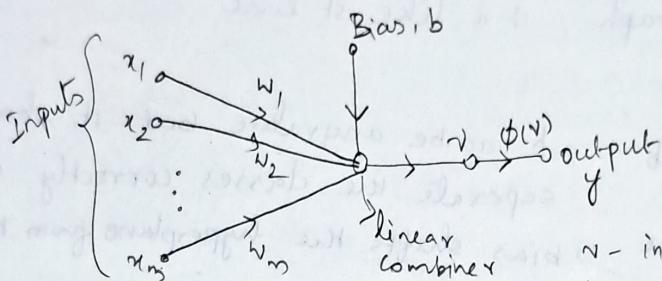


## UNIT - 2

### Single Layer Perceptron:

- Rosenblatt's perceptron / perceptron
- built around a non-linear neuron, called as McCulloch-Pitts model of a neuron.
- such a neuron model consists of a linear combiner followed by a hard limiter (performing the sigmoid function).
- signal flow graph of the perceptron



$v$  - induced local field  
 $\phi(v)$  - hard limit function  
 $w_1, w_2, \dots, w_m$  - synaptic weights  
 $x_1, x_2, \dots, x_m$  - inputs

$$\phi(v) = \begin{cases} 1, & v \geq 0 \\ -1, & v < 0 \end{cases}, \quad \phi(v) = \begin{cases} 1, & v > 0 \\ 0, & v = 0 \\ -1, & v < 0 \end{cases} \text{ signum}$$

$$v = \sum_{i=1}^m w_i x_i + b$$

**Goal of Perceptron:** To correctly classify the set of externally applied stimuli  $x_1, x_2, \dots, x_m$  into one of two classes  $C_1$  or  $C_2$ .

**Decision rule:** for the classification is

$$\text{If output of perceptron } y = +1, \quad x_1, x_2, \dots, x_m \in C_1$$

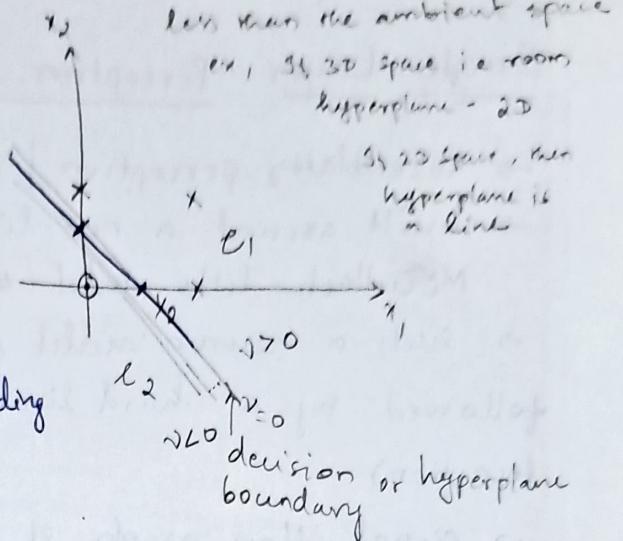
$$\quad \quad \quad " \quad \quad y = -1, \quad x_1, x_2, \dots, x_m \in C_2$$

To understand this better, we need to plot the input variables  $x_1, x_2, \dots, x_m$  on  $m$ -dimensional space, then there are two decision regions separated by a hyperplane defined by

$$\sum_{i=1}^m w_i x_i + b = 0$$

For example, consider OR gate

$x_1$	$x_2$	$v$	$y$
0	0	-1/2	0 $\rightarrow C_2$
0	1	1/2	1
1	0	1/2	1
1	1	3/2	1



what is the equation corresponding to the hyperplane.

$$w_1x_1 + w_2x_2 + b = 0$$

but from the graph, it is like st. line

$$y = mx + c$$

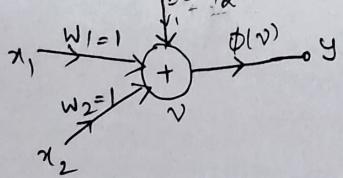
$x_2 = -x_1 + b$ ,  $b$  can be any value but it should separate the classes correctly i.e

→ what should be the value of  $b$ ? → bias shifts the hyperplane from the origin

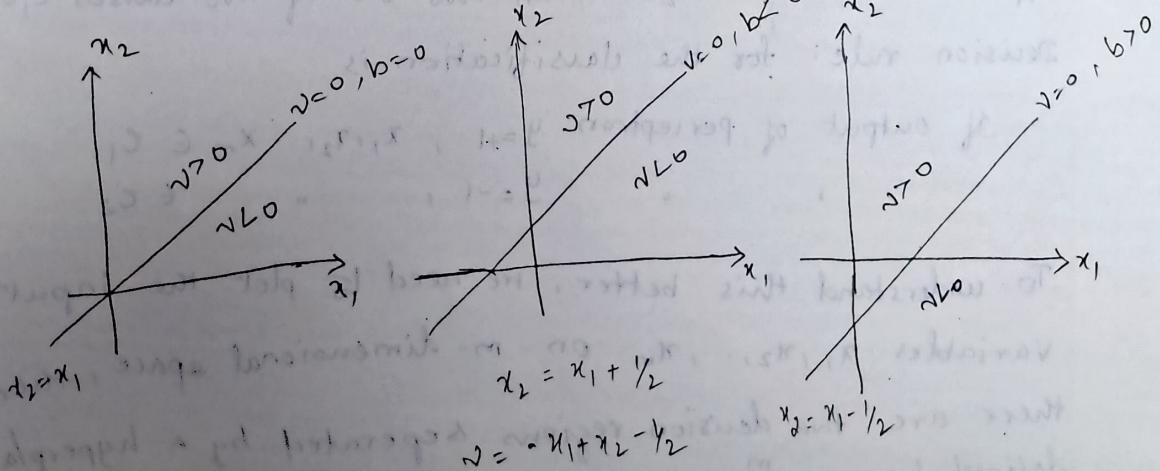
but I assume  $b = 1/2$  (it is a guess from the graph)

$$\therefore x_2 = -x_1 + 1/2 \quad \text{using this fill the table}$$

so, what is the perceptron for OR gate?



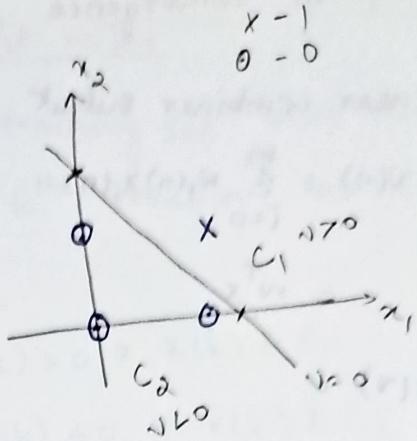
what are the different bias values & corresponding decision boundary?



## AND Gate:

$x_1$	$x_2$	$v$	$y$
0	0	-3/2	0
0	1	-1/2	0
1	0	-1/2	0
1	1	1/2	1

$\rightarrow C_1$



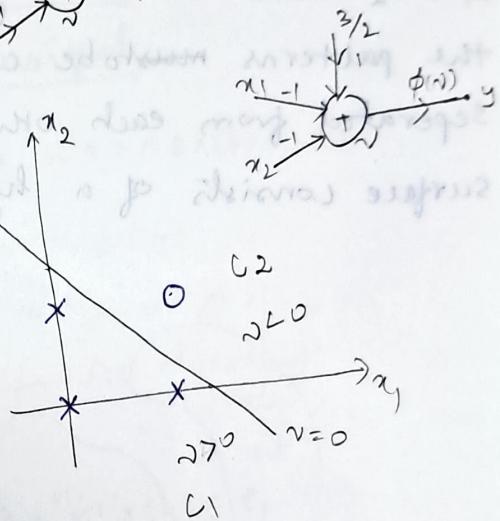
eq. of line  $x_2 = -x_1 + 3/2$ ,  $3/2$  is a guess

$$v = +x_1 + x_2 - 3/2$$

## NAND Gate:

$x_1$	$x_2$	$v$	$y$
0	0	3/2	1
0	1	1/2	1
1	0	1/2	1
1	1	-1/2	0

$\rightarrow C_2$



eq. of line  $x_2 = -x_1 + 3/2$

$$v = -x_2 + x_1 + 3/2$$

This can be applied to any boolean function.

The synaptic weights  $w_1, w_2, \dots, w_m$  of the perceptron can be adapted on an iteration-by-iteration basis.

for the adaptation we may use an error-correction rule known as the perceptron convergence algorithm.

for a given set of inputs  $x_1, x_2, \dots, x_n$  and desired output  $v$ , the perceptron converges if the error function  $E = \sum_{i=1}^n (y_i - v_i)^2$  decreases towards zero. This means that the output  $y$  approaches the desired value  $v$  as the number of iterations increases.

$\Rightarrow x$  from  $0 < x^T w$

$\Rightarrow x$  from  $0 \leq x^T w$

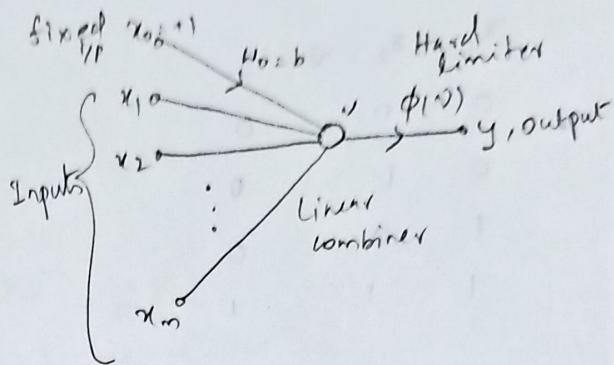
## Perception Convergence Theorem:

linear combiner output

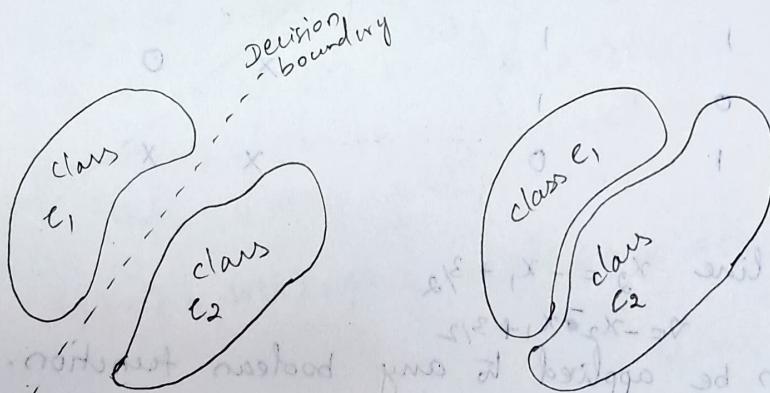
$$v(n) = \sum_{i=0}^m w_i(n)x_i(n)$$

$$= w^T x$$

$$\phi(v) =$$



For the perception to function properly, the two classes  $C_1 \times C_2$  must be linearly separable. It means that the patterns ~~must~~ to be classified must be sufficiently separated from each other to ensure that the decision surface consists of a hyperplane.



- a) a pair of linearly separable patterns
- b) a pair of non-linearly separable patterns.

