

(2)

(AWFERA)

Course

Deep Learning

Lecture 1 + 2

Topic :-

(Overview) + (Application)

"Study of "neural network" is called Deep learning."

- Tuning of neural networks
- Optimization of models
- Techniques of learning of neural networks.
- Different Types of Architecture of NN.
 - Images (CNN)
 - Sequence Data like Stock Exchange, weather (RNN)

AIRLWA

(2)

- Auto-Encoders.
- Transformers.

What is Learning?

The process through which system or individuals improve performers or knowledge over time base on experience.

P. C. Datalink
P. C. Learning & Cognition Model

Human learning Machine learning
Biological brain → Digital Brain

Data Flow
Data

Thinking → Perception → Action

AI :-

When machine (like computer, Robots) are programmed to think and make decision like humans.

(3)

- Rule base System, Expert.
- Based System, Data Learn
- Algorithms.

Sub field:-

“ML” learning focused
from previous Data.

“DL” is subset of ML
which learns from neural
networks or Diff Types of
neural networks architecture.

- (s 1960) credit.
- (1970)

Application and factors:-

Popular in different
fields like

- Financial.
- Health Care
- Law and others

Factors:-

- * Data Availability (Tera bytes
(Generate))
- * Computational Powers
- * Algorithm Advancements
 - ↓ neural networks
 - (RNN, AutoEncoders, Transformers
(GPT combined like, Bard))

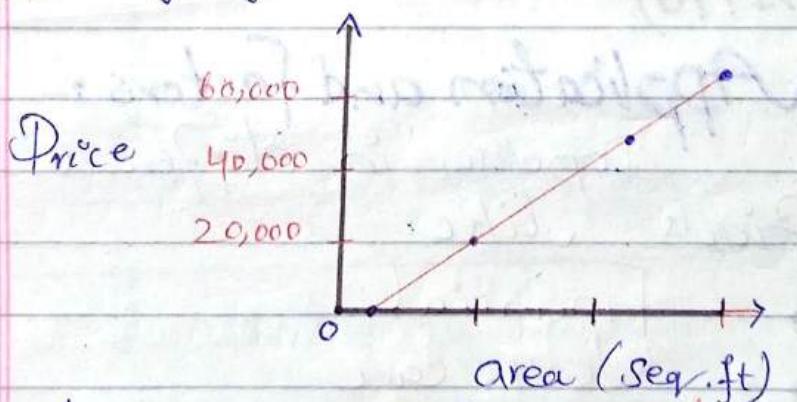
(4)

- Super-Computer's
 - Multicore-processing
 - GPUs available.
- * ↓ Pure-Mem is called Super-Computer. (Must fast speed of neural networks)

Example 8~

"House-Price" prediction
is an example of [Given by Area and Price].

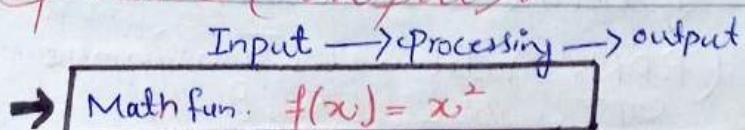
Independent variable → x -axis = Area
dependent variable → y -axis = Price. = Predicted



|| Straight Line is called Models / Neurons ||

Neurons:-

① A "neuron" in (AI) is like a simple math function. It takes inputs (data) ^{apply} math function and produce (output).



(5)

→ For Multiple features:-

like: Locality, facilities
↓
Location School, hospital
 Shopping mall

Multiple features → Better Prediction

From Simple To Complex (Automatic learning)
features Neural Networks.

Simple To Complex
features Combined Together.

Automatic learning [Layers hidden
Using in this condition]

Using of DL :- (Neural Networks)

In ML Paradigm.

- Supervised Labeled (Input, output) Prediction
- Un-Supervised Unlabeled Patterns, clusters
- Self-Supervised Raw-Data Improved Understanding

★ Application of DL :-

- Spam deletion, Self driving (Tesla)
- Cars, → face (Recognition), Speech
- Airports, glass of blind person.
- Healthcare → [diagnoses, ^{drug} disease], Translation apps)
- NLP (Siri, Google Assistant, Alexa) /
Chatbots,
- AI in Arts ⇒ Deep fake

Lectures # 3 + 4 (6)

Topics-

Neural Network's

→ Example : House Price Prediction

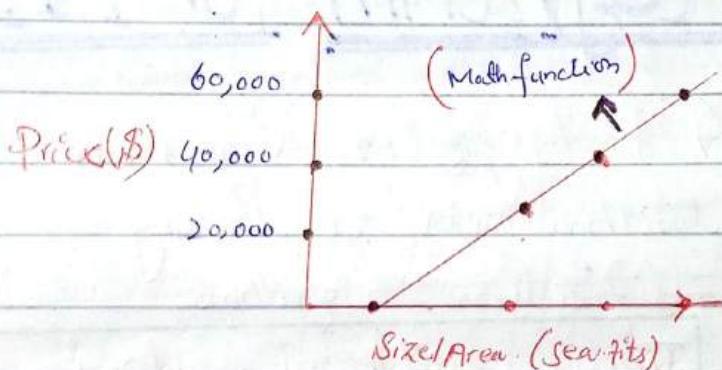
Suppose:- (Single features)

→ 6 House data Set

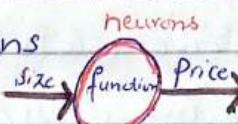
→ Input (Size), output = Price

→ how To calculate house Price

Using house Size? [linear Regression] → ML algorithm



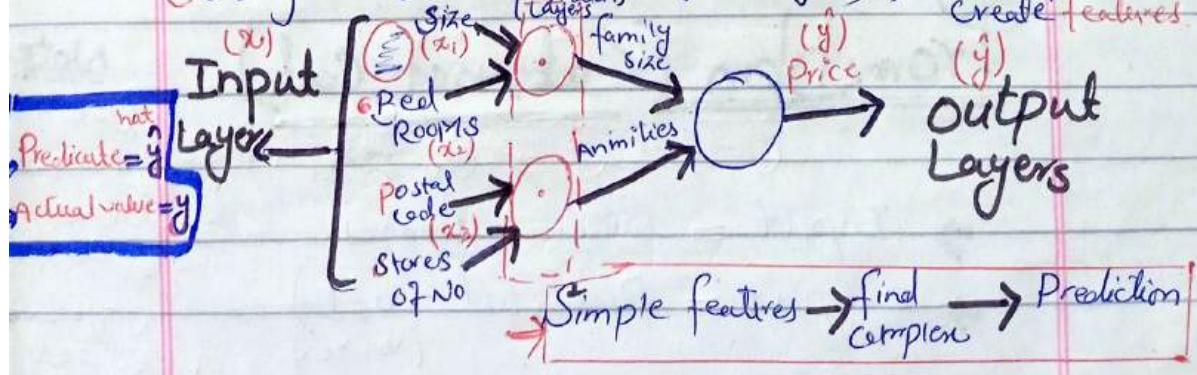
Simple neurons
is a function.



Take any
Values from Input,
output Predicted
values. (y).

For Multiple features :- (SNN)

Using many neurons arrange (Layers) find complex
create features.



