

Equivalence between weak and strong learning.

- We consider a modified version of AdaBoost
- Instead of weighted sample (D_t, S) we draw an (unweighted) i.i.d. sample S_t from D_t
- We feed S_t to the weak learner, which produces $h_t = WL(S_t)$
- Note: In practical applications, if you have an algorithm

that does support sample weights, you can use the same trick.

- The sizes of the samples S_1, S_2, \dots will be m_0

- If WL is γ -weak learner

$$\text{err}_{D_t}(\text{WL}(S_t)) \leq \gamma$$

- So if $T > \frac{\ln |S|}{2\gamma^2}$ then

AdaBoost outputs h such that $\widehat{\text{err}}_S(h) = 0$.

- Ix10 want prove that

$\text{err}_D(h)$ is small without relying on the assumption that h_1, h_2, \dots, h_T belong to a class with finite VC dimension.

• We use compression schemes.

Definition: A learning algorithm

$A: (X \times Y)^* \rightarrow Y^X$ is called

k -compressing if there exist

a mapping $\phi: \bigcup_{t=0}^k (X \times Y)^t \rightarrow Y^X$

and for any labeled sample

$S = ((x_1, y_1), (x_2, y_2), \dots, (x_m, y_m))$

there exist $r \in \{0, 1, \dots, k\}$ and indices $i_1, i_2, \dots, i_r \in \{1, 2, \dots, m\}$ such that

$$A(s) = \phi((x_{i_1}, y_{i_1}), \dots, (x_{i_r}, y_{i_r})).$$

- ϕ is called decompressor
for A

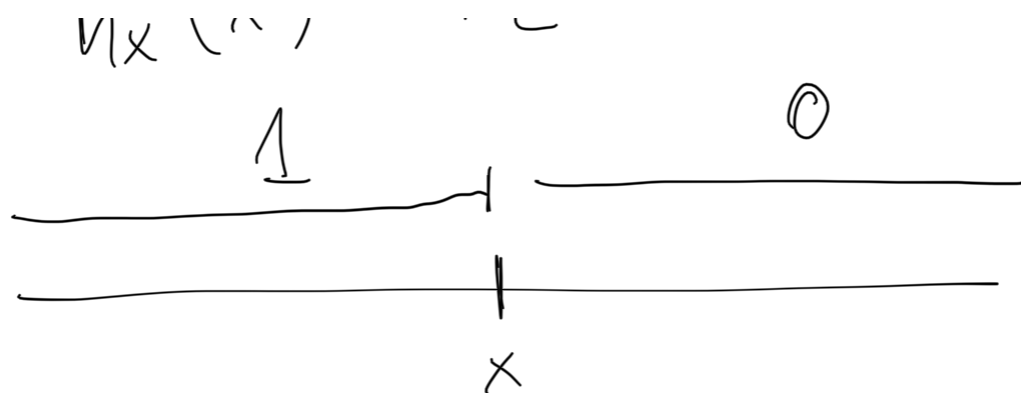
Example:

- Let $X = \mathbb{R}$ and $Y = \{0, 1\}$

- $\phi((x, y)) = h_x \in Y^X$

where

$$h_x(x') = \mathbb{1}[x' \leq x]$$



and

$$\phi(\emptyset) = h_0$$

where $h_0(x) = 0$ for all $x \in X$.

- Consider the class of thresholds

$$H = \{h_x : x \in \mathbb{R}\}$$

- There exists a particular version of ERM that is 1-compressing :

Given $S = (x_1, y_1) \dots (x_m, y_m)$

let $H_S = \{h_0, h_{x_1}, h_{x_2}, \dots, h_{x_m}\}$

The algorithm outputs

$$\hat{h} = \underset{h \in H_S}{\operatorname{argmin}} \widehat{E}R_S(h)$$

Theorem:

Let \mathbb{D} be arbitrary probability distribution over $X \times Y$.

Let A be a k -compressing algorithm.

Let S be an i.i.d. sample from \mathbb{D} of size $m > k$.

Then

$$\Pr \left[\widehat{\text{err}}_S(A(S)) = 0 \text{ and } \text{err}_D(A(S)) > \varepsilon \right] \\ \leq (1 - \varepsilon)^{m-k} \sum_{r=0}^k m^r$$

Proof:

• Let $I = (i_1, i_2, \dots, i_r)$ be
be a tuple of $r \leq k$ indices
 $i_1, i_2, \dots, i_r \in \{1, 2, \dots, m\}$

• Let $S = ((x_1, y_1), \dots, (x_m, y_m))$

• Let $S_I = ((x_{i_1}, y_{i_1}), \dots, (x_{i_r}, y_{i_r}))$

• Let ϕ be the decompressor

for A .

