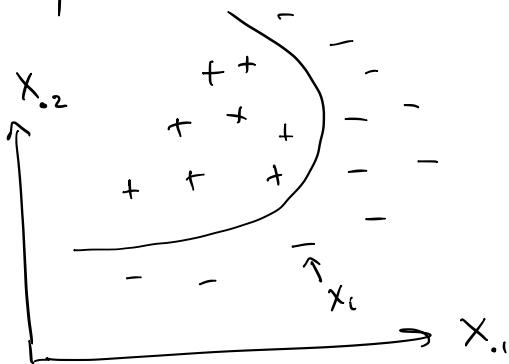


Separation happens - points can often be separated by a hyperplane in high dimensions.

If you want curvy decision boundaries, use a trick - map to a higher dim space



$$\vec{x}_i = [x_{i1}, x_{i2}] \rightarrow [x_{i1}, x_{i2}, x_{i1}^5] = x_i^{\text{new}}$$
$$f(\vec{x}_i) = \vec{\omega} \cdot \vec{x}_i \rightarrow f^{\text{new}}(\vec{x}_i) = \vec{\omega}^{\text{new}} \cdot x_i^{\text{new}}$$
$$= \omega_1 x_{i1} + \omega_2 x_{i2}$$
$$= \omega_1^{\text{new}} x_{i1} + \omega_2^{\text{new}} x_{i2} + \omega_3^{\text{new}} x_{i1}^5$$

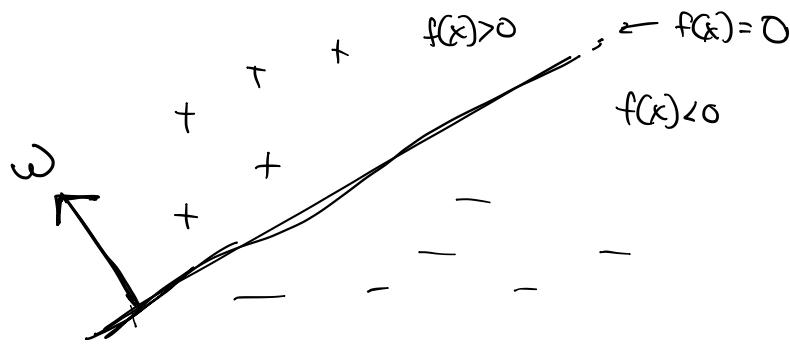
Intercepts are easy

$$\vec{x}_i = [x_{i1}, x_{i2}] \rightarrow \tilde{x}_i = [x_{i1}, x_{i2}, 1]$$

$$\begin{aligned} f(\vec{x}) &= \omega \cdot x + b \\ &= \omega_1 x_{i1} + \omega_2 x_{i2} + b \end{aligned}$$
$$\begin{aligned} \tilde{f}(x_i) &= \tilde{\omega} \cdot \tilde{x}_i \\ &= \tilde{\omega}_1 x_{i1} + \tilde{\omega}_2 x_{i2} + \underbrace{\tilde{\omega}_3 \cdot 1}_b \\ &= \omega \cdot x + b \end{aligned}$$

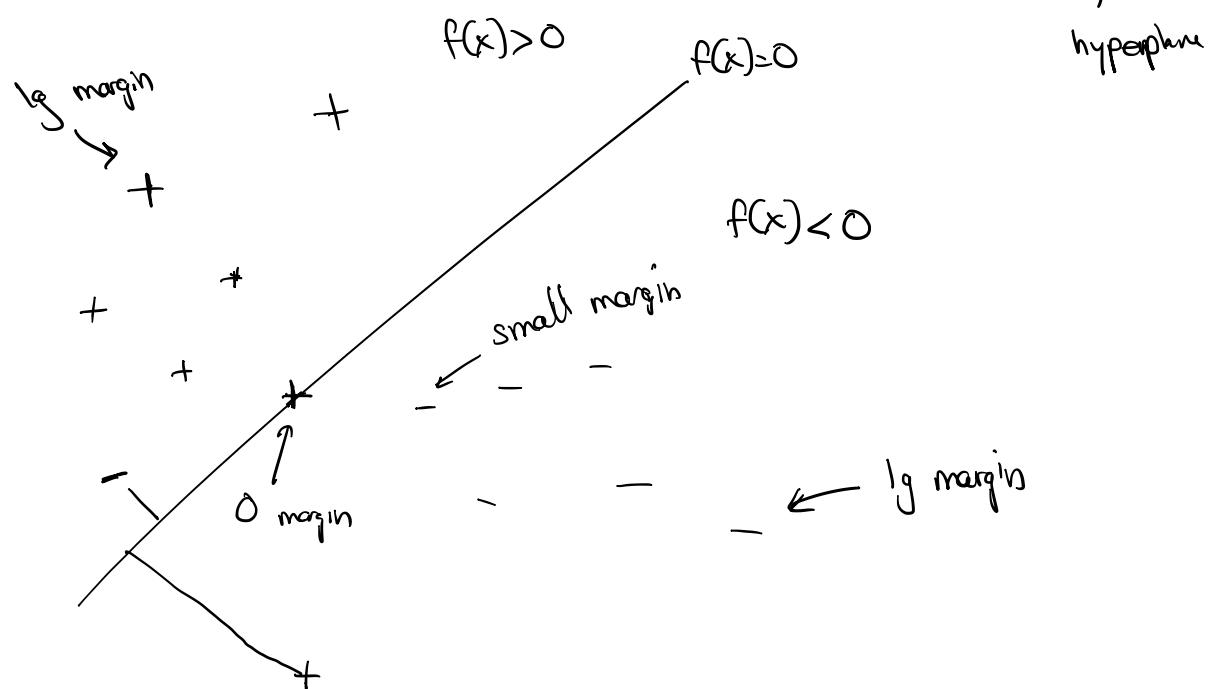
A hyperplane is represented by a vector that is perpendicular to it.

