

(1)

## Rademacher Complexity

$\mathcal{H}$ : hypothesis set  $h: \mathcal{X} \rightarrow \mathcal{Y}$   $h \in \mathcal{H}$

We can have various notions of a loss function  $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{IR}$

Given  $\mathcal{H}$  we define the losses

$$\mathcal{G} := \{ g: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{IR} \mid g(x, y) = L(h(x), y) \}$$

We are trying to associate a complexity measure with  $\mathcal{G}$  that measures how rich  $\mathcal{G}$  is.

Idea: If  $\mathcal{G}$  is very expressive it can even fit random noise.

Def

The Rademacher complexity of  $\mathcal{G}$  w.r.t samples of size  $m$  is

$$R_m(\mathcal{G}) = \mathbb{E}_{S \sim \mathcal{G}^m} \left[ \sum_{i=1}^m \mathbb{E}_{x_i \sim \mathcal{X}} [ \dots ] \right]$$

## Mc DIARMID'S INEQUALITY

let  $V$  be a set and  $f: V^m \rightarrow \mathbb{R}$  be a f.s.t.

$\forall i \exists c_i > 0 \quad \forall x_1, \dots, x_m, x_i' \in V$

$$|f(x_1, x_2, \dots, x_i, \dots, x_m) - f(x_1, \dots, x_i', \dots, x_m)| \leq c_i$$

Let  $X_1, \dots, X_m$  be random variables taking values in  $V$ .

$$\text{Then } \forall \epsilon > 0 \quad \mathbb{P}(f(X_1, \dots, X_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon) \leq e^{-2\epsilon^2/mc^2}$$

$$\text{where } c^2 = \sum_i c_i^2.$$

$$\text{If } f(x_1, \dots, x_m) = \frac{1}{m} \sum x_i \text{ & } c = \frac{1}{m} \sum_{x_i \in \{0, 1\}} x_i$$

we get the Hoeffding inequality.

We consider  $G$  to be a family of  $f^{\text{ns}}: Z \rightarrow \mathbb{R}$ .

Let  $D$  be a distribution over  $Z$  &  $g \in G$ .  
We write  $L(g) = \mathbb{E}_{Z \sim D} [g(z)]$

Given a sample  $S$  of size  $m$  drawn according to  $D$  from  $Z$  we write  $L_S(g) = \frac{1}{m} \sum_{i=1}^m g(z_i)$ .

The basic goal: we want to show that  $L_S$  &  $L$  are not too far apart with high probability for all  $g \in G$ . Rademacher complexity will give us a quantitative handle on this.

let  $\vec{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_m)$  be a family of i.i.d. random variables s.t.  $\forall i \in \{1, \dots, m\} \sigma_i = +1$  with prob.  $\frac{1}{2}$  and  $-1$  with prob.  $\frac{1}{2}$ .

def

The empirical Rademacher complexity of  $G$  with respect to the sample  $S$  of size  $m$  is

$$R_S(G) = \mathbb{E}_{\vec{\sigma}} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z_i) \right]$$

If we write  $\vec{g}_S = (g(z_1), \dots, g(z_m))$  then we can write  $R_S(G) = \mathbb{E}_{\vec{\sigma}} \left[ \sup_{g \in G} \frac{\vec{\sigma} \cdot \vec{g}_S}{m} \right]$

If this is high it means that  $G$  can correlate well with the random variables  $\sigma_i$  which means  $G$  is very expressive.

(3)

Def

The Rademacher complexity of  $G$  wrt  $S$  of size  $m$  is

$$R_m(G) = \mathbb{E}_{S \sim \mathcal{D}^m} [R_S(G)]$$

Note  $R_m$  depends on the distribution.

Thm!

