

PERCEPTRON
Learning
Algorithm

① Perception: Algorithm used for binary classification of data.

Sou. Sarath
IP UD 20

\vec{w}_{lp} : is a vector $\in \mathbb{R}^d$

y_{lp} : is a label $y \in \{-1, +1\}$

$$y_{lp} = \text{sgn}(\vec{w}_{lp}^T \vec{x}_n + b) = \text{sgn}(b + \sum w_i x_i)$$

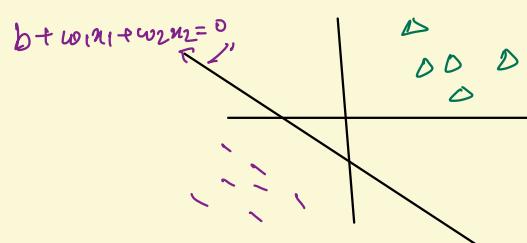
$\vec{w}_{lp} \propto \vec{x}_n$ if $\text{sgn}(\vec{w}_{lp}^T \vec{x}_n + b) \geq 0 \rightarrow \text{predict } y=1$
 $\text{sgn}(\vec{w}_{lp}^T \vec{x}_n + b) < 0 \rightarrow \text{predict } y=-1$

$b = b$ bias term
 $w_i = \text{weights}$

try to find a plane $\vec{w}_{lp}^T \vec{x}_n + b$ s.t \uparrow

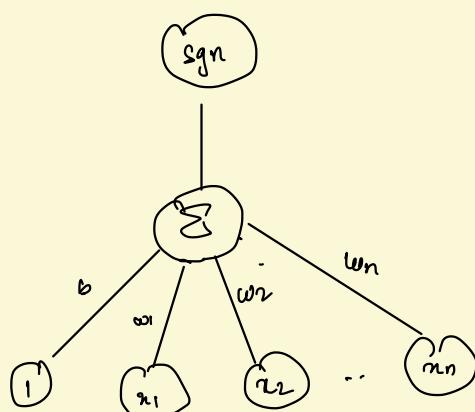
- gives a linear classifier that represents a hyperplane which separates the space into 2 half-spaces.

Goal: to find a hyperplane



* Parameters: $\vec{w} = \begin{pmatrix} b \\ w_1 \\ w_2 \end{pmatrix}$, $\eta = \text{learning rate}$ \rightarrow higher the learning rate faster is the convergence to the hyperplane.
 coefficient $\vec{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix}$

$$\vec{w}^T \vec{x} = (b, w_1, w_2) \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = \vec{w}^T \vec{x}.$$



We give training data and try to create a hyperplane using it through the algorithm so that for the new data, x , it can classify it using the hyperplane.

Check out [Youtube: Ritvik](#). Watch

The Perceptron Algorithm

- Initialize $\vec{w}_0 = \vec{0} \in \mathbb{R}^n$

- For each training example (\vec{x}_i, y_i)

$$y' = \text{sgn}(\vec{w}_0^T \vec{x}_i)$$

If $y \neq y'$

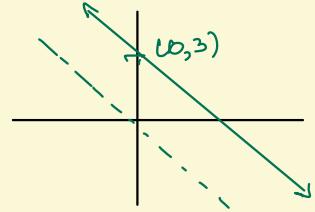
$$\text{Update } w_{t+1} \leftarrow w_t + \tau(y_i x_i)$$

. Return final Vector

Some Prerequisites:

* Hyperplane: The eqn is $w \cdot x + b$, where w is the normal to the hyperplane and b is an offset.

Example: $x + y = 3$: $y = -x + 3$

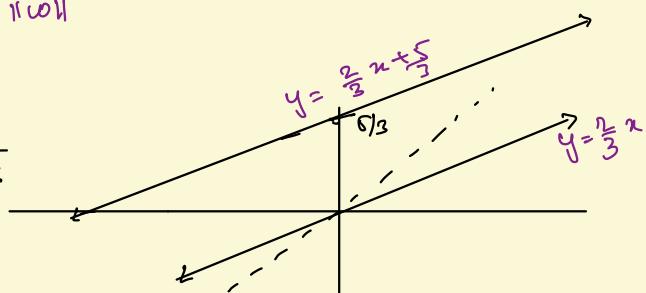


If we multiply $w \cdot x + b$ by $\frac{1}{\|w\|}$, then we have,

$$\hat{w} = \frac{w}{\|w\|}, \quad \hat{w} \cdot x + b' = 0. \quad \hat{w} = \text{unit normal to the hyperplane.}$$

