

## Understanding Perceptron Algorithm - Handwritten Notes

### Table of Contents

1. What is Perceptron Algorithm
2. What is Threshold Logic Unit
3. How the Perceptron is Trained
4. Understanding the Update Rule
5. Assumptions for convergence
6. Proof of Convergence

→ The neural networks are further developed in various architectures, and now these

Artificial neural networks (ANNs) are the very core of Deep Learning.

They are very powerful, making them ideal to tackle large and highly complex machine learning tasks such as classifying billions of images (e.g. Google Images), translating texts, recommending the best videos to watch, etc.

### simplest Artificial neural network Architecture

#### (The Perceptron)

→ The Perceptron is one of the simplest ANN Architecture, as it is based on an artificial neuron called - "threshold logic Unit (TLU)".

→ The inputs and outputs are now numbers (instead of binary on/off values) and each input connection is associated with a weight.

→ Before seeing this architecture, it would be much clear if we can understand the Perceptron Arch.

## perceptron algorithm

Input:  $\{(x_1, y_1), \dots, (x_n, y_n)\}$   $x_i \in \mathbb{R}^d$

Iteration:  $w^0 = 0 \in \mathbb{R}^d [0, \dots, 0]$

until convergence:

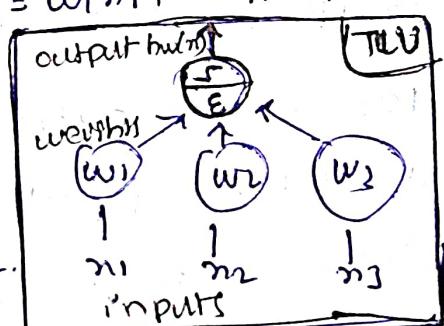
- pick  $(x_i, y_i)$  pair from the dataset
- If  $\text{sign}(w_t^T x_i) = y_i$  (do nothing)
- Else,  $w^{t+1} = w^t + x_i y_i$  update rule.
- End.

Assumption: the dataset is linearly separable.

→ The TLU computes a weighted sum of its inputs, then applies a step function to that sum and outputs the result.

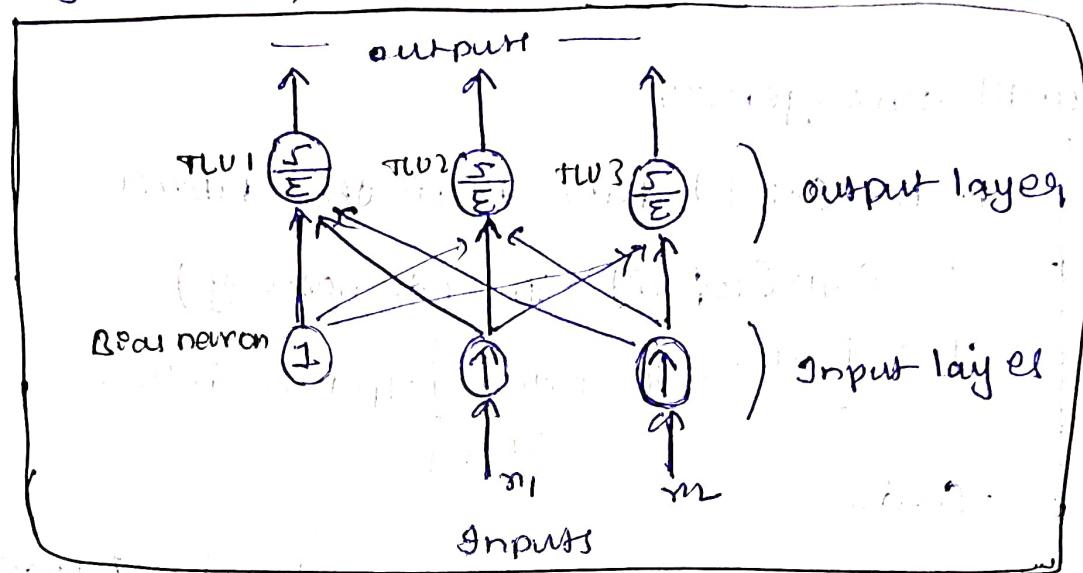
$$h_w(x) = \text{step}(z), \text{ where } z = x^T w = w_0 x_0 + w_1 x_1 + \dots + w_n x_n,$$

→ Training a TLU in this case means that finding the right values for  $w_0, w_1, \dots, w_n$ .



