


Boosting

The Weak PAC Model

- Exactly the same as PAC, but no ϵ
- Ask that $f \in \mathcal{C}, P, \delta$
algo outputs $h \in \mathcal{H}$
s.t. $\epsilon(h) \leq \frac{1}{2} - \alpha \triangleq \epsilon_0$
- Here α is a fixed constant (e.g. $\alpha = 0.01$) or even $\alpha = \frac{1}{\text{poly}(n)}$
- Algo runs in time $\text{poly}(n, \text{size}(c), \delta)$

Once upon a time...

- Every C known to be weakly learnable was known to be strongly learnable
- Every C known to be hard for strong learning was also hard for weak learning
- Q: Could weak learning \equiv strong learning?

Note: If so, needs
distribution-free
property of PAC.

\exists classes C and
specific P s.t.

- C weakly learnable
w.r.t. P
- C not strongly
learnable w.r.t. P

Boosting

- Let L be an algo for weakly learning C by \mathcal{H}
- Idea: Use L as a subroutine in an algo L' for strongly learning C by \mathcal{H}'

Attempt #1

- Run weak L ℓ times on $c, P \rightarrow$ get h_1, h_2, \dots, h_ℓ
 - Let
- $$h(x) = \underbrace{\text{MAJ}(h_1(x), \dots, h_\ell(x))}_{\text{majority vote}}$$
- Good idea?

Attempt #2

- Run L on c, P → get h_1 ,
- Run L again, but force it to do something different:

$P_1 = P$ ~~$h_1 = c$~~ $h_1 \neq c$

P_2 ~~$h_2 = c$~~ $h_2 \neq c$
shrink s grows

More formally, if

$$\varepsilon(h_1) = \varepsilon_0 < 1/2$$

Then define P_2 :

$\forall x \in X$:

If $h_1(x) = c(x)$

$$P_2(x) = \frac{1}{2(1-\varepsilon_0)} P_1(x)$$

If $h_1(x) \neq c(x)$

$$P_2(x) = \frac{1}{2\varepsilon_0} P_1(x)$$

- By design,
 $\epsilon_{P_2}(h_1) = 1/2$
- We can simulate
 P_2 from P_1, h_1 (how?)
- Run L on $P_2 \Rightarrow$
 get $h_2 \neq h_1$
- Now what?

Breaking ties

• Define P_3 :