Let  $G$  be a family of functions  $Z \rightarrow [0, 1]$ . Let  $S$  be a sample of size  $m$  drawn iid according to  $\mathcal{D}$  from  $Z$ . Then  $\forall \delta > 0$  with probability at least  $1 - \delta$ , the following holds

$$\forall g \in G \quad L(g) \leq L_S(g) + 2R_m(G) + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{m}}\right)$$

Proof

Want a bound on  $L(g) - L_S(g)$ .

Define a new random variable  $\Phi(S)$  by

$$\Phi(S) = \sup_{g \in G} (L(g) - L_S(g))$$

Idea: give a bound on  $\mathbb{E}[\Phi]$  & use a concentration inequality to argue that  $\Phi$  is close to  $\mathbb{E}[\Phi]$  with high probability.

Write  $S = (z_1, \dots, z_m)$  and consider  $\Phi$  as a function of  $(z_1, \dots, z_m)$ .

Short digression: If  $X$  is any compact set &  $f, g: X \rightarrow [0, 1]$  are continuous then

$$|\sup_{x \in X} f(x) - \sup_{x \in X} g(x)| \leq \sup_{x \in X} |f(x) - g(x)|.$$

Suppose  $\sup f$  is attained at  $x_0$  &  $\sup f > \sup g$

$$\sup f(x) - \sup g(x) = f(y) - \sup g(x) \leq f(y) - g(y) \leq \sup(f - g)$$

if  $\sup g > \sup f$  flip the roles of  $f$  and  $g$ .

④

Now what happens if we change  $z_i$  to  $z_i'$ ?

$$|\bar{\Phi}(z_1, \dots, z_m) - \bar{\Phi}(z_1, \dots, z_i', \dots, z_m)|$$

$$= \left| \sup_{g \in \mathcal{G}} [L(g) - L_{S'}(g)] - \sup_{g \in \mathcal{G}} [L(g) - L_S(g)] \right|$$

$$\leq \sup_{g \in \mathcal{G}} |L(g) - L_S(g) - L(g) + L_{S'}(g)|$$

$$= \sup_{g \in \mathcal{G}} \frac{1}{m} |g(z_i) - g(z_i')| \leq \frac{1}{m}$$

Thus we can use Mc Diarmid's inequality with  $c = \frac{1}{m}$ . Thus  $\forall \epsilon > 0$

$$P(\bar{\Phi}(S) \geq E[\bar{\Phi}(S)] + \epsilon) \leq e^{-2m\epsilon^2}$$

$$\text{Take } \epsilon = \sqrt{\frac{\log 1/\delta}{2m}} \quad c^{-2m\epsilon^2} = \delta$$

So with probability at least  $1 - \delta$

$$\bar{\Phi}(S) \leq E[\bar{\Phi}(S)] + \epsilon = \sqrt{\frac{\log 1/\delta}{2m}}$$

Now we need an upper bound on  $E[\bar{\Phi}(S)]$ .

Suppose we have a second sample

$$S' = \{z_1', \dots, z_m'\} \text{ also drawn to } \mathcal{D}$$

$$E_{S \sim \mathcal{D}^m} [\bar{\Phi}(S)] = E_{S \sim \mathcal{D}^m} \left[ \sup_{g \in \mathcal{G}} (L(g) - L_S(g)) \right] \quad \text{def of } \bar{\Phi}$$

$$= E_S \left[ \sup_{g \in \mathcal{G}} E_{S' \sim \mathcal{D}^m} (L_{S'}(g) - L_S(g)) \right] \quad \begin{matrix} L(g) = \\ E_{S' \sim \mathcal{D}^m} [L_{S'}(g)] \end{matrix}$$

$$\leq E_{S, S' \sim \mathcal{D}^m} \left[ \sup_{g \in \mathcal{G}} (L_{S'}(g) - L_S(g)) \right]$$

$$= E_{S, S' \sim \mathcal{D}^m} \left[ \sup_{g \in \mathcal{G}} \left\{ \frac{1}{m} \sum_{i=1}^m g(z_i') - g(z_i) \right\} \right] \quad \begin{matrix} \text{def of} \\ L_S \end{matrix}$$

For any RV's  
 $\sup_{i:1} E[X_i]$   
 $\leq E[\sup_i X_i]$

Now consider the following sum

$$\sum_{i=1}^m \sigma_i [g(z'_i) - g(z_i)]$$

where  $\sigma_i = \pm 1$  with equal probabilities.

