

Lecture 22

distribution-free weak learning } "boosting"
⇒ strong learning

average vs. worst case complexity

Weak vs. Strong Learning

Def. Algorithm A "weakly PAC learns" concept class \mathcal{C} if $\exists \epsilon > 0$

st. $\forall c \in \mathcal{C} + \forall$ dists \mathcal{D}

$\forall \delta > 0$

$\leftarrow (\delta = \frac{1}{4}$ or $\frac{1}{n^2}$ doesn't affect)

with prob $\geq 1 - \delta$

given examples of c

A outputs h s.t. $\Pr_{\mathcal{D}} [h(x) = c(x)] \geq \frac{1}{2} + \frac{\epsilon}{2}$

not good
compared
to
 $1 - \epsilon$ or 99%

↑
advantage
over
guessing

It was first conjectured that weak learning is easier than strong (i.e. \exists fctns that can weakly learn but not strongly learn)

Surprise!!

Can "boost" a weak learner

Thm if \mathcal{C} can be weakly learned on

any dist \mathcal{D} then \mathcal{C} can be

(strongly) learned

ie. $\forall \epsilon$

dependence on γ ?

δ ?

ϵ ?

Applications:

1) "theoretical"

- uniform distribution algorithms for poly term DNF weight- w poly threshold fctns (Boosting + KM)

low degree alg doesn't work well

- Ave case vs. worst case complexity

2) practical: "Boosting"

Freund-Schapire

Good & Bad Ideas

1) simulate weak learner several times
on same distribution & take

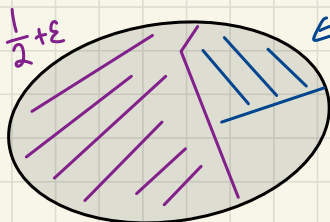
majority answer

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



$\leftarrow \frac{1}{2} + \epsilon$ of non-purple

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

3) Keep some samples on which you are ok in your filtering.

Always use majority vote on previous hypotheses to predict value of new samples.

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

Filtering Procedures:

- decide which samples to keep vs. throw out
- samples on which you guess

Correctly: needed for checking future hypotheses

incorrectly: needed for improvement

The setting

- Given labelled examples

$$(x_1, f(x_1)) \quad (x_2, f(x_2)) \quad \dots$$

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

$$f \in \mathcal{C} \quad \text{"target function"}$$

- Given weak learning alg W which weakly learns (advantage $\frac{\epsilon}{2}$) on any dist \mathcal{D}

$$\text{error} \quad \frac{1}{2} - \frac{\epsilon}{2}$$

$= \beta$

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

Plan:

1. simple "modest" accuracy boosting procedure
2. recursively use ① to drive down error

Part I: Modest Improvement

Given: oracle to f
example oracle \mathcal{D}
weak learning algorithm WL

Algorithm:

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

note: \rightarrow
now also
have
oracle
for h_1 !

create example oracle \mathcal{D}_2 :

Question - how many samples
of \mathcal{D} needed per
output sample of \mathcal{D}_2 ?

flip coin:

heads - draw examples from \mathcal{D}
until find x s.t.

$$h_1(x) = f(x)$$

output x

" h_1 correct"

tails - draw examples from \mathcal{D}
until find x s.t.

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

output x

" h_1 incorrect"

"normalize"

\mathcal{D}

to make h_1
err half
the time

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

(Algorithm conti.)

note
 $\text{err}(h_2) < \frac{1}{2}$
 \mathcal{D}_2
so $h_1 \neq h_2$

$h_2 \leftarrow$ run WL on \mathcal{D}_2 for f

Create example oracle \mathcal{D}_3 :

draw examples from \mathcal{D} until find
 x st. $h_1(x) \neq h_2(x)$

output x

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

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

on x , evaluate $h_1(x), h_2(x), h_3(x)$
& output majority answer

Error Analysis of "Modest Improvement"

define 3 error probabilities:

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

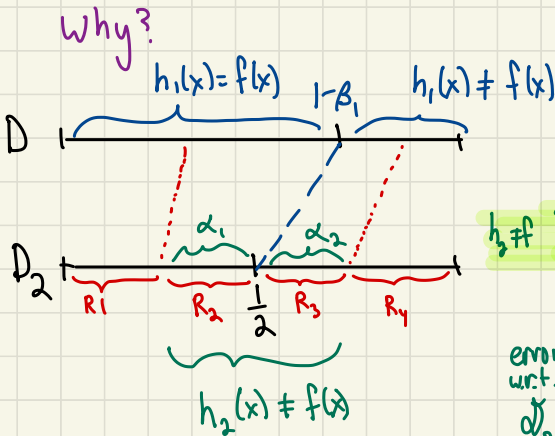
$$\beta_2 = \Pr_{\mathcal{D}_2} [h_2(x) \neq f(x)]$$

$$\beta_3 = \Pr_{\mathcal{D}_3} [h_3(x) \neq f(x)]$$

Observation:

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

" " \neq " " $D(x) = 2\beta_1 D_2(x)$



4 regions:

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

h_3 "good" in R_2, R_4

(Sort x 's so that picture works)

On reweighting between $D + D_2$: (proof of observation)

for x s.t. $h_1(x) = f(x)$:

Total wt of x s.t. $h_1(x) = f(x)$

goes from $1 - \beta_1$ to $\frac{1}{2}$

\Rightarrow relative wts of x 's stays same

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

$$\sum_{\substack{x \text{ s.t.} \\ h_1(x) = f(x)}} \underbrace{D(x) \cdot \alpha}_{D_2(x)} = \frac{1}{2}$$

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

$$\begin{aligned} \text{so, } D_2(x) &= D(x) \cdot \alpha \\ &= \frac{1}{2 \cdot (1 - \beta_1)} \cdot D(x) \end{aligned}$$

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

for x s.t. $h_1(x) \neq f(x)$:

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

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

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

More general observation:

$$\forall S, \Pr_{x \in \mathcal{D}} [x \in S] = 2(1 - \beta_1) \cdot \Pr_{x \in \mathcal{D}_2} [h_1(x) = f(x) \wedge x \in S] \\ + 2\beta_1 \cdot \Pr_{x \in \mathcal{D}_2} [h_1(x) \neq f(x) \wedge x \in S] \quad (*)$$

Bounding error of WL:

β ← error guarantee (by assumption) on WL's output

$$g(\beta) \leftarrow 3\beta^2 - 2\beta^3$$

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

note: $g(\beta) \leq \beta$ but how much better?

not always better since $g(\frac{1}{2}) = \frac{1}{2}$

Proof $\text{err}_\beta(h)$ from 2 types:

Type ① x st. $h_1(x) = h_2(x) \neq f(x)$ (both h_1, h_2 wrong)

"lost
cause"
case

so no matter whether $h_3(x)$ wrong or right
 $h = \text{maj}(h_1, h_2, h_3)$ will be wrong on x

Type ② x st. $h_1(x) \neq h_2(x)$

here $h_3(x)$ determines if h correct

defined
on
next
page
↓

$$\text{so } \text{err}_\beta(h) = \Pr_{x \in D} [h_1(x) \neq f(x) + h_2(x) \neq f(x)] \quad \left. \vphantom{\Pr_{x \in D}} \right\} 2\beta_1, \alpha_2$$
$$+ \Pr_{x \in D} [h_3(x) \neq f(x) \mid h_1(x) \neq h_2(x)] \cdot \Pr_{x \in D} [h_1(x) \neq h_2(x)] \quad \left. \vphantom{\Pr_{x \in D}} \right\} \text{def of } \beta_3$$

$$\leq \Pr_{x \in D} [h_1(x) \neq f(x) + h_2(x) \neq f(x)]$$

$$+ \beta \cdot \Pr_{x \in D} [h_1(x) \neq h_2(x)]$$

$$\left. \vphantom{\Pr_{x \in D}} \right\} \beta_3 \leq \beta$$

equation
(0)

