

16/08/22

Assessment Methodology :

CT1 - 20

CT2 - 20

Programming Assignment - 20 (Group of 2 or 3)

End Sem - 40

Machine Learning :

→ Training a machine to think like humans

2 phases :

- * Training - used to help the system learn
- * Testing - used to test if the system has learnt useful features

If machine recognizes the digit '0' as '6' or '9'
→ it has not learnt to distinguish 0, 6 or

0, 9

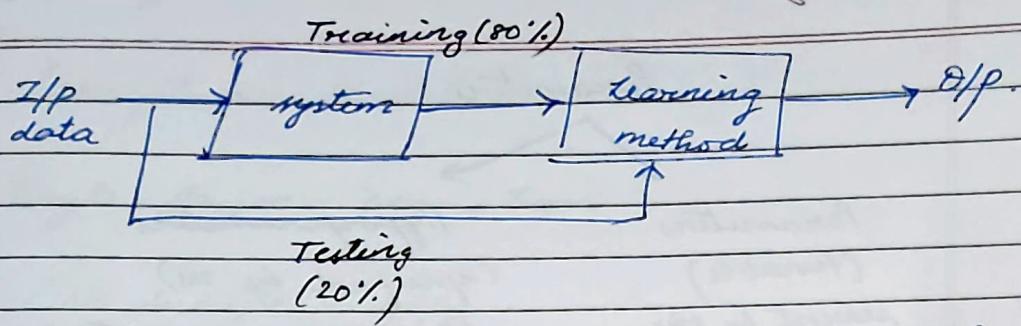
→ retrain the model with more examples.

We need to extract ~~useful~~ meaningful features from the data

Earlier → expert system (Rule based system)
we provide the rules to recognize sth
(Eg) : 7 → has a horizontal and vertical line connected at 1 point

↙

then ↗ will also be recognized as 7
so we need to give clear & explicit rules
⇒ Difficult and complex



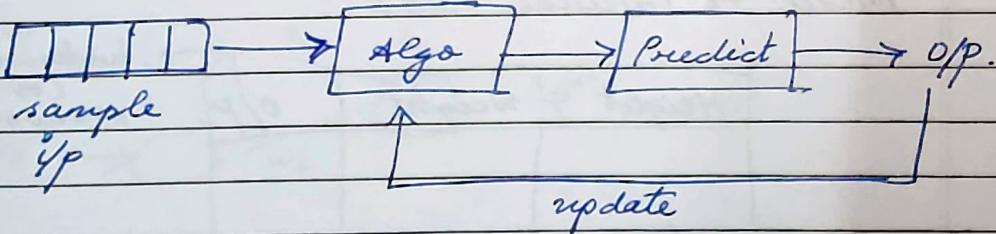
- ✓ while training, the system comes up with a learning method using the train features
- ✓ while testing, we directly get the o/p of for the test set from the method learnt by the system

If test set accuracy is ↓, we need to change the learning method

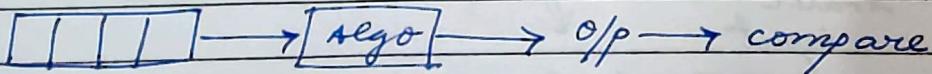
Train set - contains both data (features) & label

Test set - contains only the data (features)

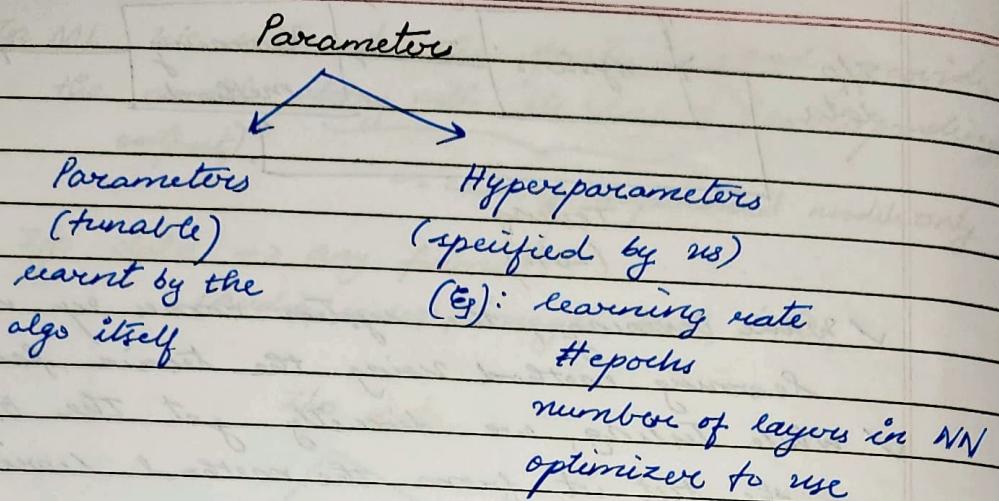
Training :



Testing :

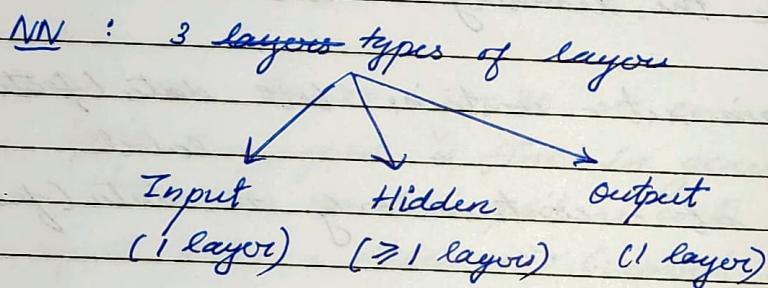


Algo → depends on the dataset and the parameters we use



Parameters → also affect accuracy

19/08/22



Model to calculate BMI :

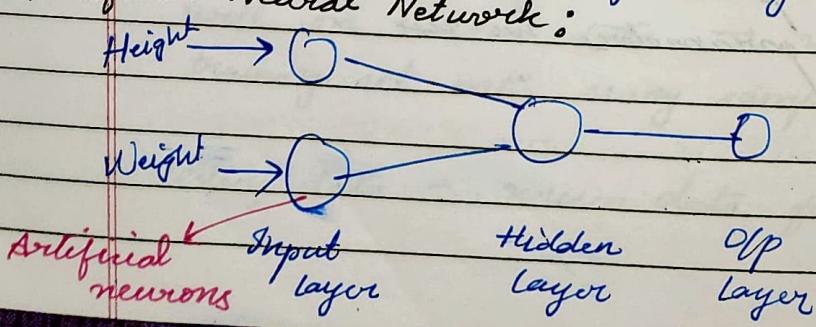
Height	Weight	O/P.	underweight (or) overweight
static input			

Input layer has 2 neurons

→ 1 for height attribute

another for weight attribute

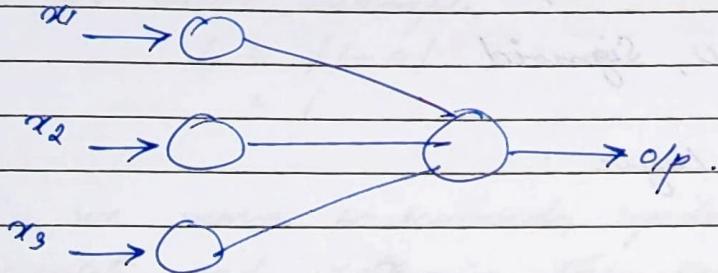
Artificial Neural Network:



Perception Model :

1. Single layer perception :

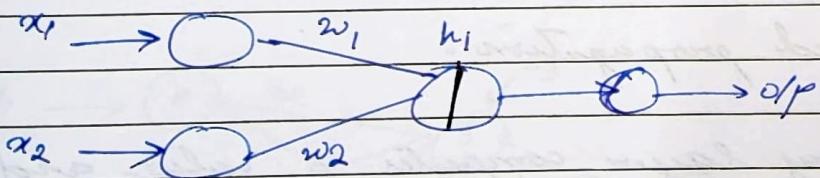
✓ No hidden layer
↙



Weights :

Responsible for transforming the info
layer layer by layer

$$\text{Output (Input)}_{\text{layer}} = \text{Input (hidden)}_{\text{layer}}$$



2 processes take place in a neuron

- ① linear operation on inputs (summing the product of inputs and weights)
- $$h_1 = \sum w_i x_i = w_1 x_1 + w_2 x_2$$

w_i are randomly initialized

if all $w_i = 0$ then our model always computes 0 only

\Rightarrow To avoid this we use bias.