Let  $C_1 \times C_2$  be 2 classes and let  $H_1 \subset C_1$  &  $H_2 \subset C_2$  be available. Given these training set, the training process involves the adjustment of weight vector  $w$  in such a way that 2 classes are linearly separable. i.e.  $\exists w$  s.t we may state

$w^T x > 0$  for every  $x \in C_1$  } given the  $H_1 \times H_2$ , training problem is to find  $w$  s.t these  
 $w^T x \leq 0$  for every  $x \in C_2$  } inequalities are satisfied  
 our choice

## Algorithm for adapting the weight vector:

① If the  $k^{\text{th}}$  member of training set  $x(k)$  is correctly classified at the  $k^{\text{th}}$  iteration, no correction is required for  $w$ .

$$w(k+1) = w(k) \quad \text{if } w^T x(k) > 0 \quad \& \quad x(k) \in C_1$$

$$w(k+1) = w(k) \quad \text{if } w^T x(k) \leq 0 \quad \& \quad x(k) \in C_2$$

② Otherwise, the weight vector of the perceptron is updated according to the following rule.  
i.e if  $x(k)$  is incorrectly classified.

$$w(k+1) = w(k) - \eta(k) x(k) \quad \text{if } w^T x > 0 \quad \& \quad x(k) \in C_2$$

$$y \text{ can be either } 1 \text{ or } 0 \quad \& \quad e(k) = y_d - y(k)$$

$$\text{If } w^T x > 0 \Rightarrow y=1 \quad \& \quad x(k) \in C_2 \Rightarrow y_d=0 \Rightarrow e(k)=0-1=-1$$

$$\text{WKT, } w(k+1) = w(k) + \eta e(k) x(k) \rightarrow \text{LMS algorithm}$$

$$= w(k) - \eta x(k)$$

$$w(k+1) = w(k) + \eta x(k) \quad \text{if } w^T x \leq 0 \quad \& \quad \underbrace{x(k)}_{y=-1/0} \in \underbrace{C_1}_{y_d=1}$$

$$e(k) = y_d - y = 1 - 0 = 1$$

where  $\eta(n)$  - learning rate parameter (can be fixed or adaptive)

Proof: If there is a weight vector  $w^*$  s.t  $\phi(w^* x(k)) = y(k)$   $\forall k$  then for any starting vector  $w$ , the perceptron learning rule will converge to a weight vector that gives the correct response for all training patterns & it will do so in finite no. of steps.

Assumptions: for ease of analysis

- i) Inputs to the perceptron originates from a linearly separable classes
- ii)  $w(1) = 0$  - initial weight
- iii)  $\eta = 1$

Now consider that perceptron incorrectly classified

$$w^T x < 0 \quad \text{& } x(k) \in C_1$$

$$\therefore w(k+1) = w(k) + x(k)$$

$$k=1, \quad w(2) = w(1) + \overset{\text{chosen}}{x(1)}$$

$$w(2) = x(1)$$

$$k=2, \quad w(3) = w(2) + x(2)$$

$$= x(1) + x(2)$$

$$k=3, \quad w(4) = w(3) + x(3)$$

$$= x(1) + x(2) + x(3)$$

In general,

$$w(k+1) = x(1) + x(2) + x(3) + \dots + x(k) \quad \stackrel{\text{①}}{\rightarrow} \quad w^T x < 0, \quad x(k) \in C_1$$

$$w(k+1) = -x(1) - x(2) - x(3) - \dots - x(k) \quad \rightarrow \quad w^T x > 0, \quad x(k) \in C_2$$

Since the classes  $C_1$  &  $C_2$  are assumed to be linearly correctly classified terms

separable,  $\exists$  a solution  $w_0$  for which  $w^T x > 0$  that can classify all inputs correctly.

for a fixed solution  $w_0$ , we may define a +ve number  $\alpha$  as

$$\alpha = \min_{x(n) \in H_1} w_0^T x(n)$$

Multiply eq. ① by  $w_0^T$

$$\Rightarrow w_0^T w(k+1) = \underbrace{w_0^T x(1)}_{\substack{\hookrightarrow \text{each term is} > \alpha, \text{ bcs } \alpha = \text{min. value}}} + \underbrace{w_0^T x(2)}_{> \alpha} + \underbrace{w_0^T x(3)}_{> \alpha} + \dots + \underbrace{w_0^T x(k)}_{> \alpha}$$

$$w_0^T w(k+1) \geq k\alpha$$

Squared norm of  $w_0^T w(k+1) \geq k\alpha$  by using Cauchy-Schwartz inequality

$$\|w_0\|^2 \|w(k+1)\|^2 \geq k^2 \alpha^2 \quad (\text{states that-})$$

$$\geq (k^2 \alpha^2) \quad k=1, 2, \dots, n$$

$$\|w(k+1)\|^2 \geq \frac{k^2 \alpha^2}{\|w_0\|^2} \quad \rightarrow \text{lower bound on } w$$

③

To find the upper bound, consider

$w(k+1) = w(k) + x(k)$   $\quad \text{for } k=1, \dots, n \quad x(k) \in H,$   
obtain the squared Euclidean norm on both sides

### NORM:

A map  $\|\cdot\|: V \rightarrow \mathbb{R}$ , where  $V$  is a vector space over  $\mathbb{R}$  s.t.

- $\|x\| \geq 0 \quad \forall x \in V; \quad \|x\| = 0$
- $\|\alpha x\| = |\alpha| \|x\| \quad \forall x \in V, \alpha \in \mathbb{R}$
- $\|x+y\| \leq \|x\| + \|y\| \quad \forall x, y \in V$

is called a norm on  $V$

$$\textcircled{2} \Rightarrow \|w(k+1)\|^2 = \|w(k)\|^2 + \|x(k)\|^2 + 2w^T(k)x(k)$$

$$\begin{aligned} \|w(k+1)\|^2 &= (w(k) + x(k))^T (w(k) + x(k)) \\ &= w^T(k)w(k) + w^T(k)x(k) + x^T(k)w(k) + x^T(k)x(k) \\ &= w^T(k)w(k) + 2w^T(k)x(k) + x^T(k)x(k) \\ &= \|w(k)\|^2 + \|x(k)\|^2 \quad \leftarrow 0 \text{ misclassification} \end{aligned}$$

$$\|w(k+1)\|^2 \leq \|w(k)\|^2 + \|x(k)\|^2$$

$$\|w(k+1)\|^2 - \|w(k)\|^2 \leq \|x(k)\|^2$$

$$\|w(k+1)\|^2 \leq \sum_{k=1}^n \|x(k)\|^2, \quad \because \|w(0)\|^2 = 0$$

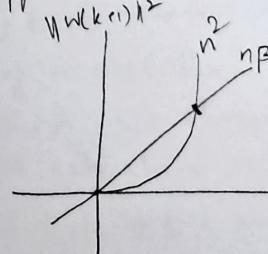
$$\leq n\beta$$

where  $\beta = \max_{x(k) \in H} \|x(k)\|^2$

$$\|w(k+1)\|^2 \leq n\beta \quad \text{--- (4)} \quad \dashrightarrow \text{upper limit}$$

$$\frac{\|x\|^2}{\|w\|^2} \leq \|w(k+1)\|^2 \leq n\beta$$

st line  
parabola



It is impossible to satisfy (3) & (4) simultaneously

but at some point those 2 coincides beyond that  
 $w(k)$  will not change at all.

$\exists n = n_{\max}$  s.t

$$\frac{m_{\max}^2 \alpha^2}{\|w_0\|^2} = n_{\max} \beta$$
$$\Rightarrow n_{\max} = \frac{\|w_0\|^2 \beta}{\alpha^2}$$

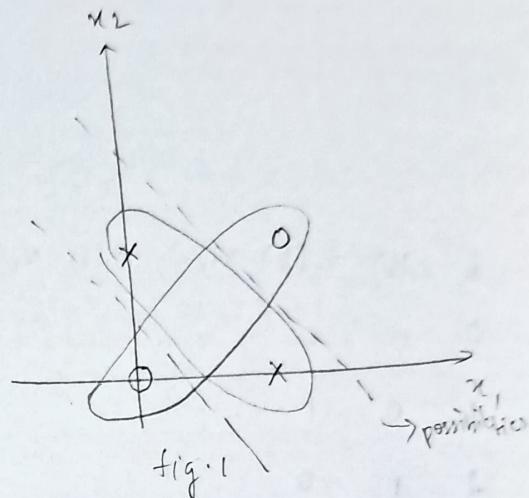
We have thus proved that for  $\eta=1$ ,  $w(0)=0$ , it gives that a solution vector  $w_0$  exists, the rule for adapting the synaptic weights of the perceptron must terminate after at most  $n_{\max}$  iterations.

## Multilayer Perceptrons (chapter 4 in Simon Haykin)

### XOR Problem

- <sup>single layer</sup> Perceptron cannot classify input patterns that are not linearly separable.
- However nonlinearly separable patterns are of common occurrence.
- This problem arises in XOR problem.

$k$	0	1	2	3
$x_1(k)$	0	0	1	1
$x_2(k)$	0	1	0	1
$d(k)$	0	1	1	0
$c_2$	$c_1$	$c_1$	$c_2$	



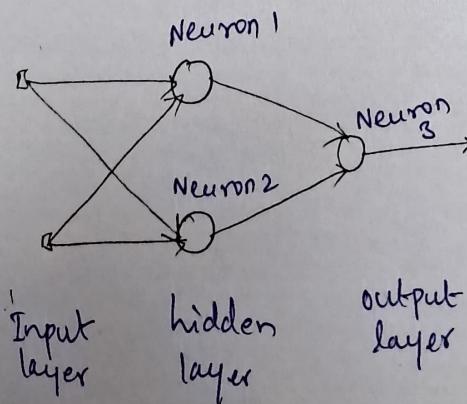
→ Generally, single neuron with 2 inputs results in a straight line for a decision boundary in the input space.

$$\begin{array}{l} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \quad \begin{array}{l} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array}$$

$c_1 y = 1$   
 $c_2 y = -1/0$

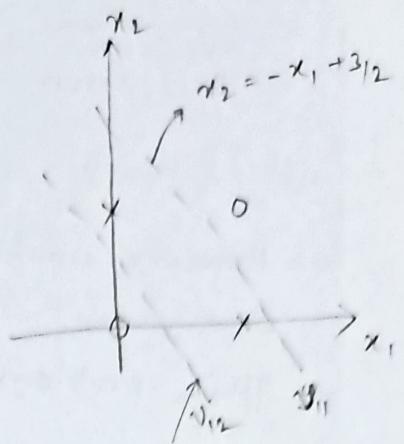
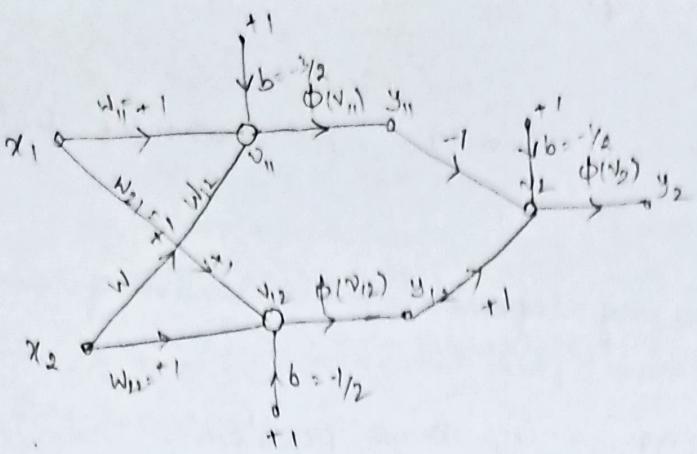
→ Fig. 1 clearly shows that we cannot construct a straight line for a decision boundary, which implies with single perceptron we cannot solve this problem.

→ We may solve XOR problem by using a single hidden layer with 2 neurons, as shown in SFG.



