



Need  $\Theta\left(\frac{\ln \frac{1}{\delta}}{\epsilon}\right)$  samples to learn  $\hat{r}$  which has error  $\leq \epsilon$  w.p.  $1-\delta$ .

Assume no pt masses in  $D$

Pf: The errors our alg makes are only in the shaded region (realizable & smallest  $\mathcal{O}$  containing all training pts).

$$\mathcal{H}_{\text{Bad}} = \left\{ \tilde{r} : D(\tilde{r}, r^*) \geq \epsilon, \tilde{r} < r^* \right\}$$

$$\begin{aligned} & \Pr[\text{selecting } \tilde{r} \in \mathcal{H}_{\text{BAD}}] \rightarrow r_{\text{max}} \text{ is the max} \\ & \leq \Pr[\text{We see no samples } \in B(r^*, 0) \setminus B(r_{\text{max}}, 0)] \\ & = (1-\epsilon)^m \leq e^{-\epsilon m} \end{aligned}$$

Recap:  $\xrightarrow{\text{realizable}} \Theta\left(\frac{\log \frac{|\mathcal{H}|}{\delta}}{\epsilon}\right)$  - PAC learnable

Finite  $\mathcal{H} \xrightarrow{\text{Non}} \Theta\left(\frac{\log \frac{|\mathcal{H}|}{\delta}}{\epsilon^2}\right)$  - Agnostically PAC learnable

Others? Finiteness isn't necessary ( $\uparrow$ , circle example)

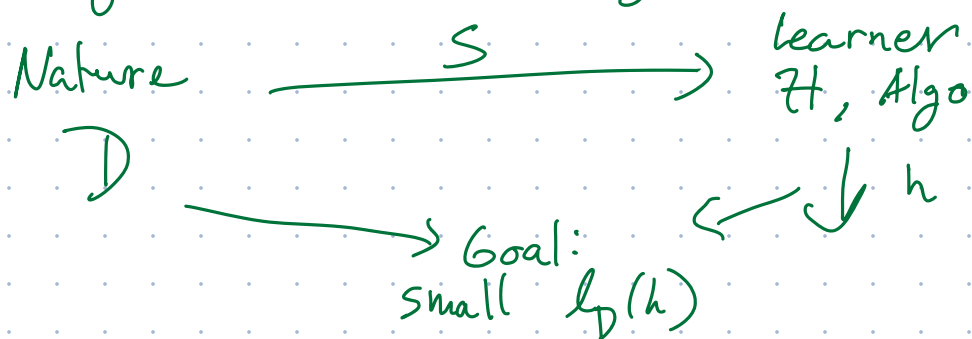
But we need info about  $D$  or  $\mathcal{H}$ ...

Our prior knowledge was that

- $D$  had a ground truth circular decision bd
- $\mathcal{H}$  was the set of circular decision bds

Q: Was this (necessary) <sup>yes</sup> or ~~just what we used?~~

Learning Problem  $\mathcal{P}(\mathcal{X}, \mathcal{Y})$  fixed.



Q: Can there be a universal alg that PAC learns every problem  $\mathcal{P}$ ?

Q': Is prior knowledge necessary for learning?

Parametric  $\mathcal{D}$   
Realizable  $\mathcal{H}$   
Small  $\min_{h \in \mathcal{H}} L_D(h)$

How much prior info should / can we assume?

- ↑ <sup>A lot</sup>
    - Learner knows  $\mathcal{D}, h^*$  (Bayes opt)  $\rightarrow \mathcal{H} = \{h^*\}$  achieves min risk
    - Learner knows a class
      - Realizable finite  $\mathcal{H} \Rightarrow$  PAC  $\tilde{m} \frac{|\mathcal{H}|}{\epsilon}$
      - Some realizable inf  $\mathcal{H} \Rightarrow$  " "
      - Nonrealizable finite  $\mathcal{H} \rightarrow$  Agnostic PAC  $\tilde{m} \frac{\ln |\mathcal{H}|}{\epsilon^2}$
  - Learner knows nothing:
  - ↓ Nothing
- No free lunch theorem:

Let  $\boxed{\text{Alg}}$  be any learning alg,  $\mathcal{Y} = \{0, 1\}$ , any  $\mathcal{X}$   
 $\forall m \leq \frac{|\mathcal{X}|}{2}$ , Then  $\exists \mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$  s.t.

