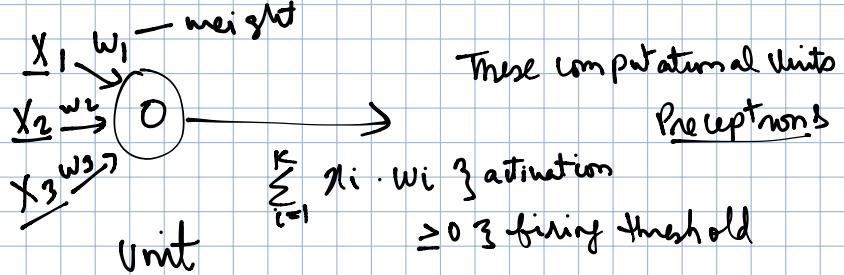


* Neural Networks

* They are computational units.

* Artificial Neural Networks

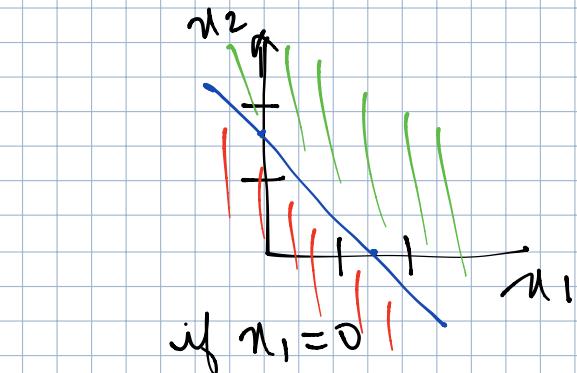


These computational units are called
Perceptrons

Ques:

$$\begin{array}{ll}
 x_1 & 1 \quad w_1 1/2 \\
 x_2 & 0 \quad w_2 3/5 \\
 x_3 & -1.5 \quad w_3 1
 \end{array}
 \quad \text{Threshold } \theta \geq 0$$

activation $\Rightarrow 1 \times \frac{1}{2} + 0 \times \frac{3}{5} - 1.5 \times 1 = -1$ output=? $y=0$
 $\Rightarrow -1 \leq \theta$, where $\theta \geq 0 \quad \therefore y=0$ (binary output)



$$\begin{aligned}
 &\text{set } n_1=0 \quad n_2=? \\
 &\hookrightarrow 0 \times \frac{1}{2} + \square \times \frac{3}{4} = \frac{3}{4} \\
 &\text{solution} \\
 &\Rightarrow n_2 \frac{1}{2} = \frac{3}{4} \\
 &n_2 = \frac{3}{2}
 \end{aligned}$$

$$\begin{aligned}
 &0,1 \quad \text{weights matter} \\
 &w_1 = 1/2 \\
 &w_2 = 3/4 \\
 &\theta = 3/4
 \end{aligned}$$

return 1
return 0

Eg. same case ↑,

$$\begin{array}{l} \{x_1 \in \{0, 1\}, w_1 = \frac{1}{2} \\ x_2 \in \{0, 1\}, w_2 = \frac{1}{2} \end{array} \quad \theta = \frac{3}{4}$$

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.

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

* 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 perceptrons

* XOR

x_1	x_2	AND	\oplus	XOR	OR	\ominus
0	0	0	1	0	0	0
0	1	0	-1	1	1	-1
1	0	0	-1	1	1	-1
1	1	1	-2	0	1	-1

$$\Rightarrow R = (A \cap B)^{*2}$$

* Perception Rule

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

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

$$\hat{y} = \left(\sum_i w_i x_i \geq 0 \right)$$

Vegetable

W. J. H. G.

linen

30 0

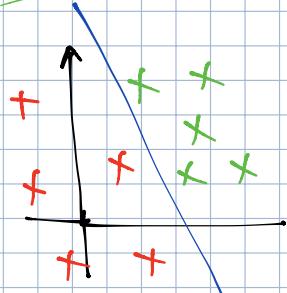
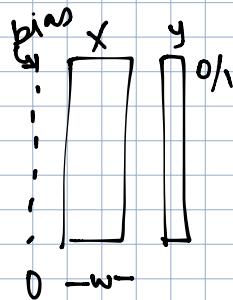
$$\begin{array}{r}
 \text{Step 1: Linearity} \\
 \begin{array}{ccccccc}
 & y & - & 5 & = & 0 & - 1 \\
 & 0 & - & 0 & - & 0 & - 0 \\
 \hline
 & - & 5 & - & 0 & - & 1
 \end{array}
 \end{array}$$

y : target

$$\vec{y} = \text{output}$$

η : learning rate

π: input



data is linearly separable in this case.

