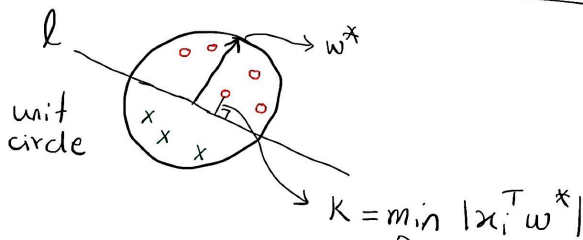


1. if the data is linearly separable, then there exists a w^* such that $\forall (x_i, y_i) \in D : y_i (x_i^T w^*) > 0$

For convenience, suppose that we re-scale our data such that $\|w^*\| = 1$ and $\|x_i\| \leq 1 \forall x_i \in D$



we define k as above. \uparrow (min distance of l to the data points)

\Rightarrow claim: perceptron algorithm, make at most $\frac{1}{k^2}$ mistakes until convergence.

proof of the claim:

from the algorithm definition we have:

$\begin{cases} y(x^T w) < 0 \rightsquigarrow \text{we need update} \\ y(x^T w^*) > 0 \rightsquigarrow w^* \text{ makes a separating line.} \end{cases}$

consider on update : $w \rightarrow w + yn$ (*)

$$(*) \Rightarrow (w + yn)^T w^* = w^T w^* + y(n^T w^*) \geq w^T w^* + \underline{K}$$

[because $y(n^T w^*) = |n^T w^*| \geq K$]

\Rightarrow after each update $w^T w^*$ increases at least K .

①

$$(*) \Rightarrow (w + yn)^T (w + yn) = w^T w + \underbrace{2y(w^T n)}_{< 0} + y^2(n^T n)$$

[because we have update]

$y^2 = 1$ and $\|n\| \leq 1$

$$\Rightarrow (w + yn)^T (w + yn) \leq w^T w + 1$$

\Rightarrow after each update $w^T w$ increases at most 1

②

from ① and ②, after i updates we have :

$$\begin{cases} w^T w^* \geq iK \\ w^T w \leq i \end{cases}$$

$$\Rightarrow i.k \leq w^T w^* = \|w\| \cos \theta \quad [\theta \text{ is angle between } w \text{ and } w^*]$$

$$\leq \|w\| = \sqrt{w^T w} \leq \sqrt{i}$$

$$\Rightarrow i.k \leq \sqrt{i} \Rightarrow i \leq \frac{1}{k^2}$$

\Rightarrow number of updates is at most $\frac{1}{k^2}$

\Rightarrow we will converge finally \checkmark

2. Suppose we have $\|\theta_{MAP}\|_2 > \|\theta_{ML}\|_2$

and we know that:

$$P(\theta_{MAP}) = \frac{1}{(2\pi)^{\frac{n+1}{2}} \times \sqrt{|\tau^2 I|}} \times \exp\left(-\frac{1}{2\tau^2} (\|\theta_{MAP}\|_2)^2\right)$$

by assumption $< \frac{1}{(2\pi)^{\frac{n+1}{2}} \sqrt{|\tau^2 I|}} \exp\left(-\frac{1}{2\tau^2} (\|\theta_{ML}\|_2)^2\right)$

$$= P(\theta_{ML})$$

$$\left[\text{because } \theta_{ML} \text{ maximizes } \right] \leq p(\theta_{ML}) \prod_{i=1}^n p(y^{(i)} | x^{(i)}, \theta_{ML})$$

there for $\|\theta_{\text{map}}\|_2 \leq \|\theta_{m_L}\|_2$ ✓

proof: assume that each sample is a d -dimensional vector like: $x = (x_1, x_2, \dots, x_d)^T$.

then if $P(y=1|X) \geq P(y=0|X) \Rightarrow \text{predict 1}$

if $P(y=0|X) > P(y=1|X) \Rightarrow \text{predict } 0$

By Bayes rule, if $\frac{P(X|y=1)P(y=1)}{P(X|y=0)P(y=0)} \geq 1$ we will predict 1 for the label. $\underbrace{\hspace{10em}}_{(*)}$

Because we have the Naive Bayes assumption, we have $P(X|y) = \prod_{i=0}^d P(x_i|y)$. So if we re-write $(*)$, we have:

$$\frac{P(y=1)}{P(y=0)} \times \prod_{i=0}^d \frac{P(x_i|y=1)}{P(x_i|y=0)} \geq 1$$

