

## Lecture 20

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

(if time) average vs. worst case complexity

Last time: algorithm that weakly learns monotone fctns  
↑  
Slightly better than random guess

What good is weak learning?

Unclear

here can only weakly learn on  
uniform distribution

ability to weakly learn on  
all distributions

⇒ ability to strongly learn  
[Schapire]

"boosting"

# Weak vs. Strong Learning

Def. Algorithm  $A$  "(strongly) PAC learns" concept class  $\mathcal{C}$

if  $\forall c \in \mathcal{C} + \forall$  dists  $\mathcal{D} \quad \forall \epsilon, \delta$

given examples of  $c$  (labelled)

with prob  $\geq 1 - \delta$

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

# 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

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  $\gamma$  ?

$\delta$  ?

$\epsilon$  ?

# 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

$$\boxed{E_x} \rightarrow (x_1, c(x_1)) (x_2, c(x_2)) \dots$$

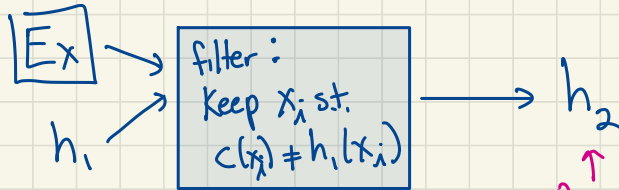
$x_i \in \mathcal{X}$

1. Run WL on  $E_x \rightarrow h_1$

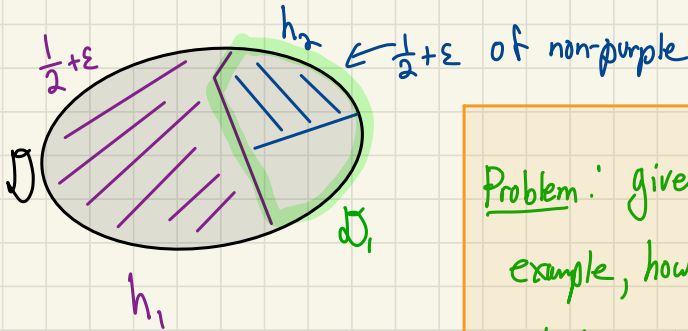


$$\Pr_{x \in \mathcal{D}} [h_1(x) = c(x)] \geq \frac{1}{2} + \gamma$$

2. use  $h_1$  to "filter"



$$\Pr_{x \in \mathcal{D}_1} [h_2(x) = c(x)] \geq \frac{1}{2} + \gamma$$



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

e.g.  $\uparrow$  if  $h_1(x) = h_2(x) = b$  seems good to output  $b$  but what if  $h_1(x) \neq h_2(x)$ ?

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  
... (lots and lots)

Filtering Procedures: (lots of these)

- 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

domain  $X$

$\mathcal{C}$  concept class:  $c \in \mathcal{C}$  maps  $X \rightarrow \{0,1\}$

• Given labelled examples

$$(x_1, c(x_1)) (x_2, c(x_2)) \dots$$

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

$c \in \mathcal{C}$  "target function"

• Given weak learning alg  $W_L$  which

weakly learns (advantage  $\gamma$ ) on any dist  $\mathcal{D}$

$$\parallel \text{error} \frac{1}{2} - \gamma$$

$$\begin{aligned} & \text{error}_{\mathcal{D}}(h) \\ &= \Pr_{x \in \mathcal{D}} [c(x) \neq h(x)] \end{aligned}$$

