

Lecture 21

finish distribution-free weak learning
 \Rightarrow strong learning

} "boosting"

average vs. worst case complexity

Last time:

Weak vs. Strong Learning

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

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 (labelled)

in "strong" learning this is $1 - \epsilon$

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

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

↑ advantage over guessing

Assume $c \in \mathcal{C}$ maps $X \rightarrow \{0, 1\}$
← domain

\mathcal{C} is concept class where s bounds size of WL-description of concepts
← parametrized by "size" n of inputs
← fctn of n

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

Strong learner

- using $\frac{\text{poly}(n, s, \log \frac{1}{\epsilon}, \log \frac{1}{\delta})}{\epsilon}$ samples & time
- Outputs hypotheses of size $\text{poly}(n, s, \log \frac{1}{\epsilon})$ (useful later)

(\star can evaluate in poly time)

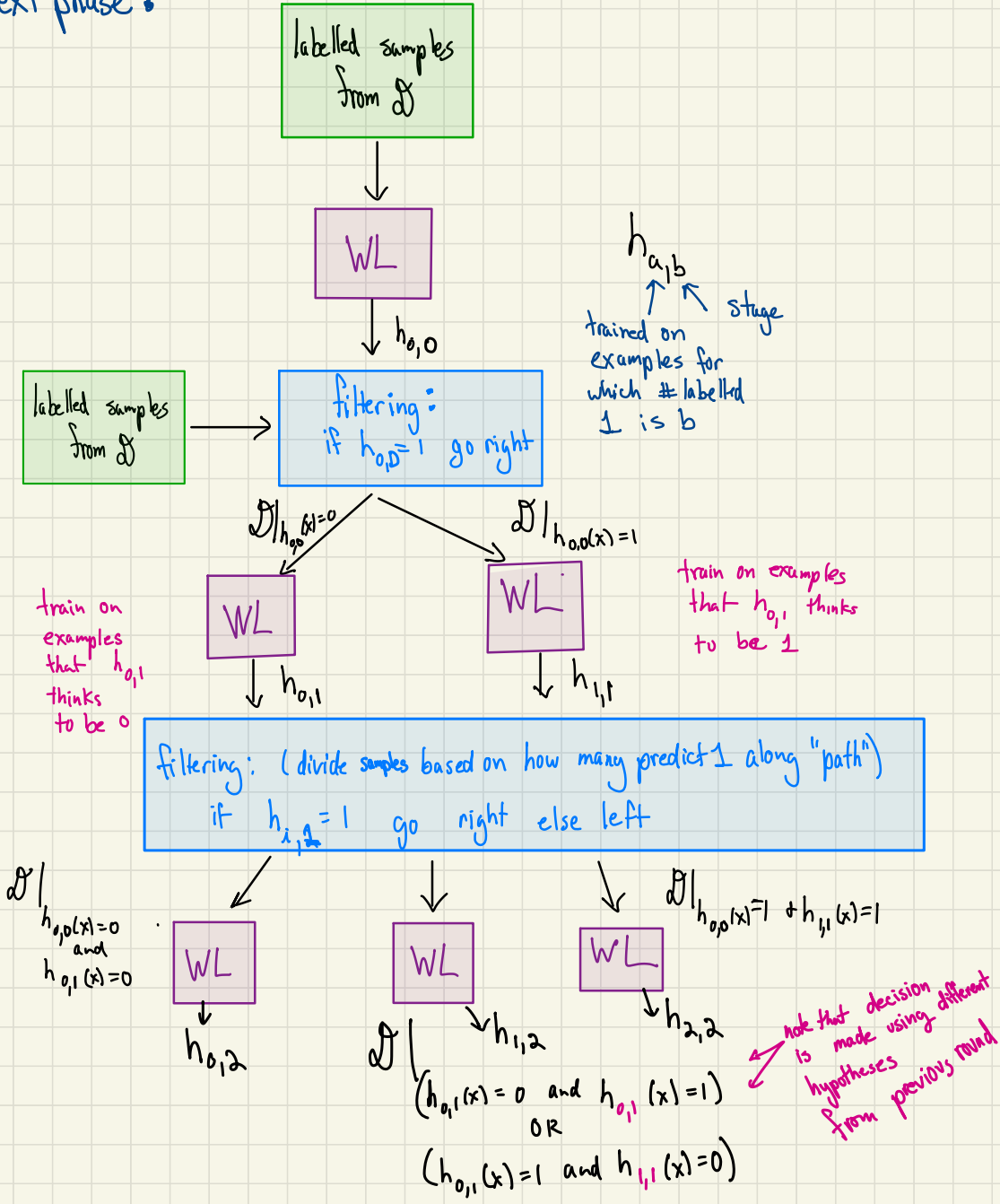
Why? given \mathcal{A} learning \mathcal{C}
use \mathcal{A} with $\epsilon_0 = 1/4$
boost \mathcal{A} to arbitrary ϵ

