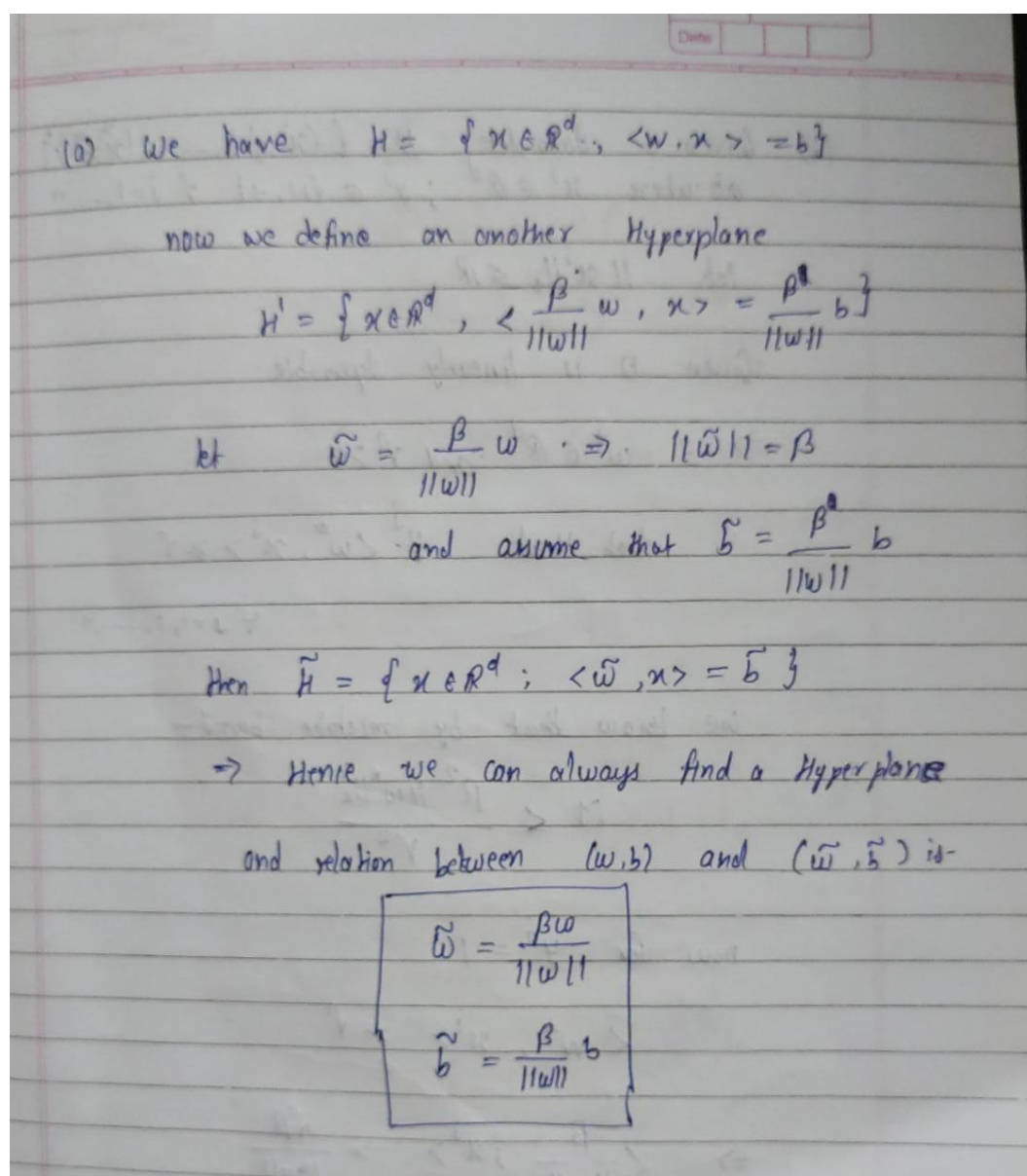


1. (a) [5 marks] Recall that a hyperplane  $H = (w, b)$  for some  $w \neq \mathbf{0} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$  is defined as  $H = \{x \in \mathbb{R}^d : \langle w, x \rangle = b\}$ . Show that for every hyperplane  $H$  and for every  $\beta > 0$ , there exists another hyperplane  $\tilde{H} = (\tilde{w}, \tilde{b})$  such that  $\|\tilde{w}\|_2 = \beta$ . Illustrate the relationship between  $(w, b)$  and  $(\tilde{w}, \tilde{b})$ . (Recall that for  $u \in \mathbb{R}^d$ ,  $\|u\|_2 = \sqrt{\langle u, u \rangle} = \sqrt{\sum_{i=1}^d |u_i|^2}$  is the  $\ell_2$  norm of  $u$ ).
- (b) [5 marks] Consider a data set  $D = \{(x^1, y^1), \dots, (x^n, y^n)\}$ , where  $x^j \in \mathbb{R}^d$ ,  $y^j \in \{+1, -1\}$ ,  $\forall j = 1, 2, \dots, n$ . Let  $\max_j \|x^j\|_2 \leq R$ . Recall that  $D$  is linearly separable if there exist  $w^* \in \mathbb{R}^d$  and  $\gamma > 0$  such that  $y^j \langle w^*, x^j \rangle \geq \gamma$ ,  $\forall j = 1, \dots, n$ . Show that if  $D$  is linearly separable, the mistake bound proved in class

$$M \leq \frac{R^2 \|w^*\|_2^2}{\gamma^2}$$

can be rewritten simply as  $M \leq \frac{R^2}{\eta^2}$ , where  $\eta > 0$  (which might be same as  $\gamma$  or different from  $\gamma$ ). (Hint: Use the previous result about hyperplane.)



1b) Given a data set  $D = \{(x^1, y^1), \dots, (x^n, y^n)\}$   
where  $x^j \in \mathbb{R}^d$ ;  $y^j \in \{+1, -1\}$   $\forall j=1, \dots, n$

$$\text{let } \|x^j\|_2 \leq R$$

Given  $D$  is linearly separable.

$$\Rightarrow \exists w^* \in \mathbb{R}^d \text{ and } \gamma > 0$$

