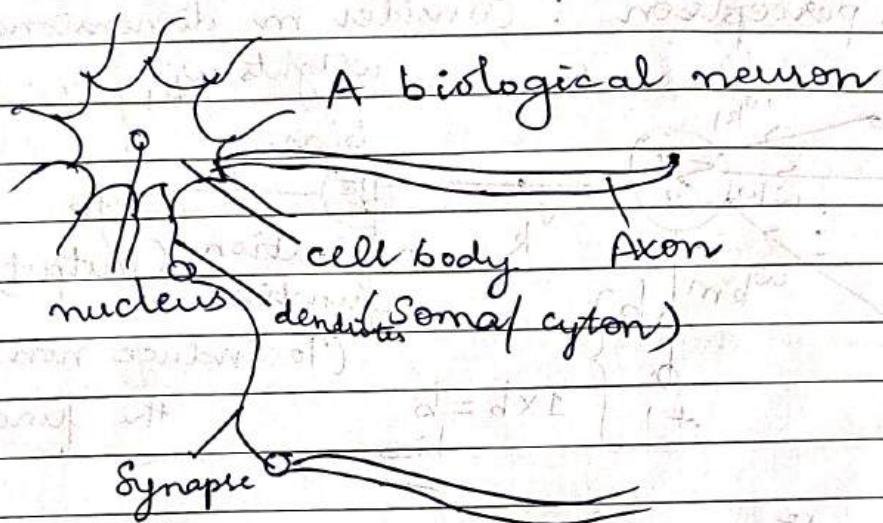


Deep Learning

Structure of neuron:

Axon long in size.

Dendron: for receiving information from synapse
cell body: for processing information (axon of another neuron)



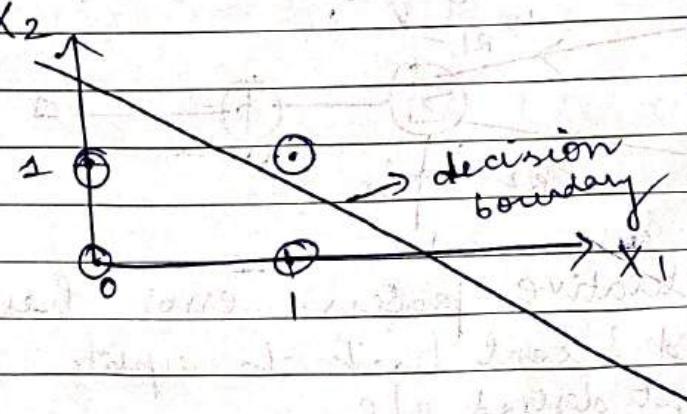
If a decision boundary can be realised to separate the distinct classes — linearly separable data.

Logic Gate AND, OR etc.

x_1	x_2	$x_1 x_2$
0	0	0
0	1	0
1	0	0
1	1	1

XOR gate
(non linearly
separable)

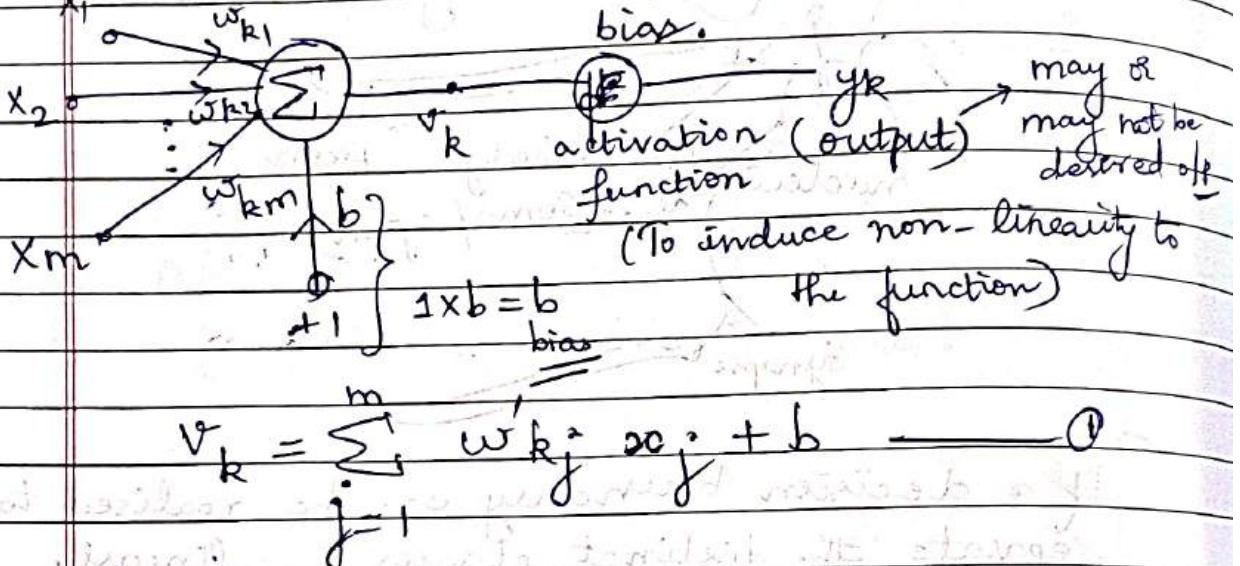
multiple decision
boundaries to
separate data.



A single perceptron is needed for solving linearly separable problems.

Multi-layer perceptron (consists of inputs & multiple hidden layers) to solve non-linearly separable problems.

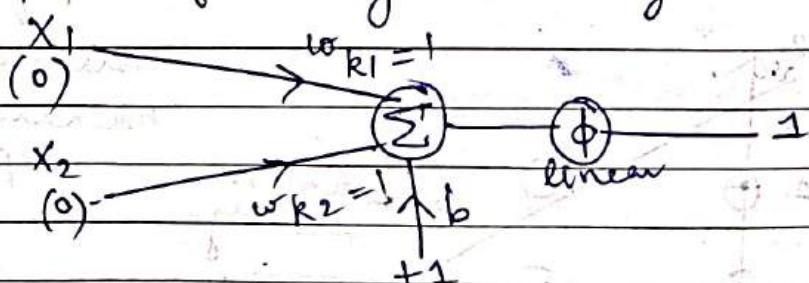
A perceptron : Consider m dimensional data, with weights w_{k1}, \dots, w_{km} with bias.



$$y_k = \phi\left(\sum_{j=1}^m w_{kj} x_j + b\right)$$

Consider 2 inputs x_1, x_2 with weights 1 each and a bias term 1 is induced.

for and gate (O/P by multiplication)



It is an iterative process. error has to be calculated & sent back to update weights to get desired O/P.

$$(0 \times 1) + 1 \rightarrow 1$$

~~+ (0 \times 1)~~

~~(0 \times 1) + 1 = 0~~

desired output = 0

so, error should be calculated.

Activation function:

1. Hardlim :-

$$\phi(v_k)$$

+1

hardlim

$$\phi(v_k)$$

$$y = \text{hardlim}(v_k)$$

$$\phi(v_k) = \begin{cases} 1 & \text{if } v_k \geq 0 \\ 0 & \text{if } v_k < 0 \end{cases}$$

$$(1) x_1$$

$$w_2 = -2$$

$$(1) x_2$$

$$w_1 = 1$$

$$\text{hardlim } y_k = 0 \text{ or } 1$$

$$v_k = -2(1) + 1(1) + 0 \text{ bias}$$

$$v_k = -1$$

$$0 \text{ or } 1$$

2. Symmetric hard limit:

$$\phi(v_k)$$

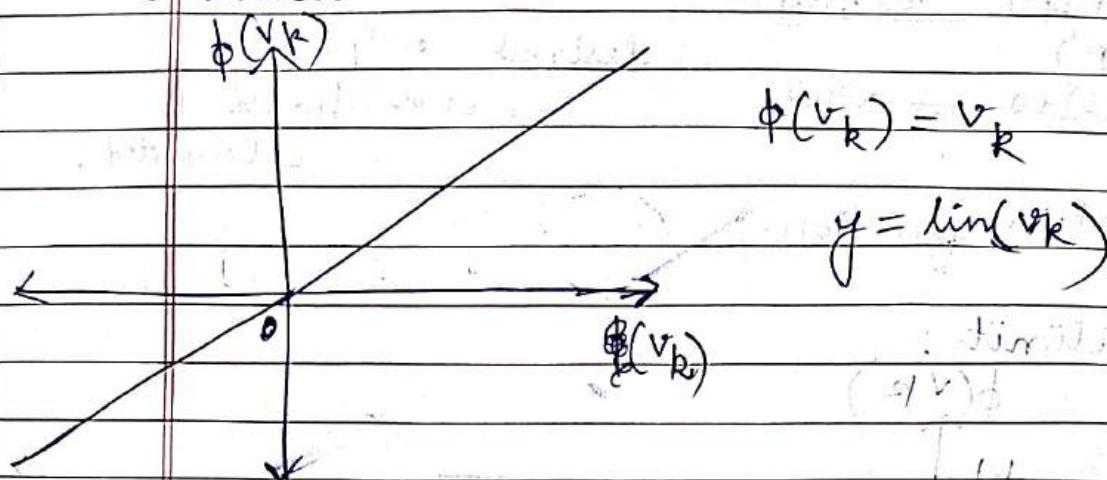
+1

-1

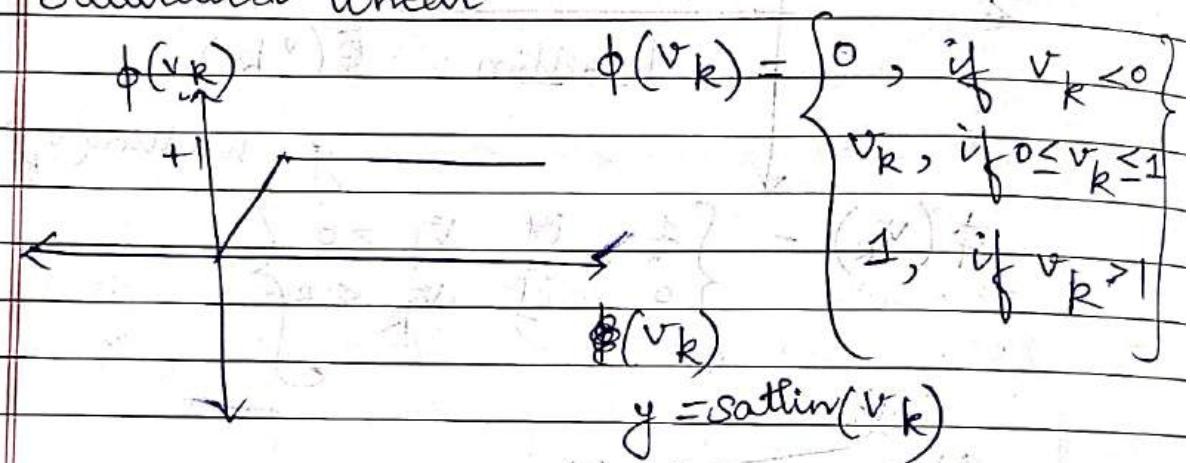
$$\phi(v_k) = \begin{cases} +1 & \text{if } v_k \geq 0 \\ -1 & \text{if } v_k < 0 \end{cases}$$

$$y = \text{hardlim}_{\text{symmetric}}(v_k)$$

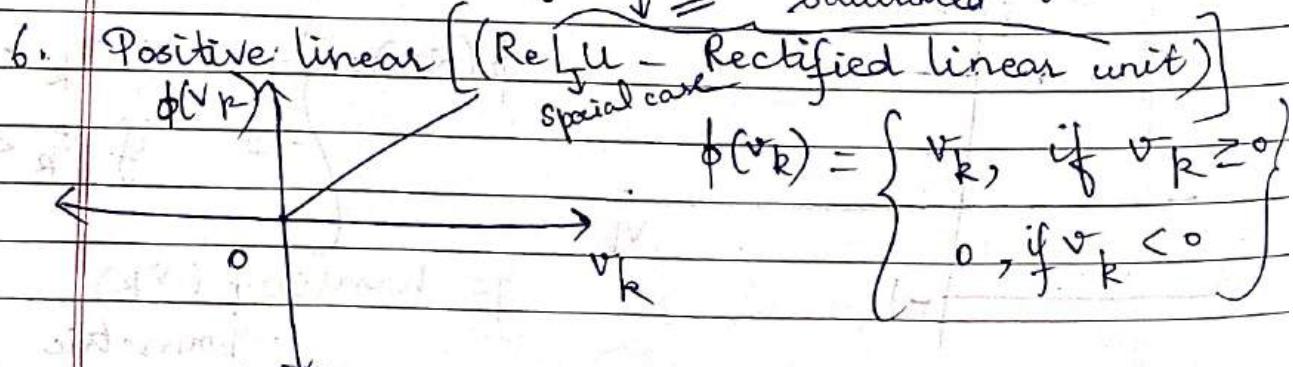
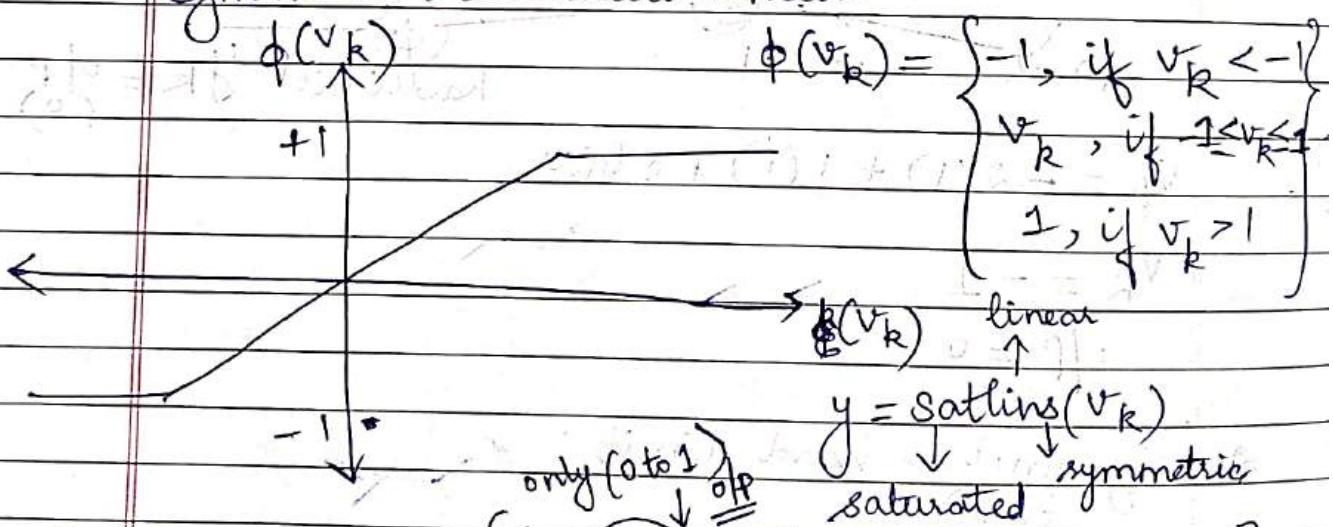
3. Linear



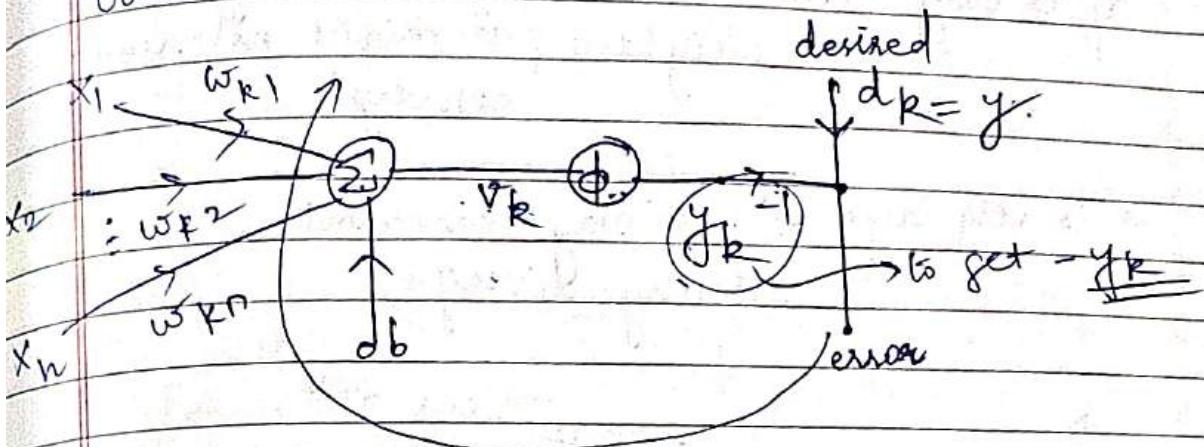
4. Saturated linear



5. Symmetric saturated linear



Gradient - descent :

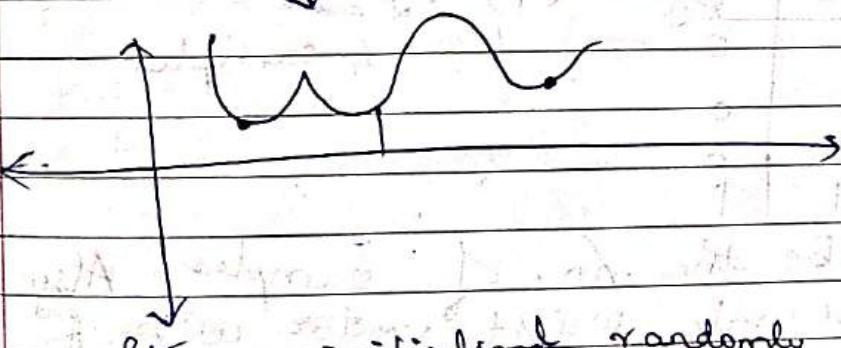
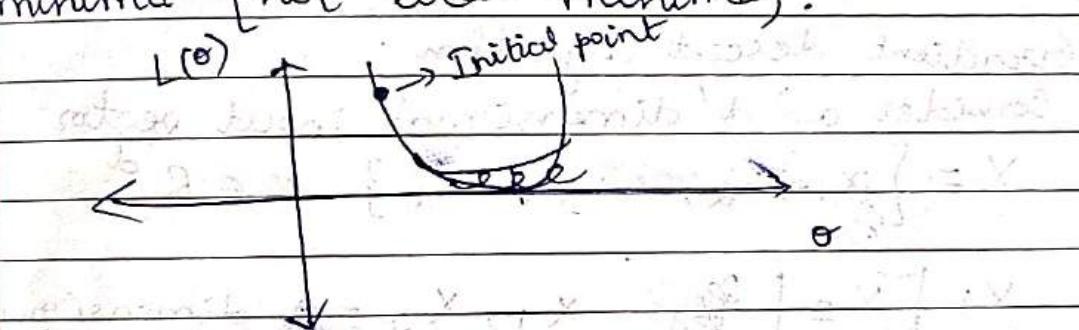


$$v_k = \sum_{i=1}^m w_i^T x + b$$

$$y_k = f(v_k) \quad \text{error} = y - y_k$$

Gradient descent algorithm :-

- Start from any point on curve with initial weights and tend to reach the only global minima (not local minima).



our objective is to minimize error ..