1.  $\exists f: \mathcal{X} \rightarrow \mathcal{Y}$  w.  $L_D(f) = 0$
2.  $P_S(L_D(\text{Alg}(S)) \geq \frac{1}{8}) \geq \frac{1}{2}$

PF idea Randomly label all of  $X$ , make dist uniform over  $X$ .  $S$  will contain  $\leq \frac{|X|}{2}$  samples, on the remaining  $\frac{|X|}{2}$  or  $\frac{1}{2}$  of samples, can only guess  $\hat{f}$ , error  $\frac{1}{2}$  (Total error  $\geq \frac{1}{4}$ )

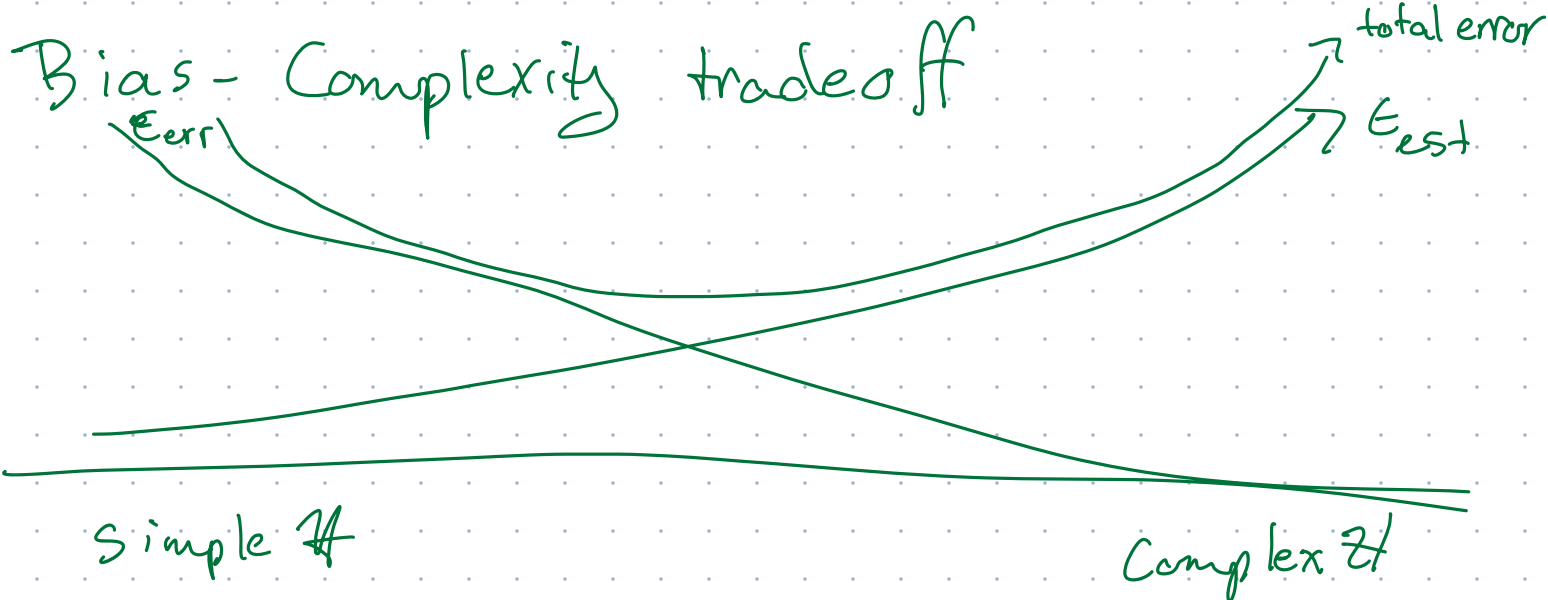
But  $f(x) = 1$ ,  $L_D(\hat{f}) = 0$ .

Corollary With infinite  $X$  and  $\mathcal{H}$  all fns:  $X \rightarrow \{0, 1\}$   
 $\mathcal{H}$  is not PAC learnable.

$$L_D(h_S) = \underbrace{(L_D(h_S) - \min_{h' \in \mathcal{H}} L_D(h'))}_{\substack{\text{ERM} \\ \uparrow \\ \text{ERM}}} + \underbrace{\min_{h' \in \mathcal{H}} L_D(h')}_{\substack{\downarrow \\ \epsilon_{\text{approx}}, \text{approx.}}}$$

- $\epsilon_{\text{est}}$ , estimation err
- generally  $\downarrow$  w  $m \uparrow$  for low bias
- Depends on complexity of  $\mathcal{H}$ ,

- $\perp$  of  $m$
- fn of  $D, \mathcal{H}$



What  $\mathcal{H}$ 's are PAC learnable?  $\rightarrow$  Vladimir Vapnik / Alexey Chervonevskis 1970