If $h_1(x) \neq h_2(x)$

$$P_3(x) = P_1(x)/Z$$

$$Z = P_1[h_1(x) \neq h_2(x)]$$

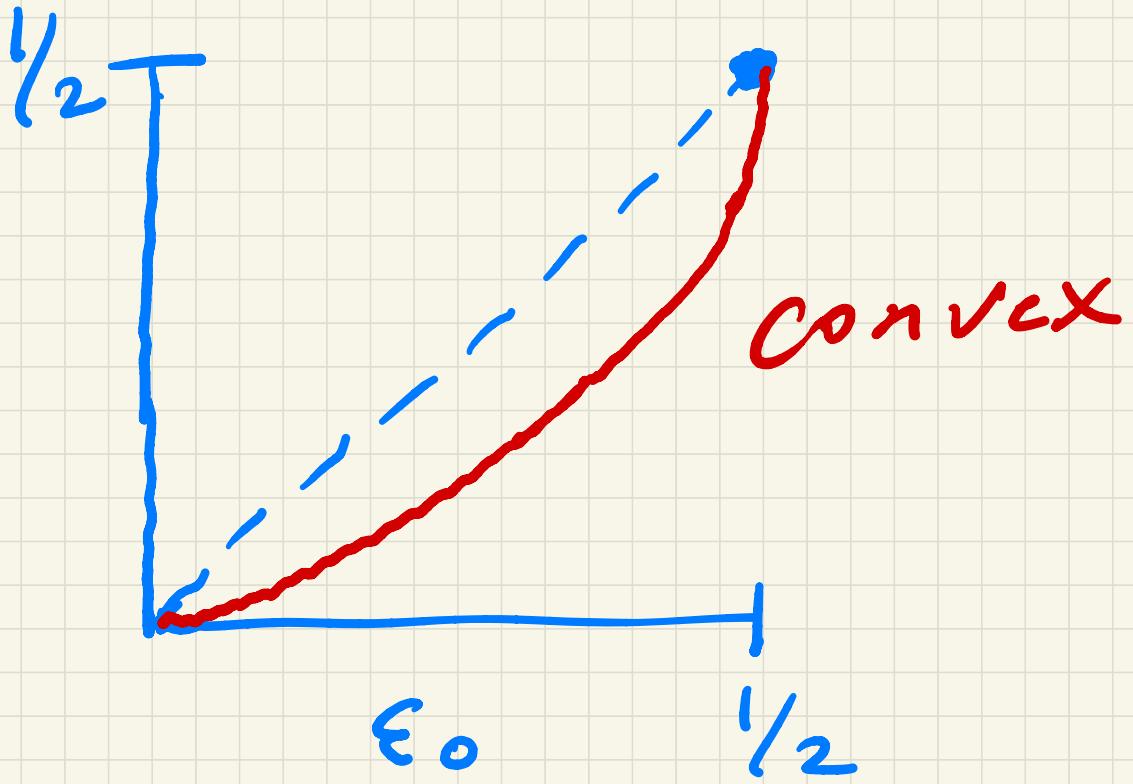
If $h_1(x) = h_2(x)$

$$P_3(x) = 0.$$

- We can simulate P_3 from P_1, h_1, h_2 (how?)
- Run L on P_3
 \Rightarrow get h_3
- Now let
$$h(x) = \text{MAJ}(h_1(x), h_2(x), h_3(x))$$
- Now what?

