

544NA

9/24/2013

LINEAR CLASSIFIER: $a_k^{\wedge} = \vec{w}^T \vec{\phi}^n$ TRAINING CRITERIA

① LINEAR MSE $\mathcal{E} = \sum_{n=1}^N |t^n - a_k^{\wedge}|^2$

$$\Rightarrow \vec{w}^T = \Phi^T T = (\Phi^T \Phi)^{-1} \Phi^T T$$

$$\vec{w} = (\Phi^T \Phi)^{-1} \Phi^T \vec{t}$$

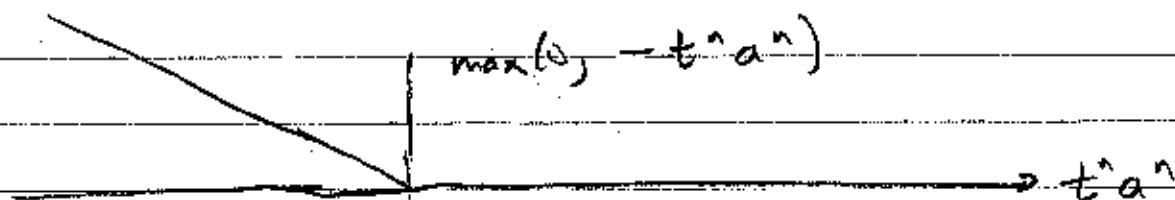
② NONLINEAR MSE $\mathcal{E} = \sum_{n=1}^N |t^n - h(a^n)|^2$

$$w_j = w_j - \eta \frac{\partial \mathcal{E}}{\partial w_j}$$

$$= w_j + \eta (t^n - h(a^n)) h'(a^n) \phi_j^n$$

③ PERCEPTRON $\mathcal{E} = \sum_{n=1}^N \max(0, -t^n a^n)$

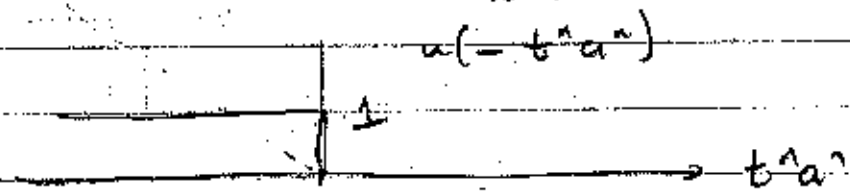
$t^n a^n \begin{cases} > 0 \\ < 0 \end{cases} \begin{cases} \text{IF } t^n, a^n \text{ SAME SIGN} \\ \text{IF OPPOSITE SIGN} \end{cases}$



TOKEN $\vec{\phi}^n$ IS MISCLASSIFIED

IFF $\text{sign}(t^n) \neq \text{sign}(a^n)$ UNIT STEP

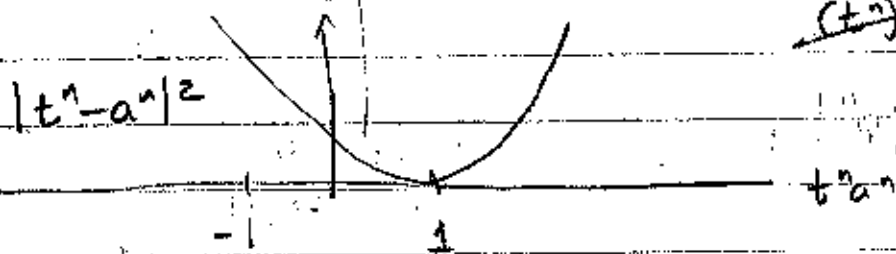
$$\Leftrightarrow (\# \text{ ERRORS}) = \sum_{n=1}^N u(-t^n a^n)$$



