

4 - Rademacher Complexity: Properties and Applications

Recall: z_i are iid data points, σ_i iid Rademacher r.v.'s

$$\neq \text{RC of } \mathcal{F}: \mathcal{R}(\mathcal{F}) = \mathbb{E} \left[\sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right]$$

$$\neq \text{ERC of } \mathcal{F}: \hat{\mathcal{R}}(\mathcal{F}) = \mathbb{E} \left[\sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \mid z_{1:n} \right]$$

- Some properties of RC:

- Monotonicity: $\mathcal{F}_1 \subseteq \mathcal{F}_2$, $\mathcal{R}(\mathcal{F}_1) \leq \mathcal{R}(\mathcal{F}_2)$

- Linear combination: $\mathcal{F}_1 + \mathcal{F}_2 = \{f_1 + f_2 : f_1 \in \mathcal{F}_1, \text{ and } f_2 \in \mathcal{F}_2\}$

$$\mathcal{R}(\mathcal{F}_1 + \mathcal{F}_2) = \mathcal{R}(\mathcal{F}_1) + \mathcal{R}(\mathcal{F}_2)$$

- Scaling w/ scalar: $\mathcal{R}(\mathcal{F} \cdot c) = |c| \cdot \mathcal{R}(\mathcal{F})$

- Convex-hull: For \mathcal{F} finite

$$\mathcal{R}(\text{convex-hull}(\mathcal{F})) = \mathcal{R}(\mathcal{F})$$

Recall:

→ **Theorem** (Generalization due to RC): Loss is bounded by \mathcal{B} , and $\mathcal{G} = \{g(z) = \ell(z, f) : f \in \mathcal{F}\}$, then w.p. $1 - \delta$,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f^*) \leq 4 \cdot \mathcal{R}(\mathcal{G}) + \mathcal{B} \sqrt{\frac{2 \log 2/\delta}{n}}$$

$$\text{Also, } \sup_{\mathcal{F}} |\hat{\mathcal{R}}(f) - \mathcal{R}(f)| \leq 2 \mathcal{R}(\mathcal{G}) + \mathcal{B} \sqrt{\frac{2 \log 2/\delta}{n}}$$

- **Theorem** (Talagrand's Contraction Principle): Let g be a L -Lipschitz cont. func. and $g \circ \mathcal{F} = \{g \circ f : f \in \mathcal{F}\}$, then

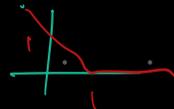
$$\mathcal{R}(g \circ \mathcal{F}) \leq L \cdot \mathcal{R}(\mathcal{F})$$

Goal: i) Relate $\mathcal{R}(g)$ to $\mathcal{R}(\mathcal{F})$.

ii) $\mathcal{R}(\mathcal{F})$ decays w/ n .

Ex (SVMs): - Data $z = (y, x)$ where $y \in \{-1, +1\}$, $f \in \mathcal{F}$
- Hinge loss: $\ell(z, f) = \max\{0, 1 - yf(x)\}$
 $\phi: s \rightarrow \max\{0, 1 - s\}$ is 1-Lip.

* RC of loss class



$$\mathcal{R}(g) \text{ where } g = \left\{ z = (y, x) \rightarrow \ell(z, f) \mid f \in \mathcal{F} \right\} \\ = \left\{ z \rightarrow \phi(y \cdot f(x)) \mid f \in \mathcal{F} \right\}$$

Define the intermediate fnc class

$$\mathcal{H} = \left\{ z = (y, x) \rightarrow y \cdot f(x) \mid f \in \mathcal{F} \right\} \\ \Rightarrow g = \phi \circ \mathcal{H} \text{ where } \phi \text{ is 1-Lip.}$$

$$\text{(by Talagrand)} \Rightarrow \mathcal{R}(g) \leq \mathcal{R}(\mathcal{H}) = \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \underbrace{\sigma_i y_i}_{\stackrel{d}{=} \sigma_i \perp x_i} f(x_i) \right]$$

