

Today:

Weak vs. Strong Learning

via "Boosting"

Def. Algorithm A weakly-PAC learns

concept class C if $\forall c \in C$

\forall dists \mathcal{D}

$\forall \epsilon, \delta > 0$

with prob $\geq 1 - \delta$

given labelled examples of $C \leftarrow (x, c(x))$

from dist \mathcal{D}

↓

A outputs h st. $\Pr_{x \in \mathcal{D}} [h(x) \neq c(x)] \leq \epsilon$

$\frac{1}{2} - \frac{\gamma}{2}$

γ is "advantage"

is weak learning easier than strong-
learning
regular PAC learning?
surprisingly "no"!

Thm if \mathcal{C} can be weakly learned
on any dist \mathcal{D} then \mathcal{C} can be
(strongly) PAC-learned
 $\forall \epsilon$

Applications:

1) "Theoretical"

- Boosting + KM \rightarrow unif dist learning algs
for poly term DNF
(better than low deg alg)
- insights into average case + worst case
complexity

2) Practical -

many boosting algs

Freund Schapire ... (many many)

history:

Schapire

Freund Schapire

:

lots more

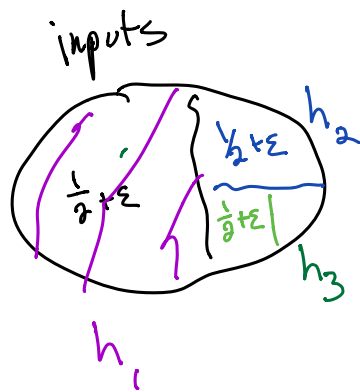
Good & Bad ideas

1) simulate weak learner many times on
same distribution & take $\begin{cases} \text{majority answer} \\ \text{best answer} \end{cases}$

\Rightarrow better confidence (δ)

but doesn't reduce error

2) **filter** out examples on which current hypothesis does well & rerun weak learner on part where we did badly



training phase:
 $x, f(x)$

later:
 given x , predict $f(x)$
 here we don't know
 if h_1, h_2 or h_3
 correct
 which section
 are we in?

- 3) keep some "good" samples in filters
 use majority vote on all hypothesis h_1, h_2, h_3
 (don't need to know "color" of x)

Filtering Procedures

- decide which samples to keep/throw
- need keep some samples on which you did well
 but weight more on those on which you need to improve

The setting:

- Given labelled examples $(x_1, f(x_1)) (x_2, f(x_2))$

$$\begin{array}{l} x_i \in \mathcal{X} \\ \text{target} \\ \text{fctn} \rightarrow f \in \mathcal{C} \end{array}$$

- Given weak learning alg WL which weakly learns $\left\{ \begin{array}{l} \text{advantage } \gamma \\ \text{error } \frac{1}{2} - \gamma/2 \end{array} \right\}$ on any dist \mathcal{D}'

$$\text{error}_{\mathcal{D}}(h) \equiv \Pr_{x \in \mathcal{D}} [f(x) \neq h(x)]$$

The plan:

- simple "modest" boosting procedure (error slightly improves)
- recursively use \uparrow to drive down error

Part I: Modest improvement

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

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

create an example oracle \mathcal{D}_2

flip coin:

normalize \mathcal{D}
st. errs
half the
time

$$\text{err}_{\mathcal{D}_2}(h_1) = \frac{1}{2}$$

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

output x

Tails - " " " " " "

" " $h_1(x) \neq f(x)$ " h_1 incorrect"

output x

how many samples?

H: ≤ 2
T: if $\geq \frac{1}{\epsilon}$
then done!