LINEAR MSE: IF $t^n \in \{-1, 1\}$, THEN

$$|t^n - a^n|^2 = 1 - 2t^n a^n + (a^n)^2$$

$$= 1 + 2t^n a^n + \frac{(t^n a^n)^2}{(t^n)^2} \rightarrow 1$$



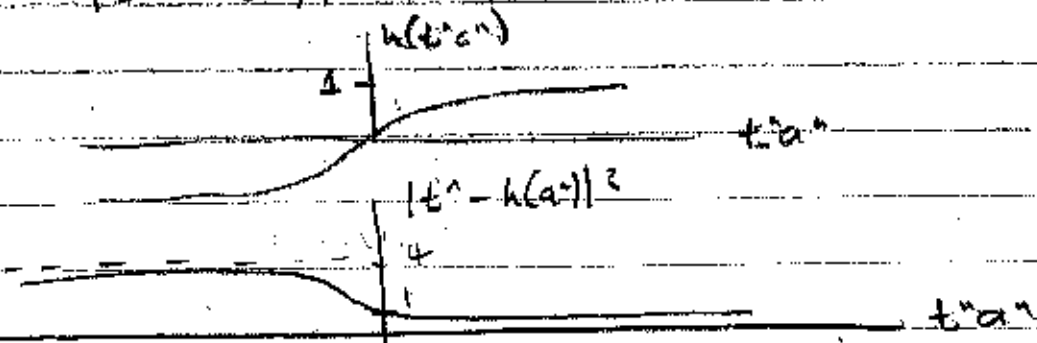
NONLINEAR MSE: $|t^n - h(a^n)|^2$

$$= 1 - 2t^n h(a^n) + h^2(a^n)$$

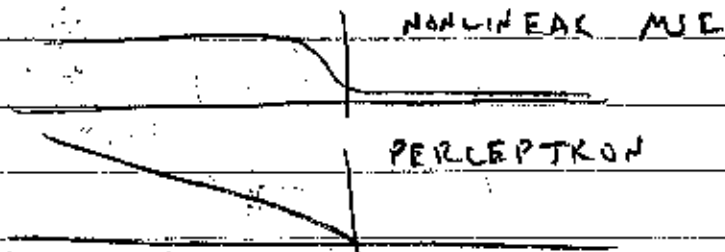
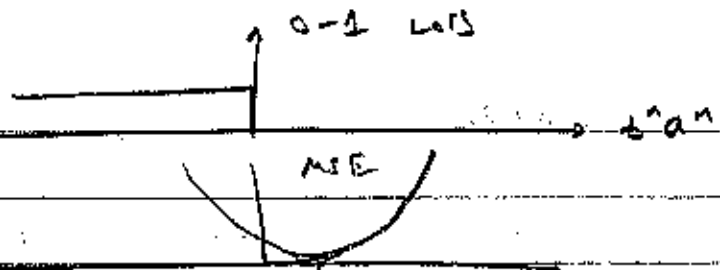
BUT $h(-a^n) = -h(a^n)$

$$h(t^n a^n) = t^n h(a^n)$$

$$|t^n - h(a^n)|^2 = 1 - 2h(t^n a^n) + h^2(t^n a^n)$$



LOSS FUNCTIONS



ADVANTAGES

DISADVANTAGES

0-1 LOSS

IT'S WHAT YOU REALLY CARE ABOUT

NOT CONVEX \Rightarrow NO GLOBAL OPTIMUM
NOT DIFFERENTIABLE \Rightarrow NO GRADIENT SOLUTION EITHER!

LINEAR MSE

CONVEX \Rightarrow HAS GLOBAL OPTIMUM

VERY GOOD TOKENS AS BAD AS VERY BAD TOKENS

NONLINEAR MSE

CLOSE TO 0-1 LOSS; DIFFERENTIABLE

"VERY BAD TOKENS" HAVE $h'(a_n^k) \approx 0 \Rightarrow$ GRADIENT DESCENT "GIVES UP"

PERCEPTRON

FIXES THIS PROBLEM; GRADIENT DESCENT IS VERY EASY!!