Idea: Filter by hypothesis value

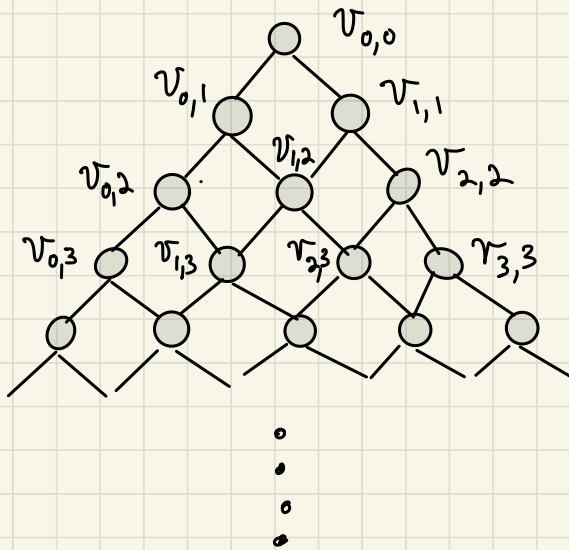
[Long Serodio]
"Martí Boost"

next phase:

assume $C: X \rightarrow \{0,1\}$

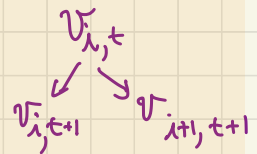


Putting together the filters:

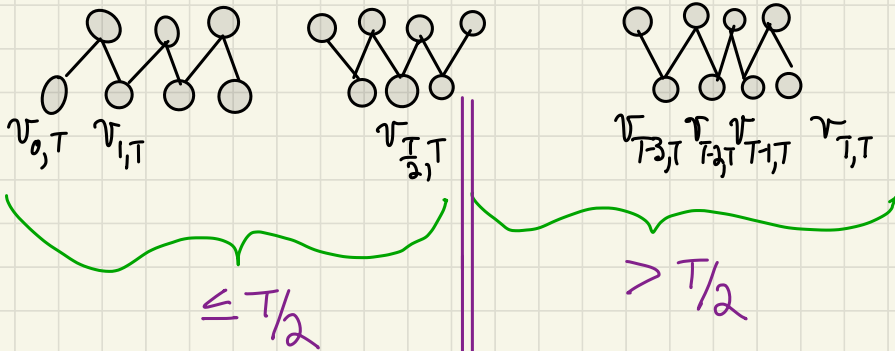


T stages
 $T+1$ layers

layer $0 \leq t \leq T$ has t nodes
 each node in layer $\leq T-1$ has 2 outgoing edges:



nodes in layer T have no outgoing edges



for new unlabelled sample, if it ends up on this side output "0"

else, output "1"

"Majority":

Node v_{ij} labelled by hypothesis h_{ij}

Routing of examples: *← ignores labels of examples*

given x , at node $v_{i,j}$

• evaluate $h_{i,j}(x)$ to get $y \in \{0,1\}$

• send on edge labelled by y , i.e. $v_{i+y, j+1}$

eventually reaches $v_{l,t}$

↑ level
of hypotheses evaluating to 1 on the way

↑ next level
Left: $y=0$
Right: $y=1$

At t^{th} stage:

branching program induces $t+1$ different distributions

$\mathcal{D}_{0,t} \mathcal{D}_{1,t} \mathcal{D}_{2,t} \dots \mathcal{D}_{t,t}$

where $\mathcal{D}_{j,t} = \mathcal{D}|_{\text{reach } v_{j,t}}$

run WL on each to get $h_{i,t} \forall 0 \leq i \leq t$

Final classifier:

given unlabeled x , route to level T

if reaches $v_{l,T}$ s.t. $l < T/2$ output 0

else output 1

End of review

Proof of correctness in special case of stronger weak learner:

Assume WL has 2-sided advantage wrt \mathcal{D} :

$$\Pr_{x \in \mathcal{D}^+} [h(x)=1] \geq \frac{1}{2} + \gamma$$

$$\Pr_{x \in \mathcal{D}^-} [h(x)=0] \geq \frac{1}{2} + \gamma$$

depends on
concept,
not hypothesis
(unknown to learner)

output of WL satisfies

note: all these
should be
labelled 1

$$\Pr_{x \in \mathcal{D}^+} [h(x)=1] \geq \frac{1}{2} + \gamma$$

$$\Pr_{x \in \mathcal{D}^-} [h(x)=0] \geq \frac{1}{2} + \gamma$$

these should
be 0

(can show "reduction" to this case)

main idea:

show $x \in \mathcal{D}^+$ follows random walk biased right
 $x \in \mathcal{D}^-$ " " " " left

Azuma's \neq for sub-Martingales

(similar to Chernoff/Hoeffding when don't have complete independence)

def "submartingale" is sequence

$Y_0 \dots Y_T$ of r.v.'s with finite

means + st. $\forall 1 \leq i \leq T$

$$E[Y_i | Y_0, Y_1, \dots, Y_{i-1}] \geq Y_{i-1}$$

note:
weaker
than
independence

Thm Azuma's \neq for submartingales

Let $0 = Y_0, \dots, Y_T$ be a submartingale

with bounded differences

$$|Y_i - Y_{i-1}| \leq c \quad \forall i$$

then

$$\Pr[Y_T \leq -\lambda] \leq e^{-\lambda^2 / 2Tc^2}$$

(will use $c=1$)

Thm let γ be $\in (0, \frac{1}{2})$

$\forall t = 0..T-1$, suppose each of $t+1$ calls to WL
on $\mathcal{D}_{i,t}$ gives $h_{i,t}$ with 2-sided
advantage γ wrt. $\mathcal{D}_{i,t}^+$
then final hypothesis has error $\leq e^{-\gamma^2 T}$

Pf.

bound error on \mathcal{D}^+ (\mathcal{D}^- is identical)

$X_t \leftarrow \left[\begin{array}{l} \text{pick } x \in \mathcal{D}^+ \\ \text{route to level } t \text{ + reach } v_{i,t} \\ \text{output } i \end{array} \right.$

} process
that used
to define
value of
 X_t

Consider $X_t \mid X_0 \dots X_{t-1}$

lets condition on x reaching $v_{i,t-1}$
then $x \in (\mathcal{D}_{i,t-1})^+$

\leftarrow we are already
conditioning on
 $x \in \mathcal{D}^+$

$$\begin{aligned} \text{so } E[X_t \mid X_0 \dots X_{t-1}] &= X_{t-1} + \Pr_{x \in (\mathcal{D}_{i,t-1})^+} [h_{i,t-1}(x) = 1] \\ &\geq X_{t-1} + \frac{1}{2} + \gamma \end{aligned}$$

here we use
that WL is distribution-free

} since $h_{i,t-1}$ has 2-sided
on $\mathcal{D}_{i,t-1}$ for any i

Define $Y_t \leftarrow X_t - t(\frac{1}{2} + \gamma)$ best known lower bound on $E[X_t]$ $Y_0 = X_0$

See next page for why we care about Y_t

note: conditioning on $Y_0 \dots Y_t$ equivalent to " " $X_0 \dots X_t$

$$\text{so } E[Y_t | Y_0 \dots Y_{t-1}] = E[X_t - t(\frac{1}{2} + \gamma) | Y_0 \dots Y_{t-1}]$$

$$= E[X_t | Y_0 \dots Y_{t-1}] - t(\frac{1}{2} + \gamma)$$

$$\geq X_{t-1} + \frac{1}{2} + \gamma - t(\frac{1}{2} + \gamma)$$

$$= X_{t-1} - (t-1)(\frac{1}{2} + \gamma) = Y_{t-1}$$

\Rightarrow sub-martingale

$$\text{also } |Y_t - Y_{t-1}| = \underbrace{|X_t - X_{t-1}|}_{0 \text{ or } 1} - \underbrace{(t - (t-1))(\frac{1}{2} + \gamma)}_{\in [\frac{1}{2}, 1]} \leq 1$$

