsto R_{n,h}$  is a subGaussian process w.r.t.  $\|\cdot\|_\infty$  on  $\mathcal{H}$ .

Pf)

$R_{n,h} - R_{n,h'} = \frac{1}{\sqrt{n}} \sum \sigma_i (h - h')(z_i)$ . Recalling that  $\sigma_i$ 's are 1-sub-Gaussian,

$$\begin{aligned} \mathbb{E} \left[ \exp(\lambda(R_{n,h} - R_{n,h'})) \mid Z^n \right] &= \prod_{i=1}^n \mathbb{E} \left[ \exp\left(\frac{\lambda \sigma_i}{\sqrt{n}} (h - h')(z_i)\right) \mid z_i \right] \leq \prod_{i=1}^n \exp\left(\frac{\lambda^2}{2n} (h - h')(z_i)^2\right) \\ &= \exp\left(\frac{\lambda^2}{2} \frac{1}{n} \sum_{i=1}^n (h(z_i) - h'(z_i))^2\right) \\ &= \exp\left(\frac{\lambda^2}{2} \|h - h'\|_{L^2(\mathbb{R}^d)}^2\right) \quad \square \end{aligned}$$

So to bound  $B_n \mathcal{H} = \frac{1}{\sqrt{n}} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{\sqrt{n}} \sum \sigma_i h(z_i) \mid Z^n \right] = \frac{1}{\sqrt{n}} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} R_{n,h} \mid Z^n \right]$ ,  
 we can bound suprema of sub-Gaussian processes.

Key Lemma Let  $X_{ij}$  be  $\sigma_i^2$ -sub-G RVs,  $j=1, \dots, N$ . Then,  $\mathbb{E} \max_{1 \leq j \leq N} X_{ij} \leq \max_{1 \leq j \leq N} \sigma_j \cdot 2\sqrt{\log N}$ ,  $N \geq 2$ .

Proposition Let  $\{V_h : h \in T\}$  be a sub-Gaussian process w.r.t. a metric  $d$  on  $T$ . Let  $D := \sup_{h, h' \in T} d(h, h')$  diam(T)

Then, for any  $\delta > 0$ ,  $\mathbb{E} \sup_{h \in T} V_h \leq 2 \mathbb{E} \sup_{\substack{d(h, h') \leq \delta \\ h, h' \in T}} (V_h - V_{h'}) + 4D \sqrt{\log N(T, d, \delta)}$

Pf) Let  $N = N(T, d, \delta)$ , and  $\{h_j\}_{j=1}^N$  be a  $\delta$ -cover of  $T$ . Fix an arbitrary  $h \in T$ .  
 There exists  $j$  st.  $d(h, h_j) \leq \delta$ . Then,

$$V_h - V_{h'} = V_h - V_{h_j} + V_{h_j} - V_{h'} \leq \sup_{\substack{r, r' \in T \\ d(r, r') \leq \delta}} (V_r - V_{r'}) + \max_{1 \leq j \leq N} |V_{h_j} - V_{h'}|$$

Given another arbitrary  $\tilde{h} \in T$ , the same bound holds for  $V_{h'} - V_{\tilde{h}}$ .

Adding the two, and taking supremum over  $h, \tilde{h} \in T$

$$\sup_{h, \tilde{h} \in T} V_h - V_{\tilde{h}} \leq 2 \sup_{\substack{r, r' \in T \\ d(r, r') \leq \delta}} (V_r - V_{r'}) + 2 \max_{1 \leq j \leq N} |V_{h_j} - V_{h'}|$$

From Lemma,  $\mathbb{E} \max_{1 \leq j \leq N} |V_{h_j} - V_{h'}| \leq 2D \sqrt{\log N}$ . □

Example (A parameter class on  $[0,1]$ ) Define  $l(\theta; z) = 1 - e^{-\theta z}$ ,  $\theta \in [0,1]$ ,  $z \in [0,1]$ .  
 $\mathcal{H} = \{l(\theta; \cdot) : \theta \in [0,1]\} \subseteq \{h : [0,1] \rightarrow \mathbb{R}\}$ .

First term  $\mathbb{E} \sup_{\|h-h'\|_{L_2(\mathbb{P})} \leq \delta} R_{n,h} - R_{n,h'} = \mathbb{E} \sup_{\|h-h'\|_{L_2(\mathbb{P})} \leq \delta} \frac{1}{n} \sum \sigma_i (h(z_i) - h'(z_i)) \leq \sqrt{n} \cdot \delta$  by Cauchy-Schwarz.