This is equivalent to taking 2 samples & swapping the values at position  $i$  or not according to a fair coin toss

$$\begin{array}{cccc} z_1 & z_2 & z_3 & \cdot \begin{array}{c} z_i \\ \boxed{z_i} \\ z'_i \end{array} \cdot \cdots z_m \\ z'_1 & z'_2 & z'_3 & \cdot \begin{array}{c} z'_i \\ z_i \\ z'_m \end{array} \end{array}$$

If we get 1 swap also keep.

So the sums will have terms flipped or not with prob  $\frac{1}{2}$ . This is what  $\sigma_i$ 's do.

This swapping process does not change the distribution of  $S$  &  $S'$  so.

$$\begin{aligned} & \mathbb{E}_{S, S' \sim \mathcal{D}^m} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m [g(z'_i) - g(z_i)] \right] \\ &= \mathbb{E}_{S, S', \sigma} \left[ \sup_{g \in G} \frac{1}{m} \left[ \sum_{i=1}^m \sigma_i [g(z'_i) - g(z_i)] \right] \right] \\ &\leq \mathbb{E}_{S, \sigma} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z'_i) \right] + \mathbb{E}_{S', \sigma} \left[ \sup_{g \in G} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z_i) \right] \\ &= 2 R_m(G). \end{aligned}$$

Thus  $\mathbb{E}[\bar{\Phi}(S)] \leq 2 R_m(G)$ .

$$\bar{\Phi}(S) \leq 2 R_m(G) + \left[ \frac{\log \frac{1}{\delta}}{2m} \right]^{\frac{1}{2}} \text{ with prob} \geq 1 - \delta$$

Cor Suppose we draw a sample of size  $m$  according to  $\mathcal{D}$ . Then for any  $\delta > 0$  with prob. at least  $1-\delta$  the following holds

$$\forall g \in \mathcal{G} \quad L(g) \leq L_S(g) + 2R_S(g) + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{m}}\right)$$

Proof Consider  $R_S(g) := \mathbb{E}_{\mathcal{D}} \left[ \sup_{\mathcal{G}} \frac{1}{m} \sum_{i=1}^m g(z_i) \right]$  as a function of  $(z_1, \dots, z_m) \in \mathcal{S}$ . Changing one of the  $z_i$  to  $z'_i$  will change  $R_S$  by  $\frac{1}{m}$  at most. From Mc Diarmid's inequality with  $c = \frac{1}{m}$  &  $\epsilon = \sqrt{\frac{\log \frac{1}{\delta}}{m}}$  we have with prob at least  $1-\delta$

$$R_S(g) \leq R_m(g) + \sqrt{\frac{\log \frac{1}{\delta}}{m}}$$

Now by the union bound the probability that this is violated or the inequality in the theorem is violated is  $2\delta$ . Thus they both hold with probability at least  $1-2\delta$ . Replacing  $\delta$  by  $\delta/2$  & absorbing constants in the  $O(\cdot)$  notation we get the result.

— x —

Now we will relate Rademacher complexity & VC dimension.

Let  $H$  be a class of functions  $X \rightarrow \{+1, -1\}$ . Let  $G$  be the associated class of 0-1 loss  $f^m$ .

Prop 1 Given a sample  $S = \{(x_1, y_1), \dots, (x_m, y_m)\} \subseteq X \times \{+1, -1\}$  let  $S'$  denote the set  $\{x_1, \dots, x_m\}$  with labels removed from  $S$ . Then

$$R_S(g) = \frac{1}{2} R_{S'}(H)$$

(7)

Proof Fix  $h \in H$  and let  $g$  be the corresponding loss function:  $X \times \{+1, -1\} \rightarrow \{0, 1\}$ .

