

Weak vs. Strong Learning

Def. Algorithm A weakly "PAC learns" concept class C

if $\forall c \in C$ + \forall dists \mathcal{D} $\exists \delta > 0$
 $\forall \epsilon, \delta > 0$ ($\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't affect)

with prob $\geq 1 - \delta$
 given examples of c

A outputs h s.t. $\Pr_{\mathcal{D}} [h(x) \neq c(x)] \leq \frac{1}{2} - \frac{\delta}{2}$
 \uparrow
advantage

It was conjectured that distribution free weak learning
 was really weaker but surprise!
 Can "boost" a weak learner

Thm if C can be weakly learned on any dist \mathcal{D} then C can be
 (strongly) learned. \Leftarrow i.e. $\forall \epsilon, \mathcal{D}_0$

Will prove for case of $\mathcal{D}_0 = \mathcal{U}$

Applications

1) "Theoretical"

- Unit dist Algorithms for poly term DNF
weight w - poly threshold fctns

} low degree
alg doesn't
work well

(Boosting + KM)

- Ave case vs. worst case complexity

2) practical - Boosting Freund-Schapire

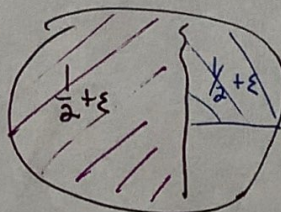
Good & Bad Ideas

- 1) simulate weak learner several times on
same distribution & take majority answer
-or-
best answer

gives better confidence

but doesn't reduce error, what if always get same answer?

- 2) filter out examples on which current hypothesis
does well & run weak learner on part where you
do badly.



Problem: given a new
example, how do you
know which section it
is in?

- 3) Keep some samples on which you are ok in filtering
 always use majority vote on all previous hypotheses
 to predict value of new samples

history: Schapire, Freund-Schapire, Impagliazzo-Servedio, Klivans

Filtering Procedures

- decide which samples to keep, which to throw out
- samples on which so far you guess correctly \leftarrow need for checking future hypotheses
 incorrectly \leftarrow need to improve on these

The setting

- Given labeled examples
 $(x_1, f(x_1)), (x_2, f(x_2)), \dots$

$$x_i \in \mathcal{X}$$

target function

$$\rightarrow f \in \mathcal{C}$$

- Given weak learning alg WL which weakly learns (advantage $\frac{\epsilon}{2}$) on any dist \mathcal{D}
 or equivalently error $\frac{1-\epsilon}{2}$
 $= \beta$ (use $\text{error}_{\mathcal{D}}(h) \equiv \Pr_{x \in \mathcal{D}}[f(x) \neq h(x)]$)

Plan:

1. simple "modest" accuracy boosting procedure
(error slightly improves)
2. recursively use ① to drive down error

Part I. Modest Improvement:

Algorithm: Given oracle to f , \mathcal{D} + weak learning alg WL

$\rightarrow h_1 \leftarrow$ run WL on \mathcal{D} for f

how also have oracle for h_1

Create example oracle \mathcal{D}_2 : how many calls to \mathcal{D} per output sample of \mathcal{D}_2 ?

"normalize" \mathcal{D} to make \mathcal{D} err half the time

flip coin:

Heads - draw examples from \mathcal{D} until find x s.t. $h_1(x) = f(x)$ " h_1 correct"

Output x

Tails - " " " " " " " " " h_1 incorrect"

Output x

so $err_{\mathcal{D}_2}(h_1) = 1/2$

note $err_{\mathcal{D}_2}(h_1) < 1/2$ $\rightarrow h_2 \leftarrow$ run WL on \mathcal{D}_2 for f
so $h_2 \neq h_1$

Create example oracle \mathcal{D}_3 : how many calls to \mathcal{D} per sample of \mathcal{D}_3 ?
draw examples from \mathcal{D} until find x s.t. $h_1(x) \neq h_2(x)$

Output x

$h_3 \leftarrow$ run WL on \mathcal{D}_3 for f

output $h \equiv \text{maj}(h_1, h_2, h_3)$

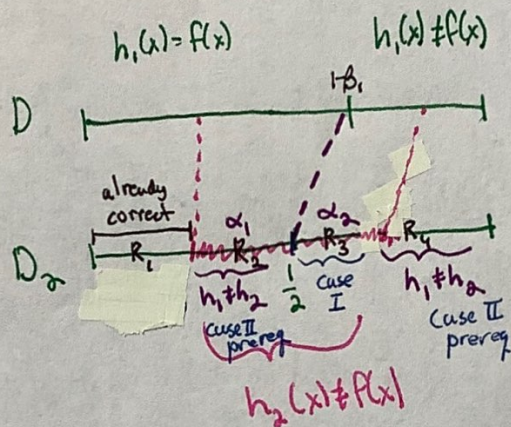
(ie. h is a concept which evaluates h_1, h_2, h_3 + outputs majority answer)

Error Analysis of "Modest Improvement": let $\beta_1 = \Pr_{D_1} [h_1(x) \neq f(x)]$ (3 error probabilities)
 $\beta_2 = \Pr_{D_2} [h_2(x) \neq f(x)]$ $\beta_3 = \Pr_{D_3} [h_3(x) \neq f(x)]$

Observation if $h_1(x) = f(x)$ then $D(x) = 2(1 - \beta_1) D_2(x)$
 " $h_1(x) \neq f(x)$ " $D(x) = 2\beta_1 D_2(x)$

why?

