

Last time: • confidence boosting (easy)

• start accuracy boosting:

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

3-stage

Today: • finish

- discuss boosting-by-filtering vs boosting over a fixed sample
- start discussion + analysis of AdaBoost (✓)

Questions?

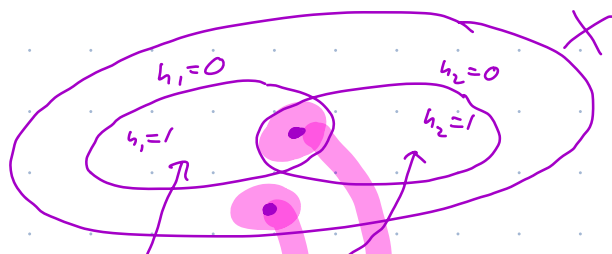
(addit.) $C \leq B$: toss coin w/ $\Pr[\text{coin} = H] = p$
 $m = \frac{O(\log(1/\delta))}{\epsilon^2}$ times:

observed freq. of H is
 $\pm \epsilon$ close to p

with prob. $\geq 1 - \delta$.

At this point, have $h_1, \& h_2$.

Need to construct h_3 , & then our final h will be
 $h = \text{MAJ}(h_1, h_2, h_3)$.



matters only on $A \cup B$ h_3 's value doesn't matter here -

So \mathcal{D}_3 should put all its wt. on regions $A \cup B$.

We take \mathcal{D}_3 to be \mathcal{D}_1 with 0 wt on x s.t. $h_1(x) = h_2(x)$
& wt on $A \cup B$ "scaled up" so total weight = 1.

Run A on $EX(c, \mathcal{D}_3)$ to get h_3 ;

$$h = \text{MAJ}(h_1, h_2, h_3).$$

How to sample?

Draw x from $EX(c, \mathcal{D})$; eval. $h_1(x), h_2(x)$
Keep x as $(\text{draw from } EX(c, \mathcal{D}_3))$ iff $h_1(x) \neq h_2(x)$,
else discard & try again.

Let $\eta := \mathcal{D}[h_1(x) \neq h_2(x)]$.

Takes, on avg, $\frac{1}{\eta}$ repetit. to get xy .

If η tiny: sad (inefficient!)...

but good news is, h_3 only affects
 h on η frac. of \mathcal{D} . So

if n tiny (as witnessed by many reps of draw from $EX(c, \mathcal{D})$ without getting $h_1 \neq h_2$)
 just take $h_3 \equiv \text{anything}$, $h' = \text{MAJ}(h_1, h_2, h_3)$

So can efficiently construct h' s.t.
 $|\text{err}_{\mathcal{D}}[h', c] - \text{err}_{\mathcal{D}}[h, c]| \leq .001$

→ "factor-1000 overhead"

We'll ignore this & pretend have bona fide h_1, h_2, h_3 from above procedure.

To show: $\mathcal{D}\{h(x) \neq c(x)\} = 35.2\%$

Divide X into 4 regions:

$$X = R_1 \cup \dots \cup R_4$$

R_1 : x s.t. $h_1(x) = h_2(x) = c(x)$	R_2 : $h_1(x) = c(x)$, $h_2(x) \neq c(x)$	R_3 : $h_1(x) = h_2(x) \neq c(x)$	R_4 : $h_1(x) \neq c(x)$, $h_2(x) = c(x)$
$\mathcal{D}_2(R_1) = .5 - p$	$\mathcal{D}_2(R_2) = p$	$\mathcal{D}_2(R_3) = .4 - p$	$\mathcal{D}_2(R_4) = .1 + p$
$\mathcal{D}_1(R_1) = \frac{6}{5}(.5 - p)$	$\mathcal{D}_1(R_2) = \frac{6}{5}p$	$\mathcal{D}_1(R_3) = \frac{4}{5}(.4 - p)$	$\mathcal{D}_1(R_4) = \frac{4}{5}(.1 + p)$