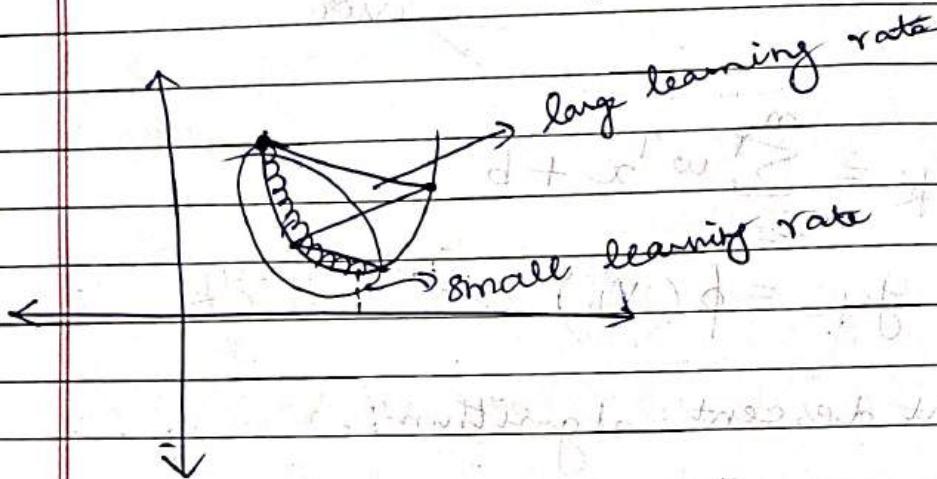
$$w_{n+1} = w_n - \eta \frac{\partial L}{\partial w}$$

learning rate / step size

how it effects the performance
(how it converges)

η is very small - no drastic change in weights (it might not reach convergence)

η is very large - weights may overshoot (it may diverge)



Initialisation of weights and learning rate affect the model performance.

Gradient descent algorithm:

Consider a 'd' dimensional input vector

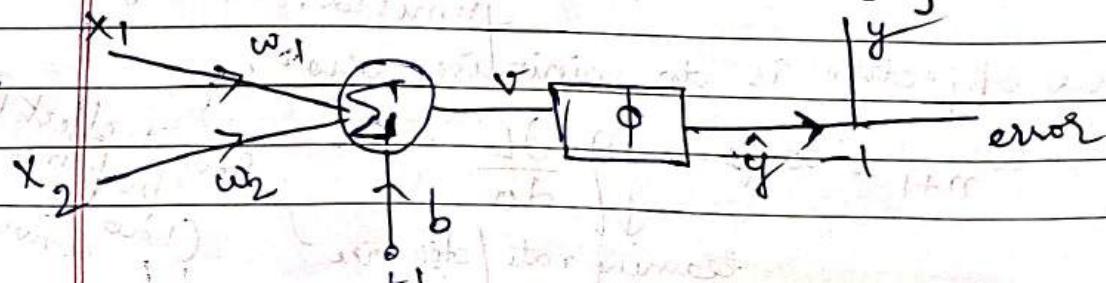
$$X = \{x_1, x_2, \dots, x_d\} \quad X \in \mathbb{R}^d$$

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

$x_1, x_2 = 2$ dimensions
4 samples

Let m be the no. of samples. Also, a d dimensional weight vector $w \in \mathbb{R}^d$

$$w = \{w_1, w_2, \dots, w_d\}$$



$y \rightarrow$ desired output

$\hat{y} \rightarrow$ obtained/predicted output

$$\text{Error } (e) = y - \hat{y}$$

$$\text{Loss } (\omega, b) = \underbrace{\frac{1}{2} \sum_i (y_i - \hat{y}_i)^2}_{\substack{\text{weight bias} \\ \text{We need to reduce error w.r.t. weight vector \& bias}}} \quad \begin{array}{l} \text{[obtained by partial derivative]} \\ \text{to reduce computational intensity} \end{array}$$

(5)

We need to reach the lowest possible point from initial guess, to do this we follow the minimization which is done by calculation of gradient.

$$\omega_{n+1} = \omega_n - \eta \frac{\partial L(\omega)}{\partial \omega}$$

$$\hat{y}_i = \phi(\omega x_i + b) \quad (\text{considering } i^{\text{th}} \text{ sample})$$

In this example $L = \text{(one sample)}$

$$1 \leq i \leq 4. \quad \boxed{\phi(v_k) = v_k}, \text{ linear activation fn}$$

Let ϕ be linear activation function.

So, that the derivative = 1.

$$\hat{y}_i = \omega x_i + b$$

$$\omega_{n+1} = \omega_n - \eta \frac{\partial L(\omega)}{\partial \omega} \quad (2) \quad b_{n+1} = b_n - \eta \frac{\partial L(b)}{\partial b} \quad (3)$$

$$\frac{\partial L(w)}{\partial w} = \frac{\partial L(w)}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w} \quad [\text{chain rule}]$$

Digression

$$\begin{aligned}\text{error} &= y - \hat{y} \\ &= 3 - 2 \\ &= 1\end{aligned}$$

$$\begin{aligned}\text{error} &= y - \hat{y} \\ &= 3 - 4 \\ &= -1\end{aligned}$$

Net error becomes 0

To avoid this ambiguity, we use the mean squared error.

$$\text{MSE} = \frac{1^2 + (-1)^2}{2} = \frac{1+1}{2} = 1 \quad (\text{net error})$$

That is why loss is mean squared error.

$$\text{Differentiate } L(w, b) = \frac{1}{2} \sum (y_i - \hat{y}_i)^2 \text{ w.r.t. } \hat{y}_i$$

$$\begin{aligned}\frac{\partial L(w)}{\partial \hat{y}_i} &= \frac{1}{2} \times 2 (y_i - \hat{y}_i) \cdot \frac{\partial}{\partial \hat{y}_i} (-\hat{y}_i) \\ &= -(y_i - \hat{y}_i) \quad \text{--- (6)}\end{aligned}$$

$$\text{Differentiate } \hat{y}_i = w x_i + b \text{ w.r.t. } w$$

$$\frac{\partial \hat{y}_i}{\partial w} = x_i \quad \text{--- (7)}$$

$$\frac{\partial L(w)}{\partial w} = -(y_i - \hat{y}_i)x_i \quad (\text{from } ④)$$

Substitute above equation in ②

$$w_{n+1} = w_n - \eta(- (y_i - \hat{y}_i))x_i$$

$$\therefore w_{n+1} = w_n + \eta(y_i - \hat{y}_i)x_i \quad ⑧$$

$$w_{n+1} = w_n + \eta \text{ error } x_i$$

$$\frac{\partial L(b)}{\partial b} = \frac{\partial L(b)}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial b} \quad ⑨$$

$$\therefore L(w, b) = \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

$$\begin{aligned} \frac{\partial L}{\partial b} &= \frac{1}{2} \times 2 (y_i - \hat{y}_i) \frac{\partial}{\partial \hat{y}_i} (-\hat{y}_i) \\ &= -(y_i - \hat{y}_i) \end{aligned} \quad ⑩$$

Differentiate $\hat{y}_i = wx_i + b$ w.r.t. b

$$\frac{\partial \hat{y}_i}{\partial b} = 1 \quad ⑪$$

from ⑨, ⑩, ⑪

$$\frac{\partial L(b)}{\partial b} = -(y_i - \hat{y}_i) \quad (\text{from } ④)$$

$$b_{n+1} = b_n - \eta \frac{\partial L(b)}{\partial b}$$

$$b_{n+1} = b_n + \eta(y_i - \hat{y}_i) \quad ⑫$$

For an AND gate, calculate the error and append overall error, plot decision boundary along with I/P (use scatter plot) for 100 iteration epochs.

I/P array

O/P array $x = np.array([$

error = ()

for no of epochs learning rate $\eta \leq 0.1$

for length (I/P)

$$y = wx + b$$

$$y = \text{hardlim}(y)$$

error

update weights

append (total error)

Decision-boundary:

$$wx + b = 0$$

$$w_1 x_1 + w_2 x_2 + b = 0$$

plot from -0.5 to 0.5

$$x_2 = -b - w_1 x_1$$

AND

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

Solve an AND problem (assuming initial weights)

$$[w_1, w_2] = [1, 2]^T, \text{ initial bias} = 1$$

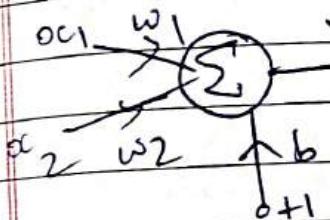
$$w_{n+1} = w_n + \gamma (y - \hat{y}_i) x_i$$

$$\text{and } b_{n+1} = b_n + \gamma (y_i - \hat{y}_i) \quad \gamma = 1$$

using gradient descent algorithm

1st iteration :-

$$(i) \quad \partial c_1 = 0, \quad \partial c_2 = 0, \quad y = 0$$



$$v = w_1 x_1 + w_2 x_2 + b$$

$$= 1(0) + 2(0) + 1$$

$$= 1$$

$$\hat{y} = \text{hardlim}(1)$$

$$= 1$$

error = desired - obtained

$$= 0 - 1$$

$$= -1$$

$$w_{n+1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1(-1) \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$w_{n+1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$b_{n+1} = 1 + 1(-1) = 0$$

$$w_{n+1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad b_{n+1} = 0$$

$$ii) \quad x_1 = 0$$

$$x_2 = 1$$

$$y = 0$$

$$v = w_1 x_1 + w_2 x_2 + b$$

$$= 1(0) + 2(1) + 0$$

$$= 2$$

$$\hat{y} = \text{hardlim}(2)$$

$$= 1$$

$$\text{error} = 0 - 1 = -1$$

$$w_{n+1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1(-1) \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$b_{n+1} = 0 + 1(-1) = -1$$

$$w_{\text{new}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b_{\text{new}} = -1$$

iii) $x_1 = 1, x_2 = 0$
 $\hat{y} = 0$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(1) + 1(0) - 1 \\ &= 0 \end{aligned}$$

$$\hat{y} = \text{hardlim}(0) = 1$$

$$\text{error} = -1$$

$$\begin{aligned} w_{n+1} &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1(-1) \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \end{aligned}$$

$$b_{n+1} = -1 + (-1)(-1) = -2$$

$$w_{\text{new}} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad b_{\text{new}} = -2$$

iv) $x_1 = 1, x_2 = 1$
 $\hat{y} = 1$

$$v = w_1 x_1 + w_2 x_2 + b$$

$$\begin{aligned} &= 0(1) + 1(1) - 2 \\ &= -1 \end{aligned}$$

$$\hat{y} = \text{hardlim}(-1) = 0$$

$$\text{error} = 1 - 0 = 1$$

$$w_{n+1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 1(1) \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$b_{n+1} = -2 + 1(1) = -1$$

$$w_{\text{new}} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, b_{\text{new}} = -1$$

2nd iteration:-

i) $x_1 = 0, x_2 = 0, \hat{y} = 0$

$$v = w_1 x_1 + w_2 x_2 + b = 1(0) + 2(0) - 1 = -1$$

$$\hat{y} = \text{hardlim}(-1) = 0$$

$$\text{error} = 0$$

No need to update weights & biases.

ii) $x_1 = 0, x_2 = 1, \hat{y} = 0$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(0) + 2(1) + -1 \\ &= 1 \end{aligned}$$

$$\hat{y} = \text{hardlim}(1) = 1$$

$$\text{error} = 0 - 1 = -1$$

$$w_{n+1} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1(-1) \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$b_{n+1} = -1 + (1)(-1) = -2$$

$$w_{\text{new}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b_{\text{new}} = -2$$

iii) $x_1 = 1, x_2 = 0, \hat{y} = 0$

$$v = w_1 x_1 + w_2 x_2 + b = 1(1) + 0(0) - 2 = -1$$

No need to update weights & biases. $\hat{y} = \text{hardlim}(-1) = 0$ error = 0 - 0 = 0

iv) $\alpha_1 = 1, \alpha_2 = 1, \hat{y} = 1$

$$v = w_1 \alpha_1 + w_2 \alpha_2 + b$$

$$= 1(1) + 1(1) - 2$$

$$= 0$$

$$\hat{y} = \text{hardlim}(0) = 1$$

$$\text{error} = 0 - 1 = 0$$

No need to update weights & biases.

however, we didn't get for step 2nd combinations. So, need to reiterate.

3rd iteration:

(i) $\alpha_1 = 0, \alpha_2 = 0, \hat{y} = 0$

$$v = w_1 \alpha_1 + w_2 \alpha_2 + b$$

$$= 0(1) + 0(1) - 2$$

$$= -2$$

$$\hat{y} = \text{hardlim}(-2) = 0$$

$$\text{error} = 0 - 0 = 0$$

No need to update weights & biases.

(ii) $\alpha_1 = 0, \alpha_2 = 1, \hat{y} = 0$

$$v = w_1 \alpha_1 + w_2 \alpha_2 + b$$

$$= 0(1) + 1(1) - 2$$

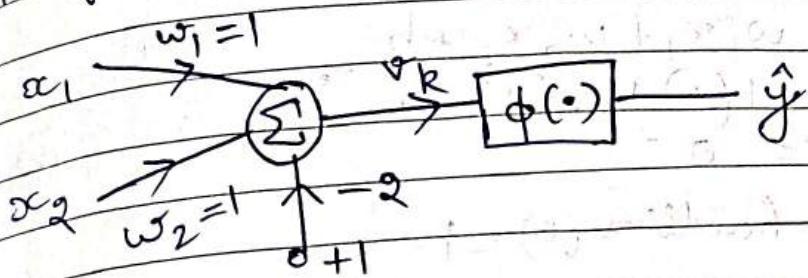
$$= -1$$

$$\hat{y} = \text{hardlim}(-1) = 0$$

$$\text{error} = 0$$

No need to update weights & biases.

