

Generalization

We want to show that

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

achieves near-optimal population loss.

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:

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

are small w.h.p.

Why is this useful?

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

Taking inf over θ , we get

$$\mathbb{E} \ell(\hat{\theta}_n; Z) \leq \inf_{\theta \in \Theta} \mathbb{E} \ell(\theta; Z) + \bar{\Delta}_n + \Delta_n$$

So if ε_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.

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

Bounded differences

Bounded differences will play a key role in showing Δ_n is small.

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

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

Part 1 We show Δ_n is concentrated around its mean w.h.p.

Define $g(z_1, \dots, z_n) := \sup_{\theta \in \Theta} \frac{1}{n} \sum l(\theta; z_i) - \mathbb{E}l(\theta; z_i)$ so that $g(z^n) = \Delta_n$. We'll apply bounded diff.

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

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

From bounded differences, $\mathbb{P}(\Delta_n - \mathbb{E}\Delta_n \geq t) \leq \exp\left(-\frac{nt^2}{2M^2}\right)$. Equivalently, $\Delta_n \leq \mathbb{E}\Delta_n + M\sqrt{\frac{2t}{n}}$ w.p. $\geq 1 - e^{-t}$

So now, it suffices to control $\mathbb{E}\Delta_n$!

We begin with concentration results for light-tailed RVs.

Def A RV X is σ^2 -subGaussian if $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq \exp\left(\frac{\sigma^2 \lambda^2}{2}\right) \quad \forall \lambda \in \mathbb{R}$
 \hookrightarrow light-tailed RV \hookrightarrow MGF of $N(0, \sigma^2)$

From Markov's inequality, $\forall \lambda \geq 0$

$$\begin{aligned} \mathbb{P}(X - \mathbb{E}X \geq t) &= \mathbb{P}(\lambda(X - \mathbb{E}X) \geq \lambda t) = \mathbb{P}(e^{\lambda(X - \mathbb{E}X)} \geq e^{\lambda t}) \leq e^{-\lambda t} \mathbb{E} e^{\lambda(X - \mathbb{E}X)} \\ &\leq \exp\left(\frac{\sigma^2 \lambda^2}{2} - \lambda t\right) \end{aligned}$$

Taking min over $\lambda \geq 0$, we get $\mathbb{P}(X - \mathbb{E}X \geq t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)$.

Similarly, we have $\mathbb{P}(X - \mathbb{E}X \leq -t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)$.

Example ε : random signs (Rademacher) is 1-sub-Gaussian.

$$\begin{aligned} \mathbb{E} e^{\lambda \varepsilon} &= \frac{1}{2}(e^{-\lambda} + e^{\lambda}) = \frac{1}{2} \left(\sum_{k=0}^{\infty} \frac{(-\lambda)^k}{k!} + \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} \right) = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} \leq 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k}}{(2k)!} \\ &\leq 1 + \sum_{k=1}^{\infty} \frac{\lambda^{2k}}{2^k k!} = e^{\frac{\lambda^2}{2}} \end{aligned}$$

Example $X \in [a, b]$, $\mathbb{E}X = 0$. Then, X is $(b-a)^2$ -sub-Gaussian.

Pf) $\mathbb{E}_X e^{\lambda X} = \mathbb{E}_X e^{\lambda(X - \mathbb{E}X + X')}$ $\stackrel{\text{Jensen}}{\leq} \mathbb{E} e^{\lambda(X - X')}$ where X' indep copy of X

Let ε be random signs indep of everything so that $X - X' \stackrel{D}{=} \varepsilon(X - X')$

$$\begin{aligned} \mathbb{E} e^{\lambda(X - X')} &= \mathbb{E}_{X, X'} \mathbb{E}_{\varepsilon} e^{\lambda \varepsilon (X - X')} \leq \mathbb{E}_{X, X'} e^{\frac{\lambda^2}{2} (X - X')^2} \text{ by Ex 2} \\ &\leq e^{\frac{\lambda^2}{2} (b-a)^2} \end{aligned}$$

Actually, we can show X is $\frac{(b-a)^2}{4}$ -sub-G.

cf. $X \in [a, b]$ is $\frac{(b-a)^2}{4}$ -sub-Gaussian.