Claim: $\sigma_i y_i \stackrel{d}{=} \sigma_i \perp x_i$

proof: For any A ,

$$\mathbb{P}(x_i \in A, \sigma_i y_i = 1)$$

$$= \mathbb{P}(x_i \in A, \sigma_i = 1, y_i = 1) + \mathbb{P}(x_i \in A, \sigma_i = -1, y_i = -1)$$

$$= \mathbb{P}(x_i \in A, \sigma_i = 1, y_i = 1) + \mathbb{P}(x_i \in A, \sigma_i = 1, y_i = -1)$$

$$= \mathbb{P}(x_i \in A, \sigma_i = 1) = \mathbb{P}(x_i \in A) \mathbb{P}(\sigma_i = 1)$$

$$= \mathbb{P}(x_i \in A) \mathbb{P}(\sigma_i y_i = 1)$$

$$= \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right]$$

$$= \mathcal{R}(\mathcal{F})$$

* If we have RC of \mathcal{F} ,

we can characterize the generalization error of SVMs.

Ex (Misclassification error):

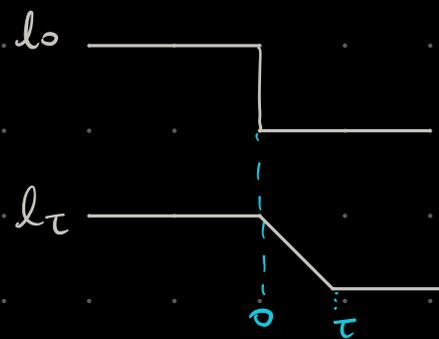
- Data $z = (y, x)$ where $y \in \{\pm 1\}$, $f \in \mathcal{F}$

- 0-1 loss: $l_0(z, f) = \mathbb{1}\{y \neq f(x)\} = \mathbb{1}\{y \cdot f(x) \leq 0\}$
 $f: \mathbb{R}^d \rightarrow \{\pm 1\}$

$$l_0(y \cdot f(x)) = \mathbb{1}\{y \cdot f(x) \leq 0\}$$

$$\text{where } l_0(s) = \begin{cases} 1 & s \leq 0 \\ 0 & s > 0 \end{cases}$$

- l_0 is not Lipschitz, we need surrogate losses (i.e. hinge)



$$l_\tau(s) = \begin{cases} 1 & s < -\tau \\ 1 - \frac{s}{\tau} & -\tau \leq s \leq \tau \\ 0 & s > \tau \end{cases}$$

$\hookrightarrow l_\tau$ is $\frac{1}{\tau}$ -Lip cont.

- What is the generalization error of $\hat{f}_\tau = \arg\min_{\mathcal{F}} \hat{R}_\tau(f)$?

Assump: $\mathbb{P}(0 \leq y \cdot f(x) \leq \tau) \leq C \cdot \tau$ for small τ .

$$\text{Ex: } y = \text{sign}(x) \quad x \sim \mathcal{N}(0, 1) \Rightarrow \mathbb{P}(0 \leq y \cdot \text{sign}(x) \leq \tau) = \sqrt{\frac{2}{\pi}} \cdot \tau + o(\tau)$$

- Define $f_\tau^* = \arg\min_{\mathcal{F}} R_\tau(f)$ (f_0^* minimizes $R_0(f)$)

- By Thm on RC, w.p. at least $1 - \delta$

$$(i) \quad R_\tau(\hat{f}_\tau) \leq R_\tau(f_\tau^*) + 4R(\mathcal{G}) + \sqrt{\frac{2 \log \frac{2}{\delta}}{n}}$$

where $\mathcal{G} = \{z \rightarrow l_\tau(z, f) \mid f \in \mathcal{F}\}$.

(ii) \Rightarrow By Talagrand's Concentration: $\mathcal{R}(f) \leq \frac{1}{\tau} \cdot \mathcal{R}(F)$

(iii) $\lambda_\tau \geq \lambda_0 \Rightarrow \mathcal{R}_\tau(f) \geq \mathcal{R}_0(f) \Rightarrow \mathcal{R}_\tau(\hat{f}_\tau) \geq \mathcal{R}_0(\hat{f}_\tau)$