$$\hat{y} = \text{sign}(f(x)) = \underset{f \text{ is hyperplane}}{\text{sign}}(\vec{\omega} \cdot \vec{x})$$



decision bdry is $\{x : \vec{\omega} \cdot \vec{x} + b = 0\}$

Margins - how far away we are from the decision boundary
in the correct direction. margin = $y \cdot f(x) = y \cdot (\vec{w} \cdot \vec{x} + b)$

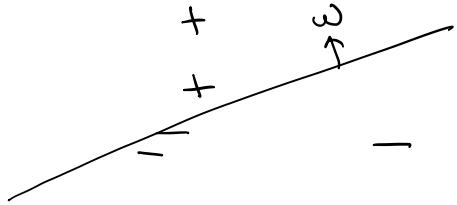


Perceptron idea

data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$

if point is correctly classified, does nothing

if point is incorrectly classified, move decision boundary towards it



$\vec{\omega}^{(t)}$ ~ model at time t

if i is correct $\text{sign}(\vec{\omega}^{(t)} \cdot \vec{x}_i) \cdot y_i = 1$

if i is misclassified $\text{sign}(\vec{\omega}^{(t)} \cdot \vec{x}_i) \cdot y_i = -1$

$$\text{and } y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i) < 0$$

want $y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i) > 0$ for all i

if i is misclassified at t

$$\vec{\omega}^{(t+1)} = \vec{\omega}^{(t)} + y_i \vec{x}_i$$

helpful because

$$y_i (\vec{\omega}^{(t+1)} \cdot \vec{x}_i) = y_i ((\vec{\omega}^{(t)} + y_i \vec{x}_i) \cdot \vec{x}_i) = y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i) + y_i (y_i \vec{x}_i \cdot \vec{x}_i)$$

want this large

$$= \underbrace{y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i)}_{\text{negative}} + \underbrace{\|\vec{x}_i\|^2}_{\text{positive}}$$

Perceptron Algorithm

input : $\{(\vec{x}_i, y_i)\}_{i=1}^n$

initialize $\vec{\omega}^{(0)} = (0, 0, \dots, 0)$

for $t=1, 2, \dots$

 locate an i such that $y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i) \leq 0$.

 if none, stop and output $\vec{\omega}^{(t)}$

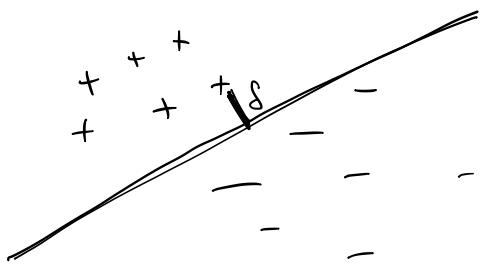
 else $\vec{\omega}^{(t+1)} = \vec{\omega}^{(t)} + y_i \vec{x}_i$

 end

end

Perceptron Convergence Bound

Assume $\{(\vec{x}_i, y_i)\}_{i=1}^n$ are separable by linear functions, so there is a good separator $\vec{\omega}$ st. $\|\vec{\omega}\|_2 = 1$ and $y_i(\vec{\omega} \cdot \vec{x}_i) \geq \delta$ for all i . Assume $\max_i \|x_i\|_2 = 1$ (normalize)