We can write  $g(x, y) = \frac{1}{2} (1 - y h(x))$

$$\begin{aligned}
 R_S(H) &= \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(x_i, y_i) \right] \\
 &= \mathbb{E}_\sigma \left[ \sup_{h \in H} \frac{1}{2m} \sum_{i=1}^m \sigma_i (1 - y_i h(x_i)) \right] \\
 &= \mathbb{E}_\sigma \left[ \frac{1}{2m} \sum_{i=1}^m \sigma_i \right] + \mathbb{E}_\sigma \left[ \sup_{h \in H} \frac{1}{2m} \sum_{i=1}^m -\sigma_i y_i h(x_i) \right] \\
 &= \frac{1}{2} \mathbb{E}_\sigma \left[ \sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i y_i h(x_i) \right] \\
 &= \frac{1}{2} \mathbb{E}_\sigma \left[ \sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \\
 &= \frac{1}{2} R_{S'}(H).
 \end{aligned}$$

LEMMA Let  $A \subseteq \mathbb{R}^m$  be a finite set of vectors with  $\|\vec{a}\| \leq 1$  for all  $\vec{a} \in A$ . Then

$$\mathbb{E}_\sigma \left[ \max_{\vec{a} \in A} \sum_{i=1}^m \sigma_i a_i \right] \leq \sqrt{2 \log |A|}$$

where  $\sigma_i$  are i.i.d. RV in  $\{+1, -1\}$  uniform.

PROOF Homework  $\rightarrow$  MASSART'S LEMMA

Cor 1 Let  $H$  be a family of functions on  $X$  taking values in  $\{+1, -1\}$  with  $\text{VC dim} = d$ . Let  $S = \{x_1, \dots, x_m\} \subseteq X$ . Then

$$R_S(H) = O\left(\sqrt{\frac{\log(m/d)}{m/d}}\right)$$

(3)

Proof Write  $A = \left\{ \frac{1}{\sqrt{m}} (h(x_1), \dots, h(x_m)) \mid h \in H \right\}$

Then  $A$  is a family of vectors in  $\mathbb{R}^m$  each of length at most 1 &  $|A| \leq \left(\frac{cm}{d}\right)^d$  [SAUER'S LEMMA]

$$\begin{aligned} \text{Thus } R_s(H) &= \mathbb{E} \left[ \sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \\ &= \frac{1}{\sqrt{m}} \mathbb{E} \left[ \sup_{\vec{a} \in A} \sum_{i=1}^m \vec{\sigma}_i \cdot \vec{a} \right] \\ &\leq \frac{1}{\sqrt{m}} \sqrt{2 \log \left( \frac{cm}{d} \right)^d} \quad [\text{MASSART'S LEMMA}] \\ &= O \left( \sqrt{\frac{\log(m/d)}{m/d}} \right). \quad \blacksquare \end{aligned}$$

Putting all these together we get

Thm 2 Let  $H$  be a family of functions on a domain  $X$ , taking values in  $\{+1, -1\}$  with VC dimension  $d$ . Let  $S$  be a sample of size  $m$  drawn iid from a fixed distribution  $D$ . Let  $\text{err}(h) = \mathbb{P}_{x \sim D}[h(x) \neq y]$  &  $\hat{\text{err}}(h)$  be the sample error. Then

$$\text{err}(h) \leq \hat{\text{err}}(h) + O\left(\sqrt{\frac{\log m/d}{m/d}}\right) + O\left(\sqrt{\frac{\log 18}{m}}\right)$$

If we use  $g$  as the 0-1 loss  $f''$  of  $h$  then  $\text{err}(h) = L(g)$  &  $\hat{\text{err}}(h) = L_S(g)$  & the result follows from Cor 1, Prop 1 & Thm 1.