- A neuron can make one neural networks and stack of different layers is called "Deep Neural Networks"

Logistic Regression

Algorithm for NN

- A Logistic Regression

Example Algorithm is a Binary

For Images Classification algorithms

Categorisation

Spam detection

Medical

diagnose

Two possible of output

- Positive class (1) = Target category
- Negative class (0) = Everything else

(Input output) pairs = (x, y)

where: x = Input, y = actual output

Problem (Cat, not Cat)

- Input = Given Images
- output = Predicted (0, 1)
(No, Yes)

Note:

Face

Introducing

probabilistic

filter

category

etc

3 Channel Using Image: (8)

⇒ Divided Images into three channels

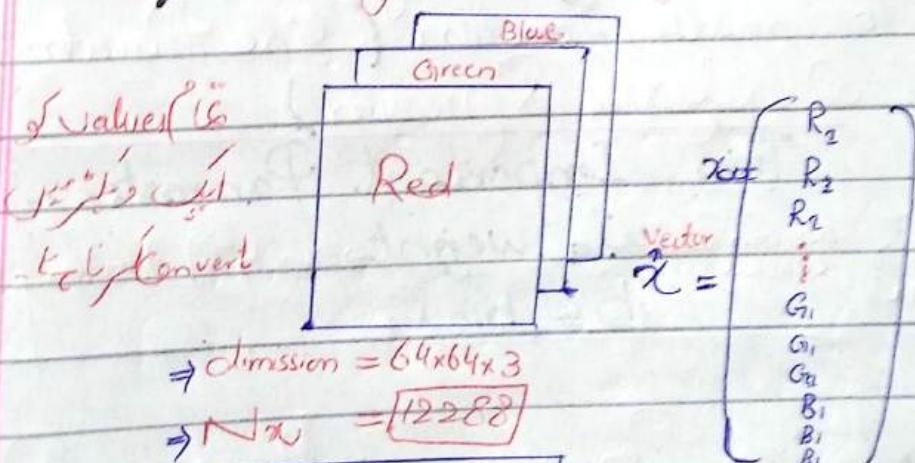
→ (Red, Green, Blue)

→ Multiple pixels create a vector and then pass to the Logistic Regression models {0,1}

→ x has M samples \rightarrow $x \in \mathbb{R}^{N_x \times M}$

→ (M ' Sample) Num of training examples / Samples.

→ 1 Images has three channels



* Logistic Regression In Math
 form:- $y \leftarrow \text{Input value} \rightarrow \hat{y} \leftarrow \text{Output value}$

→ $x = \text{input} \rightarrow x \rightarrow y \leftarrow \hat{y} \text{ hat value } \{0, 1\}$
 → $\hat{y} = \text{predicted} \rightarrow P(y=1 | x) \rightarrow \text{Probability}$
 $\rightarrow y=1 (+ve) class$
 $\rightarrow x = \text{Input value}$

* Defin Parameters :-

Calculate \hat{y} using ' x '.

→ What is probability $P(y=1)$

given ' x ' being a "cat picture".

→ Probability of Space value $\{0,1\}$

→ Logistic Regression give value b/w them

→ $(x \rightarrow \hat{y})$ ' x ' is a vector

• which is not control. So we

have Take ' x ' as Input and

Performed (maths function) and

convert $\{0,1\}$ Edding Some values

To math function (like Parameters)

of Logistic Regression

Two Important Parameters:-

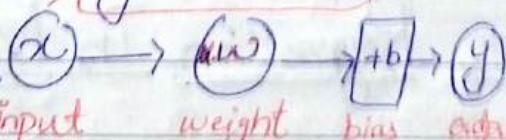
We can values control w = weight \rightarrow branch of function

b = bias

w values
can Min, Max
constant

Linear Regression

$$y = w_0 + b$$



Two operation
Using in this
case multiplication
Addition

well defined
operation
for vectors

Logistic Regression

$$y = \sigma(wx+b)$$

Applies sigmoid
function

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

Example:-

(10)

Suppose:-

vector $x = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix}$

Logistic Regression

$\rightarrow \boxed{\hat{y} = w^T x + b}$ —— \star

weights:

$$w = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}, w^T = \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$$

Transpose of w (changing Rows into columns)

Using equation \star

$$\hat{y} = w^T x + b$$

$$\hat{y} = [1 \ 2 \ 1] \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} + (-1) \quad \therefore b = -1$$

$$\begin{aligned} \hat{y} &= (2 \times 1) + (2 \times 1) + (1 \times 3) + (-1) \\ &= 2 + 2 + 3 - 1 \end{aligned}$$

$$\rightarrow \boxed{\hat{y} = 6}$$

$\{+\infty, -\infty\}$ = we just
~~(-ve, +ve, 0)~~ $\{0, 1\}$ value

\rightarrow Linear Regression Convert to Logistic Regression
 Logistic Regression Called "Sigmoid" function
 Take any numbers and convert into $\{0, 1\}$ numbers.

$$\rightarrow \boxed{\text{Sigmoid function} = \hat{y} = \sigma(w^T x + b)}$$

where,

$$\sigma = \frac{1}{1 + e^{-z}} \Rightarrow z = w^T x + b$$

$\therefore -\infty \leq z \leq \infty$

Case 1 :-

$$\rightarrow z = w^T x + b = 0$$

Using equation

$$\rightarrow \sigma = \frac{1}{1+e^{-z}}$$

$$= \frac{1}{1+e^{(-0)}}$$

$$\sigma = \frac{1}{1+1} \quad \because [e^0 = 1]$$

$$\rightarrow \boxed{\sigma = \frac{1}{2} = 0.5}$$

Case 2 :-

Larg positive Num

Using equation

$$\sigma = \frac{1}{1+e^{-z}}$$

$$= \frac{1}{1+e^{-(\text{larg Num})}}$$

$\therefore e^{-(\text{larg Num})}$

$$\therefore = \frac{1}{e^{(\text{larg Num})}}$$

$\therefore \frac{1}{e^{(\text{larg Num})}} \approx 0$
larg number approximation

$$= \frac{1}{1+0}$$

$$\rightarrow \boxed{\sigma = 1}$$

Case 3 :-

'z' is very large Negative Numbers

$$\rightarrow \sigma = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-[\text{larg-ve Num}]}}$$

$$= \frac{1}{1+e^{-[\text{larg-ve Num}]}}$$

for example,

$$\left. \begin{array}{l} \text{eg } z = (-5M) \\ \therefore e^{-5M} \approx 0 \end{array} \right\} \Rightarrow -5(5M)$$

$$= \frac{1}{1+e^{(\text{larg-ve Num})}}$$

$\therefore e^{\text{larg Num}}$ is very

$$\sigma = \frac{1}{1+\text{large Number}} \quad \text{large number} = 0$$

$$\rightarrow = \boxed{0}$$

Hence answer in Probability
Space

Lectures # 5+6

* Loss Function

How Logistic Regression Train & What Means?

جیسے کسی بھی b اور w پر
multiplied 'x' کے بعد ہے تو اسے find کر
کریں اس کا answer ہے /
 $y \approx 2, y = 2$ Representable
2 images of x weight
except 0.9, 0.8 \rightarrow (x, y) a vector

To purpose To find out
the value Actual Images
tables.