$\text{err}_{\mathcal{D}_2}(h_2) < \frac{1}{2} \rightarrow h_2 \leftarrow$ run WL on \mathcal{D}_2

so $h_1 \neq h_2$ create example oracle \mathcal{D}_3

draw examples until find x st. $h_1(x) \neq h_2(x)$
output x

if need too many samples ($> \frac{1}{\epsilon}$)
skip

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

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

ie. evaluate h_1, h_2, h_3 on x + output most common answer

Error analysis:

error $\leq \beta$ promised by WL

$$\beta_1 = \Pr_{\mathcal{D}} [h_1(x) \neq f(x)]$$

$$\beta_2 = \Pr_{\mathcal{D}} [h_2(x) \neq f(x)]$$

$$\beta_3 = \Pr_{\mathcal{D}} [h_3(x) \neq f(x)]$$

for simplicity
(worst case)

assume

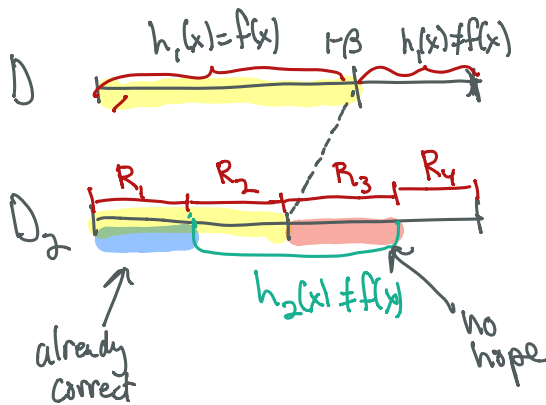
$$\beta = \beta_1 = \beta_2 = \beta_3$$

Observation

$$\text{If } h_1(x) = f(x) \text{ then } D(x) = 2(1 - \beta_1) D_2(x)$$

$$\neq \qquad \qquad \qquad = 2\beta_1 D_2(x)$$

why?



$R_1: h_1 = h_2 = f$ already maj correct
 $R_2: h_1 = f \neq h_2$
 $R_3: h_1 = h_2 \neq f$ already maj incorrect
 $R_4: h_1 \neq h_2 = f$
 $h_3 \geq \frac{1}{2} + \frac{\epsilon}{2}$ good in $R_2 + R_4$

total wt of x st. $h_1(x) = f(x)$ goes from

$1 - \beta$ to $\frac{1}{2}$

† rel wts of these x 's stay same

$$\sum_{\substack{x \text{ st} \\ h_1(x) = f(x)}} D(x) = 1 - \beta \quad \left. \vphantom{\sum} \right\} 1 - \beta = \frac{1}{2\alpha}$$

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

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

$$\dagger D(x) = 2(1-\beta) D_2(x)$$

$\beta =$ err of output of WL
guarantee

Main Lemma

$$\text{err}_\beta(h) \leq \underbrace{3\beta^2 - 2\beta^3}_{\equiv g(\beta)}$$

$$g(\beta) \ll \beta$$

idea of pf:

$\text{err}_0(h)$ has 2 types:

1) x s.t. $h_1(x) = h_2(x) \neq f(x)$ (both wrong)
so h_3 can't fix

2) x s.t. $h_1(x) \neq h_2(x)$

here h_3 decides if h correct

$$\begin{aligned} \text{(c) } \text{err}_0(h) &= \Pr_{x \in \mathcal{D}} [h_1(x) \neq f(x) \vee h_2(x) \neq f(x)] \\ &+ \underbrace{\Pr_{x \in \mathcal{D}} [h_3(x) \neq f(x) \mid h_1(x) \neq h_2(x)] \cdot \Pr_{x \in \mathcal{D}} [h_1(x) \neq h_2(x)]}_{\text{def of } \beta_3 = \text{error of WL on } \mathcal{D}_3} \end{aligned}$$

$$\alpha_1 = \Pr_{x \in \mathcal{D}_2} [h_1(x) = f(x) \vee h_2(x) \neq f(x)] \quad \text{region 2}$$

$$\alpha_2 = \Pr_{x \in \mathcal{D}_2} [h_1(x) \neq f(x) \vee h_2(x) = f(x)] \quad \text{region 3}$$

