

## Last time:

- finish proof that can efficiently improperly learn 3-term DNF
- " can't eff. properly learn 3-term DNF  
(unless  $NP \subseteq RP$ ).

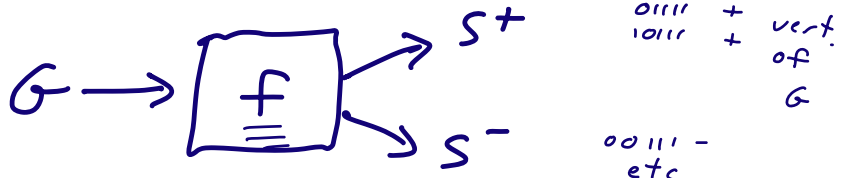
Today: - finish pf that (prove L2, completing pf about the transformation  $f$ )  
- start VC dimension unit (KV Chap. 3) / unit on sample complexity of PAC learning

- lower bound on  $\downarrow$ :  $\Omega(\text{VCdim}(\mathcal{E})/\epsilon)$
- start upper bound on  $\downarrow$ : (roughly)  $O(\text{VCdim}(\mathcal{E})/\epsilon)$

Reminder: in-class closed-book midterm on Mon.

## Questions?

To prove:



(L2) If  $\exists$  3-term DNF cons. w/  $S^+ + S^-$ ,  
then  $G$  is 3-colorable.

Pf : Let  $T_R \vee T_Y \vee T_B$  be any 3-term DNF  
cons. w/  $S^+ + S^-$ .  
 $\hookrightarrow$  every + ex is acc. by some term

Consider the col. which colors vertex  $i$  of  $G$

R  
Y  
or B

dep. on which  
of  $T_R, T_Y, T_B$  acc.  $i$   $1 \dots 101 \dots 1$

(break ties arbitrarily).

Is this col. legit?

Claim: Let  $\{i,j\}$  be any edge. Then above col.  
doesn't color both  $i$  &  $j$  R. (same for B, Y).

PF: Consider  $i,j = 1,2$  : show not both sat by  $T_R$

sat by  $T_R$   $\left\{ \begin{array}{l} \rightarrow 01111 + \text{ex } 1 \leftarrow \\ \rightarrow 10111 + \text{ex } 2 \leftarrow \\ \rightarrow 00111 - \text{ex } 12 \leftarrow \end{array} \right.$

$x_1$  not be in  $T_R$   
 $\bar{x}_1$  not in  $T_R$

wouldn't be + by  $T_R$

Similarly,  
 $x_2$  not in  $T_R$   
 $\bar{x}_2$  not in  $T_R$ .

But then if  $T_R$  acc the  $2 + \text{ex}$ ,  
it would have to acc  $i$  : contrad.



# SAMPLE COMPLEXITY OF PAC LEARNING

(For this whole unit, we only worry about, not comput. efficiency/runtime.)

Motiv.: Given  $\mathcal{C}$ , basic questions:

- is  $\mathcal{C}$  PAC learnable (from a finite # of samples)?
- if so, how many ex. needed?

Sps  $\mathcal{C}$  finite.

CHF gives  $\frac{1}{2}$  answer:  $\boxed{\frac{1}{\epsilon}(\ln |\mathcal{C}| + \ln \frac{1}{\delta})}$  ex.

is enough. (CHF for  $\mathcal{C}$  using  $\mathcal{H} = \mathcal{C}$ )

• lower bound?

• what if  $|\mathcal{C}| = \infty$ ?

Fix  $\delta = \frac{1}{10}$ .

$VCDIM(\mathcal{C})$  : whole story!

Main results of this unit:

(pretend  $\delta = \frac{1}{10}$ )

Fix any  $\mathcal{C}$ . Let  $d := VCDIM(\mathcal{C})$

→  $\textcircled{!}$  • Any PAC learning alg for  $\mathcal{C}$  must use  $\Omega\left(\frac{d}{\epsilon}\right)$  ex.  
(so if  $d = \infty$ , no finite # of ex. is enough to PAC learn)

☺ • If  $d$  finite: then  $\mathcal{C}$  is PAC learnable

(CHF for  $\mathcal{C}$  using  $\mathcal{E}$ ) with  $\approx \underline{\underline{O\left(\frac{d}{\epsilon}\right)}}$  ex.

↳ like CHF result:  $VCOIM$ , not  $\ln|\mathcal{C}|$ ,  
is "right" measure.

Lower bd on SC of PAC learning:

Thm Fix any  $\mathcal{C}$ , let  $d = VCOIM(\mathcal{C})$ .

Any PAC learner for  $\mathcal{C}$  that learns to  
error  $\epsilon$  & conf.  $\delta = \frac{1}{10}$  must use

$\Omega\left(\frac{d}{\epsilon}\right)$  ex.

