

- Last time:
- finished proof of our CHF theorem
  - Application: efficiently PAC learning LTFs over  $\mathbb{R}^7$
  - We now understand sample complexity of PAC learning :)
- 

Today: Boosting!

- • confidence boosting (easy)
- • start accuracy boosting:
  - setup, boosting framework
  - proof of concept: boosting a 40% error learner to a 35.2% error learner (three stage process)

Questions?

---

Recall PAC learning def:

• Alg.  $A$  PAC learns  $\mathcal{C}$  if  $\forall c \in \mathcal{C}, \forall \text{dist } \mathcal{D}$ ,

•  $\forall \epsilon > 0$ ,

•  $\forall \delta > 0$ ,

given  $\epsilon, \delta + EX(c, \mathcal{D})$ ,  $A$  produces  $h$  s.t.  $\text{err}_{\mathcal{D}}[h, c] \leq \epsilon$   
w.p.  $\geq 1 - \delta$ .

---

What happens if we relax the def? Turns out...  
nothing -- if  $\mathcal{C}$  is <sup>efficiently</sup> learnable under relaxed def,  
then " " " " original def!

---

Relax  $\delta$  (confidence) requirement:

Say  $\mathcal{C}$  is "low-confidence PAC learnable" if:

•  $\forall c \in \mathcal{C}, \forall \text{dist } \mathcal{D}$

•  $\forall \epsilon > 0,$

•  ~~$\forall \delta > 0,$~~  for  $\delta = 0.9,$

given  $\epsilon, \delta \in EX(c, \mathcal{D}), A$  produces  $h$  s.t.  $\text{err}_{\mathcal{D}}[h, c] \leq \epsilon$   
w.p.  $\geq 1 - \delta = 0.1$

Fact: if  $\mathcal{C}$  eff. LC-PAC learnable with  $m(\epsilon)$  examples,  
then  $\mathcal{C}$  is also eff. "strongly PAC learnable",  
in fact using  $\underline{O(\log(1/\delta) \cdot m(\epsilon/2))} + \underline{O(\log(1/\delta)/\epsilon^2)}$  ex.

Pf sketch:

The alg has 2 stages:

1) Run orig. LC-PAC learner with error param  $\epsilon/2,$   
 $C_1 = O(\log(1/\delta))$  times.

$\Pr[\text{none of } C_1 \text{ runs yields an } \epsilon/2\text{-error hyp}]$   
 $\leq (1 - 0.1)^{C_1} < \delta/2.$

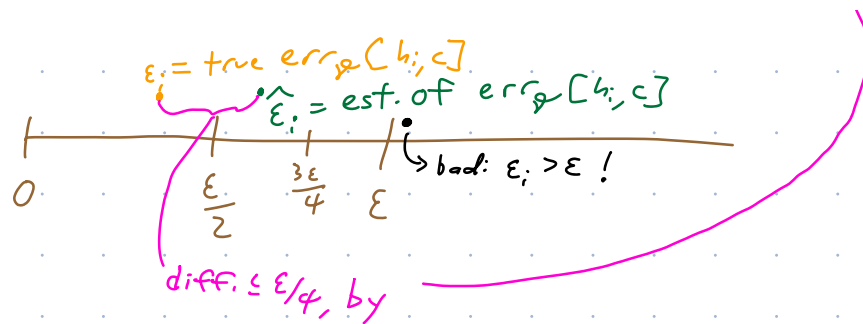
So have hyp's  $h_1, \dots, h_{C_1},$  + w.p.  $\geq 1 - \delta/2,$   
at least one  $h_i$  has  $\text{err}_{\mathcal{D}}[h_i, c] \leq \epsilon/2.$

2) Hyp. testing: Draw  $\underline{O(\frac{\log(1/\delta^2)}{\epsilon^2})}$  fresh ex.,

measure observed error of each  $h_i$  on them,  
output the  $h_i$  that does best.