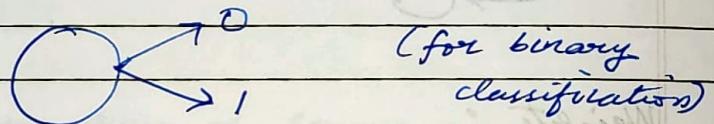
$$\therefore h_i = \sum w_i x_i + b$$

② apply activation function

$$y = f(h) \quad \begin{matrix} \rightarrow \text{Brings the values in the} \\ \downarrow \text{activation fn} \quad \text{range 0 to 1 (always!)} \end{matrix}$$

(Eg): ReLU, Sigmoid

Output layer :

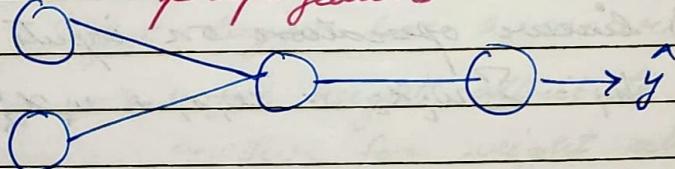


If values are $\begin{pmatrix} 0.75 \\ 0.3 \end{pmatrix} \Rightarrow$ class 0 is more probable
 $(\geq 0.5 \Rightarrow 1)$
 $(< 0.5 \Rightarrow 0)$.

Forward propagation :

- Every layer computes a value and forwards it to the next layer

forward propagation



If original label $y=0$ but predicted label $\hat{y}=1$ then we have to retrain the model by updating the weights (bias, inputs are constants!)

loss function :

Calculates the loss = difference in predicted and actual value.

$$\textcircled{1} \quad L = |y - \hat{y}|$$

For above example

$$L = |1 - 0|$$

$$= 1.$$

$\Rightarrow 100\%$ error?

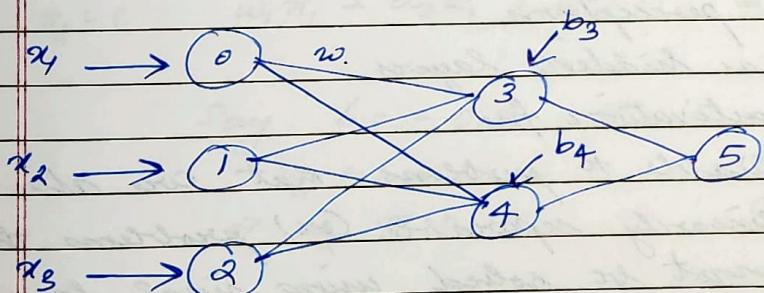
So we move backwards, update all the weights and retrain the model.

Optimizers :

Responsible for updating weights

(Eg): Adam, RMSprop, SGD (Stochastic Gradient Descent).

2-layer perceptron



bias - associated with each node

Multilayer - lower chances of error ~~less~~
of error

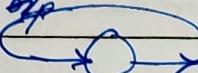
Topologies

Complete Recurrent

forward

backward

(every node in (there are
1 layer is
connected to
every node
in the
next layer).
self weights
→ self



1 forward + 1 backward prop → 1 epoch

22/08/22

ANN with > 1 hidden layer → DNN

Single layer perceptron :-

- i) No hidden layer
- ii) So, no activation function
- iii) applicable to problems that are linearly separable.

Multi layer perceptron :-

- i) It has hidden layers
- ii) Has activation fn
- iii) applicable to problems that are also non-linearly separable (or) problems that cannot be solved using single layer perceptron model

McCulloch - Pitt Algorithm (for single layer perceptron model).

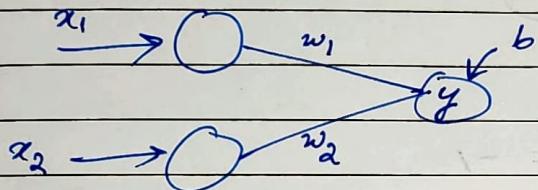
$$\text{op} \begin{cases} 1 & \text{if } w_0x + b > 0 \\ 0 & \text{if } w_0x + b \leq 0. \end{cases}$$

1. Logical AND

x_1	x_2	y	
0	0	0	$w_1x_1 + w_2x_2 + b = 0$
0	1	0	$w_1x_1 + w_2x_2 + b = 0$
1	0	0	$w_1x_1 + w_2x_2 + b = 0$
1	1	1	$w_1x_1 + w_2x_2 + b = 1.$

using
Hebb's
rule

Model



Using McCulloch-Pitt Algo

$$x_1=0, x_2=0 \quad w_1x_1 + w_2x_2 + b \leq 0 \Rightarrow b \leq 0$$

$$x_1=0, x_2=1 \quad w_1x_1 + w_2x_2 + b \leq 0 \Rightarrow w_2 + b \leq 0$$

$$x_1=1, x_2=0 \quad w_1x_1 + w_2x_2 + b \leq 0 \Rightarrow w_1 + b \leq 0$$

$$x_1=1, x_2=1 \quad w_1x_1 + w_2x_2 + b > 0 \Rightarrow w_1 + w_2 + b > 0$$

Say $b = -0.5$, then

$$w_2 + b \leq 0$$

$$\Rightarrow w_2 \leq 0.5$$

$$\text{Say } w_2 = 0.4.$$

$$w_1 + b \leq 0$$

$$\Rightarrow w_1 \leq 0.5$$

$$\text{Say } w_1 = 0.4.$$

$$w_1 + w_2 + b > 0$$

$$\Rightarrow w_1 + w_2 > 0.5.$$

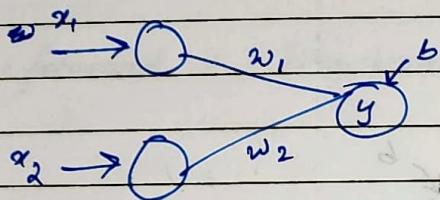
check:

$$0.4 + 0.4 - 0.5$$

$$= 0.3 > 0 \quad \checkmark$$

2. Logical OR

x_1	x_2	y		
0	0	0	$w_1x_1 + w_2x_2 + b = 0 \leq 0$	$b \leq 0$
0	1	1	$w_1x_1 + w_2x_2 + b = 1 > 0$	$w_2 + b > 0$
1	0	1	$w_1x_1 + w_2x_2 + b = 1 > 0$	$w_1 + b > 0$
1	1	1	$w_1x_1 + w_2x_2 + b = 1 > 0$	$w_1 + w_2 + b > 0$



