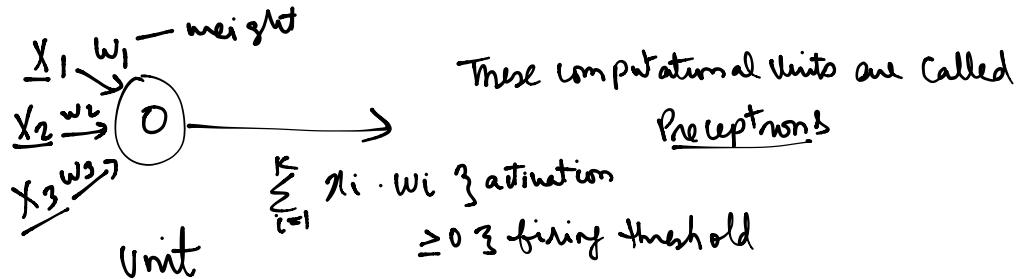


## \* Neural Networks

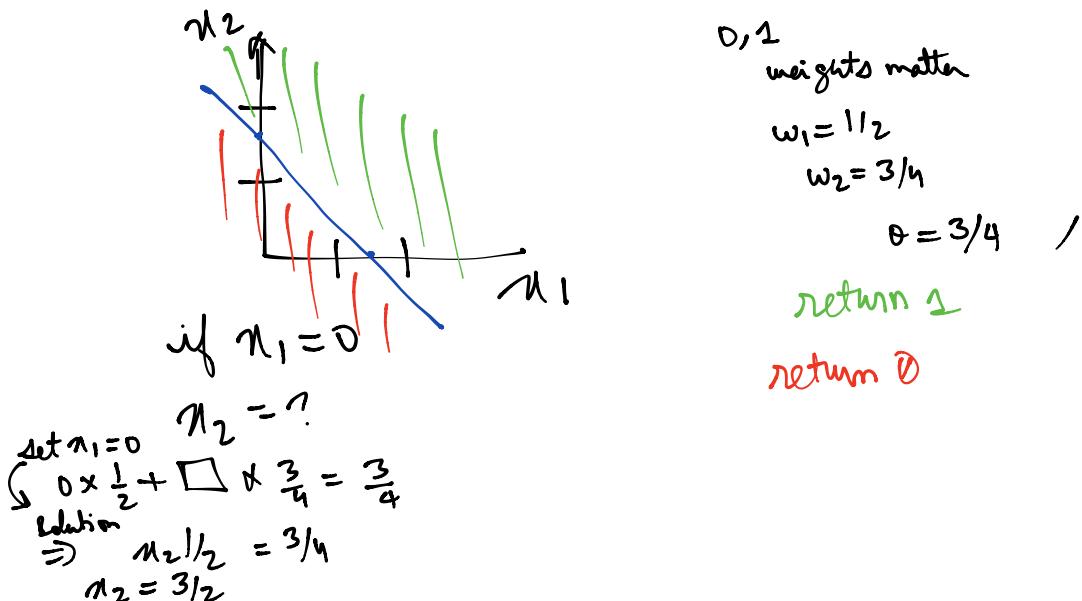
\* They are computational units.

## \* Artificial Neural Networks



Qn 3:

$$\begin{array}{ll}
 \begin{matrix} x_1 & 1 \\ x_2 & 0 \\ x_3 & -1.5 \end{matrix} & \begin{matrix} w_1 & 1/2 \\ w_2 & 3/5 \\ w_3 & 1 \end{matrix} \\
 \text{activation} \Rightarrow 1 \times \frac{1}{2} + 0 \times \frac{3}{5} - 1.5 \times 1 & \text{output} = ? \quad y = 0 \\
 \Rightarrow -1 & -1 \leq \theta, \text{ where } \theta \geq 0 \quad \therefore y = 0 \quad (\text{binary output})
 \end{array}$$



Eg. same case ↑,

$$\begin{array}{l} w_1 \in \{0, 1\} \\ w_2 \in \{0, 1\} \end{array} \quad w_1 = \frac{1}{2} \quad \theta = \frac{3}{4}$$