Type ② - Calculating $\Pr_{\mathcal{D}}[h_1(x) \neq h_2(x)]$:

Partition \mathcal{D}_2 into 2 parts:

- 1) x st. $h_1(x) = f(x)$ regions R_1, R_2
 2) " " $h_1(x) \neq f(x)$ regions R_3, R_4

R_2 : $\alpha_1 = \text{err of } h_2 \text{ wrt } \mathcal{D}_2 \text{ on part 1} = \Pr_{x \in \mathcal{D}_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$

R_3 : $\alpha_2 = \text{ " " " " " " " " } 2 = \Pr_{x \in \mathcal{D}_2} [h_1(x) \neq f(x) \wedge h_2(x) \neq f(x)]$

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

R_2 :

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

$$= 2 \cdot (1 - \beta_1) \Pr_{x \in \mathcal{D}_2} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)]$$

α_1

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

use reweighting calculations from before

need to reweight these two cases differently

And $\Pr_{x \in \mathcal{D}_2} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)] = \frac{1}{2} - \alpha_2$ ← defn of α_2

$\Pr = \frac{1}{2}$
by construction
of \mathcal{D}_2

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

Putting together:

$$\begin{aligned} \Pr_{x \in \mathcal{D}} [h_1(x) \neq h_2(x)] &= \\ & \Pr_{x \in \mathcal{D}} [h_1(x) = f(x) \wedge h_2(x) \neq f(x)] \\ & + \Pr_{x \in \mathcal{D}} [h_1(x) \neq f(x) \wedge h_2(x) = f(x)] \\ & = 2(1-\beta_1)\alpha_1 + 2\beta_1(\frac{1}{2} - \alpha_2) \end{aligned}$$

Finally:

$$\text{err}_\beta(h) \leq 2\beta_1\alpha_2 + \beta(2(1-\beta_1)\alpha_1 + 2\beta_1(\frac{1}{2} - \alpha_2))$$

⋮

assume
 $\beta = \beta_1 = \beta_2$
use $\alpha_1 + \alpha_2 = \beta_2$

$$\leq \beta^2 + 2\beta(1-\beta)(\alpha_1 + \alpha_2) \leq 3\beta^2 - 2\beta^3$$



Part II Recursive accuracy boosting

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

we want tiny error

main idea: Recursion

Algorithm: given ρ, D'

if $\rho \geq$ promised error of WL, return result of WL on D'

else:

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

define $D'_2 + D'_3$ as in "modest boost"

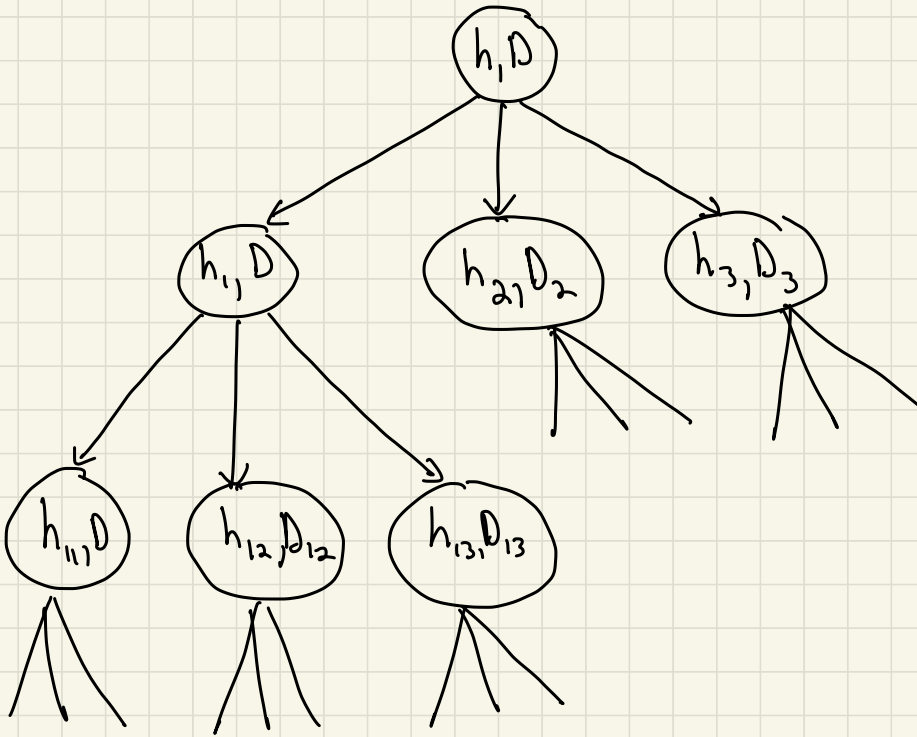
$h_1 \leftarrow \text{strong learn}(\beta, Ex(f, D'))$