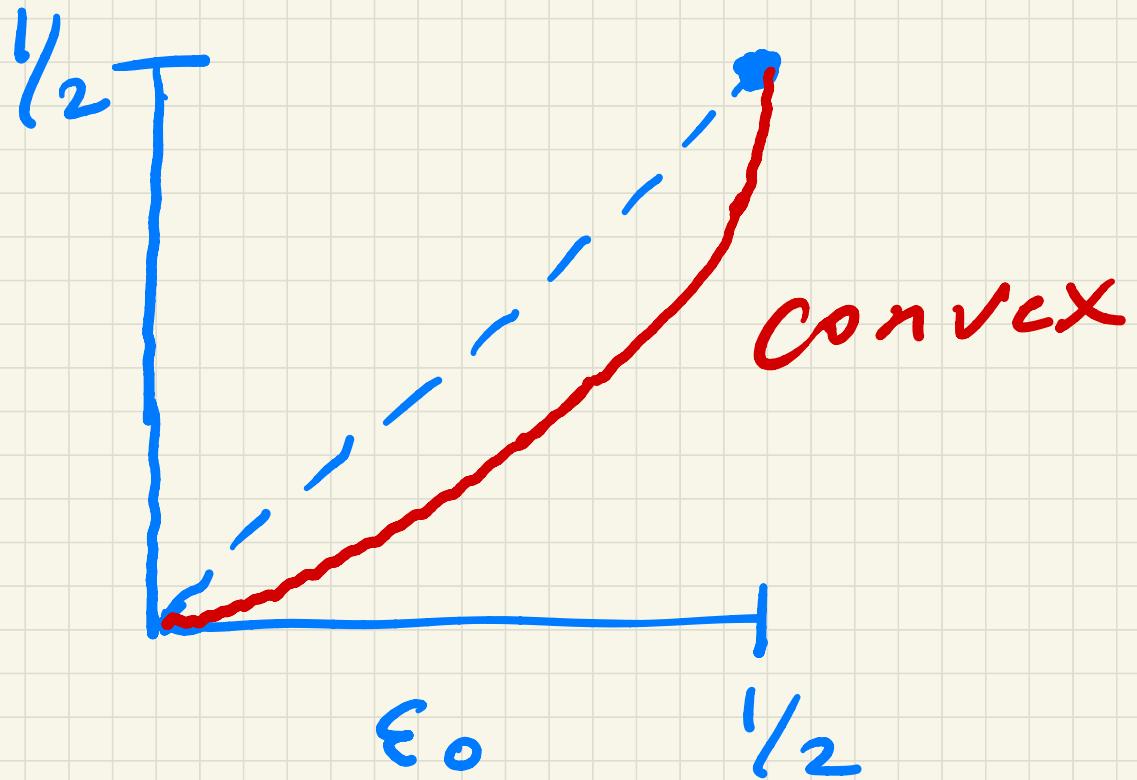
Lemma

$$\epsilon_{P_i}(h) \leq 3\epsilon_0^2 - 2\epsilon_0^3$$



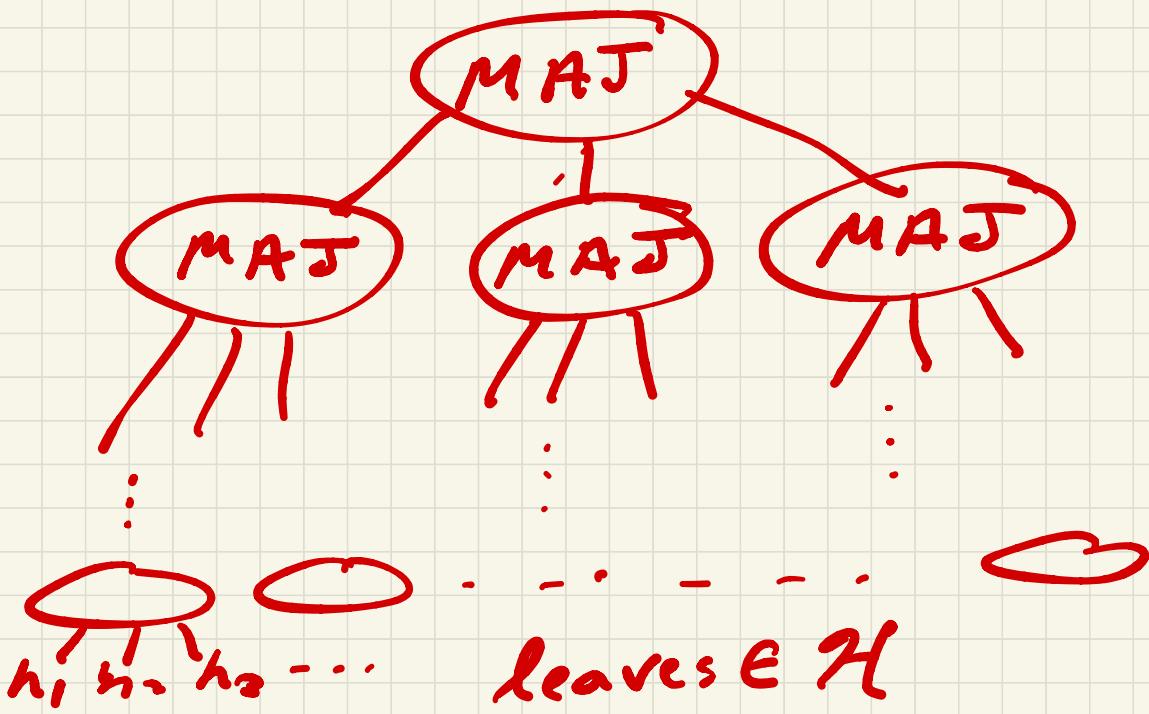
So now have a
better weak learner!