Example :-

→ Cat Images x by Trained weight
 $= 1$

→ Non-Cat Images vector x , Trained weight
 $= '0', or near To '0'$

Training:-

(13)

- ① In Neural Networks To find out Such type of Input to ^{convert} output (Labeled) Transform successfully. (**Parameters**) like weight, bias.

→ \hat{y} = Predicted Label.

→ y = Actual Label

Given → Ideal change $\hat{y} = y$

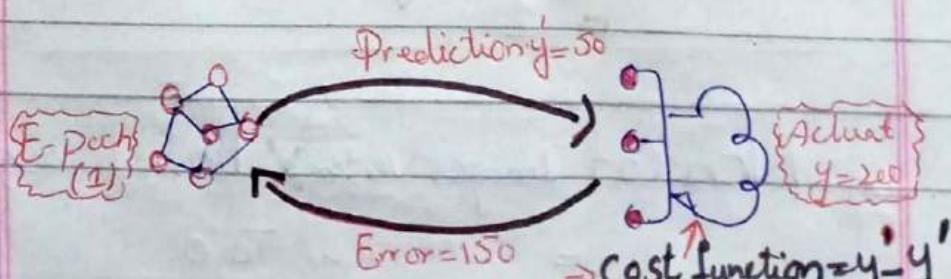
Given Cat Image $[1 = y]$ Label.

and \hat{y} to find equation

given by value "Using Special function To find weight and bias. distance"

Loss function:- (Cost)

- ② This function tell us how much distance b/w y and \hat{y} is called "Cost" Function.



Distance / Loss function

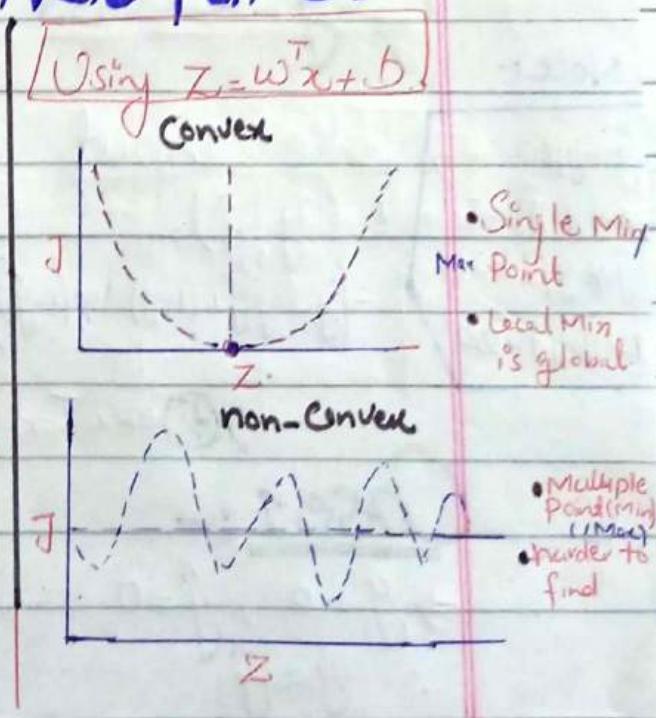
- It is given by distance b/w y and \hat{y} to find out the weight and bias whose value is "Minimum value".

- Take input (y, \hat{y})
- Given output (Some Num.)
- '0' = Min distance.
- '1' = Max distance.

- Algorithm is not Proper training. Special Category.

Convex - Non Convex function

- Distance b/w y and \hat{y}
- To choose a special type of function is Convex function
- y = Actual
- \hat{y} = Predicted
- σ = Discrepancy



(15)

Notes:-

To Neural Networks training

mosty function is used

by ("Convex Function").

→ optimin point (Min, Max)

both point use.

→ Single (one Point) is
Global Minima.

→ "non Convex" function
use multiple point (local
Global Minima).

Logistic Regression

Using function "Negative
Log-Loss function:-

Note:-

$$\hat{y} \approx y \quad (\text{only})$$

$$\Rightarrow L(y, \hat{y}) = 0$$

$$\hat{y} \neq y \quad (\text{error})$$

$$\Rightarrow L(y, \hat{y}) = \infty$$

Take Input

$$L(y, \hat{y}) =$$

$$-\left[y \log \hat{y} + (1-y) \log(1-\hat{y})\right]$$

$$= -\left[0 \log(0) + (1-0) \log(1-0)\right]$$

Putting values in eq (A)

→ \oplus evaluation

Case 1 :-

$$\Rightarrow y = 0, \hat{y} = 0$$

$$\therefore y \approx \hat{y}$$

$$= -\left[0 + (1) \log(1)\right]$$

$$= -\left[1 \log(1)\right]$$

$$= -\left[1(0)\right]$$

$$\Rightarrow = 0$$

$$\therefore \frac{10}{10} = 1 \\ \log \frac{1}{10} = 0$$

Case 2 :-

- When $y \neq y'$
- $y = 1, y' = 0$
- Putting values in equation *

$$= -[1 \cdot \log(0.1) + (1-1) \cdot \log(1-1)]$$

$$= -\left[\frac{1}{10} \log_{10}^{10^{-1}} + (0) \cdot \log(0)\right]$$

→ $\log_{10}(10^{-1}) = -1$ base and value is same

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

$$= -[-2]$$

→ $L(y, y') = 2$

① It's also called Cross-Entropy loss, binary cross-Entropy

→ To Tell us distance b/w them y and y' (To improve)

'B' & 'W' Shows,

→ To best Training Min value

② This function is defined for one value (Image) Training Set for (2 pictures)

Loss Function :-

Act as a mathematical distance function that minimize the Error b/w Predicted & Actual Labels To Enhance model performance.

Case :- Nth Value Sample

Example / Images. $m \approx 10^3$

Collective multiple images for Training.

For 'm' Sample:-

Adding all 'loss' of images and then Take Average.

m Sample = Add $L(y, y')$ for all cost function Represent by ' J '

$\hat{y} \approx y$
Using Randomly Selected by (w, b)

$$\Rightarrow J(w) = \frac{\sum_{i=1}^m L(\hat{y}_i, y_i)}{M}$$

Total Loss of Training set Samples
 $\hat{y} = \sigma(w^T x + b)$

$$J(w, b) = -\frac{1}{M} \left[y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

$$\Rightarrow -\frac{1}{M} \left\{ y \log [\sigma(w^T x + b)] + (1-y) \log [1 - \sigma(w^T x + b)] \right\}$$

Gradient Descent

(18)

- ① In this lecture we have study to optimize our parameters like (w, b).

$$\rightarrow \boxed{y \approx \hat{y}} \text{ So called}$$

"trained the models."

جسے \hat{y} find y bias, weight!

Output \rightarrow Images \rightarrow Input

Assume:- \hat{y} to y , Actual

- ② We just using only one weights for neural networks To Train it for Furthermore To Extend it for multiple weights, bias.