→ Here, a perceptron is simply composed of a single layer TLUs with each TLU connected to all the inputs. An extra bias feature is generally added ( $n_{0,1}$ ), which just outputs 1 all the time, called a bias neuron.

→ Here is a perceptron with two inputs and 3 outputs, this perceptron can classify instances simultaneously into 3 different binary classes, making it multioutput classifier.



→ So how this perceptron is trained?

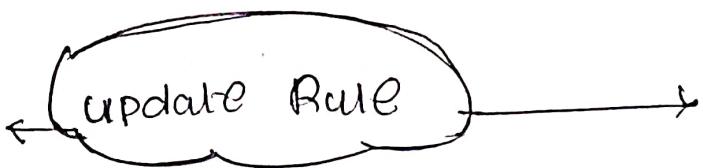
Just like update rule, we have seen before,

→ Perceptron is fed one training instance at a time, and for each instance it makes a set of predictions.

For every output neuron that produced a wrong prediction, it reinforces the connection weights from the inputs that would have contributed to correct prediction.

→ And if the training instances are linearly separable, it is said that perceptron would converge to a solution.

How?



$$w^{t+1} = w^t + \eta_i y_i$$

→ if mistake  $\left[ \text{sign}(w^T x_i) \neq y_i \right]$

\* Now, we need to understand how by convergence we will get a bet for  $w^t$ !

\* If we think about it, there are two different ways this update rule can go wrong.

mistake ①

$$\text{predicted} = 1 \quad \text{if } (w^T x_i) \geq 0$$

Actual = -1

$$y_i = -1$$

mistake ②

$$\text{predicted} = -1 \quad \text{if } \text{sign}(w^T x_i) < 0$$

Actual = +1

$$y_i = +1$$

→ say, we are using the update rule as

→ so, if we take update rule as a

$$w^{t+1} = w^t + \eta_i y_i$$

good one, then it should not make mistake for the input for which previous weight had made a mistake.

→ let's find the new weight dot product with the point where old weight made a mistake.

$$\rightarrow (\omega^{t+1})^T n_i = (\omega^t + \alpha_i y_i)^T n_i = \omega^T n_i + \alpha_i \|n_i\|^2$$

+ now, considering type I mistake,

we have,  $\begin{cases} \omega^T n_i \geq 0 \\ y_i \approx -1 \end{cases}$

Substitute them in new weight,

And we get,

$$\begin{aligned} (\omega^{t+1})^T n_i &= \underbrace{\omega^T n_i}_{\geq 0} + \underbrace{\alpha_i \|n_i\|^2}_{\geq 0} \\ &\quad \underbrace{-1}_{\text{This whole thing is negative.}} \end{aligned}$$

+ From this, what we can say is that the new weight's dot product with  $n_i$  is equal to, old weight's dot product with  $n_i$  - Something.

+ So, we are subtracting out something from old  $\omega$ 's dot product.

+ well, the old  $\omega$ 's was positive, and by subtracting something, we can say that the value is reducing.

+ Of course, it doesn't immediately become negative, but however it is moving in the right direction as we are subtracting!!

→ similarly if we consider type 2,

where,

$$(w^T_{\cdot i}) < 0$$

$$y_i = +1$$

substitute them in new weight;

And we get,

$$(w^{t+1})^T_{\cdot i} = \underbrace{w^T_{\cdot i}}_{< 0} + \underbrace{y_i \|x_i\|^2}_{> 0}$$

→ and here we are adding something.

→ overall we can say that,

update rule pushes  $w$  in the right direction

for  $x_i$ .

if dot product is -ve, when it should have been +ve,  
the update rule is subtracting something.

→ But then, the question raised that,

okay we fixed the previous mistake,

but what's the guarantee that this

new weight didn't break anything that

was previously correct!!

→ So problem is - refining  $w$  for one  $x$ ,

might affect decision for

other data points. so we

need more careful argument

for convergence.

④ Let's solve an example~

$$\begin{array}{c} \text{3} \left( \begin{bmatrix} -1 \\ 1/2 \end{bmatrix}, -1 \right) \\ \text{1} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}, +1 \right) \\ \hline \text{2} \left( \begin{bmatrix} 0 \\ -1 \end{bmatrix}, +1 \right) \end{array}$$

Is this a linearly separable dataset?