Therefore,



$$\hat{y} = \phi(x_1 + x_2 - 2)$$

OB problem

$$\begin{array}{ccc} x_1 & x_2 & \hat{y} \\ 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{array} \quad [w_1, w_2] = [1, 1]^T$$

bias = 0

1st iteration:

$$(i) x_1 = 0, x_2 = 0, \hat{y} = 0$$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(0) + 1(0) + 0 \\ &= 0 \end{aligned}$$

$$\hat{y} = \text{hardlim}(0) = 0$$

$$\text{error} = 0 - 1 = -1$$

~~No need to update weights or bias.~~

~~(ii) $x_1 = 0, x_2 = 1, \hat{y} = 0$~~

$$w_{n+1} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} + (-1) \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$b_{n+1} = 0 + (-1) = -1$$

(ii) $x_1 = 0, x_2 = 1, \hat{y} = 1$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(0) + 1(1) - 1 \\ &= 0 \end{aligned}$$

$$\hat{y} = \text{hardlim}(0) = 1$$

error = 0

(iii) $x_1 = 1, x_2 = 0, \hat{y} = 1$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(1) + 1(0) - 1 \\ &= 0 \end{aligned}$$

$$\hat{y} = \text{hardlim}(0) = 1$$

error = 0

(iv) $x_1 = 1, x_2 = 1, \hat{y} = 1$

$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(1) + 1(1) - 1 \\ &= 1 \end{aligned}$$

$$\hat{y} = \text{hardlim}(1) = 1$$

error = 0

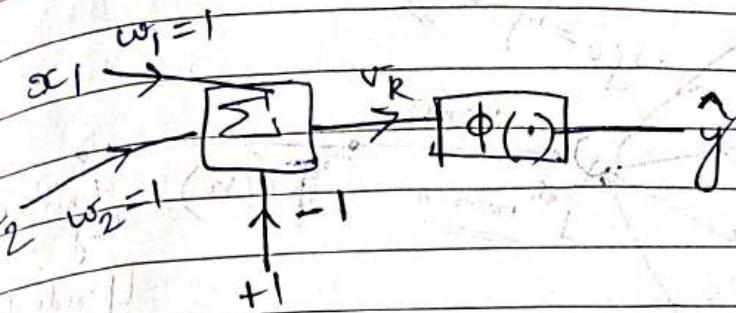
2nd iteration

(i) $x_1 = 0, x_2 = 0, \hat{y} = 0$

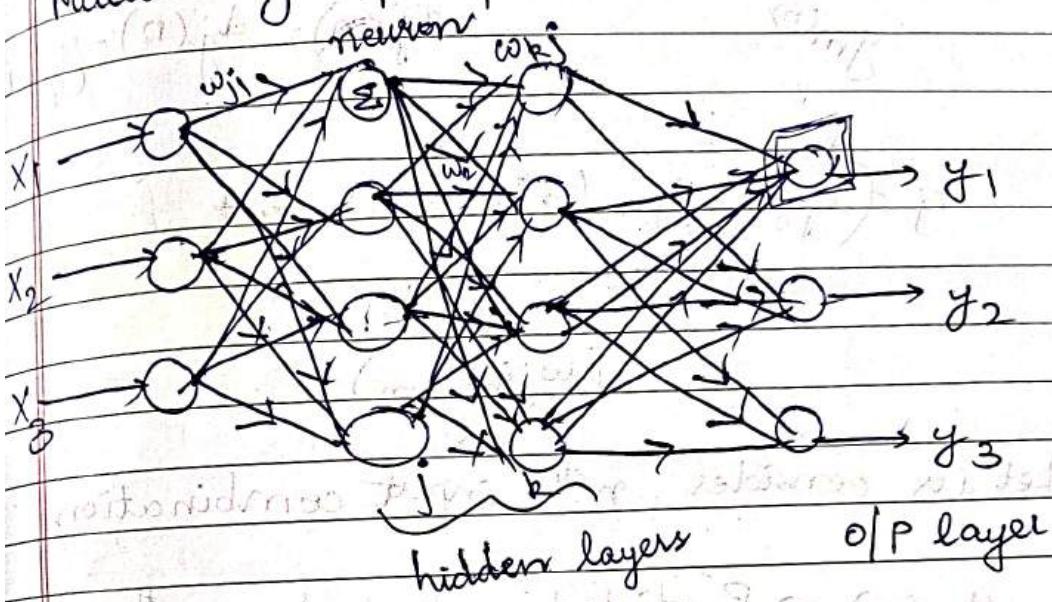
$$\begin{aligned} v &= w_1 x_1 + w_2 x_2 + b \\ &= 1(0) + 1(0) - 1 \\ &= -1 \end{aligned}$$

$$\hat{y} = \text{hardlim}(-1) = 0$$

error = 0



Multi-layer perception:



- Required for non linearly separable problems.

Ex:- XOR, XNOR gate

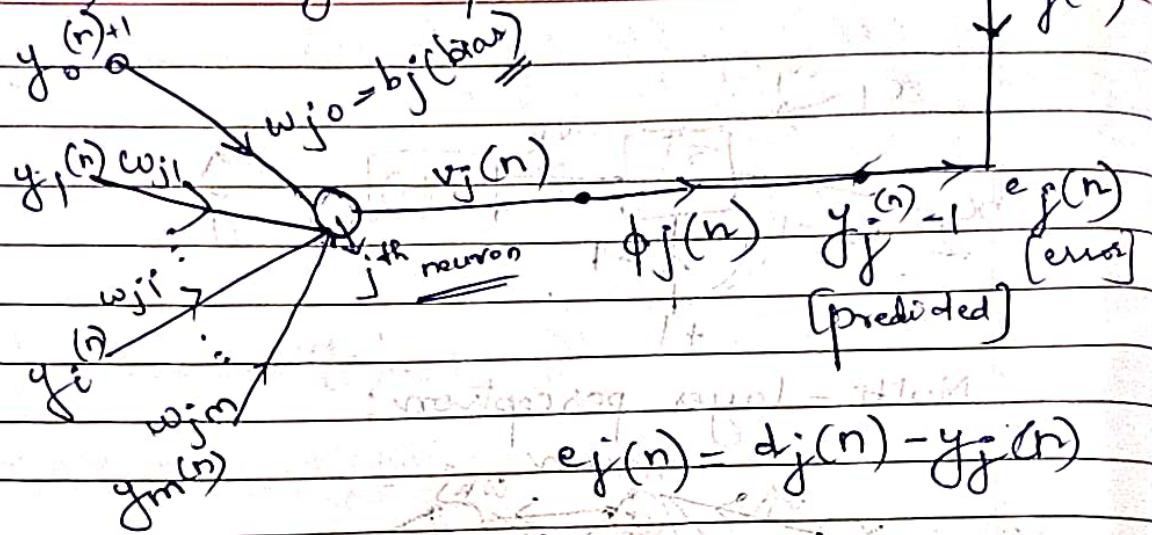
Back propagation:

Weights get updated for hidden layers and input layers by propagating the gradient of loss function using the chain rule.

Cases:

- error update for O/P neuron (If error occurs in O/P layer, it is similar to single layer perception)
- error update for hidden neuron (hidden layer receives updated weights from different layers, error from O/P layer is also received by hidden layer)

Case i) when j is o/p node -



$$v_j = (w_{j0} \times y_0) + (w_{j1} \times y_1) + \dots + (w_{j2} \times y_2) + \dots + (w_{jm} \times y_m)$$

Let us consider n^{th} input combination

$y_j^{(n)}$ → Predicted output for n^{th} input.

$d_j^{(n)}$ → Desired output

$$\text{error} = e_j^{(n)} = d_j^{(n)} - y_j^{(n)} \quad \text{--- (4)}$$

$$\text{Loss} = Z(n) = \frac{1}{q} \sum_{j \in C} e_j^{(n)^2} \quad \text{--- (1)}$$

(for all classes)

$$= \frac{1}{q} \sum_{j \in C} (d_j^{(n)} - y_j^{(n)})^2 \quad \text{--- (2)}$$

Weight updation =

$$w_{ji}^{n+1} = w_{ji}^{n} - \eta \frac{\partial Z(n)}{\partial w_{ji}(n)} \quad \text{change in weights}$$

$$w_{ji}^{n+1} = w_{ji}^{n} + \Delta w_{ji}$$

η = learning rate

$\frac{\partial Z(n)}{\partial w_{ji}(n)}$ = Gradient of weight vector

for n^{th} iteration, error found for o_i, o_j we update weights for all vectors \rightarrow [see updated w_1, w_2 , etc]

The entire vector is updated, so w_{ji}

$$\Delta w_{ji} = (-\eta) \frac{\partial Z(n)}{\partial w_{ji}(n)} \quad (3)$$

\downarrow $w_{ji}(n) \in [j^0, j^1, \dots, j^L]$

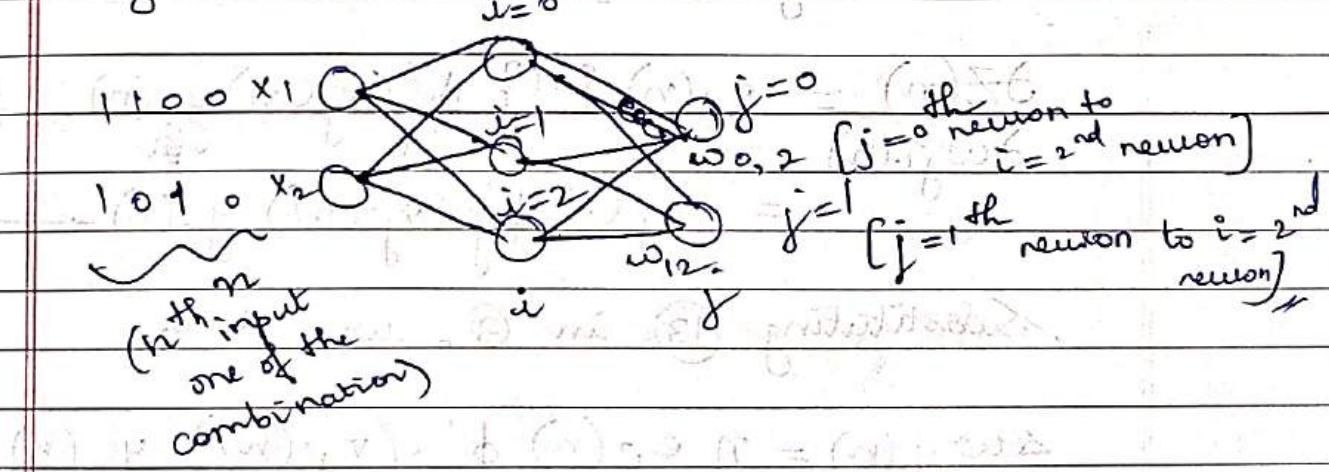
$$v_j(n) = \sum_{i=0}^m w_{ji}(n) y_i(n) \quad (5)$$

\downarrow j^0, j^1, \dots, j^L

Indicates from which layer of neuron correction is established

$$y_j(n) = \phi_j(v_j(n)) \quad (6)$$

Digression:



Continuation:

find $\frac{\partial Z(n)}{\partial w_{ji}(n)}$

$\frac{\partial Z(n)}{\partial w_{ji}(n)}$

$$\frac{\partial Z(n)}{\partial w_{ji}(n)} = \frac{\partial Z(n)}{\partial e_j(n)} \cdot \frac{\partial e_j(n)}{\partial y_j(n)} \cdot \frac{\partial y_j(n)}{\partial v_j(n)} \cdot \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

[chain rule]

From ④, $\frac{\partial Z(n)}{\partial e_j(n)} = e_j(n)$ — ⑧

From ④, $\frac{\partial e_j(n)}{\partial y_j(n)} = -1$ — ⑨

$$\frac{\partial y_j(n)}{\partial v_j(n)} \quad [e_j(n) = d_j(n) - g_j(n)]$$

From ⑥, $y_j(n) = \phi_j(v_j(n))$

$$\frac{\partial y_j(n)}{\partial v_j(n)} = \phi'_j(v_j(n))$$
 — ⑩

From ⑤, $\frac{\partial v_j(n)}{\partial w_{ji}(n)} = y_i(n)$ — ⑪

Substituting ⑧, ⑨, ⑩, ⑪, ⑫ in ⑦ we get

$$\frac{\partial Z(n)}{\partial w_{ji}(n)} = e_j(n) \{ \dots \} \cdot \phi'_j(v_j(n)) \cdot y_i(n)$$
 — ⑫

$$= -e_j(n) \phi'_j(v_j(n)) y_i(n)$$
 — ⑬

Substituting ⑬ in ③, we obtain

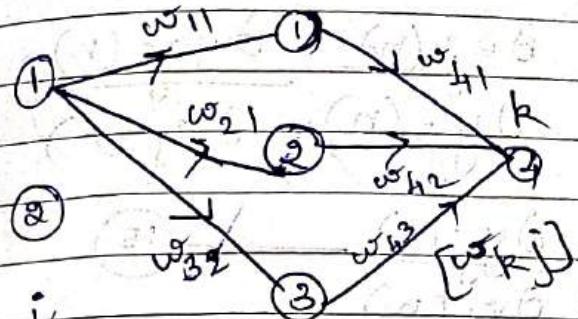
$$\Delta w_{ji}(n) = \eta e_j(n) \phi'_j(v_j(n)) y_i(n)$$

Assume local gradient 's'

$$s_j = e_j(n) \phi'_j(v_j(n)), j \text{ is output node.}$$
 — ⑭

Therefore, $\Delta w_{ji}(n) = \eta s_j y_i(n)$

Digressions :-



$n \rightarrow$ combination of

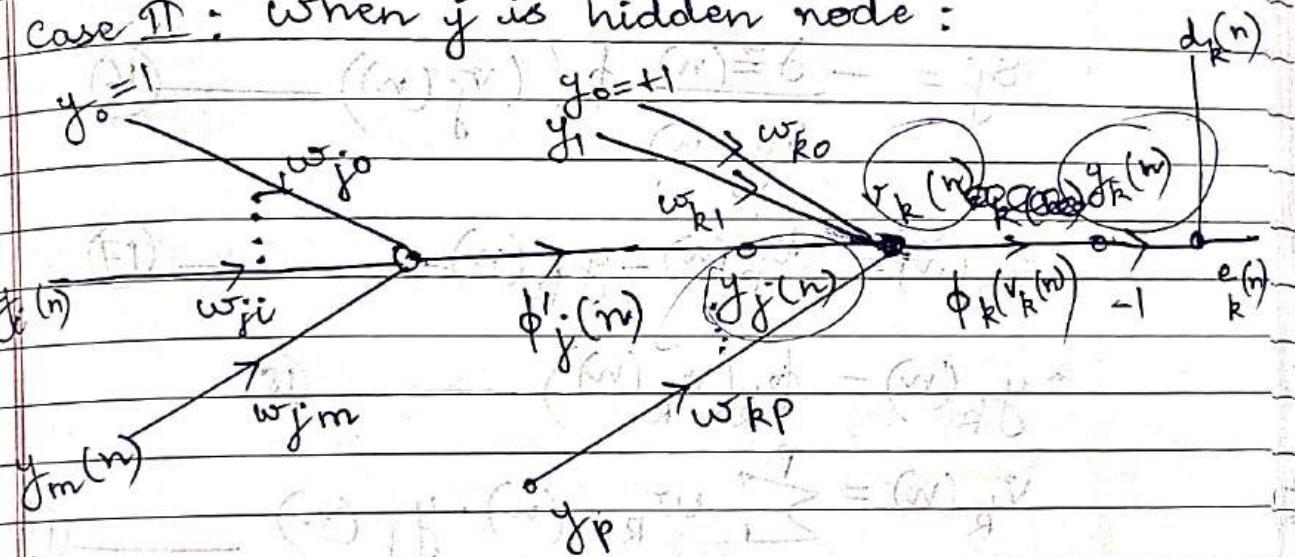
inputs
(not weight vector)

$n=1 \rightarrow \infty$

combination

$$[w_{ji}] \quad j \quad \text{etc.}$$

case II : when j is hidden node :



From (14), for output node we had -
[previously]

$$\delta_j = e_j(n) \phi'_j(v_j(n))$$

{Trying to do a backward pass.}

$$\delta_j = -e_j(n)(-1) \phi'_j(v_j(n))$$

From ⑧

$$\delta_j = -\frac{\partial z(n)}{\partial e_j(n)} (-1) \phi'_j(v_j(n))$$

From ⑨

$$\delta_j = -\frac{\partial z(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial v_j(n)} \phi'_j(v_j(n))$$

from ⑩

$$\delta_j = -\frac{\partial z(n)}{\partial y_j(n)} \cdot \frac{\partial e_j(n)}{\partial y_j(n)} \cdot \frac{\partial y_j(n)}{\partial v_j(n)}$$

$$\delta_j = -\frac{\partial z(n)}{\partial y_j(n)} \cdot \frac{\partial e_j(n)}{\partial y_j(n)} \quad ⑯$$

but $\frac{\partial y_j(n)}{\partial v_j(n)} = \phi'_j(v_j(n))$

$$\delta_j = -\frac{\partial z(n)}{\partial y_j(n)} \cdot \phi'_j(v_j(n)) \quad ⑯$$

$$e_k(n) = d_k(n) - y_k(n) \quad ⑯$$

$$y_k(n) = \phi_k(v_k(n)) \quad ⑯$$

$$v_k(n) = \sum_{j=0}^p w_{kj}(n) \cdot y_j(n) \quad ⑯$$

$$z(n) = \sum_{k \in C} c_k(n) \quad ⑯$$

(all classes)
or combination

$$\frac{\partial z(n)}{\partial y_j(n)} = \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial y_j(n)} \quad ⑯$$

(from ⑯)
output to input

$$\frac{\partial z(n)}{\partial y_j(n)} = \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial y_k(n)} \frac{\partial y_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)} \quad ⑯$$

From ⑯

$$\frac{\partial z(n)}{\partial y_i(n)} = -\sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial y_k(n)} \frac{\partial v_k(n)}{\partial v_i(n)} \frac{\partial y_i(n)}{\partial y_i(n)} \quad ⑯$$

from (18)

$$\frac{\partial Z(n)}{\partial y_j(n)} = - \sum_k e_k(n) \phi'_k(v_k(n)) \frac{\partial v_k(n)}{\partial y_j(n)} \quad (24)$$

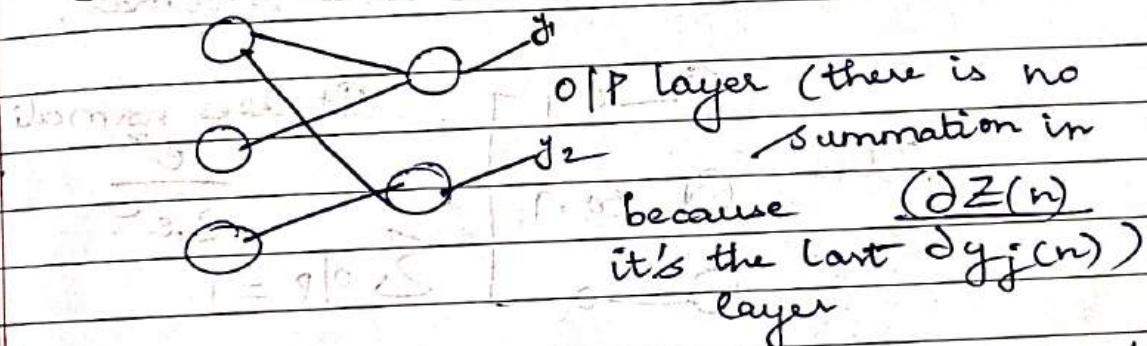
from (19)

$$\frac{\partial Z(n)}{\partial y_j(n)} = - \sum_k e_k(n) \phi'_k(v_k(n)) w_{kj}(n) \quad (25)$$

$$\therefore \frac{\partial Z(n)}{\partial y_j(n)} = - \sum_k \delta_k w_{kj}(n) \quad (26)$$

where $\delta_k = e_k(n) \phi'_k(v_k(n))$

Digression:



However, in hidden layer, there is summation in (26) because the neuron receives error correction from different layers.