Theorem: The perceptron algorithm makes at most $\frac{1}{\delta^2}$ mistakes

Define ω^* to be a "good separator" $y_i(\vec{\omega}^* \cdot \vec{x}_i) \geq \delta \forall i$ and $\|\vec{\omega}^*\| = 1$

at T iterations, show $\vec{\omega}^{(T+1)}$ is "close to" ω^* . $\Rightarrow T \leq \frac{1}{\delta^2}$

$$\text{show } \sqrt{T} \leq \frac{\vec{\omega}^* \cdot \vec{\omega}^{(T+1)}}{\|\vec{\omega}^*\| \|\vec{\omega}^{(T+1)}\|} \stackrel{\cos \leq 1}{\leq} 1.$$

$$\sqrt{T} \leq 1$$

$$\sqrt{T} \leq \frac{1}{\delta} \text{ so } T \leq \frac{1}{\delta^2}$$

Step 1:

$$(\vec{\omega}^* \cdot \vec{\omega}^{(t+1)}) - (\vec{\omega}^* \cdot \vec{\omega}^{(t)}) = \vec{\omega}^* (\omega^{(t+1)} - \omega^{(t)}) = \vec{\omega}^* \cdot y_i \vec{x}_i$$

↑
mistake
at +

$$= y_i (\vec{\omega}^* \cdot \vec{x}_i) \geq \delta$$

$$(\omega^* \cdot \omega^{(T+1)}) = (\vec{\omega}^* \cdot \vec{\omega}^{(T+1)}) - (\vec{\omega}^* \cdot \omega^{(1)})$$

↑
(0, 0, 0)

$$= \sum_{t=1}^T (\vec{\omega}^* \cdot \vec{\omega}^{(t+1)}) - (\vec{\omega}^* \cdot \vec{\omega}^{(t)}) \geq T\delta$$

Step 2: for each t

$$\|\vec{\omega}^{(t+1)}\|^2 = \|\vec{\omega}^{(t)} + y_i \vec{x}_i\|^2 = \|\vec{\omega}^{(t)}\|^2 + 2y_i (\vec{\omega}^{(t)} \cdot \vec{x}_i) + y_i^2 \|\vec{x}_i\|^2$$

$$\|\vec{\omega}^{(t+1)}\|^2 \leq \|\vec{\omega}^{(t)}\|^2 + 1.$$

$$\text{so } \|\vec{\omega}^{(T+1)}\|^2 \stackrel{\|\vec{\omega}\|=0}{\leq} T \Rightarrow \|\vec{\omega}^{(T+1)}\| \leq \sqrt{T}.$$

Step 3: combine

$$1 \geq \frac{\vec{\omega}^* \cdot \vec{\omega}^{(T+1)}}{\|\vec{\omega}^*\| \|\vec{\omega}^{(T+1)}\|} \geq \frac{T\delta}{\sqrt{T}} = \sqrt{T} \delta \quad \begin{cases} \delta \sqrt{T} \leq 1 \\ \sqrt{T} \leq 1/\delta \\ T \leq 1/\delta^2 \end{cases} \quad \square$$

Winnow Algorithm (Littlestone, 1988)

Input: $\{(x_i, y_i)\}_{i=1}^n$, parameter $\eta > 0$

Initialize: $w_j^{(0)} = \frac{1}{p}$ $j=1 \dots p$ (p features)

For $t = 0, 1, 2, \dots$

 locate i s.t. $y_i(\omega^{(t)} \cdot x_i) \leq 0$.

 if none, stop and output $\omega^{(t)}$

$$\text{else } \forall j \quad w_j^{(t+1)} = \frac{w_j^{(t)} e^{\eta y_i x_{ij}}}{Z_t}$$

end

$$Z_t \leftarrow \text{normalization}$$

$$Z_t = \sum_j w_j^{(t)} e^{\eta y_i x_{ij}}$$

idea:

reward features that agree with label:

$$\text{if } \underbrace{\text{sign}(x_{ij})}_{} = y_i \text{ then } w_j^{(t+1)} \propto w_j^{(t)} \cdot e^{\eta (+ve)}$$

punish features that disagree

$$\text{if } y_i \neq \text{sign}(x_{ij}) \text{ then } w_j^{(t+1)} \propto w_j^{(t)} e^{\eta (-ve)}$$

Winnow Algorithm (Littlestone 1988)

Input : $\{(\bar{x}_i, y_i)\}_{i=1}^n$, parameter $\eta > 0$

Initialize $\omega_j^{(0)} = \frac{1}{p}$ $j=1 \dots p$ (p features)

For $t=0, 1, 2 \dots$

Locate i s.t. $y_i(\omega^{(t)} \cdot \bar{x}_i) \leq 0$.