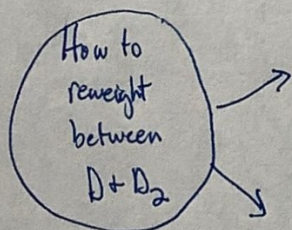
sort x 's
st.
picture
works



R_1 : $h_1 = h_2 = f$ maj correct
 R_2 : $h_1 = f \neq h_2$ h_3 helps?
 R_3 : $h_1 = h_2 \neq f$ maj incorrect
 R_4 : $h_1 \neq h_2 = f$ h_3 helps?

h_3 is "good" on $R_2 + R_4$

total wt of x st. $h_1(x) = f(x)$ goes from $1 - \beta_1$ to $\frac{1}{2}$
 + relative wts of such x 's stay same



$$\sum_{x \text{ st. } h_1(x) = f(x)} D(x) = 1 - \beta_1$$

$$\sum_{x \text{ st. } h_1(x) \neq f(x)} (D(x) \cdot \alpha) = \frac{1}{2}$$

$$1 - \beta_1 = \frac{1}{2\alpha}$$

$$\text{so } D_2(x) = D(x) \cdot \alpha = \frac{1}{2(1-\beta_1)} \cdot D(x)$$

$$+ D(x) = 2(1-\beta_1) D_2(x)$$

for x st $h_1(x) \neq f(x)$

$$\beta_1 = \frac{1}{2\alpha'}$$

$$\text{so } D_2(x) = \alpha' \cdot D(x) = \frac{1}{2\beta_1} D(x)$$

$$+ D(x) = 2\beta_1 D_2(x)$$

More general observation:

$$\Pr_{x \in S} [x \in S] = 2(1 - \beta_1) \Pr_{x \in D_2} [h_1(x) = f(x) \wedge x \in S] + 2\beta_1 \Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge x \in S] \quad (*)$$

Let β be error guaranteed by assumption on output of WL
 Let $g(\beta) = 3\beta^2 - 2\beta^3$

Main Lemma

$$\text{err}_D(h) \leq g(\beta)$$

note: $g(\beta) < \beta$ but by how much?
 $g(1/2) = 1/2$

Pf.

$\text{err}_D(h)$ comes from 2 types:

1) x st. $h_1(x) = h_2(x) \neq f(x)$ (both h_1 & h_2 are wrong)
 so no matter if $h_3(x)$ is wrong or right,
 h (which is majority) will be wrong
 of h_1, h_2, h_3

2) x st. $h_1(x) \neq h_2(x)$
 so $h_3(x)$ decides if h is correct

$$\begin{aligned} \text{so } \text{err}_D(h) &= \Pr_{x \in D} [h_1(x) \neq f(x) + h_2(x) \neq f(x)] \\ &\quad + \Pr_{x \in D} [h_3(x) \neq f(x) \mid h_1(x) \neq h_2(x)] \cdot \Pr_{x \in D} [h_1(x) \neq h_2(x)] \\ &= \Pr_{x \in D} [h_1(x) \neq f(x) + h_3(x) \neq f(x)] \\ &\quad + \underbrace{\beta_3}_{\text{def of } \beta_3} \cdot \Pr_{x \in D} [h_1(x) \neq h_2(x)] \\ &\leq \Pr_{x \in D} [h_1(x) \neq f(x) + h_2(x) \neq f(x)] + \beta \cdot \Pr_{x \in D} [h_1(x) \neq h_2(x)] \end{aligned} \quad (0)$$

Finish 2nd term:

partition D_2 into 2 equal parts:

- 1) x st. $h_1(x) = f(x)$
- 2) " " " \neq "

$$\text{let } \alpha_1 = \text{"err von part 1"} \stackrel{\text{of } h_2 \text{ (wrt } D_2)}{=} \Pr_{x \in D_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$$

$$\alpha_2 = \text{" " " 2"} \stackrel{\text{of } h_2 \text{ (wrt } D_2)}{=} \Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge h_2(x) \neq f(x)]$$

$$\alpha_1 + \alpha_2 = \beta_2$$

$$\text{" } x \in S \text{"} \stackrel{\text{of } h_2 \text{ (wrt } D_2)}{=} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$$

$$\text{Then } \Pr_{x \in D_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$$

$$= 2(1-\beta_1) \cdot \underbrace{\Pr_{x \in D_2} [h_1(x) = f(x) \wedge h_1(x) = f(x) \wedge h_2(x) \neq f(x)]}_{=\alpha_1}$$

$$+ 2\beta_1 \cdot \underbrace{\Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge h_1(x) = f(x) \wedge h_2(x) \neq f(x)]}_0$$

$$= 2 \cdot (1-\beta_1) \cdot \alpha_1$$

$$\text{And } \Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)] = \frac{1}{2} - \alpha_2 \quad \leftarrow \text{algebra + structure of } D_2$$

$$\text{so } \Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)] = 2\beta_1 \left(\frac{1}{2} - \alpha_2 \right)$$