Continuation:-

Substituting (26) in (19), we get -

$$\delta_j = \phi'_j(v_j(n)) \sum_k \delta_k w_{kj}(n) \quad (27)$$

where j is hidden node.

Note:

$$\Delta w_{ji} = \eta \delta_j y_i$$

$$\delta_j = \phi'_j(v_j(n)) e_j \quad j \text{ is O/P node.}$$

$$\delta_j = \phi'_j(v_j(n)) \sum_k S_k w_{kj} \quad j \text{ is hidden node.}$$

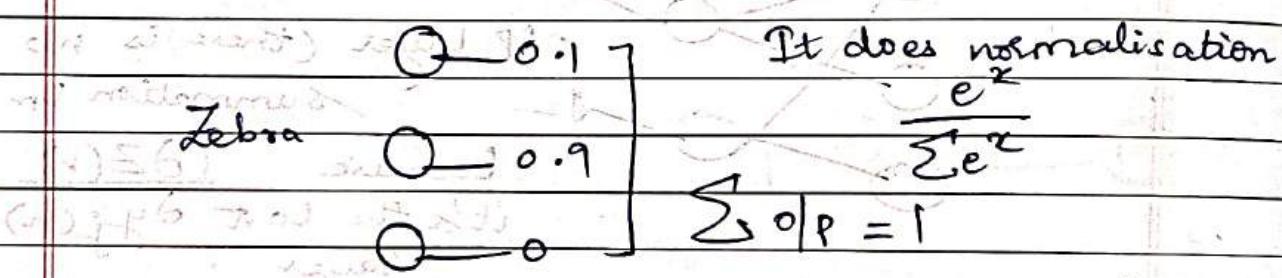
If there is one more hidden layer,
do substitution again (series of steps!)

Sigmoid - binary classification

Softmax - multiclass classification

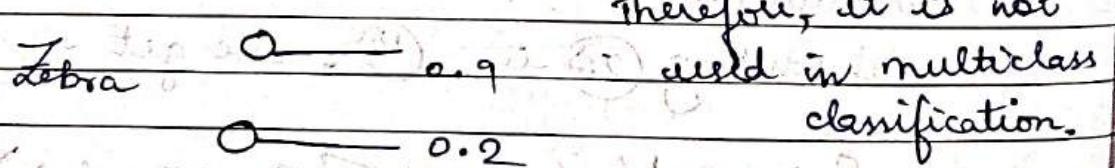
[Probability oriented]

Suppose there are three classes cat, dog, zebra & we used softmax.



If sigmoid is used in O/P all values

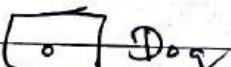
will lie from 0 to 1



Multi-label classification:-

In same image, there are n-number of

(multi-label) images.



≈ 1 cat

≈ 1 cat

Sigmoid activation
because

softmax gives probabilistic output.

How probabilities of softmax is calculated?

$$\begin{array}{l}
 \rightarrow 0.7 = 0.7 \\
 \rightarrow 0.2 = \frac{0.7 + 0.2 + 0.1 + 0.4}{14} = 7/14 \\
 \rightarrow 0.1 = 1/14 \\
 \rightarrow 0.4 = 4/14
 \end{array}$$

Activation functions:

Log-sigmoid: [Ranges from 0 to 1]

$$y_j = \phi(v_j(n)) = \frac{1}{1 + e^{-av_j(n)}} \quad *$$

$$\phi'(v_j(n)) = \frac{1}{(1 + e^{-av_j(n)})^2} \frac{\partial}{\partial v_j(n)} (1 + e^{-av_j(n)})$$

$$\phi'(v_j(n)) = \frac{a e^{-av_j(n)}}{(1 + e^{-av_j(n)})^2}$$

$$\phi'(v_j(n)) = \left\{ a \cdot 1 - \frac{a e^{-av_j(n)}}{1 + e^{-av_j(n)}} \right\} \left\{ \frac{e^{-av_j(n)}}{1 + e^{-av_j(n)}} \right\}$$

from *

$$\phi'(v_j(n)) = a \left\{ \phi(v_j(n)) \right\} \left\{ \frac{e^{-av_j(n)}}{1 + e^{-av_j(n)}} \right\}$$

Adding & subtracting 1.

$$\phi'(v_j(n)) = a \left\{ \phi(v_j(n)) \right\} \left\{ \frac{1 + e^{-av_j(n)}}{1 + e^{-av_j(n)}} - 1 \right\}$$

$$\phi'(v_j(n)) = a \left[\phi(v_j(n)) \right] \left[1 + \left[\frac{e^{-av_j(n)} - 1 - e^{-av_j(n)}}{1 + e^{-av_j(n)}} \right] \right]$$

$$\phi'(v_j(n)) = a \phi(v_j(n)) \left[1 - \frac{1}{1 + e^{-av_j(n)}} \right]$$

$$\therefore \phi'(v_j(n)) = a \phi(v_j(n)) [1 - \phi(v_j(n))] \quad \star\star$$

$$\phi'(v_j(n)) = a y_0(1 - y_0) \quad \star\star$$

② Tan hyperbolic: [Ranges from +1 to -1]

$$y_0 = \phi(v_j(n)) = a \tanh(b v_j) \quad ①$$

$$\phi'(v_j(n)) = ab \sec^2 h(b v_j)$$

$$\phi'(v_j(n)) = ab (1 - \tan^2 h(b v_j))$$

Multiply & divide by a .

$$\phi'(v_j(n)) = \frac{a^2 b}{a} (1 - \tan^2 h(b v_j))$$

$$\phi'(v_j(n)) = \frac{b}{a} (a^2 - a^2 \tan^2 h(b v_j))$$

It is of the form
 $(a^2 - b^2) = (a+b)(a-b)$

$$\phi'(v_j(n)) = \frac{b}{a} [a + a \tanh(b v_j)] [a - a \tanh(b v_j)]$$

from ①

$$\phi'(v_j(n)) = \frac{b}{a} \{a + \phi(v_j)\} \{a - \phi(v_j)\}$$

$$\phi'(v_j(n)) = \frac{b}{a} \{a + y_0\} \{a - y_0\} \quad \star\star$$

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$\left\{ e^{-\infty} = 0 \right\}$
 $\left\{ e^{\infty} = 1 \right\}$

$$x \rightarrow \infty \Rightarrow \tanh x \rightarrow 1$$

$$x \rightarrow -\infty \Rightarrow \tanh x \rightarrow -1$$

Conclusion : forward pass
 backward pass

Sigmoid activation function

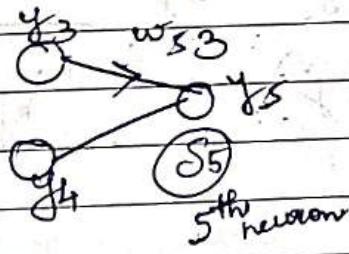
$$\Delta w_{ji} = \eta \delta_j y_i$$

$$\delta_j = \phi'(v_j(n)) e_j(n)$$

$= y_j(1-y_j)$ error ; j is output node.

$$= y_{pred}(1-y_{pred}) \quad (y_{desired} - y_{pred})$$

Note —



$$\delta_5 = y_5(1-y_5)(y_d - y_5)$$

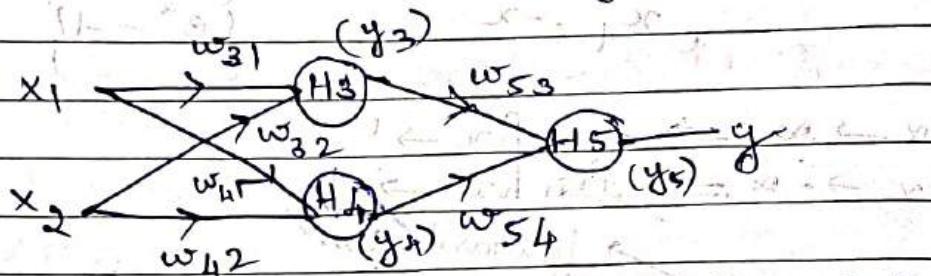
$$\Delta w_{53} = \eta \delta_5 y_3$$

$$\delta_j = \phi'(v_j(n)) \sum \delta_k w_{kj} ; j \text{ is hidden node}$$

$$\delta_j = \phi(v_j)(1-\phi(v_j)) \sum \delta_k w_{kj}$$

$$= y_j(1-y_j) \sum \delta_k w_{kj}$$

Problem on back propagation:



$$w_{31} = 0.1 \quad w_{41} = 0.2 \quad x_1 = 0.3 \text{ and } y_t = 0.5$$

$$w_{32} = 0.4 \quad w_{42} = 0.3 \quad x_2 = 0.5 \quad y_t = 0.5$$

$$w_{53} = 0.2 \quad w_{54} = 0.4 \quad \text{desired opf } y_t = 0.5 \quad \eta = 1$$

W.K.T. $w_{n+1} = w_n + \Delta w_{ji}$

and $\Delta w_{ji} = \eta \delta_j y_i$, ✓

where $\delta_j = y_j(1-y_j)$ (error); j is op neuron

$$\delta_j = y_j(1-y_j) \sum_k \delta_k w_{kj} \quad ; j \text{ is hidden neuron}$$

Forward pass:

$$H_3 = w_{31}x_1 + w_{32}x_2$$

$$= 0.1(0.3) + 0.4(0.5)$$

$$= 0.03 + 0.20$$

$$= 0.23$$

$$y_3 = \phi(H_3)$$

$$= \frac{1}{1+e^{-H_3}} - \frac{1}{1+e^{-0.23}} = 0.557$$

$$H_4 = w_{41}x_1 + w_{42}x_2$$

$$= 0.3(0.3) + 0.2(0.5)$$

$$= 0.09 + 0.1$$

$$= 0.19$$

$$y_4 = \frac{1}{1+e^{-H_4}} = \frac{1}{1+e^{-0.19}} = 0.547$$

$$\begin{aligned} H_5 &= y_3 w_{53} + y_4 w_{54} \\ &= 0.557(0.2) + 0.547(0.4) \\ &= 0.3802 \end{aligned}$$

$$y_5 = \frac{1}{1+e^{-H_5}} = \frac{1}{1+e^{-0.3802}} = 0.581$$

As we know y_3, y_4, y_5 , our O/P desired is 0.5 but we got 0.581 so we need to update the weights.

$$\text{error} = y_d - y_5 = 0.5 - 0.581 = -0.081$$

We have to minimise the error, by finding local gradients hence weight updation.

y_5 is output neuron

$$\text{from (1)} \quad \delta_5 = y_5(1-y_5)(\text{error})$$

$$\begin{aligned} &= 0.581(1-0.581)(-0.081) \\ &= 0.581(0.419)(-0.081) \\ &= -0.019 \end{aligned}$$

Similarly, calculate δ_3 and δ_4

H_3 and H_4 are hidden neurons.

from (2)

$$\delta_j = y_j(1-y_j) \sum_k \delta_k w_{kj} \quad [k \text{ is O/P neuron}]$$

$$\begin{aligned} \delta_3 &= y_3(1-y_3) \delta_5 w_{53} \\ &= 0.557(1-0.557)(-0.019)(0.2) \\ &= -0.000937 \end{aligned} \quad [\text{O/P has one neuron, so no summation required}]$$

$$\begin{aligned}\delta_4 &= y_4(1-y_4) \delta_5 w_{54} \\ &= 0.547 (1-0.547) (-0.019) (0.4) \\ &= -0.00188\end{aligned}$$

$$\therefore \Delta w_{ji} = \eta \delta_j y_i ; w_{n+1} = w_n + \Delta w_{ji}$$

$$\begin{aligned}\Delta w_{53} &= \eta \delta_5 y_3 \\ &= 1 (-0.019) 0.557 \\ &= -0.0105\end{aligned}$$

$$\begin{aligned}w_{53}(\text{new}) &= w_{53}(\text{old}) + \Delta w_{53} \\ &= 0.2 - 0.0105 = 0.1895\end{aligned}$$

Similarly,

$$\begin{aligned}\Delta w_{54} &= \eta \delta_5 y_4 \\ &= 1 \times (-0.019) (0.547) \\ &= -0.0103\end{aligned}$$

$$\begin{aligned}w_{54}(\text{new}) &= w_{54}(\text{old}) + \Delta w_{54} \\ &= 0.4 - 0.0103 \\ &= 0.3897\end{aligned}$$

In the same way, $\Delta w_{ji} = \eta \delta_j y_i$

$$\Delta w_{31} = \eta \delta_3 x_1$$

$$\begin{aligned}&= 1 (-0.000937) (0.3) \\ &= -0.000281\end{aligned}$$

$$\begin{aligned}w_{31}(\text{new}) &= w_{31}(\text{old}) + \Delta w_{31} \\ &= 0.1 - 0.000281\end{aligned}$$

$$= 0.099720$$

$$\begin{aligned}\Delta w_{32} &= \eta \delta_3 x_2 \\ &= 1 (-0.009317) 0.5 \\ &= -0.00468\end{aligned}$$

$$\begin{aligned}w_{32}(\text{new}) &= w_{32}(\text{old}) + \Delta w_{32} \\ &= 0.4 - 0.00468 = 0.39953\end{aligned}$$

$$\begin{aligned}\Delta w_{41} &= \eta \delta_4 x_1 \\ &= 1 (-0.0018)(0.3) \\ &= -0.000564\end{aligned}$$

$$\begin{aligned}w_{41}(\text{new}) &= w_{41}(\text{old}) + \Delta w_{41} \\ &= 0.3 + (-0.000564) \\ &= 0.29943\end{aligned}$$

$$\begin{aligned}\Delta w_{42} &= \eta \delta_4 x_2 \\ &= 1 (-0.00188)(0.5) \\ &= -0.00094\end{aligned}$$

$$\begin{aligned}(w_{42})_{\text{new}} &= (w_{42})_{\text{old}} + \Delta w_{42} \\ &= 0.2 + (-0.0094) \\ &= 0.19905\end{aligned}$$

Updated weights:

$$w_{81} = 0.0997 \quad w_{41} = 0.29943 \quad w_{53} = 0.1895$$

$$w_{32} = 0.39953 \quad w_{42} = 0.19905 \quad w_{54} = 0.3897$$

Note - The δ values decreases as we propagate back from output layers to hidden layers. At some point, the δ value might become zero. This is called vanishing gradient problem. Weight updation depends on

Δw_{ji} , which depends on δ , so it may not happen if gradients vanish.

More the deeper the layers become in ANN, we come across this problem.

Iteration 2 - forward pass:

$$H_3 = w_{31}x_1 + w_{32}x_2 \\ = 0.0997(0.3) + 0.3953(0.5) \\ = 0.2296$$

$$y_3 = \phi(H_3) = \frac{1}{1+e^{-H_3}} = \frac{1}{1+e^{-0.2296}} = 0.557$$

$$H_4 = w_{41}x_1 + w_{42}x_2 \\ = 0.29943(0.3) + 0.19905(0.5) \\ = 0.1893 \quad [\text{Same value because change in weight vector was small}]$$

$$y_4 = \frac{1}{1+e^{-H_4}} = \frac{1}{1+e^{-0.1893}} = 0.547$$

$$H_5 = y_3 w_{53} + y_4 w_{54} \\ = 0.557(0.1895) + 0.547(0.3897) \\ = 0.3187$$

$$y_5 = \frac{1}{1+e^{-0.3187}} = 0.579$$

$$\text{error} = y_t - y_5 \\ = 0.5 - 0.579 \\ = -0.079$$

which depends on δ , so it may not
gradient vanish.

so the layers become in ANN,
across this problem.