Assumptions:

- Each neuron is represented by a McCulloch-Pitts model, which uses threshold function as activation
- Bits 0 & 1 are represented by the levels 0 & +1 respectively.



$$v_{11} = x_1 + x_2 - \frac{3}{2}$$

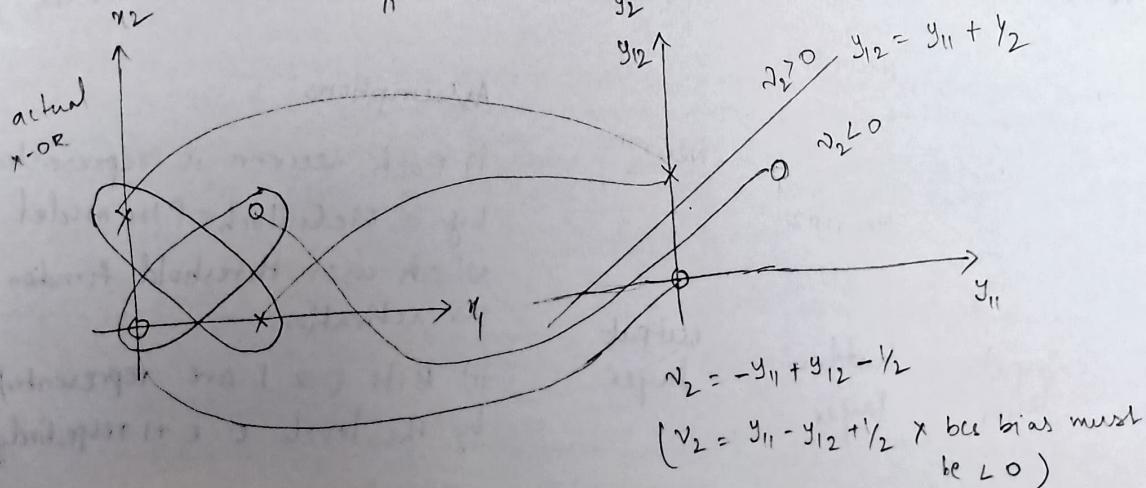
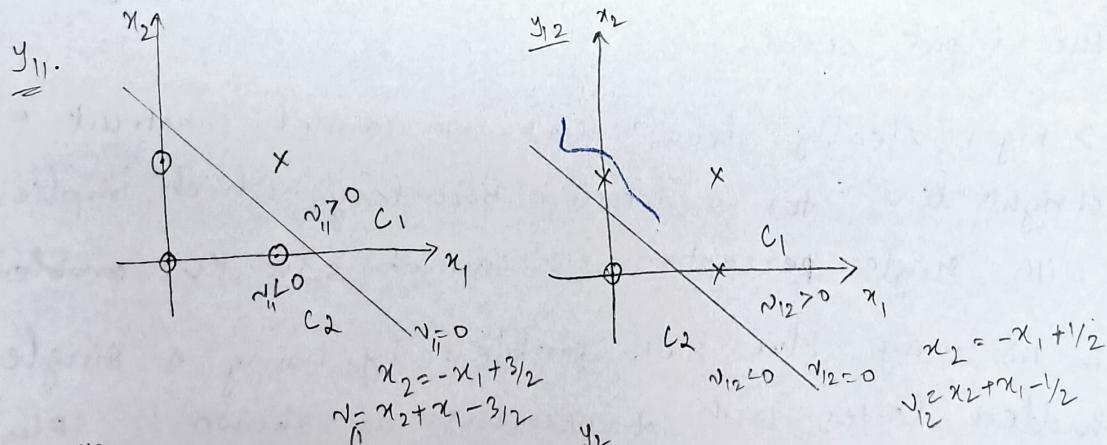
$$v_{12} = -x_1 + \frac{1}{2}$$

$$v_{12} = x_1 + x_2 - \frac{1}{2}$$

AND

OR

$k$	$x_1(k)$	$x_2(k)$	$v_{11}(k)$	$y_1(k)$	$v_{12}(k)$	$y_1(k)$	$v_{21}(k)$	$y_2(k)$
0	0	0	$-\frac{3}{2}$	0	$-\frac{1}{2}$	0	$-\frac{1}{2}$	0
1	0	1	$-\frac{1}{2}$	0	$\frac{1}{2}$	1	$\frac{1}{2}$	1
2	1	0	$-\frac{1}{2}$	0	$\frac{1}{2}$	1	$\frac{1}{2}$	1
3	1	1	$\frac{1}{2}$	1	$\frac{3}{2}$	1	$-\frac{1}{2}$	0



① AND gate:

$$\begin{array}{ccc} x_1 & x_2 & y_d \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{array} \quad \left. \begin{array}{l} x(1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ x(2) = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \\ x(3) = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ x(4) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \end{array} \right\} l_1$$

let

$$w(1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \eta = 1$$

$$x(1) \in l_1$$

$$w^T(1)x(1) = 0 \geq 0 \quad : \quad y(1)=1 \Rightarrow \text{which is correctly classified}$$

$$\therefore y_d(1)=0$$

update w as follows.

$$w(2) = w(1) - x(1) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

$$w^T(2)x(2) = 0 \geq 0 \quad y(2)=1 \Rightarrow x(2) \in l_1$$

which is incorrectly classified  $\because y_d(2)=0$

$$w(3) = w(2) - x(2) = \begin{pmatrix} 0 \\ -1 \end{pmatrix}.$$

$$w^T(3)x(3) = \begin{pmatrix} 0 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 0 \geq 0 \Rightarrow y(3)=1$$

$$\Rightarrow x(3) \in l_1$$

$$w(4) = w(3) - x(3)$$

which is incorrectly classified

$$= \begin{pmatrix} 0 \\ -1 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

$$w^T(4)x(4) = (-1 \quad -1) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = -2 < 0 \quad \therefore y(4)=0$$

$$\Rightarrow x(4) \in l_2$$

which is incorrect

(3)

$$w(5) = w(4) + x(4)$$

$$= \begin{pmatrix} 1 \\ -1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

again same result will occur.

hence start with different  $w = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$\Sigma$   $w^T(1)x(1) = (1 \ 1) \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0 \quad y(1) = 1 \notin l_2$   
hence incorrectly classified.

$$w(2) = w(1) - x(1)$$

$$= \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$\Pi$   $w^T(2)x(2) = (1 \ 1) \begin{pmatrix} 0 \\ 1 \end{pmatrix} = 1 > 0 \quad y(2) = 1 \notin l_1$   
incorrectly classified  
 $= 0$

$$w(3) = w(2) - x(2)$$

$$= \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$\text{III}$   $w^T(3)x(3) = (0 \ 0) \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0 \quad x(3) \notin l_1$   
incorrect classification.

$$w(4) = w(3) - x(3)$$

$$= \begin{pmatrix} 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$\text{IV}$   $w^T(4)x(4) = (0 \ 0) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0 \geq 0 \quad x(4) \in l_1$   
correct classif?

$$w(5) = w(4) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

but again for other pattern it fails.

$$\omega(1) = \begin{pmatrix} +\gamma_2 \\ -\gamma_2 \end{pmatrix}$$

I  $\omega^T(1)x(1) = 0 \Rightarrow x(1) \in \ell_1$  incorrect

$$\omega(2) = \omega(1) - x(1)$$

$$= \begin{pmatrix} -\gamma_2 \\ -\gamma_2 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} +\gamma_2 \\ -\gamma_2 \end{pmatrix}$$

II  $\omega^T(2)x(2) = \begin{pmatrix} +\gamma_2 \\ -\gamma_2 \end{pmatrix}^T \begin{pmatrix} 0 \\ 1 \end{pmatrix} = -\gamma_2 < 0 \Rightarrow x(2) \in \ell_2$  correct.

$$\omega(3) = \omega(2).$$

III  $\omega^T(3)x(3) = (-\gamma_2 \quad -\gamma_2) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = -\gamma_2 < 0 \Rightarrow x(3) \in \ell_2$  good.

$$\omega(4) = \omega(3)$$

IV  $\omega^T(4)x(4) = (-\gamma_2 \quad -\gamma_2) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = -\gamma_2 - \gamma_2 = -1 \Rightarrow x(4) \in \ell_2$ .

$$\omega(5) = \omega(4) + x(4)$$

$$= \begin{pmatrix} -\gamma_2 \\ -\gamma_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} \gamma_2 \\ \gamma_2 \end{pmatrix}$$

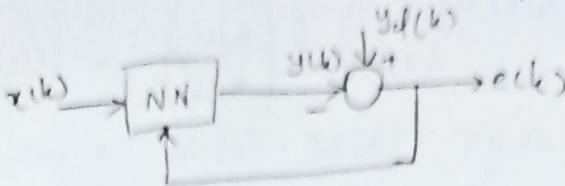
incorrect

000

$$\begin{matrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{matrix}$$

$$\begin{matrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \end{matrix}$$

Back propagation Algorithm: is applied to multilayer FF network.

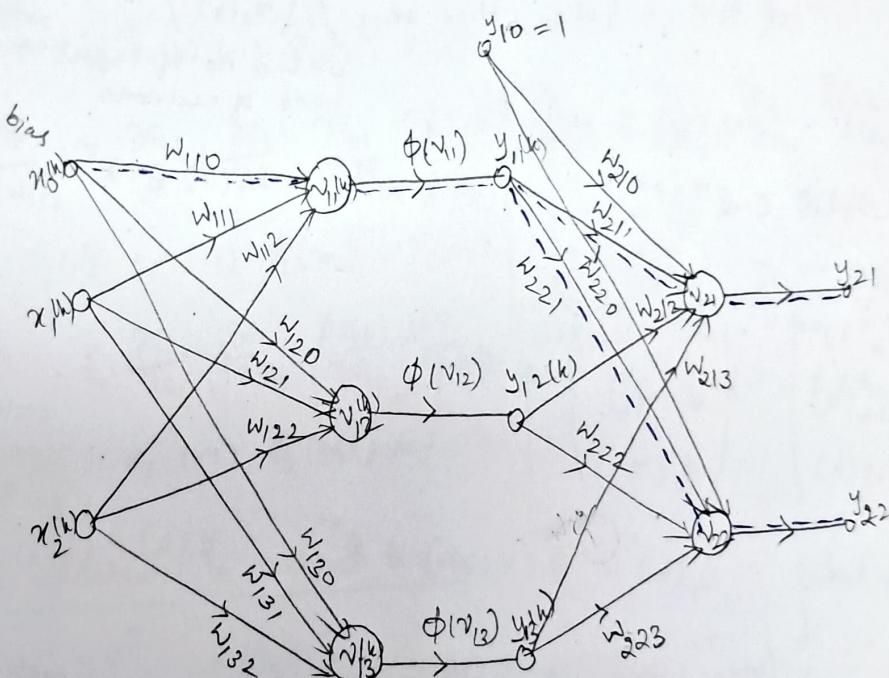


- \* uses continuous differentiable function as activation function
- \* based on error-correction learning rule

There are two parts

1. forward pass : during this information flows in forward direction & weights are fixed
2. Backward pass : error is back propagated & weights are updated.

$$m_0 : m_1 : m_2 \\ 2 : 3 : 2$$



layer:  $l=0$

index:  $j$   
i/P layer

Weij  $\Rightarrow$   $l$ -layer number

$i$  - neuron number

$j$  - input connected to

$l=1$

hidden layer

$l=2$

o/p layer

## Forward Pass

$$x(k) = \begin{pmatrix} x_0(k) \\ x_1(k) \\ x_2(k) \end{pmatrix} \quad y_0(k) \stackrel{\Delta}{=} \begin{pmatrix} 1 \\ x_0(k) \end{pmatrix} = \begin{pmatrix} 1 \\ x_1(k) \\ x_2(k) \end{pmatrix}$$

$$v_{11}(k) = w_{110} x_1 + w_{111} x_1(k) + w_{112} x_2(k) = \sum_{j=0}^2 w_{11j} x_j$$

$$v_{12}(k) = w_{120} + w_{121} x_1(k) + w_{122} x_2(k) = \sum_{j=0}^2 w_{12j} x_j$$

$$v_{13}(k) = w_{130} + w_{131} x_1(k) + w_{132} x_2(k) = \sum_{j=0}^2 w_{13j} x_j$$

$$\begin{aligned} v_1(k) &= \begin{pmatrix} v_{11} \\ v_{12} \\ v_{13} \end{pmatrix} = \begin{pmatrix} w_{110} & w_{111} & w_{112} \\ w_{120} & w_{121} & w_{122} \\ w_{130} & w_{131} & w_{132} \end{pmatrix} \begin{pmatrix} 1 \\ x_1(k) \\ x_2(k) \end{pmatrix}, \quad y_0(k) = 1 \\ &\qquad\qquad\qquad \text{3x3} \rightarrow \text{no. of input nodes} + 1 \\ &\qquad\qquad\qquad \underbrace{\text{no. of neurons}}_{\text{no. of neurons}} \end{aligned}$$