Putting together

$$\Pr_{x \in D_2} [h_1(x) \neq h_2(x)] = \Pr_{x \in D_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$$

$$+ \Pr_{x \in D_2} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)]$$

$$= 2(1-\beta_1)\alpha_1 + 2\beta_1 \left(\frac{1}{2} - \alpha_2 \right) \quad (1)$$

2nd term:

Also

$$\Pr_{x \in D} [h_1(x) \neq f(x) + h_2(x) \neq f(x)] = 2\beta_1 \alpha_2 \tag{2}$$

So putting (1) + (2) into (0):

$$\begin{aligned} \text{err}_0(h) &\leq 2\beta_1 \alpha_2 + \beta \cdot [2(1-\beta_1)\alpha_1 + 2\beta_1(\frac{1}{2}-\alpha_2)] \\ &= \beta^2(1-2\gamma_1) + 2\beta\gamma_2(1-\beta) + 2\gamma_1\beta \quad \text{assume } \beta_1 = \beta \\ &\quad \vdots \\ &\leq \beta^2 + 2\beta(1-\beta)(\gamma_1 + \gamma_2) \quad \text{assume } \beta_2 = \beta \\ &\leq 3\beta^2 - 2\beta^3 \end{aligned}$$

Part II Recursive Accuracy Boosting

one application of procedure takes error $\beta \rightarrow \leq 3\beta^2 - 2\beta^3$

but we want error \leq tiny

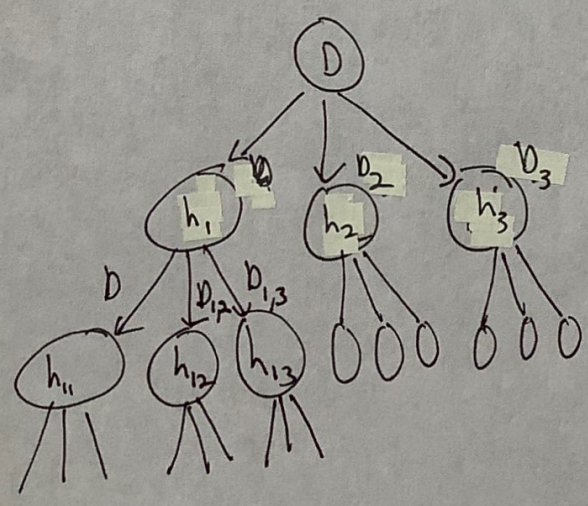
idea: recursion

Strong-learn Algorithm:
(first try)

given p, D' return result of WL on D'
if $p \geq \frac{1}{2} - \epsilon$ (error required from level below to get error $\leq \epsilon$ here)
 $\beta \leftarrow g^{-1}(p)$

get $D'_2 \text{ \& } D'_3$ as in "modest boost"

- $h_1 \leftarrow \text{strong-learn}(\beta, \text{Ex}(f, D'))$
- $h_2 \leftarrow \text{strong-learn}(\beta, \text{Ex}(f, D'_2))$
- $h_3 \leftarrow \text{strong-learn}(\beta, \text{Ex}(f, D'_3))$
- $h \leftarrow \text{maj}(h_1, h_2, h_3)$
- return h



Sample Complexity:

Main issues:

- how many recursive calls?
depth + size of recursion tree
- how many samples required to construct filters?

Depth of recursion: (assume wk advantage γ is $\frac{1}{2} \Rightarrow$ error $\beta \leq \frac{1}{4}$
but also works for $\gamma = \Omega(\frac{1}{n^c})$
i.e. $\beta \sim \frac{1}{2} = \frac{1}{n^c}$)

Claim if $\beta \leq \frac{1}{4}$, $g(\beta) \leq 3\beta^2 = \frac{1}{3}(3\beta)^2$
error decreases in k steps $\leq \frac{1}{3}(3\beta)^{2^k} \leq (\frac{3}{4})^{2^k}$

- $\Rightarrow k = \Theta(\log \log (1/\epsilon))$ depth suffices
- \Rightarrow size $3^{\log \log (1/\epsilon)} \approx \Theta(\log (1/\epsilon))$ suffices

samples?

Problem: filtering could take a while

good news!
if filtering takes a while than we've already learned.
if can't find samples st. $h_1(x) = f(x)$
output \vec{h}_1 (this can't happen though)
" " " " " " $h_1(x) \neq f(x)$
can halt + output h_1 for final answer
" " " " " " $h_1(x) \neq h_2(x)$
don't calculate h_3