

1. Project proposal due on Mar 16.

2. HW 2 out (on class website)

(latest version on Tuesday Mar. 2)

due Mar 23

Weak PAC learning: \mathcal{H} . A is a δ -weak PAC learner for \mathcal{H} , if \exists a fn $f: (0,1) \rightarrow \mathbb{N}$, for any

$\delta > 0$, for any D realizable by \mathcal{H} , w/ $m \geq f(\delta)$,
iid samples from D

A produces a classifier h , s.t. w.p. $1-\delta$:

$$\text{err}(h, D) \leq \frac{1}{2} - \delta.$$

Adaboost

if weak learner A returns classifiers from base hypothesis class B ,

for any $h^* \in H$, for any dist^(*) D_x over X ,

$$\exists f \in B. \quad \mathbb{E}_{x \sim D_x} I(f(x) \neq h^*(x)) \leq \frac{1}{2} - \gamma.$$

AdaBoost gives an algorithmic pf showing

for any $h^* \in H$. \exists distribution D_B over B .

such that $h^*(x) = \text{sign}\left(\sum_{f \in B} D_B(f) \cdot f(x)\right)$.

A more direct pf:

rewrite (*):

$$\max_{D_x} \left(\min_{f \in B} \mathbb{E}_{x \sim D_x} \frac{I(f(x) \neq h^*(x))}{2} \right) \leq \frac{1}{2} - \gamma$$

$$\Downarrow$$
$$\min_{D_x} \max_{f \in B} \mathbb{E}_{x \sim D_x} f(x) h^*(x) \geq 2\gamma.$$

A(5)

$$\min_{\substack{D_X \\ P}} \max_{\substack{D_B \in \Delta(B) \\ \mathcal{G} \text{ set of all distns on } B}} \mathbb{E}_{f \sim D_B} \mathbb{E}_{X \sim D_X} f(x) h^*(x) \quad (*)$$

Von Neumann's Minimax Theorem:

given $A \in \mathbb{R}^{m \times n}$. $\Delta^d = \{v_1 \sim v_d\} : \forall i: v_i \geq 0, \sum_i v_i = 1$
 protocol 1 modified protocol 2 modified.

$$\min_{P \in \Delta^m} \max_{\mathcal{G} \in \Delta^n} P^T A \mathcal{G} = \max_{\mathcal{G} \in \Delta^n} \min_{P \in \Delta^m} P^T A \mathcal{G}.$$

A:

		Bob		
		R	P	S
Alice	R	0	1	-1
	P	-1	0	1
	S	1	-1	0

$$\max_i \min_j A_{ij} = 1$$

protocol 1:

- Alice chooses i (row) $P \in \Delta^m$
- Bob chooses j after seeing i $\mathcal{G} \in \Delta^n$
- Alice suffer loss of A_{ij}

$$P^T A \mathcal{G} = \mathbb{E}_{i \sim P} \mathbb{E}_{j \sim \mathcal{G}} [A_{ij}] \max_j A_{ij}$$

If Alice & Bob behave optimally, what would be Alice's pay off?

$$\max_j \min_i A_{ij} = -1$$

protocol 2:

- Bob chooses j (column)
 - Alice chooses i after seeing j
 - Alice loss = A_{ij}
- $\left. \begin{array}{l} \text{min}_i A_{ij} \end{array} \right\}$

If Alice & Bob behave optimally, what would be Alice's pay off?

Applying von Neumann's theorem on (\mathcal{F}) .

$$\max_{D_B \in \Delta(B)} \min_{D_X} \mathbb{E}_{X \sim D_X} \mathbb{E}_{f \sim D_B} f(x) h^*(x) \geq 2\delta.$$

$$\max_{D_B \in \Delta(B)} \min_{\mathcal{X}} \left[\mathbb{E}_{f \sim D_B} f(x) \right] h^*(x) \geq 2\delta.$$

$\exists D_B$, for all $x \in \mathcal{X}$,

$$\sum_{f \in B} D_B(f) f(x) \cdot h^*(x) \geq 2\delta > 0$$

if we define (say B is finite)