In general,  $v_1(k) = w_1(k) y_0(k)$ ,  $y_0(k) \in \mathbb{R}^{m_0+1}$ ,  $v_1(k) \in \mathbb{R}^{m_1}$ ,  $w_1(k) \in \mathbb{R}^{m_0 \times m_0+1}$

$$\bar{y}_1 \stackrel{\Delta}{=} \begin{pmatrix} y_{11}(k) \\ y_{12}(k) \\ y_{13}(k) \end{pmatrix} = \begin{pmatrix} \phi_1(v_{11}) \\ \phi_1(v_{12}) \\ \phi_1(v_{13}) \end{pmatrix} = \phi_1 \begin{pmatrix} v_{11}(k) \\ v_{12}(k) \\ v_{13}(k) \end{pmatrix} = \phi_1(v_1(k))$$

In general,  $\bar{y}_1(k) = \phi_1(v_1(k)) \in \mathbb{R}^{m_1}$ ,  $y_1(k) = \begin{pmatrix} 1 \\ \bar{y}_1(k) \end{pmatrix}$

$$\begin{pmatrix} v_{21}(k) \\ v_{22}(k) \end{pmatrix} = \begin{pmatrix} w_{210} & w_{211} & w_{212} & w_{213} \\ w_{220} & w_{221} & w_{222} & w_{223} \end{pmatrix} \begin{pmatrix} y_{10}(k) \\ y_{11}(k) \\ y_{12}(k) \\ y_{13}(k) \end{pmatrix} = w_2(k) \begin{pmatrix} 1 \\ \bar{y}_1(k) \end{pmatrix}$$

$$v_2(k) \stackrel{\Delta}{=} w_2(k) y_1(k) \quad 4 \times 1$$

In general,

$$y_2(k) = \begin{pmatrix} y_{21}(k) \\ y_{22}(k) \end{pmatrix} = \phi_2 \begin{pmatrix} v_{21}(k) \\ v_{22}(k) \end{pmatrix} = \phi_2(v_2(k))$$

## Backward Pass

$$e_1(k) = y_{d_1}(k) - y_{s_1}(k)$$

$$e_2(k) = y_{d_2}(k) - y_{s_2}(k)$$

$$e_3(k) = y_{d_3}(k) - y_{s_3}(k) \quad 1 \leq k \leq 2$$

$$\xi(k) = \eta_2 \sum_{s=1}^2 e_s^2(k) = \frac{1}{2}(e_1^2(k) + e_2^2(k))$$

$$\begin{aligned} \frac{\partial \xi(k)}{\partial w_{210}(k)} &= \frac{\partial \xi}{\partial e_1(k)} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial w_{210}} \quad \left| \begin{array}{l} \frac{\partial \xi}{\partial w_{210}} = \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial w_{210}} \\ = e_1(k) \cdot (-1) \phi'_2(v_{21}(k)) \cdot y_{10}(k) \end{array} \right. \\ &= e_1(k) \cdot (-1) \phi'_2(v_{21}(k)) \cdot y_{10}(k) \quad \left| \begin{array}{l} = e_2(k) \cdot (-1) \phi'_2(v_{22}) \cdot y_{10}(k) \end{array} \right. \end{aligned}$$

$$\begin{aligned} \frac{\partial \xi}{\partial w_{211}} &= \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial v_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial w_{211}} \quad \left| \begin{array}{l} \frac{\partial \xi}{\partial w_{221}} = \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial v_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial w_{221}} \\ = e_2(-1) \phi'_2(v_{22}) \cdot y_{11}(k) \end{array} \right. \\ &= e_1(k) \cdot (-1) \phi'_2(v_{21}) \cdot y_{11}(k) \end{aligned}$$

$$\begin{aligned} \frac{\partial \xi}{\partial w_{212}} &= \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial w_{212}} \quad \left| \begin{array}{l} \frac{\partial \xi}{\partial w_{222}} = \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial w_{222}} \\ = e_2(-1) \phi'_2(v_{22}) \cdot y_{12}(k) \end{array} \right. \\ &= e_1(-1) \phi'_2(v_{21}) \cdot y_{12}(k) \end{aligned}$$

$$\begin{aligned} \frac{\partial \xi}{\partial w_{213}} &= \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial w_{213}} \quad \left| \begin{array}{l} \frac{\partial \xi}{\partial w_{223}} = \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial w_{223}} \\ = e_2(-1) \phi'_2(v_{22}) \cdot y_{13} \end{array} \right. \\ &= e_1(-1) \phi'_2(v_{21}) \cdot y_{13}(k) \end{aligned}$$

Weight updation in output layer

$$\begin{pmatrix} \Delta w_{210}(k) \\ \Delta w_{211}(k) \\ \Delta w_{212}(k) \\ \Delta w_{213}(k) \end{pmatrix}^T = -\eta_2 \begin{pmatrix} \frac{\partial \xi}{\partial w_{210}} \\ \frac{\partial \xi}{\partial w_{211}} \\ \frac{\partial \xi}{\partial w_{212}} \\ \frac{\partial \xi}{\partial w_{213}} \end{pmatrix}^T = \eta_2 \underbrace{e_1(k) \phi'_2(v_{21}(k))}_{\delta_{21}} \underbrace{\begin{pmatrix} y_{10}(k) \\ y_{11}(k) \\ y_{12}(k) \\ y_{13}(k) \end{pmatrix}}_T, \quad \begin{array}{l} \downarrow \\ \delta_{21} = -\frac{\partial \xi}{\partial e_1} = \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial w_{210}} \\ = e_1(-1) \phi'_2(v_{21}) \end{array}$$

local gradient

$$\begin{pmatrix} \Delta w_{220}(k) \\ \Delta w_{221}(k) \\ \Delta w_{222}(k) \\ \Delta w_{223}(k) \end{pmatrix}^T = -\eta_2 \begin{pmatrix} \frac{\partial \xi}{\partial w_{220}} \\ \frac{\partial \xi}{\partial w_{221}} \\ \frac{\partial \xi}{\partial w_{222}} \\ \frac{\partial \xi}{\partial w_{223}} \end{pmatrix}^T = \eta_2 \underbrace{e_2(k) \phi'_2(v_{22}(k))}_{\delta_{22}} \underbrace{\begin{pmatrix} y_{10} \\ y_{11} \\ y_{12} \\ y_{13} \end{pmatrix}}_T, \quad \begin{array}{l} \downarrow \\ \delta_{22} = -\frac{\partial \xi}{\partial e_2} = \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial w_{220}} \\ = e_2(-1) \phi'_2(v_{22}) \end{array}$$

$$\begin{array}{l} = e_2(k) \phi'_2(v_{22}(k)) \end{array}$$

$$\Delta w_2(k) = \begin{pmatrix} \Delta w_{210} & \Delta w_{211} & \Delta w_{212} & \Delta w_{213} \\ \Delta w_{220} & \Delta w_{221} & \Delta w_{222} & \Delta w_{223} \end{pmatrix}$$

$$= \begin{pmatrix} \eta_2 \delta_{21}(k) y_{10}(k) & \eta_2 \delta_{21}(k) y_{11}(k) & \eta_2 \delta_{21}(k) y_{12}(k) & \eta_2 \delta_{21}(k) y_{13}(k) \\ \eta_2 \delta_{22}(k) y_{10}(k) & \eta_2 \delta_{22}(k) y_{11}(k) & \eta_2 \delta_{22}(k) y_{12}(k) & \eta_2 \delta_{22}(k) y_{13}(k) \end{pmatrix}$$

$$= \eta_2 \begin{pmatrix} \delta_{21}(k) \\ \delta_{22}(k) \end{pmatrix} (y_{10}(k) \ y_{11}(k) \ y_{12}(k) \ y_{13}(k))$$

$$\text{but } \begin{pmatrix} \delta_{21} \\ \delta_{22} \end{pmatrix} = \begin{pmatrix} e_1(k) \phi'_2(v_{21}) \\ e_2(k) \phi'_2(v_{22}) \end{pmatrix}$$

$$\Delta w_2(k) = \eta_2 \begin{pmatrix} \delta_{21}(k) \\ \delta_{22}(k) \end{pmatrix} y_1^T(k)$$

$$w_2(k+1) = w_2(k) + \Delta w_2(k)$$

$$= \begin{pmatrix} e_1(k) \\ e_2(k) \end{pmatrix} \odot \begin{pmatrix} \phi'_2(v_{21}) \\ \phi'_2(v_{22}) \end{pmatrix}$$

Hadamard product

Weight update in 1st layer

$$\begin{aligned} \frac{\partial E}{\partial w_{110}} &= \frac{\partial E}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial v_{11}} \cdot \frac{\partial v_{21}}{\partial w_{110}} + \frac{\partial E}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial v_{11}} \cdot \frac{\partial v_{22}}{\partial w_{110}} \\ &= e_1(-1) \phi'_2(v_{21}) w_{211} \phi'_1(v_{11}) x_0 + e_2(-1) \cdot \phi'_2(v_{22}) \cdot w_{221} \phi'_1(v_{11}) \cdot x_0 \\ &= -\underbrace{e_1(k) \phi'_2(v_{21})}_{\delta_{21}} w_{211} \phi'_1(v_{11}) x_0 - \underbrace{e_2(k) \phi'_2(v_{22})}_{\delta_{22}} w_{221} \phi'_1(v_{11}) x_0 \\ &= -\delta_{21} w_{211} \phi'_1(v_{11}) x_0 - \delta_{22} w_{221} \phi'_1(v_{11}) x_0 \end{aligned}$$

$$\begin{aligned} \Delta w_{110} &= -\eta_1 \frac{\partial E}{\partial w_{110}} = \eta_1 \delta_{21} w_{211} \phi'_1(v_{11}) x_0 + \eta_1 \delta_{22} w_{221} \phi'_1(v_{11}) \cdot x_0 \\ \Delta w_{111} &= -\eta_1 \frac{\partial E}{\partial w_{111}} = \eta_1 \delta_{21} w_{211} \phi'_1(v_{11}) x_1 + \eta_1 \delta_{22} w_{221} \phi'_1(v_{11}) \cdot x_1 \\ \Delta w_{112} &= -\eta_1 \frac{\partial E}{\partial w_{112}} = \eta_1 \delta_{21} w_{211} \phi'_1(v_{11}) x_2 + \eta_1 \delta_{22} w_{221} \phi'_1(v_{11}) \cdot x_2 \end{aligned}$$

$$\Delta w_1 = \begin{pmatrix} \Delta w_{110} & \Delta w_{111} & \Delta w_{112} \\ \Delta w_{120} & \Delta w_{121} & \Delta w_{122} \\ \Delta w_{130} & \Delta w_{131} & \Delta w_{132} \end{pmatrix}$$

$$\Delta W_{120} = -\eta_1 \frac{\partial \xi}{\partial W_{120}}$$

$$\begin{aligned} \frac{\partial \xi}{\partial W_{120}} &= \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_2} \cdot \frac{\partial v_{21}}{\partial v_{12}} \cdot \frac{\partial v_{12}}{\partial v_{10}} + \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial v_{12}} \cdot \frac{\partial v_{12}}{\partial v_{10}} \\ &= e_1 (-1) \phi'_2(v_{21}) w_{212} \phi'_1(v_{12}) x_0 + e_2 (-1) \phi'_2(v_{22}) w_{222} \phi'_1(v_{12}) x_0 \\ &= -e_1 \phi'_2(v_{21}) w_{212} \phi'_1(v_{12}) x_0 - e_2 \phi'_2(v_{22}) w_{222} \phi'_1(v_{12}) x_0 \end{aligned}$$

$$\Delta W_{120} = \eta_1 e_1 \phi'_2(v_{21}) w_{212} \phi'_1(v_{12}) x_0 + \eta_1 e_2 \phi'_2(v_{22}) w_{222} \phi'_1(v_{12}) x_0$$

$$\Delta W_{121} = \eta_1 e_1 \phi'_2(v_{21}) w_{212} \phi'_1(v_{12}) x_1 + \eta_1 e_2 \phi'_2(v_{22}) w_{222} \phi'_1(v_{12}) x_1$$

