

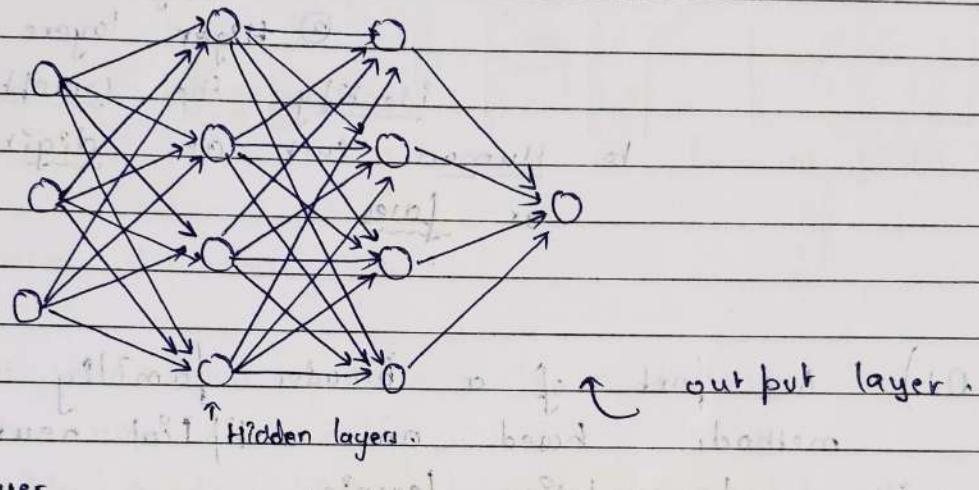
Deep learning :-

PAGE NO.:

DATE: / /

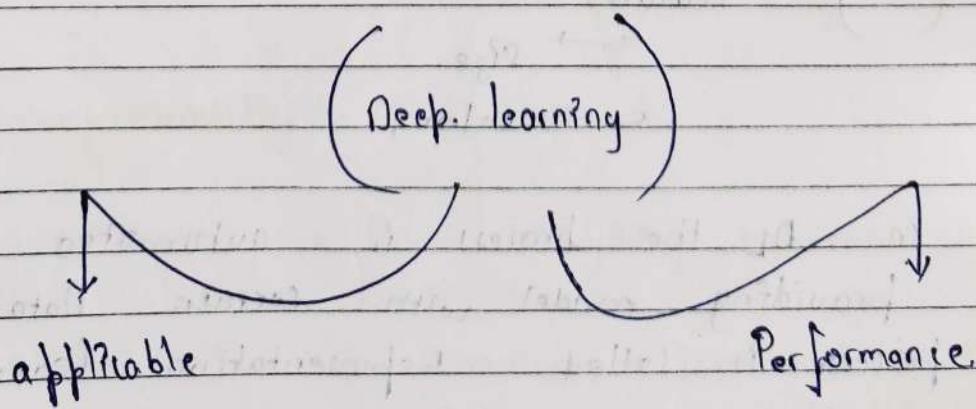
→ DL algo. attempts to draw similar conclusions as humans could by continually analyzing data with a given logical structure. called neural network.

e.g.



→ Different types of neural network:-

- ① ANN
- ② CNN → ex (works best on image data)
- ③ RNN → " " " " " " speech / text
- ④ GAN → generate thing. (text, images)



- DP → (Algo uses multi-layer).

~~multilayer~~ → progressively extract higher level features from other features.

e.g.

(image processing) ⇒ ①. lower layers may identify edges / shapes.

② Higher layers may identify the concepts relevant to Human such as digit or letters or faces.

- (DL) is part of a broader family of (ML) methods based on artificial neural network with representation learning.

- Representation learning :-

→ ex :- for Dog / cat classification
In ML you have to create features manually.

(ML) → Features.

→ size.

→ color.

But.

In DL the process is automated by providing model with correct Data.
This process is called Representation Learning.

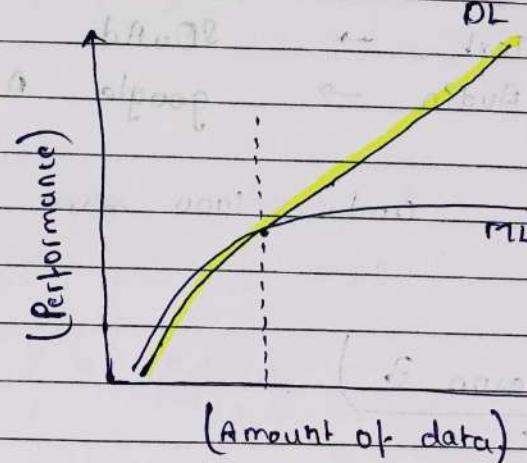
Deep Learning vs Machine Learning.

PAGE NO.:
DATE: / /

① Data Dependency :-

DL need more data
as compare to
Machine learning.

- DL is Data Hungry



② Hardware Dependency :-

{ DL me honot iyada. ~ matrix multiplication
to run note hai }

→ which needs **GPU** to run (quickly)
↑
(Powerfull)

③ Training Time :-

- In DL Training time is very high

④ Features selection :-

FL → manual feature creation.

DL → Automatically done by model.

⑤ Interpretability :-

ML → we know everything happening in each
and every step.

DL → Process of extraction of feature cannot be explained
to anyone. It is like a black box.

Data for Deep Learning → (Reason 1)

PAGE NO.:

DATE: / /

① Image → Microsoft COCO

② Video → YouTube 817 (6.1 million videos)

③ Text → SQuAD (Wikipedia 1,50,000)

④ Audio → Google AudioSet (20,00,000 sound clips)

and 1000 more datasets available //

• (Reason 2)

(DL) → Transistor → 2 years (?)
cost → half

NVIDIA → CUDA → GPU

(DL) → lot of { matrix } → CPU → low → (2010)
data operations speed

Parallel processing

GPU → reduces 10-15x time

than CPU

→ Hardware

• FPGA → fast, low power
↳ expensive, reprogrammable,ustomizable

e.g. Xilinx

• ASIC

→ TPU → Tensor Processing unit (google)
google colab
CPU option.

→ Edge TPU → Drone, watches
smart glass

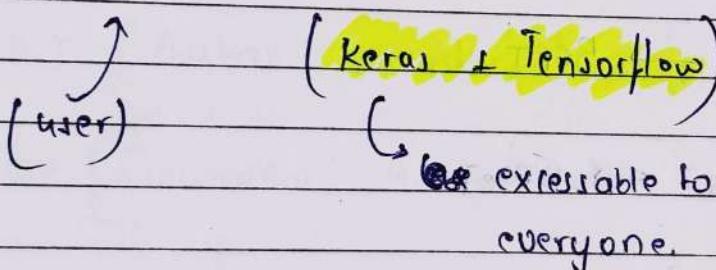
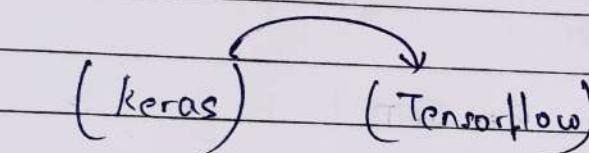
- Reason 3 → { Frameworks / Libraries }

sklearn → ML

- For Deep Learning →
 - ① Pytorch (Facebook)
 - ② Tensorflow (Google)

- Tensorflow (2015)

→ Powerful → Difficult



- Pytorch

AI researchers loved it.

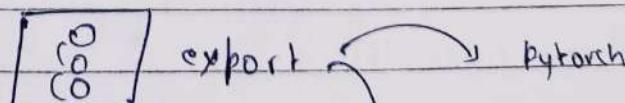
research driven.

- more industry driven

→ For conversion of code in two libraries.

Pytorch ↔ Tensorflow

⇒ Prop. draw GUI

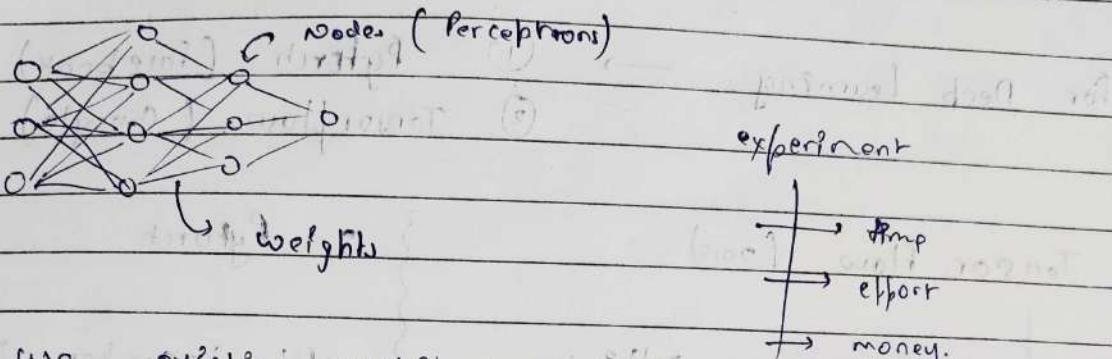


AutoML

or (Custom vision AI) or (CreateML)

• (Reason 4) → Deep learning Architectures.

→ Creating different architectures for best results



→ use existing architecture.

{ transfer learning }

- Image classification → ResNet

- Text classification → BERT

- Image segmentation → UNet

- Image dom → Pix2Pix

- Object detection → Yolo

- Speech generation → WaveNet

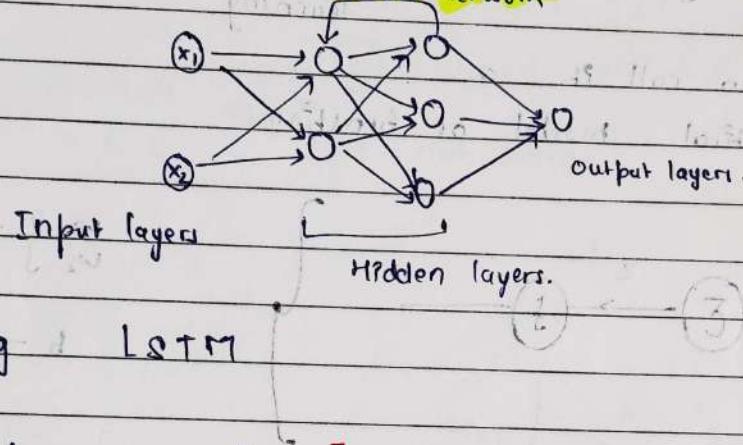
• Types of Neural Network.

1.1 [Multi layer Perceptron.] { supervised Problems }

2.2 [Convolutional Neural Network.] { Image object detection etc }

3.3 [Recurrent Neural Network]

Recurrent network



4.4 [Auto encoders.]

5.5 [Generative Adversarial Networks]

• History! (OL)

→ Applications.

(1) Game-playing agents. { Deep reinforcement learning }

(2) Virtual Assistant.

(3) Image colorization

(8) object detection

(4) Adding audio to mute videos

(5) Image caption generation.

Text translation.

(7) Pixel Restoration

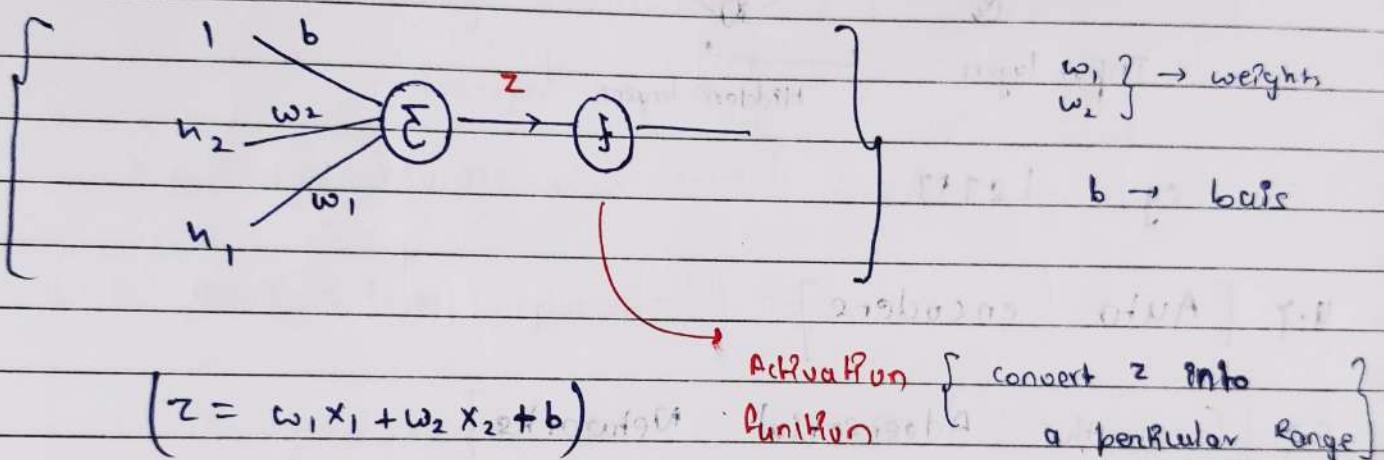
• Deep learning.

* Perceptron ??

well,

we can also call it a mathematical model or function

Algorithm → supervised
machine learning



e.g.-

IQ	CGPA	Motivation
----	------	------------

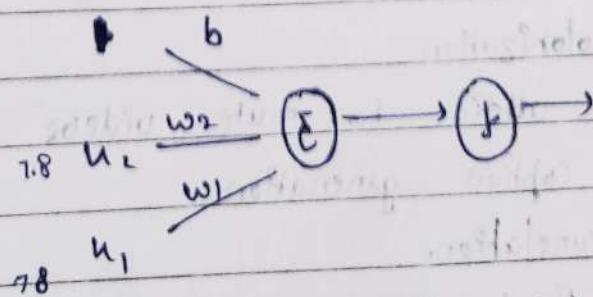
78 78 1

69 61 0

By the mathematical model of Perceptrons

$x_1 \rightarrow \text{IQ}$

$x_2 \rightarrow \text{CGPA}$



The whole purpose of training process is to

PAGE NO.:

DATE: / /

Find the value of w_1, w_2, b in training.

→ Let say by training we get $w_1 = 1, w_2 = 2$

and student with

x_0	x_1	x_2	Placed
1.00	5.1	?	

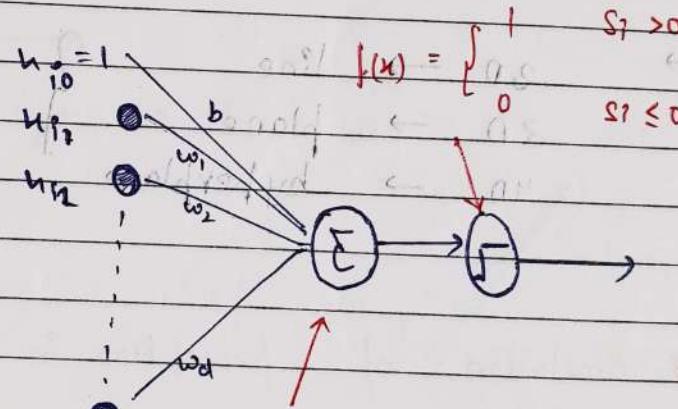
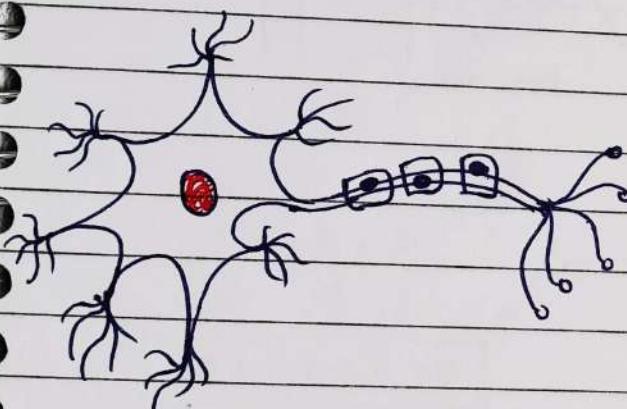
$$b = 3$$

So,

$$z = 1.00 \times 1 + 5.1 \times 2 + 3 \times 1$$

$z \geq 0 \rightarrow$ ① placed

$z \leq 0 \rightarrow$ ② not placed.

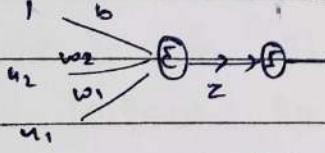


$$s_i = b x_{i0} + w_0 x_{i1} + w_1 x_{i2} + \dots + w_d x_{id}$$

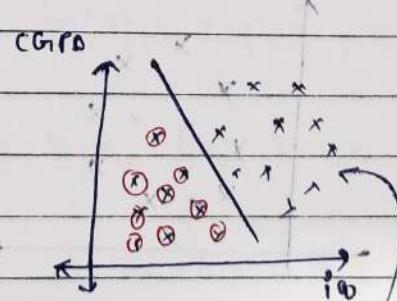
Neuron

Perceptron

• Geometric Intuition



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



$$\boxed{Ax + By + C \geq 0} \text{ region}$$

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

let say I have 3 inputs.

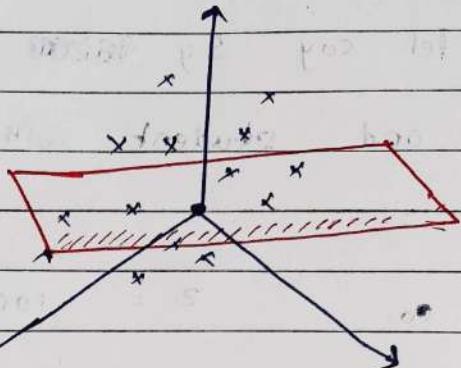
$i_0 | \text{bias} | 1^{\text{st}} \text{ input} | \text{Placed}$

$$f(z) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b \geq 0$$

$$Ax + By + Cz + b \geq 0$$

Plane. boundary

↪ How Perception acts.



- It is a **Binary classifier.**



2D → Line.

3D → plane

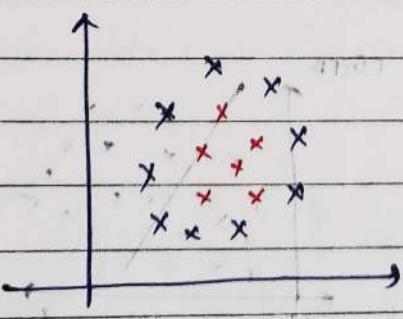
4D → hyperplane

* limitation of perception :-

- It can only classify linear or part of linear cases.

It fails in non-linearly separable cases.

e.g.



∴ (Perception fails here)

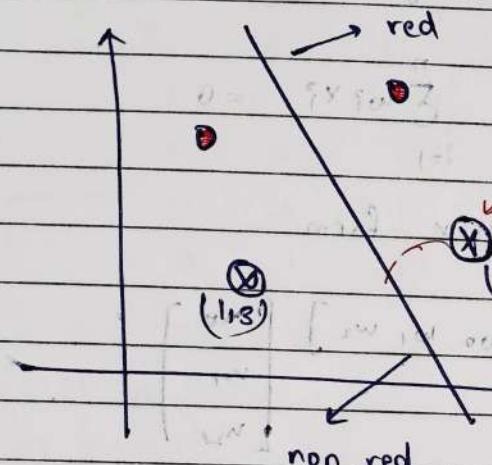
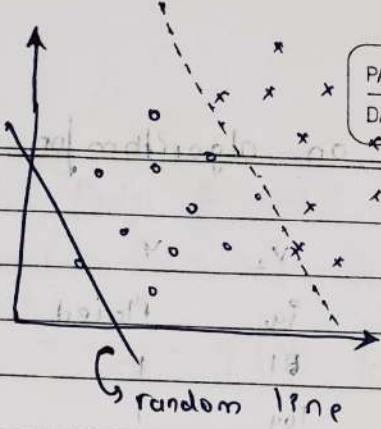
Perceptron trick.

PAGE NO.:

DATE: / /

Apply transformation

to converge from random
line to original



Trick

2 3 5
U S I

$$-2x - 2y + 4 = 0$$

$$2x + 3y + 5 = 0$$

Then line will

bring (4,5)

in non red region

or

(to move a point to
+ve region)

2 3 5

1 3 10 = 4

$$3x + 6y + 6 = 0$$

- If a negative pt. is in positive region.

DO ①

- If a +ve pt. is in (-)ve region. DO ②

Note 8- In YML do not apply these big changes instantly
instead do that.

$$\eta = \text{learning rate} = 0.01$$

$$\begin{matrix} \text{so} & 2 & 3 & 5 \\ & 0.01 & 0.03 & 0.01 \end{matrix}$$

same for (+)ve

$$2.01x + 3.03y + 5.01 = 0$$

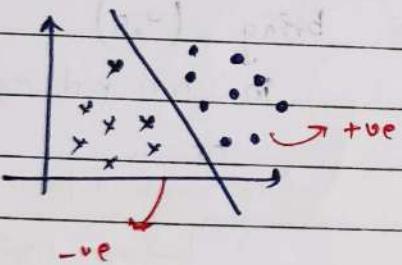
- Let's Build an algorithm for this :-

x_0	x_1	x_2	y	
CgPA	92	Placed		$AX + BY + C = 0$
1	7.5	81	1	$w_0 + w_1x_1 + w_2x_2 = 0$
1	8.9	log	0	$w_0x_0 + w_1x_1 + w_2x_2 = 0$
1	!	;	1	$\sum_{i=1}^n w_i x_i = 0$

matrix form

$$\begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = 0$$

$$\begin{cases} w_i x_i \geq 0 & \text{if } \dots \\ < 0 & \text{else} \end{cases}$$



epoch $\rightarrow 1000$, $\eta = 0.01$

for i in range (epochs) :

randomly select a row

if $x_i \in N$ and $\sum w_i x_i \geq 0$ then

$$w_{\text{new}} = w_{\text{old}} - \eta x_i$$

if $x_i \in P$ and $\sum w_i x_i < 0$ then

$$w_{\text{new}} = w_{\text{old}} + \eta x_i$$

Simplifying algo.

for i in Range(0, 1000) :

Pick random student.

$$w_{\text{new}} = w_{\text{old}} + \eta (y_i - \hat{y}_i) x_i$$

Intuition

$$y_i \quad \hat{y}_i \quad y_i - \hat{y}_i$$

$$1 \quad 1 \quad 0 \}$$

$$\rightarrow w_{\text{new}} = w_{\text{old}}$$

$$1 \quad 0 \quad 1$$

$$0 \quad 1 \quad -1$$

$y_i \rightarrow$ actual result

$\hat{y}_i \rightarrow$ model prediction

$$w_1 = w_0 + \eta x_1$$

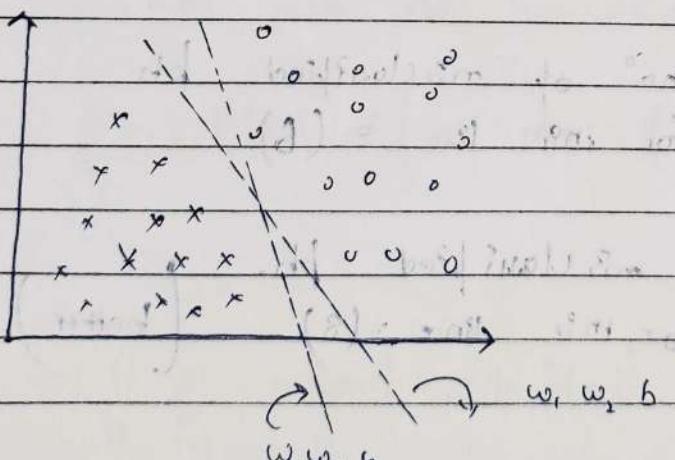
$$w_0 = w_0 - \eta x_1$$

for pt in

-ve region

→ Practical done in (DL folder)

* Problems with Perceptron



• result by m3 trick

can't be quantify.

(Kisna acha result aaya hai?)

• running code u can

w, w2, b get any of these 1/men.

→ can be different with

initialization for same d. for same 2-3 times.

will give diff. ans.

Loss function (Perceptron.) :-

The method to tell which solution is better than the other.