$$\text{such that } y^j \langle w^*, x^j \rangle \geq \gamma$$

$$\forall j=1, 2, \dots, n$$

We know that by mistake bound  $\rightarrow$

$$M \leq \frac{R^2 \|w^*\|_2^2}{\gamma^2}$$

now for  $y^j = 1$

$$\langle w^*, x^j \rangle \geq \gamma$$

$$\Rightarrow \left\langle \frac{B}{\|w\|_2}, x^j \right\rangle \geq \frac{\sqrt{B}}{\|w\|_2}$$

$$(\text{since } B > 0 \\ \|w\|_2 > 0)$$

and for  $y^i = -1$

$$\langle w^*, x^i \rangle < r$$

$$\left\langle \frac{\beta}{\|w^*\|_2} w^*, x^i \right\rangle < \frac{r\beta}{\|w^*\|_2}$$

Hence  $(\tilde{w}, \tilde{b})$  is also linearly separable.

$$\text{where } \tilde{w} = \frac{\beta w^*}{\|w^*\|_2}$$

$$\tilde{b} = \frac{r\beta}{\|w^*\|_2}$$

$$\text{now } \|\tilde{w}\|_2^2 = \left\| \frac{\beta(w^*)}{\|w^*\|_2} \right\|_2^2 = \beta^2$$

by mistake bound inequality  $\rightarrow$

$$M < \frac{R^2 \|\tilde{w}\|_2^2}{r^2}$$

$$M < \frac{R^2 \beta^2}{r^2}$$

$$\text{let } \eta = \frac{r}{\beta}$$

$$\text{then } \boxed{M \leq \frac{R^2}{\eta^2}}$$

2. Consider the perceptron learning algorithm with a starting point  $w^0 = [\theta \ \theta \ \dots \ \theta]^T$  where  $\theta \in [0, 1]$ , used to train on a linearly separable data set.

- (a) [7 marks] Find a suitable upper bound on the number of mistakes for the choice of starting point  $w^0$  given above.
- (b) [3 marks] Compare and contrast the bound you obtained in part (a) with the bound discussed in class. Explain the changes observed in the bound, and explain the dependence of the bound you obtained in part (a) on the choice of  $w^0$ .
- (c) [2 marks] Justify if your bound is tight.

Q2 (a) Let  $w^0 = [\theta \ \theta \ \dots \ \theta]^T$  where  $\theta \in [0, 1]$ .