$$\text{Say } b = -1.$$

$$w_2 + b > 0$$

$$w_2 > 1 \Rightarrow \text{Say } w_2 = 2$$

$$w_1 + b > 0$$

$$w_1 > 1 \Rightarrow \text{Say } w_1 = 2$$

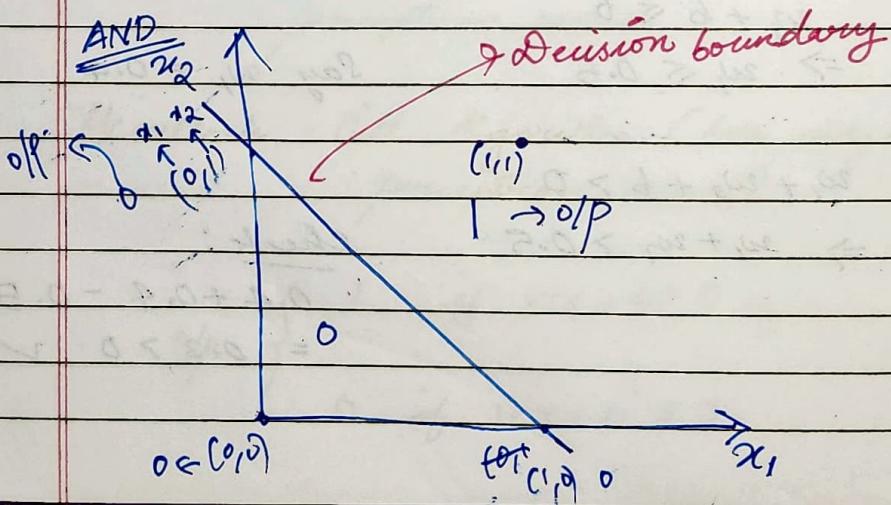
$$w_1 + w_2 + b > 0$$

$$2 + 2 - 1 > 0$$

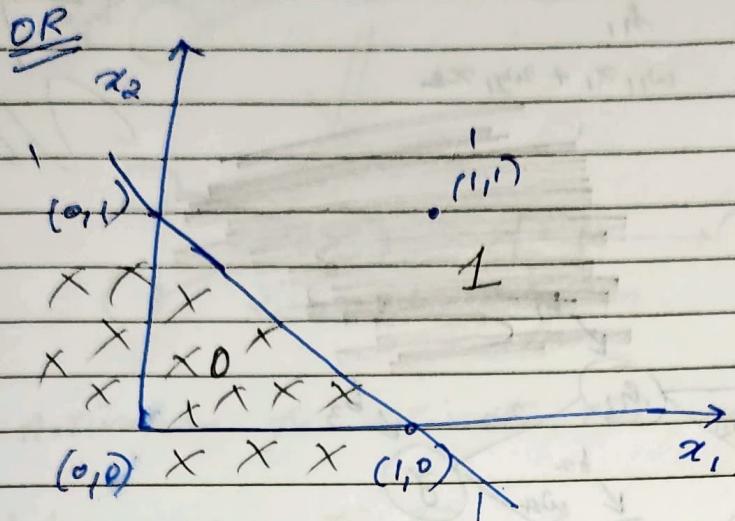
$3 > 0$ is true

\therefore One possible solution is

$$b = -1, w_1 = 2, w_2 = 2$$

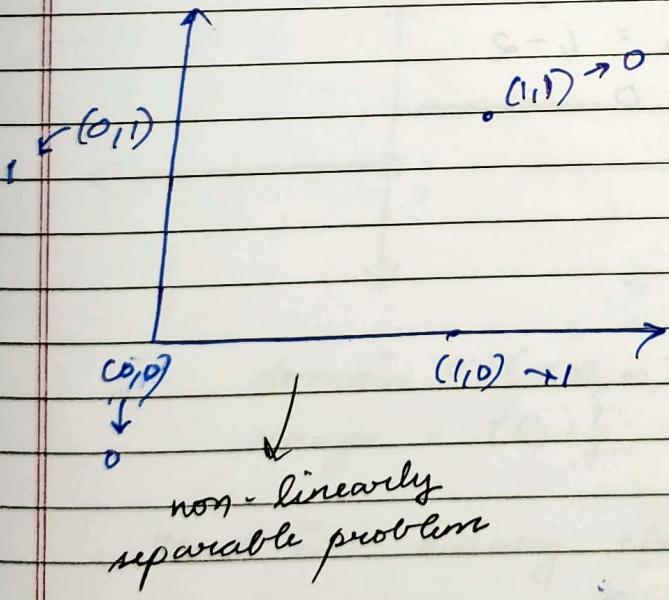


If decision boundary is a line, then problem start at hand is linearly separable.



3. XOR

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



→ No linear boundary possible

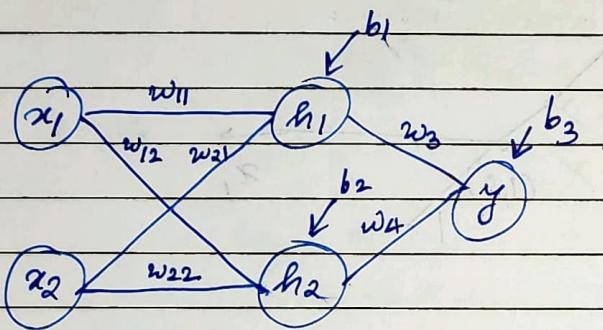
so XOR is not linearly separable & hence cannot be implemented using single layer perceptron model

non-linearly separable problem

Multi-layer perceptron model

1. XOR

x_1	x_2	y	h_1
0	0	0	$w_{11}x_1 + w_{21}x_2$
0	1	1	
1	0	1	
1	1	0	



$$h_1 \rightarrow w_{11}x_1 + w_{21}x_2 + b_1$$

$$h_2 \rightarrow w_{12}x_1 + w_{22}x_2 + b_2$$

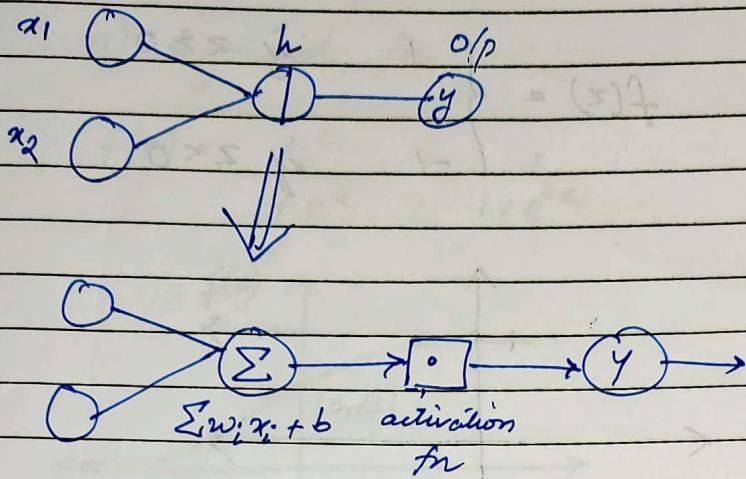
$$\text{suppose } w_{11} = w_{12} = w_{21} = w_{22} = 1$$

$$b_1, b_2 = 0, -1$$

$$w_3, w_4 = 1, -2$$

$$b_3 = 0$$

22/08/22
Activation functions



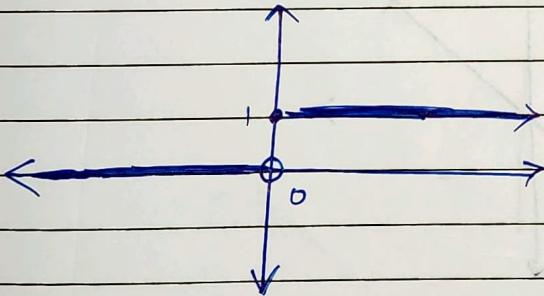
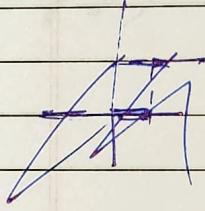
Activation fn activates the neurons.

2 types of activation fns:

1. Linear fn

a) Step function

$$f(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0. \end{cases}$$



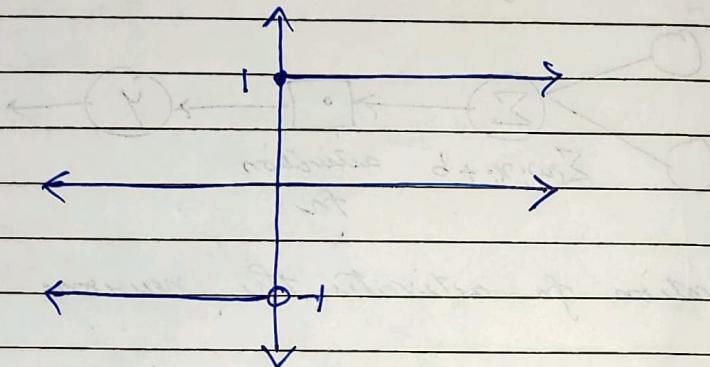
$$\text{domain} = (-\infty, \infty)$$

$$\text{range} = \{0, 1\}.$$

Used in binary classification.

b) sigmoid function :

$$f(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

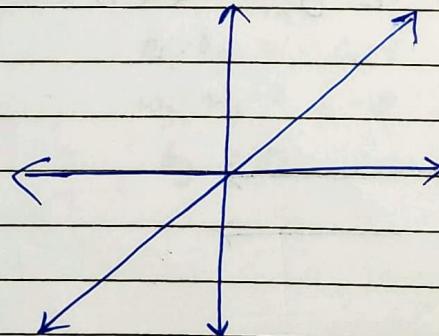


$$\text{domain} = (-\infty, \infty)$$

$$\text{range} = \{-1, 1\}.$$

c) linear function

$$f(z) = z$$



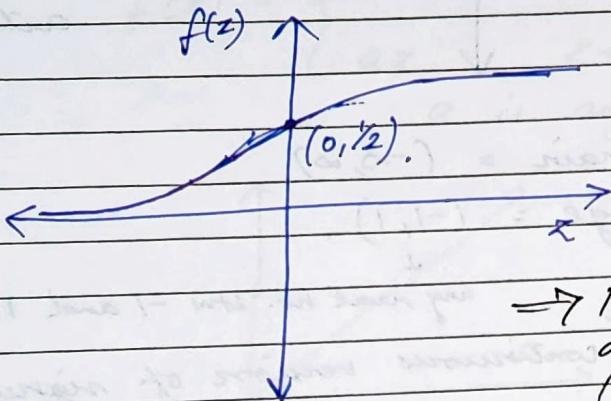
$$\text{domain} = (-\infty, \infty)$$

$$\text{range} = (-\infty, \infty).$$

2. On Non-linear function

a) Sigmoid fn

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



\rightarrow Non-zero centred activation fn
(does not touch 0).

domain = $(-\infty, \infty)$

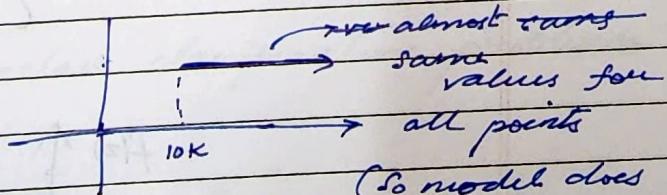
range = $[0, 1]$.

(continuous version
of step fn)

any real no. b/w 0 & 1.

Problem: Local convergence

say $z \in [10000, \infty)$.