$$\Delta W_{122} = \eta_1 \underbrace{e_1 \phi'_2(v_{21})}_{\delta_{21}} w_{212} \phi'_1(v_{12}) x_2 + \eta_1 \underbrace{e_2 \phi'_2(v_{22})}_{\delta_{22}} w_{222} \phi'_1(v_{12}) x_2$$

$$\begin{aligned} \frac{\partial \xi}{\partial W_{130}} &= \frac{\partial \xi}{\partial e_1} \cdot \frac{\partial e_1}{\partial y_{21}} \cdot \frac{\partial y_{21}}{\partial v_{21}} \cdot \frac{\partial v_{21}}{\partial v_{13}} \cdot \frac{\partial v_{13}}{\partial v_{10}} + \frac{\partial \xi}{\partial e_2} \cdot \frac{\partial e_2}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial v_{22}} \cdot \frac{\partial v_{22}}{\partial v_{13}} \cdot \frac{\partial v_{13}}{\partial v_{10}} \\ &= e_1 (-1) \phi'_2(v_{21}) w_{213} \phi'_1(v_{13}) x_0 + e_2 (-1) \phi'_2(v_{22}) w_{223} \phi'_1(v_{13}) x_0 \\ &= -e_1 \phi'_2(v_{21}) w_{213} \phi'_1(v_{13}) x_0 - e_2 \phi'_2(v_{22}) w_{223} \phi'_1(v_{13}) x_0 \end{aligned}$$

$$\begin{aligned} \Delta W_{130} &= -\eta_1 \frac{\partial \xi}{\partial W_{130}} = \eta_1 \underbrace{e_1 \phi'_2(v_{21})}_{\delta_{21}} w_{213} \phi'_1(v_{13}) x_0 + \eta_1 \underbrace{e_2 \phi'_2(v_{22})}_{\delta_{22}} w_{223} \phi'_1(v_{13}) x_0 \\ &= \eta_1 \delta_{21} w_{213} \phi'_1(v_{13}) x_0 + \eta_1 \delta_{22} w_{223} \phi'_1(v_{13}) x_0 \end{aligned}$$

$$\begin{aligned} \Delta W_1 &= \begin{cases} \eta_1 \delta_{21} w_{211} \phi'_1(v_{11}) x_0 + \eta_1 \delta_{22} w_{221} \phi'_1(v_{11}) x_0 & ( ) x_1 \\ (\eta_1 \delta_{21} w_{212} \phi'_1(v_{12}) + \eta_1 \delta_{22} w_{222} \phi'_1(v_{12})) x_0 & ( ) x_1 \\ (\eta_1 \delta_{21} w_{213} \phi'_1(v_{13}) + \eta_1 \delta_{22} w_{223} \phi'_1(v_{13})) x_0 & ( ) x_2 \end{cases} \end{aligned}$$

$$\begin{aligned} &= \eta_1 \left[ \begin{pmatrix} \delta_{21} w_{211} + \delta_{22} w_{221} \\ \delta_{21} w_{212} + \delta_{22} w_{222} \\ \delta_{21} w_{213} + \delta_{22} w_{223} \end{pmatrix} \right] \underbrace{\begin{pmatrix} \phi'_1(v_{11}) \\ \phi'_1(v_{12}) \\ \phi'_1(v_{13}) \end{pmatrix}}_{3 \times 1} \underbrace{\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}}_{2 \times 2} \cdot \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \\ &\quad \text{Hadamard product } \Rightarrow \text{diag}(\phi'_1) \downarrow \text{matrix} \end{aligned}$$

$$= \eta_1 \left[ \left( \begin{pmatrix} w_{211} & w_{221} \\ w_{212} & w_{222} \\ w_{213} & w_{223} \end{pmatrix} \begin{pmatrix} \delta_{21} \\ \delta_{22} \end{pmatrix} \right) \right] \odot \begin{pmatrix} \phi'_1(v_{11}) \\ \phi'_1(v_{12}) \\ \phi'_1(v_{13}) \end{pmatrix} \underbrace{\begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}}_{3 \times 1} \underbrace{\text{diag}(\phi'_1)}_{3 \times 1} \underbrace{\begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}}_{3 \times 1}$$

$$\delta_1(k) = \underbrace{\bar{w}_1^T(k) \delta_2(k)}_{3 \times 2 \times 1} \odot \phi'_1(v_1(k)) \xrightarrow{\text{decay } (\phi'_1)} 3 \times 2$$

$$\Delta w_1 = \eta_1 \delta_1(k) y_0^T(k), \quad y_0 = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix}$$

$$w_1(k+1) = w_1(k) + \Delta w_1(k)$$

Summary of Back propagation Algorithm:

consider  $(L+1)$  layered FFNN with  $m_0$  input nodes and  $m_1, m_2, \dots, m_L$  neurons in the subsequent layers.

Let  $\phi_1, \phi_2, \dots, \phi_L$  be the activation functions of these layers where these functions are continuously differentiable.

Let  $\{(x(k), y_d(k))\}_{k=1}^N$  be a given set of training patterns.

Then for each  $k=1, \dots, N$ .

a) forward pass:

$$v_l(k) = w_l(k) y_{l-1}(k), \quad \bar{y}_l(k) = \phi_l(v_l(k)), \quad 1 \leq l \leq L$$

$$\text{where } y_0(k) = \begin{pmatrix} 1 \\ x(k) \end{pmatrix}, \quad y_l(k) = \begin{pmatrix} 1 \\ \bar{y}_l(k) \end{pmatrix}, \quad 1 \leq l \leq L-1$$

$$x(k) \in \mathbb{R}^{m_0}, \quad \bar{y}_l(k) \in \mathbb{R}^{m_l}, \quad w_l(k) \in \mathbb{R}^{m_l \times (m_{l-1} + 1)}, \quad 1 \leq l \leq L$$

$$e(k) = d(k) - \bar{y}_L(k)$$

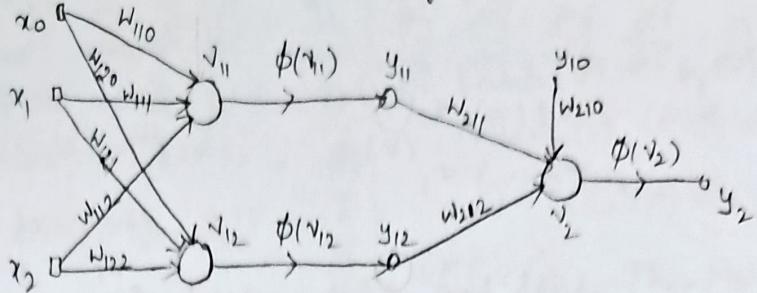
b) Backward Pass

$$\delta_L(k) = e(k) \odot \phi'_L(v_L(k))$$

$$\delta_l(k) = \left( \bar{w}_{l+1}^T(k) \delta_{l+1}(k) \right) \odot \phi'_l(v_l(k)), \quad l = L-1, L-2, \dots, 1$$

$$\Delta w_l(k) = \eta_l \delta_l(k) y_{l-1}^T(k), \quad 1 \leq l \leq L$$

$2:2:1 \rightarrow$  X-OR Gate wiring  $\phi$  - continuously  
back propagation differentiable



forward pass:

$$v_1 = \begin{matrix} w_1 \\ 2 \times 3 \end{matrix} y_0 = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

$$= \begin{pmatrix} w_{110} & w_{111} & w_{112} \\ w_{120} & w_{121} & w_{122} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

$$v_{11} = w_{110}x_0 + w_{111}x_1 + w_{112}x_2$$

$$v_{12} = w_{120}x_0 + w_{121}x_1 + w_{122}x_2$$

↓

$$\begin{pmatrix} v_{11} \\ v_{12} \end{pmatrix} = \begin{pmatrix} w_{110} & w_{111} & w_{112} \\ w_{120} & w_{121} & w_{122} \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

$$\bar{y}_1 = \phi_1(v_1) = \begin{pmatrix} y_{11} \\ y_{12} \end{pmatrix} = \begin{pmatrix} \phi(v_{11}) \\ \phi(v_{12}) \end{pmatrix}$$

$$y_0 = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix}, \quad \bar{y}_1 = \begin{pmatrix} 1 \\ \bar{y}_1 \end{pmatrix}$$

$$\bar{v}_2 = w_2 y_1 = \underbrace{\begin{pmatrix} w_{210} & w_{211} & w_{212} \end{pmatrix}}_{w_2} \begin{pmatrix} y_{10} \\ y_{11} \\ y_{12} \end{pmatrix} = w_2 \begin{pmatrix} 1 \\ \bar{y}_1 \end{pmatrix}; \quad y_2 = \begin{pmatrix} 1 \\ \bar{y}_2 \end{pmatrix}$$

$$y_2 = \phi_2(\bar{v}_2) \quad , \quad e = y_d - y_2$$

Backward pass:

$$\delta_2 = e \odot \phi'_2(\bar{v}_2(u))$$

$$\Delta w_2 = \bar{y}_2 \delta_2 y_1^T$$

$$w_2(k+1) = w_2(k) + \Delta w_2(k)$$

$$\delta_1 = (\bar{w}_2 \cdot \delta_2) \odot \phi'_1(v_1(u))$$

$$= \begin{pmatrix} w_{211} \\ w_{212} \end{pmatrix} \delta_2 \odot \phi'_1(v_{11}(u))$$

$$\Delta w_1(k) = \eta_1 \delta_1(k) y_0^T(k)$$

$$= \begin{pmatrix} \Delta w_{110} & \Delta w_{111} & \Delta w_{112} \\ \Delta w_{120} & \Delta w_{121} & \Delta w_{122} \end{pmatrix}$$

$$w_1(k+1) = w_1(k) + \Delta w_1(k)$$

$$\text{O/P layer} \quad \epsilon_1(k) = \frac{1}{2} e^2(k)$$

$$\frac{\partial \epsilon}{\partial w_{210}} = \frac{\partial \epsilon}{\partial e} \cdot \frac{\partial e}{\partial y_2} \cdot \frac{\partial y_2}{\partial v_2} \cdot \frac{\partial v_2}{\partial w_{210}}$$

$$= e(k) \cdot (-1) \phi'_2(\bar{v}_2) \cdot y_{10}$$

$$\frac{\partial \epsilon}{\partial w_{211}} = \frac{\partial \epsilon}{\partial e} \cdot \frac{\partial e}{\partial y_2} \cdot \frac{\partial y_2}{\partial v_2} \cdot \frac{\partial v_2}{\partial w_{211}}$$

$$= e(k) \cdot (-1) \phi'_2(\bar{v}_2) \cdot y_{11}$$

$$\frac{\partial \epsilon}{\partial w_{212}} = \frac{\partial \epsilon}{\partial e} \cdot \frac{\partial e}{\partial y_2} \cdot \frac{\partial y_2}{\partial v_2} \cdot \frac{\partial v_2}{\partial w_{212}}$$

$$= e \cdot (-1) \cdot \phi'_2(\bar{v}_2) \cdot y_{12}$$

$$\Delta w_2 = \begin{pmatrix} \Delta w_{210} \\ \Delta w_{211} \\ \Delta w_{212} \end{pmatrix} = -\eta_2 \begin{pmatrix} \frac{\partial \epsilon}{\partial w_{210}} \\ \frac{\partial \epsilon}{\partial w_{211}} \\ \frac{\partial \epsilon}{\partial w_{212}} \end{pmatrix}^T$$

$$= -\eta_2 e \cdot (-1) \cdot \phi'_2(\bar{v}_2) \begin{pmatrix} y_{10} \\ y_{11} \\ y_{12} \end{pmatrix}^T$$

$$= \eta_2 e \underbrace{\phi'_2(\bar{v}_2) y_1^T}_{\delta_2}$$

Weight updation 1<sup>st</sup> layer:

$$\frac{\partial E}{\partial w_{110}} = \frac{\partial E}{\partial e} \cdot \frac{\partial e}{\partial y_2} \cdot \frac{\partial y_2}{\partial v_2} \cdot \frac{\partial v_2}{\partial y_{11}} \cdot \frac{\partial y_{11}}{\partial v_{110}} \cdot \frac{\partial v_{110}}{\partial w_{110}}$$