- forward pass:

$$x_1 + w_{32} x_2 \\ (0.3) + 0.3953(0.5)$$

$$H_3 = \frac{1}{1+e^{-H_3}} = \frac{1}{1+e^{-0.2296}} = 0.557$$

$$x_1 + w_{42} x_2$$

$$y_3(0.3) + 0.19905(0.5) \\ H_3 \quad [\text{Same value because change in weight vector was small}]$$

$$e^{-H_4} = \frac{1}{1+e^{-0.1893}} = 0.547$$

$$w_{53} + y_3 w_{54}$$

$$+(0.1895) + 0.547(0.3897)$$

$$e^{-0.3187} = 0.579$$

$$-y_5 \\ -0.579 \\ 0.079$$

In XOR gate, forward & backward pass for all the possible combinations, then move on to next iteration.

$$\delta_5 = y_5(1-y_5) \text{ error}$$

$$= 0.579(1-0.579)(0.079) \\ = -0.0192$$

$$\delta_3 = y_3(1-y_3) \delta_5 w_{53}$$

$$= 0.557(1-0.557)(-0.0192)(0.1895) \\ = 0.557(0.443)(-0.0192)(0.1895) \\ = -0.0008977$$

$$\delta_4 = y_4(1-y_4) \delta_5 w_{54}$$

$$= 0.547(1-0.547)(-0.0192)(0.3897) \\ = 0.547(-0.453)(-0.0192)(0.3897) \\ = -0.001854$$

$$\Delta w_{53} = \eta \delta_5 y_3$$

$$= -0.0192(0.557) \\ = -0.0106944$$

$$w_{53}(\text{new}) = w_{53}(\text{old}) + \Delta w_{53}$$

$$= 0.1895 - 0.0106944 \\ = 0.1788056$$

$$\Delta w_{54} = \eta \delta_5 y_4$$

$$= -0.0192(0.547) \\ = -0.0105024$$

$$\begin{aligned}w_{54}(\text{new}) &= w_{54}(\text{old}) + \Delta w_{54} \\&= 0.3897 - 0.0105024 \\&= 0.3791976\end{aligned}$$

$$\begin{aligned}\Delta w_{31} &= \eta \delta_3 x_1 \\&= (-0.0008977)(0.3) \\&= -0.00026931\end{aligned}$$

$$\begin{aligned}w_{31}(\text{new}) &= w_{31}(\text{old}) + \Delta w_{31} \\&= 0.0997 - 0.00026931 \\&= 0.09943069\end{aligned}$$

$$\begin{aligned}\Delta w_{32} &= \eta \delta_3 x_2 \\&= (-0.0008977)(0.5) \\&= -0.00044885\end{aligned}$$

$$\begin{aligned}w_{32}(\text{new}) &= w_{32}(\text{old}) + \Delta w_{32} \\&= 0.39953 - 0.00044885 \\&= 0.39908115\end{aligned}$$

$$\begin{aligned}\Delta w_{41} &= \eta \delta_4 x_1 \\&= (-0.001854)(0.3) \\&= -0.0005562\end{aligned}$$

$$\begin{aligned}w_{41}(\text{new}) &= w_{41}(\text{old}) + \Delta w_{41} \\&= 0.29943 - 0.0005562 \\&= 0.2988738\end{aligned}$$

$$w_{42} = \eta s_4 x_2$$

$$= (-0.001854) (0.5)$$

$$= -0.000927$$

$$w_{42(\text{new})} = w_{42(\text{old})} + \Delta w_{42}$$

$$= 0.19905 - 0.000927$$

$$= 0.198123$$

updated weights:-

$$w_31 = 0.09943069 \quad w_41 = 0.2988738$$

$$w_32 = 0.39908115 \quad w_42 = 0.198123$$

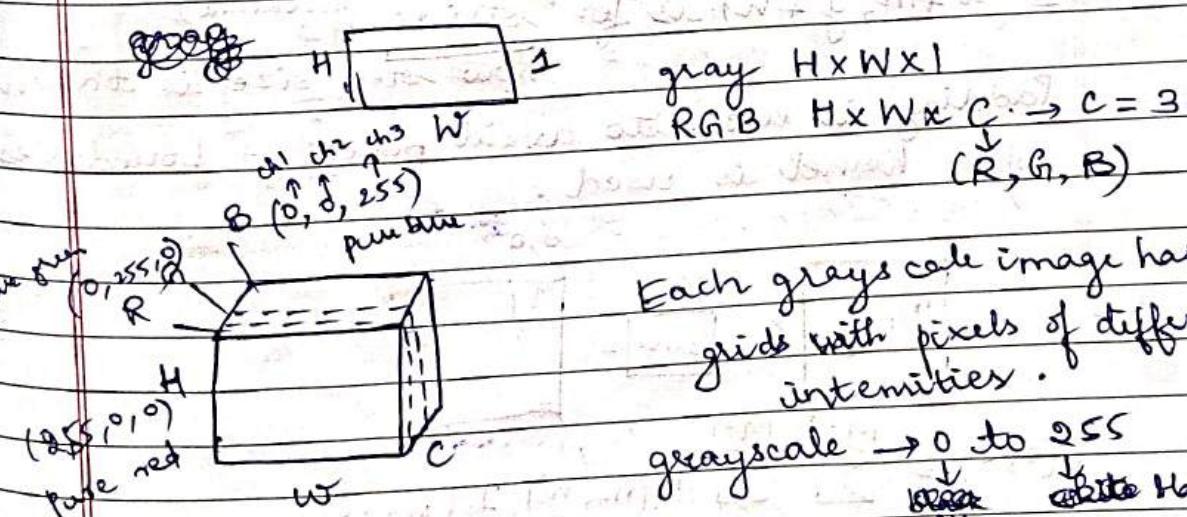
$$w_33 = 0.1988056 \quad w_43 = 0.3491976$$

ANN - classifier not feature extractor

convolution - feature extraction

Convolutional Neural networks:

A signal is 1D representation and image is 2D representation. At every time stamp, there'll be amplitudes. Images are grayscale or RGB.



Each grayscale image has grids with pixels of different intensities.

For colour images, three such channels exist with values for each pixel ranging from 0 to 255.

$H \downarrow \quad W \downarrow \quad 3 \text{ channels}$
 $32 \times 32 \times 3$

(Height & width will be same across all channels)

Feature map:-

We define a kernel which will run through the entire image matrix. This creates a feature map. This is the process of feature extraction.

$3 \times 3, 5 \times 5 \text{ or } 7 \times 7$ sized kernels are used in convolutions. 7×7 has larger field of view comparatively but it is computationally expensive.

! why can't we use 2×2 kernel?

Let $I(i, j) \rightarrow$ input image

$K(m, n) \rightarrow$ Kernel / filter

Convolution: $O(i, j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K(m, n) \cdot I(i+m, j+n)$

where $K(m, n)$ is a Kernel with dimension $M \times N$.

$\begin{bmatrix} 1, 0, 1, 1, \dots, 0, 1, P \\ 1, 0, 1, 1, \dots, 0, 1, P \\ \vdots \\ 1, 0, 1, 1, \dots, 0, 1, P \end{bmatrix}$ } generalised representation for image & Kernel

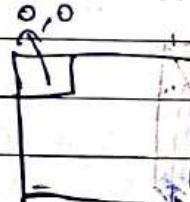
If it is image $\rightarrow P=32, Q=32$, etc.

If it is Kernel $\rightarrow Q=3, P=3$.

$I(i+m, j+n)$ is for stride: (movement)

how step size is taken

Padding is used to avoid pixels at boundaries when Kernel is used.



$$O(0,0) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K(m, n) I(m, n)$$

Assume kernel size = 3×3

$M=3, N=3$

$$o(0,0) = \sum_{m=0}^2 \sum_{n=0}^2 k(m,n) I(m,n)$$

$$m=0 \Rightarrow k(0,0) I(0,0) + k(0,1) I(0,1) + k(0,2) I(0,2)$$

$$m=1 \Rightarrow +k(1,0) I(1,0) + k(1,1) I(1,1) + k(1,2) I(1,2)$$

$$m=2 \Rightarrow +k(2,0) I(2,0) + k(2,1) I(2,1) + k(2,2) I(2,2)$$

Consider -

a	b	c	d	q	r	s
e	f	g	h	t	u	v
i	j	k	l	w	x	y
m	n	o	p			

$$K(m,n) \rightarrow M=3$$

$$N=3$$

$$o(0,0) = axq + br + cs + et + fu + gv + iw + jx + ky$$

(element-wise multiplication)

$M-1 \quad N-1$

$$\text{Similarly, } o(i,j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} k(m,n) I(i+m, j+n)$$

$$i=0, j=1 \quad M-1 \quad N-1$$

$$o(0,1) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} k(m,n) I(m, 1+n)$$

$m=0$ to 2 , $n=0$ to 2

0	1	2
0	$o(0,0)$	$o(0,1)$
1	$o(1,0)$	$o(1,1)$
2	1	2

2×2

$$o(0,1)$$

0	a	b	c	d
1	e	f	g	h
2	i	j	k	l
3	m	n	o	p

1	q	r	s
2	t	u	v
3	w	x	y

no space
downwards or
towards the
right
(convolution)

$$o(1,0) \quad o(1,1)$$



~~import numpy as np~~

~~def sigmoid(x):~~

~~return 1/(1+np.exp(-x))~~

~~def sigmoid_der(x):~~

~~return x*(1-x)~~

class NN:

~~def __init__(self, inputs):~~

~~self.inputs = inputs~~

~~self.l = len(self.inputs)~~

class Perception:

~~def __init__(self, learning_rate = 0.01, n_iters=100):~~

~~self.lr = learning_rate~~

~~self.n_iters = n_iters~~

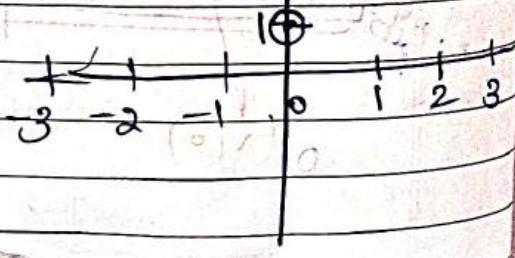
~~self.activation_func = self.sigmoid~~

~~self.weights = None~~

~~self.bias = None~~

$$x[n] = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

$$x[-n] = \begin{bmatrix} 3 & 2 & 1 \end{bmatrix}$$



Ex-

Image

$$O(0,0) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K(m,n) I(m, n)$$

Assume kernel size = 3×3

When an image is convolved with kernel, there is reduction in size. 4×4 image \star 3×3 kernel $\rightarrow 2 \times 2 \times 1$ image.
 This was the case for grayscale image (channel = 1).
 \downarrow
 $H \times W \times 1$

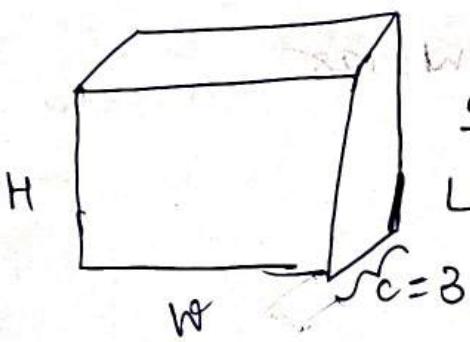
channel no. \rightarrow tensor.

For RGB channels

$$I(H \times W \times C) \star k(m \times n \times C) \Rightarrow O(H_1 \times W_1 \times k)$$

change in no. of channels

RGB



$I(H, W, 3)$

$\hookrightarrow k(m, n, 3)$

\hookrightarrow no. of channels

Depth of Kernel is decided by preceding layer

Convolution process is the same, but it happens for the three channels individually & there will be a stride. Three feature maps will be obtained, upon which element wise addition occurs.

3×3 image

a_1	b_1	c_1
d_1	e_1	f_1
g_1	h_1	i_1

R.E

a_2	b_2	c_2
d_2	e_2	f_2
g_2	h_2	i_2

G

a_3	b_3	c_3
d_3	e_3	f_3
g_3	h_3	i_3

B

If I/F layer has "64 channels, depth of kernel also has 64 channels

3×3 kernel

m_1	n_1	o_1
p_1	q_1	r_1
s_1	t_1	u_1

\downarrow

v_1	w_1
-------	-------

\downarrow

v_2	w_2
-------	-------

v_3	w_3
-------	-------



1 Kernel with depth 3 \rightarrow 1 feature map

Same kernel cannot multiply with different channels

$v_1 + v_2 + v_3$ because of different weights

$v_2 + v_3$ because of different weights

v_3 because of different weights

To increase the no. of feature maps —
we must consider more no. of kernels

one kernel
 $(3, 3, 3) \times 64 \rightarrow 64$ feature maps.
(They have different combination of weights)

Consider —

$$\text{Image} = (H \times W \times C)$$

① convolution $\rightarrow (3, 3) \times 64 \rightarrow H \times W \times 64$

depth $\rightarrow C$

Convolution $\rightarrow (3, 3) \times 128 \rightarrow H \times W \times 128$

depth = 64
(of previous layer)

No. of kernels is user defined.

Convolution $(3, 3) \times 256 \rightarrow H \times W \times 256$

depth = 128
 $3 \times 3 \times 128$

Total channels
 $= 128 \times 256$

①

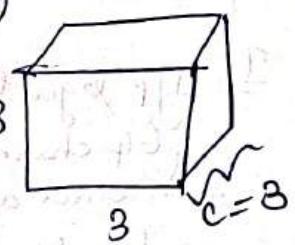


Image \rightarrow
 $H \otimes W \otimes 1$

(64 Kernels)

one image \rightarrow multiplied with 64 kernels

64 channels

IP image



{128 Kernels}

②

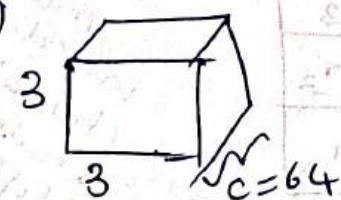


Image \rightarrow multiplied with 128 kernel

128 channels

$$o(i, j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K(m, n) I(i+m, j+n) \quad \{ \text{stride} = 1 \}$$

(stride $(i+m, j+n)$)

For a 3×3 kernel, stride $\rightarrow 1$ ✓
 2 ✓ } In general, we limit by
 3 ✓ stride 1 or 2.

By including stride,

$$o(i, j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} K(m, n) I(i+s+m, j+s+n) \quad (s = \text{stride})$$

a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y

5x5 image

1	2	3
4	5	6
7	8	9

3x3

x	x	x
x	x	x
x	x	x

3x3 image.

(stride $\rightarrow 1$)

Dimensionality of image reduced because it's entirely not concentrating on all pixels.

There are 4 diagonal & adjacent pixels

so, when kernel is multiplied entire concentration is on

pixel g, because adjacent pictures don't vary much in information.

Moving by one step, 5 overlaps with h, similarly i is

central pixel, so on

We are losing information of edge pixels hence dimensionality reduced.