Let's define $p := \mathcal{D}_2[R_2]$. Set things up s.t. $\mathcal{D}_2[h_2(x) \neq c(x)] = .4$

so means $\mathcal{D}_2(R_3) = 0.4 - p$

We set things up s.t. $\mathcal{D}_2[h_1(x) \neq c(x)] = .5$,

$\mathcal{D}_2[h_1(x) = c(x)] = .5$, so

$$\mathcal{D}_2(R_1) = .5 - p \quad + \quad \mathcal{D}_2(R_4) = .1 + p$$

Recalling how \mathcal{D}_2 was obtained from \mathcal{D}_1 , get green values

$(\mathcal{D}_3$: all wt is on R_2, R_4)

h correct on $R_2 \cup R_4$
exactly when h_3 correct

h is correct : 100% on R_1 ,
0% on R_3 ,
60% on $R_2 \cup R_4$.

Rewrite:

$$h \text{ 's error} = \mathcal{D}_1(R_3) + \frac{4}{10} (\mathcal{D}_1(R_2) + \mathcal{D}_1(R_4))$$

$$So \Pr_{x \sim \mathcal{D}_1} [h(x) \neq c(x)] =$$

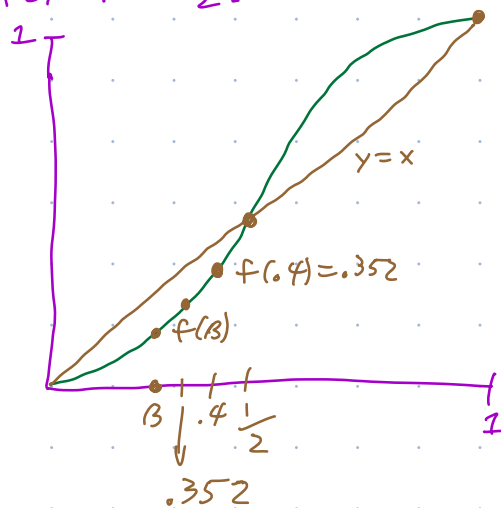
$$= \frac{4}{5}(.4 - p) + \frac{4}{10} \cdot \left(\frac{6}{5}p + \frac{4}{5}(.1 + p) \right)$$

$$= \frac{1.6}{5} + \frac{1.6}{50} + \left(-\frac{40}{50}p + \frac{24}{50}p + \frac{16}{50}p \right)$$

$$= \frac{32}{100} + \frac{3.2}{100} + 0 = 35.2\%$$

More generally: suppose weak learner only promised to have error $B < \frac{1}{2}$.

Above analysis \Rightarrow hyp h has error $\leq 3B^2 - 2B^3$;
 $< B$ for $B < \frac{1}{2}$.



Boosting - by - Filtering ^{← (what we did)} vs

Boosting over a Fixed Sample.

→ Inefficient: must discard precious data to do it !!

→ In practice: $y^i = c(x^i)$

• Draw $S = (x^1, y^1), \dots, (x^m, y^m)$
Only ever use these.

• The dist's we use: all over

• Initial dist. put wt $\frac{1}{m}$ on each data pt (\mathcal{D}).
New dist. \mathcal{D}' : new assignment of wt to points.

$$\sum_{i=1}^m \mathcal{D}'(x^i) = 1$$

• weak learning guarantee: for any \mathcal{D} (say $w_i = \mathcal{D}(x^i)$),
if run weak learner using w_1, \dots, w_m as wts,
get h s.t.

$$\sum_{i: h(x^i) \neq c(x^i)} w_i \leq \frac{1}{2} - \gamma.$$

• Goal of booster: construct a final h s.t.

$$\sum_{i: h(x^i) \neq c(x^i)} \frac{1}{m} \leq \epsilon. \quad (\text{typ. go for } 0 \text{ error on})$$

data set.)

Ada Boost: simple, powerful, adaptive
alg. to boost over a sample.

Idea of Ada Boost:

- given h_t, \mathcal{D}_t : construct \mathcal{D}_{t+1} by reweighting \mathcal{D}_t s.t. h_t has 50% acc. under \mathcal{D}_{t+1} . (scale regions \uparrow \downarrow like \mathcal{D}_2 did).

- At end: combines h_1, h_2, \dots using weighted maj. vote. Hyp's that did better: more weight. ("adaptivity")

Let's see it... 😊

Next time: analyze Ada Boost.