$x_1$	$x_2$	$= 0$	$y = 0$
0	0	$= 0$	$y = 0$
0	1	$= 0.5$	$y = 0$
1	0	$= 0.5$	$y = 0$
1	1	$= 1$	$y = 1$

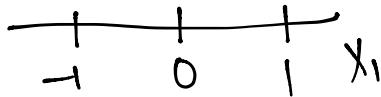
only time  $y=1$  is when  $x_1$  &  $x_2$  are True: conjunction AND Boolean

\* This gives us a way to define a logic gate!

\* we want  $y$  to be a OR LOGIC GATE.

$$\begin{array}{l} x_1 \quad x_2 \\ 0 \quad 0 \\ 0 \quad 1 \\ 1 \quad 0 \\ 1 \quad 1 \end{array} \quad \begin{array}{l} = 0 \\ = 0.5 \\ = 0.5 \\ = 1 \end{array} \quad \begin{array}{l} y = 0 \\ y = 1 \\ y = 1 \\ y = 1 \end{array} \quad \begin{array}{l} w_1 = 0.5 \\ w_2 = 1.0 \\ \theta = 0.5 \end{array}$$

\* Not (one variable)



$$w_1 = -1$$

$$\theta = 0$$

\* AND, OR, NOT are all expressible as perceptron units

\* Lets design a XOR using a network of perceptron

\* XOR

$x_1$	$x_2$	AND	$w$	XOR	OR	$\theta = 0$
0	0	0	-1	0	0	
0	1	0	-1	1	1	
1	0	0	-1	1	1	
1	1	1	-2	0	1	

OR - (AND \* 2)

\* Perception Rule

$$w_i = w_i + \Delta w_i \quad \left\{ \begin{array}{l} \text{one change} \\ \text{weight by } \Delta w_i \end{array} \right.$$

$$\Delta w_i = \eta (y - \hat{y}) x_i$$

$$\hat{y} = (\sum w_i x_i \geq 0)$$

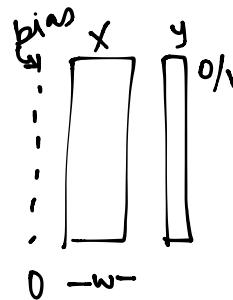
if it does  
not work  
it is not  
linearly separable

$y$ : target

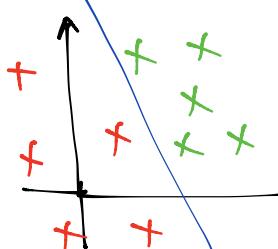
$\hat{y}$ : output

$\eta$ : learning rate

$x$ : input



$y$	$\hat{y}$	$y - \hat{y}$
0	0	0
0	1	-1
1	0	-1
1	1	0



data is linearly separable in this case.

and perception rule will find it!. The algorithm should only be run when data is linearly separable.

\* Gradient descent: a more robust algorithm for non-linear separability  
\* imagine the output is not thresholded.

$$a = \sum_i x_i w_i \quad \hat{y} = \{\hat{a} \geq 0\}$$

$$E(w) = \frac{1}{2} \sum_{(x,y)} (y - \hat{a})^2$$

use chain rule derivative

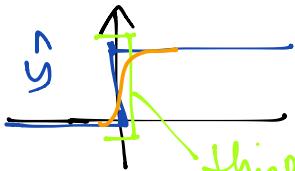
$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{(x,y) \in D} (y - a)^2 \\
 &= \sum_{(x,y) \in D} (y - a) \frac{\partial}{\partial w_i} - \sum_i x_i w_i' \\
 &= \sum_{(x,y) \in D} (y - a)(-x_i) \quad \text{looks just like perceptron rule!}
 \end{aligned}$$