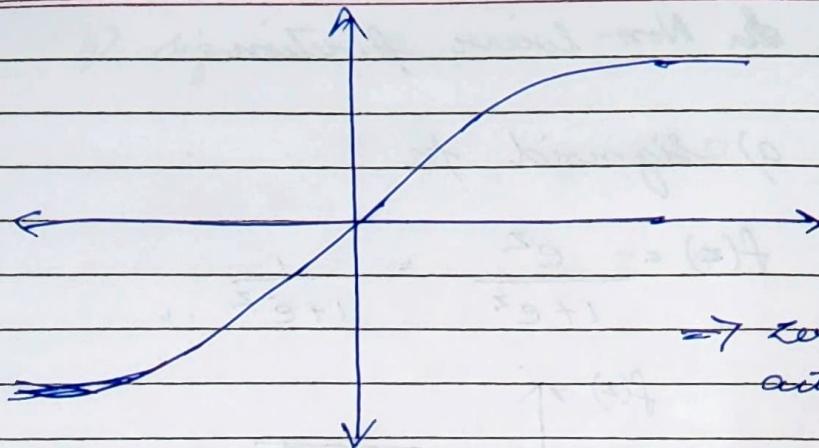


b) tanh function

(or) hyperbolic tangent function

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

so we use tanh



\Rightarrow zero-centered activation fn

domain = $(-\infty, \infty)$

range = $(-1, 1)$.

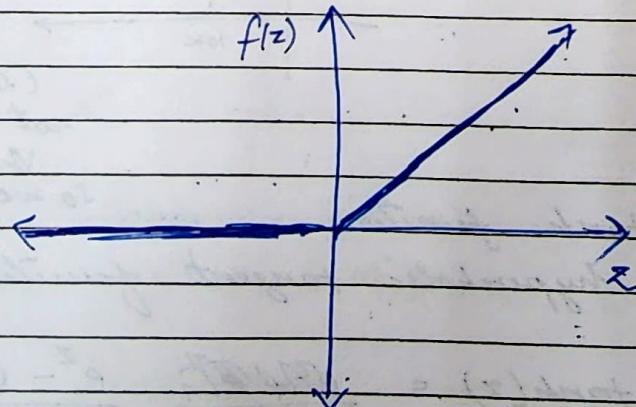
any real no. b/w -1 and 1.

(continuous version of signum).

c) ReLU (Rectified linear unit).

tanh, sigmoid graphs are non linear.
So we introduce linearity by using the
full-fn

$$f(z) = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$



domain = $(-\infty, \infty)$

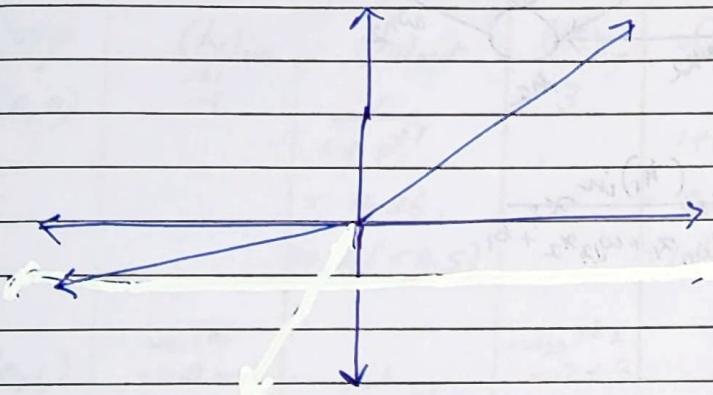
range = $[0, \infty)$.

d) Leaky ReLU

If ReLU does not give good performance we can use leaky ReLU.

$$f(z) = \begin{cases} z & , z \geq 0 \\ az & , z < 0 \end{cases}$$

a is any constant



$$\text{domain} = (-\infty, \infty)$$

$$\text{range} = (-\infty, \infty)$$

Softmax activation function

✓ Used for multi-class classification problems

$$f(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

where z_1, z_2, \dots, z_n are the n outputs

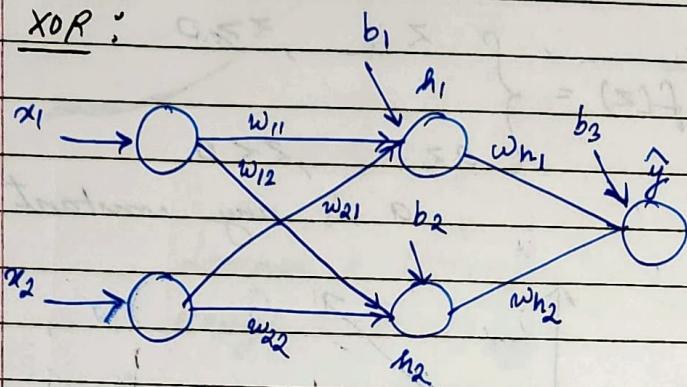
$$f(z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

⋮

$$f(z_n) = \frac{e^{z_n}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

$$f(z_i) = \frac{e^{z_i}}{\sum_{k=1}^n e^{z_k}}$$

26/08/22

XOR :

$$HP = 0/P \quad (h_1)_{in}$$

$$x_1 \quad x_2 \quad y \quad w_{11}x_1 + w_{21}x_2 + b_1$$

$$0 \quad 0 \quad 0$$

$$0 \quad 1 \quad 1$$

$$1 \quad 0 \quad 1$$

$$1 \quad 1 \quad 0$$

$$(h_1)_{in} = w_{11}x_1 + w_{21}x_2 + b_1$$

$$(h_1)_{out} = Act((h_1)_{in})$$

$$(h_2)_{in} = w_{12}x_1 + w_{22}x_2 + b_2$$

$$(h_2)_{out} = Act((h_2)_{in})$$

$$y_{in} = w_{h_1}(h_1)_{out} + w_{h_2}(h_2)_{out} + b_3$$

$$y_{out} = Act(y_{in})$$

Given $w_{11} = w_{21} = 2$

$$w_{12} = w_{22} = -2$$

$$[b_1, b_2] = [-1, 3]$$

$$w_{h_1} = w_{h_2} = 2$$

$$b_3 = -3.$$

Use sigmoid activation fn and solve

$$\begin{array}{c} \text{Input} & (h_1)_{\text{in}} & (h_1)_{\text{out}} & (h_2)_{\text{in}} & (h_2)_{\text{out}} \\ \rightarrow & t_0, 0 & b_1 & \frac{1}{1+e^{-b_1}} & \end{array}$$

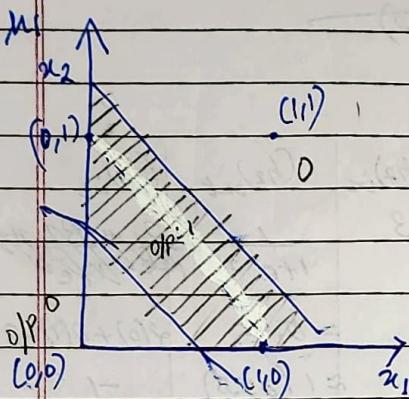
input intercept $t_0, 0$	$(h_1)_{\text{in}}$ $= b_1$ -1	$(h_1)_{\text{out}}$ $\frac{1}{1+e^{-1}}$ $= 0.26$ $\approx 0 (< 0.5)$	$(h_2)_{\text{in}}$ 3	$(h_2)_{\text{out}}$ $\frac{1}{1+e^{-3}}$ $= 0.95$ $\approx 1 (\geq 0.5)$	$(y)_{\text{in}}$ $2(0) + 2(1) - 3$ -1	y_{out} $\frac{1}{1+e^1}$ $= 0.26$ $0.$
$(0, 1)$	$\frac{w_{21} + b_1}{-2 + 1}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1	$\frac{w_{22} + b_2}{-2 + 3}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1	$\frac{2(1) + 2(1) - 3}{-1}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1

$(0, 1)$	$\frac{w_{21} + b_1}{-2 + 1}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1	$\frac{w_{22} + b_2}{-2 + 3}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1	$\frac{2(1) + 2(1) - 3}{-1}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1
----------	--	---	--	---	---------------------------------------	---

$(1, 0)$	$\frac{w_{21} + b_1}{-2 + 1}$ $= -1$	$\frac{1}{1+e^{-3}}$ $= 0.75$ ≈ 1	$\frac{w_{22} + b_2}{-2 + 3}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1	$\frac{2(1) + 2(1) - 3}{-1}$ $= 1$	$\frac{1}{1+e^{-1}}$ $= 0.75$ ≈ 1
----------	---	---	--	---	---------------------------------------	---

$(1, 1)$	$\frac{w_{21} + b_1}{-2 + 1}$ $= 2 + (-1) - 1$ $= 0$	$\frac{1}{1+e^{-3}}$ $= 0.95$ ≈ 1	$\frac{w_{22} + b_2}{-2 + 3}$ $= -1$	$\frac{1}{1+e^{-1}}$ $= 0.26$ ≈ 0	$\frac{2(1) + 2(1) - 3}{-1}$ $= -1$	$\frac{1}{1+e^{-1}}$ $= 0.26$ ≈ 0
----------	--	---	---	---	--	---

$(h_1)_{out}$	$(h_2)_{out}$	y	
0,0	0	1	0
0,1	1	1	1
1,0	1	1	1
1,1	1	0	0
		↓	↓
		OR	NAND
			$AND(h_1, h_2)$



Given $w_1 = w_{12} = w_{21} = w_{22} = 1$

$b_1, b_2 = 0, -1$

$w_{h_1}, w_{h_2} = 1, -2$

$b_3 = 0$.

