

Lecture 7

- * You can sign up for class presentations.
- * Scribe lecture notes for lectures 1-5 are online.
- * Last time
 - VC Theorem
 - Rademacher complexity

Today

Finish Rademacher complexity

Let $Z = X + Y$

F : function class $Z \rightarrow \mathbb{R}$

D : distribution over Z .

Definition (Rademacher Complexity).

Let F be a family of real-valued functions $f: Z \rightarrow \mathbb{R}$ where $Z = X + Y$. Then the Rademacher complexity $R(F)$ is defined as

$$R(F) = \frac{1}{n} \mathbb{E}_{\sigma \sim \pm 1^n} \left[\sup_{f \in F} \sum_{i=1}^n \sigma_i f(z_i) \right].$$

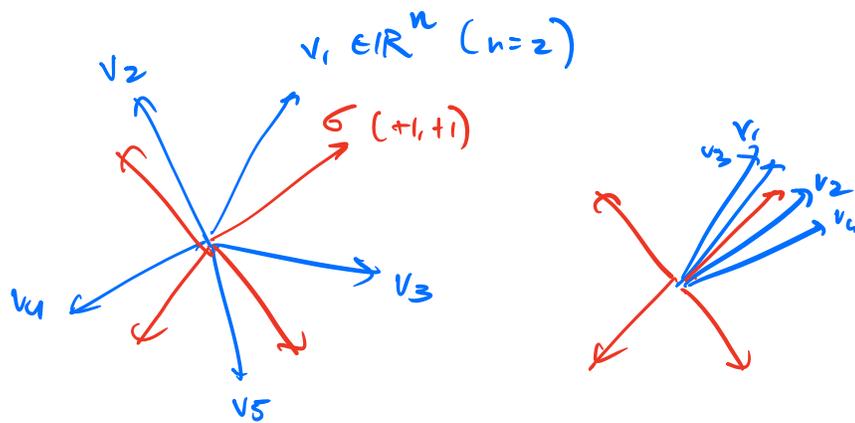
More generally, given a set of vectors $A \subset \mathbb{R}^n$, the Rademacher complexity $R(A)$ is defined as

$$R(A) = \frac{1}{n} \mathbb{E}_{\sigma \sim \pm 1^n} \sup_{a \in A} \sum_{i=1}^n \sigma_i a_i.$$

Intuition: $R(F)$ captures how well can F fit random noise.

If F can fit random noise,
 F will probably overfit.

Geometric picture



In expectation over $\sigma \sim \{\pm 1\}^n$, what is the max. inner product we can get with σ ?

How do we use Rademacher?

$$S = \{ (x_i, y_i), i \in [n] \}$$

\mathcal{H} : function from $\mathcal{X} \rightarrow \mathcal{Y}$.

$$\mathcal{H} \circ S = \{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \}$$

$\ell(h(x), y)$: instead of writing $\ell(h(x), y)$ we can write $\ell(h, z) = \ell(h(x), y)$ where $z = (x, y)$

$$\ell \circ \mathcal{H} \circ S = \{ (\ell(h, z_1), \dots, \ell(h, z_n)) : h \in \mathcal{H} \}$$

For e.g. $\mathcal{H} = \{ h_1, h_2, h_3 \}$

$$\ell \circ \mathcal{H} \circ S = \{ (\ell(h_1, z_1), \dots, \ell(h_1, z_n)), \\ (\ell(h_2, z_1), \dots, \ell(h_2, z_n)), \\ (\ell(h_3, z_1), \dots, \ell(h_3, z_n)) \}$$

Lemma (Symmetrization with Rademacher)

$$\mathbb{E}_{S \sim D^n} \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2 \mathbb{E}_{S \sim D^n} R(\ell \circ \mathcal{H} \circ S)$$