\* let's take a look further:

$$\begin{aligned}
 \Delta w_i &= \eta_v (y - \hat{y}) x_i \text{ perceptron : guarantee : finite convergence} \\
 \Delta w_i &= \eta_v (y - a) x_i \rightarrow \text{not thresholding : gradient descent : calculates} \\
 &\quad \text{activation} \quad \text{robust to data not linearly separable,} \\
 &\quad \text{move the weights in the -ve direction.} \quad \text{converge to local optima}
 \end{aligned}$$

\* So why don't we use gradient descent all the time at even  $\delta$

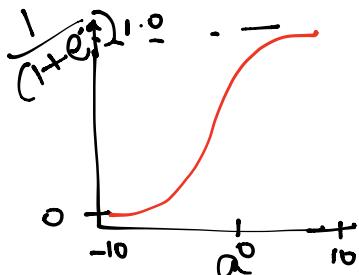
\* non-differentiable: it a discontinuous function



this part makes it discontinuous, lets make it smooth using sigmoid S

\* Sigmoid: - differentiable threshold

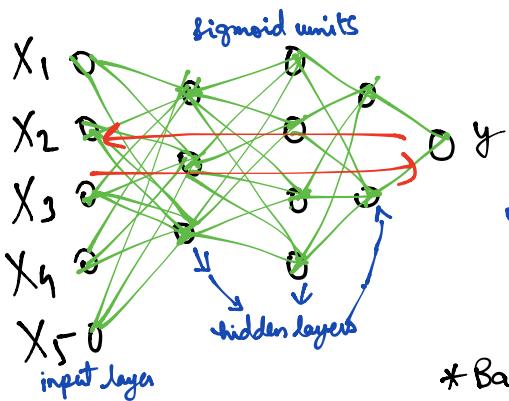
$$\begin{aligned}
 \delta(a) &= \frac{1}{1+e^{-a}} \\
 a \rightarrow \infty & \quad \delta(a) \rightarrow 0 \\
 a \rightarrow +\infty & \quad \delta(a) \rightarrow 1
 \end{aligned}$$



activation get to 0  $\Rightarrow$  as a get negative derivatives goes to 0  
when a grows in size, derivative goes toward 1. (As activation get +ve)

## \* Neural Networks

### Sigmoid units



\* This mapping from input to output is differentiable.  
in terms of weight

\* we can move all the weights from input to output  
to get a desirable result.

output \* This lead to the idea of back propagation  
+ how moving the weigh up or downwards to produce  
the out put we want.

\* Back Propagation: computing derivatives . chain rule  
with respect to all different weights of the network.

\* Many local optima: (it analogous to a perceptron since we are not doing the thresholding  
we trying to set the weights so the error is low, some time even after setting all the weights it makes  
the error worse, but infact it gets stuck.)

\* Gradient descent: minimizing errors  $\rightarrow$  optimizing weights

$\rightarrow$  gradient descent

$\rightarrow$  use advanced methods:

\* momentum descent: as we are doing gradient descent, and we get stuck in an error, we continue (bounce out)  
forward in the direction we have been working towards

\* higher order derivatives  $\rightarrow$  instead of just using variation in weights

\* randomized optimisation: apply to make things more robust

\* penalty for complexity: we just don't want to minimize errors, but we want to penalize for using a  
a system to complex. (overfitting ex: Decision Tree growing too much, then we use pruning to fill out the tree  
by barding up, limiting min-sample-split)

$\hookrightarrow$  regression overfitting (order of polynomial)

$\hookrightarrow$  size of the Decision Tree

\* Here in Neural Network  $\rightarrow$  More and more node with make mapping more complex

$\rightarrow$  more nodes

$\rightarrow$  more layers

$\rightarrow$  we can have complex network by larger numbers of the weight. So sometimes we want to penalize  
the network by keeping the numbers in a reasonable range

#### \* Restrictive Bias & Inductive Bias

- \* Restrictive Bias: tells us about the representation power of the data structure we are using  
ex in this case \* Network of Neurons and set of hypothesis you are able to consider
- \* we are restricting to a subset of models