Consider the data set  $D = \{(x^1, y^1), \dots, (x^n, y^n)\}$ , where  $x^j \in \mathbb{R}^d$ ,  $y^j \in \{+1, -1\}$ ,  $\forall j = 1, 2, \dots, n$ . Let  $\max \|x^j\|_2 \leq R$ .  
As  $D$  is linearly separable  $\exists w^* \in \mathbb{R}^d$  and  $\gamma > 0$  s.t.  
 $y^j \langle w^*, x^j \rangle \geq \gamma, \forall j = 1, 2, \dots, n$ .

Using Perceptron update,  $w^{t+1} = w^t + y^t x^t$   

$$\begin{aligned} \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle &= \langle w^*, w^{t+1} - w^t \rangle \\ &= \langle w^*, y^t x^t \rangle \\ &= y^t \langle w^*, x^t \rangle > \gamma \quad (*) \end{aligned}$$

For a total progress of  $T$  rounds, where  $M$  mistakes are made, we obtain:

$$\begin{aligned} \sum_{t=0}^T \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle &= \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{mistake is} \\ \text{made in round } t}} \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle \\ &+ \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{mistake is} \\ \text{not made in} \\ \text{round } t}} \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle \\ &= \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{mistake is} \\ \text{made in round } t}} \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle \\ &\geq M \gamma \quad [\text{As total } M \text{ mistakes} \\ &\quad \text{have been made and} \end{aligned}$$

$$\therefore \sum_{t=0}^T \langle w^*, w^{t+1} \rangle - \langle w^*, w^t \rangle > M \gamma \quad (**) \quad (*)]$$



This can be justified as when no mistake is made

$$\langle \omega^*, \omega^{t+1} \rangle - \langle \omega^*, \omega^t \rangle = 0.$$

Now,

$$\begin{aligned} \sum_{t=0}^T \langle \omega^*, \omega^{t+1} \rangle - \langle \omega^*, \omega^t \rangle &= \langle \omega^*, \omega^{T+1} \rangle - \langle \omega^*, \omega^T \rangle \\ &+ \langle \omega^*, \omega^T \rangle - \langle \omega^*, \omega^{T-1} \rangle \\ &+ \dots \\ &+ \langle \omega^*, \omega^1 \rangle - \langle \omega^*, \omega^0 \rangle \\ &= \langle \omega^*, \omega^{T+1} \rangle - \langle \omega^*, \omega^0 \rangle \\ &= \langle \omega^*, \omega^{T+1} - \omega^0 \rangle \end{aligned}$$

Using (\*\*) we get

$$\begin{aligned} M\gamma &< \langle \omega^*, \omega^{T+1} - \omega^0 \rangle \\ &\leq \|\omega^*\| \|\omega^{T+1} - \omega^0\| \quad [\text{Using Cauchy-Schwarz Inequality}] \\ \therefore M\gamma &< \|\omega^*\| [\|\omega^{T+1}\| + \|\omega^0\|] \\ &\quad [\text{By Triangle Inequality}] \end{aligned}$$

$$\|\omega^{T+1}\| + \|\omega^0\| = \|\omega^{T+1}\| + \sqrt{d} \theta$$

$$M^2 \gamma^2 < \|\omega^*\|^2 [\|\omega^{T+1}\| + \sqrt{d} \theta]^2$$

$$= \|\omega^*\|^2 [\|\omega^{T+1}\|^2 + d\theta^2 + 2\sqrt{d}\theta \|\omega^{T+1}\|]$$

$$\therefore \frac{M^2 \gamma^2}{\|\omega^*\|^2} < [\|\omega^{T+1}\|^2 + 2\sqrt{d}\theta \|\omega^{T+1}\| + d\theta^2] - (***)$$

Using Perceptron update rule, we get

$$\omega^{t+1} = \omega^t + \gamma^t x^t$$

$$\begin{aligned}
\Rightarrow \|\omega^{t+1}\|^2 &= \|\omega^t + y^t x^t\|^2 \\
&= \|\omega^t\|^2 + \|y^t x^t\|^2 + 2\langle \omega^t, y^t x^t \rangle \\
&= \|\omega^t\|^2 + \|y^t\| \|x^t\|^2 + 2y^t \langle \omega^t, x^t \rangle \\
&\leq \|\omega^t\|^2 + \|x^t\|^2 \quad \text{[PA]}
\end{aligned}$$