$$\tilde{x} = (f(x))_{f \in B} \in \{-1, +1\}^N \quad N = |B|$$

$(\tilde{x}, h^*(x))_{x \in \mathcal{X}}$ is linearly separable by

$$w = (D_B(f))_{f \in \mathcal{B}} \in \mathbb{R}^N$$

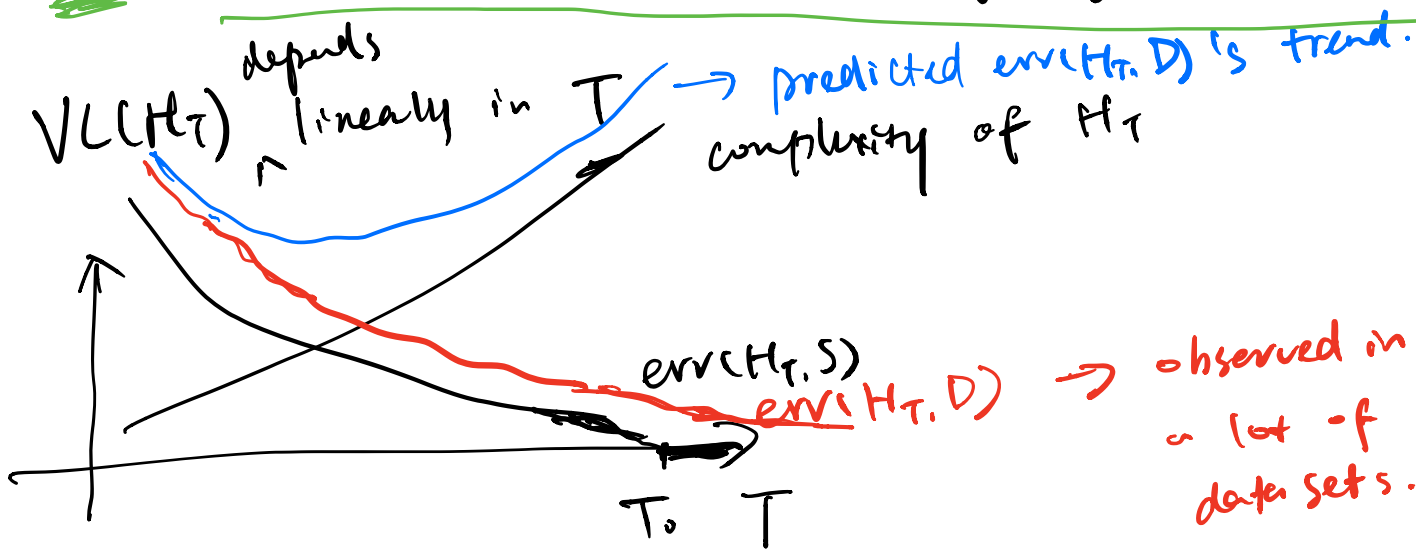
Generalization error bounds for boosting.

choose T in Ada Boost:

$$- T \uparrow. \quad \text{err}(H_T, S) \leq \exp(-2T\epsilon^2)$$

$$\text{err}(H_T, D) \leq \text{err}(H_T, S) + \sqrt{\frac{VC(H_T) T}{m}} \quad \leftarrow \epsilon = N=|\mathcal{B}|$$

$$H_T \in \left\{ \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) : \alpha_1, \dots, \alpha_T \in \mathbb{R}, \right. \\ \left. \forall t, h_t \in \mathcal{B} \right\}$$



Motivates the large margin theory for boosting.

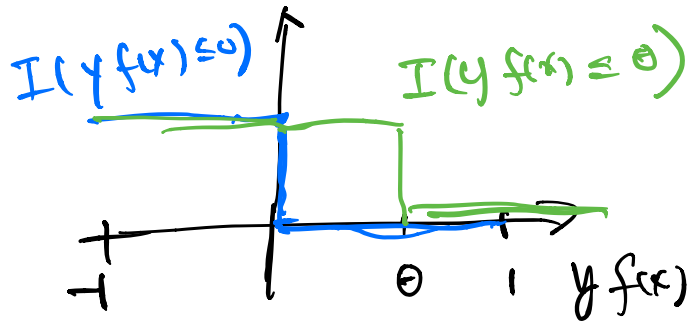
Thm: suppose \mathcal{B} is finite. $C(\mathcal{B}) = \left\{ \sum_{h \in \mathcal{B}} \alpha_h h(x) : \right.$

$\left. \sum_{h \in \mathcal{B}} |\alpha_h| \leq 1 \right\}$ is the set of voting classifier over

B. Fix margin $\theta \in (0, 1]$. Then, w.p. $1 - \delta$:
 (in iid training examples $S \sim D$): for all $f \in C(B)$:
 margin of f on (x, y) .

$$P_D (y f(x) \leq \theta) \leq P_S (|y f(x)| \leq \theta) + O\left(\frac{1}{\theta} \sqrt{\frac{\ln |B|}{m}}\right)$$

error of sign(f)



Application to Ada Boost:

$$\bar{f}_T = \frac{\sum_{t=1}^T \alpha_t h_t(x)}{\sum_{t=1}^T \alpha_t} \in C(B)$$

Let $\theta = \frac{\gamma}{2}$.