Yes, there is  $w \in \mathbb{R}^2$

$$s.t. w^T n_1 \geq 0 \Rightarrow y_1 = +1$$

$$w^T n_2 \leq 0 \Rightarrow y_2 = -1$$

Solving with perceptron as it's linearly separable.

Initially we can take the weight  $[0, 0]$

And that gives us  $w^0 = [0, 0]$

$$\begin{array}{c|c|c} w^0 T n_1 = 0 & w^0 T n_2 = 0 & w^0 T n_3 = 0 \\ \hat{y}_1 = +1 & \hat{y}_2 = +1 & \hat{y}_3 = +1 \end{array} \rightarrow \text{This is a mistake.}$$

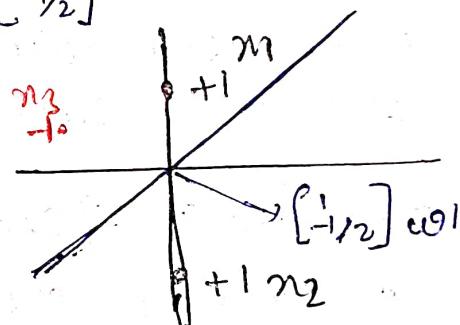
$$y_3 = -1,$$

so we find  $w^1$

$$w^1 = w^0 + n_3 y_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1/2 \end{bmatrix} \cdot -1 = \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}$$

Now with the new  $w^1$ , we can see that it made mistake for  $n_1$ .

So we find  $w^2$



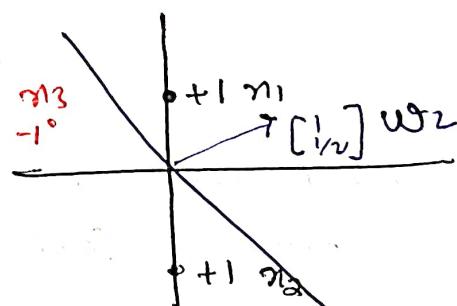
$$w^2 = w^1 + n_1 y_1 = \begin{bmatrix} 1 \\ -1/2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \cdot 1 = \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}$$

Now with the new  $w^2$ , we can see that it made mistake for  $n_2$ .

So we find  $w^3$ ,

$$w^3 = w^2 + n_2 y_2 = \begin{bmatrix} 1 \\ 1/2 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 1/2 \end{bmatrix}$$

And now  $w^3$  has become  $w_1$ , and then  $w^4$  will become  $w_2$  & it will keep looping like this & it will not have a convergence!!



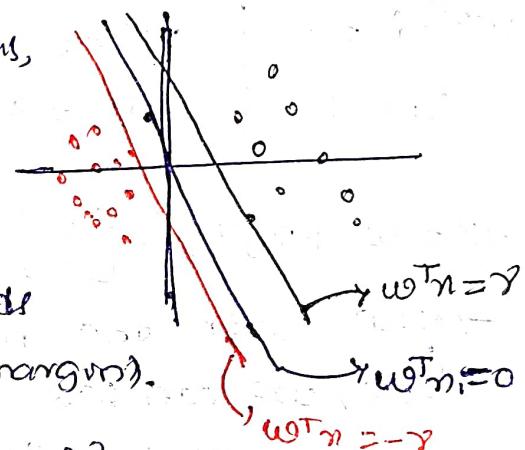
→ Now if we wonder is perceptron wrong or data wrong. And if we closely observe the data, the points actually lie on the decision boundary, i.e. are say  $y_i$   
 $w^T x_i + b$  → it can also be equal to 0.

i.e.,  $w^T x_i + b$  can be 0. so we need to make sure it can't be zero.

So the data assumption has to be changed.

## ① Linear separability with "γ" margin.

→ If we take the data like this, where there is a region with no datapoints. Then data can be modelled well with perceptron. If needed areas a small region ( $\gamma$ -margin).

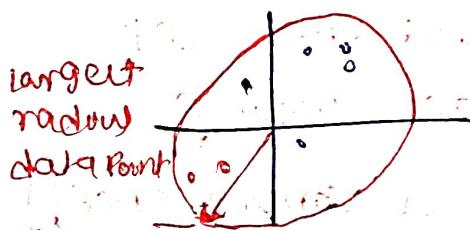


→ A dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  is linearly separable with  $\gamma$ -margin.