ReLU activation

	$w_{11}x_1 + w_{21}x_2 + b_1$	$w_{12}x_1 + w_{22}x_2 + b_2$	$w_{h_1}h_1 + w_{h_2}h_2 + b_3$	
(0,0)	0	0	-1	0
(0,1)	$1+0=1$	1	$1+0=1$	1
(1,0)	$1+0=1$	1	$1-1=0$	0
(1,1)	$1+1+0=2$	2	$1+1-1=1$	1
			$2-2+0=0$	0

$$(h_1)_{\text{out}} \quad (h_2)_{\text{out}}$$

0,0	0	0
0,1	1	0
1,0	1	0
1,1	2	1

29/08/22

 y - ground truth \hat{y} - predicted value

loss function :

1. Mean Squared Error :

$$\text{MSE} = \frac{1}{m} \sum (y_i - \hat{y}_i)^2$$

 m - number of samples y_i - original output \hat{y}_i - predicted output

2. Mean Absolute Error :

$$\text{MAE} = \frac{1}{m} \sum |y_i - \hat{y}_i|$$

Linear regression

$$y = mx + c$$

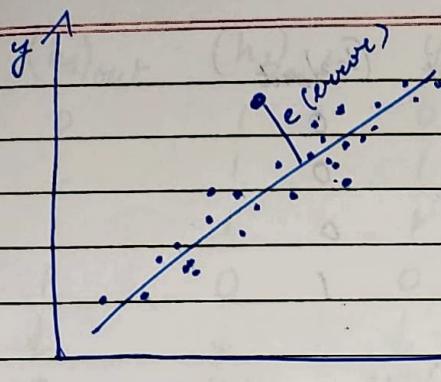
or

$$y = \alpha + \beta x$$

In DL

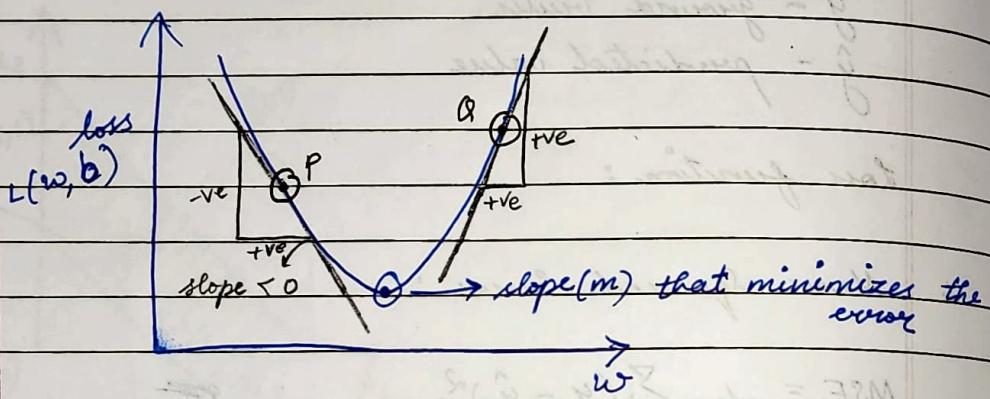
$$y = \sum w_i x_i + b$$

\downarrow \downarrow \downarrow
 w x b



$$e = (y - \hat{y})^2$$

$$= [y - (m_2 + c)]^2$$



Gradient Descent

It is an iterative optimization algorithm to find the minimum of a function.

$$\text{slope} = \frac{\Delta Y}{\Delta X}$$

$$w_{\text{new}} = w_{\text{old}} - \gamma \frac{\partial L}{\partial w}$$

γ - learning rate

for point P

\rightarrow -ve slope

$$w_{\text{new}} = w_{\text{old}} - \gamma \left(-\frac{\partial L}{\partial w} \right)$$

$$= w_{\text{old}} + \gamma \frac{\partial L}{\partial w}$$

\rightarrow magnitude of slope

For point Q

~~$$W_{\text{new}} = W_{\text{old}} - \gamma \frac{\partial L}{\partial W}$$~~

(We need to reduce W to reach & minimize error).

$$\hat{y} = mx + c$$

Error is L (Loss)

$$= \frac{1}{n} (y - \hat{y})^2$$

$$= \frac{1}{n} (y - (mx + c))^2$$

$$\frac{\partial L}{\partial m} = \frac{2}{n} (y - \hat{y})(-x) = -\cancel{2L} \frac{-2x(y - \hat{y})}{n}$$

$$\frac{\partial L}{\partial c} = \frac{2}{n} (y - \hat{y})(-1) = \cancel{2L} - \frac{2(y - \hat{y})}{n}$$

$$m_{\text{new}} = m_{\text{old}} - \gamma \frac{\partial L}{\partial m}$$

$$= m_{\text{old}} + \frac{2\gamma}{n} (y - (mx + c))$$

$$c_{\text{new}} = c_{\text{old}} - \gamma \frac{\partial L}{\partial c}$$

$$= c_{\text{old}} + \frac{2\gamma}{n} (y - (mx + c))$$

3. Binary Cross-entropy loss:

$$L = -\frac{1}{n} \sum [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

Used for classification problems

MSE

BCE

- ✓ for linear regression for logistic regression
classification problems
- ✓ does not have the has the ability to
ability to update weights update weights
- ✓ applicable only to simple applicable to complex
problems problems.

Multi-class classification

uses probabilistic distributions

3 classes

(Eg): cat - 0.9 → man cat - 1
 dog - 0.7 dog - 0
 rabbit - 0.5 rabbit - 0

↙
data point ∈ cat class

✓

inherently finds
probabilistic distribution
of predicted values

Variants of Gradient Descent :

1. Stochastic Gradient Descent

Calculate loss and update weights after for every data point

2. Mini batch GD

Divide dataset into batches and update weights for every batch

30/08/22

Binary Cross entropy Loss (BCE)

$$L = - (y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i))$$

	y	\hat{y}
1	0.95	
1	0.1	

Find MSE, BCE

y	\hat{y}	MSE	BCE
1	0.95	0.0025	$-1 \log 0.95 = 0.051$
1	0.1	0.81	$-1 \log 0.1 = 2.302$

point 1
is almost accurate
point 2 is
almost inaccurate

there is significant
difference b/w the
values allowing us to
penalize 2nd
prediction better

$$L = -[y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)]$$

1. $y=1, \hat{y}=0$

$$L = -y_i \log(\hat{y}_i) \Rightarrow \text{if target is } 1 \\ = -\log \hat{y}_i$$

2. $y=0, \hat{y}=1$

$$L = -(1-y_i) \log(1-\hat{y}_i) \Rightarrow \text{if target is } 0 \\ = -\log(1-\hat{y}_i)$$

Categorical Cross entropy

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

Backpropagation:

2 steps :

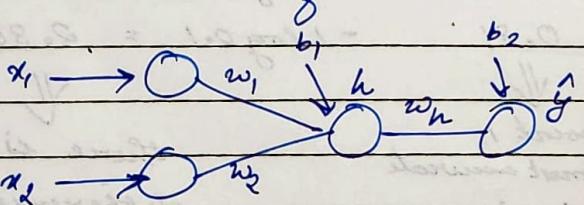
✓ weight updation

✓ chain rule of calculus

Weight updation :

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W_{\text{old}}}$$

Chain rule of Calculus



$$(1) h_{\text{in}} = w_1 x_1 + w_2 x_2 + b_1$$

$$h_{\text{out}} = \text{Act}(h_{\text{in}})$$

$$y_{in} = w_h \times h_{out} + b_2$$

$$y_{out} = \text{act}(y_{in})$$

This is the forward pass

We now calculate error

$$L = (y - \hat{y})^2$$

Now we backpropagate & update weights

$$(w_h)_{\text{new}} = (w_h)_{\text{old}} - \gamma \frac{\partial L}{\partial w_h}$$

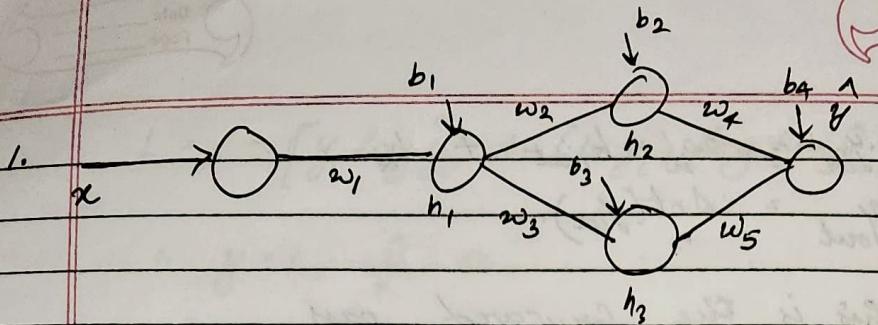
$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial y_{out}} \cdot \frac{\partial y_{out}}{\partial y_{in}} \cdot \frac{\partial y_{in}}{\partial w_h}$$

