# Lecture 17

## Boosting

# Weak Learning

def. algorithm $A$ **weakly** PAC learns concept class $C$ if $\exists \gamma > 0$ s.t.

$$\forall c \in C \quad \& \quad \forall \text{ dists } \mathcal{D},$$

given examples of $c$ according to $\mathcal{D}$

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

$\uparrow$
**advantage**

Thm if $C$ can be weakly PAC learned (on any $\mathcal{D}$) then

$C$ can be (strongly) PAC learned.

# Weak vs. Strong Learning

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

if $\forall\ c \in C$ & $\forall$ dists $\mathcal{D}$ $\qquad \exists\ \gamma > 0$

$\forall\ \not{\epsilon}, \delta > 0 \qquad \left(\delta = \frac{1}{4} \text{ or } \frac{1}{n^2} \text{ doesn't affect}\right)$

with prob $\geq 1 - \delta$

given examples of $c$

$A$ outputs $h$ s.t. $\Pr_{\mathcal{D}}\left[h(x) \neq c(x)\right] \leq \not{\epsilon}$

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

$\uparrow$
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 <u>any</u> dist $\mathcal{D}$ then $C$ can be (strongly) learned.

# Applications

1) "Theoretical"
- Unif dist Algorithms for poly term DNF
   weight $w$ - poly threshhold fctns     } low degree
                                             alg doesn't
                                             work well

   ∴  (Boosting + KM)
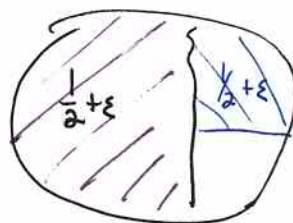
- Ave case vs. worst case

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

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 ← need for checking future hypotheses

incorrectly ← need to improve on these

## The setting

- Given labelled examples

$$(x_1, f(x_1)), (x_2, f(x_2)), \ldots$$

$$x_i \in_R \mathcal{X}$$
$$f \in \mathcal{C}$$

- Given weak learning alg $WL$ which weakly learns (advantage $\frac{\gamma}{2}$) on **any** dist $\mathcal{D}'$

# Boosting Algorithm

- Stage 0 (Initialize)

$$\mathcal{D}_0 \leftarrow \mathcal{D}$$

run WL on $\mathcal{D}_0$ to generate (whp)

$$C_1 \quad \text{s.t.} \quad \Pr_{\mathcal{D}_0}[f(x) = C_1(x)] \geq \tfrac{1}{2} + \gamma/2$$

- For $i = 1 \ldots T = O(\frac{1}{\gamma^2 \varepsilon})$ stages, stage $i$: (can stop if Majority($C_1 \ldots C_i$) correct on $\geq 1-\varepsilon$ inputs)