So can use Azuma's \neq for submartingales

$$\text{which says } P[Y_T \leq -\lambda] \leq e^{-\lambda^2/2T}$$

taking $\lambda = \gamma \cdot T$ & note

$$\Pr[Y_T \leq -\lambda] \quad \leftarrow \text{drifted } \geq \lambda \text{ more than expected to the left}$$

$$= \Pr[X_T - T \cdot (\frac{1}{2} + \gamma) \leq -\gamma \cdot T]$$

$$= \Pr[X_T \leq \frac{T}{2}] \quad \leftarrow \text{this is exactly when we err}$$

$$= \Pr_{x \in \mathcal{X}^T} [h(x) = 0] \quad \text{for output hypothesis } h$$

$$\begin{aligned} \text{So error rate} &\leq e^{-\frac{(T\gamma)^2}{2T}} \\ &= e^{-T \cdot \gamma^2 / 2} \end{aligned}$$

by defn of output hypothesis

\Rightarrow only need $T \sim O((\log 1/\epsilon) / \gamma^2)$
to get error $\leq \epsilon$

Sample / time complexity:

main issue: what if it takes a

long time to get enough samples of either \pm

for one of the weak learners?

(perhaps all the samples get filtered
& placed in other buckets?)

this will likely happen
→
for \pm samples
on $-$ side
& viceversa

Can "freeze" the node

(won't contribute much to error)

freezing rule: if $v_{i,t}$ reached by "b"-labelled samples

with $\text{prob} \leq \frac{\epsilon}{T(T+1)}$

then label $v_{i,t}$ by "1-b" & make
↑ other label

it terminal

if not frozen, then get to $v_{i,t}$ with label b
every $T(T+1)/\epsilon$ samples!

Hypothesis: if no freezing on path, predict as majority
of $h_{i,t}(x)$ on path
(same as before)

else if freeze at level t node $v_{j,t}$

any example that reaches node $v_{j,t}$
gets labelled via $v_{j,t}$'s label

Why is the error ok?

• each frozen node contributes $\leq \frac{\epsilon}{T(T+1)}$ to

error

• $\leq T(T+1)$ frozen nodes $\Rightarrow \leq \epsilon$ total error

• need to redo martingale analysis for unfrozen part
(define intermediate process to make it similar)

Boosting \Rightarrow

average case complexity insights:

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

Pf idea

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

will output consistent hypothesis of

poly size in $n (= \log \frac{1}{\epsilon}) + |c|$

\leftarrow # bits used by WL
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

why nontrivial?
output for any single input can be hardwired into ckt. \Rightarrow no "hard" inputs. how to find a set of "collectively hard" inputs?

Pf idea

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

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

Amplifying hardness

Goal: take worst case hard fctn & turn into
ave case hard fctn.

how?

Yao's xor lemma

- works for any hard fctn

- intuition:

given δ -biased coin ($\Pr[\text{heads}] = \delta < \frac{1}{2}$)

predict correctly with prob $1 - \delta$

predict parity of k -tosses correctly

with prob $\approx \frac{1}{2} + (1 - 2\delta)^k$

$\rightarrow \frac{1}{2}$ as $k \rightarrow \infty$

need to predict
all k ?

- is solving k independent copies of f

k times harder than solving 1?

not always: matrix-vector mult is $\Theta(n^2)$ time, matrix-matrix is $\Theta(n^3)$

def S is set
 f is ε -hard core on S for size g if
 \forall circuits C of size $\leq g$

$$\Pr_{x \in_n S} [C(x) = f(x)] \leq \frac{1}{2} + \frac{\varepsilon}{2}$$

Lemma Given f

define $f^{\oplus k}(x_1 \dots x_k) = f(x_1) \oplus \dots \oplus f(x_k)$

if f is ε -hardcore for set H

st. $|H| \geq \delta \cdot 2^n$ for ccts of size $g+1$

then $f^{\oplus k}$ is $\varepsilon + 2(1-\delta)^k$ -h.c.

for size g .