CB:  $\nearrow$  for any fixed  $j \in [C_1],$   
w.p.  $\geq 1 - \frac{\delta}{2C_1} \geq 1 - \delta^2,$   
my est. of  $\text{err}_{\mathcal{D}}[h_j, c]$  is  $\pm \epsilon/4$  accurate.

So UB  $\Rightarrow$  w.p.  $\geq 1 - \frac{\delta}{2},$  every est. is  $\pm \epsilon/4$ -acc.



Any  
 Bad  $h_i$  will appear to have error  $> \frac{3\varepsilon}{4}$   
 Any  $\frac{\varepsilon}{2}$ -acc.  $h_i$  will appear to have error  $< \frac{3\varepsilon}{4}$

So total fail prob. (of <sup>not</sup> outputting an  $h_i$  with  $\text{error}[h_i, c] \leq \varepsilon$ ) is  $\leq \frac{\delta}{2} + \frac{\delta}{2} = \delta$  ☺

We'll usually ignore  $\delta$  going forward ☺

What about accuracy?

Def: Alg  $A$  is a weak PAC learner for  $\mathcal{C}$  with advantage  $\gamma$  if  $\forall c \in \mathcal{C}, \forall \text{dist } \mathcal{D},$

~~$\forall \varepsilon > 0,$~~

$\forall \delta > 0,$

given  ~~$\varepsilon$~~   $\delta + EX(c, \mathcal{D}), A$  produces  $h$  s.t.  $\text{error}[h, c] \leq \frac{1}{2} + \gamma$   
 w.p.  $\geq 1 - \delta.$

(i.e. achieve acc.  $\frac{1}{2} + \gamma$ ).

Turns out: given any weak PAC learner, can "boost" acc. to achieve strong PAC learning!

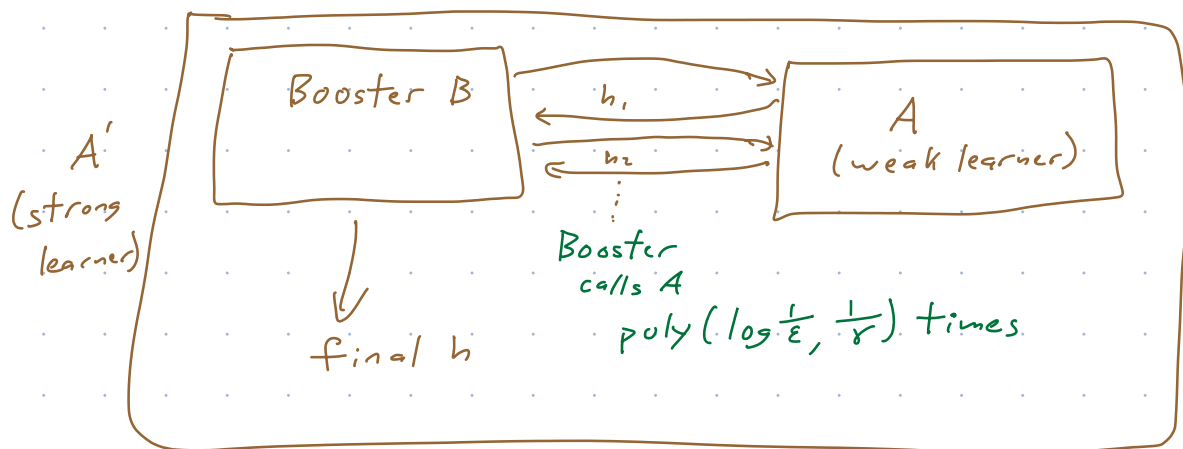
Thm: Let  $\mathcal{C}$  be any conc. class.

Sps  $A$  is an eff. weak PAC learner for  $\mathcal{C}$  with advantage  $\gamma$ .

Then there is a  $\text{poly}(\frac{1}{\gamma}, \frac{1}{\epsilon}, \log \frac{1}{\delta}, T)$ -time strong PAC learner  $A'$  for  $\mathcal{C}$  (where  $T = \text{runtime for } A$  to achieve confidence  $9/10$ , i.e.  $\delta = 0.1$ )

---

Done via efficient, explicit boosting alg.  $B$ .