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

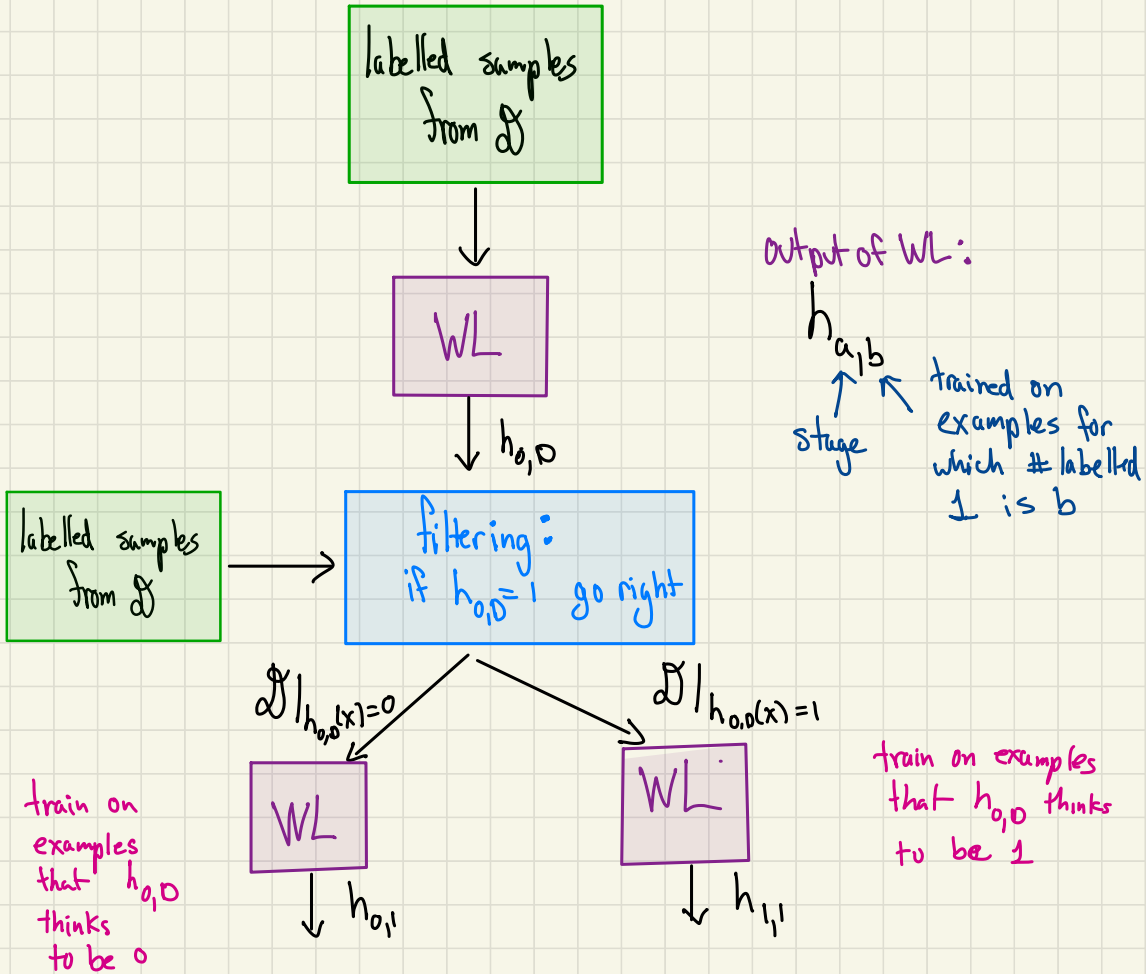
( $\star$  can evaluate in poly time)