① $P_S (y f(x) \leq \frac{\gamma}{2}) \leq \exp(-T\gamma^2)$

② complexity term $O\left(\frac{1}{\gamma} \sqrt{\frac{\ln |B|}{m}}\right)$

independent of T .

Then (more abstract version) Suppose D supported on $\{x \in \mathbb{R}^d : \|x\|_\infty \leq R_\infty\} \times \{\pm 1\}$. Fix margin value $\theta \in$

[0.1]. Then, w.p. $1-\delta$ over m samples S .

for any predictor w such that $\|w\|_1 \leq B_1$,

$$P_0 (y \langle w, x \rangle \leq 0) \leq P_S (y \langle w, x \rangle \leq \theta) + \underbrace{\text{margin of } w \text{ on } (x, y)}$$

$$O \left(\frac{B_1 R_\infty}{\theta} \sqrt{\frac{\ln d/\delta}{m}} \right).$$

$$B = \{ h_1, \dots, h_N \}$$

$$\tilde{x} = (\underbrace{h_1(x)}_{h_1(x) \in \{-1,1\}} \dots h_N(x)) \Rightarrow \|\tilde{x}\|_\infty \leq 1 := R_\infty.$$

$$(x, y) \sim \tilde{D}.$$

$$\sum_{i=1}^N \alpha_i h_i(x) \leftrightarrow \langle \alpha, \tilde{x} \rangle$$

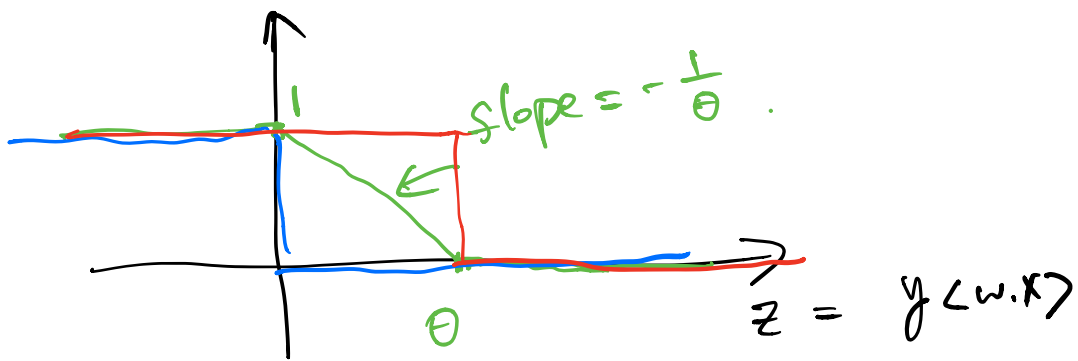
$$\|\alpha\|_1 \leq 1 \rightarrow \text{setting } B_1 = 1.$$

$$\text{choose } d = N.$$

Pf of thm :

$$\text{define ramp loss } \ell_{\theta, w}(x, y) = \phi_\theta(y \langle w, x \rangle)$$

$$\text{where } \phi_\theta(z) = \begin{cases} 1 & z \leq 0 \\ \frac{1-z}{\theta} & z \in (0, \theta) \\ 0 & z \geq \theta \end{cases}$$



We will show: w.p. $1-\delta$.

① for every $w, \|w\|_1 \leq B_1$.

$$\mathbb{E}_D \underbrace{\ell_{\theta, w}(x, y)}_{f(z)} \leq \mathbb{E}_S \underbrace{\ell_{\theta, w}(x, y)}_{f(z)} + \sqrt{\frac{\ln^2 \delta}{2m}} + 2 \cdot \text{Rad}_n(\bar{F})$$

(appears in HW2. problem 3)

where $\bar{F} = \{ \underbrace{\ell_{\theta, w}}_f : \|w\|_1 \leq B_1 \}$.

② $\mathbb{E}_D \ell_{\theta, w}(x, y) \geq \mathbb{P}_D (y \langle w, x \rangle \leq 0)$

$$\mathbb{E}_S \ell_{\theta, w}(x, y) \leq \mathbb{P}_S (y \langle w, x \rangle \leq \theta)$$

Next class:

focus on bounding

$$\text{Rad}_n(\bar{F}) \leq O\left(\frac{B_1 R_D}{\theta} \sqrt{\frac{\ln d / \delta}{m}}\right)$$

contraction inequality of Rademacher complexity