ML Loss functions

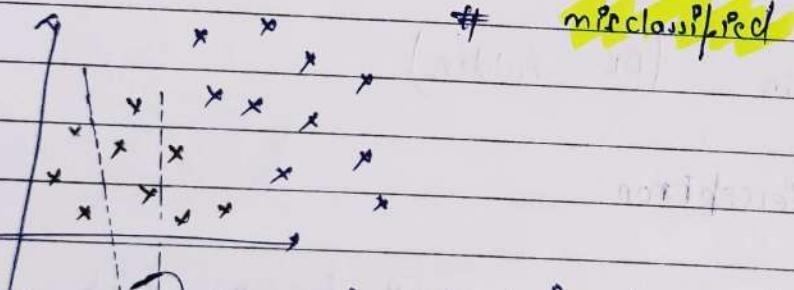
$$\text{LR} \rightarrow \text{mse} = (y_i - \hat{y}_i)^2$$

$$\begin{array}{l} \text{Log Reg.} \rightarrow \text{Log loss} \\ \text{sum} \rightarrow \text{Hinge loss} \end{array}$$

→ Let's find a loss function for Perceptron task.

$$\textcircled{1} \quad f(w_1, w_2, b) \curvearrowright \text{number} \curvearrowright \text{error.}$$

misclassified pts. (number)

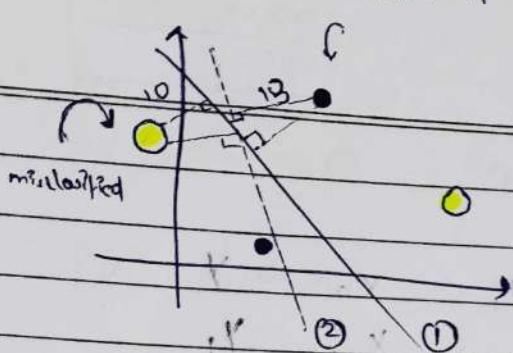


$f(w_1, w_2, b) \rightarrow$ no. of misclassified pts
for the line = (6)

no. of misclassified pts.
for the line = (3) (better)

$$\textcircled{2} \quad \text{loss function} = \text{sum of } h \text{ dist. of misclassified pt. from the line.}$$

Prev. fun. we took every misclassified pt. of the same weight.



$$\textcircled{1} \quad LF = 10 + 13$$

$$= 23$$

\textcircled{2} 00 same for other as well
man compair.

* sklearn / formula for Perceptron loss function :-

on page of stochastic gradient descent

(optimization)

$$\rightarrow L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \gamma R(w)$$

For Regularization

let take $h(x)$ for now

$$L(y_i, f(x_i)) + \max(0, -y_i f(x_i))$$

$$f(x_i) = w_1 x_1 + w_2 x_2 + b$$

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i)) \quad \text{--- (1)}$$

$\therefore (n = \text{no. of rows in data.})$

\therefore we need to find w_1, w_2, b for which value of (1) in eqn (1) is minimum

Explanation of loss function.

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

	x_1^*	x_2^*	y^*
rows	x_{11}	x_{12}	y_1
	x_{21}	x_{22}	y_2
	,	,	
	x_{i1}	x_{i2}	y_i
	x_n		

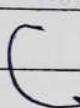
where

$$f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$$

now

$$\max(0, -y_i f(x_i))$$

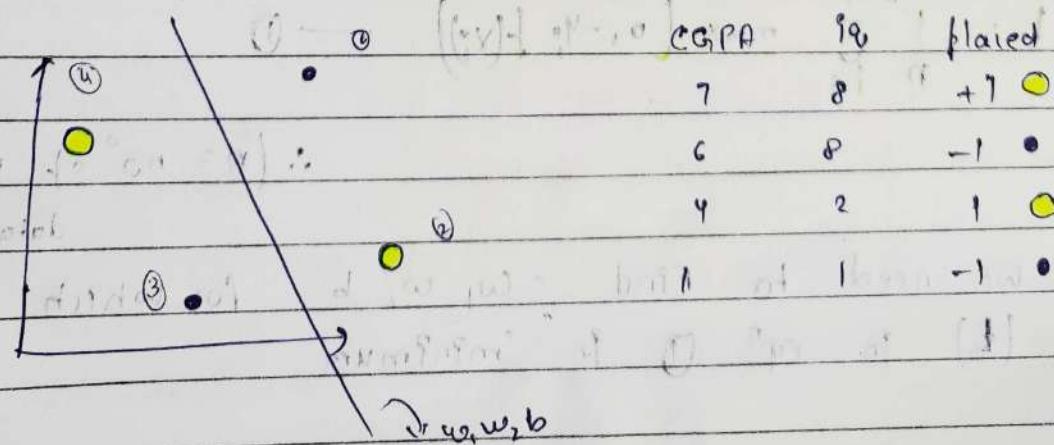
$$\begin{cases} = -y_i f(x_i) & \text{if } -y_i f(x_i) \geq 0 \\ 0 & \text{if } -y_i f(x_i) < 0 \end{cases}$$



lets do for only two fts

$$L = \frac{1}{2} [\max(0, -y_1 f(x_1)) + \max(0, -y_2 f(x_2))]$$

→ Geometric intuition.



• for pt ②

$$y_i = 1 \quad \max(0, -y_i f(x_i)).$$

$$f(x_i) = (\omega_1 x_1 + \omega_2 x_2 + b) \geq 0$$

$y_i = +ve$

$$-y_i f(x_i) \leq 0$$

$f(x) = +ve$

$$\text{so } \max(0, -y_i f(x_i)) = 0$$

↙ -ve.

• for pt ③

$$y_i = -1$$

$$\max(0, -y_i f(x_i))$$

$$f(x_i) = \omega_1 x_1 + \omega_2 x_2 + b.$$

$$f(x_i) \leq 0 \rightarrow \text{Due to pt in -ve region}$$

$$\text{so } \max(0, -y_i f(x_i)) = 0$$

$y_i < 0 = -ve$
 $f(x_i) \leq 0 = -ve$

• for pt ④

$$y_i = +1$$

$$f(x_i) = \omega_1 x_1 + \omega_2 x_2 + b \leq 0 \rightarrow \text{in +ve side of line}$$

$$\max(0, -y_i f(x_i))$$

$y_i > 0 = +ve$

$$= -y_i \cdot f(x_i)$$

$f(x_i) = +ve$

• for pt ①.

$$\max(0, -y_i f(x_i)) = -y_i f(x_i)$$

To find values for which values of w_1, w_2, b are

so that loss function becomes minimum.

$$L = \arg \min_{w_1, w_2, b} \frac{1}{n} \sum_{j=1}^n \max(0, -y_j f(x_j))$$

→ Gradient Descent :-

Start a loop.

for i in epochs :

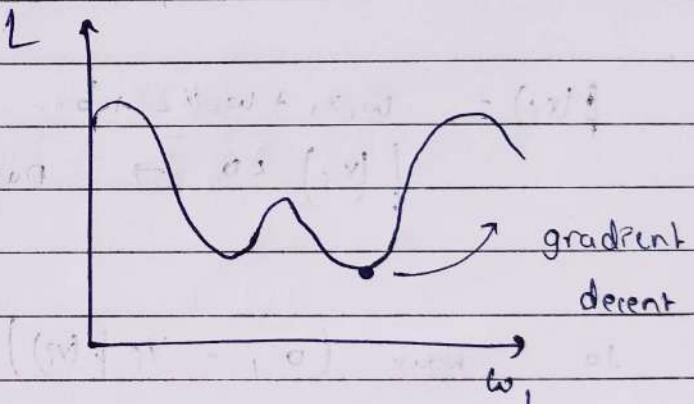
$$w_1 = w_1 + \eta \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 + \eta \frac{\partial L}{\partial w_2}$$

$$\eta = 0.01$$

$$b = b + \eta \frac{\partial L}{\partial b}$$

loop ends.



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \times \frac{\partial f(x_i)}{\partial w_1}$$

PAGE NO.:

DATE: / /

$$L = \frac{1}{n} \sum \max \left(0, -y_i f(x_i) \right)$$

$$\frac{\partial L}{\partial f(x_i)} = \begin{cases} 0 & -y_i f(x_i) \leq 0 \\ -y_i & -y_i f(x_i) > 0 \end{cases} \quad \frac{\partial f(x_i)}{\partial w_1} = x_i$$

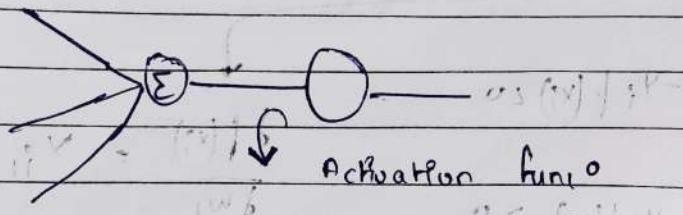
$$\textcircled{1} \quad \frac{\partial L}{\partial w_1} = \begin{cases} 0 & -y_i f(x_i) \leq 0 \\ -y_i x_i & -y_i f(x_i) > 0 \end{cases}$$

$$\textcircled{2} \quad \frac{\partial L}{\partial w_1} = \begin{cases} 0 & -y_i f(x_i) \leq 0 \\ -y_i x_i & -y_i f(x_i) > 0 \end{cases}$$

$$\textcircled{3} \quad \frac{\partial L}{\partial b} = \begin{cases} 0 & -y_i f(x_i) \leq 0 \\ -y_i & -y_i f(x_i) > 0 \end{cases}$$

Let's convert it to index form

• More about function



it is a mathematical model, which is very flexible

→ **flexibility** \Rightarrow we can change activation function (A.F)

(step) \rightarrow (sigmoid fun. o)

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

• If Activation function $= \sigma(z) = \frac{1}{1 + e^{-z}}$

then

loss function = binary cross entropy,

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

$\hat{y}_i \rightarrow$ output will probability
(check gradient descent)

so we can say.

Perceptron = logistic Regression with A.F

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

• If I am working on multiclass classification problem.

The.

$A \cdot f = \text{softmax}$

Loss funct = Categorical cross entropy.

$$L_i = \sum_{j=1}^n y_j \log(\hat{q}_j) \quad \text{--- (1)}$$

Soft max $\Rightarrow f = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$ --- (2)

let say we have 3 classes.

$$f = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

So Perceptron = Softmax regression

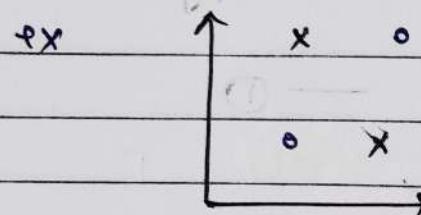
for (1) & (2)

Summary

Loss function	Activation	Output
Hinge loss	Step	perceptron - Binary classification (-1, 1)
log-loss (binary cross entropy)	sigmoid softmax	logistic regression. (0-1) - binary C.
categorical cross entropy	softmax	softmax regression → multiclass classification

- Problem with perception :-

→ It is of no use with non-linear data.
(cannot classify XOR function).



(P) This type of functions cannot be classified with perceptron

Ques 1

$$\sum_{j=1}^m \epsilon_j = 0$$

$$\epsilon_1 + \epsilon_2 + \epsilon_3 = 0$$

Ans 1

Ques 2

Ans 2

Ques 3

Ans 3

Ques 4

Ans 4

Ques 5

Ans 5

Ques 6

Ans 6

Ques 7

Ans 7

Ques 8

Ans 8

Ques 9

Ans 9

Ques 10

Ans 10

Ques 11

Ans 11

Ques 12

Ans 12

Ques 13

Ans 13

Ques 14

Ans 14

Ques 15

Ans 15

Ques 16

Ans 16

Ques 17

Ans 17

Ques 18

Ans 18

Ques 19

Ans 19

Ques 20

Ans 20

Ques 21

Ans 21

Ques 22

Ans 22

Ques 23

Ans 23

Ques 24

Ans 24

Ques 25

Ans 25

Ques 26

Ans 26

Ques 27

Ans 27

Ques 28

Ans 28

Ques 29

Ans 29

Ques 30

Ans 30

Ques 31

Ans 31

Ques 32

Ans 32

Ques 33

Ans 33

Ques 34

Ans 34

Ques 35

Ans 35

Ques 36

Ans 36

Ques 37

Ans 37

Ques 38

Ans 38

Ques 39

Ans 39

Ques 40

Ans 40

Ques 41

Ans 41

Ques 42

Ans 42

Ques 43

Ans 43

Ques 44

Ans 44

Ques 45

Ans 45

Ques 46

Ans 46

Ques 47

Ans 47

Ques 48

Ans 48

Ques 49

Ans 49

Ques 50

Ans 50

Ques 51

Ans 51

Ques 52

Ans 52

Ques 53

Ans 53

Ques 54

Ans 54

Ques 55

Ans 55

Ques 56

Ans 56

Ques 57

Ans 57

Ques 58

Ans 58

Ques 59

Ans 59

Ques 60

Ans 60

Ques 61

Ans 61

Ques 62

Ans 62

Ques 63

Ans 63

Ques 64

Ans 64

Ques 65

Ans 65

Ques 66

Ans 66

Ques 67

Ans 67

Ques 68

Ans 68

Ques 69

Ans 69

Ques 70

Ans 70

Ques 71

Ans 71

Ques 72

Ans 72

Ques 73

Ans 73

Ques 74

Ans 74

Ques 75

Ans 75

Ques 76

Ans 76

Ques 77

Ans 77

Ques 78

Ans 78

Ques 79

Ans 79

Ques 80

Ans 80

Ques 81

Ans 81

Ques 82

Ans 82

Ques 83

Ans 83

Ques 84

Ans 84

Ques 85

Ans 85

Ques 86

Ans 86

Ques 87

Ans 87

Ques 88

Ans 88

Ques 89

Ans 89

Ques 90

Ans 90

Ques 91

Ans 91

Ques 92

Ans 92

Ques 93

Ans 93

Ques 94

Ans 94

Ques 95

Ans 95

Ques 96

Ans 96

Ques 97

Ans 97

Ques 98

Ans 98

Ques 99

Ans 99

Ques 100

Ans 100

Ques 101

Ans 101

Ques 102

Ans 102

Ques 103

Ans 103

Ques 104

Ans 104

Ques 105

Ans 105

Ques 106

Ans 106

Ques 107

Ans 107

Ques 108

Ans 108

Ques 109

Ans 109

Ques 110

Ans 110

Ques 111

Ans 111

Ques 112

Ans 112

Ques 113

Ans 113

Ques 114

Ans 114

Ques 115

Ans 115

Ques 116

Ans 116

Ques 117

Ans 117

Ques 118

Ans 118

Ques 119

Ans 119

Ques 120

Ans 120

Ques 121

Ans 121

Ques 122

Ans 122

Ques 123

Ans 123

Ques 124

Ans 124

Ques 125

Ans 125

Ques 126

Ans 126

Ques 127

Ans 127

Ques 128

Ans 128

Ques 129

Ans 129

Ques 130

Ans 130

Ques 131

Ans 131

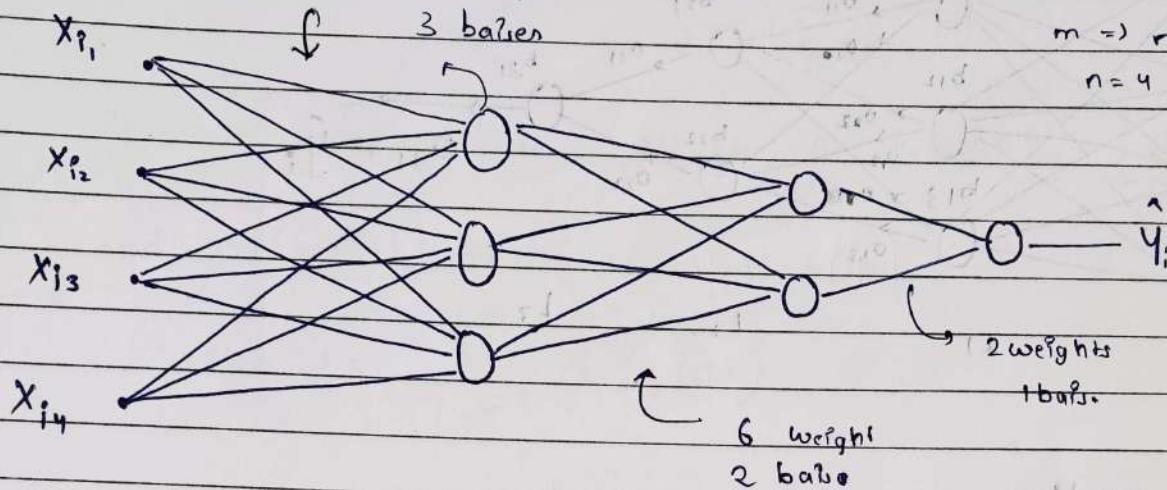
Ques 132

Ans 132

Ques 133

Multi - Layer - Perceptron (MLP)

$$4 \times 3 = 12 \text{ weights}$$



Input layer (L_0)	hidden layer (L_1)	hidden layer (L_2)	output layer (L_3)
4	3	2	1

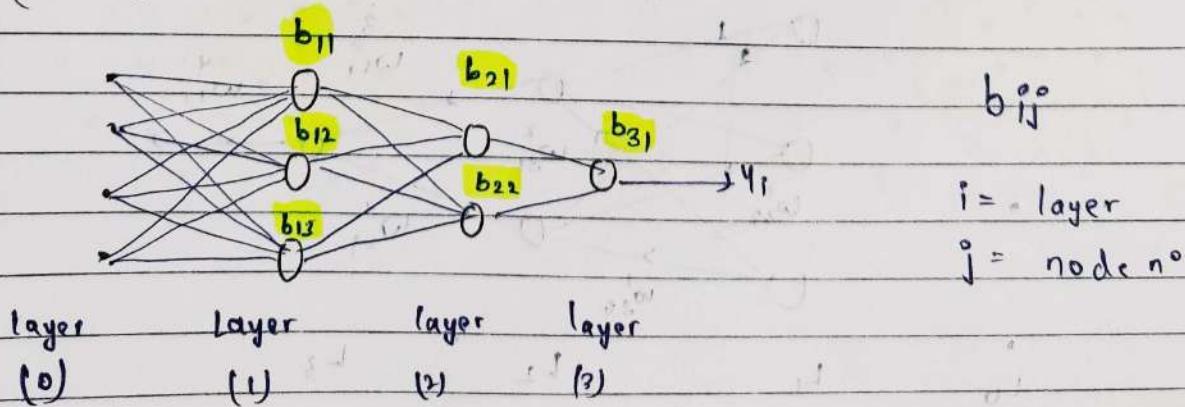
Input layer (L_0)	hidden layer (L_1)	hidden layer (L_2)	Output layer (L_3)
4	3	2	1

So In the above neural network structure there are

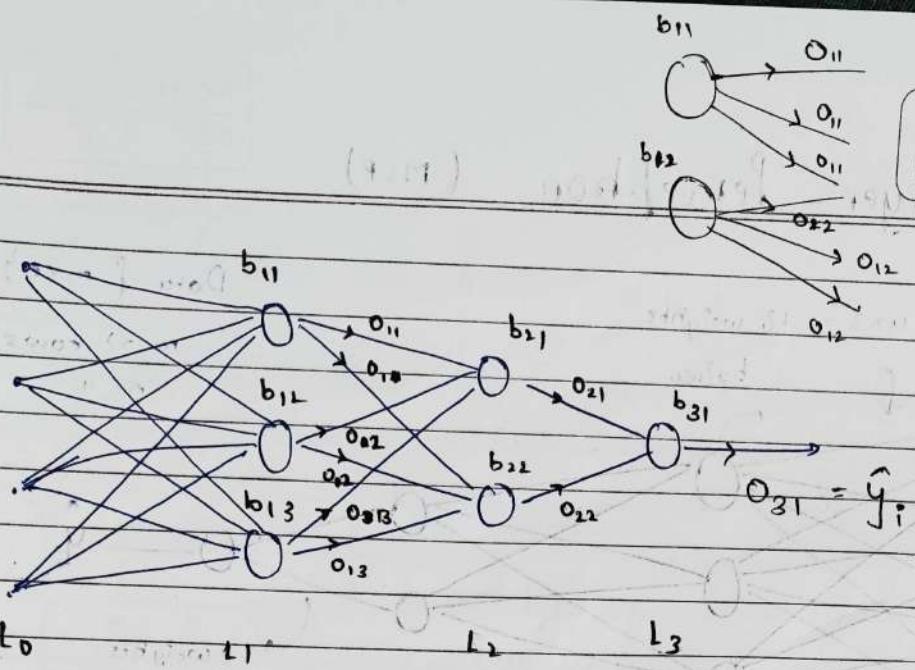
$$\text{Total weights} = (12+3) + (6+2) + (2+1)$$

(Total no. of trainable parameters.)

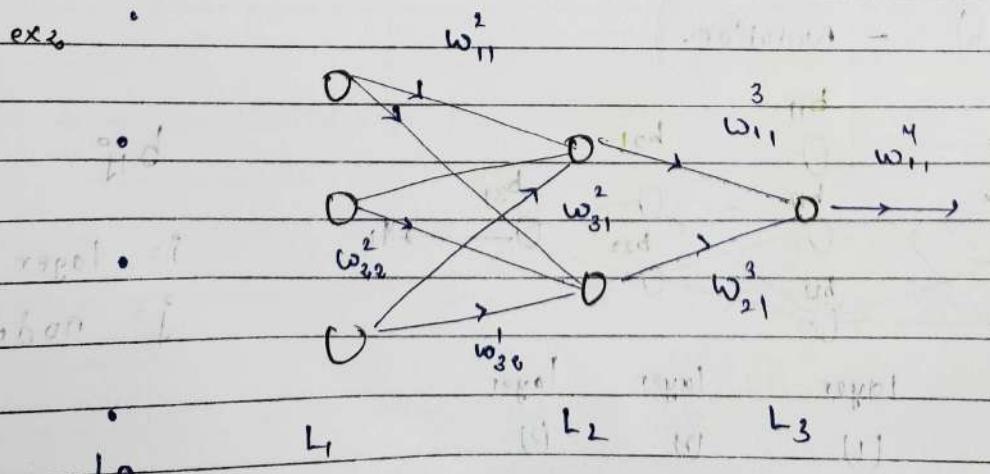
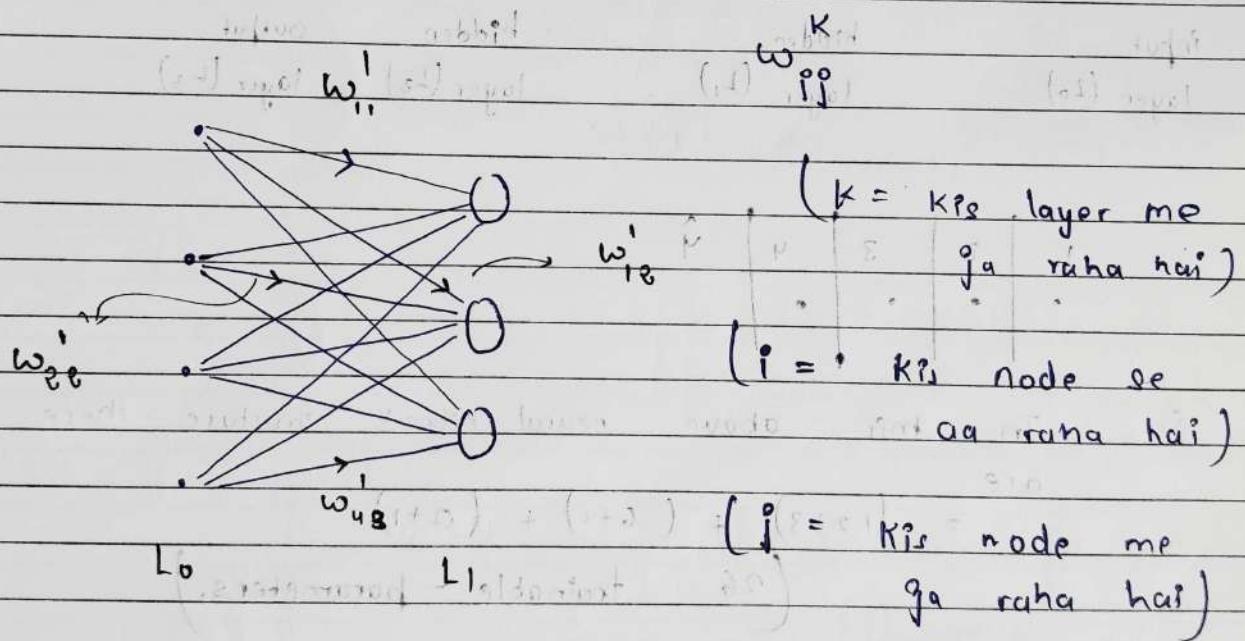
(ω, b) - Notation.