~~$(w_i)_{\text{new}}$~~ = ~~$\frac{\partial L}{\partial w_i}$~~

$$(w_i)_{\text{new}} = (w_i)_{\text{old}} - \gamma \frac{\partial L}{\partial w_i}$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial y_{out}} \times \frac{\partial y_{out}}{\partial y_{in}} \times \cancel{\frac{\partial y_{in}}{\partial w_i}} / \cancel{\frac{\partial L}{\partial w_i}}$$

$$\frac{\partial y_{in}}{\partial h_{out}} \times \frac{\partial h_{out}}{\partial h_{in}} \times \frac{\partial h_{in}}{\partial w_i}$$



$$(h_1)_{in} = w_1 x + b_1$$

$$(h_1)_{out} = \text{Act}(h_1)_{in})$$

$$(h_2)_{in} = w_2 (h_1)_{out} + b_2$$

$$(h_2)_{out} = \text{Act}(h_2)_{in})$$

$$(h_3)_{in} = w_3 (h_2)_{out} + b_3$$

$$(h_3)_{out} = \text{Act}(h_3)_{in})$$

$$y_{in} = w_4 (h_2)_{out} + w_5 (h_3)_{out} + b_4$$

$$y_{out} = \text{Act}(y_{in})$$

$$(w_2)_{new} = (w_2)_{old} - \eta \left(\frac{\partial L}{\partial y_{out}} \times \frac{\partial y_{out}}{\partial y_{in}} \times \frac{\partial y_{in}}{\partial w_2}_{old} \right)$$

$$(w_4)_{new} = (w_4)_{old} - \eta \left(\frac{\partial L}{\partial y_{out}} \cdot \frac{\partial y_{out}}{\partial y_{in}} \cdot \frac{\partial y_{in}}{\partial h_2}_{out} \cdot \frac{\partial h_2}{\partial h_2}_{in} \cdot \frac{\partial h_2}{\partial w_4}_{old} \right)$$

$$(w_1)_{new} = (w_1)_{old} - \eta \left[\frac{\partial L}{\partial y_{out}} \cdot \frac{\partial y_{out}}{\partial y_{in}} \cdot \frac{\partial y_{in}}{\partial h_2}_{out} \cdot \frac{\partial h_2}{\partial h_2}_{in} \cdot \frac{\partial h_2}{\partial w_1}_{old} + \right. \\ \left. \frac{\partial h_2}{\partial h_3}_{in} \cdot \frac{\partial h_3}{\partial h_3}_{out} \cdot \frac{\partial h_3}{\partial w_1}_{old} \right]$$

$$\left(\frac{\partial L}{\partial y_{out}} \cdot \frac{\partial y_{out}}{\partial y_{in}} \cdot \frac{\partial y_{in}}{\partial h_3}_{out} \cdot \frac{\partial h_3}{\partial h_3}_{in} \cdot \frac{\partial h_3}{\partial w_1}_{old} \cdot \frac{\partial h_1}{\partial h_3}_{out} \cdot \frac{\partial h_1}{\partial h_1}_{in} \cdot \frac{\partial h_1}{\partial w_1}_{old} \right]$$

(path via h_2 + path via h_3)

Vanishing Gradient Problem

If we use sigmoid activation for say

sigmoid \rightarrow 0 to 1

derivative \rightarrow range 0 to 0.25.

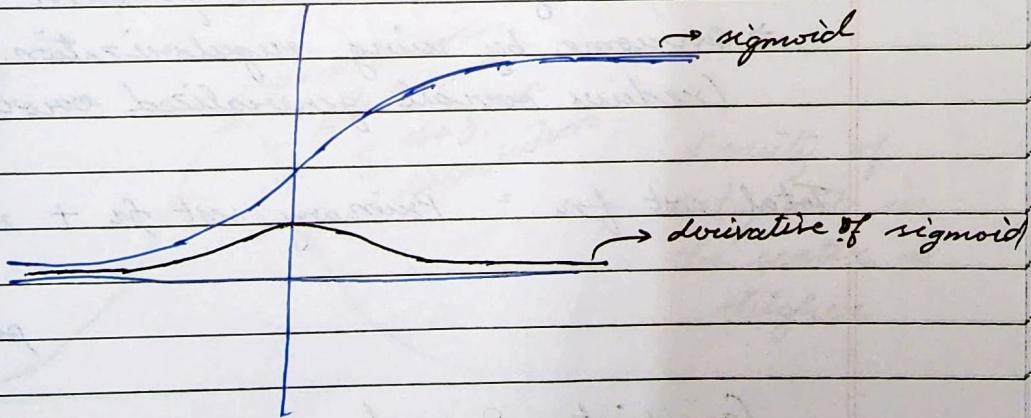
If we use a long chain rule (as we did for w_1) then

$0.25 \times 0.11 \times \dots \times 0.2$ (say) (product of all derivatives will be very small)

$$\Rightarrow (w_1)_{\text{new}} \cong (w_1)_{\text{old}}$$

\Rightarrow Model will not learn (we are stuck at a local minima)

\Rightarrow Vanishing gradient problem



02/09/22

Things to consider when creating a model

1. Cost fn (minimization problem \rightarrow min errors/loss)
2. output unit
3. hidden unit
4. architecture

Cost fn \rightarrow decided based on type of problem

Regression - statistical measures like

MSE, MAE

Classification - Binary entropy (binary classification)

Categorical entropy

(Multiclass classification)

Overfit :

Model performs well on train data
but not very well on test / unseen data

\Rightarrow overcome by using regularization.

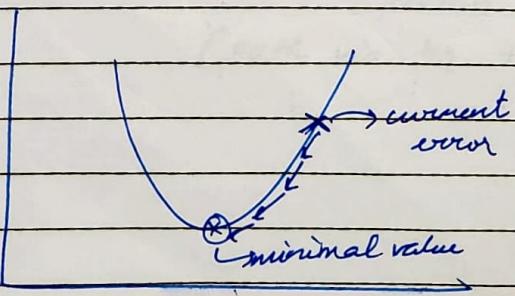
(reduces overall generalized error).

Total cost fn = Primary cost fn + regularization
 \downarrow
 penalty.

Gradient Descent

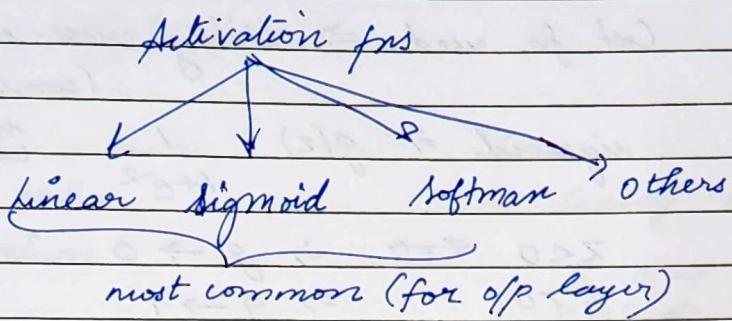
used for

✓ Reducing the error by updating weight



O/p activation need not be the same as hidden layer activations.

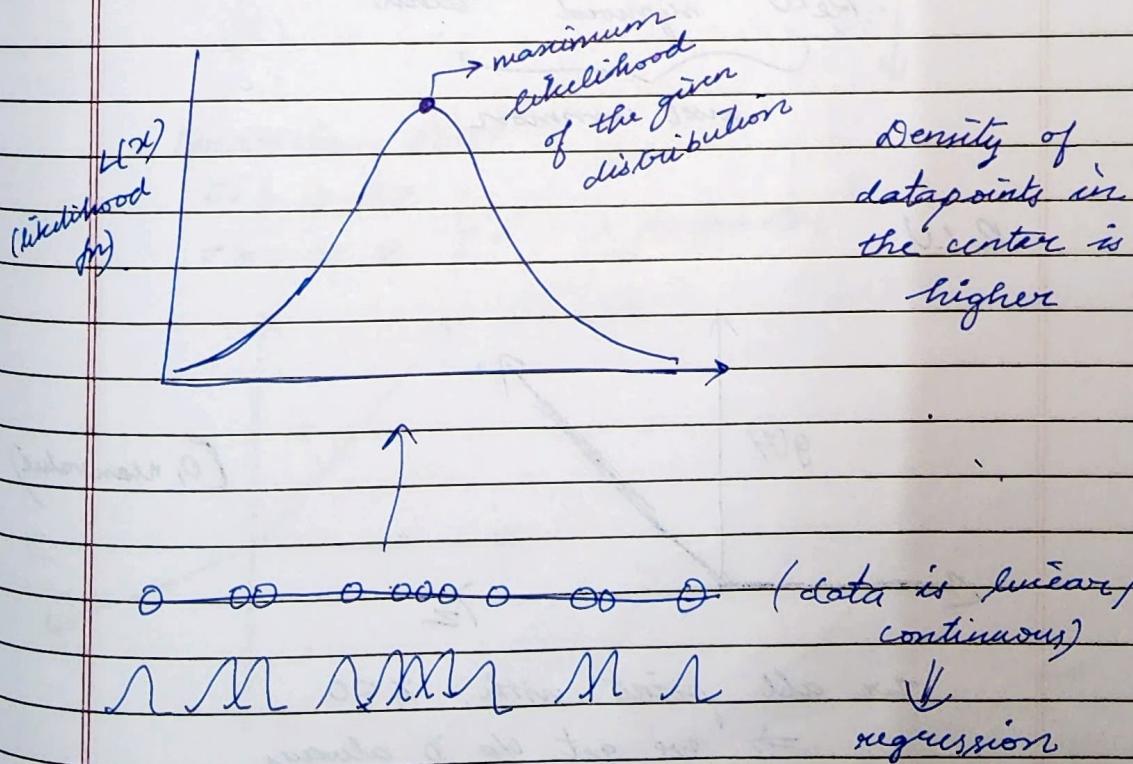
\Rightarrow mostly softmax used for o/p activation.