$$= e(-1) \phi'_2(v_2) w_{211} \phi'_1(v_{11}) x_0$$

$$\frac{\partial E}{\partial w_{111}} = e(-1) \phi'_2(v_2) w_{211} \phi'_1(v_{11}) x_1$$

$$\frac{\partial E}{\partial w_{112}} = e(-1) \phi'_2(v_2) w_{211} \phi'_1(v_{11}) x_2$$

$$\frac{\partial E}{\partial w_{120}} = \frac{\partial E}{\partial e} \cdot \frac{\partial e}{\partial y_2} \cdot \frac{\partial y_2}{\partial v_2} \cdot \frac{\partial v_2}{\partial y_{12}} \cdot \frac{\partial y_{12}}{\partial v_{12}} \cdot \frac{\partial v_{12}}{\partial w_{120}}$$

$$= e(-1) \phi'_2(v_2) v_{212} \phi'_1(v_{12}) x_0$$

$$\frac{\partial E}{\partial w_{121}} = e(-1) \phi'_2(v_2) v_{212} \phi'_1(v_{12}) x_1$$

$$\frac{\partial E}{\partial w_{122}} = e(-1) \phi'_2(v_2) v_{212} \phi'_1(v_{12}) x_2$$

$$\Delta w_1 = \begin{pmatrix} \Delta w_{110} & \Delta w_{111} & \Delta w_{112} \\ \Delta w_{120} & \Delta w_{121} & \Delta w_{122} \end{pmatrix} = \begin{pmatrix} -\eta_1 \frac{\partial E}{\partial w_{110}} & -\eta_1 \frac{\partial E}{\partial w_{111}} & -\eta_1 \frac{\partial E}{\partial w_{112}} \\ -\eta_1 \frac{\partial E}{\partial w_{120}} & -\eta_1 \frac{\partial E}{\partial w_{121}} & -\eta_1 \frac{\partial E}{\partial w_{122}} \end{pmatrix}$$

$$= \begin{pmatrix} \underbrace{\eta_1 e \phi'_2(v_2) w_{211} \phi'_1(v_{11})}_{\delta_2} x_0 & \underbrace{\eta_1 e \phi'_2(v_2) w_{211} \phi'_1(v_{11})}_{\delta_2} x_1 & \underbrace{\eta_1 e \phi'_2(v_2) w_{211} \phi'_1(v_{11})}_{\delta_2} x_2 \\ \underbrace{\eta_1 e \phi'_2(v_2) w_{212} \phi'_1(v_{12})}_{\delta_2} x_0 & \underbrace{\eta_1 e \phi'_2(v_2) w_{212} \phi'_1(v_{12})}_{\delta_2} x_1 & \underbrace{\eta_1 e \phi'_2(v_2) w_{212} \phi'_1(v_{12})}_{\delta_2} x_2 \end{pmatrix}$$

$$= \eta_1 \begin{pmatrix} \delta_2 w_{211} \phi'_1(v_{11}) & \delta_2 w_{211} \phi'_1(v_{11}) & \delta_2 w_{211} \phi'_1(v_{11}) \\ \delta_2 w_{212} \phi'_1(v_{12}) & \delta_2 w_{212} \phi'_1(v_{12}) & \delta_2 w_{212} \phi'_1(v_{12}) \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

$$= \eta_1 \delta_2 \underbrace{\left[ \begin{pmatrix} w_{211} \\ w_{212} \end{pmatrix} \odot \begin{pmatrix} \phi'_1(v_{11}) \\ \phi'_1(v_{12}) \end{pmatrix} \right]}_{\delta_1} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}_{3 \times 1} \quad \odot \rightarrow \text{Hadamard product}$$

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \odot \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$

$$\delta_1 = \delta_2 \bar{w}_2^T \odot \phi'_1(v_1)$$

$$\Delta w_1 = \eta_1 \delta_1 y_0^T \quad 2 \times 1 \times 3$$

$$= \begin{pmatrix} a_{11}b_{11} & a_{12}b_{12} \\ a_{21}b_{21} & a_{22}b_{22} \end{pmatrix}$$

① Consider 1:4:1 N/N.

$$W_1 = \begin{pmatrix} 0.2 & -0.4 \\ 0.25 & -0.4 \\ 0.4 & -0.1 \\ -0.1 & 0.4 \end{pmatrix}$$

$$\eta_1 = \eta_2 = 0.01$$

$$\phi_1(v) = \tanh(v) \quad \phi_2(v) = v.$$

$$\phi_1'(v) = 1 - \tanh^2(v) \quad \phi_2'(v) = 1$$

$$N_2 = \begin{pmatrix} -0.25 & 0.25 & 0.15 & -0.5 & 0.3 \end{pmatrix}, (\text{sech}^2(v)) \text{ or } (y(1-y))$$

Data:  $\{(\gamma_2, 1)\}$  use BPA and find the next set of wbs.

$$v_1 = W^T x \quad \text{where } x = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0.5 \end{pmatrix}$$

$$= \begin{pmatrix} \overset{\longrightarrow}{0.2} & -0.4 \\ 0.25 & -0.4 \\ 0.4 & -0.1 \\ -0.1 & 0.4 \end{pmatrix} \downarrow \begin{pmatrix} 1 \\ 0.5 \end{pmatrix} = \begin{pmatrix} 0.2 + 0.2 \\ 0.25 - 0.2 \\ 0.4 - 0.05 \\ -0.1 + 0.2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0.05 \\ 0.35 \\ 0.1 \end{pmatrix}$$

$$\bar{y}_1 = \phi_1(v_1) = \begin{pmatrix} \tanh(0) \\ \tanh(0.05) \\ \tanh(0.35) \\ \tanh(0.1) \end{pmatrix} = \begin{pmatrix} 0 \\ 0.0499 \\ 0.3363 \\ 0.0996 \end{pmatrix}$$

$$y_1 = \begin{pmatrix} 1 \\ 0.0499 \\ 0.3363 \\ 0.0996 \end{pmatrix}$$

0.461

$$v_2 = W_2 y_1 = \begin{pmatrix} -0.25 & 0.25 & 0.15 & -0.5 & 0.3 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0.0499 \\ 0.3363 \\ 0.0996 \end{pmatrix}$$

$$= -0.380785$$

$y_2 = -0.380785 \rightarrow$  linear activation function

$$e = y_d - \underline{y_2} = 1.3808$$

$$f_2 = e_1 \odot \Phi_2^T (U_2(k)).$$

1:4:1

$$= e_0 1$$

$$= 1.3808$$

$$\Delta u_2 = \eta_2 f_2 y_1^T = 0.01 \times 1.3808 \times \begin{pmatrix} 1 & 0 & 0.05 & 0.25 & 0.1 \end{pmatrix}$$

$$\Delta u_2 = \begin{bmatrix} 0.0138 & 0 & 0.00069 & 0.0048 & 0.00128 \end{bmatrix}.$$

$$f_1 = (\bar{N}_2^T f_2) \odot \Phi_1^T (U_1(k)).$$

$$\begin{pmatrix} 0.25 \\ 0.15 \\ -0.5 \\ 0.3 \end{pmatrix} \times 1.3808 \odot \begin{pmatrix} 1 & 0 & 0.05 & 0.25 & 0.1 \end{pmatrix}$$

$$\begin{pmatrix} 0.345^2 \\ 0.2071 \\ -0.6904 \\ 0.4142 \end{pmatrix} \odot \begin{pmatrix} 1 & 0.9975 \\ 0.8869 & 0.9901 \end{pmatrix}$$

$$f_1 = \begin{pmatrix} 0.345^2 \\ 0.2066 \\ -0.6123 \\ 0.4101 \end{pmatrix}.$$

$$\Delta u_1 = \eta_1 f_1 y_1^T$$

$$= 0.01 \times \begin{pmatrix} 0.345^2 \\ 0.2066 \\ -0.6123 \\ 0.4101 \end{pmatrix} (1 \quad 0.5)$$

$$= \begin{pmatrix} 0.0034 & 0.0017 \\ 0.0020 & 0.0010 \\ -0.0061 & -0.0021 \\ 0.0041 & 0.0021 \end{pmatrix}.$$

#### UE17ECE\_PG\_PO10

#### UE17ECE\_PG\_PO09

#### UE17ECE\_PG\_PO11

Become a complete professional with high integrity and ethics, with excellent professional conduct and with empathy towards the environmental and contribute to the community for sustainable development of society.

Critically evaluate the outcomes of one's actions and apply self corrective measures to improve the performance.

Contribute and communicate effectively with the society confidently, be able to write effective reports and documents by adhering to the appropriate standards, make effective presentations, give and receive clear instructions

Engage in lifelong learning with persistent scientific temper for professional advancement and effective communication of the technical information.

$$(\tanh x)^{-1} = \frac{1}{\cosh^2 x}$$

$$= \operatorname{sech}^2 x$$

4. b) Heuristic For Making the BPA perform better.

1. Sequential v/s batch mode update:

For large & redundant data, sequential processing is fast.

For smaller data & NW, B-processing is faster.

②. Maximizing information content:

Aim: Every training example presented to the BPA should be chosen on the basis that its information content is the largest possible for the task at hand.

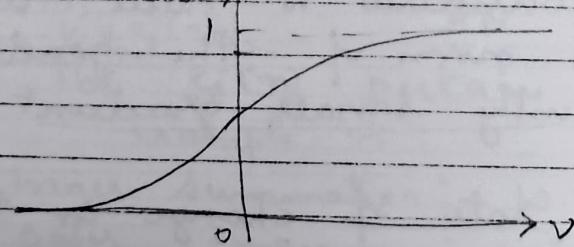
To achieve this aim ~~as~~ there are 2 ways.

①. The use of an example that results in the largest training error.

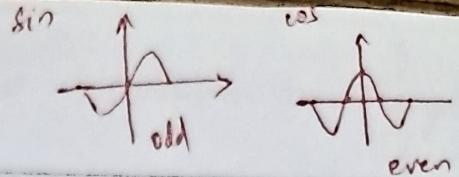
for updating w-value local gradient will be zero which is a function of error. When an example produces the larger error which leads to more larger local gradient, so faster will be the learning process.

②. The use of example that is radically different from all those previously used.

③ Activation function : antisymmetric functions are better than non-symmetric functions



→ non-symmetric



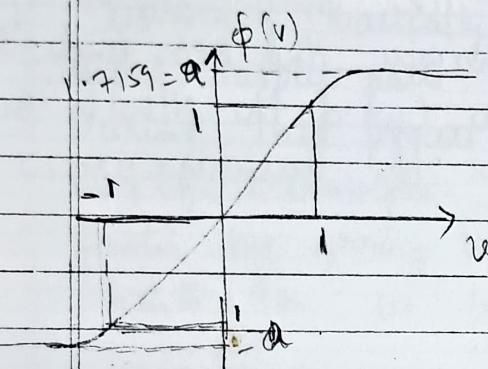
3. Activation function: choice of

As we discussed earlier activation function play an imp. role in training NN using RPA. If we use tangent hyperbolic Activation function learning will be faster than logistic. Because T.Hyperbolic is antisymmetric when as logistic non symmetric about 0.

$$\text{u } \phi(-v) = -\phi(v). \rightarrow \text{odd function (anti-symmetric)}$$

$$\phi(v) = a \tanh(bv)$$

(we can choose particular value for  $a = 1.7159$  and  $b = 2/3$ . Because the following useful opties.



$$\textcircled{1} \quad \phi(1) = 1$$

$$\phi(-1) = -1$$

\textcircled{2} At the origin the slope (i.e effective gain) of the activation fun is close to 1

$$\phi'(0) = ab. - ?$$

$$= 1.1424$$

\textcircled{3}  $\phi''(v) =$  attains its maximum value at  $v=1$ .