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

# Idea: Filter by hypothesis value

One phase of filtering

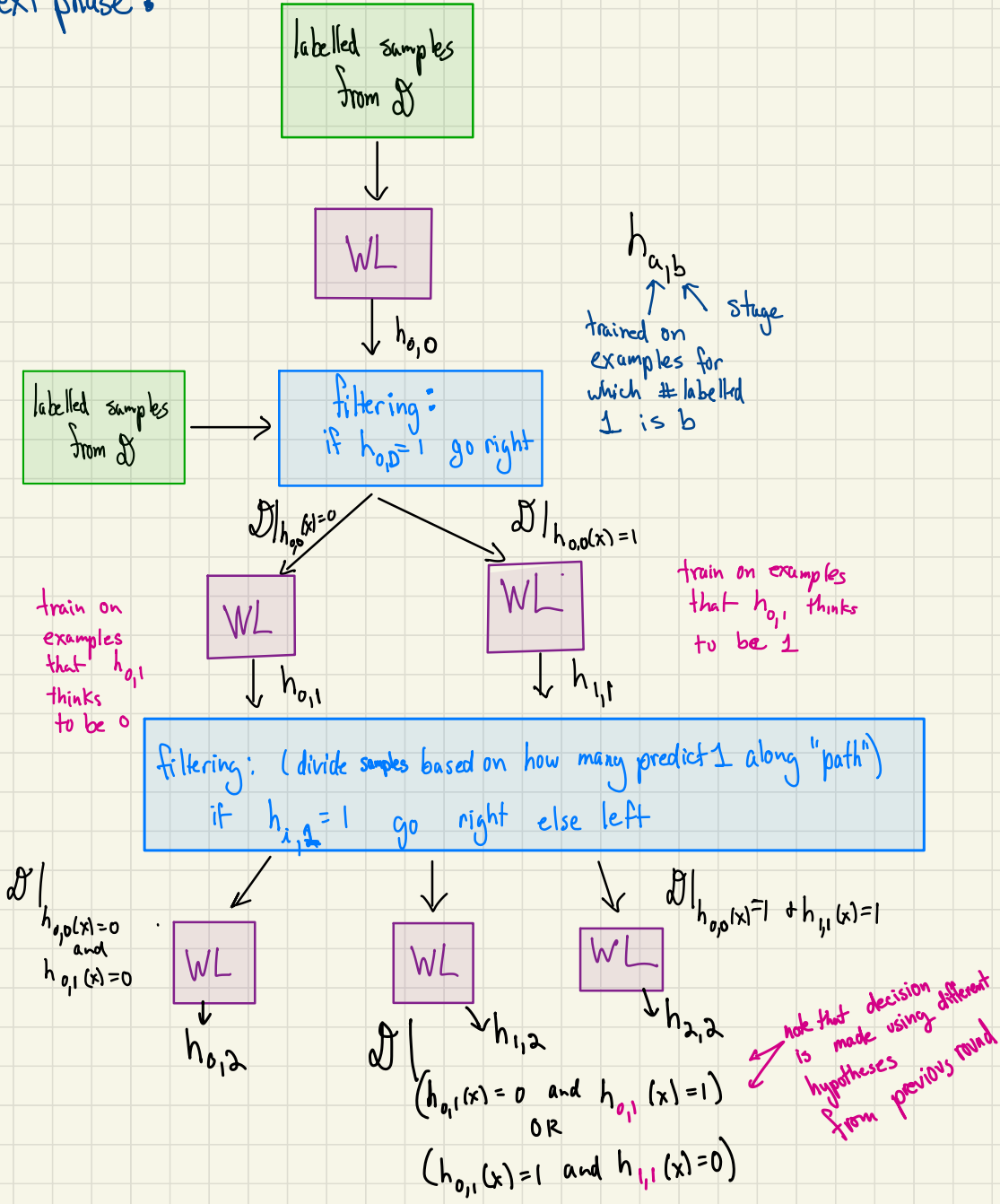
assume  $C: X \rightarrow \{0,1\}$



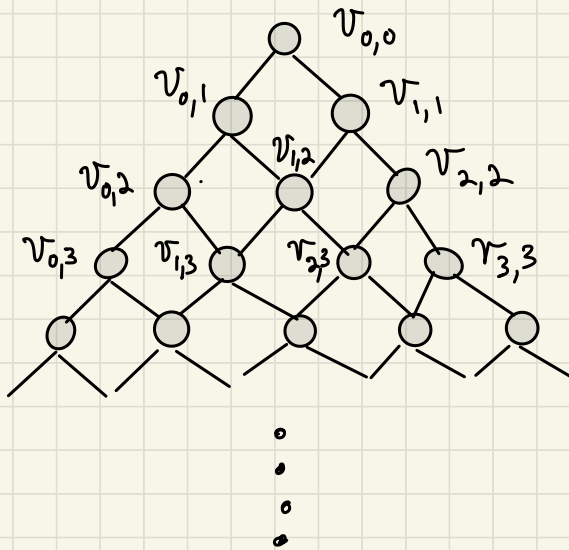
# Idea: Filter by hypothesis value (cont.)

next phase:

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

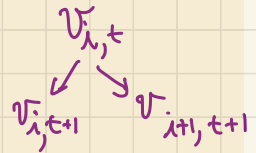


Putting together the filters:

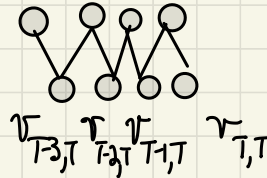
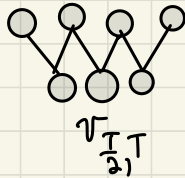
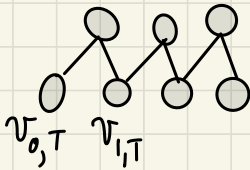