* Choosing activation fns (for o/p layer):

Depends on $P(Y/x, \theta) \Rightarrow$ distribution of o/p conditioned on inputs.

If the distribution is normal / gaussian
 \Rightarrow use linear activation



If data label is either 0 or 1 \Rightarrow use sigmoid activation
 \Downarrow
 binary classification.

Cost function \Rightarrow binary cross-entropy

$$\text{sigmoid} \Rightarrow g(z) = \frac{1}{1+e^{-z}}$$

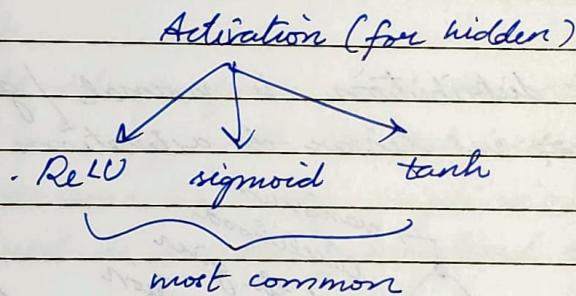
(consists of Prob of success + failure)
 (class 1) (class 0)

$$z < 0 \quad z = 0 \quad \Rightarrow g \rightarrow 0$$

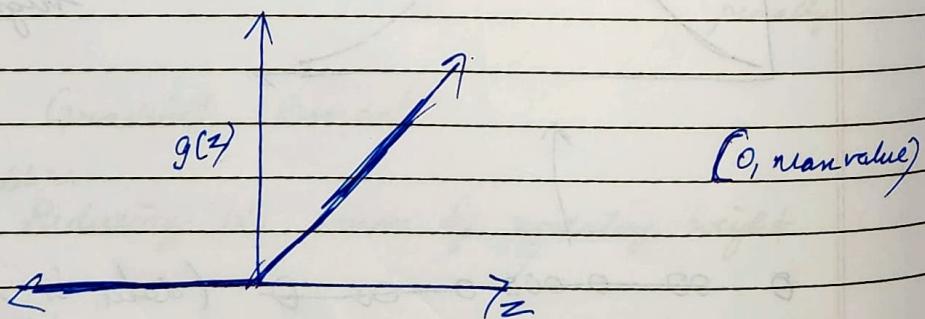
$$z > 0 \quad \Rightarrow g \rightarrow 1.$$

softmax for multiclass classification problem
 \Rightarrow most common

*



ReLU



For all points with $z < 0$

\Rightarrow we get $g(z) = 0$ always

\Rightarrow next layer always gets the $g(z) = 0$
 leads to saturation

so we use other variants like leaky ReLU

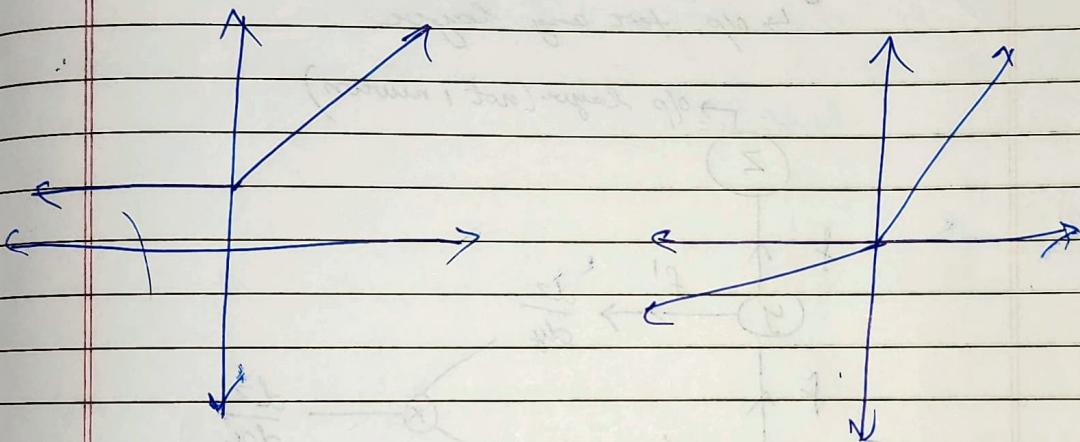
tanh $g(z) \rightarrow \text{sigmoid?}$ $\tanh \Rightarrow$

$$\varphi(g(\alpha z)) - 1$$

~~Variants
of ReLU~~Absolute value
rectification

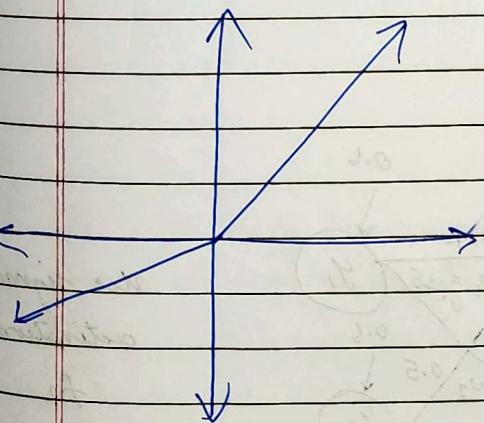
Leaky ReLU.

$$\begin{aligned} z < 0 &\Rightarrow 0.01z \\ z \geq 0 &\Rightarrow z \end{aligned}$$



Parametric ReLU

$$\begin{cases} z < 0 \Rightarrow az \\ z \geq 0 \Rightarrow z \end{cases} \quad \text{'a' is a parameter}$$



Margin

$$\max(w_1^T x + b_1, w_2^T x + b_2). \quad (\text{Are 2 sets of weights & biases?})$$

05/09/22

Initial weights - assigned randomly
 Update of weights - done using approaches/algos like Gradient Descent (Backpropagation).

2 steps:

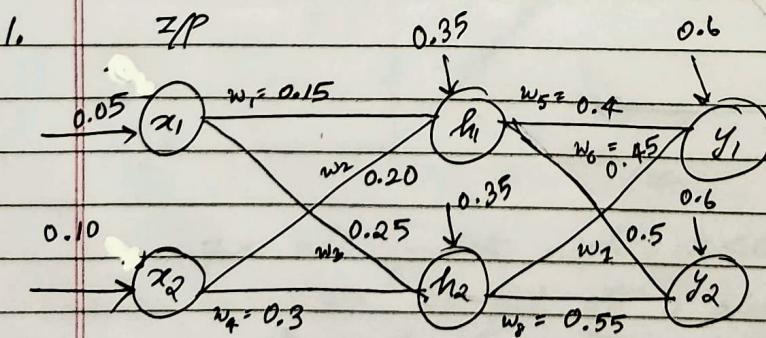
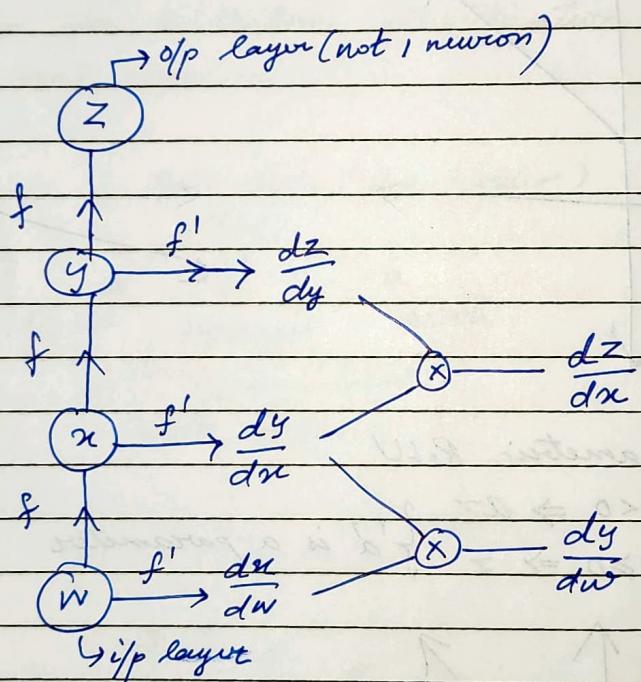
Forward propagation of weights

Backward propagation of errors.