$|v| \leq 1$  : Region is linear

$|v| > 1$  : Region is referred as saturation

$$\frac{d(\tanh)}{dx} = \frac{u^x - v^x}{\sqrt{2}} \quad \tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}, \quad \frac{d \tanh(x)}{dx} = 1 - \tanh^2(x)$$

$$\frac{d \tanh(x)}{dx} = \frac{(e^x - e^{-x})(e^x + e^{-x}) - (e^x + e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

$$= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \tanh^2(x)$$

#### ④ Target Value:

It is important to choose the target value  $y_{dj}$  within the activation function range  $a \rightarrow 1$   
for example.

Let us assume the desired output neuron is  $y_{dj}$ . If it is off by some amount  $\epsilon$ ,

$$\text{ie } y_{dj} = a - \epsilon \quad (\text{For ex. } a = 0.7159) \\ \text{or } y_{dj} = -a + \epsilon, \quad \epsilon \rightarrow +\text{ve constant}$$

(like  $y_{dj}$  cannot exceed  $a = 0.7159$  for tanh)

since  $y_{dj}$  are derived from input examples

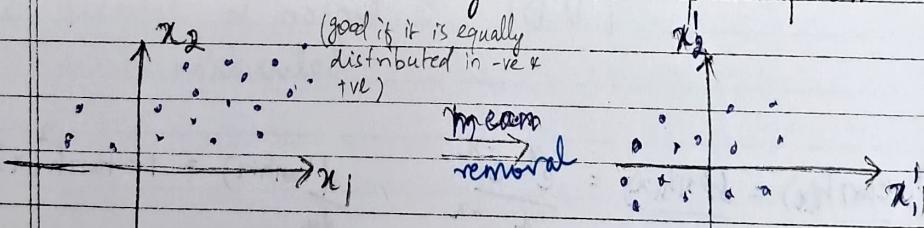
so that it will be in the linear region of activation function. Otherwise the BPA tends to drive the free parameters ( $w$ ) of the NN to infinity and thereby slow down the learning process by freezing.

#### ⑤ Normalization of inputs:

The training samples must be preprocessed before presenting to the NN.

The preprocessing steps are

① ~~mean~~ subtracting each input from mean



the data was not distributed

properly as a result after training process.

then the samples coming from inside

decorrelation (PCA)

Data are balanced

covariance

equalization (Scaling the variables)

NN not able to perform well  
covariances in  $x_1 \times x_2$  direction are equalized,

non zero mean  $\rightarrow$  large  $\Sigma$  values  $\rightarrow$  condition no is large  $\Rightarrow$  cost fun. surface will be steep

- (B) decorrelate the Training samples using PCA  
otherwise NN will be classify order

The decorrelated variables must be scaled so that their covariances are approximately equal so that

- (C) Covariances: it ensures the different system in the net learn at approximately the same speed.

- (D) Learning from hints:  
learning basically means mapping unknown input-output pairs.  
During this process it exploits the info present in the t-example. This process can be made faster by giving hints along with input examples. hints can be symmetric, invariance optis etc.

- (E) Learning rate: assign  $\eta$  in this order

$$\eta_1 > \eta_2 > \eta_3 \dots > \eta_L$$

( $\eta_1 = \eta_2 = \dots = \eta_L$ )

assign  $\eta = \alpha$

## Initialization:

A good choice for the initial values of the synaptic weights and thresholds of the N/W can be more useful in a successful N/W design. Then what is a good choice?

If I choose large<sup>initial</sup> values for  $w$ , it is highly likely that the neurons in N/W will be driven into saturation. If this happens the local gradient in BPA assume small values, which in turn will cause the learning process to slow down.

If  $w$ 's are assigned with small initial values, BPA may operate on a very flat area around the origin of the cost surface; where origin is a saddle point.

→ use of both large & small values of initial  $w$ , but it lies in between these.

To be specific, consider a task function as a f

Let the bias applied to each neuron in the network be set to zero &  $v_i = \sum_{j=1}^m w_{ij} y_j$  let ips applied to each neuron have  
 $\mu_y = E[y_j] = 0 \quad \forall i$   
 and  $\sigma_y^2 = E[(y_j - \mu_y)^2] = E[y_j^2] = 1 \rightarrow$  unit variance

i) zero mean  
ii) unit variance  
 assumed that  $E[y_i y_k] = \begin{cases} 1 & k=i \\ 0 & k \neq i \end{cases}$  iii) uncorrelated

and  $w$ 's are drawn from a uniformly distributed set of no. with '0' mean

$$\mu_w = E[w_{ij}] = 0 \quad \& \quad \sigma_w^2 = E[(w_{ij} - \mu_w)^2] = E[w_{ij}^2]$$

mean & Variance of  $v_i$  induced local field.

$$\mu_v = E[v_i] = E\left[\sum_{j=1}^m w_{ij} y_j\right] = \sum_{j=1}^m E[w_{ij}] E[y_j] = 0$$

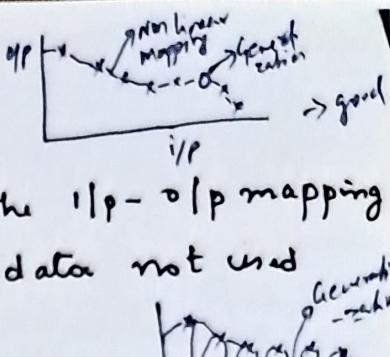
$$E[(v_i - \mu_v)^2] \leq \sigma_v^2 = E[v_i^2] = E\left[\sum_{j=1}^m w_{ij}^2 y_j^2\right] = \sum_{j=1}^m E[w_{ij}^2] = \sum_{i=1}^n E[w_{ii}^2]$$

$$\sigma_w^2 = \frac{1}{m} \sum_{j=1}^m w_{ij}^2 \quad \text{if } \sigma_v = 1$$

$m$  - no. of inputs

where  $m$  is the no. of synaptic connections per neuron

## Generalization & Cross Validation.

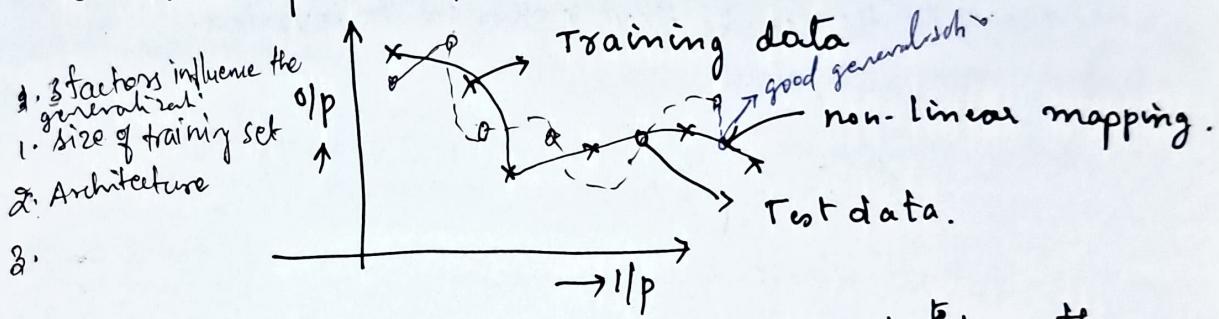


### 1) Generalization:

A N/W is said to generalize, when the I/p - o/p mapping computed by N/W is correct for test data not used in the training.

To understand the concept of generalization, we consider the N/W as sm. which performs non linear Input - output mapping.

The learning process is viewed as a curve-fitting problem. Then problem of mapping simply boils down to that of interpolation.



The training data is used to help the N/W learn the non linear input - output mapping represented by the curve.

When a test data is presented, the N/W computes the o/p as a result of interpolation performed by the N/W.

When a NN learns too many I/p - o/p samples, the N/W may end up memorizing the training set w/o trying to find the underlying function to be modelled - overfitting of our training.

When over trained the NN loses the ability to generalise.

This normally happens when more no. of H. Neurons than memory are used resulting in undesired contributions in the input space due to noise stored in it.

Generalization is influenced by.

- Size of training set and how close it is to the envt. of interest.
- Arch. of NN.
- Physical complexity of the pblm. - no control over either one of them can be fixed & other one can be varied.

Good generalization

$$\text{Size of the training set} \quad N = O\left(\frac{W}{\epsilon}\right) : O(\cdot) : \text{order of quantity.}$$

W: Total no. of free parameters (i.e. synaptic weights & biases)

$\epsilon$ : error permitted on the test data.

e.g.  $0.01 = 1\%$  of classification error.

~~Cross Validation:~~

→ generalisation is equivalent to the use of this multidimensional surface to interpolate test data. This is the motivation behind the method of RBF.

→ In the context of NN, the hidden units provide the set of functions that constitute an arbitrary basis for the input pattern when they are expanded into the hidden space. These functions are called as radial-basis functions.

Cross validation: available data set  $\xrightarrow{\text{training}}$  estimation - used to select the model  
 $\xrightarrow{\text{test}}$  validation - used to test or validate the model.  
To find multilayer perceptron with the best no. of hidden neurons & to figure out when it is best to stop training

4.15 Network Pruning Techniques  
→ Choosing a suitable topology for a NN, given 1<sup>st</sup> application is a difficult problem.

→ Usually, after a tedious trial & error process, an overfitted topology is chosen which is prone to various drawbacks like:

- ① high demand on computational resources
- ② increase in training time
- ③ non convergence of parameters
- ④ decrease in the generalization capability of a NN
- ⑤ price for HW implementation rises
- ⑥ less efficient.

If NN topology is too small, then it might not be able to learn the noise in the training data.  
To overcome this pblm, we need to find minimal topologies with best performance

There are 2 ways to achieve

- ① <sup>NN</sup> Growing approach: start with small multilayer perceptrons & then add a new neuron or a new layer of hidden neurons only when unable to reach desired specificite.
- ② <sup>NN</sup> Pruning Methods:

Generally latter method starts the training with a NN which is expected to be big enough to ensure a successful training

→ We start with a large multilayer perceptron with an adequate performance for the problem at hand, then prune it by weakening or eliminating certain synaptic weights in a selective and orderly fashion.

## Pruning Methods: deletion Regularization

### 1. deletion (Hessian based NW pruning)

→ basic idea of an analytic approach to network pruning is to use information on 2nd order derivative of error surface.

→ a local model of the error surface is constructed for analytically predicting the effect of perturbations in synaptic weights.

→ starting point in the construction of such a model is the local approximation of the cost function  $\xi_{\text{av}}$  by using a Taylor series about the operating point

$$\xi_{\text{av}}(w + \Delta w) = \xi_{\text{av}}(w) + g^T(w) \Delta w + \frac{1}{2} \Delta w^T H \Delta w + H \cdot O \cdot T \quad \rightarrow ①$$

where,  $\Delta w$  - perturbation applied to operating point  $w$

$g(w)$  - gradient vector evaluated at  $w$ .

→ To identify the set of parameters to be deleted which causes the least increase in the value of the cost function  $\xi_{\text{av}}$ .

→ To do this, we assume following approximations:

1. We assume that Extreme approximation:  $n$  parameters are deleted only after the training process has converged. In such cases  $g^T \Delta w = 0$

2. We assume that the error surface around a local minima is highly quadratic, hence  $H \cdot O \cdot T$  is neglected.

$$\therefore ① \Rightarrow \xi_{\text{av}}(w + \Delta w) - \xi_{\text{av}}(w) = \frac{1}{2} \Delta w^T H \Delta w \quad ②$$

eq. ② is the basis for the pruning procedure called optimal brain surgeon (OBS)

The goal of OBS is to set one of the synaptic weights to zero in order to minimize the incremental p.e. in fav.

Let  $w_{i(n)}$  denote this particular synaptic weight.

The elimination of this weight is equivalent to,

$$1_i^T \Delta W + w_i = 0 \quad , \quad 1_i^T \rightarrow \text{unit vector - all zeros except } i^{\text{th}} \text{ element}$$

Now the goal is

to minimize the quadratic form  $\frac{1}{2} \Delta W^T H \Delta W$  w.r.t. to incremental change in the weight vector  $\Delta W$ , subject to the constraint that  $1_i^T \Delta W + w_i = 0$  then minimize the result w.r.t. to the index  $i$ .

