

# Learning Theory: Rademacher Complexity and Infinite Hypotheses.

Friday, March 24, 2017 8:31 AM

- ▷ Last time
- Defined PAC learning model
  - Defined notion of learnability
  - Showed Learnability for finite hypotheses via uniform convergence.

- ▷ This time:
- Infinite hypotheses
  - Learnability via bounded Rademacher Complexity.

- ▷ We will still be looking at the ERM algorithm over a hypothesis space  $H$ : given dataset  $S$
- $$h_S = \operatorname{argmin}_{h \in H} \underbrace{L_S(h)}_{\text{loss on samples}} \equiv \frac{1}{m} \sum_{t=1}^m \ell(h, z_t)$$

- ▷ We want to measure performance w.r.t. true loss:  $E_{z \sim D} [\ell(h, z)] = L_D(h)$

- ▷ Our goal is to show if  $m$  is large enough
- (1) In expectation:

$$E [L_n(h)] \leq \inf L_n(h) + \epsilon$$

$$E_S [L_D(h_S)] \leq \inf_{h \in H} L_D(h) + \epsilon$$

(2) with high probability  $1 - \delta$  over  $S$

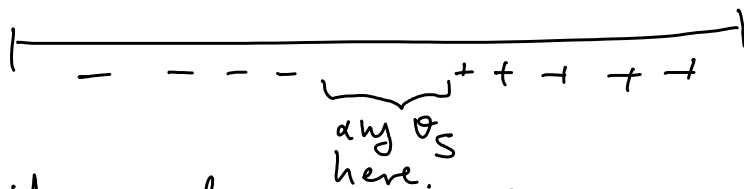
$$L_D(h_S) \leq \inf_{h \in H} L_D(h) + \epsilon.$$

▷ Last time we show an example on 1-dimension  $\theta = 1/2$



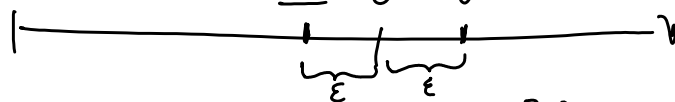
If  $H$  was all functions then we are in trouble.

- What if we restrict  $H$  to threshold functions. Then ERM picks the best threshold on sample.



- So the only misclassifications happen in the region between true  $\theta$  and  $\theta_S$ .

- But if you consider any region around  $\theta$



$$\Pr(x \in [\underline{\theta}, \bar{\theta}]) = \epsilon = P(x \in [\theta, \bar{\theta}])$$

then  $\theta_S$  has loss more than  $\epsilon$  if it either falls below  $\underline{\theta}$  or above  $\bar{\theta}$ . But it falls below  $\underline{\theta}$  if there are no  $x$ 's drawn in the region  $[\underline{\theta}, \theta]$

Each of them fails in there w.p.  $\epsilon$ .  
So  $\{\text{Prob of no } x \text{ in } [\underline{\theta}, \bar{\theta}]\} = (1-\epsilon)^m \leq e^{-\epsilon m}$

So if  $m \geq \frac{\log(2/\epsilon)}{\epsilon} \Rightarrow \{\text{above}\} \leq \frac{\delta}{2}$

Similarly  $\theta_S \in [\underline{\theta}, \bar{\theta}]$  w.p. at most  $\frac{\delta}{2}$ .

- So sample complexity:  $m(\epsilon, \delta) \leq \frac{\log(2/\epsilon)}{\epsilon} \cdot \frac{1}{\delta}$

▷ We will present a general theory of when we can make such claims for infinite hypotheses!

▷ The reason for non-learnability in the unconstrained  $H$  was that there were (the data memorization one) hypotheses that had a zero loss on sample but high loss on distribution, i.e.

$\sup_h (L_D(h) - L_S(h))$  was large!

▷ Last time we saw that if we can bound  $\sup_{h \in H} |L_D(h) - L_S(h)| \leq \epsilon$  i.e. uniform convergence, then we get learning.