By convexity of $x \mapsto e^{\lambda x}$, $e^{\lambda x} \leq \frac{b-x}{b-a} e^{\lambda a} + \frac{x-a}{b-a} e^{\lambda b}$

Take expectation on both sides. For $h = \lambda(b-a)$, $p = \frac{-a}{b-a}$, $L(h) = -hp + \log(p + pe^h)$,

$$\mathbb{E} e^{\lambda X} \leq e^{L(h)} \quad L(0) = L'(0) = 0, \quad L''(h) \leq \frac{1}{4} \quad \forall h \quad \text{so } L(h) \leq \frac{1}{8} h^2 \text{ by Taylor } \square$$

We're ready to show bdd differences now.

Thm

Let g be a function satisfying

one coordinate doesn't change f_n too much

$$|g(z_1, \dots, z_i, \dots, z_n) - g(z_1, \dots, z'_i, \dots, z_n)| \leq c_i \quad \forall 1 \leq i \leq n$$

For independent RVs Z_i 's,

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

Generalization of Hoeffding bound

Def

$\{M_i\}_{i=0}^n$ is a martingale seq w.r.t. RVs Z_1, \dots, Z_n

if M_i is (Z_1, \dots, Z_i) -measurable, $\mathbb{E}|M_i| < \infty$, and

$$\mathbb{E}[M_i | Z_{i-1}] = M_{i-1}$$

We call $\{D_i := M_i - M_{i-1}\}_{i=1}^n$ a martingale difference sequence w.r.t. Z_i

$$(\mathbb{E}[D_i | Z_{i-1}] = 0)$$

Lemma

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

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

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

Pf

$$\begin{aligned} \mathbb{E} e^{\lambda \sum D_i} &= \mathbb{E} e^{\lambda D_n} \cdot e^{\lambda \sum_{i=1}^{n-1} D_i} = \mathbb{E} \left[\mathbb{E} \left[e^{\lambda D_n} \cdot \underbrace{e^{\lambda \sum_{i=1}^{n-1} D_i}}_{Z_{i-1}\text{-measurable}} \mid Z_{i-1} \right] \right] \\ &\leq \exp\left(\frac{c_n^2 \lambda^2}{2}\right) \cdot \mathbb{E} \left[e^{\lambda \sum_{i=1}^{n-1} D_i} \right] \end{aligned}$$

By induction, we get the result. \square

Proof of bdd differences

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 $P(M_n - M_0 \geq t)$.

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

$$\leq \sup_{z_{i-1}} |\mathbb{E}_{z_{i-1}} g(z_{i-1}, z, z_{i-1}) - \mathbb{E}_{z_{i-1}} g(z_{i-1}, z', z_{i-1})| \leq c_i \quad \dots (*)$$

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

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

Part II We bound $\mathbb{E} \Delta_n$ 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 .

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

Let ε_i be iid. random signs (Rademacher RVs), indep of everything else.

From $\varepsilon_i (\ell(\theta; z_i) - \ell(\theta; z_i')) \stackrel{D}{=} \ell(\theta; z_i) - \ell(\theta; z_i')$,

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

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

$$\mathcal{R}_n \mathcal{H} := \mathbb{E} \left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(z_i) \mid z_1^n \right]$$

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

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

Collecting bounds in Parts I & II, we arrive at

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

So we conclude

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

Basic properties of Rademacher complexity:

this will be useful for HW 7.

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

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

think LP, sp 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} .

You'll show this in HW 7.

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

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.

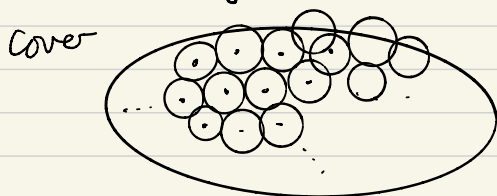
Consider a metric space (T, d) ^{any nonempty set} _{metric on T}

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

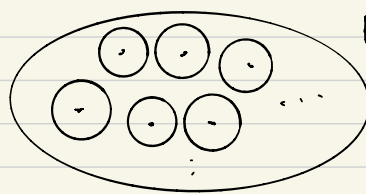
Def The ε -covering number of T is the size of the smallest ε -cover of T

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

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 δ -packing of T if $d(h_i, h_j) > \delta \quad \forall i \neq j$.

The δ -packing number of T is the size of the largest δ -packing

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

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

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

Suppose there exists 2δ -packing $\{h_1, \dots, h_M\}$ and δ -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. $d(h_i, h_k) \leq \delta, d(h_j, h_k) \leq \delta$. So $d(h_i, h_j) \leq 2\delta \times \square$.