Let $P(y=1) = p$ and $P(x_i=1|y=1) = q_i$
 $P(y=0) = 1-p$ and $P(x_i=1|y=0) = q'_i$

$$\Rightarrow P(x_i|y=1) = q_i^{x_i} (1-q_i)^{1-x_i}$$

$$P(x_i|y=0) = q'_i{}^{x_i} (1-q'_i)^{1-x_i}$$

$$\Rightarrow (*)^2 = \frac{p}{1-p} \times \prod_{i=0}^d \frac{q_i^{x_i} (1-q_i)^{1-x_i}}{q'_i{}^{x_i} (1-q'_i)^{1-x_i}} \geq 1$$

\Leftrightarrow

$$\left(\frac{p}{1-p} \times \prod_{i=0}^d \frac{1-q_i}{1-q'_i} \right) \cdot \prod_{i=0}^d \left(\frac{q_i}{q'_i} \cdot \frac{1-q'_i}{1-q_i} \right)^{x_i} \geq 1$$

$\xrightarrow{\text{log}}$

$$\text{const} + \sum_{i=0}^d x_i \log \left(\frac{q_i}{q'_i} \cdot \frac{1-q'_i}{1-q_i} \right) \geq 0$$

now let $w_i = \log\left(\frac{q_i}{q_i'} \cdot \frac{1-q_i'}{1-q_i}\right)$

\Rightarrow we have : $\text{const} + \sum_{i=0}^d x_i w_i \geq 0$

Therefore if ~~xx~~ holds, we will predict label 1 \Rightarrow The classifier is linear ✓

4. define : $\varepsilon_i(x) = y_i(x) - h(x)$

Then we would have :

$$\begin{aligned} E_n [(H_M(x) - h(x))^2] &= E_n \left[\left(\frac{1}{M} \sum_{i=1}^M y_i(x) - h(x) \right)^2 \right] \\ &= E_n \left[\left(\frac{1}{M} \sum_{i=1}^M \varepsilon_i(x) \right)^2 \right] \\ &= \frac{1}{M^2} E_n \left[\left(\sum_{i=1}^M \varepsilon_i(x) \right)^2 \right] \end{aligned}$$

by cauchy inequality we know that :

$$\begin{aligned} (1^2 + 1^2 + \dots + 1^2) (\varepsilon_1^2(x) + \dots + \varepsilon_M^2(x)) \\ \geq \left(\sum_{i=1}^M \varepsilon_i \right)^2 \end{aligned}$$


Therefore: $m \sum_{i=1}^m \epsilon_i^2(n) \geq \left(\sum_{i=1}^m \epsilon_i(n) \right)^2$

combining with \otimes , we have:

$$\begin{aligned} \frac{1}{m^2} E_n \left[\left(\sum_{i=1}^m \epsilon_i(n) \right)^2 \right] &\leq \frac{m}{m^2} E_n \left[\sum_{i=1}^m \epsilon_i^2(n) \right] \\ &= \frac{1}{m} \sum_{i=1}^m E_n \left[(y_i(n) - h(n))^2 \right] \end{aligned}$$

✓

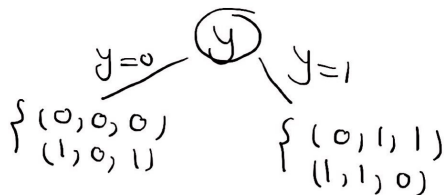
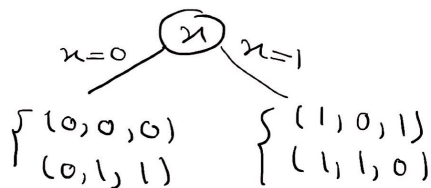
5. a) if we have m features, then we would at most 2^m different data points, and a decision tree of depth m , has at most 2^m leaves which covers all possible data. Therefore we can classify our dataset with zero training error.

b) No, consider this counterexample:
suppose we have $2 = m$ features, n and y which each can be 0 or 1. 

now suppose we have the following data:

x	y	label	
0	0	0	$(0, 0, 0)$
0	1	1	$\rightarrow (0, 1, 1)$
1	0	1	$(1, 0, 1)$
1	1	0	$(1, 1, 0)$

For a tree with depth $m-1=1$, we have these two cases:



we can see that in both choices for root node we have 50% accuracy.