▷ Now we will see that the absolute value is not really needed, i.e. we only need to have that the loss on the distribution cannot be much larger than the loss on sample, uniformly over  $H$ .  
(the other direction is not necessary, to be guaranteed uniformly)

## Learnability via Representativeness

▷ The representativeness of a sample  $S$  is defined as:

$$\text{Rep}(S, D) = \sup_{h \in H} (L_D(h) - L_S(h))$$

Lemma Representativeness implies good expected loss

$$\mathbb{E}_S [L_D(h_S)] \leq L_D(h^*) + \mathbb{E}_S [\text{Rep}(S, D)]$$

Pf

$$\mathbb{E} [L_D(h_S) - L_S(h_S)] \leq \mathbb{E}_S \left[ \sup_h (L_D(h) - L_S(h)) \right]$$

$$\mathbb{E}_S [L_D(h_S) - L_S(h_S)] \leq \mathbb{E}_S [L_S^{\sup_h (L_D(h) - L_S(h))}] \\ = \mathbb{E}_S [\text{Rep}(S, D)]$$

$$\mathbb{E}_S [L_S(h^*)] = L_D(h^*) \quad (\text{because } h^* \text{ is ind. of } S)$$

$$\Rightarrow \mathbb{E}_S [L_D(h_S)] \leq \mathbb{E}_S [L_S(h_S)] + \mathbb{E}_S [\text{Rep}(S, D)] \\ \left( \begin{array}{l} \text{since} \\ h_S \text{ max} \\ L_S \end{array} \right) \leq \mathbb{E}_S [L_S(h^*)] + \mathbb{E}_S [\text{Rep}(S, D)] \\ = L_D(h^*) + \mathbb{E}_S [\text{Rep}(S, D)] \quad \square$$

It also implies an easy h.p. bound:

- By Markov's ineq.

$$P(L_D(h_S) - L_D(h^*) > \frac{\mathbb{E}[L_D(h_S) - L_D(h^*)]}{\delta}) \leq \delta$$

$\Rightarrow$  w.p.  $1 - \delta$

$$L_D(h_S) \leq L_D(h^*) + \frac{\mathbb{E}[\text{Rep}(S, D)]}{\delta}$$

- Using McDiarmid's Inequality we can actually get: w.p.  $1 - \delta$

$$L_D(h_S) \leq L_D(h^*) + \mathbb{E}[\text{Rep}(S, D)] + 5\sqrt{\frac{\log(2/\delta)}{m}}$$

(Look at Chapter 26 of "Understanding Machine Learning")

## Bounding Representativeness

- Intuitively what if I wanted to bound

representativeness by using sample data.

• Split into train set  $S_1$ , validation set  $S_2$

• Use: 
$$\sup_h (L_{S_1}(h) - L_{S_2}(h))$$

as a proxy for representativeness.

• Split sample representativeness is exactly equal to: how much larger than the training error, can my test error be when I run ERM on train?

• Now if  $|S_1| = |S_2| = \frac{m}{2}$  then:

$$\sup_h (L_{S_1}(h) - L_{S_2}(h)) = \frac{2}{m} \sup_h \sum_{t=1}^m \epsilon_t \ell(h, z_t)$$

$$\text{where } \epsilon_t = \begin{cases} 1 & \text{if } t \in S_1 \\ -1 & \text{if } t \in S_2 \end{cases}$$

• Then the Rademacher complexity is equal to the Expected train-test split error gap, over a random split, where each data point is in  $S_1$  or  $S_2$  w.p.  $\frac{1}{2}$ .

$$\mathcal{R}(H, S) = \frac{2}{m} \mathbb{E}_{\vec{\epsilon}} \left[ \sup_{h \in H} \sum_{t=1}^m \epsilon_t \ell(h, z_t) \right]$$

- It captures: If I randomly split my data and run ERM on one set how good is the test error.