PAGE NO.:
DATE: / /



w - rotation.



MLP (intuition)

PAGE NO.:

DATE: / /

classification

$$\Sigma \cdot S \xrightarrow{\text{logistic}} A + P = \text{sigmoid}$$

] for MLP

for this perception we will get a probability.

$$P(x)$$

and

$$(1 - P(x))$$

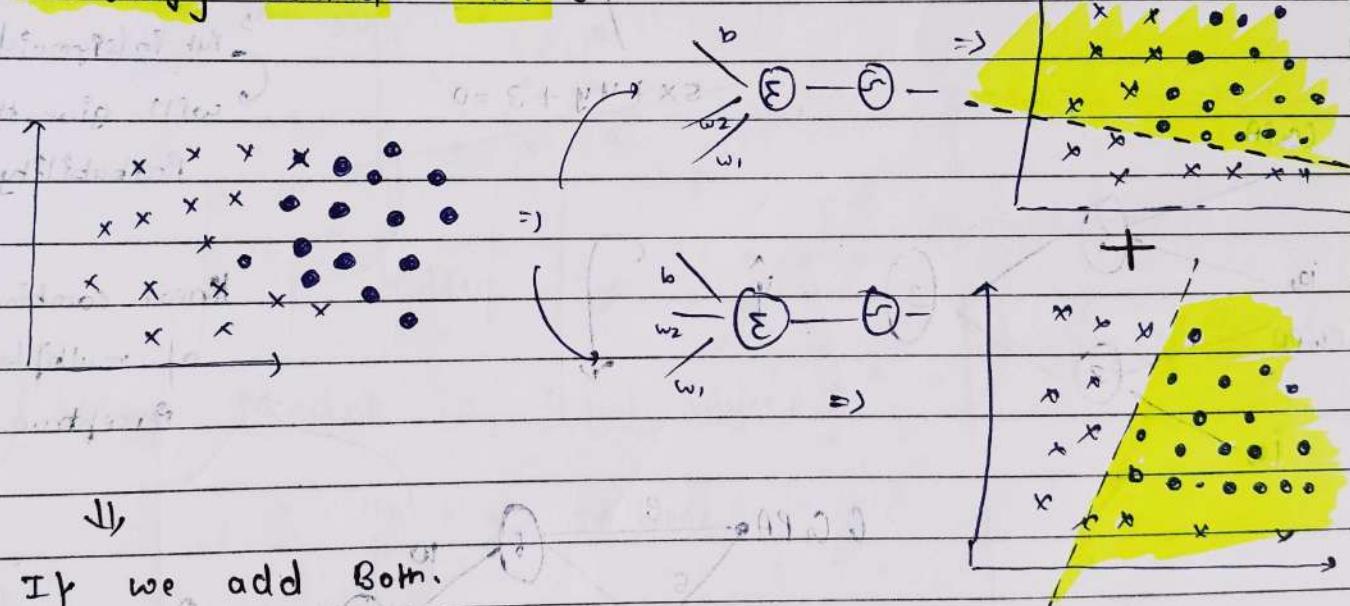
→ event not to happen.

to happen

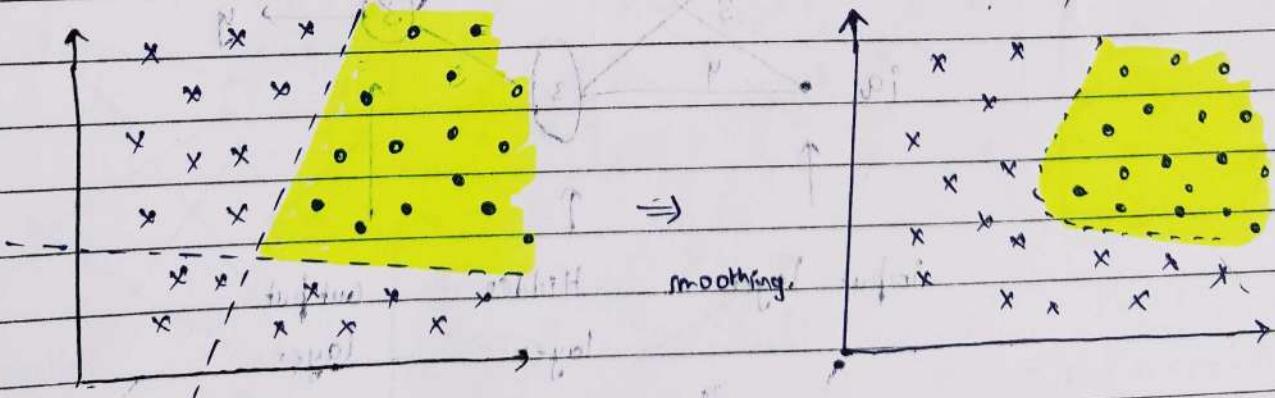
$$\text{Sigmoid} = \frac{1}{1 + e^{-z}} \Rightarrow (0 \dots 1)$$

* Idiology behind MLP :-

(Classification of data)



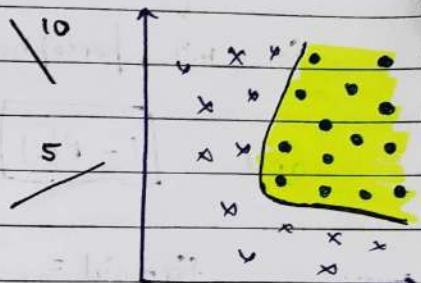
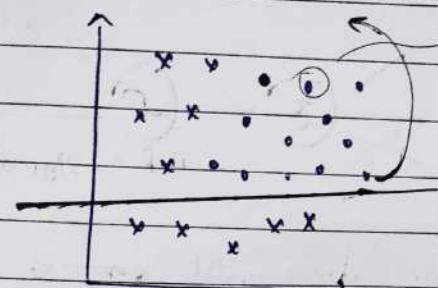
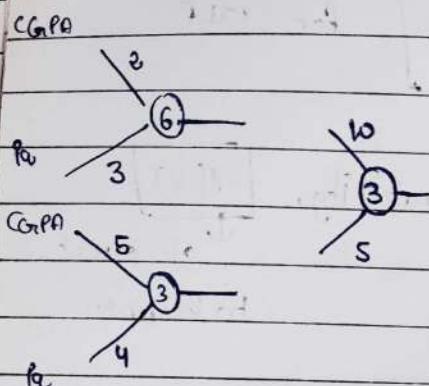
If we add Both.



smoothing.

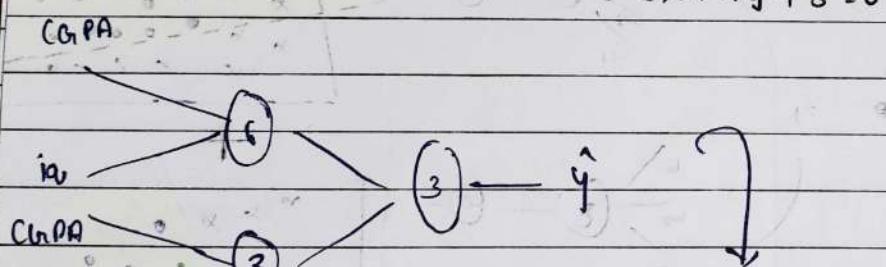
$$2x + 3y + 6 = 0$$

$$P = 0.6$$

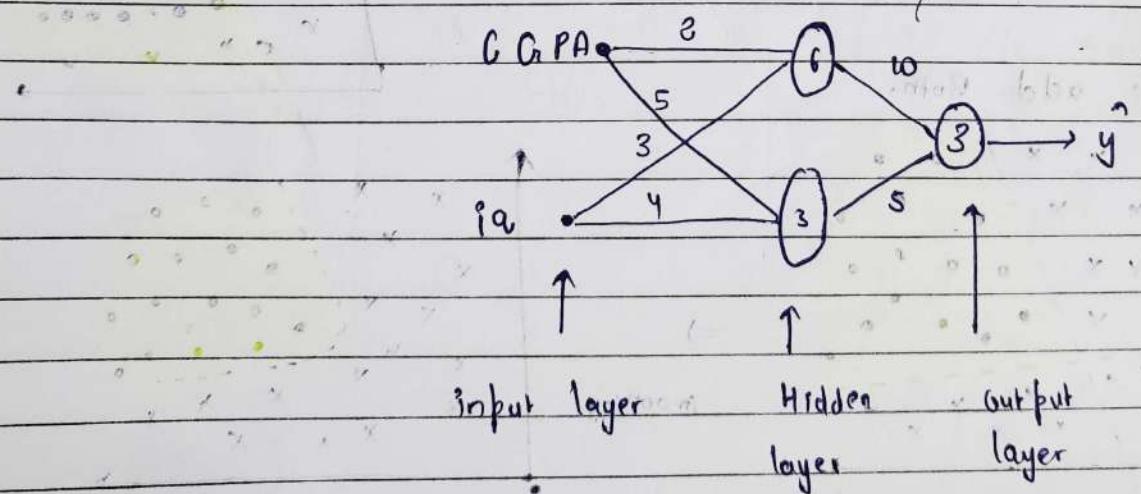


$$Z = 0.6 \times 10 + 0.8 \times 5$$

$\therefore Z = 10 + 3$
 Put in (sigmoid)
 will give the
 Probability



linear combination
of multiple
perceptrons

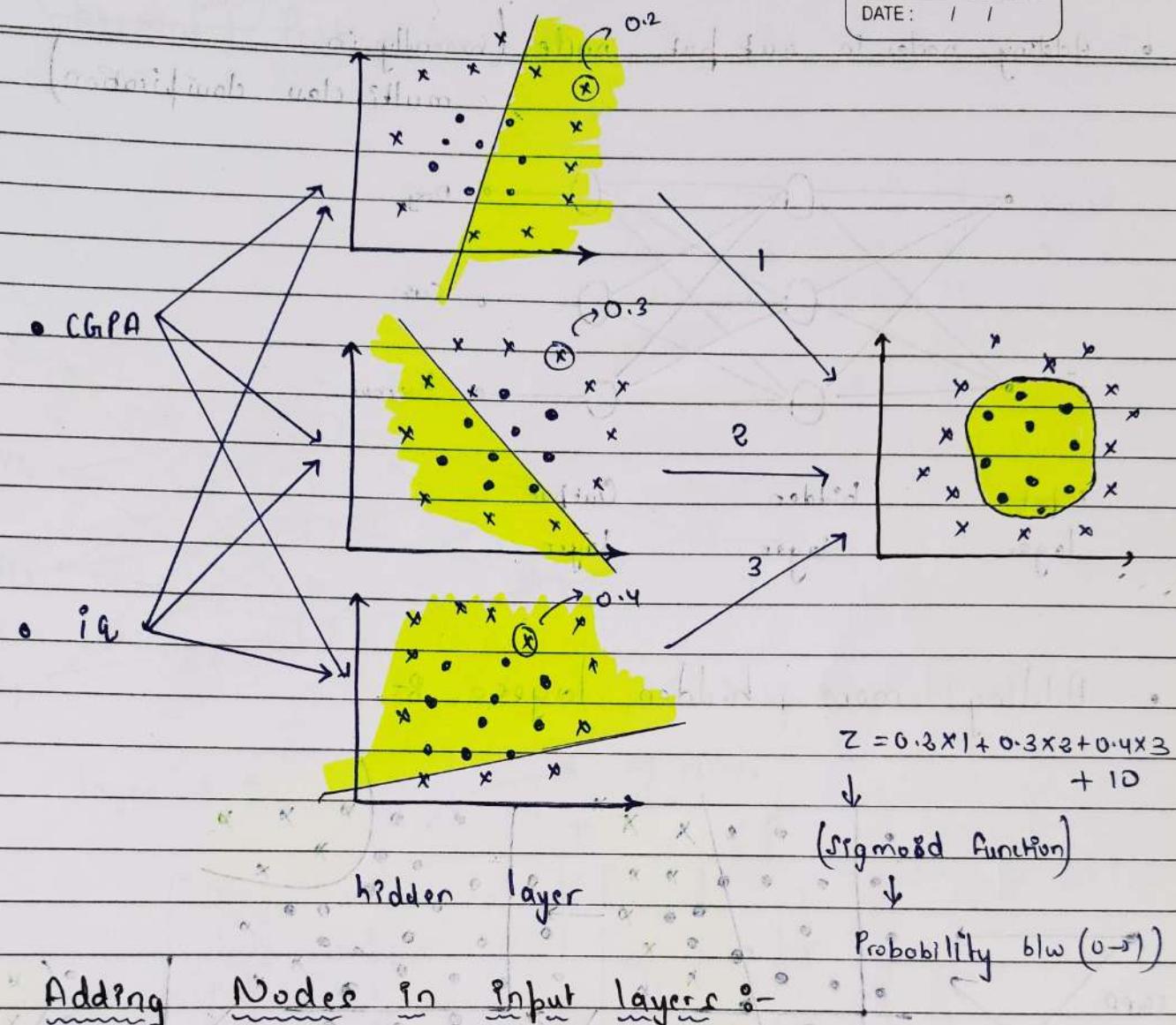


output of all nodes will probabilities.

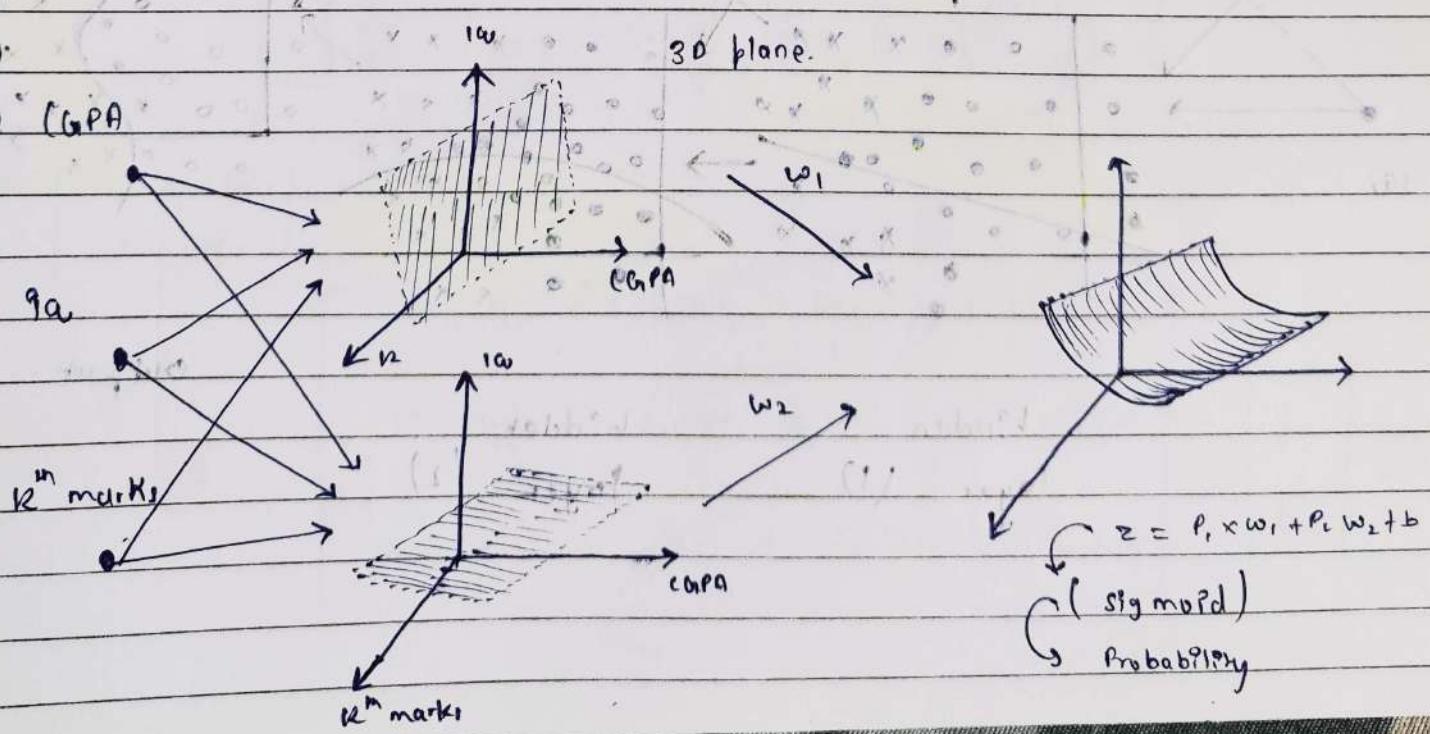
- Adding Nodes in hidden layers.

PAGE NO.:

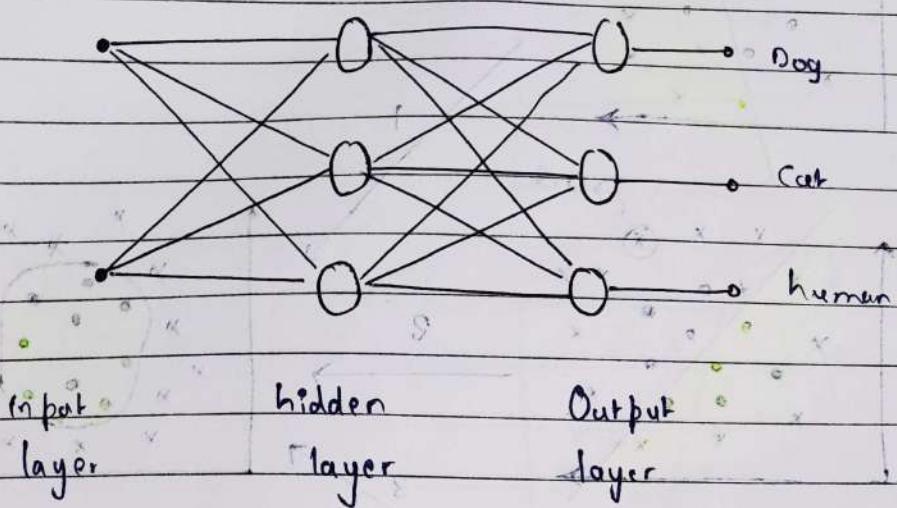
DATE: / /



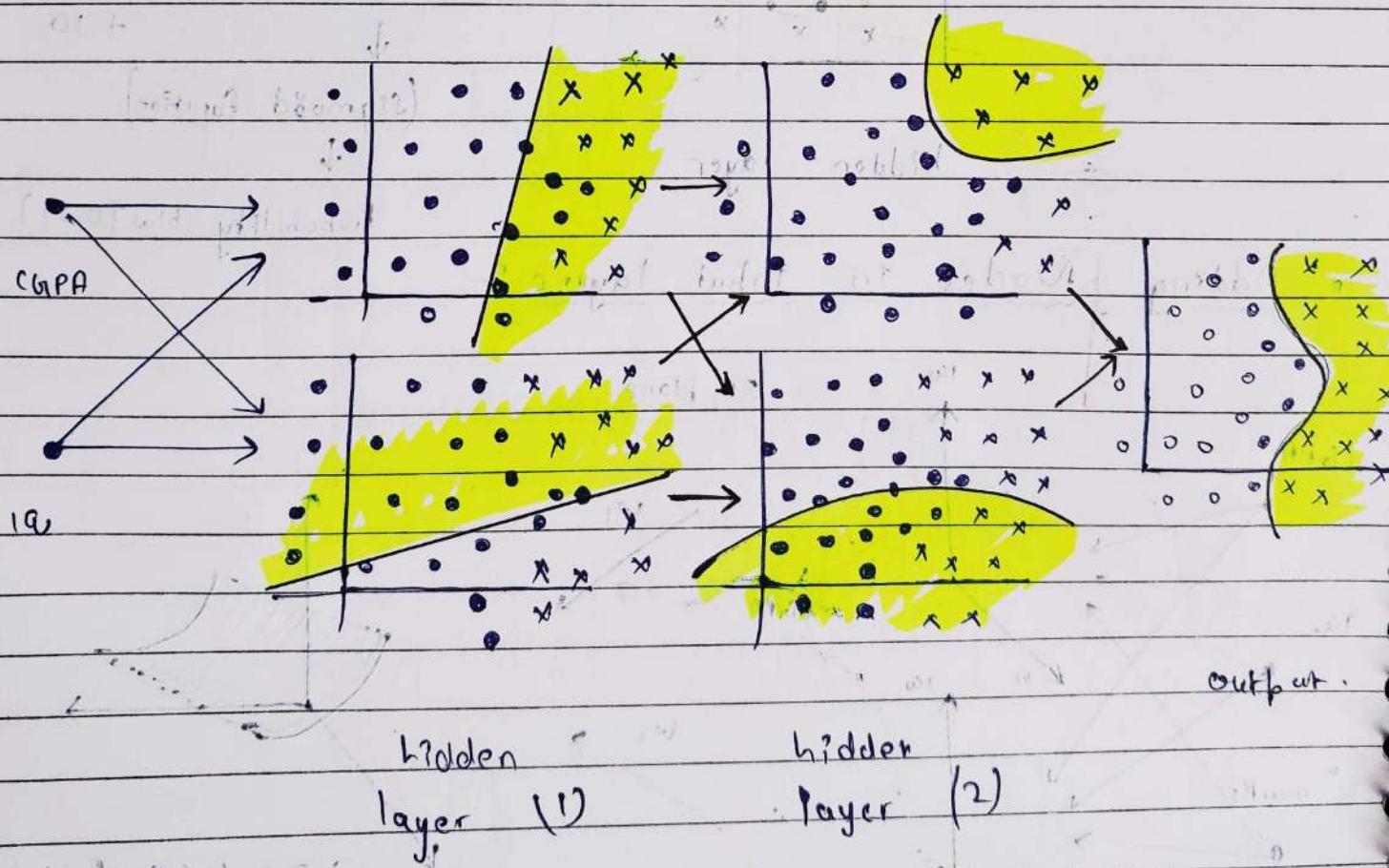
- Adding Nodes in input layers :-



- Adding nodes in output node. (generally in multi-class classification)



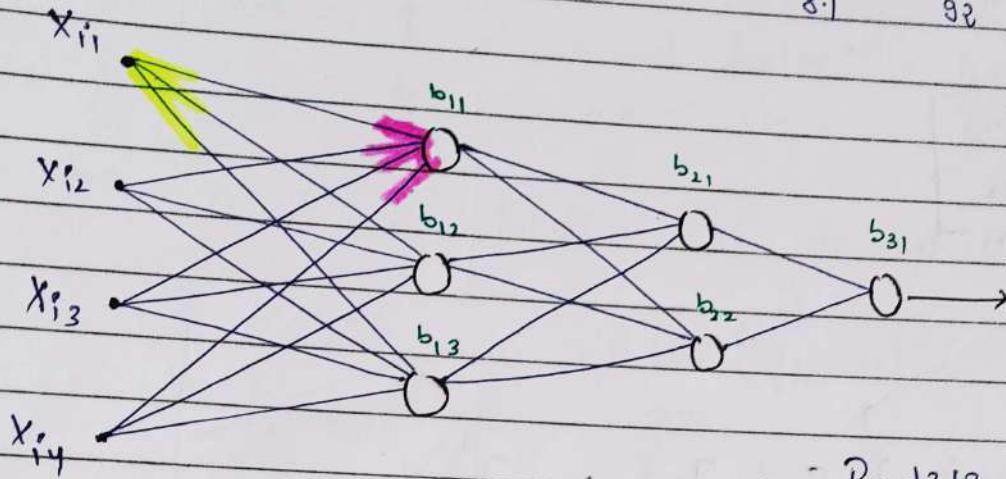
- Adding more hidden layers. 8-



forward Propagation :-

PAGE NO.: / /
DATE: / /

CGPA	19	10 th	12 th	blamed
T-2	72	69	81	*
8.1	92	75	76	1



$$\text{Prediction} \Rightarrow \sigma(\mathbf{w}^T \mathbf{x} + b) \Rightarrow$$

layer 1 #

[all entering] \mathbf{x}_{i1}

T

$$\begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix} \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

ux₁ 3x₁

[all entering]
b₁₁ node

$$4x_3 \rightarrow 3x_4$$

$$\begin{bmatrix} w_{11}^1 x_{i1} + w_{21}^1 x_{i2} + w_{31}^1 x_{i3} + w_{41}^1 x_{i4} \\ w_{12}^1 x_{i1} + w_{22}^1 x_{i2} + w_{32}^1 x_{i3} + w_{42}^1 x_{i4} \\ w_{13}^1 x_{i1} + w_{23}^1 x_{i2} + w_{33}^1 x_{i3} + w_{43}^1 x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$\sigma \left(\begin{bmatrix} w_{11}' x_{11} + w_{21}' x_{12} + w_{31}' x_{13} + w_{41}' x_{14} + b_{11} \\ w_{12}' x_{11} + w_{22}' x_{12} + w_{32}' x_{13} + w_{42}' x_{14} + b_{12} \\ w_{13}' x_{11} + w_{23}' x_{12} + w_{33}' x_{13} + w_{43}' x_{14} + b_{13} \end{bmatrix} \right) = \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix}$$

layer (2) #

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix} \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$$\sigma \left(\begin{bmatrix} w_{11}^2 o_{11} + w_{21}^2 o_{12} + w_{31}^2 o_{13} + b_{21} \\ w_{12}^2 o_{11} + w_{22}^2 o_{12} + w_{32}^2 o_{13} + b_{22} \end{bmatrix} \right) = \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix}$$

layer (3) #

$$\begin{bmatrix} w_{11}^3 \\ w_{21}^3 \end{bmatrix}^T \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix} + [b_{31}]$$

$$\sigma \left(\begin{bmatrix} w_{11}^2 o_{21} + w_{21}^2 o_{22} + b_{31} \end{bmatrix} \right) = \hat{y}_1$$

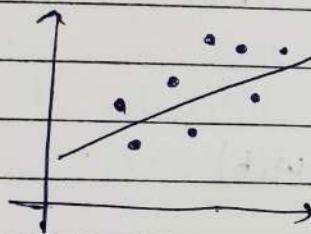
Loss function :-

PAGE NO.:
DATE: / /

→ loss function is a method of evaluating how well your algorithm is modelling your dataset.