So now let's recurse:



- Need $\text{unloglog}(\ell/\epsilon)$ levels to go from $\epsilon_0 \rightarrow \epsilon$

- Final hypothesis:



A More Practical
Formulation:

Adaboost



Setup

- Given data

$$S = \langle x_1, y_1 \rangle, \dots, \langle x_n, y_n \rangle$$
$$y_i \in \{+1, -1\}$$

- Goal: draw error on
 $S \rightarrow O$ using
weak learner

- Initial distribution

P_1 : uniform on S

- Subsequent P_t will be
reweightings on S

Adaboost

For $t = 1, 2, \dots, T$:

- run weak algo on P_t
 \Rightarrow get $h_t \in \mathcal{H}$

- choose $\alpha_t > 0$

- define f_i :

$$P_{t+1}(i) = \frac{P_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

$\overbrace{\quad\quad\quad}^{Z_t = \text{normalization}}$
 Z_t for P_{t+1}

Final classifier:

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

So \mathcal{H}' = linear combos
over \mathcal{H} .

Some Notation

- Define

$$\epsilon_t = P_t [h_t(x_i) \neq y_{i,t}]$$

- Allowing ϵ_t to vary
(including $\epsilon_t = 1/2$)

- Can easily measure

- Define $\pi_t = 1/2 - \epsilon_t$

- Advantage over $1/2$

Adaboost: Analysis

$$y_i \neq H(x_i) = \text{sign}\left(\sum_t \alpha_t h_t(x_i)\right)$$

$$\Rightarrow y_i \sum_t \alpha_t h_t(x_i) \leq 0$$

$$\Rightarrow e^{-y_i \sum_t \alpha_t h_t(x_i)} \geq 1$$

So

$$\frac{1}{m} \sum_i I[H(x_i) \neq y_i] \xrightarrow{\text{error of } H \text{ on } S}$$

$$\leq \frac{1}{m} \sum_i e^{-y_i \sum_t \alpha_t h_t(x_i)}$$

Focus on bounding this (*)

For any fixed i , recall:

$$P_{t+1}(i) = \frac{P_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

Rewrite:

$$Z_t = \frac{P_t(i) e^{-\alpha_t y_i h_t(x_i)}}{\underbrace{P_{t+1}(i)}_{=}}$$

Let's define θ_i :

$$\underbrace{P_{T+1}(i)}_{\text{not used anyway}} = \frac{1}{m} = P_i(i)$$

Next: show that

$$(*) = \prod_{t=1}^T Z_t$$

$$\prod Z_t = Z_T \cdot Z_{T-1} \cdots Z_1 =$$

$$P_T(i) e^{-\alpha_T y_i h_T(x_i)} / P_{T+1}(i)$$

$$\times P_{T-1}(i) e^{-\alpha_{T-1} y_i h_{T-1}(x_i)} / P_T(i)$$

$$\times P_{T-2}(i) e^{-\alpha_{T-2} y_i h_{T-2}(x_i)} / P_{T-1}(i)$$

⋮

$$\times P_1(i) e^{-\alpha_1 y_i h_1(x_i)} / P_2(i)$$

$$= e^{-y_i \sum_t \alpha_t h_t(x_i)}$$

This holds $\forall i$, so get
same if we average:

$$\prod_t Z_t = \frac{1}{m} \sum_i e^{-y_i \sum a_t h_t(x_i)}$$

$= (*)$ from before!

\geq error on S



Next: choose a_t 's

to make Z_t 's as

small as possible.

Adaboost

For $t = 1, 2, \dots, T$:

- run weak algo on P_t
 \Rightarrow get $h_t \in \mathcal{H}$

- choose $\alpha_t > 0$

- define f_i :

$$P_{t+1}(i) = \frac{P_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

$\overbrace{\quad\quad\quad}^{Z_t = \text{normalization}}$
for P_{t+1}

$$Z_t = \sum_i P_t(i) e^{-\alpha_t y_i h_t(x_i)}$$

= normalization for P_{t+1}

$$= \sum_i P_t(i) e^{-\alpha_t}$$

$$i : h_t(x_i) = y_i$$

$$+ \sum_{i : h_t(x_i) \neq y_i} P_t(i) e^{\alpha_t}$$

$$= (1 - \varepsilon_t) e^{-\alpha_t}$$

$$+ \varepsilon_t e^{\alpha_t}$$

If we choose

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1-\varepsilon_t}{\varepsilon_t} \right)$$