GRADIENT DESCENT :

$$\varepsilon = \sum_{n=1}^N \max(0, -t^n a^n)$$

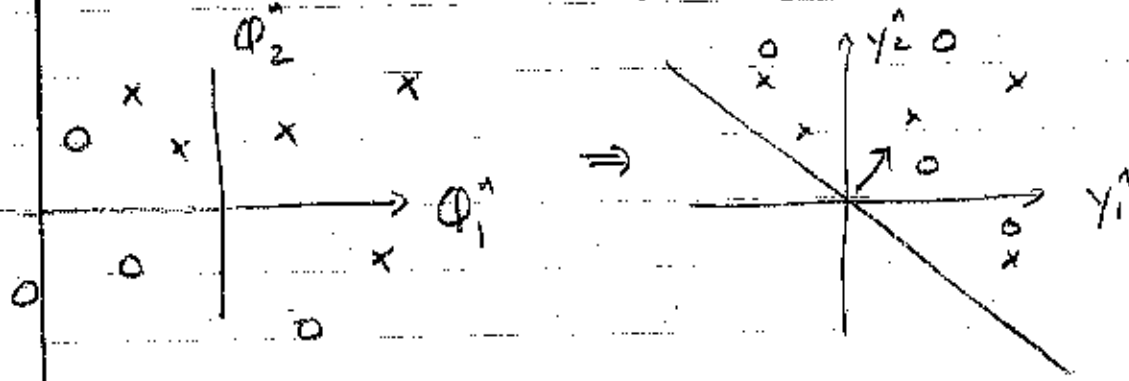
$$= \sum_{n=1}^N \max\left(0, -t^n \sum_{j=1}^M w_j \phi_j\right)$$

$$\frac{\partial}{\partial w_j} \max(0, -t^n a^n) = \begin{cases} 0 & t^n a^n \geq 0 \\ -t^n \phi_j^n & t^n a^n < 0 \end{cases}$$

$$\nabla_{\vec{w}} \varepsilon = - \sum_{n: t^n a^n < 0} t^n \vec{\phi}^n$$

GEOMETRIC INTERPRETATION

$$\text{LET } \vec{y}^n = t^n \vec{\phi}^n = \begin{cases} \vec{\phi}^n & \text{IF } t^n = +1 \\ -\vec{\phi}^n & \text{IF } t^n = -1 \end{cases}$$



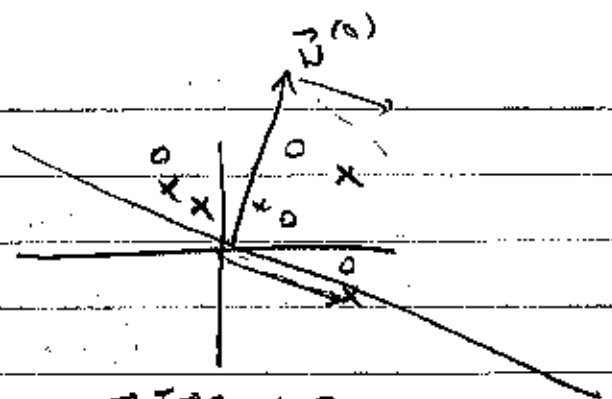
o: $t^n = -1$ ZERO ERROR:

x: $t^n = +1$ FIND \vec{w} SUCH THAT

$$\vec{w}^T \vec{y}^n > 0 \quad \text{FOR ALL } n$$

PERCEPTRON :