? If $L.F \rightarrow \uparrow$ model_perform = poor
if $L.F \rightarrow \downarrow$ model_perform = great.

→ It is a function of parameters.



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

$L(m, b) =$ By changing (m, b) to make $L(m, b)$ min

to make $L(m, b)$ min

→ Why is Loss function important?

[You can't improve what you can't measure]

→ Peter Drucker

→ L.F is the eye on model ($\frac{dL}{d\theta}$) tells you where to go to correct your performance.

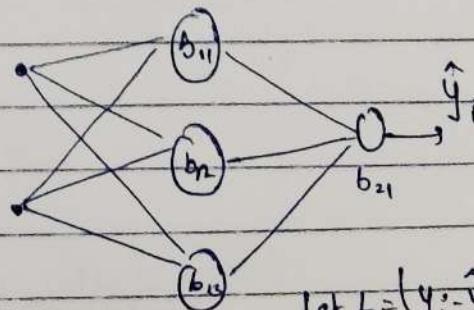
→ loss function in Deep learning?

Gpa	ia	Package (UPA)
-----	----	---------------

7.1	83	3.2
-----	----	-----

8.5	91	4.5
-----	----	-----

7.1	83	3.2
8.5	91	4.5
8.5	91	4.5



Starting with random (w, b)

→ we get a output \hat{y}_i
but it is in $(L \cdot F)$

$$\text{let } L \cdot F = (y_i - \hat{y}_i)^2 \text{ gives a no}$$

according to this, adjust the value of (w, b)
→ with the help of (gradient descent)

And the method will be repeated with other
row with multiple epochs.

And at end we will get min (w, b)

Loss function in DL

object detection
→ focal loss

Regression
- mse
- mae
- Huber loss

classification
- binary cross entropy
- categorical cross entropy
- entro

Auto encoder
- KL divergence

GAN
- discriminator loss
- Min Max gen loss

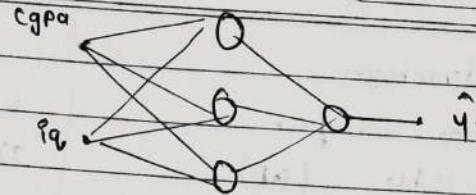
embedding
- Triplet loss

Cost Function (U/I) Loss function

PAGE NO.: / /
DATE: / /

Cgpa | y_i | \hat{y}_i | Package | Prediction

6.3	100	6.3	6.1
7.1	91	4.1	4
8.5	83	3.5	3.7
9.2	102	7.2	7



Loss fun. \rightarrow Single training eg.

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

$$= (6.3 - 6.1)^2$$

(cost function)

\rightarrow Is calculated for the whole Batch.

$$L.F. = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$L.F. = \frac{1}{4} ((6.1 - 6.3)^2 + (4.1 - 4)^2 + (3.5 - 3.7)^2 + (7.2 - 7)^2) = C.F.$$

\rightarrow Types of Loss functions :-

(1) mean-squared-error (mse) :-

$$L.F. = (y_i - \hat{y}_i)^2 \quad \text{why square}$$

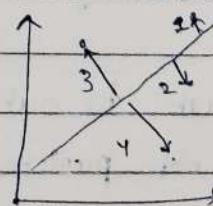
lets say point

$$y_i - \hat{y}_i = 1 \Rightarrow L.F. = 1$$

$$2 \Rightarrow L.F. = 4$$

$$3 \Rightarrow L.F. = 9$$

$$4 \Rightarrow L.F. = 16$$



as distance b/w pt and line matters

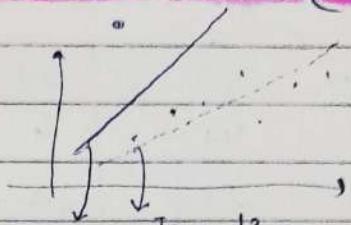
\therefore (further pt will be consider greater loss)

Advantages

- 1) easy to interpret
- 2) Differentiable (GD)
- 3) 1 local minima

Disadvantages

- 1) error isn't (squared) \rightarrow diff
- 2) Robust to outliers (Not)



Note

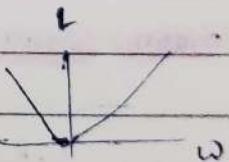
- for applying MSE in DL
- last node ~~the~~ activation function should be ["linear"]

$$\text{Cost Func} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

(2) Mean Absolute Error :-

$$L.F = |y_i - \hat{y}_i|$$

$$C.F = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



Advantage

- 1) Robust to outliers
(as not punishes extra to outliers)

Disadvantage

- 1) Not Differentiable

(so differentiation

is not possible)

Conclusion \Rightarrow more outliers \Rightarrow MSE

PAGE NO.:
DATE: / /

less outliers = MAE

(3) Huber loss :-

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq 8 \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

mse
huber
mae
mse
mse
mae
mse
mse

Note works best when there are very much outliers.

(4) Binary cross entropy :- (logloss)

→ classification

cepa | pa | placement

→ Two classes.

{ } ;

$$\text{Loss function} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Note

y → actual

\hat{y} → predicted

Activation function can only be

$$\text{(out. function} = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i) \right] \text{ (sigmoid)}$$

(5) Categorical Cross Entropy: [used in softmax Regression]

→ Multidate classification

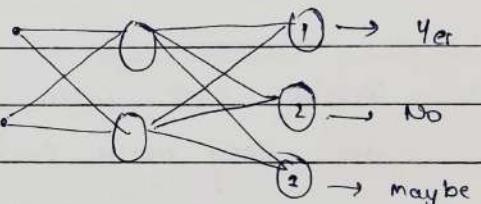
cgpa | placed | one-Hot encoding
80 Yes 1 0 0
60 No 0 1 0
70 maybe 0 0 1

$$\left\{ L = - \sum_{j=1}^k y_j \log(\hat{y}_j) \right\}$$

(K = no. of classes in data)

→ Activation function

should be (soft max)



⇒ Sparse categorical Cross Entropy.

cgpa | placed | Encoding

→ No one Hot encoding is done

→ Faster than CCE

7 70 Yes 1

8 80 No 2

⇒ CCE & SCCE

6 60 maybe 3

are same

But SCCE is faster when you have many output categories.

$$\left\{ L.P = - \sum_{i=1}^k y_i \log(\hat{y}_i) \right\} \quad \left\{ C = - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k y_{ij} \log(\hat{y}_{ij}) \right\}$$

Summary

→ Reg → mse
 } out pers — mae

→ Classification → binary → bce
 } multi → cces
 } scc

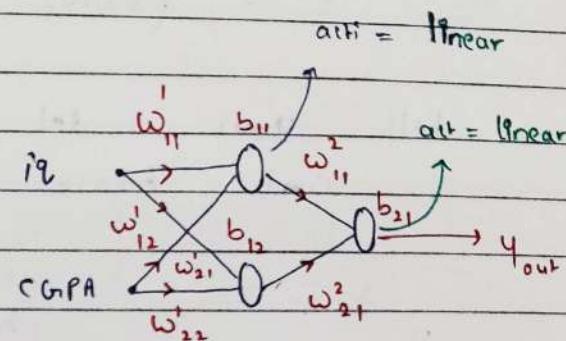
Back propagation :-

PAGE NO.:

DATE: / /

- Algo that helps to find the current value of weights and biases.
- forward propagation
- Gradient descent.

\rightarrow	ia	CGPA	LPA
80	8	3	
60	9	5	
70	5	8	
120	7	11	



→ Steps.

(1) Init w, b

↪ You can give random values.

↪ or you can initialize with $w \rightarrow 1, b \rightarrow 0$

(2) Select a point (row) → students.

(3) Predict (lpa) → forward prop [Dot product]

	ia	CGPA	LPA
①	80	8	3

→ let say prediction comes to be (18 lpa)

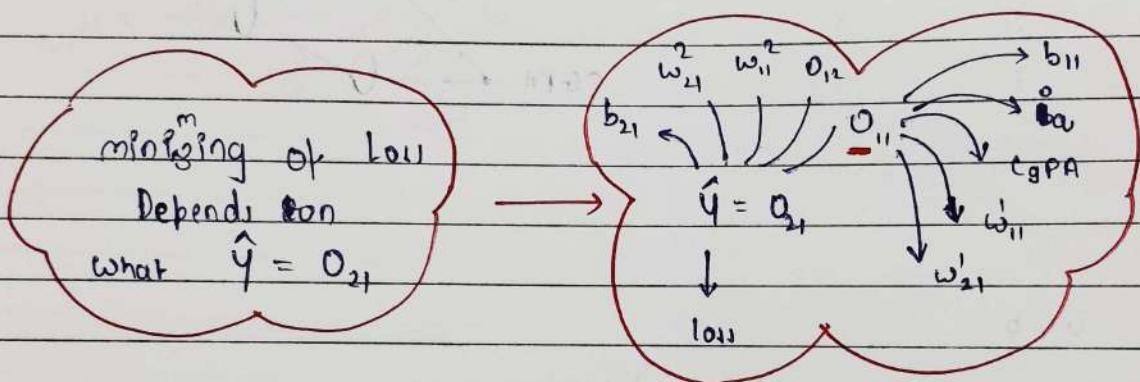
{ Big weights / Biases
are not correct. }

(3). choose a loss function

$$\text{for linear } L = (y - \hat{y})^2 \quad (\text{mse})$$

$$L = (3 - 18)^2 = 225 \quad \text{error}$$

tell your model to adjust the weights
and biases to reduce it.



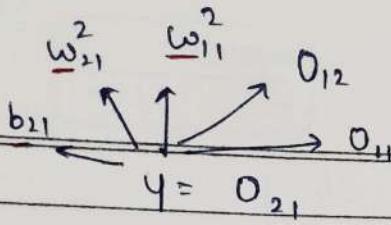
complex backward dependencies

→ to going backward and changing all parameters
is called Back propagation

(4). Update weights | B&B -

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

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



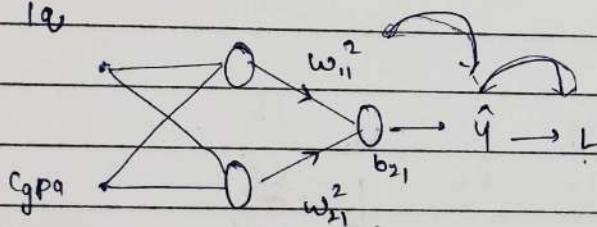
$$\check{w}_{11}^2_{\text{new}} = \check{w}_{11}^2_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{11}^2} \right) \quad \begin{matrix} \text{needs} \\ \text{to bp} \\ \text{calculated} \end{matrix}$$

$$\left(\check{w}_{21}^2_{\text{new}} = \check{w}_{21}^2_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{21}^2} \right) \right) \quad \left(b_{21}^{\text{new}} = b_{21}^{\text{old}} - \eta \left(\frac{\partial L}{\partial b_{21}} \right) \right)$$

Target
value

$$\left\{ \frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}} \right\} \quad \left\{ \frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}} \right\} \quad \left\{ \frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{21}^2}, \frac{\partial L}{\partial b_{21}} \right\}$$

i.e.



$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{q}} \times \frac{\partial \hat{q}}{\partial w_{11}^2}$$

\therefore (chain rule.)

$$\rightarrow \left\{ \frac{\partial L}{\partial \hat{q}} = \frac{\partial}{\partial \hat{q}} (Y - \hat{q})^2 = -2(Y - \hat{q}) \right\}$$

$$\left\{ \frac{\partial \hat{q}}{\partial w_{11}^2} = \frac{\partial}{\partial w_{11}^2} (O_{11} w_{11}^2 + O_{12} w_{21}^2 + b_{21}) \right\}$$

$$= O_{11}$$

So

$$\left\{ \frac{\partial L}{\partial w_{11}^2} = -2(Y - \hat{q}) O_{11} \right\}$$

Similarly,

$$\left\{ \frac{\partial L}{\partial w_{12}^2} = -2(Y - \hat{q}) O_{12} \right\}$$

$L \rightarrow \hat{q} \rightarrow s$

PAGE NO.: _____
DATE: / / / /

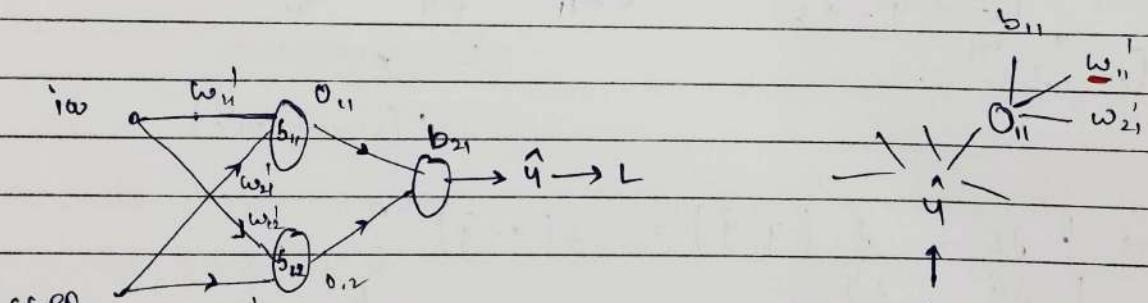
$$\frac{\partial L}{\partial \hat{q}} = \frac{d}{d\hat{q}} (4 - \hat{q})^2 = -2(4 - \hat{q})$$

$$\frac{\partial \hat{q}}{\partial b_{21}} = \frac{d}{db_{21}} [0_{11} \omega_{11}^2 + 0_{12} \omega_{21}^2 + b_{21}] = 1$$

$$\boxed{\frac{\partial L}{\partial b_{21}} = -2(4 - \hat{q})}$$

→ we have calculated all the 3 derivatives in first block.

② $\left(\frac{\partial L}{\partial \omega_{11}}, \frac{\partial L}{\partial \omega_{21}}, \frac{\partial L}{\partial b_{11}} \right)$



$$\frac{\partial L}{\partial \omega_{11}} = \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial \omega_{11}}$$

$$\frac{\partial L}{\partial \hat{q}} = -2(4 - \hat{q})$$

$$\frac{\partial \hat{q}}{\partial \omega_{11}} = \frac{d}{d\omega_{11}} (0_{11} \omega_{11}^2 + 0_{12} \omega_{21}^2 + b_{21})$$

$$= \omega_{11}^2$$

$$\frac{\partial q_1}{\partial \omega_{11}} = \frac{d}{d\omega_{11}} (\omega_{11}^2 (f\omega) + \omega_{21}^2 (CPA) + b_{21}) = (9a)$$

$$\left(\frac{dL}{d\omega_{11}^1} = -2(4-\hat{q}) \cdot \omega_{11}^2 \cdot (ia) \right)$$

$$\left(\frac{dL}{d\omega_{21}^1} = -2(4-\hat{q}) \cdot \omega_{21}^2 \cdot (CGPA) \right)$$

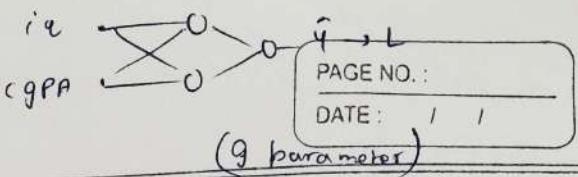
$$\left(\frac{dL}{db_{12}} = -2(4-q) \cdot \omega_{12}^2 \cdot (1) \right)$$

(3) $\left(\frac{dL}{d\omega_{12}^1}, \frac{dL}{d\omega_{22}^1}, \frac{dL}{db_{12}} \right)$

$$\frac{dL}{d\omega_{12}^1} = \frac{dL}{dG} \cdot \frac{d\hat{q}}{dO_{12}} \cdot \frac{dO_{12}}{d\hat{\omega}_{12}} = -2(4-\hat{q}) \cdot \omega_{21}^2 \cdot (ia)$$

$$\frac{dL}{d\omega_{22}^1} = \frac{dL}{d\hat{q}} \cdot \frac{d\hat{q}}{dO_{12}} \cdot \frac{dO_{12}}{d\omega_{22}^1} = -2(4-\hat{q}) \cdot \omega_{21}^2 \cdot (CGPA)$$

$$\frac{dL}{db_{12}} = \frac{dL}{d\hat{q}} \cdot \frac{d\hat{q}}{dO_{12}} \cdot \frac{dO_{12}}{db_{12}} = -2(4-\hat{q}) \cdot \omega_{21}^2 \cdot (1)$$



loop → (100 times) epochs.

for i in Range (4) :

$$(\omega \rightarrow 1) \quad (b \rightarrow 0)$$

(a) 1st student \rightarrow forward prop \rightarrow predict (LPA)

(b) Loss calculate (mse) →

(c) adjust all weights and bias.

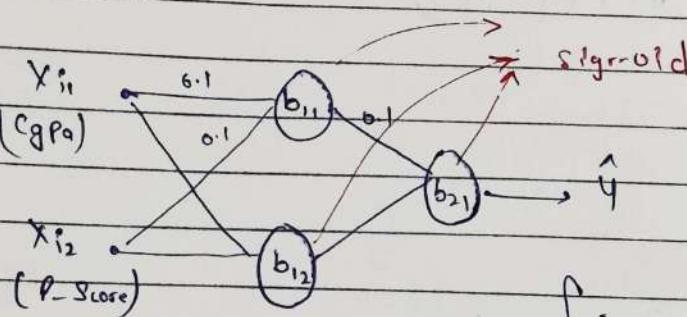
$$\left(\omega/b \right)_{\text{new}} = \left(\omega/b \right)_{\text{old}} + \eta \frac{\partial L}{\partial \omega_{\text{old}}} \quad \boxed{\downarrow}$$

→ Calculate avg loss for the epoch (9 times)

Back Propagation - Classification

PAGE NO.:

DATE: / /

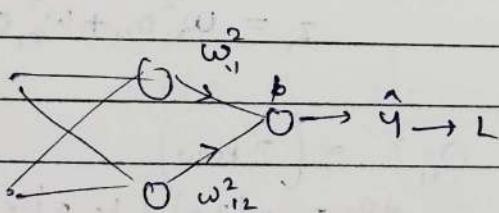


All three nodes have activation func. = (sigmoid.)

$$\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Cg Pa	P-score	Placed.
8	8	0
9	7	1
6	4	0
3	6	0
2	1	0

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_{11}^2} = -(4-9) O_n$$



$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_{21}^2} = -(4-9) O_n$$

$$\frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial b_{21}} = -(4-9)$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{d}{d\hat{y}} [-y \log(\hat{y}) - (1-y) \log(1-\hat{y})]$$

$$= -\frac{y}{\hat{y}} - \frac{(1-y)}{1-\hat{y}} = -\frac{y + y\hat{y} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})} = \frac{-(4-9)}{\hat{y}(1-\hat{y})}$$

$$\hat{y} = \sigma(z)$$

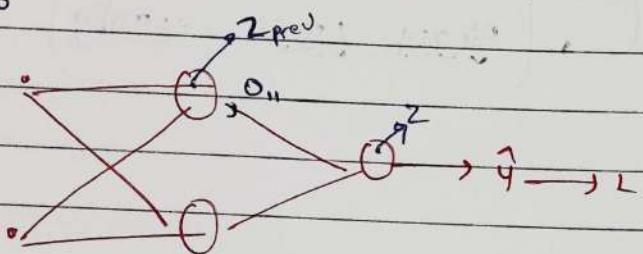
PAGE NO.: / /
DATE: / /

$$\frac{\partial \hat{y}}{\partial z} = \frac{d}{dz} (\sigma(z)) = \sigma(z) (1 - \sigma(z)) = \hat{y} (1 - \hat{y})$$

$$\rightarrow \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} = -(y - \hat{y})$$

$$z = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21}$$

→ now



$$\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{21}}, \frac{\partial L}{\partial b_{21}}$$

$$(L \rightarrow \hat{y} \rightarrow z_p \rightarrow o_{11} \rightarrow z_{prev} \rightarrow w_{11})$$

$$z_p = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21}$$

$$o_{11} = \sigma(z_{prev})$$

$$\delta o_{11} = \sigma'(z_p) [1 - \sigma(z_p)]$$

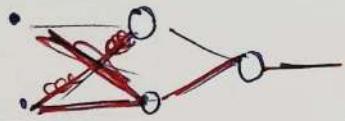
$$\left(\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_p} \cdot \frac{\partial z_p}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial z_{prev}} \cdot \frac{\partial z_{prev}}{\partial w_{11}} \right)$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $(y - \hat{y}) \quad \omega_{11}^2 \quad o_{11}(1 - o_{11}) \quad x_{11}$

$$\frac{\partial L}{\partial w_{11}} = -(y - \hat{y}) \omega_{11}^2 o_{11} (1 - o_{11}) x_{11}$$

$$\frac{\partial L}{\partial w_{21}} = -(y - \hat{y}) \omega_{21}^2 o_{12} (1 - o_{12}) x_{12}$$

$$\frac{\partial L}{\partial b_{21}} = -(y - \hat{y}) \omega_{21}^2 o_{12} (1 - o_{12})$$



$$\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_{12}}$$

PAGE NO.:

DATE: / /

$$\frac{\partial^2 V}{\partial \theta_{12}} = \omega_{11}^2 \theta_{11} + \omega_{21}^2 \theta_{12} + b_{12}$$

$$\frac{\partial L}{\partial w_{12}} = \underbrace{\frac{\partial L}{\partial \dot{y}}}_{-(y-\hat{y})} \times \underbrace{\frac{\partial \dot{y}}{\partial \theta_{12}}}_{\omega_{21}^2} \times \underbrace{\frac{\partial \theta_{12}}{\partial \theta_{12}}}_{1} \times \underbrace{\frac{\partial \theta_{12}}{\partial \dot{\theta}_{12}}}_{x_{i_1}} \times \underbrace{\frac{\partial \dot{\theta}_{12}}{\partial w_{12}}}_{\theta_{12}(1-\theta_{12})}$$

$$L \rightarrow \dot{y} \rightarrow \ddot{y} \rightarrow \theta_{12}$$

$$\downarrow \\ w_{12} \leftarrow \ddot{y}$$

$$\boxed{\frac{\partial L}{\partial w_{12}} = -(y-\hat{y}) \omega_{21}^2 \theta_{12} (1-\theta_{12}) x_{i_1}}$$

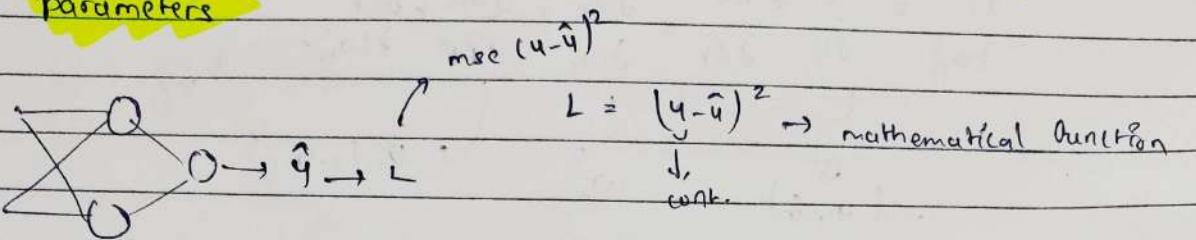
$$\boxed{\frac{\partial L}{\partial w_{22}} = -(y-\hat{y}) \omega_{21}^2 \theta_{12} (1-\theta_{12}) x_{i_2}}$$

$$\boxed{\frac{\partial L}{\partial b_{12}} = -(y-\hat{y}) \omega_{21}^2 \theta_{12} (1-\theta_{12})}$$

\rightarrow code is done as well.