$\left(\frac{d-1}{32\epsilon}\right)$

Pf: First: give a  $\mathcal{D}$  s.t.  $\Omega(d)$  ex needed  
to achieve  $\epsilon = \frac{1}{8}, \delta = \frac{1}{8}$ .

Let  $S = \{x^1, \dots, x^d\}$  be shuff. by  $\mathcal{C}$ .

Let  $\mathcal{D} = \text{unif over } S$ . (0 wt elsewhere).

Suppose alg  $A$  only gets  $\frac{d}{2}$  calls to  $EX(\mathcal{C}, \mathcal{D})$ .

After calls,  $A$  has seen labels of  $\leq \frac{d}{2}$  ex  
not " " "  $\geq \frac{d}{2}$  " .

Let target  $c \in \mathcal{C}$  be chosen by picking it unif.  
from the  $2^d$  concepts that shatter  $S$ .

So each of the  $2^d$  poss. lab. of  $S$  = ly likely.

Hence

$$\mathbb{E}[\text{error of } A\text{'s hyp}] \geq \frac{1}{2} \cdot \frac{d}{2} \cdot \frac{1}{d} \leftarrow \begin{array}{l} \text{error on each unseen pt} \\ \downarrow \quad \downarrow \rightarrow \frac{d}{2} \text{ unseen pts} \\ \leftarrow \text{prob. of each unseen pt} \end{array}$$

$$= \frac{1}{4}.$$

Let  $\text{acc} = 1 - \text{error}$ ;  $\mathbb{E}[\text{acc}] \leq \frac{3}{4}$   
 $\leq \frac{3}{4} \cdot \frac{7}{6} = \frac{7}{8}$

Markov:  $\Pr[\text{acc} \geq \frac{7}{8}] \leq \frac{6}{7}.$

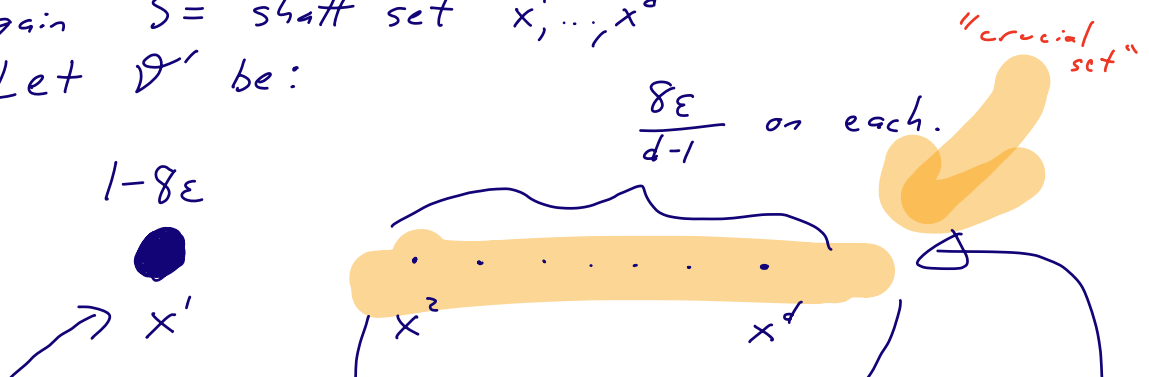
So  $\Pr[\text{error} \geq \frac{1}{8}] \geq \frac{1}{7}$ . Failed to  $(\epsilon = \frac{1}{8}, \delta = \frac{1}{8})$ -PAC learn.

Markov: if  $X \geq 0$  is a r.v.,  
 $\mu = \mathbb{E}[X]$ , for any  $k \geq 1$   
 then  $\Pr[X \geq k \cdot \mu] \leq \frac{1}{k}.$

Let's prove whole thing: give a dist  $\mathcal{D}'$   
 s.t.  $\Omega(d/\epsilon)$  ex. to get error  $\epsilon$  + conf.  $1 - \delta = \frac{9}{10}$ .

Again  $S =$  small set  $x^1, \dots, x^d$

Let  $\mathcal{D}'$  be:



$8\epsilon$  total mass

Idea:  $1-8\epsilon$  acc. very easy to achieve;  
getting  $\frac{7}{8}$  of last  $8\epsilon$  is slow:  
on avg,  $\frac{1}{8\epsilon}$  draws needed to get 1 ex from  
so "slowdown" of  $\Omega(\frac{1}{\epsilon})$  on top of earlier arg.

---

Sps learning alg  $A$  only gets  $\frac{d-1}{2}$  ex from  
crucial set  $x_1^2, \dots, x^d$ .

Let  $c \in \mathcal{C}$  (target) be rand. chosen s.t. each  
of the  $2^{d-1}$  lab. of crucial set = ly likely.