Meaning of this Algorithms

Difficult Sounding Name

- ③ Gradient \rightarrow Slop, Derivative \rightarrow "Rate of If we change Independent variable change"
 $'dx'$ to calculate how much

changing in dependent variable. $'dy'$

- ④ Descent \rightarrow (down words)

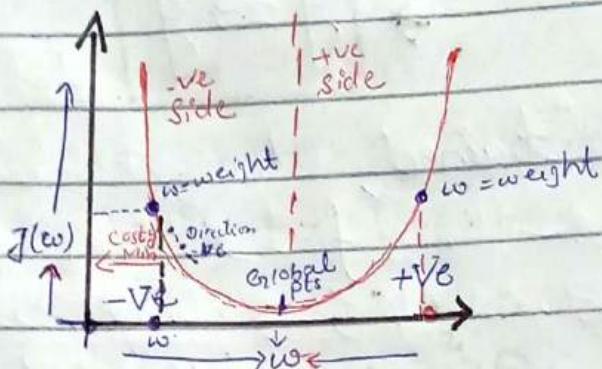
Opposite words Ascent (اوجاد)

Descent \rightarrow Gradient \rightarrow loss (loss) / weight (weight) \rightarrow minimize

Diagram:-

• We have Calculate.

$J(w) \rightarrow$ Cost function
of all training Samples.



$$\Rightarrow w = w - \text{change in loss per unit change in } w \quad (\text{updated as})$$

$$\Rightarrow w = w - dJ(w)$$

$$\therefore dJ(w) = 0$$

$$\Rightarrow w = w \quad \text{updated value two (to)} \\ \text{three times so}$$

the gradient descent Algorithms
is **Converge**. (so weight trained)

Case:- • To choose a Random value
of 'w'.

$w = w$ • w is optimize neural networks

There is not
changed using
gradient descent

Train (Means Cost Minimum)

$y^{(i)} \approx y^{(n)}$ for all Samples

(2c)

Case 1 :- [Negative Side]

Suppose we choose ' w ' randomly

$$\Rightarrow \frac{dJ(w)}{dw} < 0 \text{ (-ve Gradient)}$$

- when we change small w
slope is negative

$$w = w - \frac{dJ(w)}{dw}$$

$$= w - (-\text{ve negative Num})$$

$$\Rightarrow w = w + \text{sum Num}$$

- some w value changing (0.1)

$$w = w + (-\text{ve number})$$

$$= w + (0.1)$$

$$\Rightarrow w = w + 0.1$$

when $w = \text{Minimum} = 0$

$$w = w - (-0)$$

$$\Rightarrow w = w \text{ enshow less}$$

- minimized when w unchanged

$$\Rightarrow J(w) \approx 0$$

Case 2 :- (Positive Side)

- Suppose we choose w Randomly

$$\Rightarrow \frac{dJ(w)}{dw} > 0 \text{ (+ve slope)}$$

- Unit Cost J with.r.t

Unit change in w .

$$\Rightarrow w = w - \frac{dJ(w)}{dw} \rightarrow (*)$$

$$\Rightarrow w = w - (+\text{ve num}) \Rightarrow w = w - (\text{ve})$$

Example

Min

5 - 1

= 4

→ Change cost positive still
Same value hold so w is
more small.

Loss is

$$L(w) = L(\bar{w})$$

given dist

→ Similarly moving w small pts

Gradient near to Converg Point (Min)

Scent give

Gradient change
 w to optimize

$$w = w - (\text{very small num}) \approx 0$$

$$\rightarrow w = w - \alpha \quad (\text{cost is minimum})$$

• Same things is use for
multiple w and b bias.

from equation \star

$$\rightarrow w = w - \alpha \frac{\partial J(w)}{\partial w}$$

• α is a "Learning Rate"

we cannot control from Data?

how much Jumbled (w) changing

→ "Learning Rate α " = Gradient descent Rate
of change of descent
control

→ Jumping control Value never
miss any value Sequence Jumping
To control very Small Num

Xbywadjust (0.8) up To Soon (changed PTo P)

Too Low	Too High	Optimal
• Slow progress and Time to convergence	• overshooting failure To converge	• balance steps and smooth converge to Min

Lectures No # 7 + 8

Topic:- Part-1

Forward & Backward Propagation

* Important Concepts:-

Which methods we used in NN To Calculate and Computation in Groups.

Divided into Phases:-

- Forward Propagation.
- Backward Propagation.

Forward Propagation

To find
inference
loss

- ① The Process of Passing input values through a Computational Graph To obtain an output.

Commonly used in NN To calculate Prediction.

we have find loss in predicted
and actual values in NN.

Math forms:-

- The function we are analyzing

$$\rightarrow J(a, b, c) = 3(a + bc)^2$$

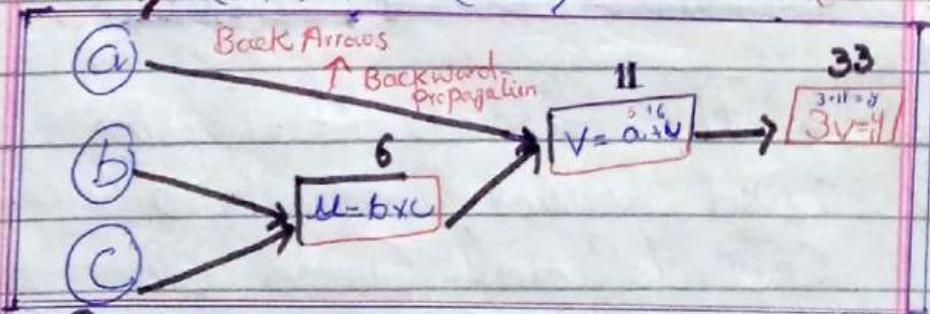
- This function involves basic arithmetic operations and can be broken into multiple steps for computation.

Simple Equation Examples:-

we using computational graphs nodes

Input Parameter

$$\rightarrow J(a, b, c) = 3(a + bc) \therefore b \cdot c = u, (a + bc) = v$$



Suppose:-

We have Suppose Random values

$$a = 5, B = 3, C = 2$$

Back-Ward Propagation

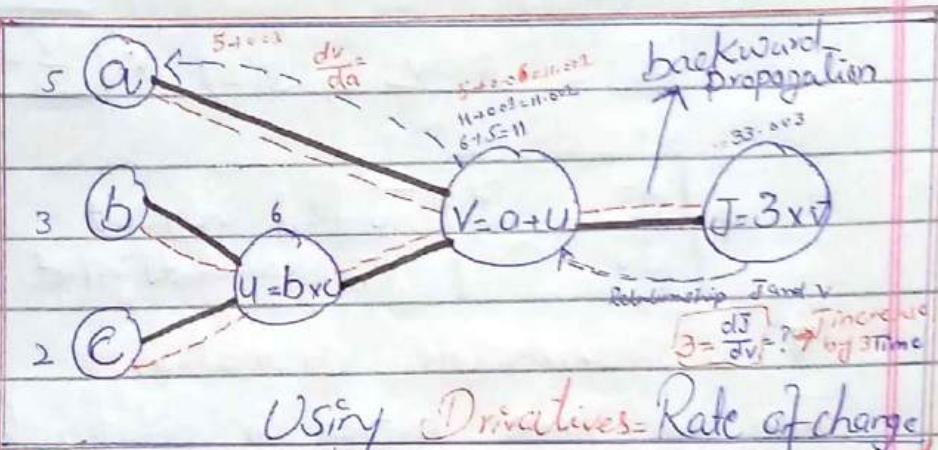
(64)

How finds changes in Input
affects the loss (Input var & output variable)

Goals: Adjust Parameters to
minimize the loss.