Thm

$$\mathbb{E}_S [\text{Rep}(D, S)] \leq \mathbb{E}_S [R(H, S)]$$

Pf

- For every fixed  $h$ : true loss of  $h$ ,  $L_D(h)$ , can be seen as expected loss of a newly sampled dataset  $S'$ , i.e.  $L_D(h) = \mathbb{E}_{S'} [L_{S'}(h)]$

$$\begin{aligned} \bullet \text{ So: } \text{Rep}(S, D) &= \sup_{h \in H} (L_D(h) - L_S(h)) \\ &= \sup_{h \in H} (\mathbb{E}_{S'} [L_{S'}(h)] - L_S(h)) \\ &\leq \mathbb{E}_{S'} \left[ \sup_{h \in H} (L_{S'}(h) - L_S(h)) \right] \end{aligned}$$

$$\bullet \text{ So: } \mathbb{E} [\text{Rep}(S, D)] \leq \mathbb{E} \left[ \sup_{h \in H} (L_{S'}(h) - L_S(h)) \right]$$

$$E_S [R_e(S, D)] \leq E_{S, S'} \left[ \underbrace{\sup_{h \in H} (L_{S'}(h) - L_S(h))}_{\text{split sample}} \right]$$

representativeness  
on a sample of  
size  $2m$ .

- Since:  $S, S'$  contain iid data I can swap any data point between them:

$$E_{S, S'} \left[ \sup_{h \in H} \frac{1}{m} \sum_{t=1}^T (\ell(h, z'_t) - \ell(h, z_t)) \right]$$

||

$$E_{S, S'} \left[ \sup_{h \in H} \frac{1}{m} \sum_t \epsilon_t (\ell(h, z'_t) - \ell(h, z_t)) \right]$$

For any  $\epsilon_t \in \{-1, 1\}$  (since this just corresponds to swapping  $z_t$  and  $z'_t$  in the samples. But  $z_t, z'_t$  are iid.)

|| if we draw  $\epsilon_t = \begin{cases} 1 & \text{w.p. } 1/2 \\ -1 & \text{w.p. } 1/2 \end{cases}$

$$E_{S, S', \epsilon} \left[ \sup_h \frac{1}{m} \sum_t \epsilon_t (\ell(h, z'_t) - \ell(h, z_t)) \right]$$

||

$$E_{S, S', \epsilon} \left[ \sup_h \frac{1}{m} \sum_t \epsilon_t \ell(h, z'_t) + \sup_h \frac{1}{m} \sum_t (-\epsilon_t) \ell(h, z_t) \right]$$



$$\begin{aligned}
& E_{S, S', \vec{c}} \left[ \sup_h \frac{1}{m} \sum_t c_t \ell(h, z'_t) + \sup_h \frac{1}{m} \sum_t (-c_t) \ell(h, z_t) \right] \\
& \quad \parallel \\
& E_{S', \vec{c}} \left[ \sup_h \frac{1}{m} \sum_t c_t \ell(h, z'_t) \right] + E_{S, \vec{c}} \left[ \sup_h \frac{1}{m} \sum_t (-c_t) \ell(h, z_t) \right] \\
& \quad \parallel \text{ since } c_t, -c_t \text{ are identically} \\
& \quad \text{dist and } S, S' \text{ too} \\
& 2 E_{S, \vec{c}} \left[ \sup_h \frac{1}{m} \sum_t c_t \ell(h, z_t) \right] \\
& \quad \parallel \\
& E_S [ R(H, S) ]
\end{aligned}$$

□

## Bounding the Rademacher complexity

▷ We want to bound for any data set of size  $m$ :

$$\frac{2}{m} \mathbb{E}_{\vec{c}} \left[ \sup_{h \in H} \sum_{t=1}^m c_t \ell(h, z_t) \right]$$

▷ Nice property of Rademacher:

- for any fixed  $h$ , variables:

$$y_t = g_t l(h, z_t) \quad \text{are mean zero.}$$

- So  $\frac{1}{m} \sum_{t=1}^m g_t l(h, z_t)$  is the sum of

mean zero random vars.

- By Hoeffding's Inequality we expect this to be of order  $\sqrt{\frac{\log 1/\delta}{m}}$  w.p.  $1-\delta$

- So then by a union bound we expect the  $\sup_h \frac{1}{m} \sum_t g_t l(h, z_t)$  to be of order

$$\sqrt{\frac{\log \#H}{m}} \quad \text{w.p. } 1-\delta. \Rightarrow \text{Expected value} = O\left(\sqrt{\frac{\log \#H}{m}}\right)$$

- This works for finite  $\#H$ . What about infinite?

- Well we don't really need to take union bound over all  $\#H$ . Just over all possible distinct functions  $H_S$  that can occur on a sample  $S$  of size  $m$ .