→ If  $\exists w^T \in \mathbb{R}^d$  s.t.,  $(w^T x_i) y_i \geq \gamma \forall i$   
 for some  $\gamma > 0$ .

→ For ease of simplification or proof, we can consider few harmless assumptions such as

## ② Radial Assumption



$\forall i (\text{every point}) \in D,$

$\|x_i\|_2 \leq R \text{ for some } R > 0$

→ For a largest point, we are away all other points. We draw that circle.

g) without loss of generality, assume  $\|w^*\| = 1$ ,

$w^*$  → The actual best line that separates the data, we know that there will be  $w^*$ , but we can't be sure all norm ( $\|w^*\|$ ) will be equal to 1.

however, we can say that there will exist a  $w^*$ , which will have norm = 1.

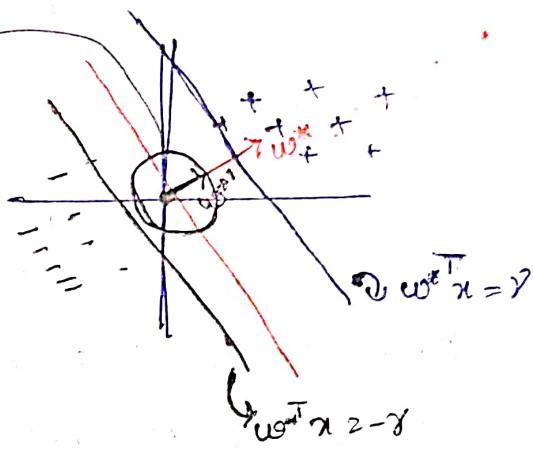
we have,

$$w^{*T} n = \gamma$$

considering  
 $\|w^*\| = 1$

divide by  $\|w^*\|$

$$\frac{w^{*T} n}{\|w^*\|} = \frac{\gamma}{\|w^*\|}$$



$$\left( \frac{w^{*T}}{\|w^*\|} \right) \cdot n = \gamma$$

$$\therefore \gamma = \frac{\gamma}{\|w^*\|}$$

→ considered  $w^*$

$$(w^*)^T n = \gamma'$$

↳ Here we just scaled the  $w^*$ .

But the advantage is that, the new  $w^*$  we have now has norm = 1.

↳ If you give me a dataset, with  $w^*$  some  $\gamma$ , say  $\gamma=10$ .

Now I can create another  $w^*$  by rescaling the  $w^*$ , to make sure it has length 1, and the  $\gamma$  will rescale accordingly. Then I will have a  $w^{*1}$  which also linearly separates data with a different  $\gamma'$ , so this assumption holds ✓

## ① Proof of Perceptron

→ The way we are gonna prove that this algorithm converges is by following idea,

48 we say that perceptron algorithm cannot make more than a certain mistakes.

### Idea of convergence.

→ If we think about it, it makes sense that if the mistakes are gonna be infinite, then it will never converge. So, if we can prove that the mistakes are gonna be finite, then we can say it converged.

### Analysis of mistakes of Perceptron

→ We know that, an update in the perceptron happens only when a mistake happens.

→ Let's say  $w^l$  is the current guess, and then a mistake happens w.r.t  $(x_i, y)$ .

Then,  $w^{l+1} = w^l + \eta \cdot y$

Let's understand the length of this update.

$$\begin{aligned}
 \|w^{l+1}\|^2 &= \|w^l + \eta \cdot y\|^2 = (w^l + \eta y)^T (w^l + \eta y) \\
 &= \|w^l\|^2 + 2 \cdot (w^l \cdot \eta y) + \|\eta y\|^2 \\
 &\quad \text{So, } \underbrace{(w^l \cdot \eta y)}_{\text{dot product}} \text{ is assumed to be } \leq 0 \text{ (because mistake)} \\
 &\leq \|w^l\|^2 + \|\eta y\|^2 \leq R^2
 \end{aligned}$$

→ If we observe, why we said and term  $(2 \cdot (\omega^T_n) y) \leq 0$ ,  
 it is because that,  $\omega^T_n$  is the prediction ( $-1/+1$ ),  
 and  $y$  is actual ( $-1/+1$ ).

And considering update happens only on mistake,  
 then both have to be opposite value in case,  
 and the total value will turn out to be negative.