→ The Why. (Back propagation).

→ Loss function is a function of all trainable parameters



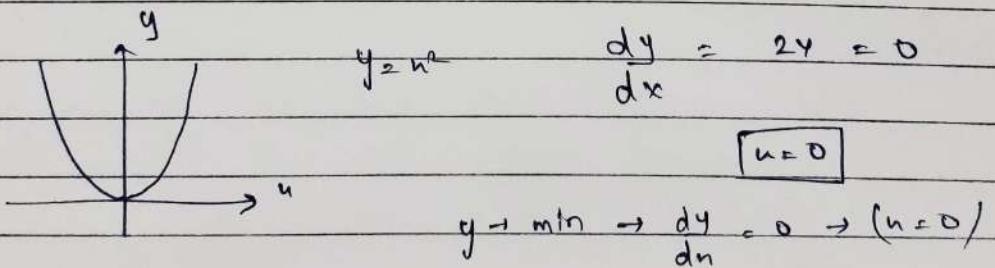
$$L(?) \rightarrow L(\hat{y})$$

$$\hat{y} = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{12}$$

$$\hat{y} = w_{11}^2 [w_{11}' x_{11} + w_{21}' x_{12} + b_{11}] + w_{21}^2 [w_{12}' x_{11} + w_{22}' x_{12} + b_{12}] + b_{12}$$

$\hat{y} \rightarrow$ depends on all trainable parameters

→ Concept of minima :-

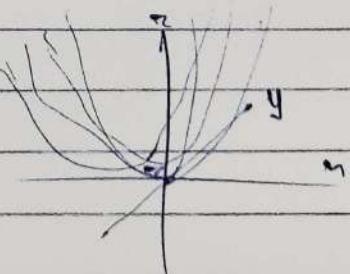


$$\text{Let } z = u^2 + y^2$$

$$(u, y) \rightarrow z \min$$

$$\frac{\partial z}{\partial u} = 0 \quad 2u = 0 \Rightarrow (u = 0)$$

$$\frac{\partial z}{\partial y} = 0 \quad 2y = 0 \Rightarrow (y = 0)$$



now

 $L \rightarrow$ function of g parameters g dimensional \rightarrow minima.

$$\left\{ \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_g} = 0 \right\}$$

- Back Propagation Assumption.

$$w_{\text{new}} = w_{\text{old}} - \gamma \left[\frac{\partial L}{\partial w_i} \right]$$

↓ small learning rate
make the transition
smooth.

MLI Memorization.

PAGE NO.:

DATE: / /

A technique used to primarily speed up computer program by storing the results of expensive function calls and returning the cached result when the same input occurs again.

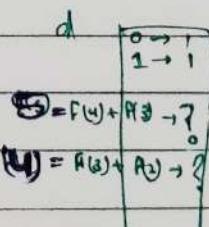
e.g. Fibonacci Series

```
→ def fib(n):  
    if n == 1 or n == 0:  
        return 1  
    else:  
        return fib(n-1) + fib(n-2)
```

But this code is very inefficient, Time to calc.
Big Fibonacci no. will reach ∞

Instead

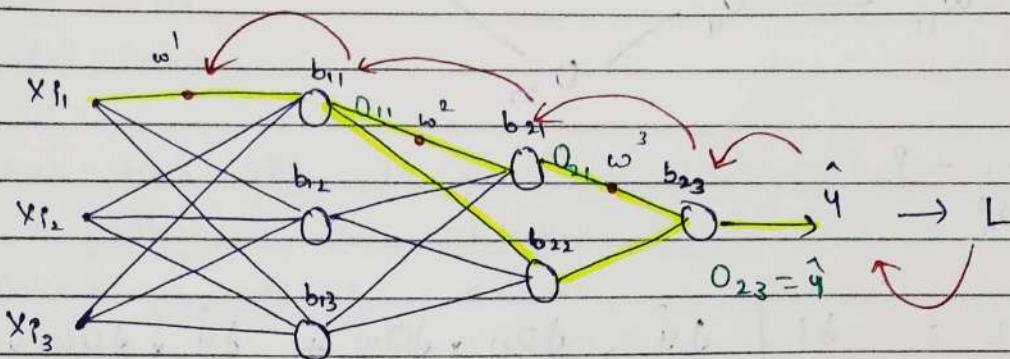
```
→ def fib(n, d):  
    if n in d:  
        return d[n]  
  
    else:  
        d[n] = fib(n-1) + fib(n-2)  
    return d[n]
```



```
→ d = {0: 1, 1: 1}  
fib(s, d) — call
```

• Small topic of Backpropagation

→ How complicated it can be w' in 3 layer in B.P.



chain rule.

$$(L \rightarrow \hat{y} \rightarrow O_{21} \rightarrow w_{11})$$

All the weights of second layer can be calculated
by m?

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{11}}$$

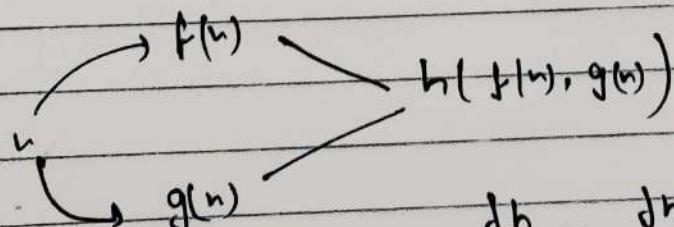
$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}} \end{array} \right.$$

→ for 1st layer:

$$L \rightarrow \hat{y} \rightarrow O_{21} \rightarrow O_{11} \rightarrow w_{11}$$

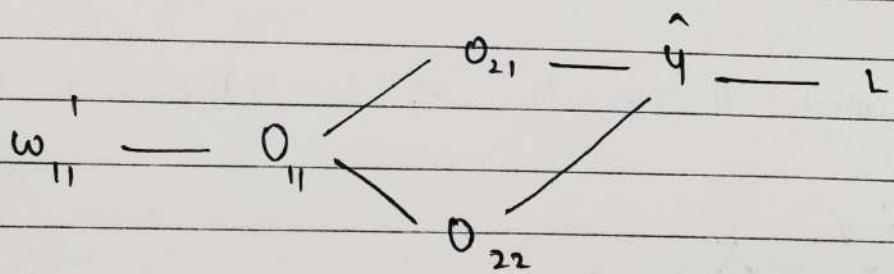
Additional complexity is there due to two path.

→ concept



$$\frac{\partial h}{\partial n} = \frac{\partial h}{\partial f(n)} \cdot \frac{\partial f(n)}{\partial n} + \frac{\partial h}{\partial g(n)} \cdot \frac{\partial g(n)}{\partial n}$$

The path is.



$$\frac{dL}{dw_{11}} = ?$$

$$\frac{dL}{dw_{11}} = \frac{dL}{d\hat{y}} \left[\frac{d\hat{y}}{dO_{11}} \times \frac{dO_{11}}{dw_{11}} + \frac{d\hat{y}}{dO_{22}} \times \frac{dO_{22}}{dw_{11}} \right]$$

→ expression complicated and complicated.
So, the derivatives which are already calculated
are stored so that it need not to
be calculated again.

Bulk prop → (chain diff rule) + memoization



maths

Keras.

Types of Gradient Descent.

- Batch
- Stochastic
- mini Batch.

→ These three differ in how much data we use to compute the gradient of the objective function

→ Depending on the amount of data, we make a trade-off b/w the accuracy of the parameter update and time it takes to perform an update.

Batch G.D (Vanilla)

- For i in range nb-epoches:
Param.grad = evaluate_grad(loss-func,
data, params)

$$\text{Params} = \text{Params} - \text{learning_Rate} * \text{Param.grad}$$

- entire dataset \rightarrow update.

$$\text{no.of epochs} = \text{no.of update.}$$

Stochastic G.D

for i in range (nb-epoches):
np.random.shuffle(data)
for example in data:
Param.grad = evaluate_gradient
(loss-func, example
, params)

Params = Params - learning_rate *
Paramgrad

Note. → In stochastic G.D. The freq of weights update is very much higher.

If, $\text{Batch size} = n$



Batch G.D

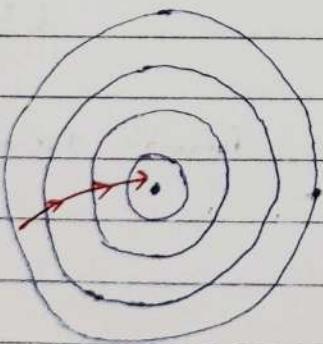
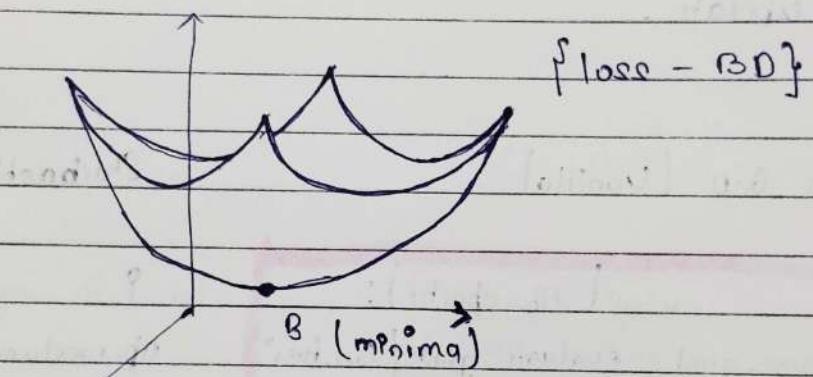
(faster)
output

Batch size = 1



Stochastic G.D

(converge faster)



Batch G.D



Stochastic G.D.

- This behaviour helps algo to move out of local minima.
- Do not converges exactly.

Mini Batch Gradient Descent :-

- Batch Size $\Rightarrow (1 \leftrightarrow n)$
(in between)

($n = \text{no. of rows}$)

BGD \longleftrightarrow SGD

320 rows

in every epoch
10 batches

batches

32

no updates

Fast :-

Batch > mini-Batch \Rightarrow smooth

Convergence

Batch < mini-Batch \Rightarrow stochastic

• Vanishing Gradient Problem in ANN. (with code)

PAGE NO.: / /
DATE: / /

Thing to remember

1. $0.1 \times 0.1 \times 0.1 \times 0.1 \Rightarrow 0.0001$

{ no smaller than 1 multiplied by
each other given even smaller no
more }

2. VGP]

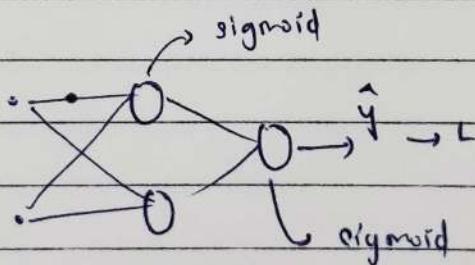
only found in Deep NN.
(many layers)

3 for VGP to occur

A-Func should be [sigmoid | tanh]

→ Intuition of VGP in Deep ANN.

→ Back prop.



$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

let's say we need to find.

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial \hat{y}_j} \times \frac{\partial \hat{y}_j}{\partial z_i} + \frac{\partial L}{\partial o_{ii}} \times \frac{\partial o_{ii}}{\partial w_{ij}}$$

$\downarrow \quad \uparrow$
 $0 < \hat{y} < 1 \quad 0 < o < 1$

let $w_0 = 1$

L.R

$$w_n = 1 - 0.01 \times 0.0001$$

$$w_n = 0.9999.$$

PAGE NO.:

DATE: / /

did not
change

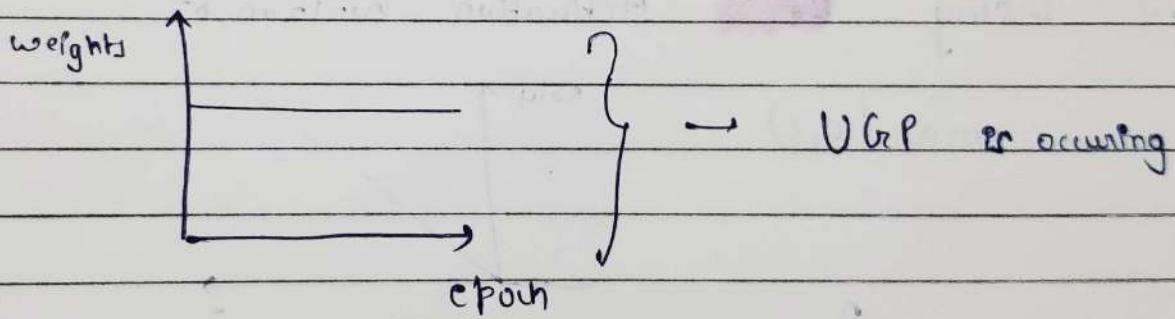
* *

So in very deep NN the initial layer weights will not change very much and then the model will not converge

How to recognize?

If loss after every epoch \rightarrow no change \rightarrow UGr

if graph of weights (W_k) epoch

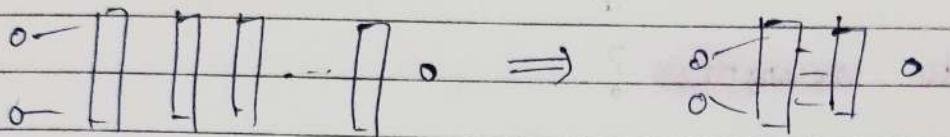


weight is not changing after every epoch

→ How to handle vanishing gradient problem.

1.7 Reduce model complexity

(if data do not consist of non linear complexity)

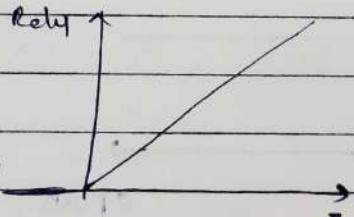


(not applicable)

$\frac{\partial L}{\partial w}$ bigger

3.7 Using ReLU Activation Function :-

$$\max(0, z)$$



$$\begin{cases} -1000 & z \leq 0 \\ \psi(z) & z = 0 \\ 1000 & z \geq 0 \end{cases}$$

→ so derivative in ReLU
 $\in (0, 1)$

$$\rightarrow |x_1 x_1 \dots| = 1 \quad \left\{ \begin{array}{l} \text{VGP don't} \\ \text{occur.} \end{array} \right\}$$

→ if "0" then
 use leaky relu

3.7 Proper weights init
→ (Glorot)
→ (Xavier)

(4) Batch normalization. \Rightarrow (layer)

(5) Residual network.

• How to improve a neural network.

1. Fine tuning of NN Hyperparameters.

- hidden layers
- neurons per layer
- learning rate
- optimizers
- Batch size.
- Activation fun.
- epochs.

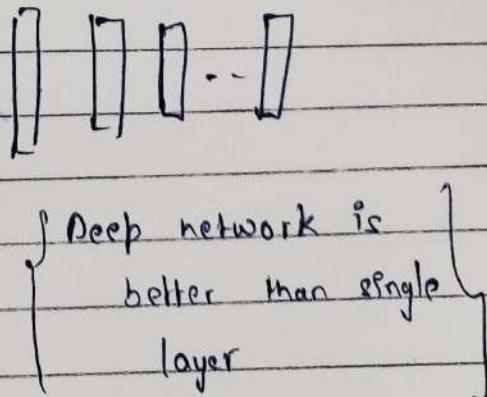
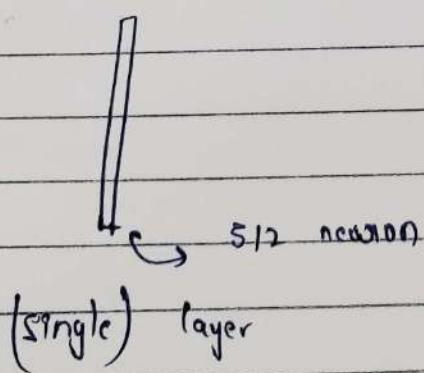
2. - Vanishing / Exploding gradient

- Not enough data

- Slow training

- overfitting

• no^o of Hidden layers



- How much layers \Rightarrow { until you get to over [31 layers]}

- Ht / trai of layers.



Deep network.

- first few layer recognize lines.
 - after that few layer recognize pattern like eyes / nose / cars.
 - The often that layers recognizes full faces -
- That's why multilayer network is better.

no. of neurons.

→ must be sufficient per layer. (not less)

Batch size.

- 3 approaches

- Small Batch size (8 to 32) } - slow
generalize better results. But give results.

- Large Batch size (approx 8000) } - fast.

To make it give good results, one thing
you can do is →

→ warming up the learning rate.

- initially take small values of (L.R)

- after every epoch increase (L.R)

→ This process will give good result in large batch size.

↳ { Try this technique first. }

- no of epochs.

→ Take as big as you want.

→ But apply early stopping

↓
Keras →

* [• when results are not getting better
or changing then keras stops
the training → applying early stopping]

Keras → (call back)

Problem with Neural net work.

Vanishing /
exploding gradients

- weights init
- Activation function
- Batch normalization
| for exploding
- gradient clipping

over fitting

- L_1 and L_2 reg
- Dropouts.

Not enough
data

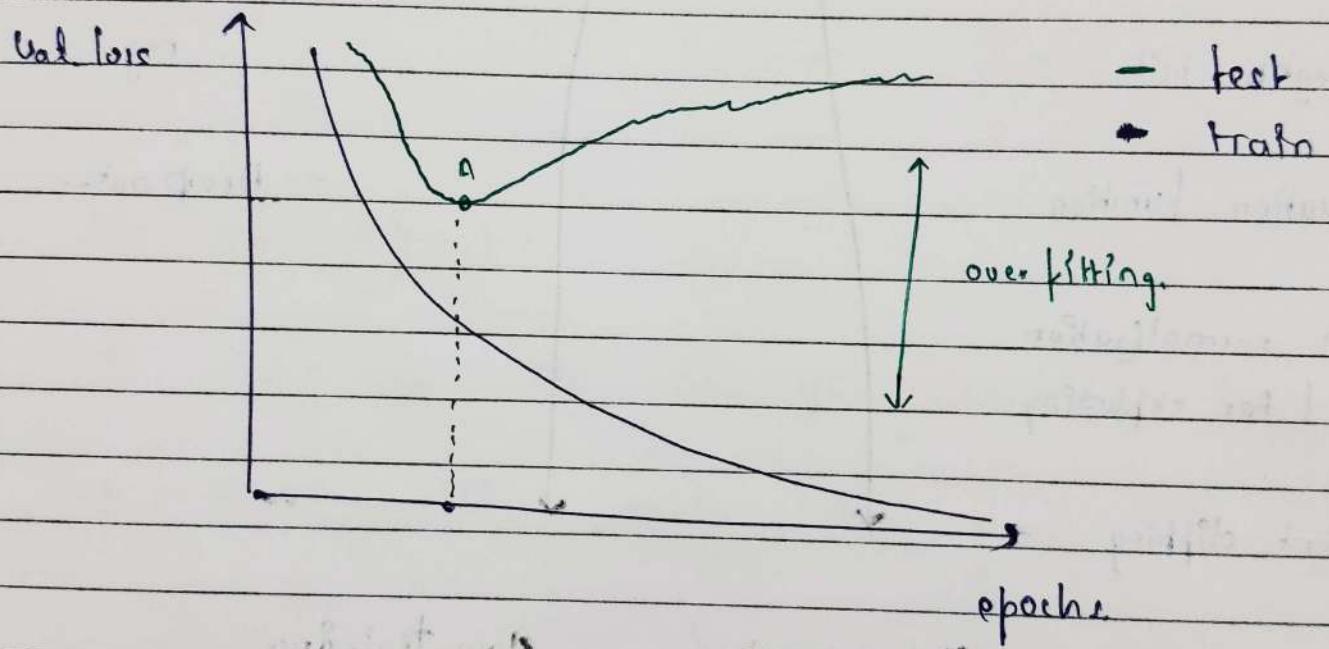
- Transfer learning
- Optimizers

- unsupervised
pretraining

- learning
rate
scheduler

* Early Stopping

- By this concept you can stop the training ~~process~~, where overfitting is started.
- By this concept you can first take As big no^o as you want



- after this point overfitting starts increasing.
- so the epoch where "A" is occurring the training stops.

=> Add in code.

Add in code.

- call back = Early stopping (

monitor = "val loss",

min_delta = 0.0001,

patience = 20,

verbose = 1

mode = "auto"

baseline = None

restore_best_weights = False

)

{ Patience → no° of epochs the algo have to wait
after overfitting starts. }

Then add.

- model.fit (x_train, y_train, validation_data = (x_test, y-test))

, epoch = 3500, call back = callback)

Improving Neural Network Performance.

1. Overfitting

- Dropout layers
- Regularization L_1 & L_2
- Early stopping ✓

2. Normalization

- Normalizing inputs ✓
- Batch Normalizing
- Normalizing Activation

3. Vanishing gradients.

- Activation function
- weight initialization

4. Gradient checking and clipping

5. Optimizers.

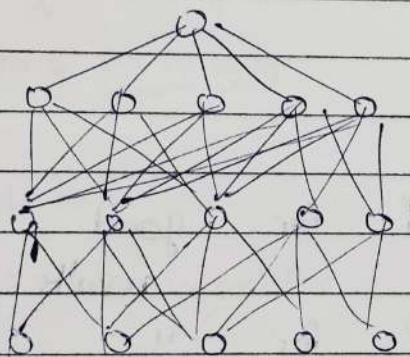
- Dropout layer :-

- N.N is prone to overfit
- To reduce that we have few techniques.
 - Add more data
 - Reduce complexity
 - Early stopping
 - Regularization $\rightarrow L_1, L_2$
 - Dropout layers.