How to concentrate on α ? \rightarrow Padding
edge pixels

Let's consider - Original Image: (H, W, c)
Kernel: $(K \times K)$

Output Image: (H_1, W_1, D_1)

$$\text{where, } H_1 = \frac{(H + 2P - K)}{S} + 1$$

$$W_1 = \frac{(W + 2P - K)}{S} + 1$$

$D_1 = \text{No. of kernels}$

$$(5 \times 5 \text{ image}) \quad H_1 = \left[\frac{5 + 2(0) - 3}{1} \right] + 1 = 3$$

$$(3 \times 3 \text{ kernel})$$

$$(1 \leftarrow \text{stride}) \quad W_1 = \left[\frac{5 + 0 - 3}{1} \right] + 1 = 3$$

Stride $\rightarrow 2$ Kernel (3×3)

	a	b	c	d	e
	f	g	h	i	j
	k	l	m	n	o
	p	q	r	s	t
	u	v	w	x	y

Convolution disadvantage -

- 1) Less field of view
[Not covering entire kernel in a single step]

Original image: $224 \times 224 \times 3$

Kernel $\rightarrow 3 \times 3$

Stride $\rightarrow 1$

padding $\rightarrow 0$

$$H_1 = \frac{(224 - 3)}{1} + 1 = 222$$

$$W_1 = \frac{(224 - 3)}{1} + 1 = 222$$

$\Rightarrow 222 \times 222 \times 1$

How α ? \rightarrow Padding
2 pixels

Original Image: (H, W, C)
 $I = (k \times k)$

Image: (H_1, W_1, D_1)

$$H_1 = \frac{H + 2P - K}{S} + 1 \quad P \rightarrow \text{Padding}$$

$$W_1 = \frac{W + 2P - K}{S} + 1 \quad S \rightarrow \text{stride}$$

D_1 = No. of kernels

$$\left[\begin{matrix} 0 & -3 \\ 1 & \end{matrix} \right] + 1 = 2 \cdot 3$$

$$\left[\begin{matrix} 0 & -3 \\ 1 & \end{matrix} \right] + 1 = 2 \cdot 3$$

Input: 3×3

Output: 2×2

Convolution disadvantage -
Less field of view
(Not covering entire kernel
in a single sketch)

$$\left[\begin{matrix} 224 & -3 \\ 1 & \end{matrix} \right] + 1 = 222$$

$$\left[\begin{matrix} 224 & -3 \\ 1 & \end{matrix} \right] + 1 = 222$$

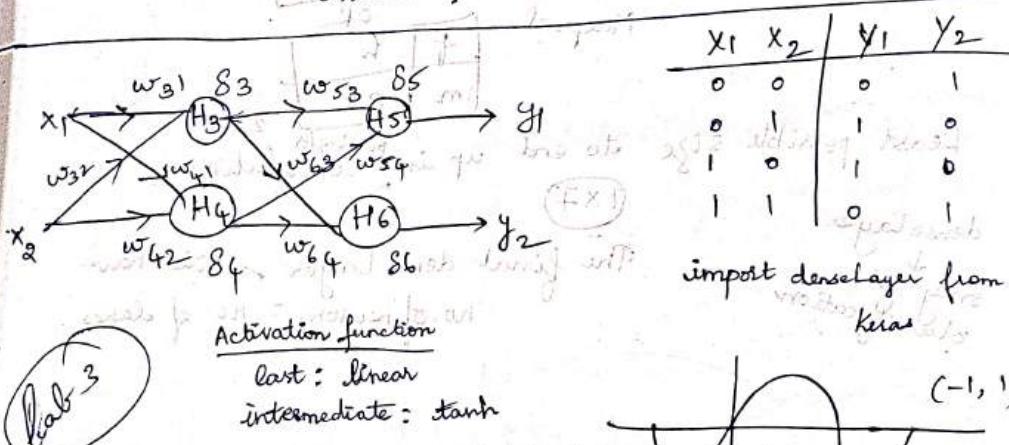
$$\Rightarrow P \rightarrow 222 \times 222 \times 1$$

When stride is increased from 1 to 2, image dimension reduced to half (less computation)

Padding $\begin{cases} \text{Valid padding} \\ \text{Same padding} \end{cases}$

No. of pixels to be added/padded depends on kernel size.

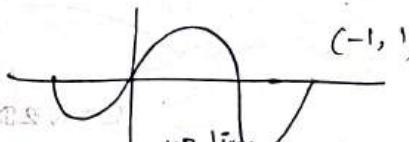
Convolution \rightarrow flip the kernel & correlate it,
↓
DL convolution is just correlation!



Activation function

Last: linear
intermediate: tanh

import denseslayer from
keras



train
for 500,
1000,
15000
epochs

$$y = 1 + \sin 2\pi x$$

not fit enough to teach function

for test we want if next x value bigger \rightarrow positive zero - interpolation will be

various factors causing the function to change

Pooling

average
max

- Reducing image size.

$224 \times 224 \times 3$

Conv2D (64, 3x3,)

MaxPool (2x2)

↓
every 4 pixels - 1 pixel

Avg pool

a	b	c	d
e	f	g	h
i	j	k	l
m	n	o	p

4×4

$a+b+t$	$c+d+g+h$
$e+f+g+h$	$i+j+k+l$
$m+n+o+p$	$t+u+v+w$
4	4

2×2

Maxpool

f	h
m	p

2×2

Least possible size to end up in convolution is

7×7

The final denseLayer should have
no. of neurons = no. of classes

dense layer

only for classification

Conv2D ();

MaxPool ;

flatten / global average pooling

Work on MNIST / Cifar 10

1. Download dataset from keras

2. Split into x-train, y-train, x-test, y-test
one hot encoding
not required for sparse categorical cross entropy

3. Use categorical-cross entropy

3. Use to_categorical to perform label encoding
one hot encoding

Consider		R_0	R_1	R_2	R_3	R_4
0	Image 1	1	0	0	0	e class 0
1	Image 2	0	1	0	0	e class 1
2	Image 3	0	0	1	0	e class 2
3	Image 4	0	0	0	1	e class 3

4. Normalise the image

- i) use dense layers directly & classify the images by reshape [flatten] without feature extraction.

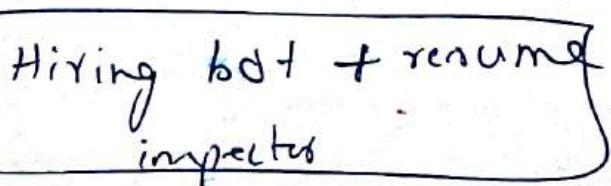
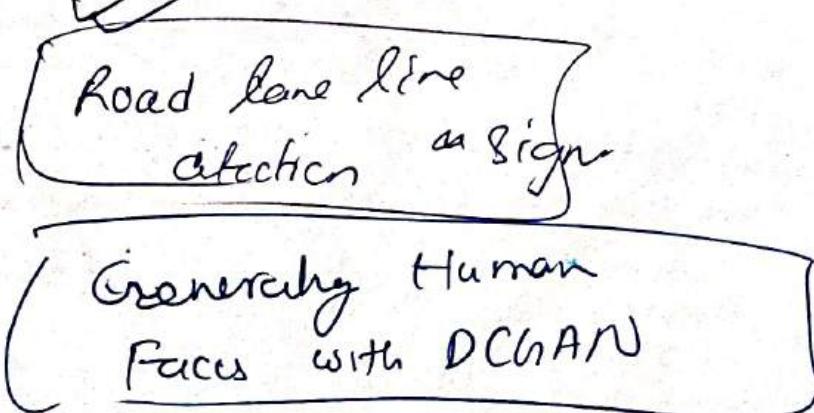
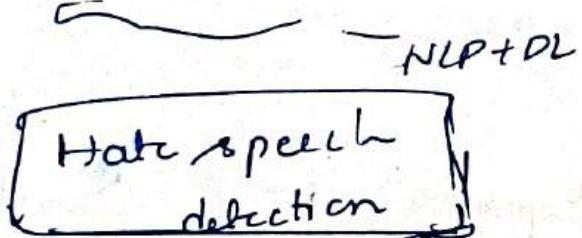
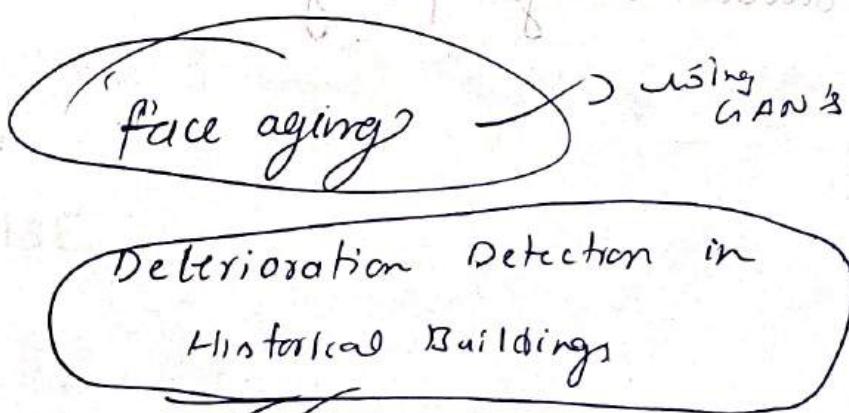
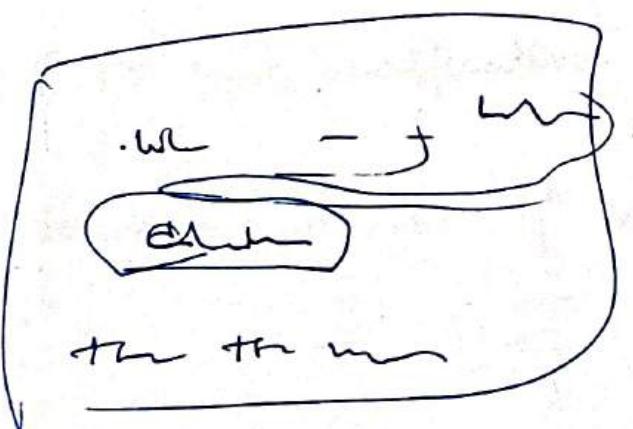
Dense

:

Dense(10)

check accuracy & (model.predict) & print some images.

- ii) Use convolutional layer, dense layer, maxpool for classifications.



Padding

2 types:

1. Valid padding - No pixels are pad
2. Same padding - Pixels are pad such that image size is preserved.

MSE good choice for regression (not classification)

Batch normalisation - used to maintain zero average and unit variance.

Average Pooling: Reduce the image size by half by averaging pixel values.

Maximum Pooling: Take the maximum value from a pixel range.

Global average pooling:

Binary cross entropy :-

$$\text{MSE: } J = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2$$

↓
desired predicted

In binary classification:-

Desired values $\rightarrow 0$ and 1

Predicted values \rightarrow range from $[0, 1]$

$\hat{y} > 0.5$ (class 1)

$\hat{y} < 0.5$ (class 0)

Sigmoid activation function brings o/p in range $[0, 1]$.

Rule :-

Loss J :

① In case of correct classification, (y_i & \hat{y}_i are same)
then $J = 0$.

② For misclassification, loss J must be very high.

↓

However, in case of MSE

$$J = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2$$

Let $y_j = 0$, and obtained $\hat{y}_j = 0.99$ { then

$\hat{y}_j = 1$ { the

penalising
of loss
is not
much

(Not much
different)

③ Always $J \geq 0$ (Always true for MSE).

Binary cross entropy $BCE = -[y \ln \hat{y} + (1-y) \ln (1-\hat{y})]$

$y \rightarrow$ desired
 $\hat{y} \rightarrow$ predicted

case(i) $y=0 ; \hat{y}=0 / y=1 ; \hat{y}=1$ [correct classification]

$$\begin{aligned} BCE &= -[0 + \ln 1] & y=1 & \hat{y}=1 \\ &= -\ln 1 & = -[\ln 1 + 0] \\ &= 0 & = -\ln 1 \\ & & & = 0 \end{aligned}$$

case(ii) $y=0 ; \hat{y} \approx 1$

$$BCE = -[0 \ln 1 + (1-0) \ln(1-1)]$$

$$= -[0 + \ln 0]$$

$$= -[-\infty]$$

output from $= \infty$ since T need a well-defined value for ③

$$y=0 ; \hat{y}=0.99$$

$$BCE = -[0 \ln 0.99 + (1-0) \ln(1-0.99)]$$

$$= -[\ln 0.01] = -[-4.63]$$

$$= 4.63$$

case(iii) clearly visible from (ii)

$$y \ln \hat{y} + (1-y) \ln (1-\hat{y})$$

\rightarrow desired

\rightarrow predicted

$y=1; \hat{y}=1$ [correct classification]

$$\begin{aligned} & y=1 \quad \hat{y}=1 \\ & = -[\ln 1 + 0] \\ & = -\ln 1 \\ & = 0 \end{aligned}$$

$$+ (1-0) \ln(1-1)$$

$\ln 0$

$$0.99 + (1-0.99) \ln(1-0.99)$$

$$0.01 = -[-4.63]$$

$$= 4.63$$

from (ii) $\Rightarrow 0.5$ $\Rightarrow 0.5$

for N examples:

$$BCE = \frac{1}{N} \sum_{i=1}^N [y_i \ln \hat{y}_i + (1-y_i) \ln (1-\hat{y}_i)]$$

BCE is the modification of categorical cross entropy.

Categorical cross entropy : (Multi-class classification)
(softmax activation function)

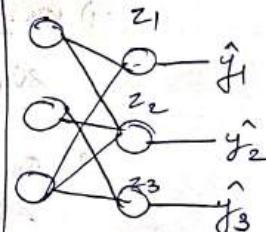
$$J = - \sum_{i=1}^n y_i \ln \hat{y}_i ; n \rightarrow \text{No. of classes}$$

$$\text{if } n=2$$

$$= -[y_1 \ln \hat{y}_1 + y_2 \ln \hat{y}_2]$$

$$\boxed{\hat{y}_2 = (1-\hat{y}_1)} \quad \boxed{\hat{y}_2 \neq (1-\hat{y}_2)}$$

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum e^{z_j}}$$



$$\hat{y}_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_1 + \hat{y}_2 + \hat{y}_3 = 1$$

sigmoid

✓ scores for two classes
positive in case of multilabel classification
let's say

$$\begin{aligned} z_1 &\rightarrow \hat{y}_1 = [0, 1] & 0.7 &\checkmark \\ z_2 &\rightarrow \hat{y}_2 = [0, 1] & 0.2 &\checkmark \\ z_3 &\rightarrow \hat{y}_3 = [0, 1] & 0.9 &\checkmark \end{aligned}$$

summation $\Rightarrow 1$

For binary classification, we use sigmoid function
& not softmax, we can use softmax too but we
are increasing the complexity (no. of parameters) of network
by additional neurons

Dense (1)

↓
sigmoid with BCE

but for softmax, we need to calculate probabilities & hence
need 2 neurons.

Dense (2)

0. 1111
0. 1100
0. 1110

Lab Practicals

Exp-1 MNIST dataset has 10 classes, make it to 2
classes. ($y=0, y=1$)

1) In the last layer, use dense (1) with
sigmoid activation function & loss as binary
cross entropy.

2) In the last layer, use dense (2) with
softmax activation function & loss as
categorical cross entropy.

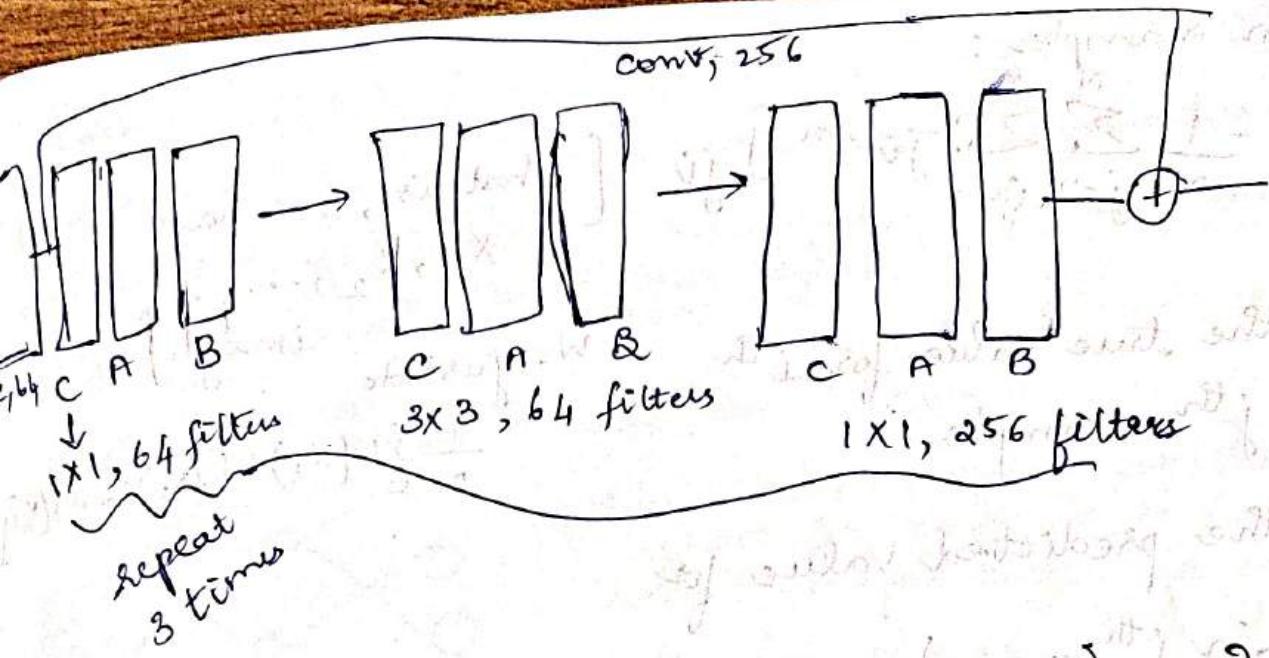
Exp-2: Develop a 50 layer residual network with
with 3×3 , 64, 1
kernel size stride
kernel

CIFAR10
dataset

1x1 layer

Conv ($3 \times 3, 64, 1$)

OR
Conv2D ($64, 3, 1$)



$\text{Conv2D}(\text{kernel size} = 3, \text{filter} = \text{filter}, \text{stride} = 1)$

$\text{Conv2D}(\text{kernel size} = 3, \text{filter} = \text{filter}, \text{stride} = 1) \times 3$

$\text{Conv2D}(\text{kernel size} = 1, \text{filter} = \text{filter}, \text{stride} = 1)$

$\text{Conv2D}(\text{kernel size} = 1, \text{filter} = \text{filter})$

Deep Residual network for image recognition

Image \rightarrow x \leftarrow cat, dog, horse

$y_i = [0 \ 1 \ 0]$ Image is of dog.

$\hat{y}_i = [0.1 \ 0.2 \ 0.7]$ But model predicted is as horse.