Basic idea of a booster: to learn  $c$  given  $EX(c, \mathcal{D})$  + weak learner  $A$ :

- booster runs  $A$  multiple times using a seq. of different distributions  $EX(c, \mathcal{D}_1), EX(c, \mathcal{D}_2), \dots$  to get hyp's  $h_1, h_2, \dots$

- booster combines  $h_1, h_2, \dots$  to create final  $h$ .

---

To describe a boosting alg, must describe:

- what are  $\mathcal{D}_1, \mathcal{D}_2, \dots$ ? (weak learner)

Cleverly designed dist's to force  $A$  to

"give some new info".

- How can you run  $A$  on  $EX(c, \mathcal{D}_i)$  when we only have  $EX(c, \mathcal{D})$ ?

Few ways:

- filter examples; or
  - work over explicit sample.
- how to combine  $h_i$ 's? Maj vote.  
(why)
  - $\otimes$  Does it work? Does it give final  $h$  s.t.  
 $\Pr_{x \sim \mathcal{D}} [h(x) \neq c(x)] \leq \epsilon$ ?

Yes - takes a pf.

---

Let's do an example pf of concept:  
3-stage boosting.

Let  $A$  be a weak PAC learner for some  $C$   
with  $\gamma = 0.1$  (acc of  $h$  is 60%; error  $\leq 40\%$ )

We'll run  $A$  3 times, & get a hyp. w/ error  $\leq 35.2\%$ .

2 simplifying assumptions:  $\nearrow$  (ignore  $\delta$ )  
assume every run of  $A$  gives a hyp. with error  
exactly 40% on whatever dist. was used.

---

Notation:

write " $\mathcal{D}[E]$ " for  $\Pr_{x \sim \mathcal{D}}[E]$ , e.g.

$\mathcal{D}[h(x)=c(x)]$  for  $\Pr_{x \sim \mathcal{D}}[h(x)=c(x)]$

Simple 3-stage booster:

[Have w.l.  $A$ ;  $EX(c, \mathcal{D})$ ]

1) Run  $A$  on  $EX(c, \mathcal{D})$  to get  $h_1$ .  
Have  $\mathcal{D}_1[h_1(x) \neq c(x)] = 0.4$

2) One idea: let  $\mathcal{D}_2$  be  $\mathcal{D}_1$  restricted to those  $x$  s.t.  $h_1(x) \neq c(x)$ .

Problem: the hyp  $h_1(x)$  has 100% acc on this  $\mathcal{D}_2$ , so w.l. could give you that... no new info !!

60% of  $\mathcal{D}$ :  
pts  $x$  s.t.  
 $h_1(x) = c(x)$

40% of  $\mathcal{D}$ :  
pts  $x$  s.t.  
 $h_1(x) \neq c(x)$

Key idea: • decrease wt of

by  $5/6$ , so now 50% of  $\mathcal{D}_2$ .

This is our  $\mathcal{D}_2$

• increase wt of 50% of  $\mathcal{D}_2$

by  $5/4$ , so now

Q: how to sim. access to  $EX(c, \mathcal{D}_2)$  given  $h_1$  &  $EX(c, \mathcal{D})$ ?

1 way: when  $EX(c, \mathcal{D})$  has  $h_1(x) = c(x)$ ,  
Keep w.p.  $\frac{2}{3}$  & try again w.p.  $\frac{1}{3}$   
 $\rightarrow p = \frac{2}{3}$

Another way: toss fair coin to decide between  
 $h_1 = c$  or  $h_1 \neq c$   
 $\rightarrow p = .6$   $\rightarrow p = 0.4$

& call  $EX(c, \mathcal{D})$  till you get an example of  
the type you want.

---

Both ways efficient:

given  $\mathbb{D} \sim \Pr[\mathbb{D} = H] = p,$

$$\mathbb{E}\{\# \text{ tosses till get } H\} = \frac{1}{p}.$$

---