Second term It's easy to check  $\theta \mapsto l(\theta; z)$  is 1-Lipschitz  $\forall z \in [0,1]$ . From above example,

$$N(\mathcal{H}, \|\cdot\|_{L_2(\mathbb{P})}, \delta) \leq N([0,1], |\cdot|, \delta) \leq \frac{1}{\delta} + 1, \quad \mathcal{V} = \sup_{\theta \in [0,1]} \frac{1}{n} \sum (1 - e^{-\theta z_i})^2 \leq 1$$

$$R_n \mathcal{H} = \mathbb{E} \left[ \sup_{\theta \in [0,1]} \frac{1}{n} \sum \sigma_i (1 - e^{-\theta z_i}) \mid Z_i^n \right] = \frac{1}{n} \mathbb{E} \sup_{h \in \mathcal{H}} R_{n,h}$$

$$\leq \frac{1}{\sqrt{n}} \cdot \left( 2\delta \sqrt{n} + 4 \sqrt{\log\left(\frac{1}{\delta} + 1\right)} \right) \text{ for any } \delta$$

$$= \frac{2}{\sqrt{n}} \inf_{\delta \in (0, \frac{1}{2})} \left( \sqrt{n} \delta + 2 \sqrt{\log\left(\frac{1}{\delta} + 1\right)} \right)$$

Setting  $\delta = \frac{1}{4\sqrt{n}}$ , we get  $R_n \mathcal{H} \leq \sqrt{\frac{\log n}{n}}$   $\square$ .

We now use a more refined argument that allows a tighter bound on the supremum.

Theorem (Radley's entropy integral) Let  $\{V_h : h \in \mathcal{T}\}$  be a sub-Gaussian process w.r.t.  $d$  on  $\mathcal{T}$ .

For any  $\delta \in [0, \mathcal{D}]$ ,

$$\mathbb{E} \sup_{h \in \mathcal{T}} V_h \leq \mathbb{E} \left[ \sup_{h, h' \in \mathcal{T}} V_h - V_{h'} \right] + 32 \int_{\delta}^{\mathcal{D}} \sqrt{\log N(\mathcal{T}, d, \varepsilon)} d\varepsilon.$$

Proof: Setting  $\delta = 0$ ,  $\mathbb{E} \sup_{h \in \mathcal{T}} V_h \leq 32 \int_0^{\mathcal{D}} \sqrt{\log N(\mathcal{T}, d, \varepsilon)} d\varepsilon$ ,  $N(\mathcal{T}, d, \delta) = 0 \forall \delta > \mathcal{D}$ .

PF)

We start with inequality from before:  $\sup_{h, h' \in \mathcal{T}} V_h - V_{h'} \leq 2 \sup_{r, r' \in \mathcal{T}} (V_r - V_{r'}) + 2 \max_{1 \leq j \leq n} |V_{h_j} - V_{h'_j}|$  *\* keep this written*

Instead of bounding the last term via Lemma, we use a chaining argument.

Recall that  $\mathcal{U} := \{h_j\}_{j=1}^n$  was a  $\delta$ -cover of  $\mathcal{T}$ .

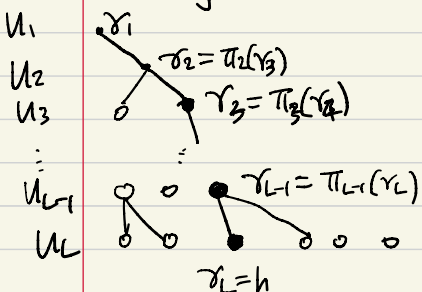
For each  $m$ , define  $\mathcal{U}_m :=$  minimal  $(\mathcal{D} \cdot 2^m)$ -cover of  $\mathcal{U}_m$  (we allow elements of  $\mathcal{T}$ ).

For  $L = \lceil \log_2 \mathcal{D}/\delta \rceil$ ,  $2^{-L} \leq \delta$ , so set  $\mathcal{U}_L = \mathcal{U}$ .

By def,  $|\mathcal{U}_m| \leq N(\mathcal{T}, d, \mathcal{D} \cdot 2^m)$ .