$T$  stages  
 $T+1$  layers

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



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 $t^{\text{th}}$ stage:

branching program induces  $t+1$  different distributions

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

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

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

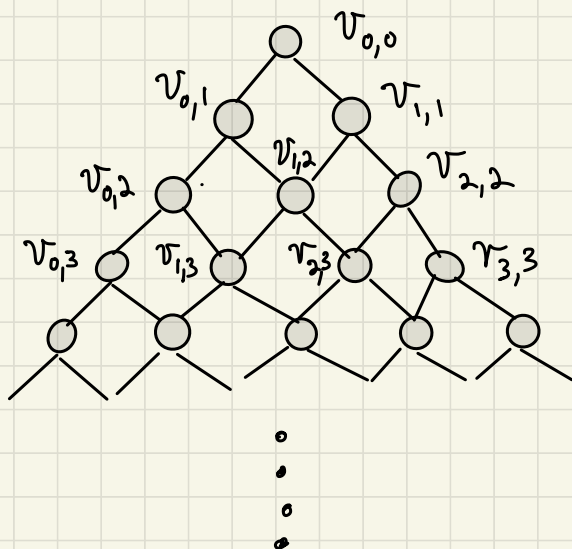
## Final classifier:

given unlabeled  $x$ , route to level  $T$

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

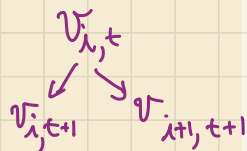
else output 1

Putting together the filters:

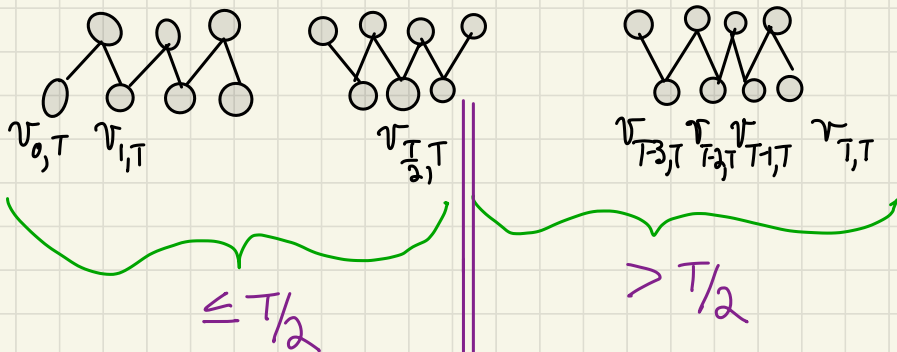


$T$  stages  
 $T+1$  layers

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  $h_{ij}$

Proof of correctness in special case:

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

$$\text{let } \mathcal{D}^+ \leftarrow \mathcal{D} \mid C(x)=1$$

$$\mathcal{D}^- \leftarrow \mathcal{D} \mid C(x)=0$$

output of WL satisfies

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

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

note: all these  
should be  
labelled 1

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.  $\forall 1 \leq i \leq T$

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

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  $\gamma$  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  $\gamma$  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)

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

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

Condition on  $X$  reaching  $\mathcal{V}_{i,t-1}$

then  $X \in (\mathcal{D}_{i,t-1}^+)^+$

$$\text{so } E[X_t \mid X_0 \dots X_{t-1}] = X_{t-1} + \Pr_{X \in (\mathcal{D}_{i,t-1}^+)^+} [h_{i,t-1}(X) = 1]$$

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

since  $h_{i,t-1}$  has 2-sided advantage

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

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

$$\begin{aligned}
 \text{so } E[Y_t | Y_0 \dots Y_{t-1}] &= E[X_t - t(\frac{1}{2} + \gamma) | Y_0 \dots Y_{t-1}] \\
 &= E[X_t | Y_0 \dots Y_{t-1}] - t(\frac{1}{2} + \gamma) \\
 &\geq X_{t-1} + \frac{1}{2} + \gamma - t(\frac{1}{2} + \gamma) \\
 &= X_{t-1} - (t-1)(\frac{1}{2} + \gamma) = Y_{t-1} \\
 &\Rightarrow \text{sub-martingale}
 \end{aligned}$$

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

So can use Azuma's  $\neq$  for submartingales

$$\text{which says } P_0[Y_T \leq -\lambda] \leq e^{-\lambda^2/2T}$$

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

$$\Pr[Y_T \leq -\lambda]$$

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

$$= \Pr[X_T \leq \frac{T}{2}]$$

$$= \Pr[h(x) = 0] \text{ for output hypothesis } h$$

$x \in \mathcal{X}^T$

by defn of output hypothesis

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

$\Rightarrow$  only need  $T \sim \log \frac{1}{\gamma \epsilon}$

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

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

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

$\Rightarrow$  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