$$b' = \frac{b}{\|w\|} = \text{distance of hyperplane from origin}$$

(ii) $-2x + 3y = 5$: $y = \frac{2}{3}x + \frac{5}{3}$



Given a hyperplane, $w \cdot x + b$

In 2D, $w_1 x_1 + w_2 x_2 + w_3 x_3 + b$

Distance of point y from $w \cdot x + b$ is: $\frac{|w \cdot y + b|}{\|w\|}$ i.e.

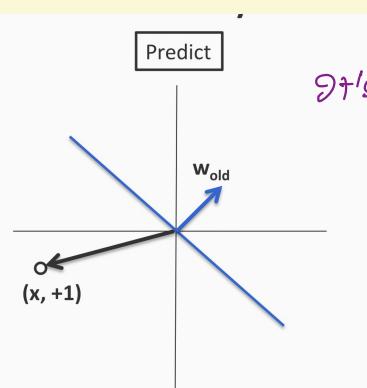
$$\frac{w_1 y_1 + w_2 y_2 + w_3 y_3 + b}{\sqrt{w_1^2 + w_2^2 + w_3^2}}$$

In (ii) dist. $(0, 0) = \frac{-2(0) + 3(0) + 5}{\sqrt{2^2 + 3^2}} = \frac{5}{\sqrt{13}}$

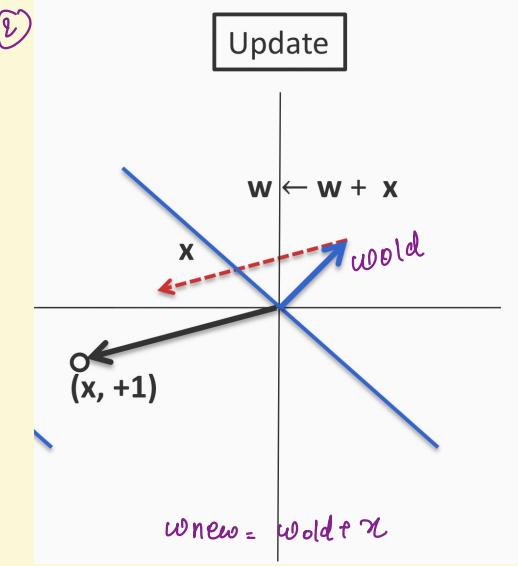
Idea Behind Perceptron Update:

Look at slide 19.

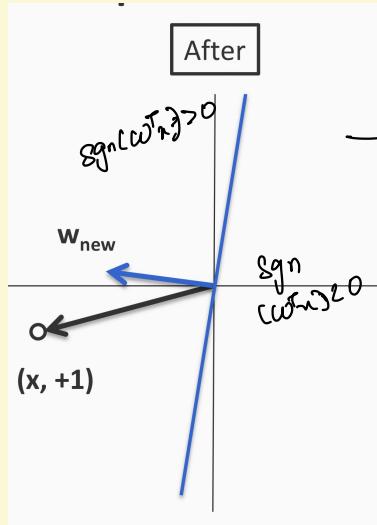
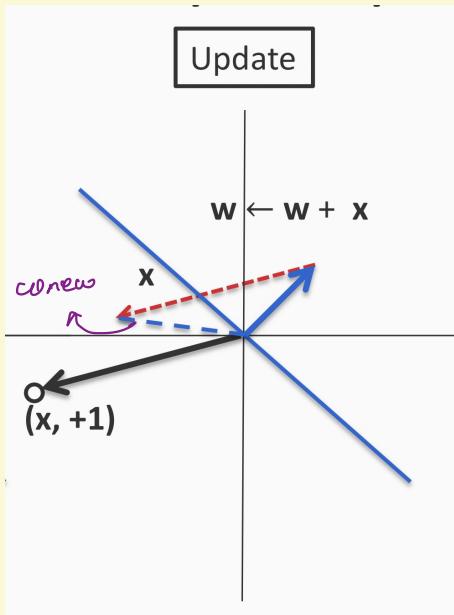
①



gt's +1 but predicted -1



\therefore find a hyperplane $w^{(0)}$ that w_{new} is a normal to that hyperplane