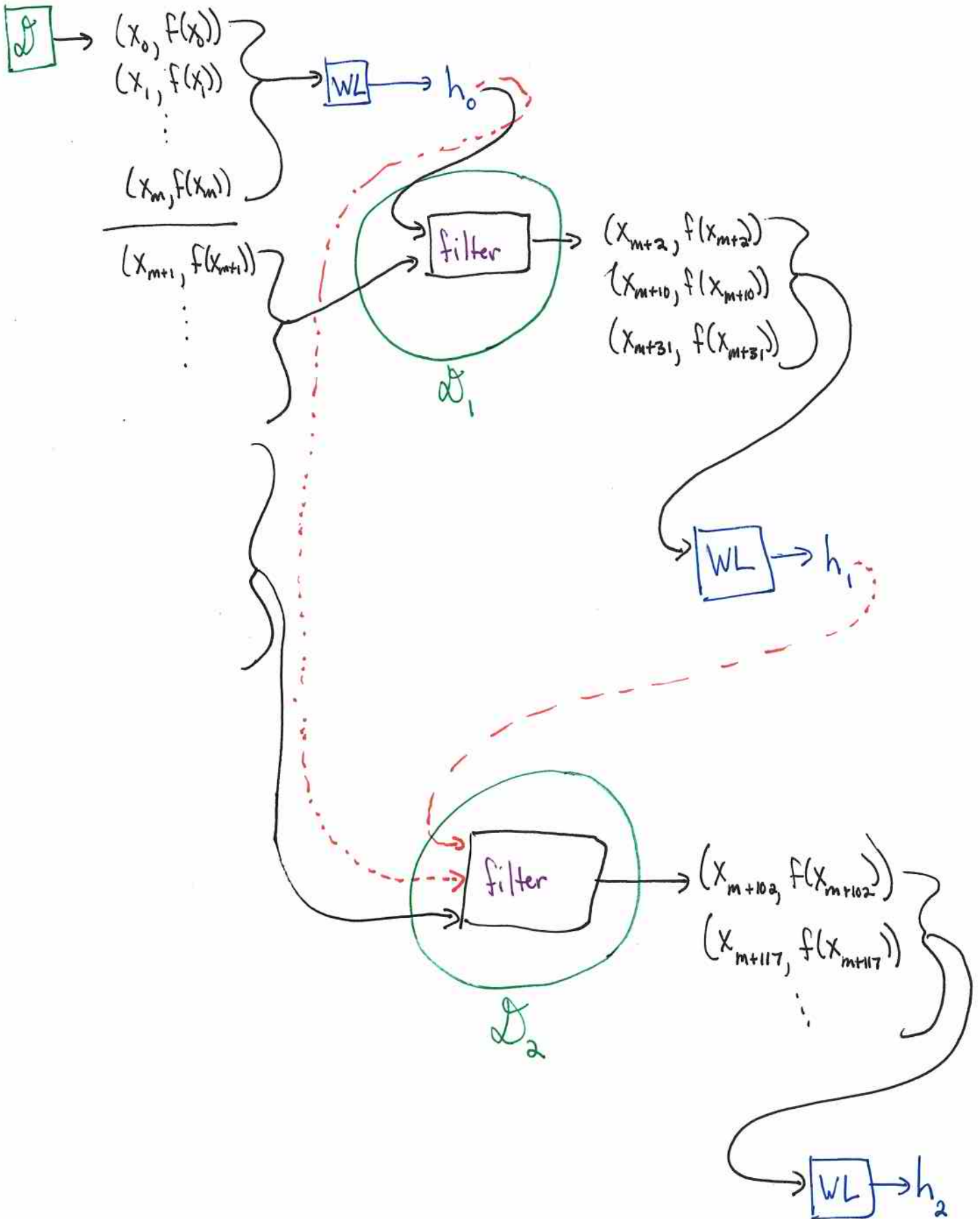    (1) Construct $\mathcal{D}_i$ via "filtering procedure":

    { favor pts on which maj of $C_1 \ldots C_i$ don't do well

    but also keep some other points }

    Will specify soon

    (2) run WL on examples from $\mathcal{D}_i$ to output

    $$C_{i+1} \quad \text{s.t.} \quad \Pr_{\mathcal{D}_i}[f(x) = C_{i+1}(x)] \geq \tfrac{1}{2} + \tfrac{\gamma}{2}$$