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

Then

$$\begin{aligned} & \Pr_{x \in \Omega} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)] \\ &= 2 \cdot (1 - \beta_1) \underbrace{\Pr_{x \in \Omega_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]}_{\alpha_1} \\ &= 2(1 - \beta_1)\alpha_1 \end{aligned}$$

$$\Pr_{x \in \Omega_2} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)] = \frac{1}{2} - \alpha_2$$

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

Putting together:

$$\Pr_{x \in \Omega} [h_1(x) \neq h_2(x)] = 2(1 - \beta_1)\alpha_1 + 2\beta_1 \left(\frac{1}{2} - \alpha_2\right)$$

$$\text{Also } \Pr_{x \in \Omega} [\text{both } h_1, h_2 \text{ wrong}] = \underbrace{2\beta_1}_{\substack{\text{fix} \\ D_2 \rightarrow D_1}} \underbrace{\alpha_2}_{\text{def}} \quad (2)$$

(1) + (2) put into previous

$$\text{err}_b(h) \leq 2\beta_1\alpha_2 + \beta [2(1 - \beta)\alpha_1 + 2\beta \left(\frac{1}{2} - \alpha_2\right)]$$

⋮

$$\leq 3\beta^2 - 2\beta^3 \quad \square$$

$$g(y) = 3y^2 - 2y^3$$

Part II Recursive Accuracy Boosting

Strong learning Algorithm:

given ρ, D'

if ρ can be achieved directly from WL
then just call WL & return result

$\beta \leftarrow g^{-1}(\rho)$ error from level below
required to get error ρ

{ get $D'_2 + D'_3$ as in "modest boost" }

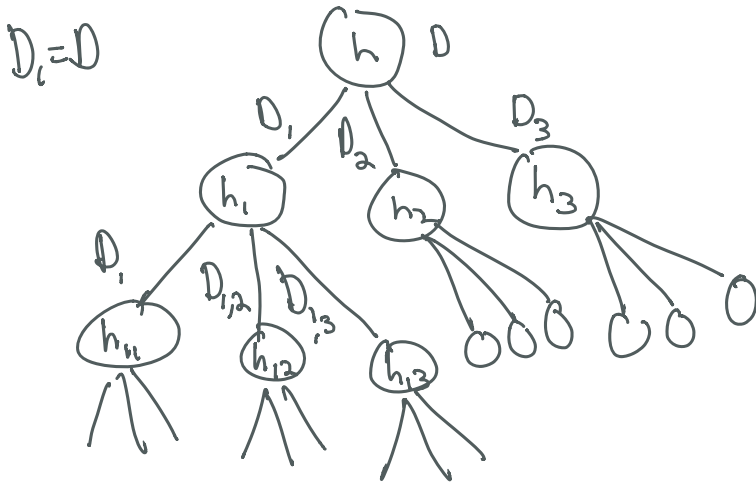
$h_1 \leftarrow \text{strong learn}(\beta, E_x(f, D'))$

$h_2 \leftarrow \text{strong learn}(\beta, E_x(f, D'_2))$

$h_3 \leftarrow \text{strong learn}(\beta, E_x(f, D'_3))$

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

return h



Sample Complexity:

how many recursive calls?

depth + size of recursion tree

how many samples to construct filters?

depth of recursion:

also does improve as long as $\beta \approx \frac{1}{2} - \frac{0(\epsilon)}{n^c}$ gives poly depth

if $\beta \leq 1/4$ then $g(\beta) = 3\beta^2 - 2\beta^3 \leq 3\beta^2$

in k steps $\leq \frac{1}{3} (3\beta)^{2^k} \leq \left(\frac{3}{4}\right)^{2^k}$

$\Rightarrow k = \Theta(\log \log \frac{1}{\epsilon})$ depth suffices

size is $O(3^{\log \log \frac{1}{\epsilon}}) \approx \Theta(\log \frac{1}{\epsilon})$