→ After the update - new hyperplane whose normal is w_{new}

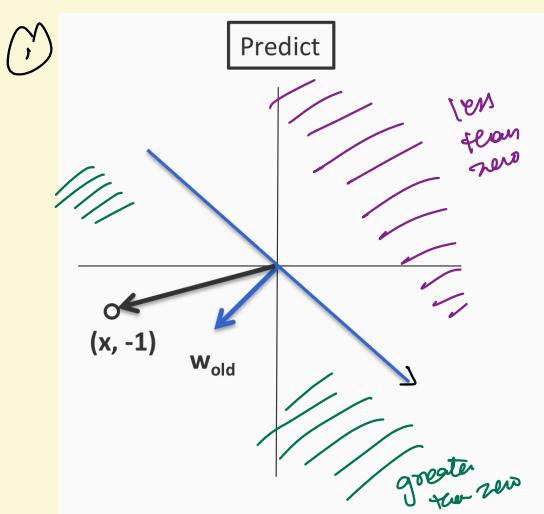
IMP

In this eg, since the label of x is $+1$, we want it to be ∞ the right of hyperplane. But in the begining it was on the left of the hyperplane (w_{old}). AIM: is to find a hyperplane

so that x lies right to the hyperplane.

After update we found a new hyperplane, s.t. the point x is to the right of the hyperplane.

Another Case Mistake on a negative eg.

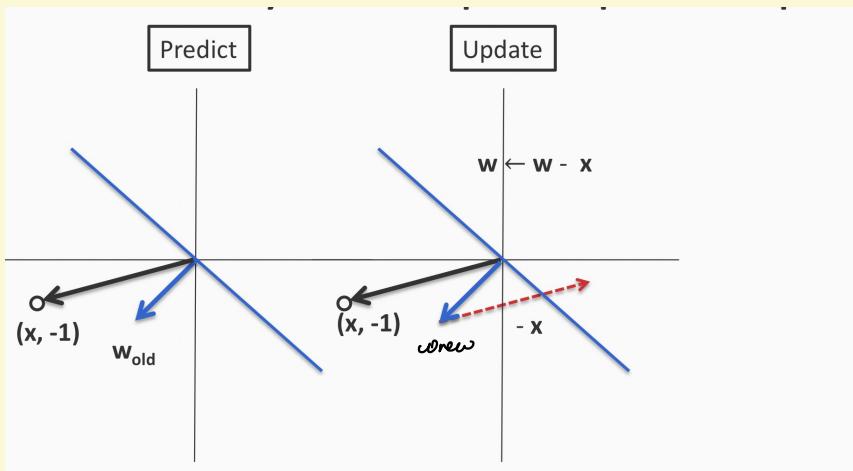


For this hyperplane, the normal is w_{old} . \therefore To the left of hyperplane $\text{sgn}(w^T x) > 0$ & right $\text{sgn}(w^T x) < 0$

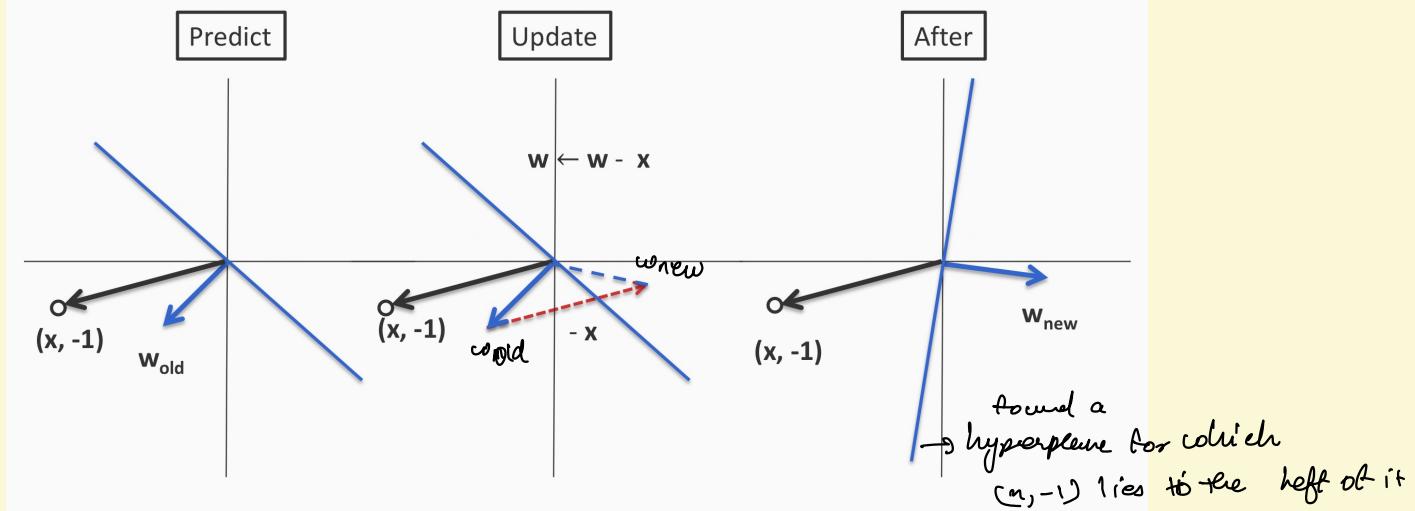
Given: $(x, -1)$ But this lies in the wrong side of the hyperplane.