Def (Restrict  $\mathcal{H}$  to  $C$ )

$\mathcal{H}$  is a class of fns  $X \rightarrow \{0, 1\}$  &  $C = (c_1, \dots, c_m) \in X$

The restriction of  $\mathcal{H}$  to  $C$  is

$$\mathcal{H}_C \triangleq \{ \tilde{h} : C \rightarrow \{0, 1\} \mid h \in \mathcal{H} \}$$

$$c_i \mapsto h(c_i)$$

[  $h$  can be represented by vec  $(h(c_1), \dots, h(c_m))$  ]

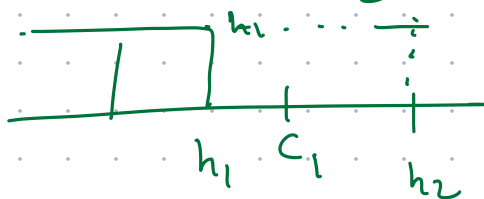
Def [Shattering]

$\mathcal{H}$  shatters  $C \subseteq X$  if  $\mathcal{H}_C$  is the set of all fns from  $C \rightarrow \{0, 1\}$ ,  $|\mathcal{H}_C| = 2^{|C|}$ .

Ex  $\mathcal{H} = \{ h_a : \mathbb{R} \rightarrow \mathbb{R} = \mathbb{I}(x \leq a) \}$  (all threshold fns)

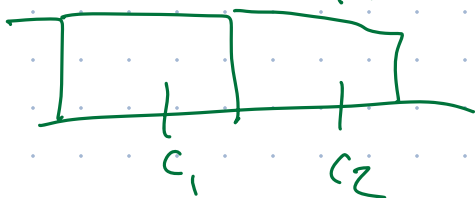
$$C = \{ c_1 \}$$

$\mathcal{H}$  shatters  $C$



Ex2  $C' = \{ c_1, c_2 \}$

$\mathcal{H}$  doesn't shatter  $C'$



Can't get  
 $h(c_1) = 0$   
 $h(c_2) = 1$

# Def VC dimension

The VC dimension of  $\mathcal{H}$ ,  $VC(\mathcal{H})$ , is the maximum-sized set  $C$  which is shattered by  $\mathcal{H}$ . (No larger sets shattered by  $\mathcal{H}$ ).

To Prove  $VCdim(\mathcal{H}) = d$ , - find a witness set  $|C| = d$  which can be shattered

• & argue no  $|C'| > d$  can be shattered.

Ex: Thresholds  $\uparrow$

$$VCdim \geq 1$$

$= 1$  (NTS  $\nexists$  a set of size 2 shatterable)

Ex Intervals  $\mathcal{H} = \{ \Pi_{a \leq x \leq b} : a \leq b \in \mathbb{R} \}$

$$VC(\text{Intervals}) = 2$$

Ex: Axis align rectangle  $\mathcal{H} = \{ \Pi_{\substack{a_1 \leq x_1 \leq b_1 \\ a_2 \leq x_2 \leq b_2}} : \substack{a_1 \leq b_1 \\ a_2 \leq b_2} \}$

$$VC(\text{Rect}) = 4$$

\*  $c_4$

\*  $c_1$

\*  $c_2$

\*  $c_3$

# Parameters  
is often  
 $\hat{=}$  VCdim

But not  
always!!!

Ex: Finite  $\mathcal{H}$

Shattering requires  $|\mathcal{H}| \geq 2^{|C|}$

$$VC(\mathcal{H}) = |C| \leq \lg_2 |\mathcal{H}|$$

# Thm 1 Fund Thm of Statistical Learning

Let  $\mathcal{H}$  be a hypothesis class from  $\mathcal{X}$  to  $\{0,1\}$ ,  
l the 0-1 loss fn. Then the following are  
equivalent:

- (1)  $\mathcal{H}$  has uniform convergence
- (2) Any ERM rule successfully agnostically PAC learns  $\mathcal{H}$
- (3)  $\mathcal{H}$  is agnostically PAC learnable
- (4)  $\mathcal{H}$  is PAC learnable
- (5) Any ERM rule PAC learns  $\mathcal{H}$
- (6)  $\mathcal{H}$  has finite VC dimension

(1  $\rightarrow$  2) Last lec

2  $\rightarrow$  3 Yup

3  $\rightarrow$  4 Yup

2  $\rightarrow$  5 Yup

2  $\rightarrow$  6 ( $\neg 6 \rightarrow \neg 2$  by NFL)

Remains to show 6  $\rightarrow$  1