To find
loss and
Relationship
between use
this loss
to minimize

Here red dashed rows represent
Backpropagation.



Using Derivatives = Rate of change

Case 1:-

Suppose:- We have find

Relationship if 'a' unit change
then affect in 'J':

• In basic mathematics that

Relationship has find "derivative"

$\left(\frac{dJ}{da}\right)$ Independent variable 'a' ke
change karna sey depend

Variable J mai kia farq aata
hein.

→ So we have $\frac{dJ}{dv} = 3 \text{ Time}$

Since when $v = 11$ $s(v) = j$

$$\Rightarrow \text{unit change } v = 11 + 0.01 \Rightarrow 3(11.001) = j \\ \Rightarrow 11.002 \Rightarrow 33.002 = j$$

Note:-
 Suppose:
 Unit change
 independent
 variable mai
 karna or
 dependent
 variable
 mai kia
 Asar hota
 ny.



→ It's means 'V' unite changes

Then J '3' time change.

Now Taking derivatives

$$\frac{dJ}{dv} = 3v$$

$$\frac{dJ}{dv} = \frac{d}{dv}(3v)$$

$$\frac{dJ}{dv} = \frac{3}{dv} \frac{dv}{dv} \Rightarrow 3(1) \quad \therefore \boxed{\frac{dv}{dv} = 1}$$

→ $\boxed{\frac{dJ}{dv} = 3}$ To find independent v
To dependent J affect.

'V' is intermediate calculation.

Then we have interest in Relationship b/w 'a' w.r.t 'J'.

Case:2

→ We have find 'a' independent variable

Unit change then value of changed in 'v' affect. ($\frac{dv}{da}$)

→ Suppose $a = 5 + 0.01 = 5.001$

$$\rightarrow v = a + u \Rightarrow v = 5.001 + 6 = \boxed{11.001}$$

→ Unit change main itni hi value change
agri gi itni change 'a' main kithi.

Na nahi multiple (3, 4, 5) sey nota hy

$$\therefore \frac{dv}{da} = 1$$

$$\Rightarrow \frac{dv}{da} = \frac{d}{da}(a+u) \Rightarrow \frac{dv}{da} = \frac{da}{da} + \frac{du}{da} \quad \therefore \text{constant } \frac{du}{da} = 0$$

$$\Rightarrow \boxed{\frac{dv}{da} = 1 + 0}$$

Case 3:-

(26)

→ find Relational 'a' and 'j'

• We have find values
using Chain Rules.

$$\rightarrow \boxed{\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}}$$

Putting values.

$$\rightarrow \frac{dJ}{da} = 3 \cdot 1 = \boxed{3} \text{ see in graph.}$$

• Unit change 'a' Resultant
in 3 unit of J .

'a' is weight and J is
cost function To minimize it.

cost function when changed w .

• aK Relation $\frac{dJ}{dw} = ?$ Asa
button/function aajak

ha loss num kam sey kam kiske gay.

changing in ' w ' update y ^{True} loss has

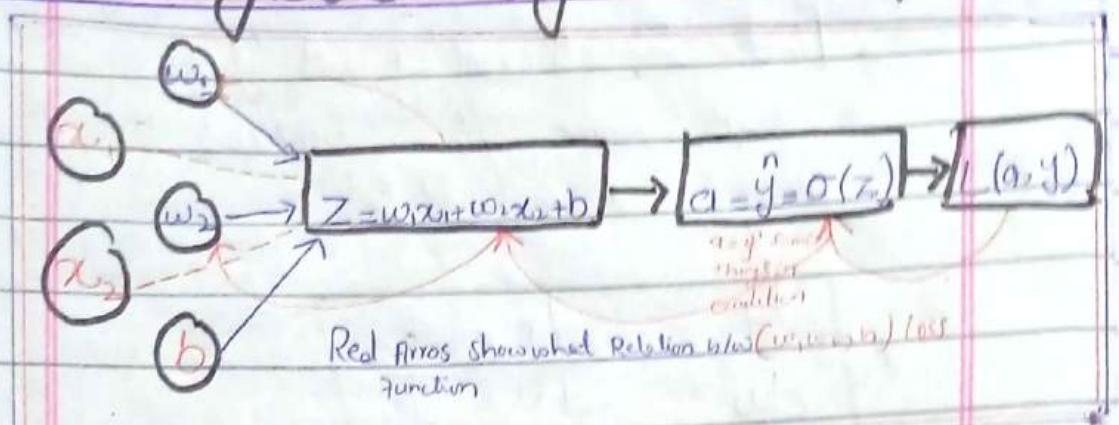
minimize the Logistic Regression Trained

we have $y \approx g$ find it some
weight To models predicted.

Video Part 2:-

Backward & Forward Propagation

Logistic Regression Graph



• We have Calculations in Logistic Regression :-

Note:-

control

w_1, w_2 - weight
not control

x_1, x_2 - pixel pic

$\rightarrow D$ = bias

→ Three para used

optimize To
loss Minimize

$$\hat{y} \approx \hat{g}$$

$$Z = w^T x + b \quad w^T \text{ is written}$$

$$\hat{y} = \sigma(Z) \rightarrow \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

So

$$Z = w_1x_1 + w_2x_2 + b$$

- Propag-
(Forward)
forward/
Propagated
- Infuse/
infuse/
Table, (P)
- Loss calcul
(cause diff
distance)

Case 1:- Loss function to
Relation $\hat{y} = a, P$ in backward

Propagation

$$\frac{dL(a, y)}{da} = \frac{-y}{a} \frac{(1-a)}{(-a)}$$

'a' Unit change in quantities

$$\left[\frac{-y}{a} \frac{(1-y)}{(1-a)} \right]$$

Case 2:-

We To calculate unit change of

Z then how much change in a ,

\hat{y} hat

$$\frac{da}{dz} = a(1-a)$$

calculation of math
provided by notes

clear
concepts

Case 3:-

How have To Calculate Relationship

b/w z , and w .

$$\rightarrow \frac{dz}{dw_1} = z_1, \quad \frac{dz}{dw_2} = z_2$$

Now:-

Now we calculate Relationship

b/w ' w ' and Loss (in backward

Propagation in Regr (Logistic Regression)

then we have Super power means

how much change (' w ') weight and

loss increase or decrease (Min, Max.)

then Neural Network Trained.

Using Chain Rule:-

$$\Rightarrow = \frac{d.z}{dw_1} \cdot \frac{da}{dz} \cdot \frac{dL}{da}$$

$$\Rightarrow = x_1(a-y)$$

This show that w_1 change (0.01)

then loss change $x_1(a-y)$ To

weight better and loss (minimize)

$$y^{(n)} \approx y^{(n)}$$

We calculate
Trained NN.