Proof:

$$\begin{aligned} & \mathbb{E}_{S \sim D^n} \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \\ & \leq \mathbb{E}_{S, S'} \sup_{h \in \mathcal{H}} \frac{1}{n} \left(\sum_{i=1}^n (\ell(h, z_i) - \ell(h, z_i')) \right) \\ & = \mathbb{E}_{S, S', \sigma_{i:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \left(\sum_{i=1}^n \sigma_i (\ell(h, z_i) - \ell(h, z_i')) \right) \\ & \leq \mathbb{E}_S \mathbb{E}_{\sigma_{i:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(h, z_i) \\ & \quad + \mathbb{E}_{S'} \mathbb{E}_{\sigma_{i:n}} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (-\sigma_i) \ell(h, z_i') \\ & \mathbb{E}_S \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2 \mathbb{E}_{S, \sigma_{i:n}} \sup_{h \in \mathcal{H}} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i \ell(h, z_i) \right) \\ & = 2 \mathbb{E}_S R(\ell \circ \mathcal{H} \circ S) \quad \blacksquare \end{aligned}$$

Theorem (Excess risk bounds using Rademacher)

Assume that for all z and $h \in \mathcal{H}$ we have that $|l(h, z)| \leq c$. Then with probability at least $(1-\delta)$ over $S \sim D^n$,

$$(1) \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2 \mathbb{E}_S [R(l \cdot \mathbb{H} \circ S)] + c \sqrt{\frac{2 \log(1/\delta)}{n}}$$

$$(2) \sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2 R(l \cdot \mathbb{H} \circ S) + 3c \sqrt{\frac{2 \log(2/\delta)}{n}}$$

(3) for any $h^* \in \mathcal{H}$,

$$R(\text{ERM}_{\mathcal{H}}(S)) - R(h^*) \leq 2 R(l \cdot \mathbb{H} \circ S) + 4c \sqrt{\frac{2 \log(4/\delta)}{n}}$$

(this holds for $h^* = \arg \min_{h \in \mathcal{H}} R(h)$)

Proof

We will keep using McDiarmid's.

Part (i)

Note that $\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h))$ satisfies the

bounded differences property with constant $\frac{2c}{n}$.

(changing any (x_i, y_i) changes the loss by at most $\frac{2c}{n}$).

\therefore Using McDiarmid's,

$$\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq \mathbb{E} \left(\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \right) + \varepsilon$$

$$\text{w.p. } 1 - \exp \left(- \frac{2\varepsilon^2}{n (2c/n)^2} \right)$$

$$= 1 - \exp \left(- \frac{n\varepsilon^2}{2c^2} \right)$$

$$\text{Choose } \varepsilon = c \sqrt{\frac{2 \log(1/\delta)}{n}}$$

We get, w.p. $(1 - \delta)$,

$$\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq \mathbb{E} \left(\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \right) + c \sqrt{\frac{2 \log(1/\delta)}{n}}.$$

Now use symmetrization lemma, and result follows. □

Part (2)

$$R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S}) = \mathbb{E}_{\sigma_{i:n}} \left(\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(h, z_i) \right)$$

this also satisfies bounded differences with constant $\frac{2c}{n}$ (swapping σ_i by σ_i' changes by $\leq \frac{2c}{n}$)

w.p. $1 - \delta$,

$$\therefore R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S}) \geq \mathbb{E}_{S'} (R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S}')) - c \sqrt{\frac{2 \log(1/\delta)}{n}}.$$

$$\Rightarrow \mathbb{E}_{S'} (R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S})) \leq R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S}) + c \sqrt{\frac{2 \log(1/\delta)}{n}}.$$

Now set $\delta = \delta/2$, w.p. $1 - \frac{\delta}{2}$,

$$\mathbb{E}_{S'} (R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S})) \leq R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S}) + c \sqrt{\frac{2 \log(2/\delta)}{n}}.$$

$$\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2 \mathbb{E}_{S'} (R(\mathcal{L} \circ \mathcal{H} \circ \mathcal{S})) + c \sqrt{\frac{2 \log(2/\delta)}{n}}$$

(By part (1)).

Result follows by doing a union bound.

w.p. $(1-\delta)$,

$$\sup_{h \in \mathcal{H}} (R(h) - \hat{R}_S(h)) \leq 2R(\log \frac{1}{\delta}) + 3c \sqrt{\frac{2 \log(2/\delta)}{n}}.$$

Part 3

Let $h_S = \text{ERM}_{\mathcal{H}}(S)$.