So, we are calculating categorical cross entropy.

$$L = - \sum_{i=1}^n y_i \ln \hat{y}_i$$

$$\therefore L = - [(0 \times \ln 0.1) + (1 \times \ln 0.2) + (0 \times \ln 0.3)]$$

$$= 1.6$$

For n examples:

$$L = -\frac{1}{N} \sum_{j=1}^N \sum_{i=1}^n y_{ji} \ln \hat{y}_{ji}$$
 [That is, we have
 x_1, x_2, \dots, x_n
images]

y_{ji} is the true value for i^{th} class in j^{th} example.

\hat{y}_{ji} is the predicted value for i^{th} class in j^{th} example.

Batch Normalization:

Normalization :- Since images have pixels in the range of $[0, 255]$, we divide each pixel by 255 so that the pixels can be brought in the range of $[0, 1]$.

While images are being trained, they will be passed as batches to neural networks.

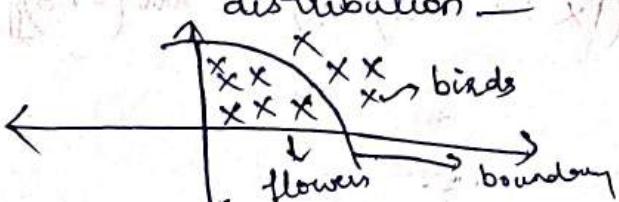
Consider that 100 flowers & birds have been split into 2 batches.

Batch 1

red rose } all
hibiscus } red flowers
green parrot

:

for batch 1 images,
distribution —

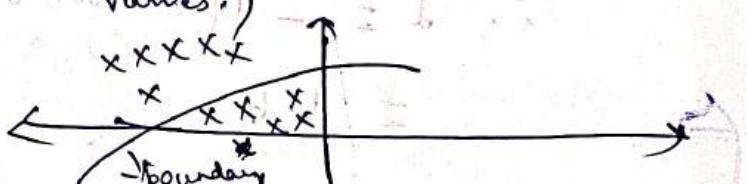


Batch 2

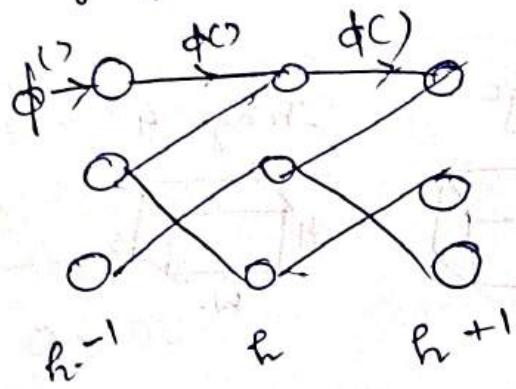
color parrot
white rose

yellow sunflower (other yellow
flowers)

(due to colour variation, distribution
varies.)



Since samples are distributing at different places & different boundaries, model will take so much time to learn the patterns, therefore batch normalization is very much significant. (for convolutional layer)



For dense layers:-

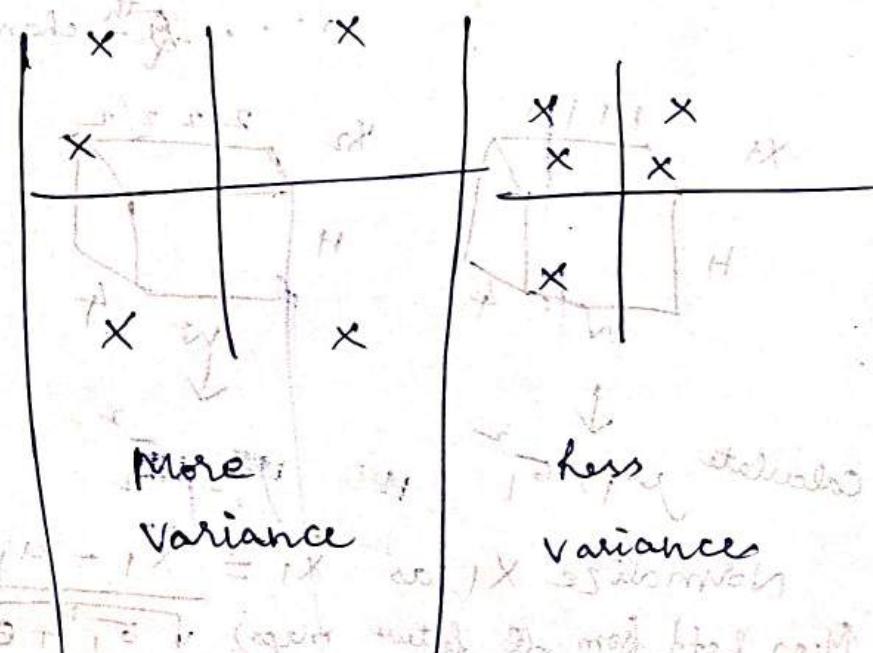
while backpropagation, weights are updated for each layer, so weights might differ for every layer, so when passing batch of images without normalization, there will be difference in distribution.

The shift in the distribution for same inputs due to multiplication with different weights for each layer is called covariate shift.

In order to suppress covariant shift, we use batch normalization where we strive to make mean zero & unit variance.

(All samples will get distributed at origin.)

Variance \rightarrow Spread of data
Samples



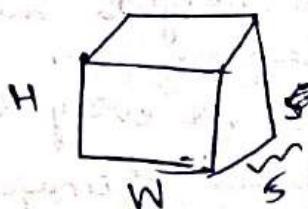
Consider:-

(HxW) Image 1 Image 2 Image 3 Image 4

Batch size = 4

No. of kernels = 5

of Image 1



$H \times W \times S$
channels

Image 2

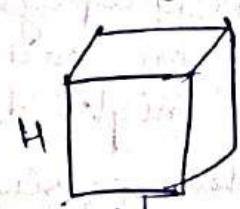


Image 3

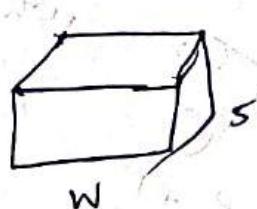
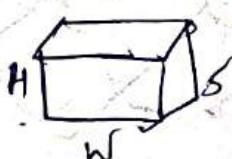


Image 4



each feature map has 5 channels \rightarrow due to
5 kernels

That is, $X \in \mathbb{R}^{4 \times H \times W \times 5}$

Image 1

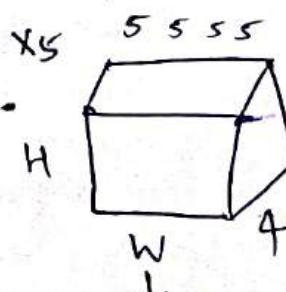
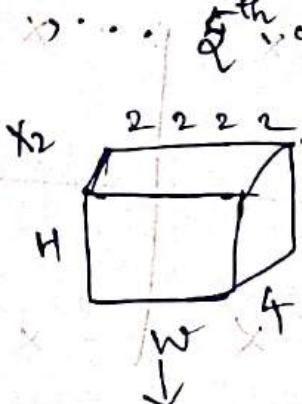
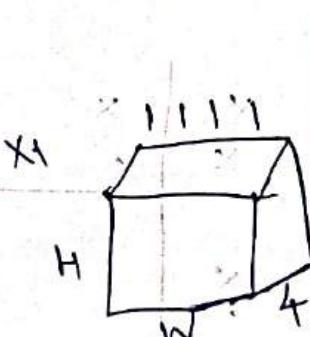


Image 4



Separation channelwise :- (Bring together respective feature maps)

Bring all the 1st channels together (which is due to kernel 1)



Calculate μ_1, σ_1^2

μ_2, σ_2^2

μ_5, σ_5^2

Normalize X_1 as $\hat{X}_1 = \frac{X_1 - \mu_1}{\sqrt{\sigma_1^2 + \epsilon}}$ & do for all images
(Mean std from all feature maps) & for every pixel

where $\epsilon = 10^{-5}$ because if in case standard deviation becomes zero, $\hat{x} \rightarrow \infty$
 ~ to avoid it, ϵ is added.

This process is called batch normalization.

(Elements of BN)

Same for both dense layers & convolution layers.

1. Obtain μ, σ^2 for all feature maps.
2. Do \hat{x} for all pixels

Consider two channels of images -

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

batch size = 2 images

each of $2 \times 2 \times 2$

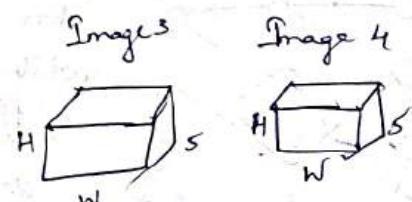
$$\mu = \frac{1+2+3+4+5+6+7+8}{8} = 4.5$$

$$\sigma^2 = \frac{(1-\mu)^2 + (2-\mu)^2 + (3-\mu)^2 + \dots + (8-\mu)^2}{8}$$

$$\begin{bmatrix} \frac{1-\mu}{\sqrt{\sigma^2+\epsilon}} & \frac{2-\mu}{\sqrt{\sigma^2+\epsilon}} \\ \frac{3-\mu}{\sqrt{\sigma^2+\epsilon}} & \frac{4-\mu}{\sqrt{\sigma^2+\epsilon}} \end{bmatrix} \quad \begin{bmatrix} \frac{5-\mu}{\sqrt{\sigma^2+\epsilon}} & \frac{6-\mu}{\sqrt{\sigma^2+\epsilon}} \\ \frac{7-\mu}{\sqrt{\sigma^2+\epsilon}} & \frac{8-\mu}{\sqrt{\sigma^2+\epsilon}} \end{bmatrix}$$

$$\begin{array}{|c|c|} \hline 3 & 6 \\ \hline 8 & 7 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 5 & 1 \\ \hline 4 & 8 \\ \hline \end{array}$$

$$3 \times 3 \times 2$$



a map has 5 channels \Rightarrow due to 5 kernels

$1 \times W \times C$ (width and depth and height of the image) \Rightarrow 5 channels

$W \times 5$ \Rightarrow 5 channels

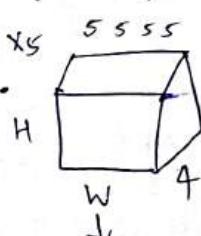
image 4

applies at where each channel may add up

respective feature maps

together (which is due to kernel 1)

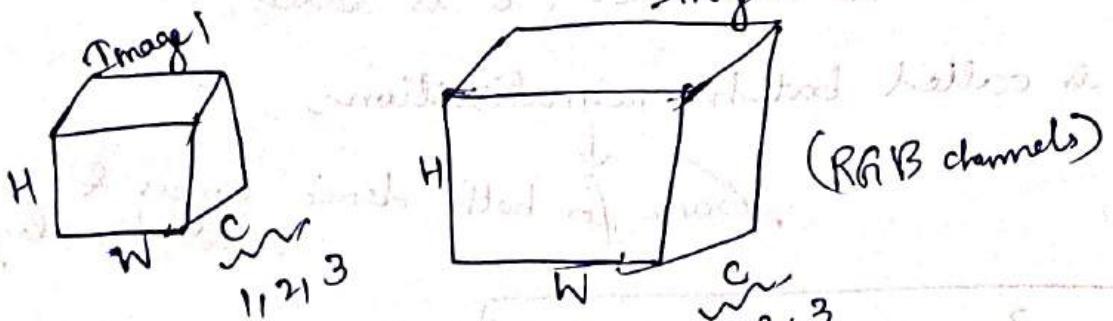
add together from every image,



$$\mu_5, \sigma^2_5$$

& do for all images
& for every output

Batch normalisation



Bringing together of channels:-

(1, 1)

(2, 2)

(3, 3)

calculate μ_1, σ_1^2 , μ_2, σ_2^2 , μ_3, σ_3^2

Perform normalisation on every pixel of respective feature map.

$$\hat{x} = \frac{x - \mu_1}{\sqrt{\sigma_1^2 + \epsilon}}$$

Let $X \in \mathbb{R}^{B \times H \times W \times C}$
 $N \rightarrow$ no. of images in one batch
 $H \rightarrow$ height
 $W \rightarrow$ width
 $C \rightarrow$ channels
 $B \rightarrow$ total no. of images before dividing into one particular batch

$$\text{Mean} = \mu_c = \frac{1}{N \cdot H \cdot W} \sum_{i=0}^{N-1} \sum_{j=0}^{H-1} \sum_{k=0}^{W-1} x_{ijk}$$

for channel #1 \rightarrow 2 feature maps

Consider :-

1	2
3	4

5	6
7	8

$$\mu_c = \frac{1+2+3+4+5+6+7+8}{2 \times 2 \times 2}$$

$$\mu_c = \frac{1}{HW} \sum_{j=1}^H \sum_{k=1}^W x_{ijk} \quad (\text{Global average pooling})$$

$$\sigma_c^2 = \frac{1}{NHW} \sum_{i=1}^N \sum_{j=1}^H \sum_{k=1}^W (x_{ij} - \mu_c)^2$$

$$\text{Normalization} \quad \hat{x}_i = \frac{x_i - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}}$$

Example: $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 3 & 4 \\ 5 & 6 \end{bmatrix} \approx \begin{bmatrix} 2 & 6 \\ 8 & 10 \end{bmatrix}$

Image 1: channel 1 $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$

channel 2 $\begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$

Image 2: channel 1 $\begin{bmatrix} 9 & 10 \\ 11 & 12 \end{bmatrix}$

channel 2 $\begin{bmatrix} 13 & 14 \\ 15 & 16 \end{bmatrix}$

$$X = \begin{bmatrix} \text{Image 1} & \begin{bmatrix} \text{channel 1} & \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} & \text{channel 2} & \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \end{bmatrix} \\ \text{Image 2} & \begin{bmatrix} 9 & 10 \\ 11 & 12 \end{bmatrix} & \begin{bmatrix} 13 & 14 \\ 15 & 16 \end{bmatrix} \end{bmatrix}$$

$$\mu_1 = \frac{1+2+3+4+9+10+11+12}{2 \times 2 \times 2}$$

$$\sigma_1^2 = \frac{52}{8} = \frac{26}{4} = \frac{13}{2} = 6.5$$

$$\sigma_1^2 = (1-6.5)^2 + (2-6.5)^2 + (3-6.5)^2 + (4-6.5)^2 + (9-6.5)^2 + (10-6.5)^2 + (11-6.5)^2 + (12-6.5)^2 / 8$$

$$\sqrt{\sigma_1^2} = \sqrt{\frac{17.25}{17.25}} = 4.15$$

$$\mu_2 = \frac{5+6+7+8+13+14+15+16}{8} = 10.5$$

$$\sigma_2^2 = \frac{(5-10.5)^2 + (6-10.5)^2 + (7-10.5)^2 + (8-10.5)^2 + (13-10.5)^2 + (14-10.5)^2 + (15-10.5)^2 + (16-10.5)^2}{8}$$

$$\sigma_2^2 = 17.25$$

$$\sqrt{\sigma_2^2} = 4.15$$

$$\hat{X} = \left[\begin{array}{cc} \frac{1-6.5}{4.15} & \frac{2-6.5}{4.15} \\ \frac{3-6.5}{4.15} & \frac{4-6.5}{4.15} \\ \frac{9-6.5}{4.15} & \frac{10-6.5}{4.15} \\ \frac{11-6.5}{4.15} & \frac{12-6.5}{4.15} \end{array} \right] \left[\begin{array}{cc} \frac{5-10.5}{4.15} & \frac{6-10.5}{4.15} \\ \frac{7-10.5}{4.15} & \frac{8-10.5}{4.15} \\ \frac{13-10.5}{4.15} & \frac{14-10.5}{4.15} \\ \frac{15-10.5}{4.15} & \frac{16-10.5}{4.15} \end{array} \right]$$

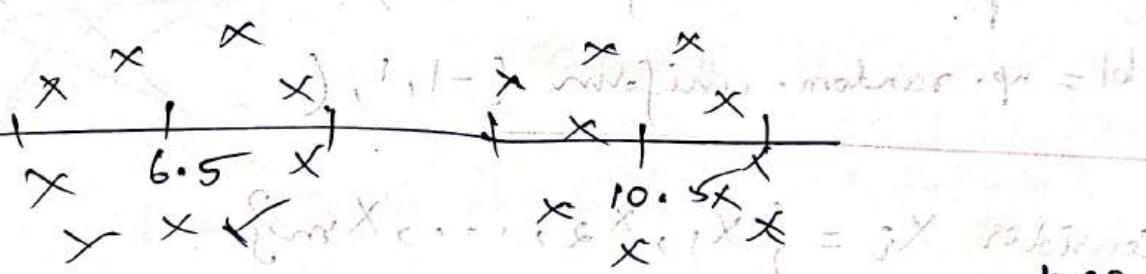
$$\hat{X} = \left[\begin{array}{cc} -1.32 & -1.08 \\ -0.23 & -0.6 \\ 0.6 & 0.84 \\ 0.618 & 1.32 \end{array} \right] \left[\begin{array}{cc} -1.3253 & -1.0842 \\ -0.843 & -0.602 \\ 0.602 & 0.843 \\ 1.084 & 1.3253 \end{array} \right]$$

μ after normalization $\rightarrow 0$

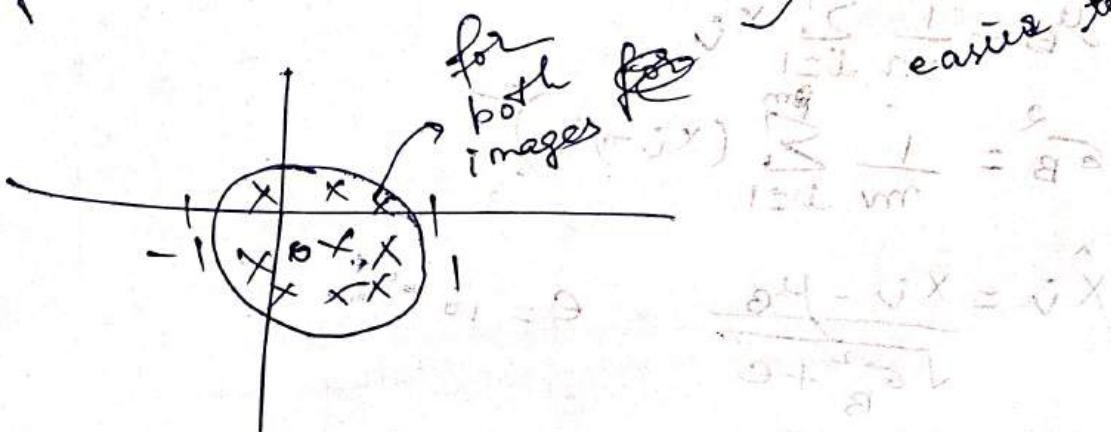
Variance after normalisation $\rightarrow 1$

$$\Rightarrow \frac{(-1.32)^2 + (1.08)^2 + (-0.84)^2 + (-0.6)^2 + (0.6)^2 + (0.84)^2 + (1.08)^2 + (1.32)^2}{8} = 1$$

Initially



After normalisation



(68) Pixel distribution
doesn't differ,
easier training.