If none, stop and output $\omega^{(t)}$

$$\text{else } \forall i, \omega_j^{(t+1)} = \frac{\omega_j^{(t)}}{Z_t} e^{\eta y_i x_{ij}}$$

$Z_t \uparrow$

normalization

$$Z_t = \sum_j \omega_j^{(t)} e^{\eta y_i x_{ij}}$$

End

Idea:

reward features that agree with the label

→ if $\text{sign}(x_{ij}) = y_i$ then $\omega_j^{(t+1)} \propto \omega_j^{(t)} \cdot e^{\eta (+ve)}$

→ if $\text{sign}(x_{ij}) \neq y_i$ then $\omega_j^{(t+1)} \propto \omega_j^{(t)} \cdot e^{\eta (-ve)}$

punish features that disagree with y_i

Winnow Convergence Bound

Assume $\max_i \|\mathbf{x}_i\|_\infty = 1$ (normalize)

Assume there is ω^* s.t. $y_i (\vec{\omega}^* \cdot \vec{x}_i) \geq \delta \quad \forall i$
 separates

and $\omega_j^* \geq 0 \quad \forall j$ and $\|\omega^*\|_1 = 1$

Theorem The Winnow Alg makes at most

$$T \leq \frac{\ln p}{n\delta + \ln \left(\frac{e^n}{e^{-n}} \right)} \text{ mistakes.}$$

proof idea: Show $\vec{\omega}^{(t)}$ is closer to $\vec{\omega}^*$ at each iteration in terms of KL divergence.

$$KL(\vec{a} \parallel \vec{b}) = \sum_i a_i \log \left(\frac{a_i}{b_i} \right)$$

Why is it a "distance"? If $\vec{a} = \vec{b}$, $\log \left(\frac{a_i}{b_i} \right) = 0$
 so $KL(\vec{a}, \vec{b}) = 0$

Turns out $KL(\vec{a}, \vec{b}) \geq 0$

$$\Phi_t = KL(\omega^* \parallel \omega^{(t)}) = \sum_j \omega_j^* \ln \left(\frac{\omega_j^*}{\omega_j^{(t)}} \right)$$

$$\Phi_t - \Phi_{t+1} = \sum_j \omega_j^* \ln \left(\frac{\omega_j^*}{\omega_j^{(t)}} \right) - \sum_j \omega_j^* \ln \left(\frac{\omega_j^*}{\omega_j^{(t+1)}} \right)$$

$$= \sum_j \omega_j^* \ln \left(\frac{\omega_j^{(t+1)}}{\omega_j^{(t)}} \right)$$

$$\stackrel{\text{alg}}{=} \sum_j \omega_j^* \ln \left(\frac{e^{\eta y_i x_{ij}}}{Z_t} \right)$$

$$= \sum_j \omega_j^* \underbrace{\ln(e^{\eta y_i x_{ij}})} - \sum_j \omega_j^* \ln Z_t$$

$$= \sum_j \omega_j^* \eta y_i x_{ij} - \ln Z_t \sum_j \omega_j^*$$

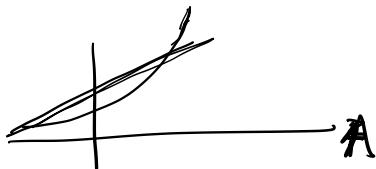
$$= \underbrace{\eta y_i}_{\text{VI} \leftarrow \text{def } \omega^*} \overrightarrow{\omega} \cdot \overrightarrow{x}_i - \ln Z_t$$

$$\geq \eta f - \ln(Z_t)$$

$$Z_t = \sum_j w_j^{(t)} e^{\eta y_i x_{ij}}$$

Use a bound of an exponential by a linear function

$$e^{\eta A} \leq \left(\frac{1+A}{2}\right) e^{\eta} + \left(\frac{1-A}{2}\right) e^{-\eta} \quad \text{holds } -1 \leq A \leq 1$$



$$Z_t \leq \sum_j w_j^{(t)} \left(\frac{1 + y_i x_{ij}}{2} \right) e^{\eta} + w_j^{(t)} \left(\frac{1 - y_i x_{ij}}{2} \right) e^{-\eta}$$

$$A = y_i x_{ij}$$

$$= \frac{e^{\eta} + e^{-\eta}}{2} \sum_j w_j^{(t)} \quad 1$$

$$+ \frac{e^{\eta} - e^{-\eta}}{2} \sum_j w_j^{(t)} y_i x_{ij}$$

$\underbrace{y_i (\overline{w}^{(t)} \cdot \bar{x}_i)}$

All \leftarrow is misclassified at t

$$\leq \frac{e^{\eta} + e^{-\eta}}{2}$$

$$-\ln(Z_t) \geq -\ln\left(\frac{e^{\eta} + e^{-\eta}}{2}\right)$$

Winnow Convergence Bound

Assume $\max_i \|x_i\|_\infty = 1$ (normalize)

