

(1)

Learning Bound based on VC dimension

$$\text{Recall } \text{II}_{\mathcal{H}}(m) \leq \frac{ed}{\sigma} \left(\frac{em}{d}\right)^d$$

Thm

Let  $\mathcal{H}$  be a hypothesis set with VC dimension  $d$ . Let  $L$  be a consistent learner for  $\mathcal{H}$  that outputs  $h \in \mathcal{H}$ . Then  $L$  is a PAC-learning algorithm for  $\mathcal{H}$  provided it is given as input a sample of size  $m$  drawn from  $X$  with distribution  $D$

$$\text{where } m \geq K_0 \left[ \frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{d}{\epsilon} \log \frac{1}{\delta} \right]$$

for some universal constant  $K_0$ .

Proof

Let  $h, h'$  be two hypotheses from  $\mathcal{H}$ . We define

$$(h \Delta h')(x) = \begin{cases} 1 & \text{if } h(x) \neq h'(x) \\ 0 & \text{if } h(x) = h'(x) \end{cases}$$

If  $h$  is the "real" hypothesis we use  $h'$

$$\text{then } \text{err}(h') = \mathbb{P}_{x \sim D} [(h \Delta h')(x) = 1]$$

Define  $\Delta_h(\mathcal{H}) = \{ \cancel{h \Delta h} \mid h' \in \mathcal{H} \}$

$$\text{Claim } \text{VC}(\mathcal{H}) = \text{VC}(\Delta_h(\mathcal{H}))$$

Fix any finite  $S \subseteq X$ . Let  $k \in \text{II}_{\mathcal{H}}(S)$

Map  $k$  to  $\hat{k} := \hat{k} \Delta h \in \text{II}_{\Delta_h(\mathcal{H})}(S)$

This map is a bijection (routine case check)

so  $|\text{II}_{\mathcal{H}}(S)| = |\text{II}_{\Delta_h(\mathcal{H})}(S)|$  for any  $S$ .

Thus if  $\mathcal{H}$  can shatter  $S$  so can  $\Delta_h(\mathcal{H})$  & vice versa. Thus we have  $\text{VC}(\mathcal{H}) = \text{VC}(\Delta_h(\mathcal{H}))$ .

(2)

Now we define the bad hypotheses wrt  $\epsilon > 0$

$$\Delta_{h, \epsilon}(H) = \{\hat{h} \in \Delta_h(H) \mid \mathbb{P}_{x \sim \mathcal{D}} [\hat{h}(x) = 1] \geq \epsilon\}$$

If we pick a hypothesis  $h'$  &  $h' \in \Delta_{h, \epsilon}(H)$   
then  $h'$  will give an unacceptable error.

We say  $S$  is an  $\epsilon$ -net for  $\Delta_h(H)$  if  
 $\forall h \in \Delta_{h, \epsilon}(H) \exists x \in S \text{ s.t. } \hat{h}(x) = 1$ .

Now we are in the realizable case so  
it is always true that some hypothesis  
exists to give an empirical error (loss,  
risk) of 0 & we have a consistent learning  
algorithm so we find this hypothesis. Thus  
our algorithm would never choose a  
concept in  $\Delta_{h, \epsilon}(H)$  if  $S$  is an  $\epsilon$ -net.

→ So we need to bound the probability  
that  $S$  is not an  $\epsilon$ -net.

Draw the sample in 2 stages of size  $m$   
will view the first sample as training  
& the second as testing.

$S_1$ : first sample of size  $m$ .

Let  $A$  be the event:  $S_1$  is not an  $\epsilon$ -net.

If  $A$  happens then ~~exists s.t.~~  $\exists \tilde{h} \in \Delta_{h, \epsilon}(H)$   
s.t.  $\forall x \in S_1 \tilde{h}(x) = 0$ .

Fix  $\tilde{h}$  & draw a second sample  $S_2$  of size  $m$ .

Now we will aim to get a lower  
bound on the number of elements of  $S_2$   
s.t.  $\tilde{h}(x) = 1$ .

(3)

Let  $X_i$  be the random variable

$$X_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ element of } S_2 \text{ satisfies } \tilde{h}(x)=1 \\ 0 & \text{otherwise} \end{cases}$$

$$X = \sum_{i=1}^m X_i$$

$$\underset{x \sim \mathcal{D}}{\mathbb{P}} [X_i = 1] \geq \epsilon \text{ since } \tilde{h} \in \Delta_{h, \epsilon}(\mathcal{D}).$$

$$\text{So } \mathbb{E}[X] \geq m\epsilon$$

$$\mathbb{P}[X < \epsilon^{m/2}] \leq \mathbb{P}[|X - \mathbb{E}[X]| > \frac{\mathbb{E}[X]}{2}]$$

So we can use Chernoff bound

$$\mathbb{P}[|X - \mu| > \alpha \mu] \leq 2 \exp[-\mu \alpha^2 / 3]$$

& obtain

$$\mathbb{P}[X < \epsilon^{m/2}] \leq 2e^{-\epsilon m / 12}$$

So if  $\epsilon m \geq 24$  the probability that at least  $\epsilon^{m/2}$  points in  $S_2$  satisfy  $\tilde{h}(x)=1$  is at least  $\frac{1}{2}$ .

Now we define a new event  $B$ :

A sample  $S = S_1 \cup S_2$  of size  $2m$  is drawn with  $|S_1| = |S_2| = m$  according to  $\mathcal{D}$ , and there exists  $\tilde{h} \in \Pi_{\Delta_{h, \epsilon}(\mathcal{D})}(S)$  such that

$$|\{x \in S \mid \tilde{h}(x)=1\}| \geq \epsilon^{m/2} \text{ & } \tilde{h}(x)=0 \text{ for } x \in S_1.$$

Note  $B$  subsumes  $A$  so  $\mathbb{P}[A \cap B] = \mathbb{P}[B]$

& hence  $\mathbb{P}[B] = \mathbb{P}[B|A] \cdot \mathbb{P}[A]$

Thus we have calculated  $\mathbb{P}[B|A]$  already it is  $\frac{1}{2} \geq \frac{1}{2}$ .

Thus  $\mathbb{P}[A] \leq 2 \mathbb{P}[B]$

(4)

But we can bound  $IP[B]$  by pure combinatorics.

We are given  $2^m$  balls out of which  $r (\geq \epsilon^m/2)$  are red & the rest are black.

If we divide the set into two equal disjoint subsets what is the probability that the first set has no red balls and the second set has all of them?

$$\binom{m}{r} / \binom{2^m}{r} \leq \left(\frac{1}{2}\right)^r$$

Thus  $IP[A] \leq IP[B] \leq$

$$\leq 2 \cdot |\Pi_{\Delta_{h(2)}}(s)| \cdot 2^{-\epsilon^m/2}$$

UNION BOUND

$$\leq 2 \cdot |\Pi_{\Delta_h(2)}(s)| \cdot 2^{-\epsilon^m/2}$$

$$\leq 2 \cdot \left(\frac{2^m}{d}\right)^d \cdot 2^{-\epsilon^m/2}$$

Now collecting factors ind. of  $d$  &  $m$  into  $K_0$  & setting this bound =  $\delta$  & solving for  $m$  we get that if  $m \geq K_0 \left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{d}{\epsilon} \log \frac{1}{\delta}\right)$

$$IP[A] < \delta$$