As above, w.p.  $\geq \frac{1}{8}$ , alg  $A$  has error rate  $\geq \frac{1}{8}$  on crucial  
set. These have prob.  $8\epsilon$  under  $\mathcal{D}'$ , so  
if  $A$  gets  $\leq \frac{d-1}{2}$  ex from crucial set, then  
w.p.  $\geq \frac{1}{8}$ ,  $A$  has error  $\geq \epsilon (= \frac{1}{8} \cdot 8\epsilon)$ .

Each call to  $EX(c, \mathcal{D})$  hits crucial set w.p.  $8\epsilon$ .

Sps  $A$  makes only  $\frac{d-1}{32\epsilon}$  calls to  $EX(c, \mathcal{D})$ .

Then  $\mathbb{E}\{\# \text{ times hit crucial set}\} = \frac{d-1}{32\epsilon} \cdot 8\epsilon = \frac{d-1}{4}$ .

Mult. CB ( $r=1$ ,  $p_m = \frac{d-1}{4}$ ):

$$\Pr\{X \geq (1+r)p_m\} \leq e^{-\frac{d-1}{12}}$$

Can assume wlog  $d \geq 100$ : then  $\downarrow \leq \frac{1}{100}$ .

So w.p.  $\geq \frac{99}{100} \cdot \frac{1}{8} \geq \frac{1}{9}$ ,

error  $\geq \epsilon$ .



## Upper Bound on Sample Complexity of PAC Learning

Main result:  $\forall \mathcal{D}$ , any hyp  $h \in \mathcal{C}$   
which is consistent with  $d = \text{VCdim}(\mathcal{C})$

$$\approx \frac{d}{\epsilon} \cdot \log \frac{1}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta}$$

many rand ex from  $EX(\mathcal{C}, \mathcal{D})$  is, w.p.  $\geq 1 - \delta$ ,  
 $\epsilon$ -accurate, i.e.  $er_{\mathcal{D}}(h, \mathcal{C}) \leq \epsilon$ .

$\hookrightarrow$  CHF for  $\mathcal{C}$  using  $\mathcal{H} = \mathcal{C}$  works, if you take  $\sim$  ex.

Hi level outline of pf:

a #

① Setup: Ponder  $\text{VCdim}(\mathcal{C})$ .

We'll consider a function

"growth function of  $\mathcal{C}$ "

$\hookrightarrow$  to get more info.

② Combinatorics arg: AMAZING THEOREM about how the growth fn behaves, for any  $\mathcal{C}$ .

---

③ Learning argument: a lot like CHF thm, but using AMAZING THM to control # ways conc in  $\mathcal{C}$  can label a data set.

---

### ① Setup

Recall  $\text{VC DIM}(\mathcal{C}) =$  size of largest shatt. set.

Another POV on  $\text{VC DIM}(\mathcal{C})$ :

Say  $S \subseteq X$ .

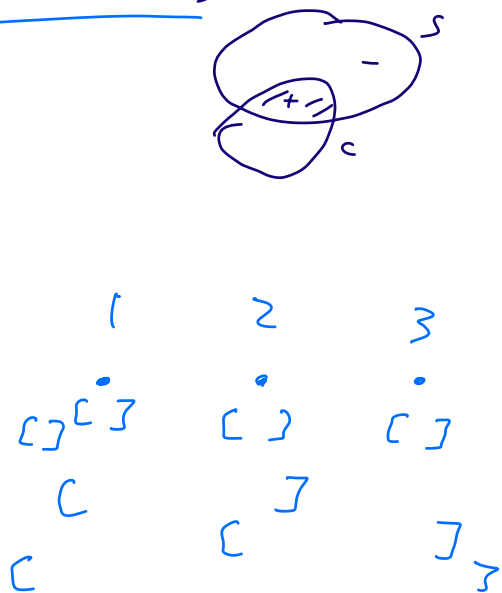
Define  $\Pi_{\mathcal{C}}(S) =$  all labelings of  $S$  induced by concepts in  $\mathcal{C}$ .

Equiv,  $\Pi_{\mathcal{C}}(S) = \{ \underline{c \cap S} : c \in \mathcal{C} \}$

Ex:  $S = \{1, 2, 3\}$

$\mathcal{C} =$  intervals of  $\mathbb{R}$

$\text{VC DIM}(\mathcal{C}) = 2$ .





$$\Pi_e(S) = \{\emptyset, \{1\}, \{2\}, \{3\}, \{1,2\}, \{2,3\}, \{1,2,3\}\}$$

7  
subsets

Note:  $\mathcal{C}$  shatters  $S$



$$|\Pi_e(S)| = 2^{|S|}$$

Def:  $\Pi_e(m) =$  growth function of  $\mathcal{C}$

$$= \max_{\substack{S \subseteq X, \\ |S|=m}} |\Pi_e(S)|$$

For  $\mathcal{C} =$  intervals,  $\Pi_e(3) = 7$

---