For each  $m$ , define  $\pi_m : \mathcal{U} \rightarrow \mathcal{U}_m$ ,  $\pi_m(h) = \operatorname{argmin}_{\tilde{h} \in \mathcal{U}_m} d(h, \tilde{h})$

Using this, we construct a chain from any  $h \in \mathcal{U}$ .  $\gamma_{m-1} = \pi_{m-1}(\gamma_m)$



$$V_h - V_{r_1} = \sum_{m=2}^L V_{r_m} - V_{r_{m-1}} \quad \text{and}$$

$$\mathbb{E} |V_h - V_{r_1}| \leq \sum_{m=2}^L \sup_{r \in \mathcal{U}_m} |V_r - V_{\pi_{m-1}(r)}| \quad //$$

Similarly, for any other  $\tilde{h} \in \mathcal{T}$ , we have same bound with  $\tilde{\gamma}_m$ 's.

We arrive at  $|V_n - V_n^*| = |V_{r_1} - V_{r_1^*} + V_n - V_{r_1} + V_{r_1^*} - V_n^*|$   
 $\leq |V_{r_1} - V_{r_1^*}| + \underbrace{|V_n - V_{r_1}| + |V_{r_1^*} - V_n^*|}_{\text{bound via chaining}}$   
 $\leq \max_{r_1, r_1^* \in U_1} |V_{r_1} - V_{r_1^*}| + 2 \sum_{m=2}^L \max_{r \in U_m} |V_r - V_{\pi_{m-1}(r)}|$

From Lemma,  $\mathbb{E} \max_{r_1, r_1^* \in U_1} |V_{r_1} - V_{r_1^*}| \leq 2D \sqrt{\log N(T, d, \frac{D}{2})}$ . and  
 since  $\max_{r \in U_m} d(r, \pi_{m-1}(r)) \leq D \cdot 2^{-(m-1)}$ , and  $|U_m| \leq N(T, d, D \cdot 2^{-m})$ , we have  
 $\mathbb{E} \max_{r \in U_m} |V_r - V_{\pi_{m-1}(r)}| \leq 2D 2^{-(m-1)} \sqrt{\log N(T, d, D \cdot 2^{-m})}$

Conclude that  $\mathbb{E} \sup_{h, h^* \in T} |V_h - V_h^*| \leq 4 \sum_{m=1}^L D \cdot 2^{-(m-1)} \sqrt{\log N(T, d, D \cdot 2^{-m})}$

Since  $s \mapsto \log N(T, d, s)$  is dec,  $D \cdot 2^{-m} \sqrt{\log N(T, d, D \cdot 2^{-m})} \leq 2 \int_{D \cdot 2^{-(m+1)}}^{D \cdot 2^{-m}} \sqrt{\log N(T, d, \epsilon)} d\epsilon$

$\Rightarrow 2 \mathbb{E} \sup_{h, h^* \in T} |V_h - V_h^*| \leq 32 \int_{\epsilon/4}^D \sqrt{\log N(T, d, \epsilon)} d\epsilon$

Combining with  $*$ , we get the result.  $\square$

Example

Recall that for  $l(\theta; z) = 1 - e^{-\theta z}$ ,  $\theta, z \in [0, 1]$ ,  $\mathcal{R}_n \mathcal{H} \leq \sqrt{\frac{\log n}{n}}$ .  
 Let's use Dudley's entropy integral.

$\mathcal{R}_n \mathcal{H} \leq \frac{32}{\sqrt{n}} \int_0^1 \sqrt{\log(1 + \frac{1}{\epsilon})} d\epsilon \leq \frac{32}{\sqrt{n}} \int_0^1 \sqrt{\log \frac{2}{\epsilon}} d\epsilon$ ,  $u = \sqrt{\log \frac{2}{\epsilon}} \Rightarrow \epsilon = 2e^{-u^2}$   
 $= \frac{32}{\sqrt{n}} \int_0^{\sqrt{\log 2}} 4u^2 e^{-u^2} du$ ,  $d\epsilon = -4u e^{-u^2} du$   
 $= \frac{C}{\sqrt{n}} \cdot (-u e^{-u^2} \Big|_0^{\sqrt{\log 2}} + \int_0^{\sqrt{\log 2}} e^{-u^2} du) = \frac{C}{\sqrt{n}}$   $\hookrightarrow$  No  $\log n$  factor!