Hence  $2(\omega^T_n) y \leq 0$ , ————— ignored

so we can upper bound whole thing as,

$$\|\omega^{t+1}\|^2 \leq \|\omega^t\|^2 + R^2$$

→ This means that the new weight length  
 will be utmost previous length +  $R^2$ ,

→ same applies for previous length.

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

..... If we continue we get ((for 1 mistake))

$$\|\omega^{t+1}\|^2 \leq \underbrace{\|\omega^0\|^2}_{=0} + tR^2$$

∴  $\boxed{\|\omega^{t+1}\|^2 \leq tR^2}$  ————— (1)

→ No or mistakes

→ So we say that the max. length it  
 can grow up 't' times  $R^2$ , and  $t$  is the  
 no. of updates/mistakes.

→ From the other side, to understand about  $(w^{l+1})$  quantity, let's use  $w^*$ .  
we have,

$$w^{l+1} = w^l + n \cdot y$$

$$(w^{l+1})^T \cdot w^* = (w^l + n \cdot y)^T \cdot w^*$$

$$= w^{lT} w^* + \underbrace{(w^{*T} n) \cdot y}$$

$\geq \gamma$  (based on  
 $w^*$  assumption)

$$(w^{l+1})^T \cdot w^* \geq w^{lT} w^* + \gamma$$

→ we can take same argument for  $w^l$ ,

$$(w^{l+1})^T \cdot w^* \geq (w^{l-1T} w^* + \gamma) + \gamma$$

..... etc. 1 mistake.

$$(w^{l+1})^T \cdot w^* \geq w^{0T} w^* + l\gamma$$

$\downarrow$   
0 (as  $w^0 = 0$ )

$$\boxed{(w^{l+1})^T \cdot w^* \geq l\gamma}$$

②

+ It says that if we make mistake,  
the dot product with  $w^*$  is increasing.

Now, if we try to make use of these two eqn,  
first we can recall the Cauchy-Schwartz inequality,

for any  $x^T y$

$$\left\| \left( \frac{x^T y}{\|y\|^2} \right) y \right\|^2 \leq \|x\|^2$$

$$(x^T y)^2 \cdot \frac{\|y\|^2}{\|y\|^4} \leq \|x\|^2$$

$$\Rightarrow (x^T y)^2 \leq \|x\|^2 \cdot \|y\|^2$$

From eq(2),

$$\lambda y \leq (w^{l+1})^T \cdot w^*, \text{ apply square.}$$

$$\lambda^2 y^2 \leq [(w^{l+1})^T \cdot w^*]^2 \leq \|w^{l+1}\|^2 \cdot \|w^*\|^2$$

From eqs  
1

$$\lambda^2 y^2 \leq \|w^{l+1}\|^2$$

$$\therefore \|w^{l+1}\|^2 \geq \lambda^2 y^2$$

From eq(3)

$$\lambda^2 y^2 \leq \|w^{l+1}\|^2 \leq \lambda R^2$$

↑  
From eq(3)

↑  
From eq(1)

$$\lambda^2 y^2 \leq \lambda R^2$$

Radius Margin Bound.

$$\lambda \leq R^2/y^2$$

$\lambda$  — no. of mistakes.

$$\rightarrow \text{So, by observing } \boxed{\lambda \leq R^2/\gamma^2},$$

this tells us that  $\lambda$  has to be at most  $R^2/\gamma^2$ ,

and at the first assumption, we said  $\gamma > 0$ .  
And so we clearly say  $\lambda$  is gonna be finite.

$\star$  And  $\lambda$  is the number of mistakes, so we can finally say that the no. of mistakes is a finite quantity  $\lambda$ .  
Hence the "Perceptron will be converged."

$\rightarrow$  If we just take linearly separable assumption without considering  $\gamma$ , i.e.  $\gamma = 0$ .

then  $\lambda \leq R^2/\gamma^2 \leq \infty$  (infinity).

In this case it is not certain. hence it can't be converged.



### Perceptron

$$\text{no. of mistakes}(\lambda) \leq \frac{R^2}{\gamma^2}$$

- under assumption,

$$\rightarrow \|w_i\|^2 \leq R^2$$

$\rightarrow$  Dataset - linearly separated with margin  $\gamma$ .

$$(w^T w_i) y_i \geq \gamma \quad \forall i, \gamma > 0$$