and perceptron 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.

$$a = \sum_i x_i w_i \quad g = \{a > 0\}$$

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

use chain rule derivative

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

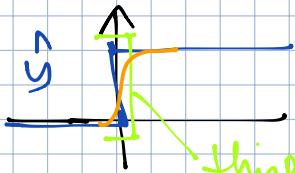
* let's take a look further:

$$\Delta w_i = \gamma_i (y - \hat{y}) x_i \text{ perceptron : guarantee : finite convergence}$$

$\Delta w_i = \gamma_i (y - a) x_i$ → not thresholding : gradient descent : calculates
robust to data not linearly separable,
move the weights in the -ve direction.

* So why don't we use gradient descent all the time at even δ

* non-differentiable : it a discontinuous function



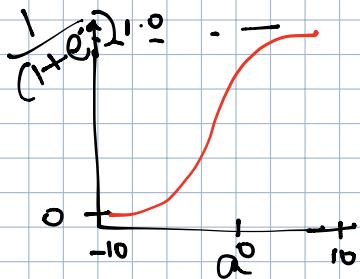
this part makes it discontinuous, lets make it smooth using sigmoid S

* Sigmoid : - differentiable threshold

$$f(a) = \frac{1}{1+e^{-a}}$$

$$a \rightarrow \infty \quad f(a) \rightarrow 0$$

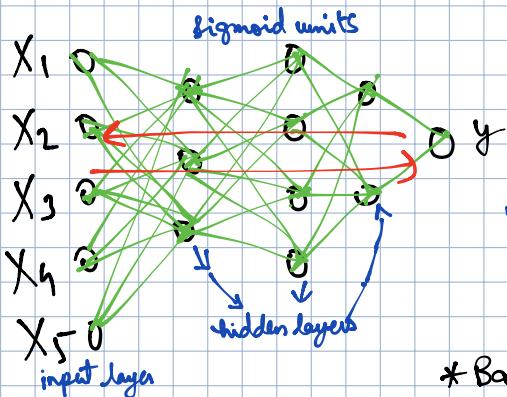
$$a \rightarrow +\infty \quad f(a) \rightarrow 1$$



activation get to 0 → as a get negative derivatives goes to 0
when a grows in size, derivative goes towards 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 optimizers : 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 system to complex . (overfitting ex: Decision Tree growing too much, then we used pruning to filling out the tree by backtracking , limiting min-sample-split)

\hookrightarrow regression overfitting (order of polynomial)

\hookrightarrow size of the Decision Tree

* Here in Neural Network \rightarrow More and more mode 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

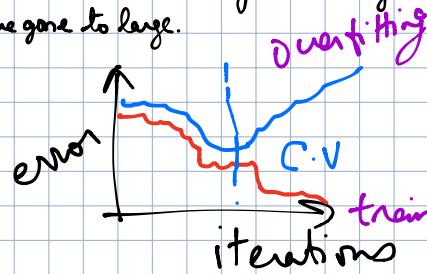
* Restrictions:

- 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 here 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/intrinsic bias of neural networks

Quiz 1:

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

Layered network weights

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

$$[2, -1]$$

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

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

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

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

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

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

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

$$\text{activation} : 1$$

$[-12, 1]$ hidden layer

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

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

$[-25] //$

Quiz 3:

$$[\begin{bmatrix} \text{input}, \text{input} \end{bmatrix}, \\ \begin{bmatrix} [3, 2], [-1, 4], [3, -5] \end{bmatrix}, \\ \begin{bmatrix} [1, 2, -1] \end{bmatrix}]$$

$$[\begin{array}{c} \text{Input} \\ \begin{bmatrix} [a, b] * [3, 2] \\ [-1, 4] \\ [3, -5] \end{bmatrix} \\ \begin{bmatrix} [3a, 2b] \\ [-a, 4b] \\ [3a, -5b] \end{bmatrix} \end{array} \quad \begin{array}{c} \text{Output} \\ * [1, 2, -1] = [3a, 2b] \\ [-2a, 8b] \\ [-3a, 5b] \end{array}]$$

