

---

---

---

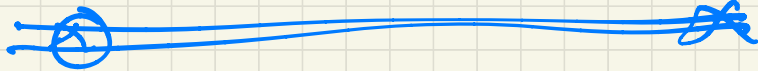
---

---



Consistency,  
Compression,  
and  
Learning:  
The Finite  $H$  Case

# Recipe for (PAC) Learning:



1. Design algo  $L$  that finds  $h \in \mathcal{H}$  consistent with sample  $S$   
( $\hat{\epsilon}_S(L(h)) = 0$ )
2. Analyze how big  $m = |S|$  must be s.t.  $\epsilon(h) \leq \epsilon$ .

We will show:

• This recipe **always** works

• Answer to 2. is **independent\*** of  $L$ -consistency is all that matters.

Let's warm up with the case of **finite  $\mathcal{H}$ .**

# Notation

$$\varepsilon(h) \stackrel{\text{a}}{=} P_{x \sim p} [h(x) \neq c(x)]$$

true error

$$\hat{\varepsilon}_S(h) \stackrel{\text{a}}{=} \frac{1}{m} \sum_i I[h(x_i) \neq y_i]$$

where  $S = \{ \langle x_0, y_0 \rangle, \dots, \langle x_m, y_m \rangle \}$

training error

~~Q: When does  $\hat{\varepsilon}_S(h) = 0$~~

imply  $\varepsilon(h)$  small?

• Fix  $\epsilon > 0$

• Call  $h \in \mathcal{H}$   $\epsilon$ -bad if  $\epsilon(h) \geq \epsilon$

•  $\forall$  fixed  $\epsilon$ -bad  $h$ :

$$P_{r_S}[\hat{\epsilon}_S(h) = 0] \leq (1-\epsilon)^m$$



indep.

$$P_{r_S}[\text{any } \epsilon\text{-bad } h \in \mathcal{H} \text{ has } \hat{\epsilon}(h) = 0]$$

$$\leq |\mathcal{H}| (1-\epsilon)^m$$

union bound

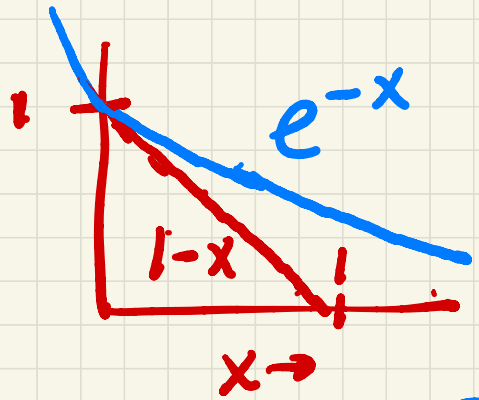
# Algebra:

$$|H|(1-\varepsilon)^m \leq$$

$$|H|e^{-\varepsilon m} \Rightarrow$$

set  $\leq \delta$  & solve:

$$m \geq \frac{1}{\varepsilon} \ln \frac{|H|}{\delta}$$



$$\therefore 1-x \leq e^{-x}$$

"complexity"

of  $\mathcal{H}$

$$: \ln |H|$$

# bits <sup>ss</sup> needed  
to describe  $h \in \mathcal{H}$

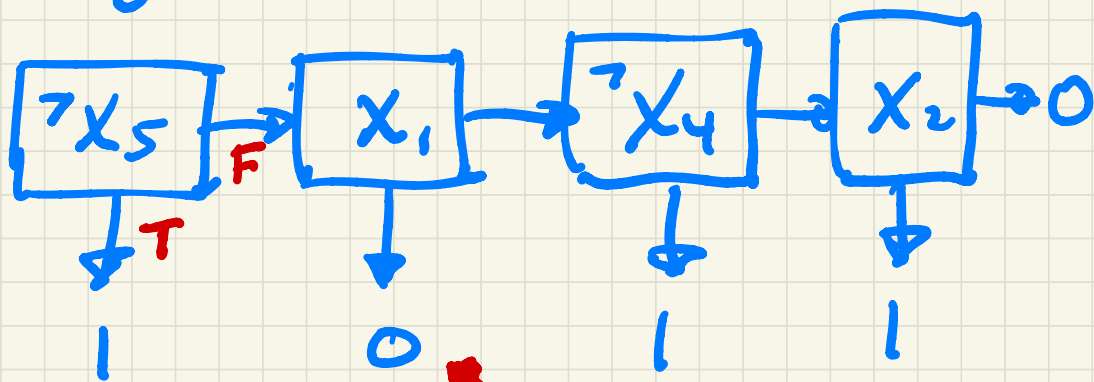
- Immediately applies to & simplifies PAC analyses for rectangles, conjunctions, 3CNF

- Any consistent  $h \in H$  suffices



Another application:  
decision lists over  $\{0,1\}^n$ .