As  $y^t \in \{1, -1\}$  and whenever update occurs (mistake has occurred)  $y^t \hat{y}^t \leq 0$  (as  $y^t \neq \hat{y}^t$ )

$$\Rightarrow y^t \text{sign}(\langle \omega^t, x^t \rangle) \leq 0$$

$$\Rightarrow y^t \langle \omega^t, x^t \rangle \leq 0.$$

$$\begin{aligned}
\therefore \|\omega^{t+1}\|^2 - \|\omega^t\|^2 &\leq \|x^t\|^2 \\
&\leq R \quad \left[ \|x^t\|^2 \leq R, \right. \\
&\quad \left. - (x \times x) \quad \forall t = 1, 2, \dots, n \right]
\end{aligned}$$

$$\begin{aligned}
\therefore \sum_{t=0}^T (\|\omega^{t+1}\|^2 - \|\omega^t\|^2) \\
&= \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{mistake} \\ \text{at round } t}} (\|\omega^{t+1}\|^2 - \|\omega^t\|^2) + \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{no mistake} \\ \text{at round } t}} (\|\omega^{t+1}\|^2 - \|\omega^t\|^2)
\end{aligned}$$

$$= \sum_{\substack{t \in \{0, \dots, T\} \\ t: \text{mistake} \\ \text{at round } t}} (\|\omega^{t+1}\|^2 - \|\omega^t\|^2)$$

$$\leq M R^2 \quad \text{[Using (***)]}$$

$$\therefore \|\omega^{T+1}\|^2 - \|\omega^0\|^2 \leq MR^2$$

$$\Rightarrow \|\omega^{T+1}\|^2 \leq MR^2 + d\theta^2$$

Using (\*\*\*)

$$\frac{M^2 R^2}{\|\omega^*\|^2} < [MR^2 + d\theta^2 + 2\sqrt{d\theta}\sqrt{MR^2 + d\theta^2} + d\theta^2]$$

$$= [MR^2 + 2d\theta^2 + 2\sqrt{d\theta}\sqrt{MR^2 + d\theta^2}]$$

considering  $M > 1$ ,

$$MR^2 + d\theta^2 \leq M^2 R^2 + M^2 d\theta^2$$

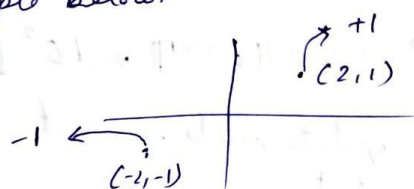
$$\text{and } 2d\theta^2 \leq 2Md\theta^2$$

$$\therefore \frac{M^2 \gamma^2}{\|w^*\|^2} < [MR^2 + 2Md\theta^2 + 2M\sqrt{d}\theta\sqrt{R^2+d\theta^2}]$$

$$\Rightarrow M < \frac{\|w^*\|^2}{\gamma^2} [R^2 + 2d\theta^2 + 2\theta\sqrt{d(R^2+d\theta^2)}]$$

(b) The bound obtained is greater than the one obtained in class for zero weight initialization. The bound increases with increase in  $d, R$  and  $\theta$ . For,  $\theta = 0$ , we obtain the same bound  $\frac{\|w^*\|^2}{\gamma^2} R^2$ . On the other hand when  $\theta = 1$ , the bound becomes  $\frac{\|w^*\|^2}{\gamma^2} [R^2 + 2d + 2\sqrt{d(R^2+d)}]$ .

(c) The bound is not tight as we can see from the example below:



The hyperplane with  $w^* = (1, 1)$  and  $\gamma = 1/2$  separates the data

For perceptron with  $w = (\frac{1}{2}, \frac{1}{2})$  no mistakes are made.

But the bound obtained is given by

$$\begin{aligned} & \frac{\|w^*\|^2}{\gamma^2} [R^2 + 2d\theta^2 + 2\theta\sqrt{d(R^2+d\theta^2)}] \\ &= \frac{2}{1/4} \left[ 5 + 2 \times 2 \times \frac{1}{4} + 2 \times \frac{1}{2} \sqrt{2(5 + 2 \times \frac{1}{4})} \right] \\ &= 74.533 \end{aligned}$$