

Hoeffding's Inequality:

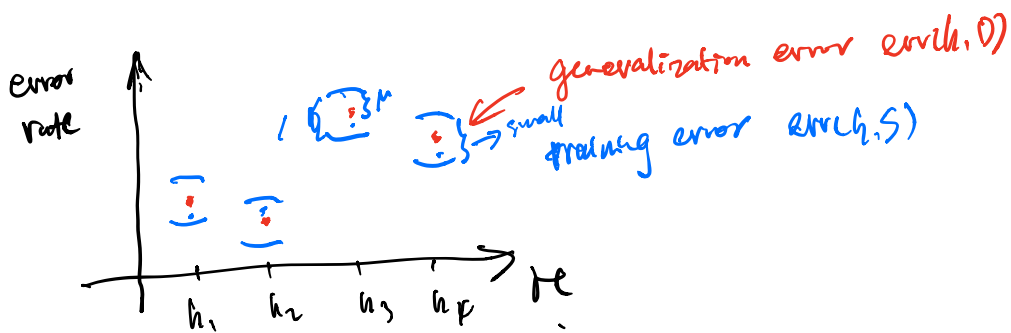
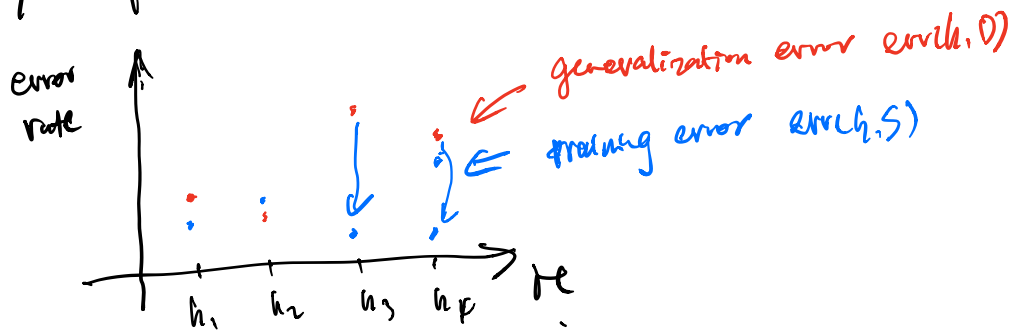
fix  $h$ .  $S =$  dataset of size  $m$  drawn i.i.d from  $D$ . (independent from  $h$ ), then:

$$P(|\text{err}(h, S) - \text{err}(h, D)| \geq \sqrt{\frac{\ln 2/\delta}{2m}}) \leq \delta.$$

Agnostic PAC  $\forall \epsilon, D$   $\min_{h \in H} \text{err}(h, D)$  may or may not be 0.

ERM input  $S$ : return  $\hat{h} = \underset{h \in H}{\text{argmin}} \text{err}(h, S)$ .

Analysis of ERM:



$$\forall h: |\text{err}(h, S) - \text{err}(h, D)| \leq \mu \quad (*)$$

$$V = \min_{h \in H} \text{err}(h, D) \quad h^* = \underset{h \in H}{\text{argmin}} \text{err}(h, D)$$

$$\text{err}(\hat{h}, S) \leq \text{err}(h^*, S) \quad (\text{optimality of ERM})$$

$$\leq \nu + \mu \quad (*)$$

$$\text{err}(\hat{h}, D) \stackrel{?}{\leq} \text{err}(\hat{h}, S) + \mu$$

$$\leq (\nu + \mu) + \mu = \nu + 2\mu.$$

Next: find  $\mu$ . (w.p.  $1 - \delta$ ) with high probability  $(1 - \delta)$

$$(*) = \bigcap_{h \in H} \{ |\text{err}(h, S) - \text{err}(h, D)| \leq \mu \}$$

$$\overline{(*)} = \bigcup_{h \in H} \{ |\text{err}(h, S) - \text{err}(h, D)| > \mu \}$$

$$P(\overline{(*)}) \leq \sum_{h \in H} P(\{ |\text{err}(h, S) - \text{err}(h, D)| > \mu \})$$

$$\stackrel{\text{Hoeffding}}{\leq} \sum_{h \in H} 2e^{-2m\mu^2} = |H| \cdot 2e^{-2m\mu^2}$$

$$\text{choose } \mu \text{ s.t. } |H| \cdot 2e^{-2m\mu^2} = \delta$$

$$\mu = \sqrt{\frac{\ln \frac{2|H|}{\delta}}{2m}}$$

In summary: w.p.  $1 - \delta$ .