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

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

Pf) ①: Let $\{h_j\}_{j=1}^N$ be a δ -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 $\|h\|_{L^2(\mathbb{P}_n)} := \sqrt{\frac{1}{n} \sum h(z_i)^2}$.
 Assume $|l(\theta; z) - l(\theta'; z)| \leq L(z) \|\theta - \theta'\|$ for some norm $\|\cdot\|$ on \mathbb{R}^d .

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

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

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

Example (Rademacher process) Consider $R_{n,h} := \frac{1}{\sqrt{n}} \sum \varepsilon_i h(z_i)$ where ε_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)
 $R_{n,h} - R_{n,h'} = \frac{1}{\sqrt{n}} \sum \varepsilon_i (h - h')(z_i)$. Recalling that ε_i 's are 1-sub-Gaussian,
 $\mathbb{E} \left[\exp(\lambda(R_{n,h} - R_{n,h'})) \mid Z^n \right] = \prod_{i=1}^n \mathbb{E} \left[\exp\left(\frac{\lambda \varepsilon_i}{\sqrt{n}} (h - h')(z_i)\right) \mid z_i \right] \leq \prod_{i=1}^n \exp\left(\frac{\lambda^2}{2n} (h - h')(z_i)^2\right)$
 $= \exp\left(\frac{\lambda^2}{2} \cdot \frac{1}{n} \sum_{i=1}^n (h(z_i) - h'(z_i))^2\right)$
 $= \exp\left(\frac{\lambda^2}{2} \|h - h'\|_{L^2(\mathbb{P}_n)}^2\right) \quad \square$

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

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

Proposition Let $\{V_h : h \in \mathcal{T}\}$ be a sub-Gaussian process w.r.t. a metric d on \mathcal{T} . Let $D := \sup_{h, h' \in \mathcal{T}} d(h, h')$ diam(\mathcal{T})
 Then, for any $\delta > 0$, $\mathbb{E} \sup_{h \in \mathcal{T}} V_h \leq 2 \mathbb{E} \sup_{\substack{d(h, h') \leq \delta \\ h, h' \in \mathcal{T}}} (V_h - V_{h'}) + 4D \sqrt{\log N(\mathcal{T}, d, \delta)}$

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

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

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

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

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

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

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

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

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

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

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

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 $\delta \in [0, D]$,

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

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

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 $U := \{h_j : j=1, \dots, n\}$ was a δ -cover of \mathcal{T} .

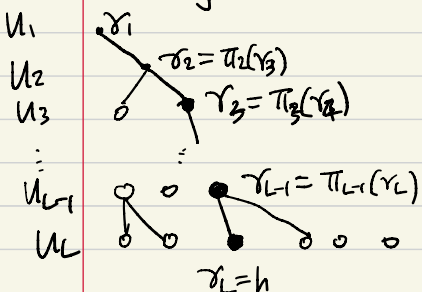
For each m , define $U_m :=$ minimal $(D \cdot 2^{-m})$ -cover of U_m (we allow elements of \mathcal{T}).

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

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

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

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



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

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

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.

$$\begin{aligned} \mathcal{R}_n \mathcal{H} &\leq \frac{32}{\sqrt{n}} \int_0^1 \sqrt{\log \left(1 + \frac{1}{\epsilon}\right)} d\epsilon \leq \frac{32}{\sqrt{n}} \int_0^1 \sqrt{\log \frac{2}{\epsilon}} d\epsilon, \quad u = \sqrt{\log \frac{2}{\epsilon}} \Rightarrow \epsilon = 2e^{-u^2} \\ &= \frac{32}{\sqrt{n}} \int_0^{\sqrt{\log 2}} 4u^2 e^{-u^2} du \quad d\epsilon = -4u e^{-u^2} du \\ &= \frac{C}{\sqrt{n}} \cdot \left(-u e^{-u^2} \Big|_0^{\sqrt{\log 2}} + \int_0^{\sqrt{\log 2}} e^{-u^2} du \right) = \frac{C}{\sqrt{n}} \quad \leftarrow \text{No } \log n \text{ factor!} \end{aligned}$$

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 $\Theta \subseteq r\mathbb{B}$, $N(\Theta, \|\cdot\|, \epsilon) \leq \left(1 + \frac{2r}{\epsilon}\right)^d$

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

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

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

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

Take a ε -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

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

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

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