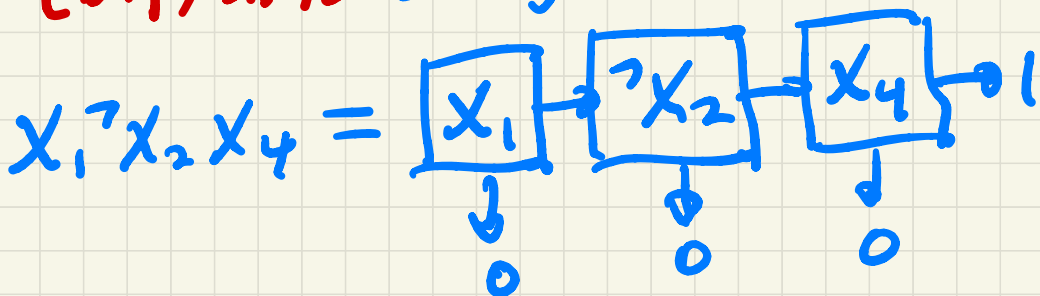
e.g.  $C \in \mathcal{C}$  given by:



$$C(01101) = 1$$

$$C(10011) = 0$$

Contains conjunctions:



# Consistent algo:

- $S \leftarrow$  all examples
- $h \leftarrow$  empty list
- while ( $S \neq \emptyset$ ):

-  $\forall z, b$ :

$$S_{z,b} \leftarrow \{ \langle x, y \rangle \in S : z = 1 \wedge x \text{ and } y = b \}$$

- find  $\max |S_{z,b}|$  s.t.

$$S_{z,\neg b} = \emptyset$$

- append  $\rightarrow \boxed{z}$  to  $h$

$\downarrow$   
 $b$

-  $S \leftarrow S - S_{z,b}$

Theorem Decision lists  
are PAC learnable  
(by  $\mathcal{H}$ =decision lists).

# Learning & Compression

- For  $h \in \mathcal{H}$ , let  $l(h)$  be the # bits needed to describe  $h$ .
- In general,  $l(h) \approx \text{poly}(n)$   
(dim. of  $X$ )  
 $\Rightarrow |\mathcal{H}| \approx 2^{\text{poly}(n)}$

- Suppose we allow

$$l(h) \approx \text{poly}(n, m)$$

$|S| \nearrow$

Good or bad idea?

If we allow

$l(h) \sim n \cdot m$  (linear  
in both)

then we can just

encode/memorize  $S$ !

( $\mathcal{H}$  = lists of  $\langle x, y \rangle$ )

~~—————~~

So linear dependence  
on  $m$  goes too far.

Let's try  $l(h) = c \cdot m^\alpha$

includes  $\nearrow$   
dep. on  $n$   $\alpha \neq 1$

$$|Z| e^{-\epsilon m} = 2^{c m^\alpha} e^{-\epsilon m}$$

$$\leq e^{c m^\alpha - \epsilon m}, \text{ set } \leq \delta:$$

$$c m^\alpha - \epsilon m \leq \ln(\delta)$$

$$\epsilon m \geq \ln(1/\delta) + c m^\alpha$$

Satisfied if:

$$m \geq \frac{2}{\epsilon} \ln(1/\delta) \quad \&$$

$$m \geq \frac{2 c m^\alpha}{\epsilon}$$

$$m \geq \frac{2cm^\alpha}{\varepsilon}$$

$$m^{1-\alpha} \geq \frac{2c}{\varepsilon}$$

$$m \geq \left(\frac{2c}{\varepsilon}\right)^{\frac{1}{1-\alpha}} \quad \leftarrow \text{blows up as } \alpha \rightarrow 1!$$

$\alpha = 0$ :  $m \sim 2c/\varepsilon$ , original bound

$\alpha = 1/2$ :  $m \sim \left(\frac{2c}{\varepsilon}\right)^2$

$\alpha = 1$ :  $m \sim \infty$

So not just consistency  
but even the slightest

compression of  $S$

yields PAC learning  
(with larger  $m$ ).