(\*) happens w.p.  $1 - \delta$

$$\Rightarrow \text{err}(\hat{h}, D) \leq \min_{h \in H} \text{err}(h, D) + 2 \cdot \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2m}}$$

here  $m$   
 $\leq \frac{1}{\delta}$

Theorem: ERM agnostic PAC learns  $\mathcal{H}$ , with a simple complexity function of  $f(\epsilon, \delta) = \frac{2}{\epsilon^2} (\ln|\mathcal{H}| + \ln \frac{1}{\delta})$

$$2 \sqrt{\frac{\ln|\mathcal{H}| + \ln \frac{1}{\delta}}{2m}} = \epsilon$$

consistency  
 $\frac{1}{\epsilon^2} (\ln|\mathcal{H}| + \ln \frac{1}{\delta})$

$$\frac{\ln|\mathcal{H}| + \ln \frac{1}{\delta}}{2m} = \frac{\epsilon^2}{4}$$

$$m = \frac{2}{\epsilon^2} (\ln|\mathcal{H}| + \ln \frac{1}{\delta})$$

Infinite hypothesis classes; Vapnik Chervonienkis Theory (VC)

$|\mathcal{H}| = \infty \Rightarrow \mathcal{H}$  is not PAC learnable?

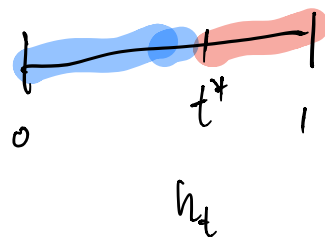
$|\mathcal{H}| < \infty \Rightarrow \mathcal{H}$  is PAC learnable by ERM.  $\forall$

Ex Infinite class can be PAC learnable

$$\mathcal{X} = [0, 1]$$

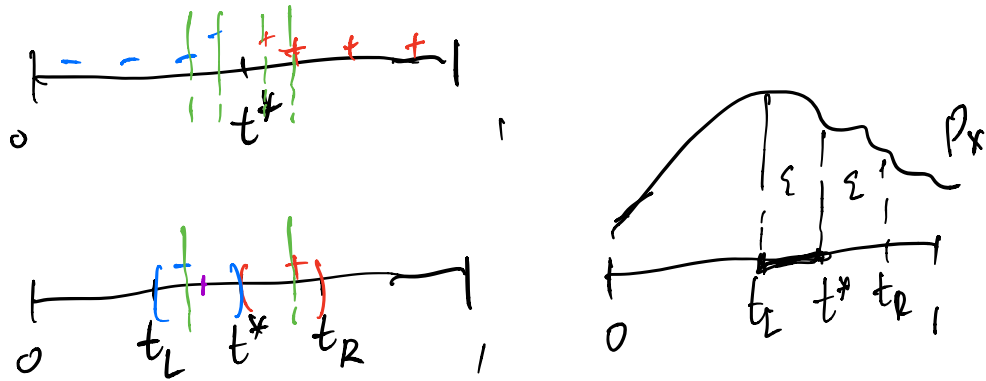
$$\mathcal{Y} = \{-1, 1\}$$

$$\mathcal{H} = \{h_t = 2I(x > t) - 1 : t \in [0, 1]\}$$



$\mathcal{D}$  realizable by  $h_{t^*}$   $X \sim \text{uniform}([0, 1])$ .

How does the consistency alg. do?