Assume there is ω^* s.t. $y_i (\underbrace{\bar{\omega}^* \cdot \vec{x}_i}_{\text{separated}}) \geq \delta \quad \forall i$

and $\omega_j^* \geq 0 \quad \forall j$ and $\|\omega^*\|_1 = 1$

Theorem the Winnow Alg makes at most

$$T \leq \frac{\ln p}{n\delta + \ln(\frac{e^n + e^{-n}}{2})} \quad \text{mistakes.}$$

$$\text{So far: } \underline{\Phi}_t - \bar{\Phi}_{t+1} \geq \gamma \delta - \ln(z_t) \geq \gamma \delta - \ln\left(\frac{e^n + e^{-n}}{2}\right)$$

$$\underline{\Phi}_0 - \bar{\Phi}_T = \sum_{t=0}^{T-1} (\underline{\Phi}_t - \bar{\Phi}_{t+1}) \geq T \Delta \quad \text{lower bound}$$

$$\underline{\Phi}_0 - \bar{\Phi}_T \leq \underline{\Phi}_0 = \text{KL}(\bar{\omega}^* \parallel \bar{\omega}_0) = \sum_j \omega_j^* \ln\left(\frac{\omega_j^*}{1/p}\right)$$

↑
KL divergence
is always non negative

$$= \sum_j \omega_j^* \ln(p \omega_j^*) \stackrel{\omega_j^* \in [0, 1]}{\leq} \sum_j \omega_j^* \ln(1) \leq \ln p \sum_j \omega_j^*$$

$$\text{So } T \Delta \leq \ln p \Rightarrow T \leq \frac{\ln p}{\Delta} \quad \text{II}$$

Winnow Convergence Bound

Assume $\max_i \|\mathbf{x}_i\|_\infty = 1$ (normalize)

Assume there is ω^* s.t. $y_i (\underbrace{\bar{\omega}^* \cdot \bar{\mathbf{x}}_i}_{\text{separates}}) \geq \delta \quad \forall i$

and $\omega_j^* \geq 0 \quad \forall j$ and $\|\omega^*\|_1 = 1$

Theorem The Winnow Alg makes at most

$$T \leq \frac{\ln p}{n\delta + \ln \left(\frac{2}{e^\eta + e^{-\eta}} \right)} \text{ mistakes.}$$

What value for η ? Set it to minimize bound

$$\text{bound}(\eta) = \frac{\ln p}{\Lambda(\eta)}$$

$$\frac{d(\text{bound}(\eta))}{d\eta} = 0$$

$$\eta = \frac{1}{2} \ln \left(\frac{1+\delta}{1-\delta} \right)$$

Corollary: For Winnow if $\eta = \frac{1}{2} \ln \left(\frac{1+\delta}{1-\delta} \right)$ then

$$\# \text{mistakes} = T \leq \frac{\ln p}{n\delta + \ln \left(\frac{2}{e^\eta + e^{-\eta}} \right)} = \frac{2 \ln p}{\delta^2}.$$

Corollary : For Winnow if $\eta = \frac{1}{2} \ln\left(\frac{1+\gamma}{1-\gamma}\right)$ then

$$\# \text{mistakes} = T \leq \frac{\ln p}{n\gamma + \ln\left(\frac{2}{e^{\eta} + e^{-\eta}}\right)} = \frac{2 \ln p}{\gamma^2}$$

how to compare with perceptron?

perceptron

$$T \leq \frac{1}{\gamma^2}$$

Winnow

$$T \leq \frac{2 \ln p}{\gamma^2}$$

not a big deal
very similar

additive

$$\tilde{\omega}^{(t+1)} = \tilde{\omega}^{(t)} + \dots$$

$$\|\tilde{x}_i\|_2 \leq 1 \quad \forall i$$

$$\|\tilde{\omega}^*\|_2 = 1$$

SVM-like

multiplicative

$$\tilde{\omega}^{(t+1)} = \tilde{\omega}^{(t)}.$$

$$\|x_i\|_\infty \leq 1 \quad \forall i$$

$$\|\tilde{\omega}^*\|_1 \leq 1$$

AdaBoost