then

$$Z_t = (1-\varepsilon_t) \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}}$$

$$+ \varepsilon_t \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}}$$

$$= 2 \sqrt{\varepsilon_t (1-\varepsilon_t)}$$

so error on $S \leq$

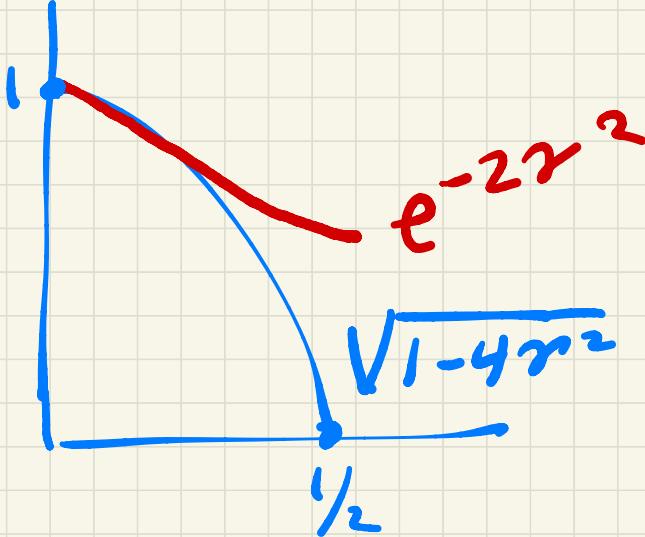
$$\frac{\pi}{t} Z_t = \frac{\pi}{t} 2 \sqrt{\varepsilon_t (1-\varepsilon_t)} \leq 1$$

$$\sqrt{4\varepsilon(1-\varepsilon)} =$$

$$\sqrt{4(\frac{1}{2}-\gamma)(\frac{1}{2}+\gamma)} =$$

$$\sqrt{4(\frac{1}{4}-\gamma^2)} =$$

$$\sqrt{1-4\gamma^2} \leq e^{-2\gamma^2}$$



$$S_0 \frac{\pi}{t} z_t \leq \pi e^{-2\sum_t \gamma_t^2}$$

$$= e^{-2 \sum_t \gamma_t^2}$$

- exponential decay in $\sum_t \gamma_t^2$
- adaptive in γ_t 's

- If all $\gamma_t \geq \gamma > 0$:
error on $S \leq e^{-2T\gamma^2}$

- Set $\epsilon < \gamma_m$ (consistency):

$$T > \frac{1}{2\gamma^2} \ln(m)$$

- $VCD = VCD(\mathcal{H}) \cdot T = dT$

• So

$$T > \frac{1}{2\gamma^2} \ln(dT)$$

suffices.

Weak = Strong via AdaBoost

- Suppose \mathcal{C} weakly PAC by \mathcal{H} with $\varepsilon_0 = \gamma_2 - \gamma$
- If we make

$$e^{-2T\gamma^2} < 1/m$$

we get a consistent

$$H(x) = \text{sign}\left(\sum_t \alpha_t h_t(x)\right)$$

- If $VCD(\mathcal{H}) = d$ then
 VCD of these $\leq Td$

So:

- Choose $m \approx \frac{Td}{\varepsilon} \log(1/\delta)$

- Then solve

$$e^{-2Tr^2} < 1/m$$

$$-2Tr^2 < \ln(1/m)$$

$$T > \frac{1}{2r^2} \ln(m)$$

$$T > \frac{1}{2r^2} \ln\left(\frac{Td}{\varepsilon} \ln(1/\delta)\right)$$

With these choices for
 m and T , with
prob. $\geq 1 - \delta$, $\varepsilon(H) \leq \varepsilon$.

\therefore Weak PAC \Rightarrow PAC.