- output $C = MAJ(C_1 \ldots C_T)$

$\mathcal{D} \rightarrow$

$(x_0, f(x_0))$
$(x_1, f(x_1))$
$\vdots$
$(x_m, f(x_m))$

$\boxed{WL} \rightarrow h_0$

$(x_{m+1}, f(x_{m+1}))$
$\vdots$

filter $\rightarrow$

$(x_{m+2}, f(x_{m+2}))$
$(x_{m+10}, f(x_{m+10}))$
$(x_{m+31}, f(x_{m+31}))$

$\mathcal{D}_1$

$\boxed{WL} \rightarrow h_1$

filter $\rightarrow$

$(x_{m+102}, f(x_{m+102}))$
$(x_{m+117}, f(x_{m+117}))$
$\vdots$

$\mathcal{D}_2$

$\boxed{WL} \rightarrow h_2$

## Filtering procedure

Given new example $x, f(x)$ from example oracle

- if majority of $c_1 \ldots c_i$ wrong, keep it

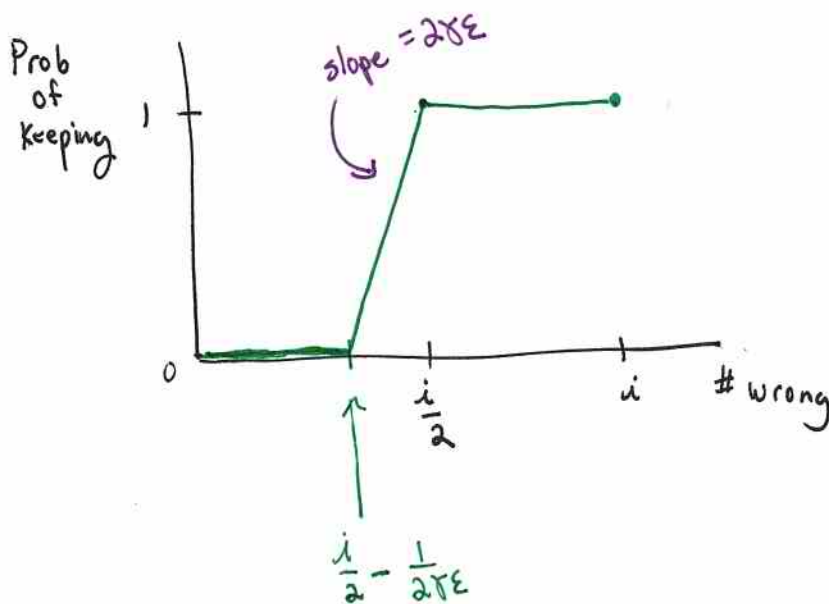  ie. $\geq \frac{i}{2}$

- if large majority right, then discard

  ie. #right − #wrong $> \frac{1}{\gamma\varepsilon}$

  or #wrong $\leq \frac{i}{2} - \frac{1}{2\gamma\varepsilon}$

- else #right − #wrong $= \frac{\alpha}{\gamma\varepsilon}$ for $0 < \alpha < 1$

  #wrong − #right $= \frac{-\alpha}{\gamma\varepsilon}$

  So keep with prob $= 1 - \alpha$

Need to show:

1) Output is has nontrivial agreement with $f$

2) # samples needed not too bad

why could it be bad?
if throw out lots of samples, might
need to wait a long time before WL
can give an output,
but if throw out too many samples then
you already have a good hypothesis!

↑

will stop if $Maj(C_1 \dots C_i)$ correct on $\geq 1-\varepsilon$ fraction
of inputs

o.w. $Maj(C_1 \dots C_i)$ incorrect on $> \varepsilon$ fraction

so filtering procedure outputs
sample with prob $\geq \varepsilon$
(+ in expectation, every $1/\varepsilon$ samples
of $\mathcal{D}$ at least one makes
it thru the filtering
system)

⟹ filtering slows down sample
collection by $\leq O(1/\varepsilon)$

So lets focus on ①

# Notation

- $R_c(x) = \begin{cases} +1 & \text{if} \quad f(x) = c(x) \\ -1 & \text{if} \quad f(x) \neq c(x) \end{cases}$

  "is c correct on x?"