Part (2) bounds $R(\hat{h}) - \hat{R}(\hat{h}) \leq (\cdot)$
For (2) we need bound on
 $\hat{R}(h) - R(h) \leq (\cdot)$
 \therefore Part (2) doesn't work.

$$R(h_S) - R(h^*) = \underbrace{R(h_S) - \hat{R}_S(h_S)}_{(1)} + \underbrace{\hat{R}_S(h_S) - \hat{R}_S(h^*)}_{(2)} - R(h^*) \leq 0$$

(1) is bounded by part (2)

For (2) we use Hoeffding's (McDiarmid's with constant $\frac{2c}{n}$)

By using this, we get that w.p. $1 - \frac{\delta}{2}$,

$$\hat{R}_S(h^*) - R(h^*) \leq c \sqrt{\frac{2 \log(2/\delta)}{n}}.$$

Combining (1) & (2),

$$R(h_S) - R(h^*) \leq 2R(\log \frac{1}{\delta}) + 4c \sqrt{\frac{2 \log(4/\delta)}{n}}.$$

Take aways.

- Could be much better than VC bound: Rademacher complexity takes the data distribution into account, and could give tighter bounds than worst case VC bounds.
- Data-dependent bound. (3) in our Theorem is a data-dependent bound, we use a training set S both for learning a hypothesis from \mathcal{H} , and for estimating its generalization error (we can check if we are overfitting).

Rademacher calculus.

Claim (Translation & Scaling)

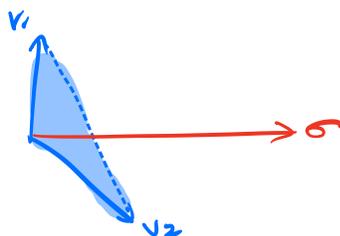
Let $A' = \{ \phi a + v, a \in A \}$

Then $R(A') = \phi R(A)$.

Proof [Exercise]

(can also show that

$$R(\text{convex hull of } A) = R(A).$$



Lemma (Massart Lemma) Let $A = \{v_1, \dots, v_m\}$ be a finite set of vectors in \mathbb{R}^n . Let $\bar{v} = \frac{1}{m} \sum_{i=1}^m v_i$.

Then

$$R(A) \leq \max_i \|v_i - \bar{v}\|_2 \cdot \sqrt{2 \log(m)}$$

Proof [Exercise]

Hint: First, by translation invariance, we can take $\bar{v} = 0$ without loss of generality.

Then use the max. of sub-gaussian result from last time. \square

Note: This gives a bound for finite hypothesis classes.

Lemma (Contraction Lemma)

For each $i \in [m]$, let $\phi_i: \mathbb{R} \rightarrow \mathbb{R}$ be a ρ -Lipschitz function i.e. $|\phi_i(x) - \phi_i(y)| \leq \rho |x - y| \forall x, y \in \mathbb{R}$.

For any $a \in \mathbb{R}^n$ define $\phi(a) \in \mathbb{R}^n$ as

$\phi(a) = (\phi_1(a_1), \dots, \phi_m(a_m))$. For a

set A , let $\phi \circ A = \{\phi(a) : a \in A\}$. Then

$$R(\phi \circ A) \leq \rho R(A).$$

Proof. Refer to book. •

Rademacher complexity of linear classes.

$$\mathcal{H}_1 = \{ hw(x) = \langle w, x \rangle : \|w\|_1 \leq B_1 \}$$

$$\mathcal{H}_2 = \{ hw(x) = \langle w, x \rangle : \|w\|_2 \leq B_2 \}$$

Lemma (l₂ bounded linear predictor)

Let $S = (x_1, \dots, x_n)$. Define $\mathcal{H}_2 \circ S = \{ \langle w, x_1 \rangle, \dots, \langle w, x_n \rangle : \|w\|_2 \leq B_2 \}$.