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

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

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

return h



issues:

- how many recursive calls?

depth & size of recursion tree

- how many samples to construct filtered distributions?

samples:

problem ... filtering can take a while

but good news!

if it takes a while to find good samples,
then we've already learned well!

e.g. to find samples s.t. $h_1(x) = f(x)$

more than half should satisfy
(in any case if few such samples,
can output T_{h_1} to get good
approx. to f)

to find samples s.t. $h_1(x) \neq f(x)$:

if can't find, then h_1 is good
approx to f

to find samples s.t. $h_1(x) \neq h_2(x)$:

if always agree, then don't need h_3

depth of recursion:

assume WL advantage δ is $\frac{1}{2} \Rightarrow \beta \leq \frac{1}{4}$

(but also works if $\delta = \Omega(\frac{1}{n^c})$)

Claim if $\beta \leq \frac{1}{4}$, $g(\beta) \leq 3\beta^2 = \frac{1}{3}(3\beta)^2$

error decrease in depth k is down to

$$\leq \frac{1}{3}(3\beta)^{2^k}$$

Important Consequences:

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

\Rightarrow size $3^{\log \log \frac{1}{\epsilon}} \approx \Theta(\log \frac{1}{\epsilon})$
 $\times S$ \leftarrow description of weak learning hypothesis

$$= O(S \log \frac{1}{\epsilon})$$

suffices to describe circuit

\mathcal{C} is concept class where s bound size of concepts

Thm \mathcal{C} learnable $\Rightarrow \exists$ efficient algorithm

• using $\frac{\text{poly}(n, s, \log \frac{1}{\epsilon}, \log \frac{1}{\delta})}{\epsilon}$ samples & time

• Outputs hypotheses of size

$$\text{poly}(n, s, \log \frac{1}{\epsilon})$$

($\&$ can evaluate in poly time)

Why? given \mathcal{A} learning \mathcal{C}

use \mathcal{A} with $\epsilon_0 = 1/4$

boost \mathcal{A} to arbitrary ϵ

Corr \mathcal{C} learnable \Rightarrow all concepts $c \in \mathcal{C}$
have poly sized ckt

Pf idea

$\forall c \in \mathcal{C}$, use $\delta < \frac{1}{2^n}$ (e.g. no error)

will output consistent hypothesis of

poly size in $n + |c|$

\leftarrow # bits to describe c

that is poly time
evaluable.

\Rightarrow poly sized ckt.



Thm. Suppose f cannot be computed by poly sized ckt's. Then there is a sequence of distributions $\{\mathcal{D}_n^*\}_{n=1}^\infty$ st. f is "average-case hard" on $\{\mathcal{D}_n^*\}_{n=1}^\infty$

↑
 no poly sized ckt gets f right more than $\frac{1}{2} + \frac{1}{\text{poly}(n)}$ of time

↑
 need hard dist for each input size to make it well defined

Pf idea

if not, f can be weakly-learned by poly sized ckt's

\Rightarrow can strongly learn in size $\text{poly}(\log \frac{1}{\epsilon}, \dots)$
 \Rightarrow can learn f with 0 error
 \Rightarrow f computable by poly-sized ckt's