To solve this constrained - optimization problem, we first construct Lagrangian

$$S = \frac{1}{2} \Delta W^T H \Delta W - \lambda (1_i^T \Delta W + w_i) \quad \text{--- (3)}$$

$\lambda$  - Lagrange multiplier  $\frac{\partial S}{\partial \lambda} = b$

$$\frac{\partial S}{\partial \Delta W} = \frac{1}{2} \cdot 2 H \Delta W - \lambda 1_i^T = 0 \Rightarrow \Delta W = H^{-1} \lambda 1_i^T$$

$$\boxed{\Delta W = \lambda H^{-1} 1_i^T}$$

$$\frac{\partial S}{\partial \lambda} = -(1_i^T \Delta W + w_i)$$

$$\frac{\partial S}{\partial \lambda} = - (1_i^T H^{-1} 1_i^T \lambda + w_i) = 0$$

$$1_i^T H^{-1} 1_i^T \lambda = -w_i$$

$$\lambda = \underbrace{(1_i^T H^{-1} 1_i^T)^{-1}}_{\text{scalar}} w_i$$

$$= \frac{-w_i}{1_i^T H^{-1} 1_i^T} \quad \text{or} \quad \frac{-w_i}{(H^{-1})_{ii}}$$

$\hookrightarrow$   $i^{\text{th}}$  element

of this inverse

$$\Rightarrow \Delta W = \frac{-w_i}{1_i^T H^{-1} 1_i^T} \cdot H^{-1} 1_i^T = \frac{-w_i}{(H^{-1})_{ii}} \cdot H^{-1} 1_i^T$$

$\lambda_i \rightarrow 1$  can be removed  
 sub  $\lambda_i \propto \Delta w$  in Eq ③ to obtain  $s_{\min}$ .  
 $\lambda_i = \lambda$  in the following derivations.

$$\begin{aligned}
 &= \frac{1}{2} (\lambda_i H^{-1} z_i)^T H(w) (\lambda_i H^{-1} z_i) - \lambda_i (z_i^T \lambda_i H^{-1} z_i + w_i) \\
 &= \frac{1}{2} \lambda_i^2 z_i^T H^{-1} H^{-1} z_i - \lambda_i (z_i^T z_i + \lambda_i^2 z_i^T H^{-1} z_i + w_i) \quad (\because H^T = H) \\
 &= \frac{1}{2} \lambda_i^2 z_i^T H^{-1} z_i - \lambda_i (\lambda_i (H^{-1})_{ii} + w_i) \\
 &= \frac{1}{2} \lambda_i^2 (H^{-1})_{ii} - \lambda_i^2 (H^{-1})_{ii} - \lambda_i w_i \\
 &= -\frac{1}{2} \lambda_i^2 (H^{-1})_{ii} - \lambda_i w_i \\
 &= -\frac{1}{2} \lambda_i^2 (H^{-1})_{ii} - \lambda_i (-\lambda_i (H^{-1})_{ii}) \\
 &= -\frac{1}{2} \lambda_i^2 (H^{-1})_{ii} + \lambda_i^2 (H^{-1})_{ii}
 \end{aligned}$$

$$\begin{aligned}
 s_{\min} &= \frac{1}{2} \lambda_i^2 (H^{-1})_{ii} \\
 &= \gamma_2 \left( \frac{-w_i}{(H^{-1})_{ii}} \right)^2 (H^{-1})_{ii}
 \end{aligned}$$

$$s_{\min} = \gamma_2 \frac{w_i^2}{(H^{-1})_{ii}}$$

Saliency

$$s_i = \frac{1}{2} \frac{w_i^2}{(H^{-1})_{ii}}$$

where  $H^*$  - Hessian matrix.

Lagrangian  $s$ , optimized w.r.t.  $\Delta w$  subject to the constraint that the  $i$ th synaptic weight  $w_i$  be eliminated is called saliency of  $w_i$ .

$\rightarrow s_i$  represents the rise in the MSE or performance measure. In OBS,  $w_i$  corresponds to smallest  $s_i$  is selected for deletion.

small  $w_i \rightarrow$  small effect on MSE, but  $s_i$  inversely  $\propto$   $w_i$  due to  $H^{-1}$ ,  $\therefore$  small  $H^{-1}$  may have rise in MSE,  $\therefore$  this is a trade off.

# Supervised learning viewed as an optimization problem

- ① Optimization problem means finding the best solution.
- ② Error is highly nonlinear function of  $w$ .  
 → Using Taylor series we expand  $E_{av}(w)$  at the current point on the error surface  $w(n)$  as  

$$E_{av}(w(n) + \Delta w(n)) = E_{av}(w(n)) + \frac{\partial E_{av}}{\partial w(n)} \frac{\Delta w(n)}{1!}$$

... continued.

$$+ \frac{\Delta w}{2} \frac{\partial^2 E_{av}}{\partial w(n)^2} \cdot \frac{\Delta w(n)}{2!} + \text{higher order terms}$$

→ sub.  $\lambda$  &  $\Delta w$  in

$$S = \frac{1}{2} (-\lambda u_i H^{-1})^T H (-\lambda u_i H^{-1}) + \frac{w_i}{u_i^T H u_i} (u_i^T - \lambda u_i H^{-1} + w_i)$$

$$s_i = \frac{w_i}{2[H^{-1}]_{ii}}$$

$H^{-1}$  - Hessian matrix

small  $w_i \rightarrow$  small effect on MSE  
but  $s_i$  inversely  $\propto$   $H^{-1}$ ;  $\therefore$  small  $H^{-1}$  may have  $\propto$  in MSE,  $\therefore$  this is a trade off.

belong to row  
prior to previous row  
prob → Lagrangian  $S$ , optimized  $w_i$  to  $\Delta w$  subject to the constraint that

the  $i$ th synaptic weight  $w_i$  be eliminated, is called the saliency of  $w_i$ .

→  $s_i$  represents the use in the MSE or performance measure. In OBS,  $w_i$  corresponds to smallest  $s_i$  is selected for deletion.

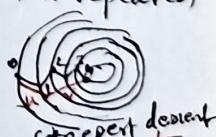
- 1st order optimization Method: rate of convergence , by adding

→ Draw back in SGD is slow convergence , by adding  $\alpha$ -momentum might help , but adding  $\alpha$  more terms will make the training process little complicated. ( $\because \Delta w$  depends on  $\Delta w$ )

→ This is bcs SGD takes many steps to reach minimum. The method repeatedly moves in the steepest direction , consequently the repeated steps are far to each other , i.e in zig-zag manner.  
 $\therefore$  it is slower

→ To avoid this , we go for another tech. CGM , in

which we consider conjugate vector  $v$  which will not take any steps , when 'u' has taken , by which it neglects multiple steps.

  
 = steepest descent  
 = conjugate gradient

## 4.16 Conjugate Gradient Method

→ it is a second-order optimization method.

→ let us assume minimization of quadratic function

$$f(x) = \frac{1}{2} x^T A x - b^T x + c.$$

Where  $x$  is a  $n \times 1$  parameter vector

$A$  :  $n \times n$  symmetric positive definite matrix

$b$  :  $n \times 1$  vector

$c$  : constant

minimization of quadratic fun.  $f(x)$  is achieved by

$$f'(x) = 0.$$

$$\Rightarrow A x^* - b = 0$$

$$x^* = A^{-1} b.$$

Thus minimizing of  $f(x)$  and solving linear sys  
of eqns  $Ax^* = b$  are equivalent problems.

Definition:  
Note for a given matrix  $A$ , we say that a set of  
non zero vectors  $[s_0 \ s_1 \ \dots \ s_{(n-1)}]$  is  $A$ -conjugate.  
If the following condition is satisfied.  
 $s_i^T (A s_j) = 0 \quad \forall i, j \text{ s.t } i \neq j.$

If  $A$  = Identity matrix, then it implies notion of orthogonality

For a given set of A-conjugate vectors  $s(0), s(1), \dots, s(w-1)$ , the 3 corresponding conjugate direction method for unconstrained minimization if the quadratic error function  $f(x)$  is defined by

$$x(n+1) = x(n) + \eta(n)s(n) \quad n = 0, 1, \dots, w-1.$$

where  $x(0)$  is an arbitrary starting vector

$\eta(n)$  : scalar defined by

$$f(x(n) + \eta(n)s(n)) = \min_{\eta} f(x(n) + \eta s(n)).$$

The procedure of choosing  $\eta$  so as to minimize the function  $f(x(n) + \eta s(n))$  for some fixed  $n$  is referred to as

line search which represents 1-dimensional minimization problem

$$f(x(n) + \eta(n)s(n)) = \frac{1}{2} [x(n) + \eta(n)s(n)]^T A [x(n) + \eta(n)s(n)] - b^T [x(n) + \eta(n)s(n)] + c.$$

$$f(x(n) + \eta(n)s(n)) = \frac{1}{2} [x(n)^T A x(n) + \eta(n) s(n)^T A x(n) + \eta(n) s(n)^T A s(n) + \eta^2(n) s(n)^T A s(n)] - b^T x(n) - \eta(n) b^T s(n) + c$$

$$\frac{\partial f(x)}{\partial \eta} = 0 + \frac{1}{2} s(n)^T A x(n) + \frac{1}{2} x(n)^T A s(n) + \eta(n) s(n)^T A s(n) - b^T s(n)$$

$$0 = s(n)^T A x(n) - b^T s(n) + \eta(n) s(n)^T A s(n)$$

$$\eta(n) s(n)^T A s(n) = -s(n)^T \{A x - b\} \quad A = A^T = A^{-1}$$

$$\eta(n) = \frac{-s(n)^T \{A x - b\}}{s(n)^T A s(n)}$$

$$x(n) - x^* = \epsilon$$

$$\eta(n) = \frac{-s(n)^T A \epsilon(n)}{s(n)^T A s(n)}$$

$$\text{where } x^* = A^{-1}b$$

starting from an arbitrary point  $x(0)$  the conjugate direction method is guaranteed to find optimal sol<sup>n</sup>  $x^*$  of  $f(x)=0$  in at most  $N$  iterations.

At successive iterations, the CG method minimizes the quadratic function  $f(x)$  over a progressively expanding linear vector space that eventually includes the global minimum of  $f(x)$ , i.e. for each iteration  $x(n+1)$  minimizes  $f(x)$  that passes through  $x(0)$ .

For the CG method to work we require the availability of set A-conjugate vectors  $s(0), s(1), \dots, s(N-1)$

It is determined in a sequential manner at current step.

→ Define residual  $r(m) = b - Ax(m)$  as the steepest descent direction  
 linear combination  
 $\wedge L.C. \text{ of } r(m) \times s(m-1) \Rightarrow s(m) = r(m) + \beta(m)s(m-1)$

scaling factor to be determined  
 where  $\beta(m)$  is chosen s.t., A-conjugacy & persists.

$$\Rightarrow s^T(m-1) A s(m) = s^T(m-1) A r(m) + \beta(m) s^T(m-1) A s(m-1).$$

$$0 = s^T(m-1) A r(m) + \beta(m) s^T(m-1) A s(m-1)$$

$$\Rightarrow \beta(m) = -\frac{s^T(m-1) A r(m)}{s^T(m-1) A s(m-1)}.$$

option To avoid the dependency on  $A$  to generate s'vectors we consider

$$x^* = A^{-1}b$$

$$e(m) = x(m) - x^*$$

Polar-Ribiere formula

$$\beta(m) = \frac{r^T(m)[r(m) - r(m-1)]}{s^T(m-1) s(m-1)}$$

Fletcher-Reeves formula

$$\beta(m) = \frac{r^T(m) r(m)}{s^T(m-1) s(m-1)}$$

Given  
 $x(0), s(0)$

find  $r(0)$

$$① x(0) = b - Ax(0)$$

$$\beta(0) = -\frac{s^T(0) A e(0)}{s^T(0) A s(0)}$$

$$x(1) = x(0) + \eta(0) \beta(0)$$

$$② x(1) = b - Ax(1)$$

$$\beta(1) = -\frac{s^T(1) A e(1)}{s^T(1) A s(1)}$$

$$s(1) = r(1) + \beta(1) s(0)$$

$$\eta(1) = -s^T(1) A e(1)$$