$9 + 5x^2 = 1$ if normalized data

B P

```

import numpy as np
import matplotlib.pyplot as plt
x = np.array([[0, 0], [0, 1], [1, 0], [1, 1]])
y = np.array([0, 1, 1, 0])

```

lr = 0.5

n_epochs = 10000

except [EOF]

input_nodes = 2

hidden_nodes = 2

output_nodes = 1

w1 = np.random.uniform(-1, 1, (input_nodes, hidden_nodes))

w2 = np.random.uniform(-1, 1, (hidden_nodes, output_nodes))

b1 = np.random.uniform(-1, 1, (1, hidden_nodes))

Consider $x_i = \{x_1, x_2, \dots, x_m\}$

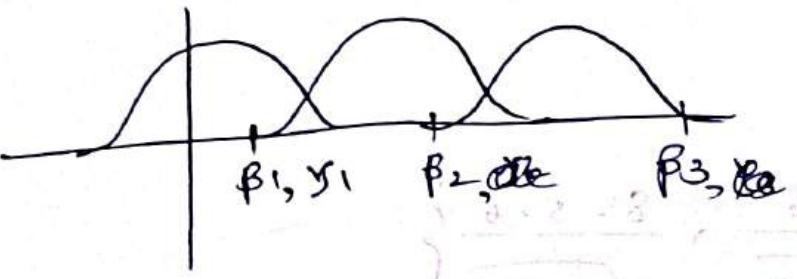
$$y_i = \sigma_{\beta}(\mathbf{w}_1 x_i + b_1)$$

$$\text{Then, } \mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

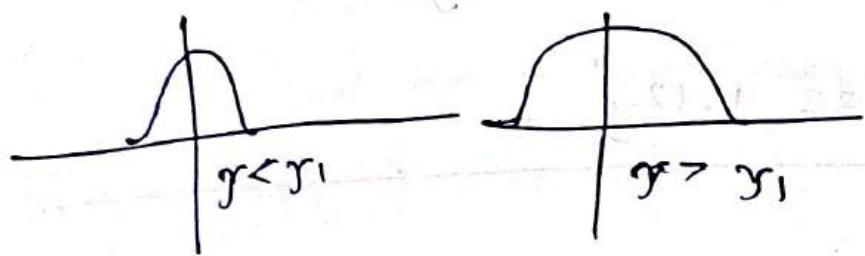
$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \epsilon = 10^{-5}$$

Batch Normalisation $y_i = \gamma \hat{x}_i + \beta$

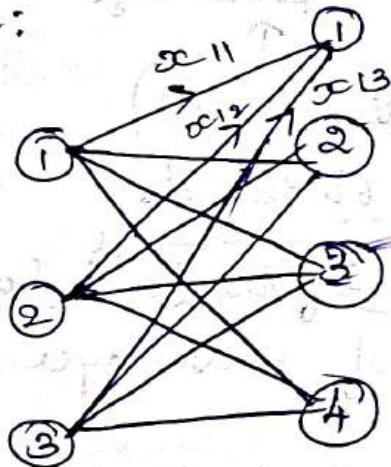


As β increases, the distribution shifts towards right.



As γ increases, (or spread) variance will increase.

Consider:



$$x = \begin{cases} x_{11} & x_{21} & x_{31} & x_{41} \\ x_{12} & x_{22} & x_{32} & x_{42} \\ x_{13} & x_{23} & x_{32} & x_{43} \end{cases}$$

$$x_{11} = 4$$

$x_{12} = 5$ for 1 neuron in hidden layer
 $x_{13} = 8$

$$S \rightarrow x = \{4, 5, 8\}$$

$$\hat{x}_i = ?$$

$$\bar{x}_x = \frac{4+5+8}{3} = \frac{17}{3} = 5.67$$

$$\sigma^2 = \frac{(4-5.67)^2 + (5-5.67)^2 + (8-5.67)^2}{3}$$

$$= \frac{(1.67)^2 + (0.67)^2 + (2.33)^2}{3} = \frac{8.66}{3}$$

$$= 2.89$$

$$\sigma = 1.7$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\hat{x}_i = \left\{ \frac{4-5.6}{\sqrt{2.89}}, \frac{5-5.6}{\sqrt{2.89}}, \frac{8-5.6}{\sqrt{2.89}} \right\}$$

$$\hat{x}_i = \{-0.8, -0.32, 1.12\}$$

Depthwise convolution:

$$x \in \mathbb{R}^{B \times H \times W \times C}$$

$$\text{image} = H \times W \times 32$$

$$\text{Want to do } \text{Conv2D}(64, (3, 3))$$

kernel size

\uparrow

with 64 kernels

Depth of this $\overset{\text{kernel}}{\text{with 64}}$

that of prev. layer = 32

Each kernel has size $3 \times 3 \times 32$ weights of 32 depth.

$$\therefore \text{Parameters} = (3 \times 3 \times 32) \times 64$$

$$= 18432$$

for 1 kernels



$$= (K \times K \times C_i) C_o$$

no. of feature maps in per layer
 no. of feature maps we seek to apply

$C_i \rightarrow \text{in channels}$

$C_o \rightarrow \text{o/p channels}$

To avoid 18432 parameters, we move to depthwise convolutions.

Depthwise

kernel $\rightarrow K \times K$ (No depth)
one kernel for each channel
 $= K \times K \times 1$

here, we don't decide no. of kernels.

feature map $X = H \times W \times 32 \rightarrow$ so directly 32 kernels will be applied.

Also, we don't introduce the parameter no. of output channels.

$\therefore \text{Parameters} = [K \times K \times C_{in}]$

C_{in} here = 32 (No. of kernels)

$$= 3 \times 3 \times 32$$

$$= 288$$

Therefore, no. of parameters reduced.

Separate multiplication with each feature map.

Advantages: less no. of parameters

Not taking cross-correlation b/w channels.

Disadvantages \rightarrow Intermix of channel information not

considered.

1st channel
with 1st channel
2nd channel
with 2nd channel
so on.

Normal convolution disadvantage \rightarrow more no. of parameters \rightarrow High interchannel interactions

Separable convolutions:

\rightarrow Depthwise convolution + pointwise convolution

\rightarrow Incorporates cross channel interaction.

Pointwise convolution

$|X| \rightarrow$ Kernel size (1, 1) with depth = previous feature map/
No. of channels in

Consider an image $H \times W \times C_{in}$

$$= 64 \times 64 \times 32$$

↓ separable convolution with 64 filters

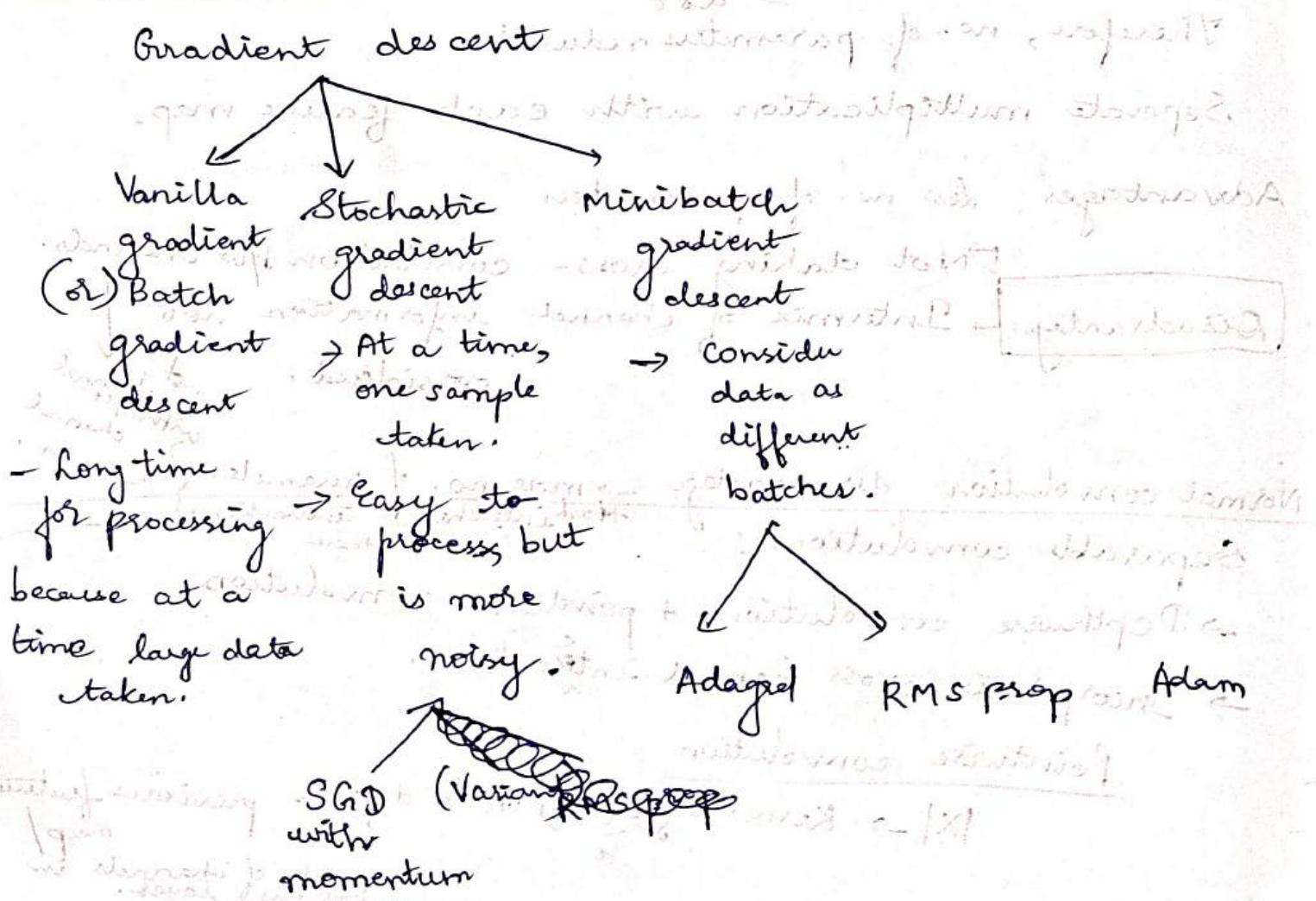
(kernel size) $\times C_{in}$ 1×1 kernel, Depth of previous layer
 $3 \times 3 \times 32$ + $1 \times 1 \times 32 \times 64$ $\xrightarrow{64 \text{ filters needed}}$

Depthwise Pointwise
 $= 288 + 2048$

$$= (K \times K \times C_{in}) + (1 \times 1 \times C_{in} \times C_{out})$$

which results in output depth (previous no. of channels) of $H \times W \times C_{out}$ (same padding)

Reduces parameter count.

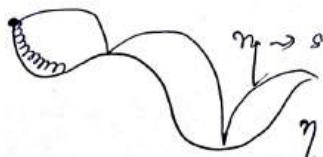


$\times W \times C_{in}$
 $\times 64 \times 32$
 separable convolution with 64 filters
 kernel size $\times C_{in}$
 $\times 3 \times 32$ + 1×1 kernel depth \times layer
 Depthwise $288 + 2048$
 $n) + (1 \times 1 \times C_{in} \times C_{out})$
 (previous depth \times no. of channels)
 in output of $H \times W \times C_{out}$
 (same padding)
 data count.

minibatch gradient descent
 consider data as different batches.

Adam RMS prop Adam
 momentum

Consider



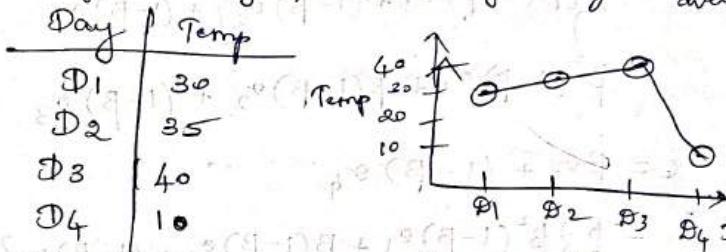
$\eta \rightarrow$ small it may reach local minima not global minima

$\eta \rightarrow$ large it may diverge.

saddle point \rightarrow where slope is zero.
(gradient)

No weight updation.
Due to this, we introduce momentum.

Exponential weighted average / Exponentially moving weighted average



there is a large dip.
how will models learn?

Therefore, we give weights to each day.

$$v_t = \beta v_{t-1} + (1 - \beta) e_t$$

β = 0.9 $v_{t-1} = 0$

previous temperature

$$v_3 = 0.9(6.2) + (1 - 0.9)40 = 5.58 + 0.1 \times 40 = 5.58 + 4 = 9.58$$

no sudden fall because previous day given

$$v_1 = 0.9(0) + (1 - 0.9)(30) = 0.1 \times 30 = 3$$

$$v_4 = 0.9(9.58) + (1 - 0.9)10 = 8.622 + 1$$

$$v_2 = 0.9(3) + (1 - 0.9)(35) = 2.7 + 0.1 \times 35 = 6.2$$

$$= 9.622$$

Algorithm:-

$$v_t = \beta v_{t-1} + (1-\beta) \sigma_t$$

$$v_{t-1} = 0$$

$$v_1 = (1-\beta) \sigma_1$$

$$v_2 = \beta v_1 + (1-\beta) \sigma_2$$

$$v_2 = \beta(1-\beta) \sigma_1 + (1-\beta) \sigma_2$$

$$v_3 = \beta v_2 + (1-\beta) \sigma_3$$

$$= \beta \{ \beta(1-\beta) \sigma_1 + (1-\beta) \sigma_2 \} + (1-\beta) \sigma_3$$

$$= \beta^2(1-\beta) \sigma_1 + \beta(1-\beta) \sigma_2 + (1-\beta) \sigma_3$$

$$v_4 = \beta v_3 + (1-\beta) \sigma_4$$

$$= \beta \{ \beta^2(1-\beta) \sigma_1 + \beta(1-\beta) \sigma_2 + (1-\beta) \sigma_3 \} + (1-\beta) \sigma_4$$

$$= \beta^3(1-\beta) \sigma_1 + \beta^2(1-\beta) \sigma_2 + \beta(1-\beta) \sigma_3 +$$

$$\beta(1-\beta) \sigma_4$$

$$v_t = (1-\beta) \{ \beta^3 \sigma_1 + \beta^2 \sigma_2 + \beta \sigma_3 + \sigma_4 \}$$

$$\text{Therefore, } v_t = (1-\beta) \{ \beta^{t-1} \sigma_1 + \beta^{t-2} \sigma_2 + \dots +$$

$$\beta \sigma_t \}$$

The no. of previous values we shall consider is given by
the formula $\rightarrow \frac{1}{1-\beta}$ (depends on the value of β)

$$\beta \text{ small } \rightarrow \beta = 0.5 \rightarrow \frac{1}{1-\frac{1}{2}} = 2$$

$$\beta \text{ large } \rightarrow \beta = 0.9 \rightarrow \frac{1}{1-0.9} = \frac{1}{0.1} = 10$$

~~start~~

iLoc[0] \rightarrow first element

How to give value of α ? we incorporate previous gradients to ensure proper convergence & not diverging the model with varying weight updation.

Prevent
large gradient \rightarrow large
case weight

Generate time series data & plot the variations with EWMA.