Aim: To rotate the hyperplane so that $(x, -1)$ lies on the left of the hyperplane. \therefore the label is -1 .

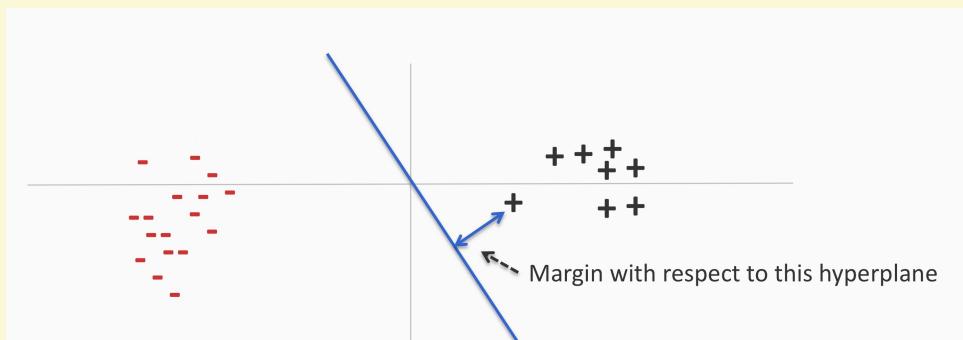
(2)



(3)

Convergence:

- ① If the data is linearly separable, perceptron algorithm will converge.
- ② **Cycling them:**
If the training data isn't separable, then learning algo. will eventually repeat the same set of weights & enter an ∞ loop.
- ③ **Margin them:**
margin of a hyperplane is the dist. b/w the hyperplane & data point nearest to it.
margin of a dataset (δ): is the max. margin for that dataset using any weight vector.

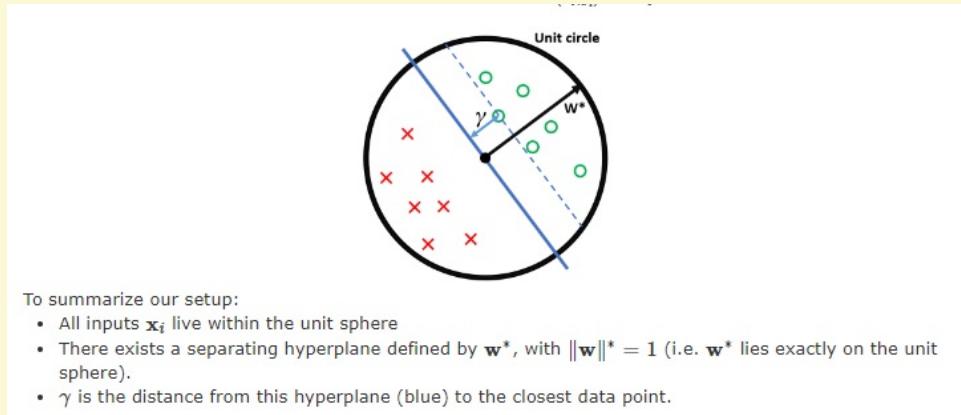


Convergence:

stmt: If a data set is nearly separable, then perceptron will find a separating hyperplane. (If a data set is linearly separable, it'll loop forever.)