• Then

$$\begin{aligned} \Pr[\widehat{\text{err}}_S(\phi(S_I)) \text{ and } \text{err}_D(\phi(S_I)) > \varepsilon] \\ \leq (1 - \varepsilon)^{m - |I|} \\ \leq (1 - \varepsilon)^{m - k} \end{aligned}$$

• Union bound over all choices of I of size $|I| \leq k$

$$\begin{aligned} \Pr[\exists I : \widehat{\text{err}}_S(\phi(S_I)) \text{ and } \text{err}_D(\phi(S_I)) > \varepsilon] \\ \leq (1 - \varepsilon)^{m - k} \sum_{t=1}^k m^t \end{aligned}$$

• $A(s) = \phi(S_I)$ for some I

of size $|I| \leq k$.



Corollary: Let $\mathcal{D}, A, k, \epsilon, m$ be as before. With probability at least $1 - \delta$, $\widehat{\text{err}}_S(A(S)) = 0$ then

$$\text{err}_{\mathcal{D}}(A(S)) \leq \frac{(k+1) \ln m + \ln(1/\delta)}{m-k}$$

Proof:

$$\sum_{r=0}^k m^r \leq m^{(k+1)}$$

$$(1 - \epsilon)^{m-k} \leq e^{-\epsilon(m-k)}$$

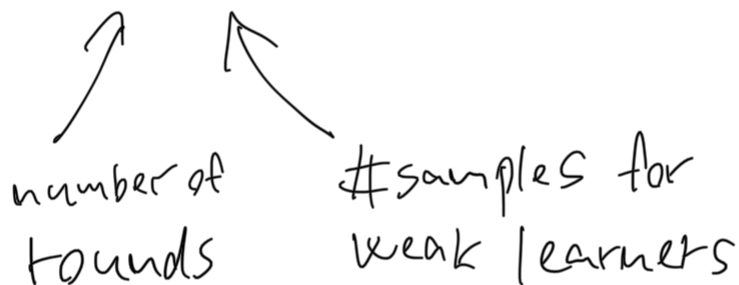
Solve

$$m^{k+1} \cdot e^{-\epsilon(m-k)} = \delta$$

for ϵ .

~~QED~~

- The sampling version of Ada Boost is a Time-compressing algorithm.



- Once the samples S_1, S_2, \dots, S_T are selected h_1, h_2, \dots, h_T are determined. In turn, these determine d_1, d_2, \dots, d_T and

$$h = \sum_{t=1}^T d_t h_t$$

Claim:

Let $\varepsilon, \delta \in (0, 1)$. Let $\eta \in (0, \frac{1}{2})$.

Let m_0 be a positive integer.

There exist positive integers T, m such that

$$T > \frac{\ln m}{2\eta^2}$$

$$\frac{(Tm_0 + 1) \ln m + \ln(T/\delta)}{m - Tm_0} < \varepsilon.$$

Proof:

• For any m , we consider

$$T := T(m) = \left\lceil \frac{\ln m}{\eta^2} \right\rceil + 1$$

$$1 - \frac{1}{2\beta^2}$$

• Clearly the first inequality is satisfied.

• The second inequality becomes

$$\frac{(T(m)m_0 + 1) \ln(m) + \ln(2/\delta)}{m - T(m)m_0} < \epsilon$$

• Clearly for large enough m the inequality holds.

~~□~~

Theorem: Let $H \subseteq \{-1, 1\}^X$ be a class of predictors.

Let $\eta \in (0, \frac{1}{2})$.

H is η -weakly PAC learnable if and only if H is PAC learnable.

Proof:

- \Leftarrow is trivial
- \Rightarrow : Suppose A is a η -weak learner for H .
- Let m_0 be sample size such such that for any distribution Q and any i.i.d. sample S' from Q labeled by any $h \in H$

$$\text{err}_{Q,h}(A(S)) \leq \eta.$$

with probability at least $1/2$.

- Let T, m be positive integers

such that

$$T > \frac{\ln m}{2\eta^2}$$

$$\frac{(Tm_0 + 1) \ln m + \ln(2/\delta)}{m - Tm_0} < \epsilon$$

(see previous lemma.)

• Let $L = \lceil \log_2(2T/\delta) \rceil$

• We run AdaBoost with subsampling and check

$$\text{err}_{D_t}(A(S_t)) \leq \eta$$

• If $\text{err}_{D_t}(A(S_t)) > \eta$, we resample S_t afresh and repeat.

... ..

- We repeat up to L times if necessary.
- If we fail after L attempts, Ada Boost fails (e.g. output all-zero classifier).
- The failure probability is at most

$$T \cdot 2^{-L} \leq T \cdot 2^{-\log_2(2T/\delta)}$$

$$= T \cdot \frac{\delta}{2T}$$

$$= \frac{\delta}{2}$$

- If the algorithm does not fail, it outputs

$$h = \sum_{t=1}^T \alpha_t h_t$$

with $\overline{\text{err}}_S(h) = 0$.

- With probability at least $1 - \frac{\delta}{2}$ the generalization error of h satisfies

$$\text{err}_D(h) \leq \frac{(Tm_0 + 1) \ln m + \ln(2/\delta)}{m - Tm_0}$$

$\leq \epsilon$

Compression theorem.

- By union bound, the algorithm fails to output h with $\widehat{\text{err}}_D(h) \leq \epsilon$ with probability at most

$$\frac{\delta}{2} + \frac{\delta}{2} = \delta$$