Similarly this calculation is done
for (w_2) and Biase (b)

(29)

- ④ we have Read Gradient Descent Algorithms.

$$\Rightarrow w_i = w_i - \alpha \frac{dL}{dw_i}$$

Apply gradient Descent iterative
and final (w , weight)

Backward propagation is used
To update weights.

Important Points:-

- When we have Training NN using both Technique like forward and backward propagation
- But when we have Trained NN there is no need To forward/backward propagation. Just (Model deployment)
using Forward propagation and find Predicted variable. (inference).

Standard Fully Connect NN:-

- ④ First NN which has diff

Given Break Throughs. This NN
Architecture (1918) *Yan Yan Kurn*
introduce it. We Read all
Basic Concepts in Previous
lectures with Using examples.

Content:-

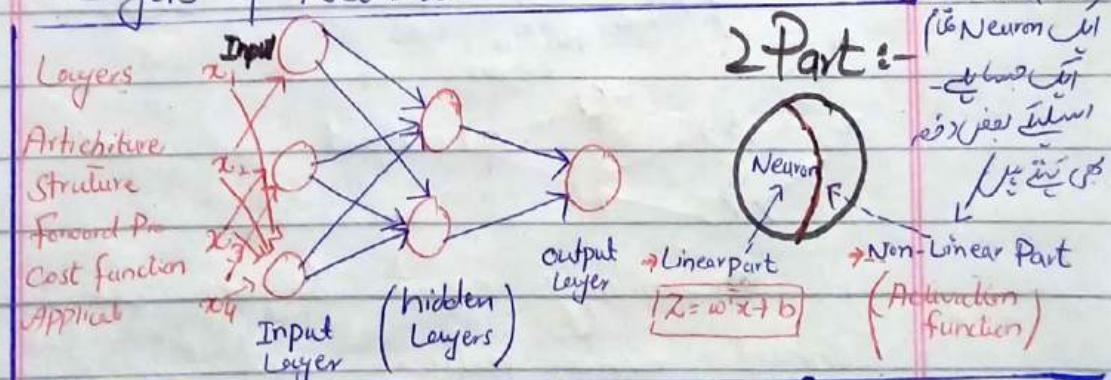
fully Connected Neural

Networks (Dense) Consists of following.

- Layers → Cost/Loss function
- Activation Function → Application
- Structure / Architecture
- Forward Propagation

Example:- Important things different

Layers of Neurons.



Input Layer → Hidden Layers → out-put Layer

Accepts Raw Features

More than one hidden layers in NN.

- to find pattern (Complex)

depending upon which type of problem has solved

Select Activation predicted

AlgorithmsActivation

Binary classification
Regression Algorithms
Continuous value
(Scores,)

Sigmoid/Tanh
ReLU/Variant
Softmax Groups
(Multi-class) programs

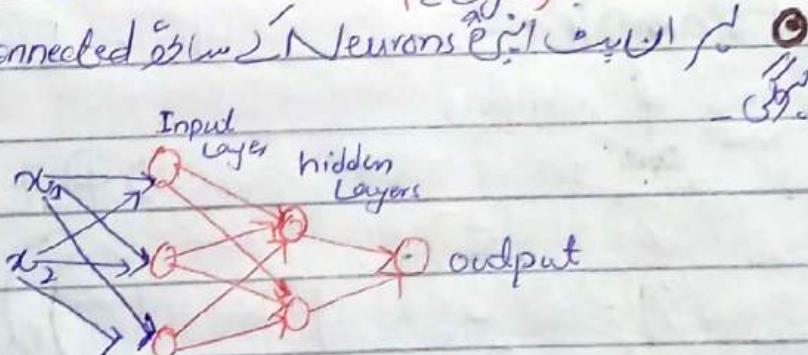
ReLU Explain
Explained

Similarly :-

- Some Case in Input Layers we have used (ReLU) activation
- Primarily start layers used and output depending upon problems Relevant activation function.

Structure :-

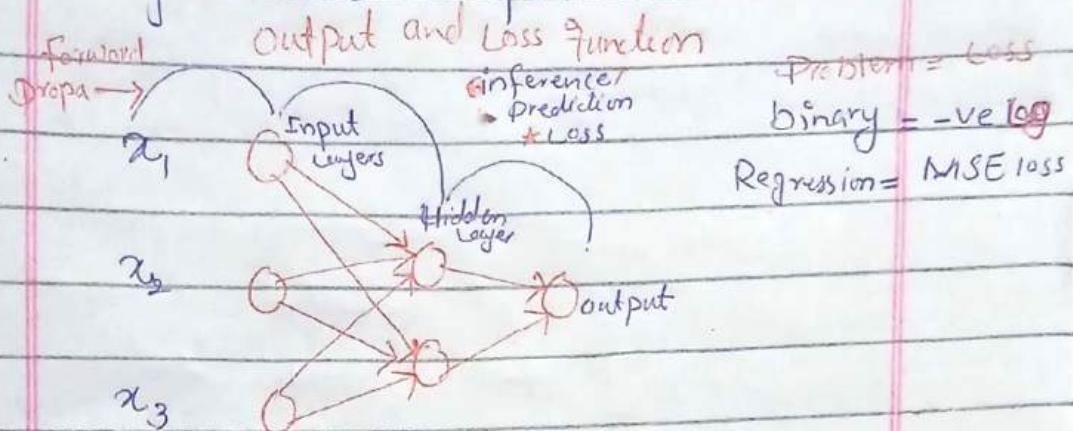
- It is also called (Fully-Connected neural Networks) or (dense NN)
(Connected between Neurons of all layers)



- Similarly next layers we connect every neuron and same in next layers is connected every neurons of layers.

Forward Propagation :-

If we moving in forward way is called forward calculate.



Equation forms:-

hidden Layer Equation

$$\rightarrow Z^{(2)} = W^{(2)} \cdot a^{(1-2)} + b^{(2)}$$

$$\rightarrow a^{(2)} = g(Z^{(2)}) \quad \because g = \text{Activation function}$$

Note: where, (hidden neuron \times input features)
dimensions $\rightarrow W^{(2)} \in \mathbb{R}^{(n \times m)}$: is the weight matrix for layer 2

3×3 Matrix
Rows col.

$a^{(2-2)}$: is the activation output from

$b^{(2)}$: the previous layers

$b^{(2)}$: is the bias

$g([z]^{(2)})$: is the activation

Rows \times column
Structure Data

Applied (SNN)

Application:- Examples

Computer Vision Problems:-

- Handwritten digits k^o pics
holi by or classifier krtey ky
digits (0 To 9) kon so digits

→ Speech
Recognition

→ NLP
(sentiment
+ve, -ve)

→ Finance
forward, Scam
detecting

→ Real Estate Analytics:-
④ House Price Prediction

(Lecture No# 09)

(33)

Topic:- Key - Activation Function

- The most popular activation function is as:

→ Sigmoid → Tanh

→ ReLU → SoftMax.