Formally:

- Let  $H_S$  be all possible distinct hypotheses

on  $S$ , i.e.  $\forall h \in H, \exists h' \in H_S$  s.t.

$$\forall z_t \in S: l(h', z_t) = l(h, z_t)$$

i.e. I cannot distinguish  $h, h'$  by just looking at their outputs on  $S$ .

Example:

$S_0$



If  $H$  is the space of threshold functions then how many distinct hypotheses can I create on  $m$  samples?

Answer: At most  $(m+1)!$

So then by the above reasoning I should expect:

$$R(H, S) = O\left(\sqrt{\frac{\log(m+1)}{m}}\right)$$

Even if  $H$  is unbounded!

Growth Rate | The growth rate of a hypothesis space  $H$  is the largest size of  $H_S$  for any data set of size  $m$ , as a function of  $m$ :

$$\tau_H(m) = \sup_{S: |S| \leq m} |H_S|$$

#  $S: |S| \leq m$

Lemma For any  $H$ :

$$R(H, S) \leq 2 \sqrt{\frac{2 \log(\tau_H(m))}{m}}$$

Pf

$$R(H, S) = 2 E_{\vec{G}} \left[ \sup_{h \in H} \frac{1}{m} \sum_{t=1}^m G_t \ell(h, z_t) \right]$$

$$= 2 E_{\vec{G}} \left[ \sup_{h \in H_S} \frac{1}{m} \sum_{t=1}^m G_t \ell(h, z_t) \right]$$

(We use the following in-expectation analogue of the Hoeffding Analysis)  
Massart's Lemma Let  $H = \{h_1, \dots, h_N\}$  be a finite hypothesis space. Then

$$E_{\vec{G}} \left[ \max_{h \in H} \frac{1}{m} \sum_{t=1}^m G_t \ell(h, z_t) \right] \leq \sqrt{\frac{2 \log N}{m}}$$

Pf

$$E_{\vec{G}} \left[ \max_{h \in H} \frac{1}{m} \sum_{t=1}^m G_t \ell(h, z_t) \right] =$$

$$\frac{1}{\lambda} E_{\vec{G}} \left[ \max_{h \in H} \log \left( \exp \left( \frac{\lambda}{m} \sum_{t=1}^m G_t \ell(h, z_t) \right) \right) \right]$$

11

.. 6 -  $h \in H$

$$\frac{1}{2} E_{\vec{\sigma}} \left[ \log \max_{h \in H} \exp \left( \frac{1}{m} \sum_t \sigma_t \ell(h, z_t) \right) \right]$$

$\wedge$  Jensen's

$$\frac{1}{2} \log E_{\vec{\sigma}} \left[ \max_{h \in H} \exp \left( \frac{1}{m} \sum_t \sigma_t \ell(h, z_t) \right) \right]$$

$\wedge$

$$\frac{1}{2} \log E_{\vec{\sigma}} \left[ \sum_{h \in H} \exp \left( \frac{1}{m} \sum_t \sigma_t \ell(h, z_t) \right) \right]$$

$\parallel$

$$\frac{1}{2} \log \sum_{h \in H} E_{\vec{\sigma}} \left[ \prod_t \exp \left( \frac{1}{m} \sigma_t \ell(h, z_t) \right) \right]$$

$\parallel$

$$\frac{1}{2} \log \sum_{h \in H} \prod_t E_{\sigma_t} \left[ \exp \left( \frac{1}{m} \sigma_t \ell(h, z_t) \right) \right]$$

$\wedge$

$$\frac{1}{2} \log \sum \prod^m \exp \left( - \frac{1}{2} \right)$$

$$\frac{1}{2} \log \sum_{u \in H} \prod_{t=1}^u \exp\left(-\frac{1^t}{2m^2}\right)$$

$$\frac{1}{2} \log \sum_{u \in H} \exp\left(\frac{1^2}{2m}\right)$$

$$\frac{1}{2} \log N \exp\left(\frac{1^2}{2m}\right)$$

$$\frac{1}{2} \log N + \frac{1}{2m} = \frac{1}{\sqrt{\frac{2m}{\log N}}} \sqrt{\frac{\log N}{2m}}$$

$$\leq \sqrt{\frac{2 \log N}{m}}$$