Observation: if  $\exists$  training example in  $[t_L, t^*]$ , and  
 $\exists$  training example in  $[t^*, t_R]$ ,  
 then  $\hat{h}$ : consistent classifier will have  
 $\text{err}(\hat{h}, D) \leq \epsilon$ .

want:

$$P(E_L \cap E_R) \geq 1 - \delta$$

$$\Leftrightarrow P(\bar{E}_L \cup \bar{E}_R) \leq \delta$$

$$P(\bar{E}_L) = P(\text{all training examples are outside } [t_L, t^*])$$

$$= \prod_{i=1}^m P(x_i \text{ outside } [t_L, t^*])$$

$$= (1 - \epsilon)^m$$

$$P(\bar{E}_R) \leq (1-\epsilon)^m$$

$$P(\bar{E}_L \vee \bar{E}_R) \leq \underline{2 \cdot (1-\epsilon)^m} \quad \begin{array}{l} \text{want this to be } \leq \delta. \\ (-x \leq e^{-x}) \end{array}$$

$$\Leftrightarrow 2 \cdot e^{-m\epsilon} \leq \delta$$

$$\Leftrightarrow m \geq \frac{1}{\epsilon} \ln \frac{2}{\delta}$$

Thm: the consistency algorithm learns  $\mathcal{H}$  (threshold class) with a sample complexity function of  $f(\epsilon, \delta) = \frac{1}{\epsilon} \ln \frac{2}{\delta}$ .

### VC dimension

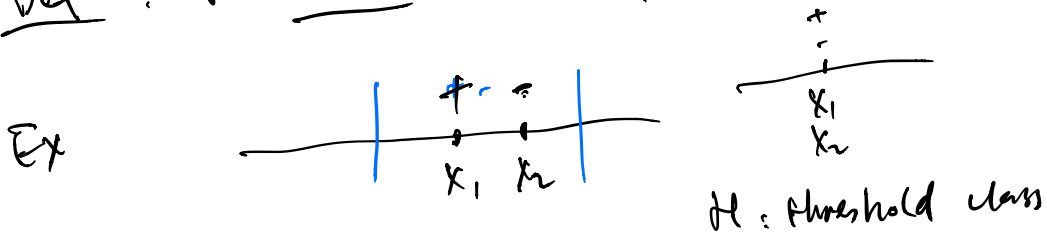
characterizing the complexity / expressiveness of hypothesis classes.

Def: hypothesis class  $\mathcal{H}$ ,  $\subseteq (\mathcal{X} \rightarrow \{+1, -1\})$  sequence of unlabeled examples  $S = (x_1, \dots, x_n)$ , define the projection of  $\mathcal{H}$  on  $S$

as 
$$\Pi_{\mathcal{H}}(S) = \left\{ \underline{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}} \right\}$$

$$|\Pi_{\mathcal{H}}(S)| \stackrel{?}{\leq} 2^n \quad \begin{array}{l} \mathcal{H} \\ \{+1, -1\} \end{array}$$

Def:  $\mathcal{H}$  shatters  $S$  if  $|\Pi_{\mathcal{H}}(S)| = 2^n$ .



$S = \{x_1\}$  ✓

$S = \{x_1, x_2\}$  ✗

Def: The VC dimension of  $\mathcal{H}$  (abbrev.  $VC(\mathcal{H})$ ).

is  $\max \{ n \in \mathbb{N} : \mathcal{H} \text{ can shatter } n \text{ points} \}$ .

$VC(\mathcal{H}_{\text{threshold}}) = 1$

Lemma  $d \in \mathbb{N}$ .  
 $VC(\mathcal{H}) = d \Leftrightarrow \left\{ \begin{array}{l} \mathcal{H} \text{ can shatter } d \text{ points, and} \\ \mathcal{H} \text{ cannot shatter } d+1 \text{ points.} \end{array} \right.$

$\Leftrightarrow \left\{ \begin{array}{l} \exists d \text{ points shatterable by } \mathcal{H}, \text{ and} \\ \nexists d+1 \text{ points shatterable by } \mathcal{H} \end{array} \right.$

Comment if  $\forall n$ . can find  $n$  points shatterable by  $\mathcal{H}$ ,  
 $VC(\mathcal{H})$  is defined as  $\infty$ .

∇ dataset of size 2 (x<sub>1</sub>, x<sub>2</sub>)

