

## LAST TIME: Online Learning in Mistake-Bounded Model

For  $t=1, 2, \dots$

- Receive  $x_t$
- Output  $\hat{y}_t \in \{0, 1\}$
- Receive  $y_t \in \{0, 1\}$  mistake if  $\hat{y}_t \neq y_t$ .

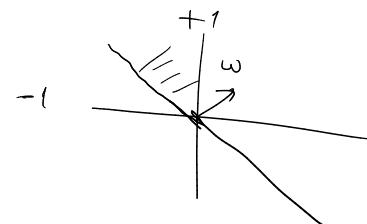
GOAL: #MISTAKES bounded even as  $t \rightarrow \infty$ .

(this requires that  $y_t = c(x_t)$  for some  $c \in C$ , where  $C$  is some concept class.)

PERCEPTRON: (Rosenblatt - 50s).

$$f: \mathbb{R}^n \rightarrow \{-1, 1\}$$

$$f_{\underline{\omega}}(x) = \begin{cases} +1 & \text{if } \underline{\omega} \cdot x \geq 0 \\ -1 & \text{if } \underline{\omega} \cdot x < 0. \end{cases}$$



Perceptron Algorithm:

$$\text{Set } \underline{\omega}_0 = \underline{0} \in \mathbb{R}^n$$

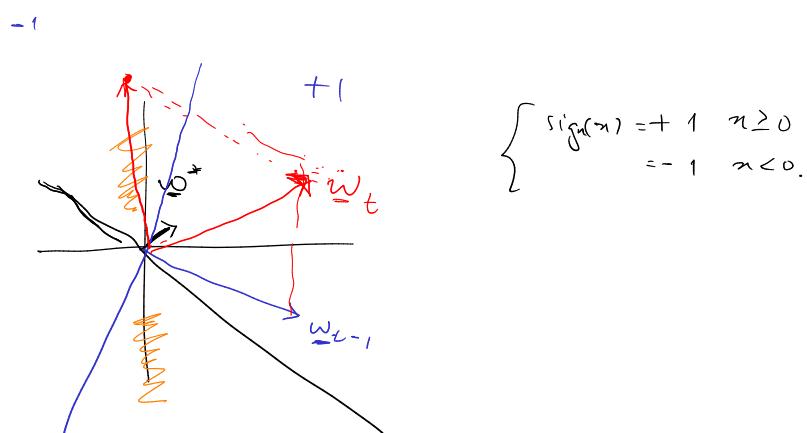
for  $t=1, 2, \dots$  do

receive  $x_t$

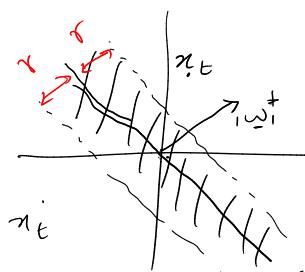
$$\hat{y}_t = \text{sign}(\underline{\omega}_{t-1} \cdot x_t)$$

if  $\hat{y}_t \neq y_t$ :

$$\underline{\omega}_t = \underline{\omega}_{t-1} + y_t x_t$$



Theorem: Suppose  $\|x_t\| \leq D \forall t$  (bounded) and  $\underline{\omega}^*$  defines the true linear separator, i.e.  $y_t = \text{sign}(\underline{\omega}^* \cdot x_t) \forall t$ . Furthermore suppose  $\|\underline{\omega}^*\| = 1$  and that  $|\underline{\omega}^* \cdot x_t| \geq \gamma > 0$  for all  $t$ ,  
 then the total #MISTAKES of perceptron  $\leq O(D^2/\gamma^2)$ .



← margin around the separator in which there is no data.

Lemma: Let  $m_t$  be the mistakes made up to time  $t$

$$(i) \quad \underline{\omega}_t \cdot \underline{\omega}^* \geq r \cdot m_t$$

$$(ii) \quad \|\underline{\omega}_t\|^2 \leq m_t \cdot D^2.$$

Proof By induction

$$(i) \quad \underline{\omega}_0 \cdot \underline{\omega}^* = 0 \geq r \cdot 0 \quad (\text{base case}).$$

if no mistake at time  $t$ ,  $\underline{\omega}_{t+1} = \underline{\omega}_t$  &  $m_{t+1} = m_t$  so clearly holds

else  $\underline{\omega}_{t+1} = \underline{\omega}_t + y_t \cdot \underline{x}_t$

$$\begin{aligned} \therefore \underline{\omega}_{t+1} \cdot \underline{\omega}^* &= \underline{\omega}_t \cdot \underline{\omega}^* + \text{sign}(\underline{\omega}_t \cdot \underline{x}_t)(\underline{\omega}^* \cdot \underline{x}_t) \\ &= \underline{\omega}_t \cdot \underline{\omega}^* + |y_t \cdot \underline{x}_t| \\ &\geq \underline{\omega}_t \cdot \underline{\omega}^* + r \geq (m_t + 1)r. \end{aligned}$$

$$|\gamma| = \text{sign}(y)$$

$$(ii) \quad \|\underline{\omega}_0\|^2 = 0 \leq 0 \cdot D^2.$$

$$\|\underline{\omega}_{t+1}\|^2 = \|\underline{\omega}_t\|^2 \leq (m_t)D^2 = (m_{t+1})D^2 \quad \text{if no mistake at time } t.$$

$$\begin{aligned} \text{else } \|\underline{\omega}_{t+1}\|^2 &= \|\underline{\omega}_t\|^2 + y_t^2 \|\underline{x}_t\|^2 + 2y_t(\underline{\omega}_t \cdot \underline{x}_t) \leq \|\underline{\omega}_t\|^2 + \|\underline{x}_t\|^2 \\ &\leq m_t D^2 + D^2 \\ &= (m_{t+1})D^2. \end{aligned}$$