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

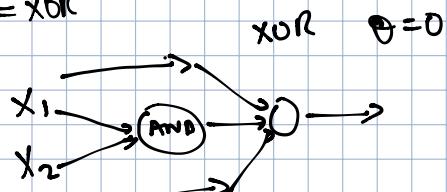
Design a XOR Network

$\Theta = 1$

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

x_1	x_2	w	Act	OR
*	*			



OR gates $\begin{bmatrix} 3, 2 \\ 4, 5 \end{bmatrix}$
Network

$\begin{bmatrix} 2, 5 \\ 1, 1 \end{bmatrix}$

input

$$\begin{bmatrix} 0, 0 \\ 4, 5 \\ 2, 5 \end{bmatrix} \begin{bmatrix} [3, 2] \\ [4, 5] \\ [2, 5] \end{bmatrix} = \begin{bmatrix} 6, 0 \\ [0, 5] \\ [0, 5] \end{bmatrix} \checkmark$$

$$\begin{bmatrix} 0, 1 \\ 4, 5 \\ 2, 5 \end{bmatrix} \begin{bmatrix} [3, 2] \\ [4, 5] \\ [2, 5] \end{bmatrix} = \begin{bmatrix} [0, 2] \\ [0, 5] \\ [0, 5] \end{bmatrix} = \begin{bmatrix} 0, 12 \\ [0, 5] \\ [0, 5] \end{bmatrix}$$

$$\begin{bmatrix} 1, 0 \\ 4, 5 \\ 2, 5 \end{bmatrix} \begin{bmatrix} [3, 2] \\ [4, 5] \\ [2, 5] \end{bmatrix} = \begin{bmatrix} [3, 0] \\ [4, 0] \\ [2, 0] \end{bmatrix} = \begin{bmatrix} 9, 0 \\ [4, 0] \\ [2, 0] \end{bmatrix}$$

$$\begin{bmatrix} 1, 1 \\ 4, 5 \\ 2, 5 \end{bmatrix} \begin{bmatrix} [3, 2] \\ [4, 5] \\ [2, 5] \end{bmatrix} = \begin{bmatrix} [3, 2] \\ [4, 5] \\ [2, 5] \end{bmatrix} = \begin{bmatrix} [9, 12] \\ [4, 12] \\ [2, 5] \end{bmatrix}$$

Output

OR gate

$$[0, 0] * [1, 1] = 0 \quad y = 0 \quad \theta = 8 : \text{OR}$$

$$[0, 12] * [1, 1] = 12 \quad y = 1$$

$$[9, 0] * [1, 1] = 9 \quad y = 1$$

$$[9, 12] * [1, 1] = 21 \quad y = 1$$

$$\cdot \quad \theta = \text{length } x_{out} - 1 \quad @ [0, 1] \\ [1, 0]$$

AND $\theta = 20$

$$\theta = (x_{1\text{out}} + x_{2\text{out}}) - 1 \text{ @ } [1, 1]$$

output

$$[0, 0] * [1, 1] = 0 \quad y = 0$$

$$[0, 12] * [1, 1] = 12 \quad y = 0$$

$$[9, 0] * [1, 1] = 9 \quad y = 0$$

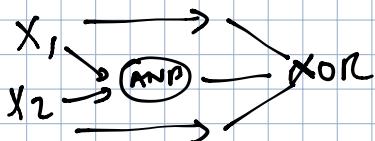
$$[9, 12] * [1, 1] = 21 \quad y = 1 \quad \text{for AND gate if output } y=1 \text{ then } y=y+2 \Rightarrow: \\ y = -2$$

N.W truth table to generate XOR

XOR:
 $\theta = 0$

Input	OR	AND	OR-AND	XOR
[0, 0]	0	0	0	0
[0, 1]	1	0	0	1
[1, 0]	1	0	0	1
[1, 1]	2	-2	1	0

XOR:



when Input

a) [0, 0] = $x_1 = 0, x_2 = 0$

b) [0, 1] = $x_1 = 0, x_2 = 1$

$$\text{XOR} = 1$$

c) [1, 0] = $x_1 = 1, x_2 = 0$

$$\text{XOR} = 1$$

d) [1, 1] = $x_1 = 1, x_2 = 1, \text{AND} = -2$

$$= x_1 + x_2 - \text{AND} = 1 + 1 - 2 = 0$$

($\therefore X \text{ ORL} = 0$)