(iv) $|\mathcal{R}_\tau(f) - \mathcal{R}_0(f)| \leq \mathbb{P}(0 \leq Y \cdot f(x) \leq \tau) \leq C \cdot \tau$

$$\Rightarrow \mathcal{R}_\tau(\hat{f}_\tau^*) \leq \mathcal{R}_\tau(f_0^*) \leq \mathcal{R}_0(f_0^*) + C \cdot \tau$$

Then, by (i-iv), w.p. $1-\delta$

$$\mathcal{R}_0(\hat{f}_\tau) - \mathcal{R}_0(f_0^*) \leq C \cdot \tau + \frac{4}{\tau} \mathcal{R}(F) + \sqrt{\frac{2 \log 2/\delta}{n}}$$

- Typically, $\mathcal{R}(F) = O\left(\frac{1}{n}\right)$. (this is next)

$$\mathcal{R}_0(\hat{f}_\tau) - \mathcal{R}_0(f_0^*) \lesssim \tau + \frac{1}{\tau n}$$

By optimizing over τ , $\lesssim \frac{1}{n^{1/4}}$ if $\tau = O\left(\frac{1}{n^{1/4}}\right)$

- We get generalization error $\leq O(n^{-1/4})$ by solving the relaxation.

- RC of Constrained Linear Models

- Goal: $\mathcal{R}(F) = O\left(\frac{1}{\sqrt{n}}\right)$

Theorem (RC of Linear Models): Let $F = \{f(x) = \langle x, \theta \rangle : \|\theta\| \leq r\}$

Then, i) $\hat{\mathcal{R}}(F) \leq \frac{r}{n} \sqrt{\sum_{i=1}^n \|x_i\|^2}$

ii) If $\mathbb{E}[\|x_i\|^2] \leq \kappa^2$, then $\mathcal{R}(F) \leq \frac{r \cdot \kappa}{\sqrt{n}}$

Remarks: 1- If we combine this bound w/ previous examples, we achieve generalization.

2 - Notice $\kappa = O(\sqrt{n})$ so $\mathcal{R}(\mathcal{F}) \leq r \sqrt{\frac{d}{n}}$.

$$\begin{aligned} \text{Proof: } \textit{i)} \quad \hat{\mathcal{R}}(\mathcal{F}) &= \mathbb{E} \left[\sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \mid x_{1:n} \right] \\ &= \mathbb{E} \left[\sup_{\|\theta\| \leq r} \frac{1}{n} \sum_{i=1}^n \sigma_i \langle x_i, \theta \rangle \mid x_{1:n} \right] \\ &= \mathbb{E} \left[\sup_{\|\theta\| \leq r} \left\langle \frac{1}{n} \sum_{i=1}^n \sigma_i x_i, \theta \right\rangle \mid x_{1:n} \right] \end{aligned}$$

$$\sup_{\|\theta\| \leq r} \langle \theta, v \rangle = r \cdot \|v\|$$

(by Jensen's Ineq.)

$$\begin{aligned} &= r \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \sigma_i x_i \right\| \mid x_{1:n} \right] \\ &\leq \frac{r}{n} \mathbb{E} \left[\left\| \sum_{i=1}^n \sigma_i x_i \right\|^2 \mid x_{1:n} \right]^{\frac{1}{2}} \\ &= \frac{r}{n} \mathbb{E} \left[\sum_{i=1}^n \|x_i\|^2 + \sum_{i \neq j} \sigma_i \sigma_j \langle x_i, x_j \rangle \mid x_{1:n} \right]^{\frac{1}{2}} \\ &= \frac{r}{n} \sqrt{\sum_{i=1}^n \|x_i\|^2} \end{aligned}$$

$$\textit{ii)} \quad \mathcal{R}(\mathcal{F}) = \mathbb{E} \hat{\mathcal{R}}(\mathcal{F}) \leq \frac{r}{n} \mathbb{E} \sqrt{\sum_{i=1}^n \|x_i\|^2}$$

(by Jensen's Ineq.)

$$\begin{aligned} &\leq \frac{r}{n} \sqrt{\mathbb{E} \sum_{i=1}^n \|x_i\|^2} \\ &\leq \frac{r}{n} \sqrt{\sum_{i=1}^n \mathbb{E} [\|x_i\|^2]} \\ &\leq \frac{r \cdot K}{\sqrt{n}} \end{aligned}$$