→ ReLU Variants

① Leaky ReLU ② Parameteric

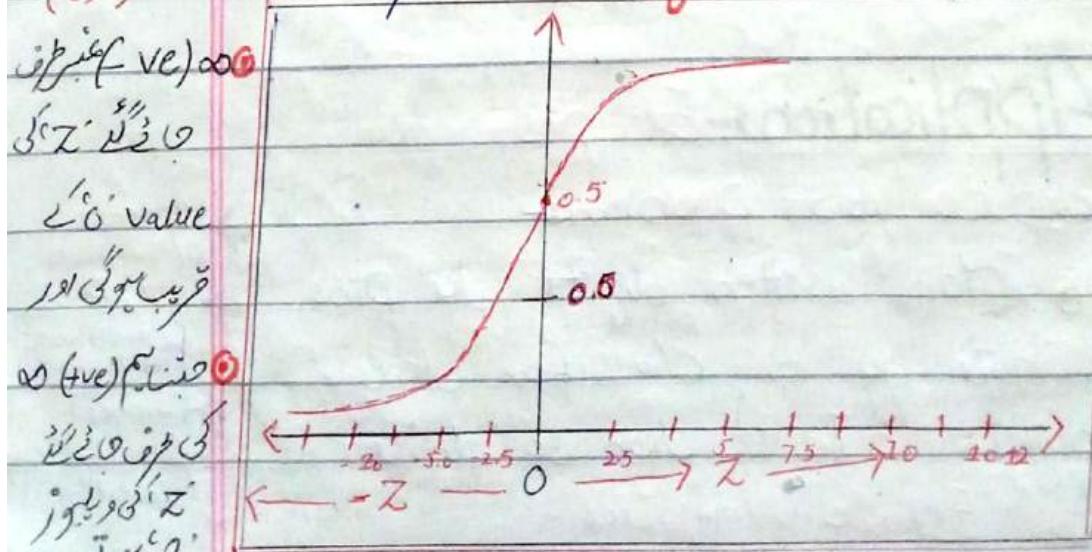
③ Exponential ReLU (ERU)

→ Sigmoid Activation Functions:-

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

- Sigmoid maps any input value z to range b/w 0 and 1, which is particularly useful for output probabilities in binary classification.

(0,1) Classification (usage)



Main Issue

① Computational
Expensive
Calculate e^z

every Time
in Neurons.

② e^{-z} is costly
function.

③ Usually not used
in hidden layers

Main Issues: (Cons)

- very Expensive for Computational power
- Non-zero Centralized output
It means if value = 0 then output value = 0 both '0' will not signified. This activation function may be suitable for no use neural network.
- Jatti converges by. (Gradient descent)

→ Always used in output layer
not used in hidden layers

Usage in Binary classification.

Note:- Particularly Problem Create (**Vanishing**

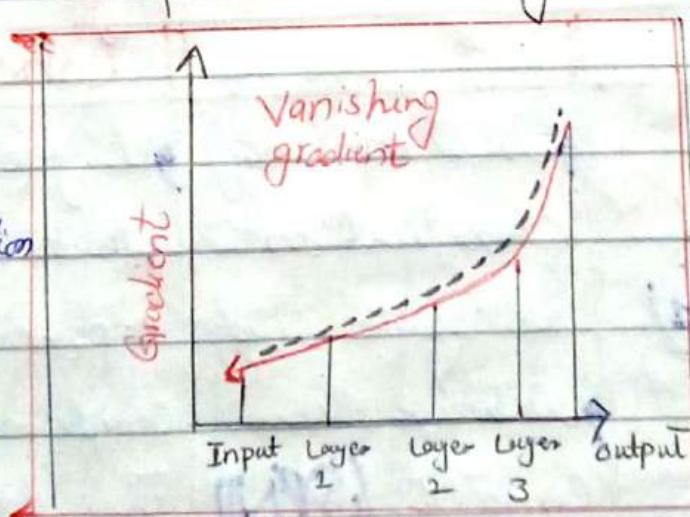
Gradient) A network learning issue where deeper layers struggle to update weight due to exponentially shrinking gradients (Mostly used when multiple hidden layers).

derivative
t_i to t_j
 $w_i - w_j$

backpropagation

$\frac{dw}{dz} = \frac{dL}{dz}$

$w = w - \alpha \frac{dL}{dz}$
not updated by neurons

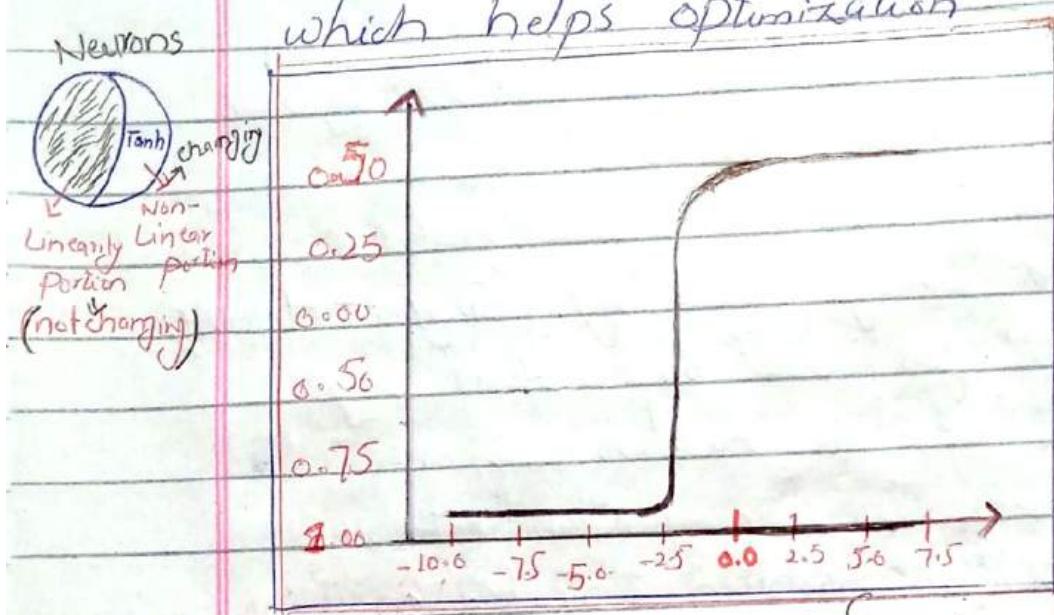


It is faced by vanishing gradient.
So not used in hidden layers.