Example Lipschitz functions  $|l(\theta; z) - l(\theta'; z)| \leq L \|\theta - \theta'\|$ ,  $\mathcal{H} = \{l(\theta; \cdot) : \theta \in \Theta\}$

Recall:  $N(\mathcal{H}, \|\cdot\|_{\mathcal{B}(\mathbb{R}^d)}, \epsilon \cdot L) \leq N(\Theta, \|\cdot\|, \epsilon)$ . If  $\Theta \subseteq r\mathbb{B}$ ,  $N(\Theta, \|\cdot\|, \epsilon) \leq (\frac{2r}{\epsilon})^d$

$\mathcal{R}_n \mathcal{H} \leq \frac{32}{\sqrt{n}} \int_0^{rL} \sqrt{\log N(\mathcal{H}, \|\cdot\|_{\mathcal{B}(\mathbb{R}^d)}, \epsilon)} d\epsilon \leq \frac{32L}{\sqrt{n}} \int_0^r \sqrt{\log N(\Theta, \|\cdot\|, \epsilon)} d\epsilon$   
 $\leq 32L \sqrt{\frac{d}{n}} \int_0^r \sqrt{\log(1 + \frac{2r}{\epsilon})} d\epsilon \leq L \cdot r \cdot \sqrt{\frac{d}{n}}$

Combining this with previous concentration result, for  $l(\theta; z) \in [0, M]$ , we have

$\mathbb{E} l(\hat{\theta}_n; z) \leq \inf_{\theta \in \Theta} \mathbb{E} l(\theta; z) + CLr \sqrt{\frac{d}{n}} + C \sqrt{\frac{L}{n}}$   $u.p. \geq 1 - e^{-t}$

Comment on measurability issues. Outer measures.

# ULLN

what if we just want to show  $\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum \ell(\theta; z_i) - \mathbb{E} \ell(\theta; z) \right| \xrightarrow{P} 0$  ?

## Theorem

Let  $H$  be an envelope function for  $\mathcal{H}$ :  $\forall h \in \mathcal{H}, |h| \leq H$ . Let  $\mathbb{E}|H(z)| < \infty$ , and define truncated version of  $\mathcal{H}$ :  $\mathcal{H}_M := \{ \underbrace{h \mathbb{1}\{|h| \leq M\}}_{=: h_M} : h \in \mathcal{H} \}$ .

If  $n \cdot \log N(\mathcal{H}_M, \|\cdot\|_{L_2(\mathcal{P}_n)}, \varepsilon) \xrightarrow{P} 0$  then  $\sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E} h(z) \right| \xrightarrow{P} 0$  for all fixed  $\varepsilon > 0, M < \infty$ .

Pf) From symmetrization,  $\mathbb{E} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E} h(z) \right| \leq 2 \mathbb{E} \mathcal{R}_n \mathcal{H}$   
 $\leq 2 \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum \tilde{O}_i (h(z_i) - h_M(z_i)) + 2 \mathbb{E} \mathcal{R}_n \mathcal{H}_M$   
 $\leq 2 \mathbb{E} H(z) \mathbb{1}\{H(z) > M\} + 2 \mathbb{E} \mathcal{R}_n \mathcal{H}_M$

Take a  $\varepsilon$ -cover  $\mathcal{H}_{M,\varepsilon}$  of  $\mathcal{H}_M$  in  $\|\cdot\|_{L_2(\mathcal{P}_n)}$ .  $\mathcal{R}_n \mathcal{H}_M \leq \mathcal{R}_n \mathcal{H}_{M,\varepsilon} + \varepsilon$

Now, note that since  $\sup_{h \in \mathcal{H}} \|h\|_{L_2(\mathcal{P}_n)} \leq M$ , Lemma gives

$$\sqrt{n} \mathcal{R}_n \mathcal{H}_{M,\varepsilon} \leq 2M \sqrt{\log N(\mathcal{H}_M, \|\cdot\|_{L_2(\mathcal{P}_n)}, \varepsilon)} \Rightarrow \mathcal{R}_n \mathcal{H}_{M,\varepsilon} \xrightarrow{P} 0$$

Same bound holds for  $\sup_{h \in \mathcal{H}} |\mathbb{E} h(z) - \frac{1}{n} \sum h(z_i)|$ .

$$\text{So } \mathbb{E} \sup_{h \in \mathcal{H}} \left| \frac{1}{n} \sum h(z_i) - \mathbb{E} h(z) \right| \leq 4 \mathbb{E} H(z) \mathbb{1}\{H(z) > M\} + 4 \mathbb{E} \mathcal{R}_n \mathcal{H}_{M,\varepsilon} + \varepsilon$$

Take  $n \rightarrow \infty$ , then let  $\varepsilon \downarrow 0, M \uparrow \infty$ . MCT gives the result.  $\square$