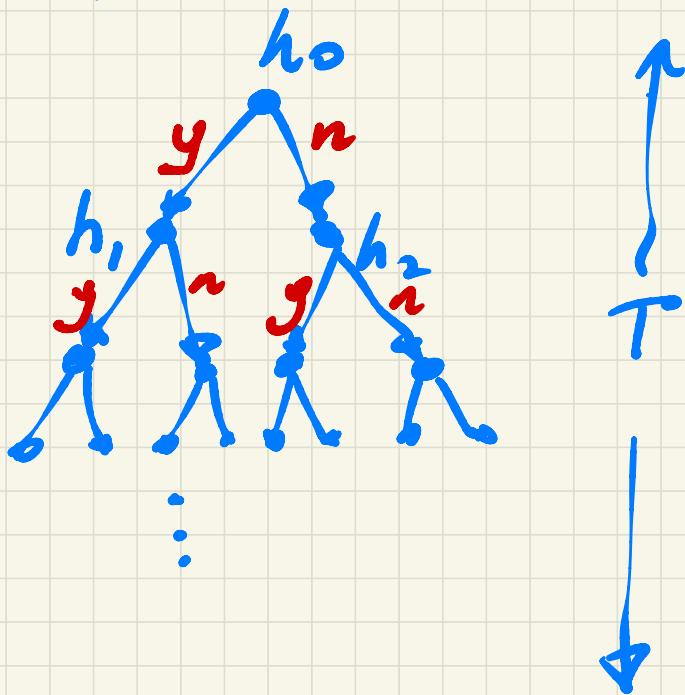


# Application to Adaptive ML/Crowdsourcing

- Imagine we have data  $S_{train}, S_{test}$
- We give  $S_{train}$  to crowd
- We keep  $S_{test}$  to validate
- Now  $\mathcal{H}$  not fixed in advance and not finite!
- How big should  $m = |S_{test}|$  be?

# A Modified Leaderboard

- When given model  $h$ :
  - improvement  $< \alpha$  on  $S_{\text{test}}$ :  
answer **no**, nothing else
  - improvement  $\geq \alpha$  on  $S_{\text{test}}$ :  
**yes**, update best error



• Observation: at most  $\frac{1}{\alpha}$   
y's on any path

$$\therefore \text{total size of tree} \\ = |\mathcal{H}| \leq \binom{T}{\frac{1}{\alpha}} \leq T^{1/\alpha}$$

$$\text{Solve: } T^{1/\alpha} e^{-\epsilon^2 m} \leq \delta$$

$$\epsilon^2 m \geq \left(\frac{1}{\alpha}\right) \ln\left(\frac{T}{\delta}\right)$$

$$m \geq \left(\frac{1}{\epsilon^2 \alpha}\right) \ln\left(\frac{T}{\delta}\right)$$

One more variation.

So far we have  
shown (for  $n$  large)

$$|\hat{\epsilon}_S(h) - \epsilon(h)| =$$

$$|0 - \epsilon(h)| = \epsilon(h) \leq \epsilon$$

for all consistent  $h$ .

What about all  
the other  $h \in \mathcal{H}$ ?

# Chernoff bounds

- Consider biased coin with  $\Pr[\text{heads}] = p$
- Flip  $n$  times, let  $\tilde{p} =$  fraction of heads

Then  $\forall \gamma > 0$ :

$$\Pr[|\tilde{p} - p| \geq \gamma] \leq 2e^{-n\gamma^2/3}$$

→ 0 exponentially fast

- A "concentration inequality"

- $\forall$  fixed  $h \in \mathcal{H}$ :

$$\Pr[|\hat{\mathbb{E}}_S(h) - \mathbb{E}(h)| \geq \varepsilon] \leq \text{blah}$$

- Prob. any

$$h \in \mathcal{H} \text{ has } \leq |\mathcal{H}| \cdot \text{blah}$$

$$| \cdot | \geq \varepsilon$$

- $\leq \delta$  if

$$m \sim \frac{1}{\varepsilon^2} \ln \frac{|\mathcal{H}|}{\delta}$$

"uniform convergence"

WYSIWYG

Fine.

But what if

$H$  is

infinite ???