Argument: Suppose $\exists w^* \text{ s.t. } y_i(w^T x_i) \geq 0 \quad \forall (x_i, y_i) \in D$
 $\text{s.t. } \|w^*\| = 1 \quad \text{and} \quad \|x_i\| \leq 1. \quad \forall x_i \in D.$

$$\text{margin: } \gamma = \min_{(x_i, y_i) \in D} |x_i^T w^*|$$



// These assumptions essentially mean that data is linearly separable.

\therefore If all the above holds, then the perceptron will make at most $(1/\gamma^2)$ mistakes.

Suppose we have a binary classification dataset with n dimensional inputs.

If the data is separable, ...

...then the Perceptron algorithm will find a separating hyperplane after making a finite number of mistakes

Proof

After t mistakes, we make 2 claims.

$$\text{Claim 1: } w^* \cdot w_t \geq t\gamma$$

$$\text{Claim 2: } \|w_t\|^2 \leq t$$

Claim 1 Proof:

$$\begin{aligned} w^* \cdot w_{t+1} &= w^* \cdot (w_t + y_t x_t) \\ &= w^* \cdot w_t + y_t (w^* \cdot x_t) \stackrel{\downarrow}{\geq} 0 \geq w^* \cdot w_t + \gamma \end{aligned}$$

$$\therefore w^* \cdot w_{t+1} \geq w^* \cdot w_t + \gamma$$

$$w_0 = 0 \Rightarrow w^* \cdot w_1 \geq \gamma, \quad w^* \cdot w_2 \geq 2\gamma, \dots \quad w^* \cdot w_t \geq t\gamma$$

Hence $w^*^\top w_t \geq t\gamma$

Claim 2 Proof:

$$\|w_{t+1}\|^2 = \|w_t + y_i x_i\|^2 \quad (\text{In case of update})$$

$$= \|w_t\|^2 + 2y_i (w_t^\top x_i) + \|x_i\|^2$$

Observe $2y_i (w_t^\top x_i) < 0$ (\because it's the reason update was made!)

and $\|x_i\|^2 \leq 1$

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + 1$$

$$w_0 = 0 \Rightarrow \|w_1\|^2 \leq 1 \quad \|w_2\|^2 \leq 2 \quad \dots \|w_t\|^2 \leq t$$

$$\therefore t\gamma \leq w_t^\top w^* = \|w_t\| \|w^*\| \cos(\theta)$$

(Definition of γ)
 $\theta = \text{angle b/w } w^* \text{ & } w_t$

$$\leq \|w_t\|$$

$$\leq \sqrt{t}$$

$$\Rightarrow t\gamma \leq \sqrt{t} \Rightarrow \sqrt{t} \leq \frac{1}{\gamma}$$

$$\Rightarrow t \leq \frac{1}{\gamma^2}$$

\therefore # of mistakes t is bounded above by a constant.

The Algorithm:

Perceptron Algorithm

```
1 Initialize  $\vec{w} = \vec{0}$                                 // Initialize  $\vec{w}$ .  $\vec{w} = \vec{0}$  misclassifies everything.  
2 while TRUE do                                     // Keep looping  
3    $m = 0$                                          // Count the number of misclassifications,  $m$   
4   for  $(x_i, y_i) \in D$  do                         // Loop over each (data, label) pair in the dataset,  $D$   
5     if  $y_i(\vec{w}^T \cdot \vec{x}_i) \leq 0$  then          // If the pair  $(\vec{x}_i, y_i)$  is misclassified  
6        $\vec{w} \leftarrow \vec{w} + y_i \vec{x}$                   // Update the weight vector  $\vec{w}$   
7        $m \leftarrow m + 1$                             // Counter the number of misclassification  
8     end if  
9   end for  
10  if  $m = 0$  then                               // If the most recent  $\vec{w}$  gave 0 misclassifications  
11    break                                         // Break out of the while-loop  
12  end if  
13 end while                                       // Otherwise, keep looping!
```

Line 5 checks if predicted sign & actual sign are equal.
Line 10 checks if there is no update in that cycle, it means the hyperplane has been found.