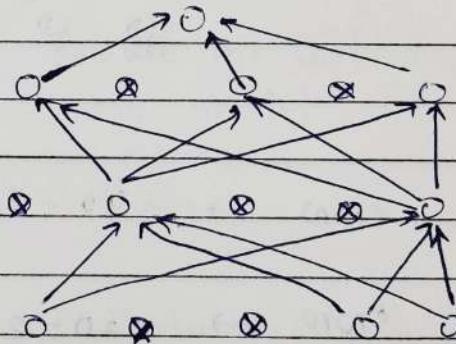
- Dropout layers \Rightarrow

$P = \text{dropout ratio}$

Tells
How many nodes will be turned off
 $P = 0.2 \Rightarrow 20\%$



(N.N)



some of nodes are dropped

- in every epoch we randomly remove some nodes in the layers.
- that means in every epoch we are changing architecture of neural network which helps in reducing overfitting

$$\left. \begin{array}{lll} P = 0.2 & P = 0.5 & P = 0.75 \end{array} \right\}$$

overfitting

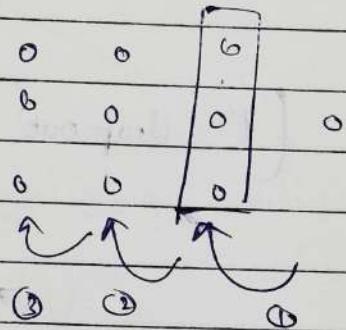
<--

underfitting

Practical Tips.

1. If overfitting then $\uparrow P$, if underfitting then $\downarrow P$

2. Try adding dropout to the last layer first
then go feather to initial layers.



3. CNN \rightarrow 40-80% $\Rightarrow P$ for good results

RNN \rightarrow 20-30% $\Rightarrow P$ for "

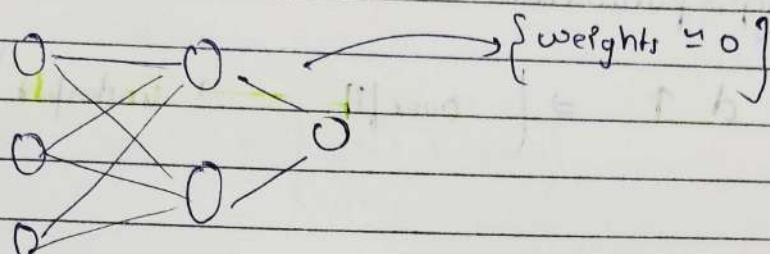
Drawback.

1. Convergence \rightarrow delay \rightarrow Training gets slow

2. loss function value changes

• Regularization (L_1 & L_2)

- When to use. → when model is overfitting



(ways to solve overfitting)

Adding more data

Reduce the complexity
of model

- Add rows

- Dropout

- Data augmentation

Early stopping.

→ In ANN → $(w | b)$

↳ loss (min)

→ Regularization $\xrightarrow{L_1}$
 $\xrightarrow{L_2}$
 $\xrightarrow{\text{L1+L2}}$

$$C = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i) + (\text{Penalty term})$$

$$C = L + \left[\frac{\lambda}{2n} \sum_{i=1}^k \|w_i\|^2 \right]$$

Cost Loss Penalty Term

for $(w_1 \rightarrow w_{10})$

(λ = Regularization factor)

PAGE NO.:

DATE: / /

$$C = L + \frac{\lambda}{2n} [w_1^2 + w_2^2 + \dots + w_{10}^2]$$

λ = Hyperparameter

\rightarrow As $\lambda \cdot \lambda \uparrow \Rightarrow$ (Overfit \rightarrow Underfit)

$$\Rightarrow C = L + \lambda \sum_{i=1}^k \|w_i\|^2 \longleftrightarrow L_2 \text{ (mostly used)}$$

$$C = L + \frac{\lambda}{2n} \sum_{i=1}^n \|w_i\| \rightarrow L_1$$

\rightarrow most accurate representation :

$$C = \sum_{i=1}^m L(y_i, \hat{y}_i) + \boxed{\sum_{l=1}^L \sum_{i=1}^n \sum_{j=1}^q \|w_{ij}\|^2}$$

\rightarrow mathematical intuition

$$w_j = w_0 - \eta \frac{dL}{dw_0}$$

L = loss function.

Add Penalty to apply regularization

$$L = L + \frac{\lambda}{2n} \sum_{i=1}^n \|w_i\|^2$$

\downarrow
n=1 for single node

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial w_0} + \frac{1}{\lambda} (\beta w_0)$$

$$= \frac{\partial L}{\partial w_0} + \frac{1}{\lambda} w_0$$

$$w_n = w_0 - \eta \left| \frac{\partial L}{\partial w_0} + \frac{1}{\lambda} w_0 \right|$$

$$[w_n = \underbrace{(1 - \eta \frac{1}{\lambda}) w_0}_{\text{weight}} - \eta \frac{\partial L}{\partial w_0}]$$

↳ (weight decay)

$L_2 \rightarrow$ weight decay.

Activation function (Brief)

$$a = g(w_1x_1 + w_2x_2 + \dots + b)$$

g = activation function → Sigmoid, relu, tanh.

- If $g(z) = z$ is linear the neuron perform linear regression or classification.

- In $g(z) =$ taken non-linear to make non-linear regression or classification.

→ Ideal activation func :-

1} Should be non-linear eg: $\sigma(z) = \frac{1}{1+e^{-z}}$

2} It should be differentiable.

exception \Rightarrow relu not differentiable at "0"

3} Computationally inexpensive.

↳ derivative \rightarrow simple
 \rightarrow easy fast.

4} Zero centered \rightarrow zero centred.

↳ normalized.

eg: \tanh

{ Proven that model converges fast when normalized data is given in input

→ If (values) of activation function \Rightarrow $k_1 \{ \text{mean} = 0 \}$
then data is normalized.

PAGE NO.:
DATE: / /

* Non-saturating :- (considered) * (good)

→ Saturating function = which bound their output values in a range.

Saturating

sigmoid $\rightarrow (0, 1)$

tanh $\rightarrow (-1, 1)$

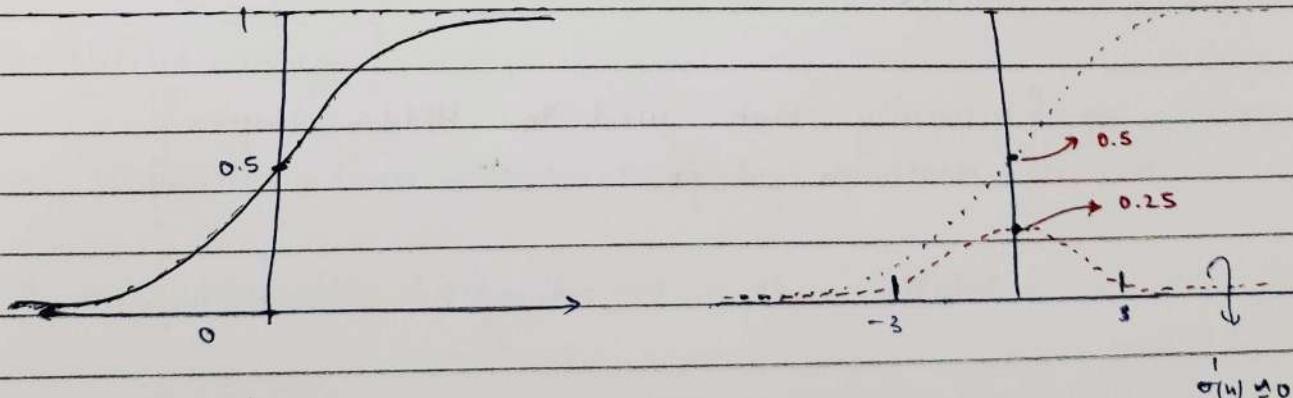
non-saturating

relu = $\max(0, x)$

Problem \Rightarrow vanishing gradient
descents.

1

Sigmoid Act. Function :-



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

$$\sigma(u) \approx 0$$

→ Advantages

(1) $[0, 1] \rightarrow$ probability \rightarrow output layer \rightarrow Binary-class

(2) Non-linear \rightarrow (non-linear data) \rightarrow good option

(3) Differentiable \rightarrow Back prop. $\rightarrow \frac{\partial L}{\partial w}$

→ Disadvantages.

1) Saturating func' $\rightarrow (-\infty, \infty) \rightarrow [0,1]$



Input Output,
 $(-\infty, \infty) \rightarrow [0,1]$
Vanishing gradient problem

2) Non zero centered

0 → output not either $\rightarrow +ve$
0 → normalized. $\rightarrow -ve$

0

0

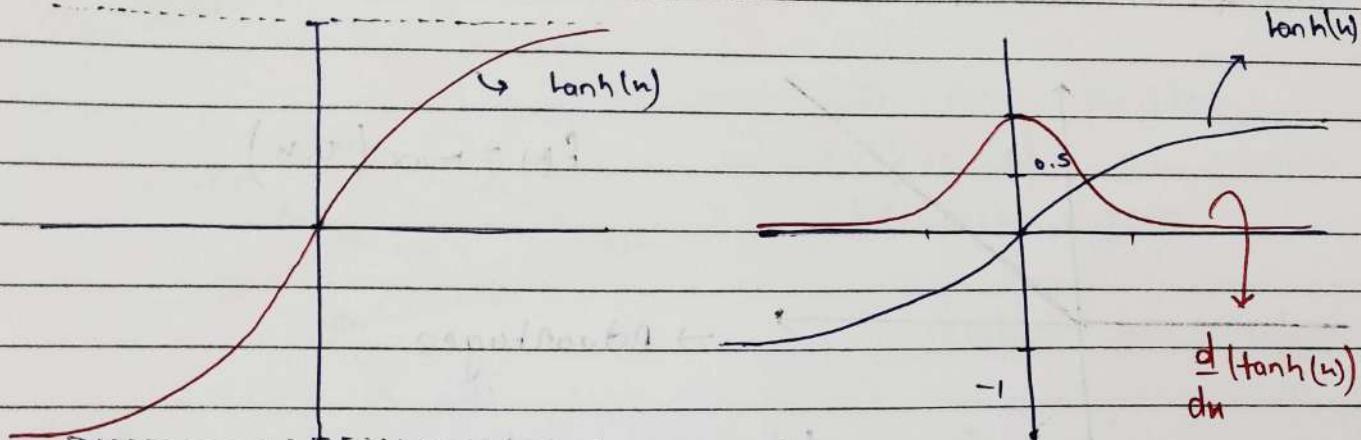
which leads to (Training slow)

{ Convergence problem }

3) $\sigma(u) = \frac{1}{1+e^{-u}} \rightarrow$ (computationally expensive)

{ generally not used in Hidden layers }
now a days.

(2) Tanh Activation function.



$$f(u) = \tanh(u) = \left(\frac{e^u - e^{-u}}{e^u + e^{-u}} \right) \quad f'(u) = \left(1 - \tanh^2(u) \right)$$

→ **Advantage.**

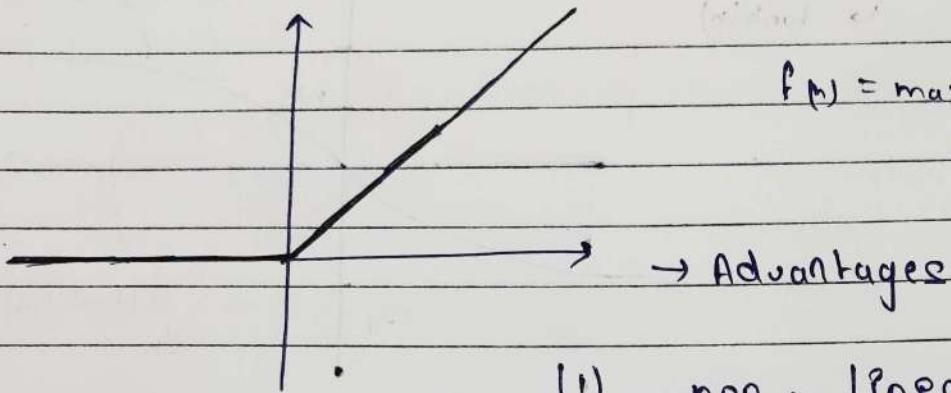
- (1) non linear
- (2) differentiable.
- (3) zero centred. → training faster than sigmoid.

→ **D3. Advantages.**

→ saturating function → vanishing gradient problem.

→ computationally expensive → slow training

Rely Activation funi.



$$f(u) = \max(0, u)$$

→ Advantages

(1) non-linear

(2) Not saturated in +ve region

(3) computationally in-expensive

(4) converges → faster → $\begin{cases} \text{sigmoid} \\ \tanh \end{cases}$

→ Disadvantages

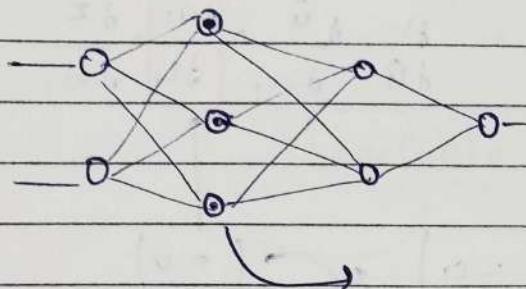
(1) non differentiable at "0" \Rightarrow $\begin{cases} \text{But} \\ \text{in code} \\ \text{we take} \end{cases} \Rightarrow f'(u) = \begin{cases} 0, & u < 0 \\ 1, & u > 0 \end{cases}$

(2) Non zero centered →

But this problem is solved

By Batch normalization

Dying ReLU Problem



→ In this problem.

output of a neuron becomes zero

irrespective of input

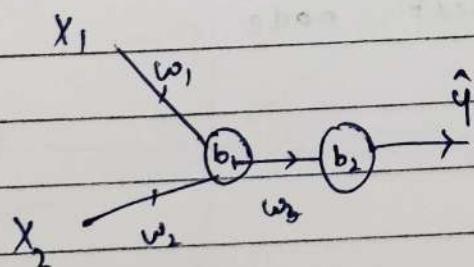
→ neuron's output become zero forever

→ If IT becomes case with

$$\left\{ \begin{array}{l} \text{no. of neuron} \\ \text{output = 0} \end{array} \right\} > 50\%$$

Then patterns will not be recognized.

→ Let's see why ?

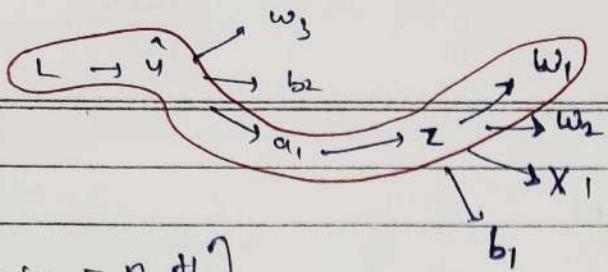


$$q_1 = \max(0, z_1) =$$

$$\text{if, } z_1 = w_1x_1 + w_2x_2 + b_1 < 0$$

$$\text{then } (a_1 = 0) = \frac{da_1}{dz_1} = 0$$

now



$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} + \boxed{\frac{\partial a_1}{\partial z}} \cdot \frac{\partial z}{\partial w_1}$$

$$\text{so } \left(\frac{\partial L}{\partial w_1} = 0 \right) \curvearrowleft \left(= 0 \right)$$

$$\text{similarly } \left(\frac{\partial L}{\partial w_2} = 0 \right)$$

→ So weights will not get updated.

→ reasons why

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

(1) learning rate is very high (η)

then $0 > z > 0$ in next node

can be

(2) High (-)ve bias.

Solution of Dying Relu.

PAGE NO.:

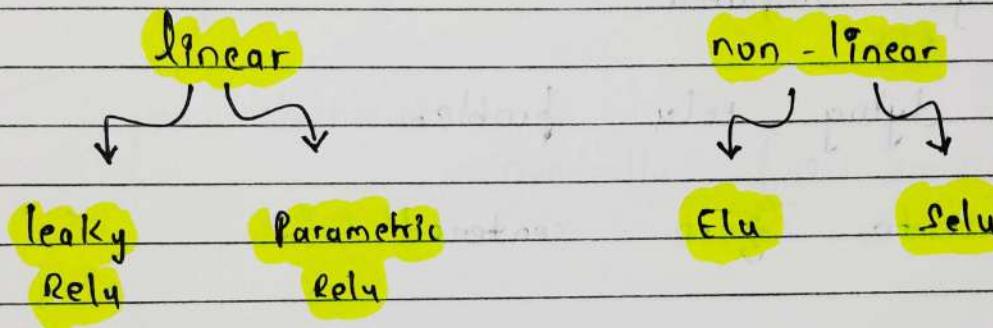
DATE: / /

→ Set low learning rate

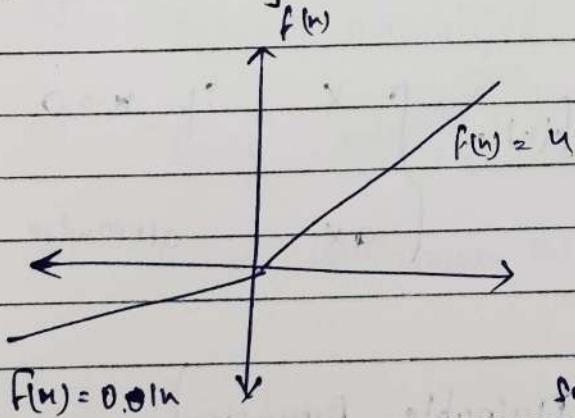
→ bias → +ve value → (0.0)

→ Don't use relu → use variants

Variants of Relu



① Leaky Relu



$$f(z) = \max(0.01z, z)$$

$$f(z) = \begin{cases} z & z \geq 0 \\ \frac{z}{100} & z < 0 \end{cases}$$

$$f'(z) = \begin{cases} 1 & z > 0 \\ 0.01 & z \leq 0 \end{cases}$$

does not become zero.

$$\therefore w_n = w_0 - \eta \left(\frac{d L}{d w_n} \right)$$

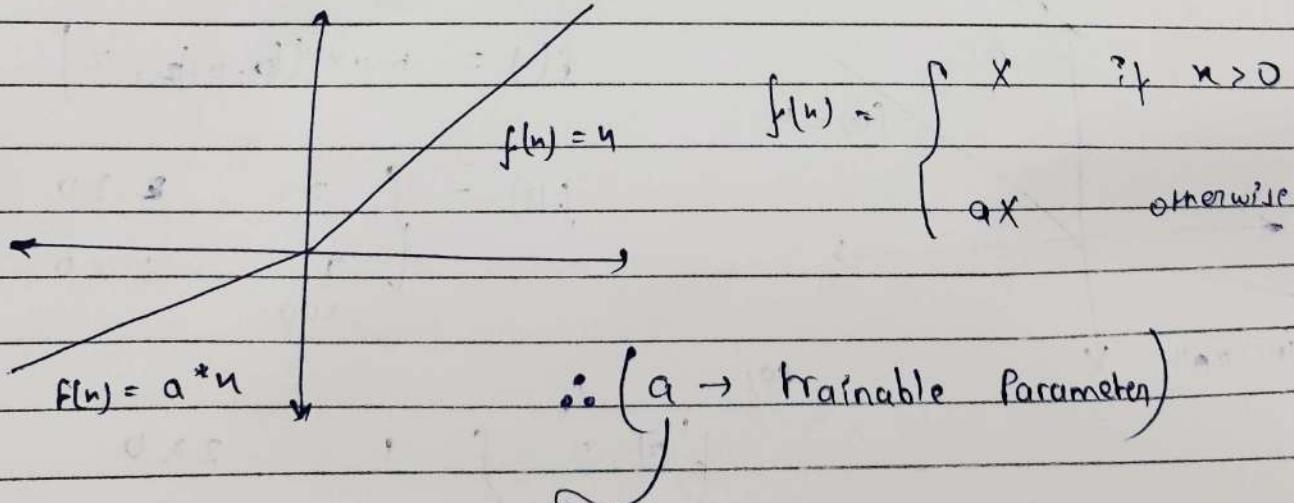
PAGE NO.: _____
DATE: / /

will not get zero
will always get small changes.
, so Dying relu will not occur.

• Advantages

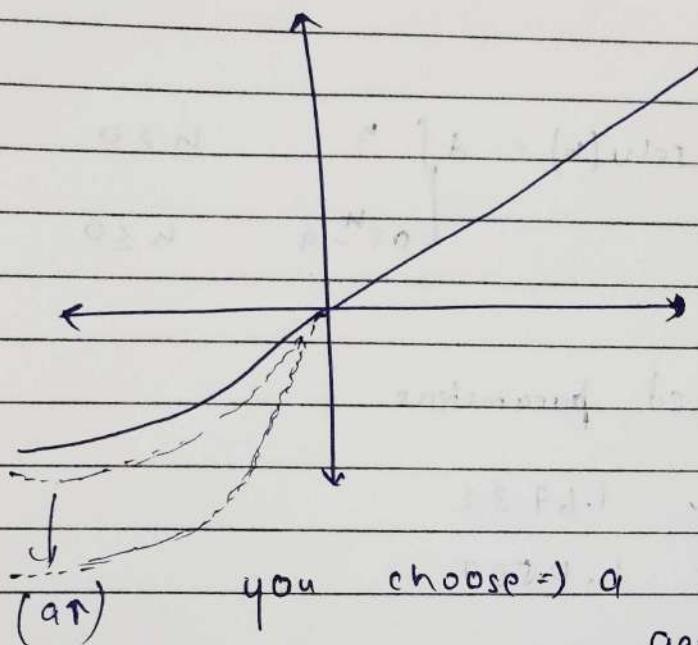
- non saturated. → unbounded
- easily computed
- No dying relu problem
- closer to zero centered.

(2) Parametric Relu



Best value gets while training

Elu - Exponential Linear Unit.



$$\text{ELU}(u) = \begin{cases} u & u > 0 \\ a(e^u - 1) & u \leq 0 \end{cases}$$

$$\text{ELU}(u) = \begin{cases} u & u > 0 \\ \text{ELU}(u) + a & u \leq 0 \end{cases}$$

you choose $\Rightarrow a$
generally (0.1 to 0.3)

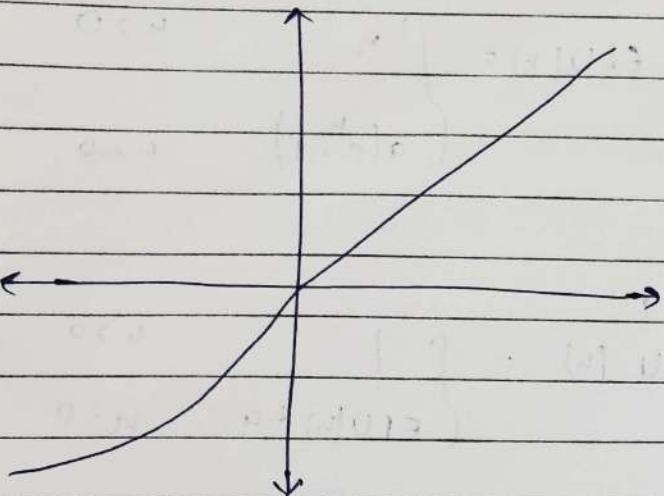
Advantages

- (1) close to zero, centered \rightarrow converges faster.
- (2) Better generalized
- (3) Dying relu (X)
- (4) always continuous as well as differentiable.

disadvantages

- (1) computationally expensive.

(4) Selu - Scaled Exponential Linear Unit



$$\text{selu}(u) = \begin{cases} \alpha \int u & u > 0 \\ \alpha e^u - \alpha & u \leq 0 \end{cases}$$

Fixed parameters

$$\alpha \approx 1.6732$$

$$\beta \approx 1.0507$$

$$\text{selu}(u) = \begin{cases} \beta u & u > 0 \\ \alpha e^u - \alpha & u \leq 0 \end{cases}$$

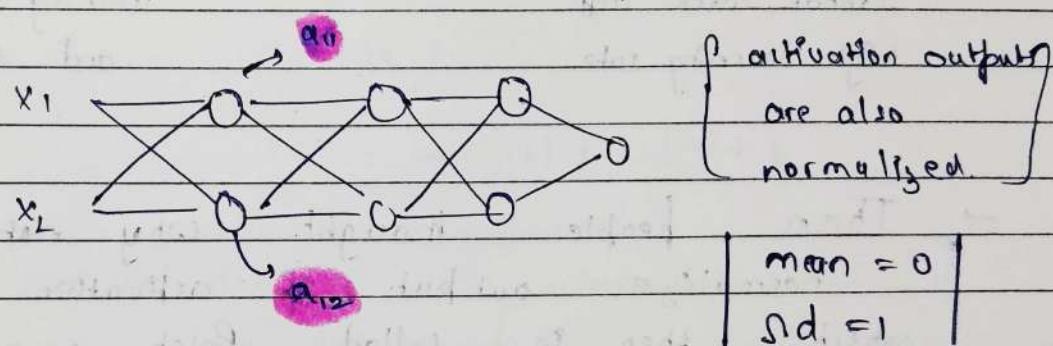
Advantage.