#### \* Neural Networks

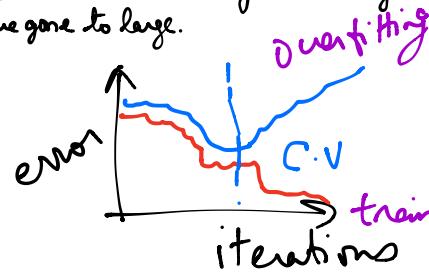
##### \* Restrictions:

- 1) we started with linear planes - Perceptrons
- 2) Sigmoids: much more complex  
hence not much of restriction

→ it represents → Boolean function: using network of threshold-like units  
→ continuous function? → no discontinuity: we can do it with single hidden  
→ Arbitrary: stitch together - two hidden layers

#### \* There is a worry of OVERFITTING!

→ we limit it by using CV, to decide how many hidden layers to use, how many nodes to use, when to stop training when weights have gone to large.



error keeps dropping in neural networks as we perform more iterations

#### \* Preference Bias:

- \* Something about the algorithm we are to learn, why we prefer 1 over the other
- \* ex: Decision Trees: we prefer nodes near the top with high information gain, longer trees, trees that were shorter to deeper

#### \* Neural Network: How do we start the algorithm: do we start with 0 or 1?

- \* lets initialize the weights to something, we can't update something which is undefined.
  - \* small random value: since we run the algorithm multiple times, what if the algorithm gets stuck, we don't want the algorithm to get stuck there again. It gives some variability which helps avoiding local minima.
  - \* why we start with small values → if weights get too big, it sometimes lead to overfitting (high complexity)
  - \* So Preference Bias → small values (low complexity), we prefer simpler explanation to complex.
  - \* we prefer correct answers, the simpler explanation is preferred Occam's Razor.
- \* Occam's Razor: entities should not be multiplied unnecessarily, given Neural Networks! → there is often a lot of unnecessary multiplication. We are not doing any better at fitting data. we should not multiply further and make it more complex. choose the simple solution.

\* we get better generalization with simpler solutions in Supervised Learning.

\* ReCap

\* Perceptrons  $\rightarrow$  threshold units

\* networks can produce any Boolean function

\* perceptron rule — finite time for linearly separable

\* general differentiable rule  $\rightarrow$  backpropagation & gradient descent

\* performance restriction bias of neural networks

Quiz 1:

Network inputs:  $[1, 2, 3]$

layered network weights

$[1, 1, -5] \ \& \ [3, -4, 2]$

$[2, -1]$

$$x_1 w_1 \Rightarrow 1 * 1 = 1$$

$$x_2 w_2 \Rightarrow 2 * 1 = 2$$

$$x_3 w_3 \Rightarrow 3 * -5 = -15$$

$$\text{activation} = -12$$

$$x_1 w_1 \Rightarrow 1 * 3 = 3$$

$$x_2 w_2 \Rightarrow 2 * -4 = -8$$

$$x_3 w_3 \Rightarrow 3 * 2 = 6$$

$$\text{activation} : 1$$

$[-12, 1]$  hidden layer

$$x_1 w_1 \Rightarrow -12 * 2 = -24$$

$$x_2 w_2 \Rightarrow 1 * -1 = -1$$

$[-25] //$

Quiz 3:

$$[\text{input}, \text{input}],$$

$$[[3, 2], [-1, 4], [3, -5]],$$

$$[[1, 2, -1]]$$

]

Input	Output	Operation
$[a, b] * [3, 2]$	$= [3a, 2b]$	$* [1, 2, -1] = [3a, 2b]$
$[-1, 4]$	$[-a, 4b]$	$[-2a, 8b]$
$[3, -5]$	$[3a, -5b]$	$[-3a, 5b]$

$$= [-2a, 15b] //$$

Design a XOR Network

								$\theta = 1$	
$x_1$	$x_1$	$w$	Activ	AND	OR	OR-AND	XOR		
0	0	.	.	0	0	0	0		
1	0	.	.	0	1	1	1		
0	1	.	.	0	1	1	1		
1	1	.	.	1	1	0	0		

design OR & AND then OR-AND = XOR       $[3, 2, 4]$