INITIALIZE: $\vec{w} = \sum \vec{y}^n$

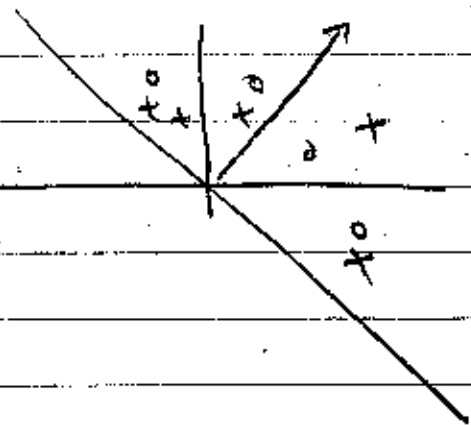
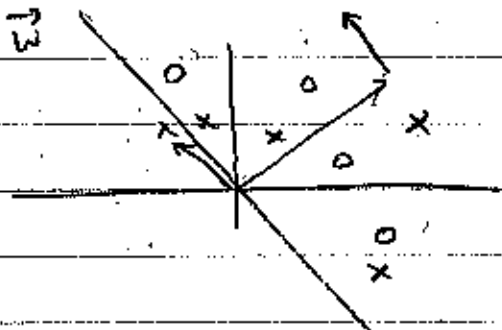


ITERATE

(a) FIND THE \vec{y}^n S.T. $\vec{w}^T \vec{y}^n < 0$

(b) ADD THEM TO \vec{w}

$$\vec{w} = \vec{w} + \sum_{\substack{n: \\ \vec{w}^T \vec{y}^n < 0}} \vec{y}^n$$



TERMINATE

WHEN ALL \vec{y}^n
CORRECTLY CLASSIFIED

CONVERGENCE =

- \vec{w} KEEPS GETTING BIGGER
- PROCESS TERMINATES IF DATA ARE LINEARLY SEPARABLE

• IF NOT, SET A THRESHOLD w_{max}
STOP WHEN $\|\vec{w}\|$ PASSES w_{max}

PROOF OF CONVERGENCE

LINEARLY SEPARABLE DATA

$$\Leftrightarrow \exists \alpha \text{ SUCH THAT } \text{sign}(\alpha^T \vec{\phi}^n) = t^n$$

$$\Leftrightarrow \hat{w}^T \vec{y}^n > 0 \text{ FOR ALL } n,$$

$$\|\hat{w}\| = 1$$

BUT NOTICE, AT THE T ITERATION,

$$\vec{w}_{(T+1)} = \vec{w}_{(T)} + \sum_{n: \hat{w}_{(T)}^T \vec{y}^n < 0} \vec{y}^n \quad \text{!} < 0 !!$$

$$\|\vec{w}_{(T+1)}\|^2 = \|\vec{w}_{(T)}\|^2 + \left\| \sum_{\hat{w}_{(T)}^T \vec{y}^n < 0} \vec{y}^n \right\|^2 + 2 \vec{w}_{(T)}^T \sum_{\hat{w}_{(T)}^T \vec{y}^n < 0} \vec{y}^n$$

$$\leq \|\vec{w}_{(T)}\|^2 + \|\vec{z}\|^2$$

SO $\|\vec{w}_{(T)}\| \leq \sqrt{T} \|\vec{z}\|, \quad \vec{z} = \sum_{n: \hat{w}_{(T)}^T \vec{y}^n < 0} \vec{y}^n$

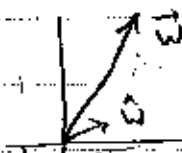
$T = \# \text{ ITERATIONS THAT HAVE ANY ERRORS}$

BUT NOTICE $\vec{w}_{(T)} = \sum_{n \in T} M^n \vec{y}^n$

$M^n = \# \text{ TIMES } \vec{y}^n \text{ HAS BEEN MISCLASSIFIED}$

$$\hat{w}^T \vec{w}_{(T)} = \sum_{n=1}^N M^n (\hat{w}^T \vec{y}^n) \geq N \tau \min_n (\hat{w}^T \vec{y}^n)$$

BUT $\hat{w}^T \vec{w}_{(T)} = \|\vec{w}_{(T)}\| \cos \theta$
 $\leq \|\vec{w}_{(T)}\|$



SO

$$N \tau \min_n (\hat{w}^T \vec{y}^n) \leq \sqrt{T} \|\vec{z}\| \quad \therefore \sqrt{T} \leq \frac{\|\vec{z}\|}{N \min_n (\hat{w}^T \vec{y}^n)}$$