- self normalizing \rightarrow activation output \rightarrow normalized
 \downarrow
mean = 0
- converges fast

- Batch normalization.

BN makes training of (DNN) faster and stable

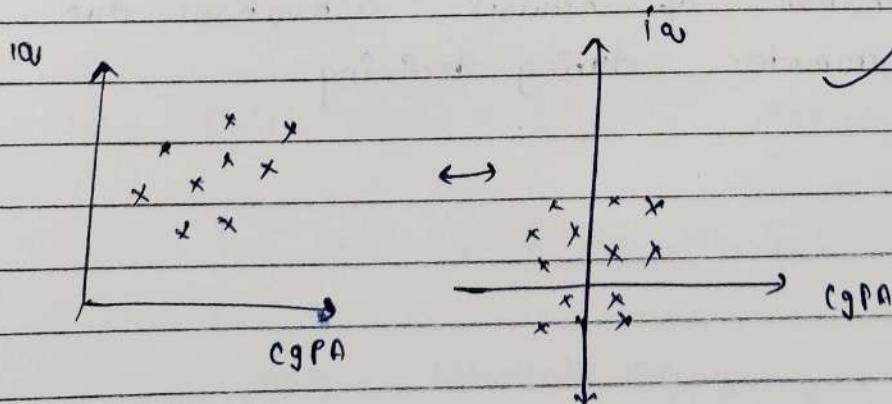
- It consists of normalizing activation vectors from hidden layers using the mean and variance of the current
- This activation normalization is done right before (or right after) the non linear function



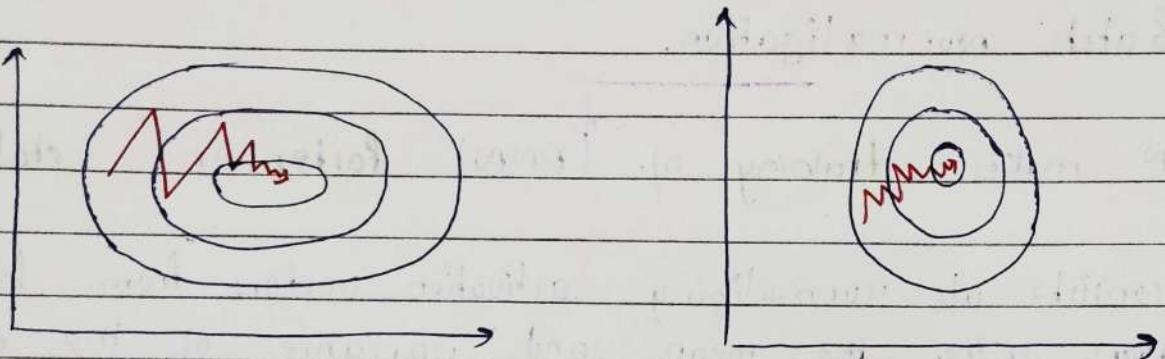
- Why Norm Batch?

cgPa	\bar{P}_a	Plotted.
7	70	1
8	80	0
9	90	1

$$\begin{cases} \mu = 0 \\ \sigma = 1 \end{cases}$$



(Before Normalization.) (after normalization)



(Before normalization)

cannot work with
high learning rate

(with normalization)

training is faster
and stable.

→ Then people thought why not to normalize output of activation units which then is called Batch normalization.

another reason

- Internal covariate shift

→ we define internal covariate shift as the change in the distribution of network activations due to network parameters during training

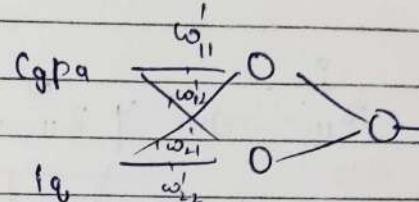
• The How?

- Mini-Batch Gradient Descent.
- layer by layer update.

Cgpa | iω | placed

8.9	100	,
6.2	89	0
7.7	91	0
9.1	76	1
1	1	1
1	1	1
1	1	1

normalize $\begin{cases} \mu=0 \\ \sigma=1 \end{cases}$



$$z_{11} = w_1(\text{Cgpa}) + w_2(i\omega) + b$$

→ step by step calculation.

(1) $z_{11} \rightarrow (z_{11})$ ^{normalized} $\rightarrow g(z_{11}) \rightarrow a_{11}$ ✓

(2) $z_{11} \rightarrow g(z_{11}) = a_{11} \rightarrow (a_{11})$ ^{normalized}

$$\therefore \left(\frac{z_{11} - \mu}{\sigma} \right)^N$$

→ take batch size = (n) assume.

- for each node we will get a (z_{11}) value corresponding to each row.
- Then calculate $(m=4)$ Batch size

$$\mu_B = \frac{1}{m} \sum_{i=1}^m z_{11}^i$$

$$\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_{11}^i - \mu_B)^2}$$

Then

calculate

$$\hat{z}_{11} = \frac{z_{11} - \mu_B}{\sigma_B + \delta}$$

(normalized)

error term

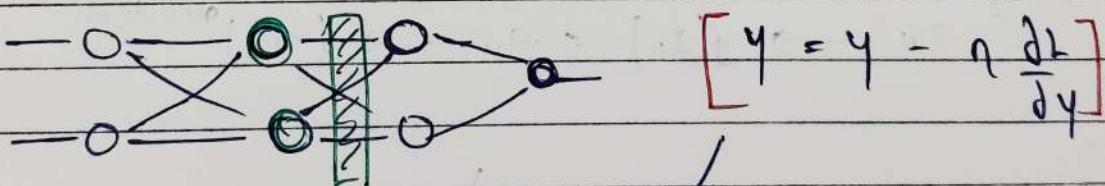
• Step (2)

$$\hat{z}_{11} = \gamma z_{11}^B + \beta$$

learnable parameters

$z_{11} \rightarrow z_{11}^B \rightarrow z_{11} \rightarrow g(z_{11}) = o_{11} \quad (\gamma=1) \quad (\beta=0)$

→ Now in ANN During training.



gamma will also be updated
during training by
G.O

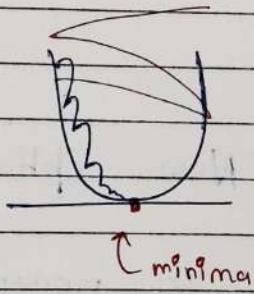
Optimizers

- Batch G.O
- stochastic G.O
- mini batch G.O

→ challenges with G.O

(1) learning rate →

$$w_n = w_0 - \eta \frac{dL}{dw}$$

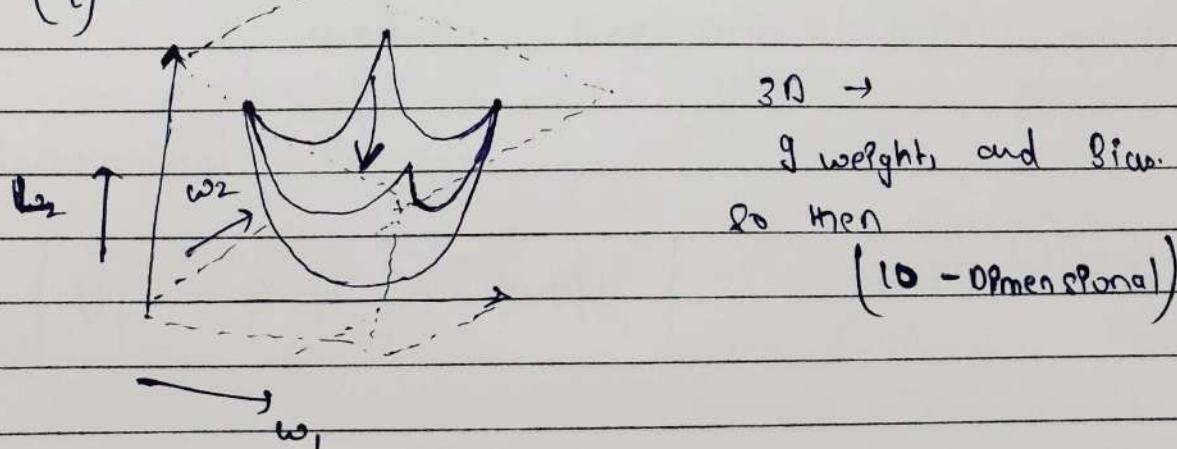


{ It is crucial to decide right }
learning rate

— Solution.

• learning rate scheduling or → problem (pre-defined)

(2) (η) is same for both directions or

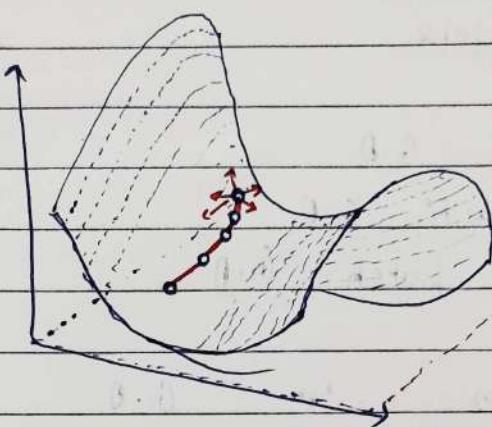


(3) local minima [Problem]

5) Saddle point

$$\frac{\partial L}{\partial w} = 0$$

$$w_p = w_0 - \eta \frac{\partial L}{\partial w}$$



- New optimizers.

1) momentum

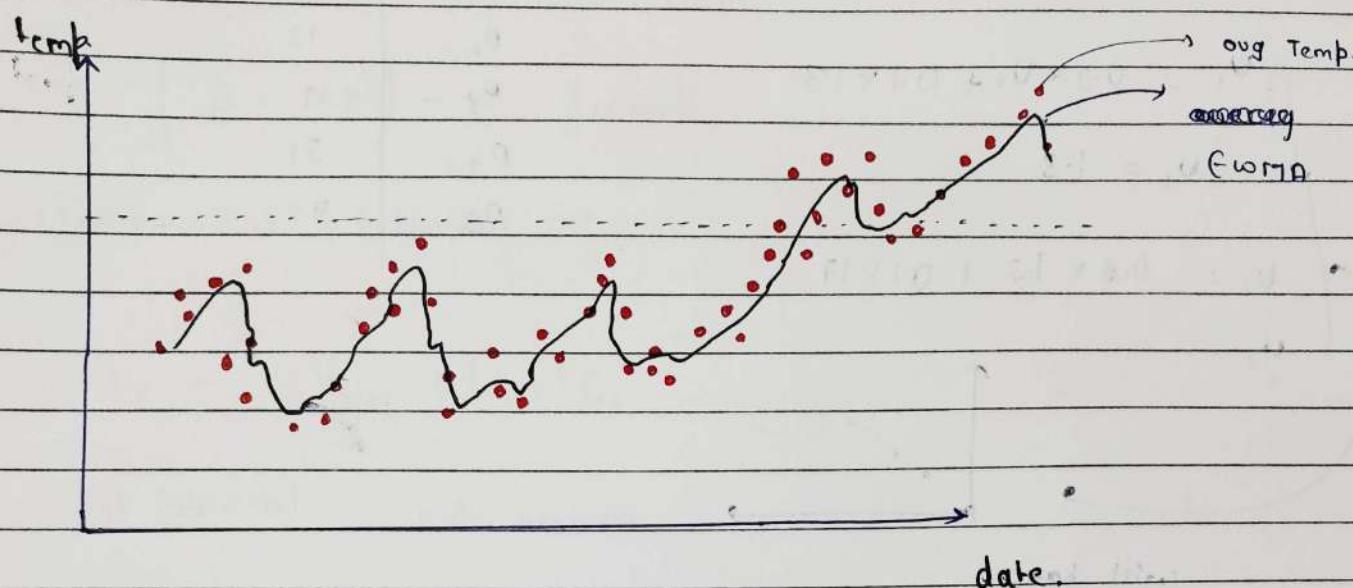
2) Ada grad

3) NAG

4) RMS Prop

5) Adam (Popular)

- Exponentially weighted moving average :- (EWMA)



→ used where ?

- Time series.
- Financial forecasting
- Signal processing
- Deep learning.

→ EWMA is used in time series pattern recognition.

- Mathematical Average.

$$U_t = \beta U_{t-1} + (1-\beta) \Theta_t$$

$\beta \rightarrow$ constant.

$$(U_0 = \Theta_0) \quad \text{generally } (\beta = 0.9)$$

index	temp(Θ)
O_1	25
O_2	13

O_3	17
-------	----

O_4	31
O_5	43

$$\beta = 0.9$$

temp. (θ)

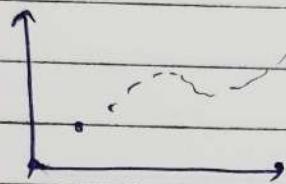
D_1	25
D_2	13
D_3	17
D_4	31
D_5	43

$$U_1 = 0.9 \times 0.9 + 0.1 \times 13$$

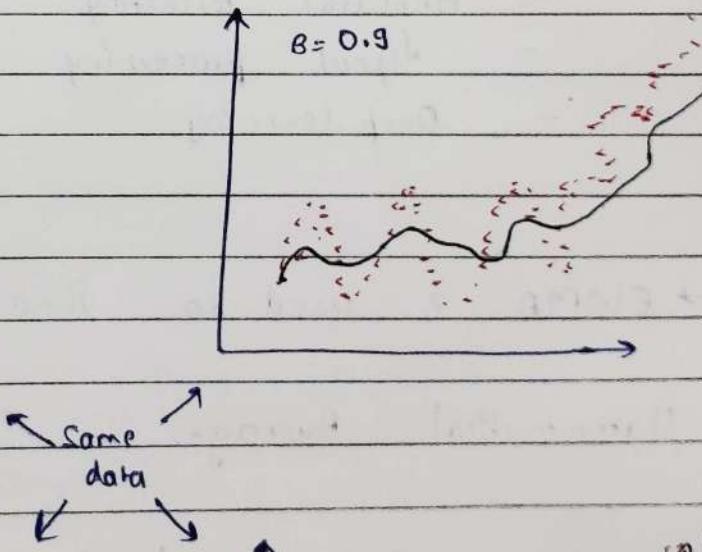
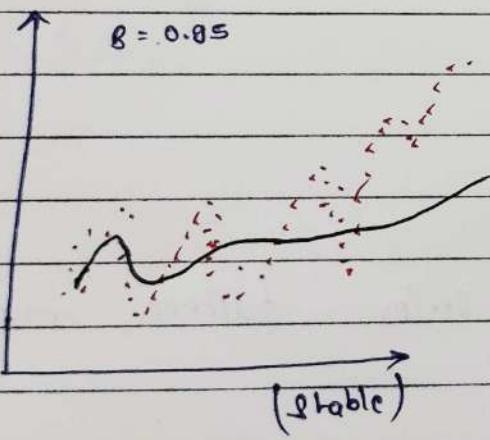
$$U_1 = 1.3$$

$$U_2 = 0.9 \times 1.3 + 0.1 \times 17$$

U_1

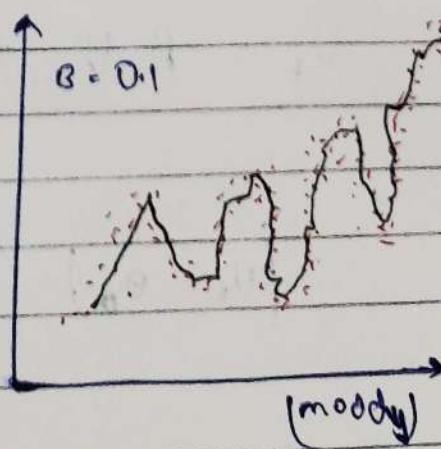
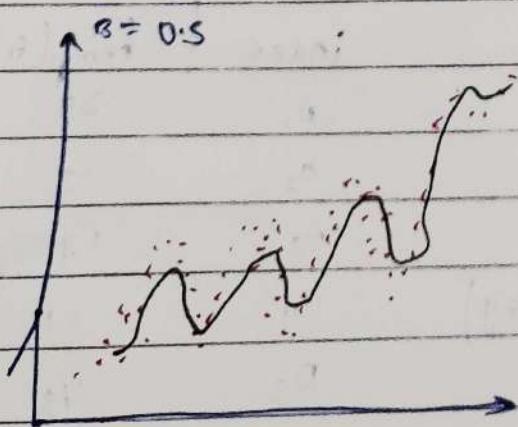


will be
formed by these values.



(stable)

Same
data



(mood)

→ As (B) is increased value of (EWMA) depends on more previous data.

$$(B = 0.9) \rightarrow \text{generally.}$$

• Mathematical intuition →

$$U_t = B U_{t-1} + (1-B) O_t$$

let ($U_0 = 0$) not necessary

$$U_1 = (1-B) O_1$$

$$\begin{aligned} U_2 &= BU_1 + (1-B) O_2 \\ &= B(1-B) O_1 + (1-B) O_2 \end{aligned}$$

$$U_3 = BU_2 + (1-B) O_3$$

$$= B^2(1-B) O_1 + B(1-B) O_2 + (1-B) O_3$$

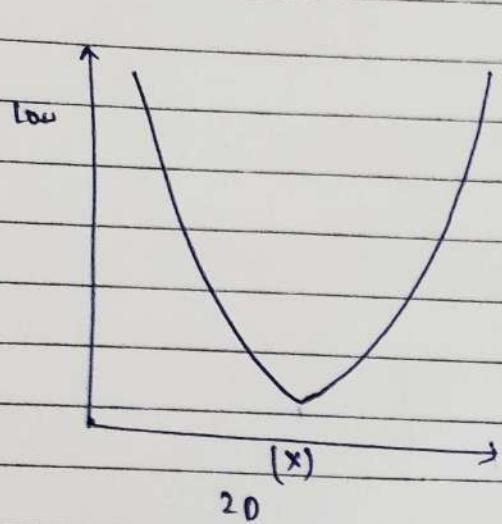
$$U_4 = BU_3 + (1-B) O_4$$

$$U_4 = (1-B) \left(\underbrace{B^3 O_1}_\uparrow + \underbrace{B^2 O_2}_\uparrow + \underbrace{BO_3}_\uparrow + O_4 \right)_\uparrow$$

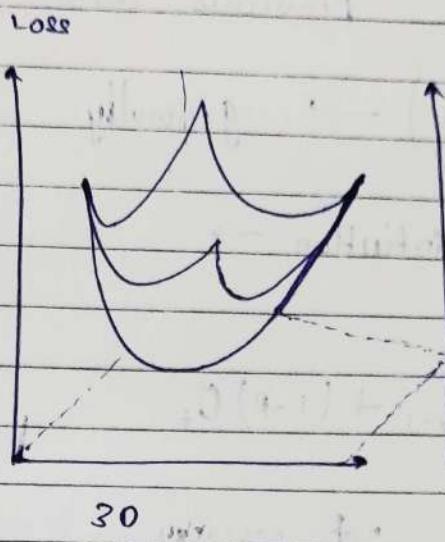
$$(B^3 < B^2 < B)$$

• smaller values are multiplied by older values.

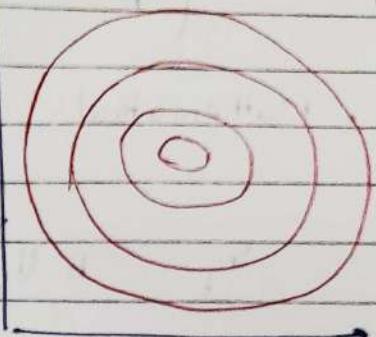
- SGD with momentum.



[momentum]

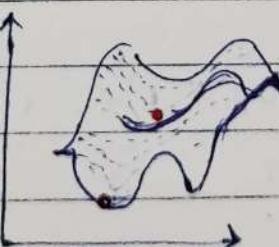
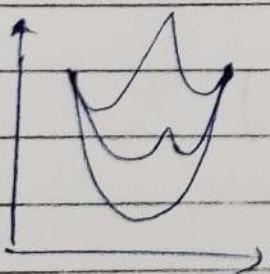
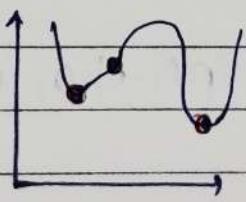
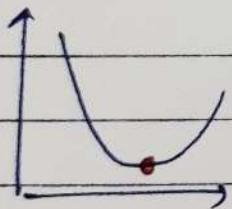


[Visualization
of Loss]



Contour

- convex (Vs) non convex

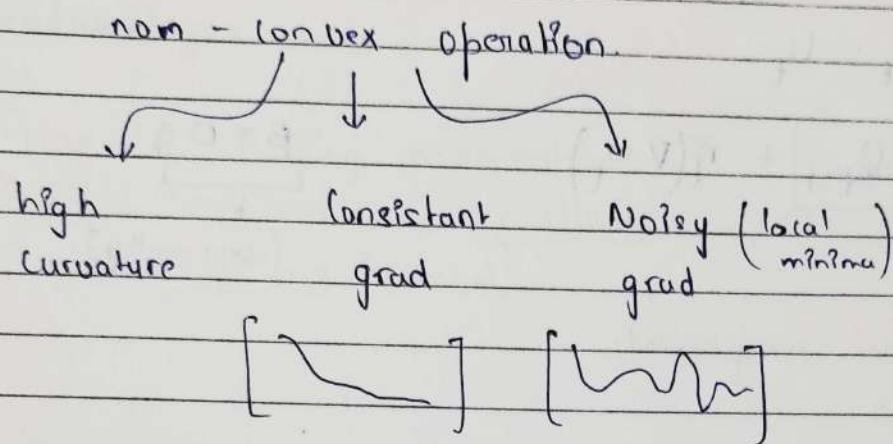


- { 1) Local minima
- 2) Saddle point
- 3) High curvature.

SGD cannot navigate
all these three.

Momentum optimization - why?

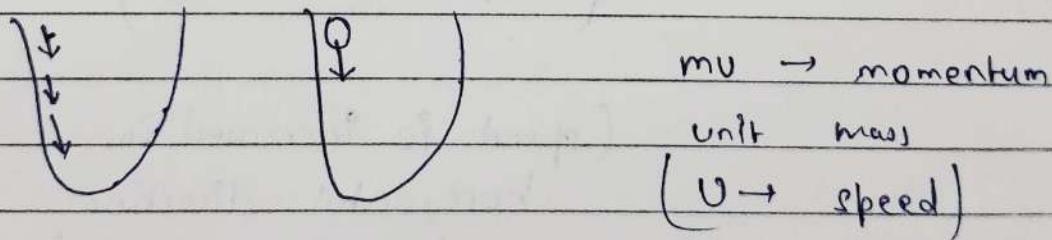
- momentum is good at navigating.



→ momentum helps in solving all these three issues.

→ (Intuition)

If previous gradients are saying you to go in a direction then you will increase speed in that direction



→ (Mathematics.)

$$w_{t+1} = w_t - \eta [\Delta w_t]$$

Sume.

$$w_{\text{new}} = w_0 - \eta \frac{\partial L}{\partial w}$$

$$w_{t+1} = w_t - \eta \nabla w_t \quad \left. \right\} \text{in momentum}$$

$$w_{t+1} = w_t - v_t$$

$$w_{t+1} = w_t - v_t$$

$$\therefore v_t = \boxed{\beta \cdot v_{t-1}} + \eta (\nabla w_t)$$

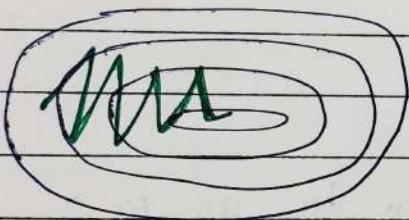
\downarrow

$$\underbrace{\beta}_{t} = 0.9$$

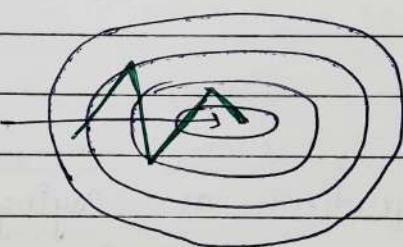
{ Generally. }

momentum component.