□

- RC of Finite Fnc Classes

Thm (Massart's Finite Lemma): Let z_1, \dots, z_n are iid and \mathcal{F} is finite. If $\sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(z_i)^2 \leq \nu^2$ a.s., then

$$\hat{\mathcal{R}}(\mathcal{F}) \leq \nu \sqrt{\frac{2 \log |\mathcal{F}|}{n}}$$

Proof: σ is a Rademacher r.v. w/ MGF $M_{\sigma}(t) = \mathbb{E} e^{\sigma t} = \frac{e^t + e^{-t}}{2} = \cosh(t)$.

$$\text{For } t \geq 0, \exp\{t \cdot \hat{\mathcal{R}}(\mathcal{F})\} = \exp\left\{t \cdot \mathbb{E} \left[\sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \mid z_{1:n} \right]\right\}$$

$z \rightarrow e^{tz}$ is convex. By Jensen's $\leq \mathbb{E} \left[\exp\left\{t \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i)\right\} \mid z_{1:n} \right]$

$z \rightarrow e^{tz}$ is monotone $= \mathbb{E} \left[\sup_{\mathcal{F}} \exp\left\{t \cdot \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i)\right\} \mid z_{1:n} \right]$

\mathcal{F} is finite: $\sup_{\mathcal{F}} \leq \sum_{\mathcal{F}}$ $\leq \mathbb{E} \left[\sum_{f \in \mathcal{F}} \exp\left\{\frac{t}{n} \sum_{i=1}^n \sigma_i f(z_i)\right\} \mid z_{1:n} \right]$

$$= \sum_{f \in \mathcal{F}} \prod_{i=1}^n M_{\sigma_i} \left(\frac{t}{n} f(z_i) \right) \mid z_{1:n}$$

$= \cosh$

$$\leq \sum_{f \in \mathcal{F}} \prod_{i=1}^n \exp\left\{\frac{t^2}{2n^2} f(z_i)^2\right\}$$

$$= \sum_{f \in \mathcal{F}} \exp\left\{\frac{t^2}{2n} \underbrace{\frac{1}{n} \sum_{i=1}^n f(z_i)^2}_{\leq \nu^2}\right\}$$

$$\exp\{t \cdot \hat{\mathcal{R}}(\mathcal{F})\} \leq |\mathcal{F}| \exp\left\{\frac{t^2}{2n} \nu^2\right\}$$

Exercise: $\frac{x^2}{2} \geq \log \cosh(x)$

$$\Rightarrow M_{\sigma}(t) \leq \exp\{t^2/2\}$$

Take log, divide by $t \Rightarrow \hat{R}(\mathcal{F}) \leq \frac{\log |\mathcal{F}|}{t} + \frac{tk^2}{2n}$

Optimize over t
 $(t = k^{-1} \sqrt{2 \log |\mathcal{F}| \cdot n}) \Rightarrow \hat{R}(\mathcal{F}) \leq k \cdot \sqrt{\frac{2 \log |\mathcal{F}|}{n}}$ □

Remarks: Recall that w.p. $1-\delta$,

$$R(\hat{f}) - R(f_*) \leq 4R(\mathcal{G}) + B \sqrt{\frac{2 \log 2/\delta}{n}}$$

where $\mathcal{G} = \{z \rightarrow \ell(z, f) \mid f \in \mathcal{F}\}$

ℓ is bdd by $B \Rightarrow \sup_{\mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(z_i, f)^2 \leq B^2 := k^2$

\Rightarrow By MFL

$$\hat{R}(\mathcal{G}) \leq B \sqrt{\frac{2 \log |\mathcal{F}|}{n}} \quad (\text{since } |\mathcal{G}| \leq |\mathcal{F}|)$$

\Rightarrow Generalization bound: w.p. $1-\delta$

- by MFL + RC: $4B \sqrt{\frac{2 \log |\mathcal{F}|}{n}} + B \sqrt{\frac{2 \log 2/\delta}{n}}$

- by union bound: $B \sqrt{\frac{2 \log |\mathcal{F}|}{n} + \frac{2 \log 1/\delta}{n}}$

The same order.