Then $R(\mathcal{H}_2 \circ S) \leq B_2 \frac{\max_i \|x_i\|_2}{\sqrt{n}}$

Proof By Cauchy-Schwarz: $\langle w, x \rangle \leq \|w\|_2 \|x\|_2$

$$\begin{aligned} \therefore n R(\mathcal{H}_2 \circ S) &= \mathbb{E}_\sigma \left[\sup_{a \in \mathcal{H}_2 \circ S} \sum_{i=1}^n \sigma_i a_i \right] \\ &= \mathbb{E}_\sigma \left[\sup_{w: \|w\|_2 \leq B_2} \sum_{i=1}^n \sigma_i \langle w, x_i \rangle \right] \\ &= \mathbb{E}_\sigma \left[\sup_{w: \|w\|_2 \leq B_2} \langle w, \sum_{i=1}^n \sigma_i x_i \rangle \right] \\ &\leq B_2 \cdot \mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i x_i \right\|_2 \right] \quad \text{--- (1)} \end{aligned}$$

Using Jensen's,

$$\mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i x_i \right\|_2 \right] = \mathbb{E}_\sigma \left[\left(\left\| \sum_{i=1}^n \sigma_i x_i \right\|_2^2 \right)^{1/2} \right]$$

$$\leq \left(\mathbb{E}_\sigma \left\| \sum_{i=1}^n \sigma_i t_i \right\|_2^2 \right)^{1/2} \quad - (2)$$

$$\mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i t_i \right\|_2^2 \right] = \mathbb{E}_\sigma \left[\sum_{i,j} \sigma_i \sigma_j \langle t_i, t_j \rangle \right]$$

Since σ_i are independent $\mathbb{E}_\sigma [\sigma_i \sigma_j] = 0$ ($\neq i, j$)

$$\begin{aligned} \mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i t_i \right\|_2^2 \right] &= \sum_{i=1}^n \|t_i\|_2^2 \\ &\leq n \max_i \|t_i\|_2^2 \quad - (3) \end{aligned}$$

Combining (1), (2) & (3) gives bound. \blacksquare

Lemma (L₁ bounded linear model)

Let $S = (x_1, \dots, x_n)$ where $t_i \in \mathbb{R}^d$ $\forall i \in [n]$.

Then $R(\mathcal{H}_1 \circ S) \leq B_1 \max_i \|x_i\|_\infty \sqrt{\frac{2 \log(2d)}{n}}$.

Proof By Holder's inequality $\langle w, v \rangle \leq \|w\|_1 \|v\|_\infty$

$$n R(\mathcal{H}_1 \circ S) = \mathbb{E}_\sigma \left[\sup_{a \in \mathcal{H}_1 \circ S} \sum_{i=1}^n \sigma_i a_i \right]$$

$$= \mathbb{E}_\sigma \left[\sup_{w: \|w\|_1 \leq B_1} \sum_{i=1}^n \sigma_i \langle w, x_i \rangle \right]$$

$$= \mathbb{E}_\sigma \left[\sup_{w: \|w\|_1 \leq B_1} \langle w, \sum_{i=1}^n \sigma_i t_i \rangle \right]$$

$$\leq \mathbb{E}_\sigma \left[B_1 \left\| \sum_{i=1}^n \sigma_i t_i \right\|_\infty \right]$$

$$= B, \mathbb{E}_G \left[\max_{j \in [d]} \left| \sum_{i=1}^n \sigma_i(x_i)_j \right| \right]$$

each term $\sigma_i(x_i)_j$ is
 $|x_i)_j|$ sub-Gaussian

since $|x_i)_j| \leq \max_i \|x_i\|_\infty$,

\therefore each term $\sigma_i(x_i)_j$ is $\max_i \|x_i\|_\infty$
 sub-Gaussian

the sum $\sum_{i=1}^n \sigma_i(x_i)_j$ is $\left(\sum_{i=1}^n (\max_i \|x_i\|_\infty)^2 \right)^{1/2}$
 sub-Gaussian $\leq \sqrt{n} \cdot \max_i \|x_i\|_\infty$

By bound for mat. of sub-Gaussian, including
 negations to take care of ± 1 ,

$$n \mathbb{R}(\mathcal{H}, \sigma) \leq B, \sqrt{n} \max_i \|x_i\|_\infty \sqrt{2 \log(2d)}$$

~~■~~