Vanishing Gradient

★ Tanh Activation Function:-

- Is a Scaled version of Sigmoid Providing outputs in the Range of $(-1, +1)$ making zero-centred which helps optimization



Tanh Activation function

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

Also: $\rightarrow 2\sigma(2z) - 1$ is very similarly To Sigmoid but some changing

Pros:-

- It is a zero-centred
- Its Range $(-1, 1)$ output.
- any problem which is shown in (LVE, +ve)

Class Then output, layer used Tanh

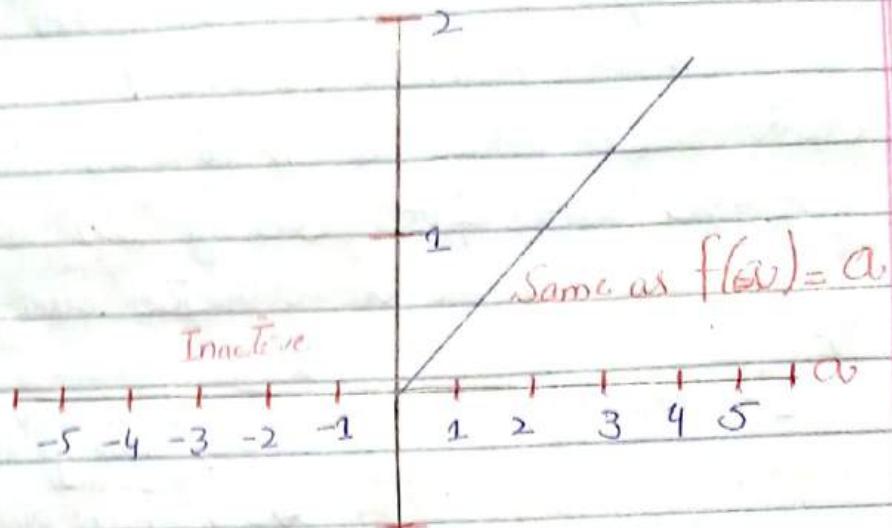
Cons:-

- Used in binary classification we have study Machine learning concepts (SFNN) loss.
- Using (e^{-z}) exponent which is very expensive for Computational power in output layers
- Some Recurrent hidden layers used it.
- Vanishing gradient

★ ReLU (Rectified Linear Unit) 36

- This most popular contemporary used activation function in (NN) (mostly used nowadays)

Rectified Linear Unit
→ $\text{ReLU}(a) = \max(0, a)$



{ If we give input as a neurons
+ve value = 0
-ve value = -A's
then give output as Same as
Input.

(ReLU) output the input value of a for all (+ve) values, and zero (0') for all (-ve) negative values. It's widely used in hidden layers To avoid the vanishing gradient problem.

Cons:-

Died Dying ReLU
problem

Neurons

in die'

if input
variables

more neurons
or no learning

Pros :-

It is more

efficient in

computational

power

Just compare use in hidden

number of both layers

max find (1000) layers present

no calculation easily used (ReLU)

deactivated on (+ve) function (NN complex)

numbers (Neural dropout provided) no vanishing gradient

★ RELU Variations-

① Leaky ReLU :-

$$g(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

Point : $z > 0$

All variation

changed To (-ve)

$$\rightarrow \text{Leaky ReLU}(z) = \max(\alpha z, z)$$

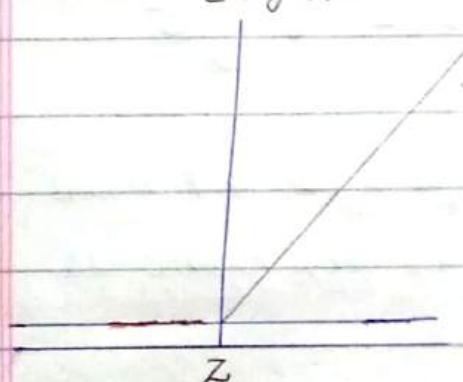
Portion and (+ve)
as it's as. for all.
not changed

we multiple by (-ve) number with
very small number ($\alpha = 0.01$)

Leak ReLU allows a small -ve step
for -ve values of z . preventing "dying ReLU"
Problem where neurons become inactive during

Training.

Leaky ReLU



* Example = -10

$$\begin{aligned} \text{ReLU Return} &= 0 \\ -10 &\rightarrow -10 \times 0.01 \\ &= -0.1 \end{aligned}$$

\Rightarrow dying ReLU Problem
Solved.

(e.g.: Speech, Images, language)

The value of
 α is fixed to $\frac{1}{2}$ $\rightarrow \alpha = 0.01$

② Parametric ReLU - Exponential ReLU

Value is αz is similar To
also a Leaky ReLU, but instead
of using a fixed value
with learn α , it introduce
data by a learnable parametric
key by α for (-ve) input.

$$\alpha = \text{Data learned}$$

It is also (-ve) numbers
address. Using exponential function

$$g(z) = \begin{cases} z & z > 0 \\ \alpha(e^z - 1) & z \leq 0 \end{cases}$$

it is ensure
output is centered

It is exponential expensive
for Training Time

ReLU Variations

Leaky ReLU

$$f(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

Point: 27

MI variation

Biased To +ve

Leaky ReLU is better than standard ReLU because it's not as flat and not exploded.

We multiply by (+ve) numbers with very small numbers.

Leaky ReLU allows a small +ve slope for -ve values of z preventing dying ReLU problem when neurons become inactive during training.

Training.

Leaky ReLU

Example: -10

$$\text{ReLU Return} = -10 \rightarrow -10 \times 0.01 = -0.1$$

Dying ReLU problem solved.

ions: much more stable

The value of α at fixed z : $\alpha = 0.01$

Parametric ReLUs: Exponential ReLU

is similar to Leaky ReLU, but instead of using a fixed value for α , it introduces a learnable parameter α .

A learnable parameter for (-ve) input.

$$g(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

It is also (-ve) number addressed using exponential.

For $z > 0$ it is e^z . For $z \leq 0$ output is zero.

It is exponential exp.

For Training Time

① RELU Variations-

A Leaky ReLU :-

$$g(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

Point : $z > 0$

All variation

changed To (-ve)

$$\text{Leaky ReLU}(z) = \max(\alpha z, z)$$

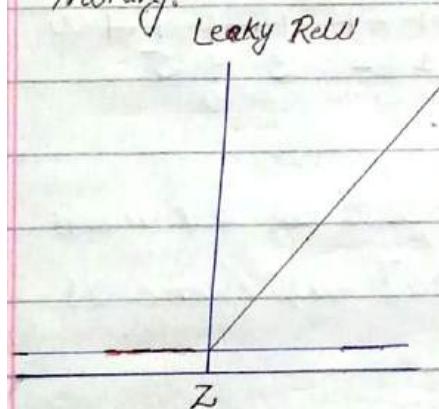
Portion and (+ve)
as it's as. for all.
not changed

we multiple by (-ve) number with
very small number ($\alpha = 0.01$)

Leak ReLU allows a small -ve step
for -ve values of z . preventing "dying ReLU"

Problem where neurons become inactive during

Training.



* Example = -10

ReLU Return = 0

$$-10 \rightarrow 10 \times 0.01$$

$$= -0.1$$

\Rightarrow dying ReLU problem
Solved.

(Cons: Speech, Images, language)

The value of

$$\alpha \text{ is fixed to } \frac{1}{2} \Rightarrow \alpha = 0.01$$

② Parametric ReLU :-

is a similar To

Reaky ReLU, but instead

of using a fixed value

for α , it introduce

a learnable parametric

for (-ve) input.

$$g(z) = \begin{cases} z & z > 0 \\ \alpha z & z \leq 0 \end{cases}$$

③ Exponential ReLU

It is also (-ve) number

address. Using exponential function.

$$g(z) = \begin{cases} z & z > 0 \\ \alpha(e^z - 1) & z \leq 0 \end{cases}$$

it is ensure
output is centred

It is exponential expensive

for Training Time