Loop.

$$y = f(x, \theta)$$

↳ o/p for any layer



Target FP 2/4 (expected)

T

$$\begin{matrix} y_1 & y_2 \\ 0.01 & 0.99 \end{matrix}$$

$$(h_1)_{in} = (h_1)_{in} = w_1 x_1 + w_2 x_2 + b_1 \\ = 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35$$

$$(h_1)_{in} - t_0 =$$

$$(h_2)_{out} = \frac{1}{1+e^{-}}$$

We got $E_{\text{total}} = 0.2984$ (MSE).

\Rightarrow not good enough
so we backpropagate and update weights

* Stopping conditions (for training a model):

- ✓ If error in current iteration is much better than previous iterations (convergence occurs)
- ✓ If max iterations have reached
 - \Rightarrow will not give the best model (weights might still have to be updated).

06/09/22

Updation:

$$W = W - \alpha \Delta$$

\downarrow \hookrightarrow gradient
learning rate

$$E_{\text{total}} = \frac{1}{2} (T_1 - (Y_1)_{\text{out}})^2 + \frac{1}{2} (T_2 - (Y_2)_{\text{out}})^2$$

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_5} &= \frac{\partial E_{\text{total}}}{\partial T_1} \times \frac{\partial T_1}{\partial (Y_1)_{\text{out}}} \times \frac{\partial (Y_1)_{\text{out}}}{\partial w_5} \\ &= \frac{\partial E_{\text{total}}}{\partial (Y_1)_{\text{out}}} \times \frac{\partial (Y_1)_{\text{out}}}{\partial (Y_1)_{\text{in}}} \times \frac{\partial (Y_1)_{\text{in}}}{\partial w_5} \\ &= \left[\frac{1}{2} \times 2 \times (T_1 - (Y_1)_{\text{out}}) \times -1 \right] \left[\frac{-1}{(1 + e^{-y_1})^2} \cdot e^{-y_1} \cdot -1 \right] \\ &\quad \left[\infty (H_1)_{\text{out}} \right] - \\ &= + \left((Y_1)_{\text{out}} - T_1 \right) \left(\frac{e^{-y_1}}{(1 + e^{-y_1})^2} \right) (H_1)_{\text{out}} \end{aligned}$$

$$\alpha = 0.5 = -0.01$$

$$T_1 = 0.01$$

$$T_2 = 0.99$$

$$(y_1)_{\text{out}} = 0.7513$$

$$(H_1)_{\text{out}} = 0.592$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = ((y_1)_{\text{out}} - T_1) \left(\frac{1}{1+e^{-y_1}} \cdot \frac{e^{-y_1}}{1+e^{-y_1}} \right) (H_1)_{\text{out}}$$

$$= ((y_1)_{\text{out}} - T_1) / ((y_1)_{\text{out}} (1 - (y_1)_{\text{out}})) H_1_{\text{out}}$$

$$= (0.7513 - 0.01) (0.7513 \times (1 - 0.7513)) \times 0.592$$

$$= 0.08216.$$

$$(w_5)_{\text{new}} = (w_5)_{\text{old}} - \alpha \cdot \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$= 0.4 - 0.5 \times 0.08216$$

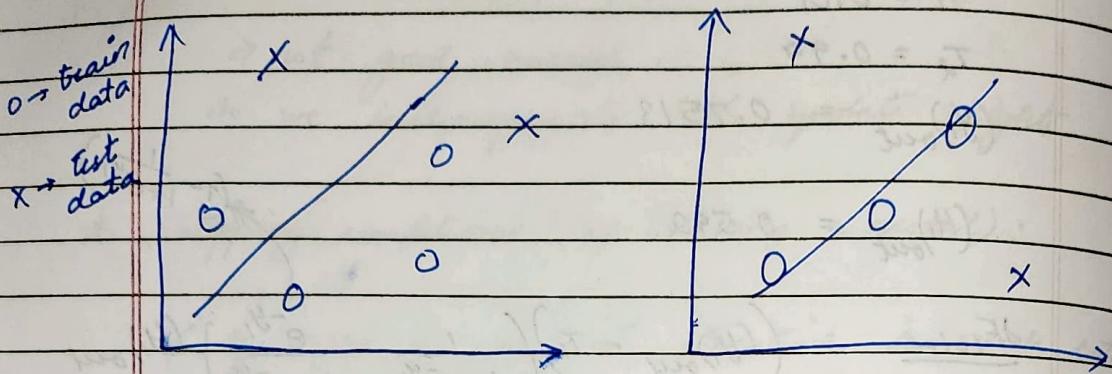
$$= 0.4 - 0.04108$$

$$= 0.35892$$

Repeat this for all w

After updation of weights, check the error
(by performing forward propagation) to see
if it is acceptable
Else repeat.

Regularization

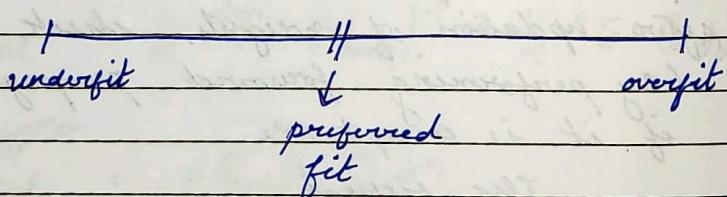


⇒ Poor performance
on both train &
test set
⇒ underfitting

⇒ good performance on
train set but poor
performance on test
set ⇒ Overfitting
Use regularization

Total cost $J(\theta)$ = primary cost + Regularization parameter (penalty).

We need to find a central point



Methods :

✓ $L^2 \Rightarrow$ penalty = $(\text{slope})^2$

✓ $L^1 \Rightarrow$ penalty = $|\text{slope}|$.

✓ Data augmentation

* flipping

* rotation

* warping (introducing noise/distortion)

✓ Early stopping (stop training when error is minimal or convergence occurs)

Even if # samples in dat train data is low, overfitting can occur

\Rightarrow use more data or data augmentation (if more data is not available).

Be careful when using data augmentation

(Ex): $9 \xrightarrow[\text{by}]{\text{rotate}} b$

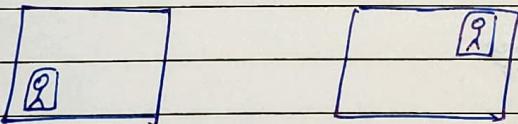
$d \xrightarrow[\text{vertically}]{\text{flip}} b$

} data is itself changed
(not preferred leads to misclassification)

$9 \xrightarrow{\text{blur}} ?$ (acceptable).

(cont. of Reg. methods) ✓ Parameter sharing

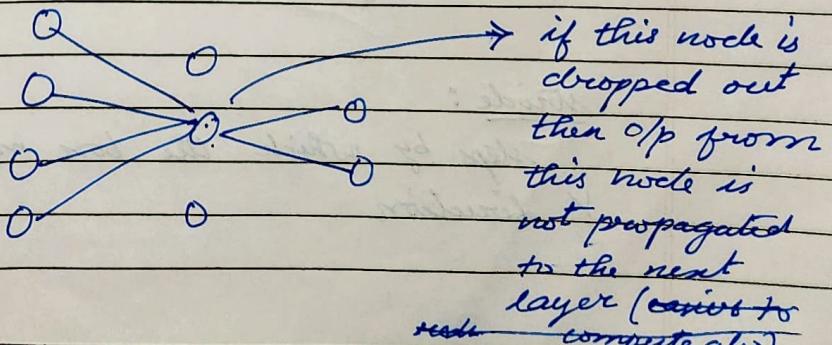
Consider 2 images in which the same person is present at different locations



To recognize the person, we can use the same parameters (can be saved & shared).

✓ Dropout

Randomly omit nodes



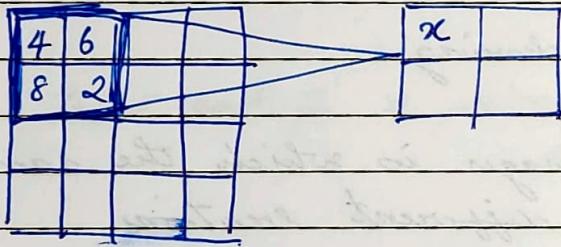
Leads to reduction in computation
 If one node causes overfitting then
 dropping it will help prevent dropout
 since we don't know which node is
 responsible for dropout, we randomly
 drop out some % of the nodes.

✓ Boosting & Bagging

Ensemble - Train many models & combine
 the results from them

09/09/22

~~stride~~



Max Pooling :

$$x = \max(4, 6, 8, 2)$$

$$= 8$$

Avg. Pooling :

$$x = \frac{4+6+8+2}{4}$$

$$= 5.$$

stride :

steps by which the box moves in X and Y direction