Proof of Theorem:

$$m_t \cdot r \leq \underbrace{\underline{\omega}_t \cdot \underline{\omega}^*}_{\text{Cauchy-Schwarz}} \leq \|\underline{\omega}_t\| \cdot \|\underline{\omega}^*\| \leq \sqrt{m_t} \cdot D.$$

$$\Rightarrow m_t \leq \frac{D^2}{r^2} \quad \forall t.$$

### HALFSPACES / LINEAR THRESHOLD FUNCTIONS

- $VC(LTF_n) = n+1$
- Linear programming (or SVMs) → get a consistent learner. } doesn't require a margin assumption.

## MONOTONE DISJUNCTIONS :

$$f(x) = \begin{cases} 1 & \text{if } x_i, \vee x_{i_2}, \dots, \vee x_{i_k} \equiv \text{TRUE.} \\ 0 & \text{otherwise} \end{cases}$$

$$f(x) \equiv x_1 \vee \dots \vee x_k$$

$$f(\underline{x}) \equiv \text{sign} \left( \underbrace{\sum_{i=1}^k x_i - \frac{1}{2}} \right).$$

(Perception Theorem) :

$$\|\tilde{\omega}^*\|^2 = k + \frac{1}{4}$$

could scale  $\tilde{\omega}^*$  by  $\frac{1}{\sqrt{k+\frac{1}{4}}}$  to get  $\omega^*$  which is a unit vector.

$$(i) \|\underline{x}_t\|^2 \leq n, D \approx \sqrt{n+1} \quad \underline{x}_t \in \{0, 1\}^n. \quad D = \Theta(\sqrt{n})$$

$$(ii) \gamma \approx \frac{1}{\sqrt{k+\frac{1}{4}}} \left(\frac{1}{2}\right) = \Theta\left(\frac{1}{\sqrt{k}}\right). \quad \gamma = \Theta\left(\frac{1}{\sqrt{k}}\right).$$

$$\text{Perception mistake bound} = O\left(\frac{D^2}{\gamma^2}\right) = O(nk).$$

GOAL : Mistake bound of  $O(k \log n)$ .

• Sample complexity for learning disjunctions on  $k$  literals =  $\text{poly}(k, \log n)$ .

• Online learning  $\Rightarrow$  PAC learning.

( if we had  $\Theta\left(\frac{1}{\varepsilon} \log \frac{M}{\delta}\right)$  steps without making a mistake we could find a good hypothesis )

$$\# \text{time steps to simulate} \approx O\left(\frac{M}{\varepsilon} \log \frac{M}{\delta}\right).$$

# Littlestone's Algorithm / WINNOW / Weighted Majority Algorithm.

1.  $\underline{w}^0 = (1, \dots, 1) \in \mathbb{R}^n$  (Learning DISJUNCTIONS)
2. For  $t=1, 2, \dots$ , do [Littlestone '89 : Learning where irrelevant variables abound..]
- $\hat{y}_t = \mathbf{1}(\underline{w}^{t-1} \cdot \underline{x}_t \geq \frac{n}{2})$
- if MISTAKE.
- if  $\hat{y}_t = 1$  and  $y_t = 0$ , ] (A)  
set  $w_i^t = 0$  for all  $i$  s.t.  $(x_t)_i = 1$
- else ( $\hat{y}_t = 0$  &  $y_t = 1$ ) ] (B)  
set  $w_i^t = 2 \cdot w_i^{t-1}$  for all  $i$  s.t.  $(x_t)_i = 1$ .

Case (A): If  $i$  s.t.  $(x_t)_i = 1$ , then " $x_i$ " cannot be in target disjunction.

Case (B): If weight not enough to cross majority, doubles all the weights for  $i$  s.t.  $(x_t)_i = 1$ .

Claim 1: For each  $i, t$ ;  $w_i^t \leq n$

Proof: Weight increase only when  $\hat{y}_t = 0$  &  $y_t = 1$ .

for any  $i$  s.t.  $(x_t)_i = 1$ , it must be that  $(w_i^{t-1}) < y_2$ .

otherwise  $w_i^{t-1} \cdot x_t \geq \frac{n}{2}$ . and  $\hat{y}_t$  would have been 1.

if you double A, we still have  $w_i^t \leq n$ . ■

Claim 2: # Promotions  $\equiv P$  (no. of steps of type (B))

# eliminations  $\equiv E$  (no. of steps of type (A))

$$E \leq P + 2.$$

Proof:

$$0 \leq \sum_i w_i^t \leq n \quad + \quad \underbrace{\frac{P \cdot n}{2}}_{\substack{\text{every promotion} \\ \text{at time} \\ 0}} - \underbrace{E \cdot \frac{n}{2}}_{\substack{\text{every elimination} \\ \text{step decreases total} \\ \text{wt by at most } \frac{n}{2}}}$$

every elimination step decreases total wt by at least  $\frac{n}{2}$ . ■

Claim 3:  $P \leq k \log_2 n \quad \forall t$ . ■

Proof: Each promotion step doubles the weight of at least one variable in the target disjunction.

$$\text{Mistakes} = P + E \leq 2P + 2 \leq 2k \log_2 n + 2. \quad \square$$