- $N_i(x) = \displaystyle\sum_{1 \leq j \leq i} R_{c_j}(x)$

  after iteration $i$, how many $c$'s correct? (#right - #wrong)

- $M_i(x) = \begin{cases} 1 & \text{if} \quad N_i(x) \leq 0 \\ 0 & \text{if} \quad N_i(x) \geq \frac{1}{\varepsilon \gamma} \\ 1 - \varepsilon \cdot \gamma \cdot N_i(x) & \text{o.w.} \end{cases}$

  prob of keeping $x$ in filtering (after stage $i$)

  note — all "wrong" $x$ included in $M$
  also some "right" $x$ included

Note that new distribution on samples is proportional to $M_i$:

- $D_{M_i}(x) \equiv \dfrac{M_i(x)}{\displaystyle\sum_x M_i(x)}$

  distribution induced by $M$

  note $\underline{\qquad} \quad D_{M_i}(x) = \mathcal{D}_i$

  $\displaystyle\sum_x M_i(x)$ includes all "wrong" $x$ but $\left.\begin{array}{c} \text{upper} \\ \text{bounds} \\ \# \\ \text{wrong} \\ x \end{array}\right\}$

  also $x$ for which maj that isn't overwhelming are correct

How correct are we wrt. $D_{M_i}$?

- $\mathrm{Adv}_c(M_i) = \displaystyle\sum_x R_c(x) M_i(x)$

  "Advantage" of $c$ on $M$
  $\sim \Pr[\text{correct}] - \Pr[\text{incorrect}]$
  $= 2 \cdot \Pr[\text{correct}] - 1$

- $\Pr_{x \in D_{M_i}}[c(x) = f(x)] = \dfrac{1}{2} + \dfrac{\mathrm{Adv}_c(M_i)}{2 \cdot \underbrace{\sum_x M_i(x)}_{\gamma/2}}$

**Note:**

if $\sum M_i(x) \geq \varepsilon \, 2^n$

$Adv_c(M_i) \geq \gamma \cdot \varepsilon \cdot 2^n$

convert claim about WL ⇒ claim about advantage

ie. if have $\gamma$ advantage on output of WL & still almost (r) wrong on lots of inputs then new advantage is pretty good

if not, then you are done

---

**Begin Proof**

For input $x$

let $A_i(x) \leftarrow \sum_{0 \leq j \leq i-1} R_{c_{j+1}}(x) \, M_j(x)$
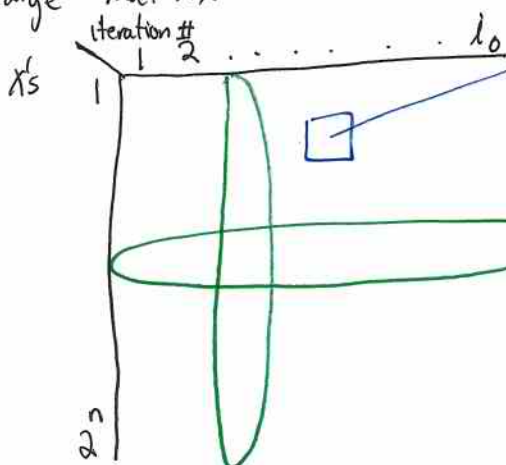
Strange —
indices don't match
$c_1 \cdots c_j$ define $\mathcal{D}_j$
but $c_{j+1}$ learned from WL on $\mathcal{D}_j$

**Claim** $A_i(x) \leq \frac{1}{\varepsilon \gamma} + \frac{\varepsilon \gamma}{2} \cdot i$

· bounds advantage per input

· only helps after $\frac{1}{\varepsilon \gamma}$ rounds

Plan for use of claim:

Consider large matrix:

iteration #

$x$'s

$(k,j)^{th}$ Entry: $R_{c_{j+1}}(x_k) \, M_j(x_k)$

$x$'s row sum $= \sum_{0 \leq j \leq i_0} R_{c_{j+1}}(x) \, M_j(x) = A_{i_0+1}(x)$

$j^{th}$ col sum $= \sum_x R_{c_{j+1}}(x) \, M_j(x)$

$= Adv_{c_{j+1}}(M_j) \leq \geq \gamma \sum_x M_j(x)$ else algorithm stops