Margin Perceptron

6. Given a training set D with a separation margin γ_0 , the original perceptron algorithm predicts a mistake when $y\mathbf{w}^{(n)T}\mathbf{x} < 0$. As we have discussed that the perceptron algorithm converges to a linear classifier that can perfectly separate D but does not necessarily achieve the maximum margin. The **margin perceptron algorithm** extends presented perceptron Algorithm to approximately maximize the margin in the perceptron algorithm, where it is considered to be a mistake when $\frac{y\mathbf{w}^{(n)T}\mathbf{x}}{\|\mathbf{w}^{(n)}\|} < \frac{\gamma}{2}$, where $\gamma > 0$ is a parameter. Prove that the number mistakes made by the margin perceptron algorithm is at most $\frac{8}{\gamma_0^2}$ if $\gamma \leq \gamma_0$.

We normalize all the examples (x_i^o, y_i^o) so that $\|x_i^o\| \leq 1$ #^o

Why Many in Perceptron Algorithm? We iterate through the data using perceptron algorithm updating the weights not only on mistakes but also on examples x that the current hypothesis gets correct by a margin $< \frac{\gamma}{2}$ [So if an ex is with w margin $\frac{\gamma}{2}$ then it's incorrect]

Predict rule $\frac{w_t^T x_i^o}{\|w_t\|} \geq \frac{\gamma}{2}$, Negative if $\frac{w_t^T x_i^o}{\|w_t\|} < \frac{-\gamma}{2}$

$\frac{y_i^o w^{(n)}{}^T x_i^o}{\|w^{(n)}\|} < \frac{\gamma}{2}$ then it's a mistake

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + \|x_i^o\|^2 + 2y_i^o w_t^T x_i^o$$

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + \|x_i^o\|^2 + 2y_i^o w_t^T x_i^o$$

$$\leq \|w_t\|^2 + 2w_t^T x_i^o + \|x_i^o\|^2$$

$$\leq \|w_t\|^2 \left(1 + \frac{1}{\|w_t\|} \frac{w_t^T x_i^o}{\|w_t\|} + \frac{\|x_i^o\|^2}{\|w_t\|^2} \right)$$

$$\|w_{t+1}\| \leq \|w_t\| \left(1 + \underbrace{\frac{2}{\|w_t\|} \frac{\gamma}{2} + \frac{1}{\|w_t\|^2}}_{\geq 0} \right)^{1/2}$$

(For $x > 0 \quad \sqrt{1+x} < 1 + \frac{x}{2}$) Verify by MVT

$$\|w_{t+1}\| \leq \|w_t\| \left(1 + \frac{\gamma}{2\|w_t\|} + \frac{1}{2\|w_t\|^2} \right)$$

$$\text{Hence } \|w_{t+1}\| \leq \|w_t\| + \frac{1}{2\|w_t\|} + \frac{\gamma}{2}$$

If the classification is $\frac{y_i \mathbf{w}_t^\top \mathbf{x}_i}{\|\mathbf{w}_t\|} < \frac{\gamma}{2}$ it's a mistake

$$\frac{\|\mathbf{w}_t\|}{y_i \mathbf{w}_t^\top \mathbf{x}_i} > \frac{2}{\gamma}$$

If for some eg we take $y_i \mathbf{w}_t^\top \mathbf{x}_i = 1$ then, $\|\mathbf{w}_t\| > \frac{2}{\gamma}$
 $(\frac{1}{\|\mathbf{w}_t\|} < \frac{\gamma}{2})$

$$\|\mathbf{w}_{t+1}\| \leq \|\mathbf{w}_t\| + \frac{\gamma}{4} + \frac{\gamma}{2} = \|\mathbf{w}_t\| + \frac{3\gamma}{4}$$

Repeating it recursively after m mistakes

$$\|\mathbf{w}_{m+1}\| \leq \|\mathbf{w}_t\| + \frac{3M\gamma}{4} < \frac{2}{\gamma} + \frac{8M\gamma}{4}$$

Also for a perceptron algo recall $M\gamma \leq \|\mathbf{w}_{m+1}\|$

$$\text{So, } M\gamma \leq \|\mathbf{w}_{m+1}\| < \frac{2}{\gamma} + \frac{8M\gamma}{4}$$

$$M\gamma \leq \frac{2}{\gamma} + \frac{8M\gamma}{4} \Rightarrow \frac{M\gamma}{4} \leq \frac{2}{\gamma} \Rightarrow \boxed{M \leq \frac{8}{\gamma^2}}$$