(difference.)



(SGP)



(SGP with momentum)

(speed is increased in
horizontal direction
for faster convergence)

• Effect of (β)

$$w_{t+1} = w_t - v_t$$

$$v_t = \beta \cdot v_{t-1} + \eta \nabla w_t$$

- Advantage of momentum.

- SGD with momentum converges faster than normal SGD
- has the ability to come out of local minima

- disadvantages

- due to having momentum if surpasses passes the Global minima, oscillates and then converges to minima (global)

- NAG - Nesterov Accelerated Gradient

- some disadvantages in momentum were tried to solve in NAG. (momentum problem)

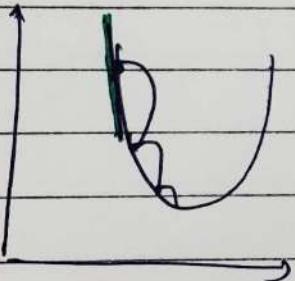
→ (Momentum.)

$$\omega_{t+1} = \omega_t - u_t$$

where

$$[u_t = \beta \cdot u_{t-1} + \eta \nabla \omega_t]$$

$$\omega_{t+1} = \omega_t - (\underbrace{\beta u_{t-1} + \eta \nabla \omega_t}_{\downarrow})$$



Jumps will be decided by these two components.

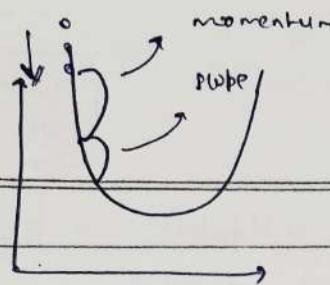
- (1) fast velocity
- (2) Gradient at that pt

→ (NAG)

In this we first calculate momentum then take first jump

then at that point we calculate slope then take second jump \perp to the slope.

NAG



PAGE NO.:

DATE: / /

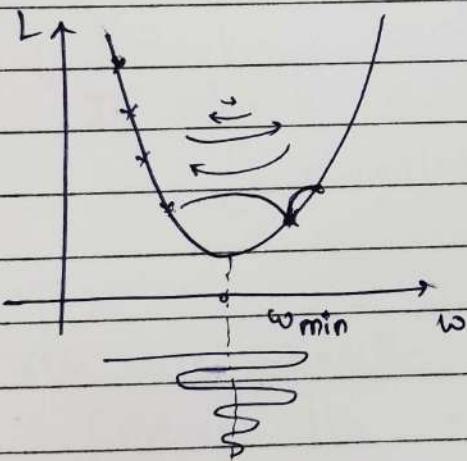
$$\omega_{ea} = \omega_t - \beta v_{t-1} \rightarrow \text{first update due to previous velocity}$$

$$v_t = \beta v_{t-1} + \eta \nabla \omega_{ea}$$

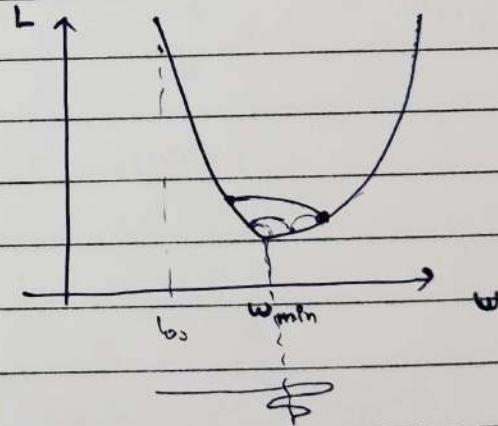
$$\Rightarrow v_t = (\omega_t - \omega_{ea}) + \eta \nabla \omega_{ea}$$

$$\omega_{t+1} = \omega_t - v_b$$

→ Geometric Intuition :-



(momentum)



(NAG)

- disadvantage.

- can stuck in local minima

Keras code

```
tf.keras.optimizers.SGD(
```

```
    learning_rate = 0.01, momentum = 0.0, nesterov = False  
, name = "SGD", **kwargs)
```

SGD	{	momentum	}	NAG
	{	momentum = 0.9	{	momentum = 0.9
		nesterov = False		nesterov = True.

AdaGrad - Adaptive Gradient.

- learning rate is not fixed
∴ (adapted as per situation)
- where AdaGrad performs good.

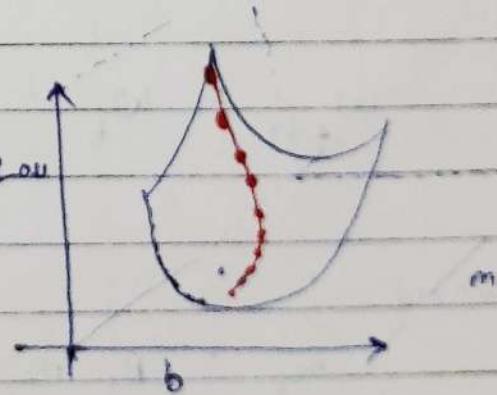
(1) features are (sparse) → mostly ('0)
then [AdaGrad] is usefull

eg.	Data set	10	cgpa	from iit	pakago
		,	,	0	
		:	{	:	
		:		0	
				0	
					0
					0
					0
					0
					0

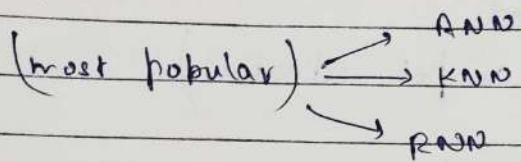
In sparse Dataset

[elongated bowl problem occurs]

→ disadvantage :-
[misses the minima]



Adam → Adaptive moment Estimation



- SGD | BGD | MB GD works on
- momentum ↗
- NAG ↗ (momentum)
- Ada grad. ↗
- RMS prop. ↗ (learning rate decay)

Adam → [works on] (momentum + learning rate decay)

• Mathematical formulation :-

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * m_t$$

Where

$$[m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla w_t] \rightarrow \text{momentum}$$

$$[v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla w_t)^2] \rightarrow \text{Adagrad}$$

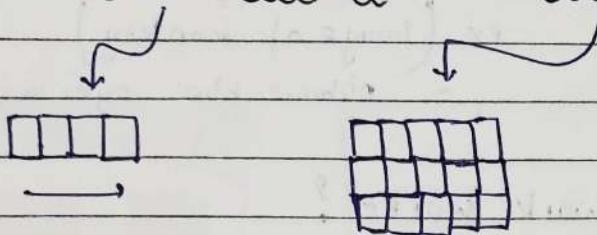
But correction

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t} \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\left. \begin{array}{l} \beta_1 = 0.9 \\ \beta_2 = 0.99 \end{array} \right\} \text{changeable}$$

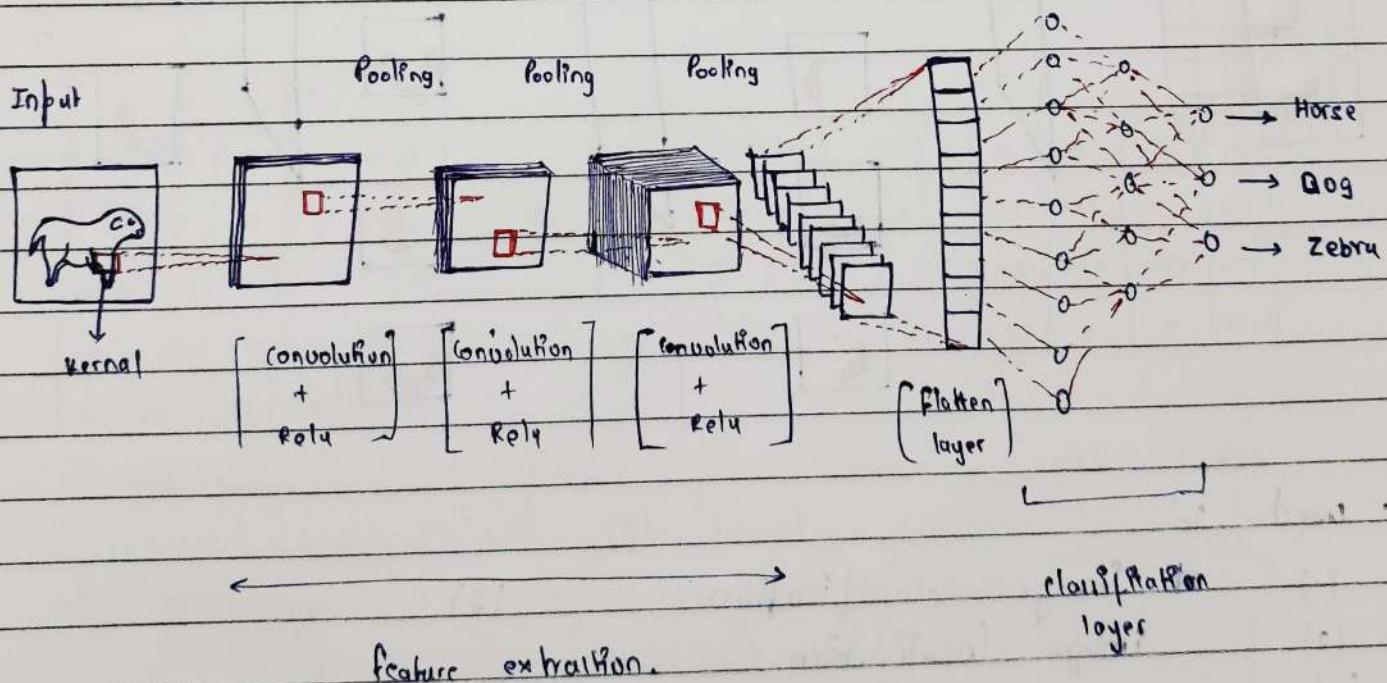
Convolutional neural network.

- CNN are a special kind of neural network for processing data that has a known grid-like topology like time series data (1D) or images (2D).



ANN \rightarrow matrix multiplication

CNN \rightarrow convolutional operation.



- Components
- Convolution
- Pooling
- FC layer - ANN

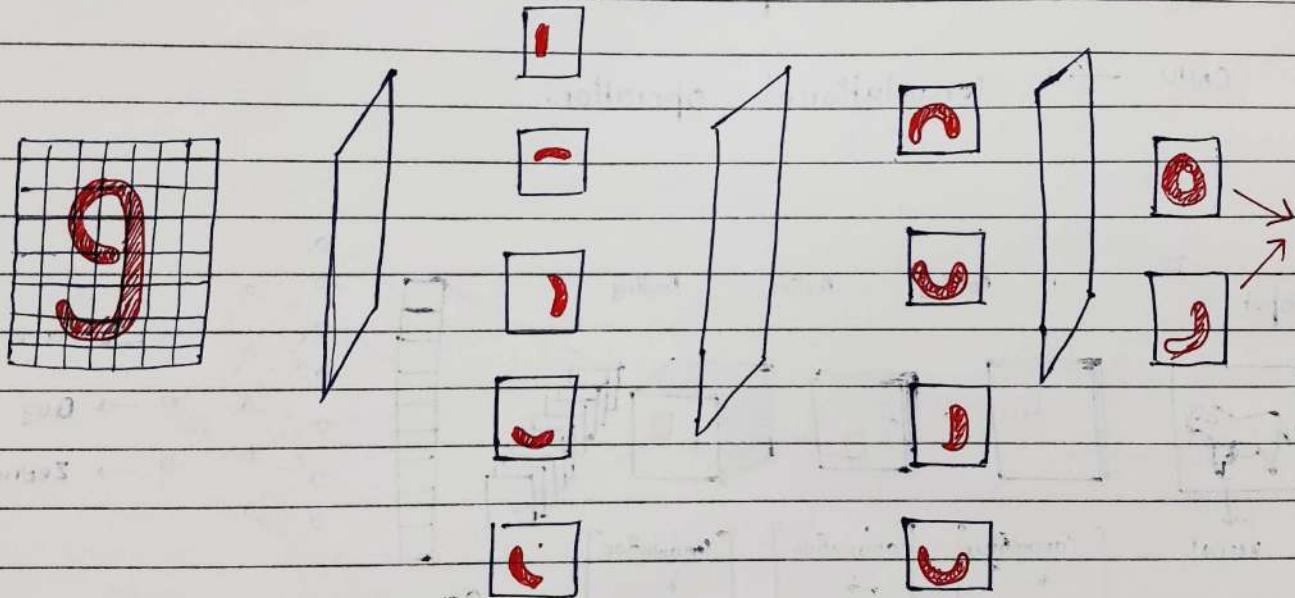
- why not ANN?

1.2 High computational cost

3.2 Overfitting

3.7 loss of imp info like spatial arrangement of pixels
ex (Image of monkey)
- Distance b/w eyes and nose

- How convolution neural network works?



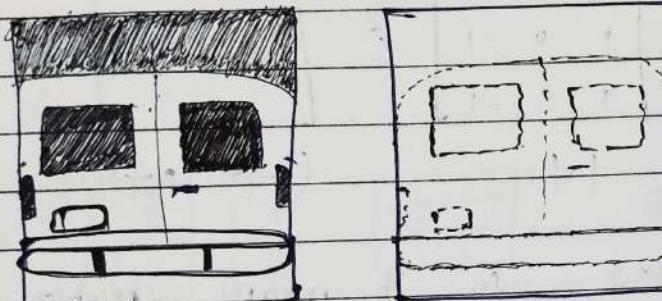
- used in

- (1) Image classification
- (2) Image localization
- (3) Object detection
- (4) Face detection and recognition
- (5) Image segmentation
- (6) Image resolution
- (7) Image colourization.

- Convolutional operation

Two Types of Images \rightarrow Black & white
 \rightarrow RGB

- Edge detection (Convolution operation)



Vertical and Horizontal
edge detection

- How algorithm Detects edges.

$$\begin{array}{|c|c|c|c|c|c|} \hline
 0 & 0 & 0 & 0 & 0 & 0 \\ \hline
 0 & 0 & 0 & 0 & 0 & 0 \\ \hline
 0 & 0 & 0 & 0 & 0 & 0 \\ \hline
 255 & 255 & 255 & 255 & 255 & 255 \\ \hline
 255 & 255 & 255 & 255 & 255 & 255 \\ \hline
 255 & 255 & 255 & 255 & 255 & 255 \\ \hline
 \end{array}
 *
 \begin{array}{|c|c|c|} \hline
 -1 & -1 & -1 \\ \hline
 0 & 0 & 0 \\ \hline
 1 & 1 & 1 \\ \hline
 \end{array}
 =
 \begin{array}{|c|c|c|c|} \hline
 0 & 0 & 0 & 0 \\ \hline
 255 & 255 & 255 & 255 \\ \hline
 255 & 255 & 255 & 255 \\ \hline
 0 & 0 & 0 & 0 \\ \hline
 \end{array}$$

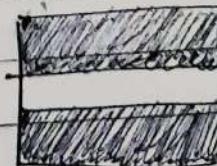
6x6 image.

Filter | Kernel

$$\begin{aligned}
 & -1 \times 0 + -1 \times 0 + -1 \times 0 \\
 & + 0 \times 0 + 0 \times 0 + 0 \times 0 \\
 & + 1 \times 0 + 1 \times 0 + 1 \times 0
 \end{aligned}$$



Horizontal edge
detector



(edge detected)

• Padding

• Problem with convolution,

$$\begin{bmatrix} 7 & 2 & 3 & 3 & 8 \\ 4 & 5 & 3 & 8 & 4 \\ 3 & 3 & 2 & 8 & 4 \\ 2 & 8 & 7 & 2 & 7 \\ 5 & 4 & 4 & 5 & 4 \end{bmatrix}_{n \times n} * \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}_{f \times f} = \begin{bmatrix} \dots \end{bmatrix}_{(n-f+1) \times (n-f+1)}$$

(1) Feature map is of 3×3 , so it means feature map is reduced, we lose some information.

(2) When convolution is taking place, top the inner value of images is having more say and outer corner is having less say.

→ To solve this

We want that feature image should not reduce in size.

$$n \times n * f \times f = (n-f+1) \times (n-f+1)$$

$$(n-f+1) = (n) \rightarrow \text{condition}$$

$$n-3+1 = 5$$

$$(n=7)$$

We need to convert input image into (7×7) image.

0	0	0	0	0	0	0
0	7	2	3	3	8	0
0	4	6	8	7	1	0
0	3	2	1	9	8	0
0	2	4	8	6	3	0
0	5	4	4	2	1	0
0	0	0	0	0	0	0

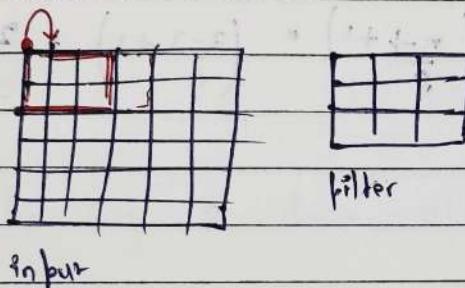
Padding.

PAGE NO.:

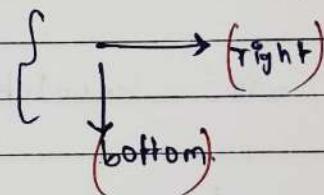
DATE: / /

If all values are "0" in padding then called zero padding.

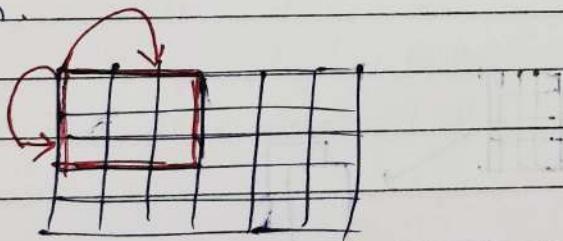
• stride :-



stride = no of block moved by the filter towards right & bottom



→ If stride = (2,2)
then.



Formula.

$$n \times n \times f \times f \rightarrow (n-f+1) \times (n-f+1)$$

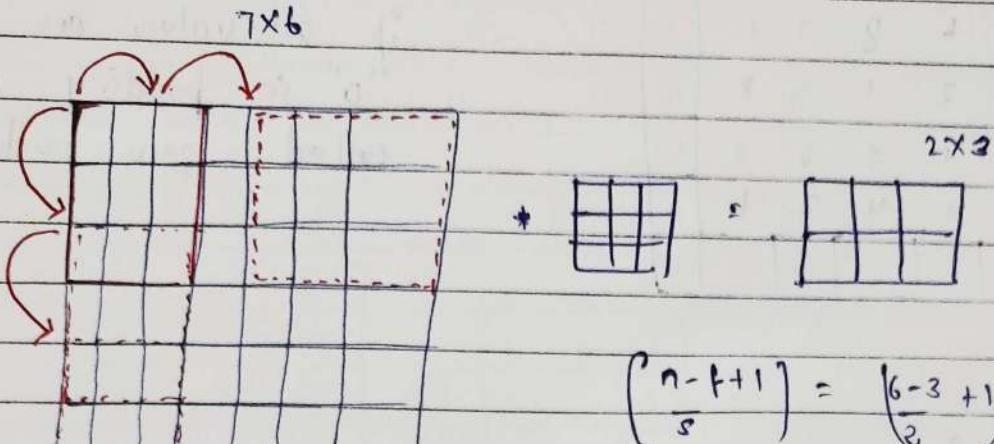
Applied stride

$$(n-f+1) \rightarrow \left[\frac{n-f+1}{s} \right]$$

Applied padding.

$$\text{also } (n-f+1) \rightarrow \left[\frac{n-f+2p+1}{s} \right]$$

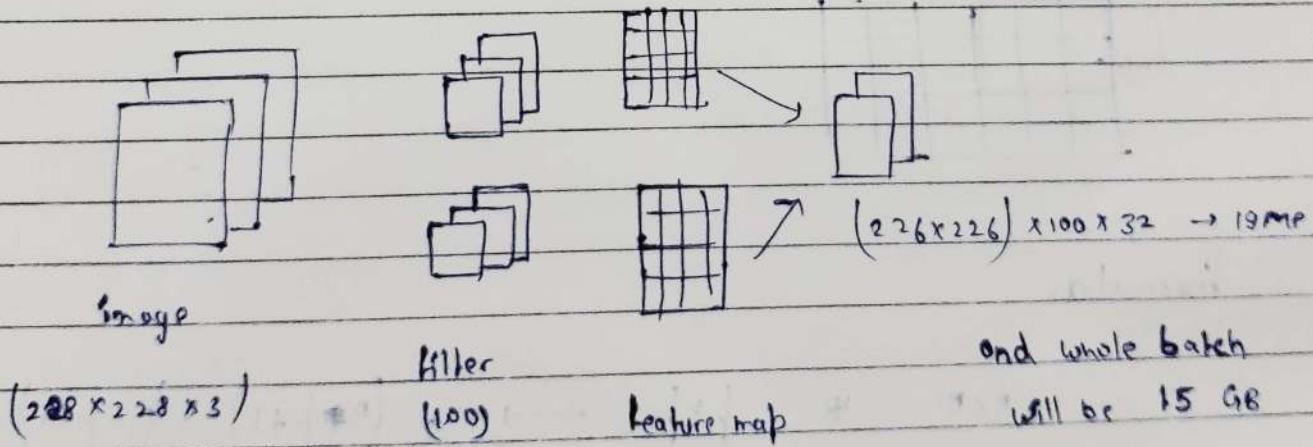
• special case :-



Pooling :-

→ Problem with convolution

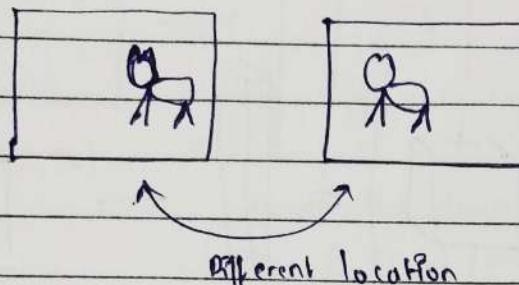
(i) memory issue.



• so we have to reduce the size of feature map

→ Translation variance. - (problem with convolution)

→ Convolution is location dependent.



Pooling

↓
Down scale your
feature map

↓
which makes it independent
to location

→ Pooling is done after convolution.

- ↪ Types
- (1) max pooling
 - (2) min "
 - (3) Avg "
 - (4) L_2 Pooling
 - (5) global Pooling

→ Calculation of image size after pooling.

$$(H \times W) \xrightarrow{\text{pooling}} (k \times k) \longrightarrow \left(\frac{H-k+2P}{s} + 1 \right) \times \left(\frac{W-k+2P}{s} + 1 \right)$$

$k = 2$ generally

$s = \text{stride}$

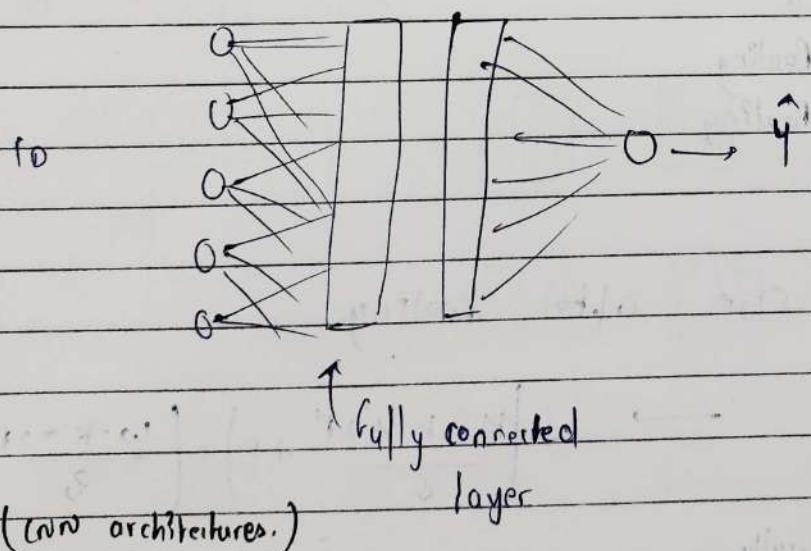
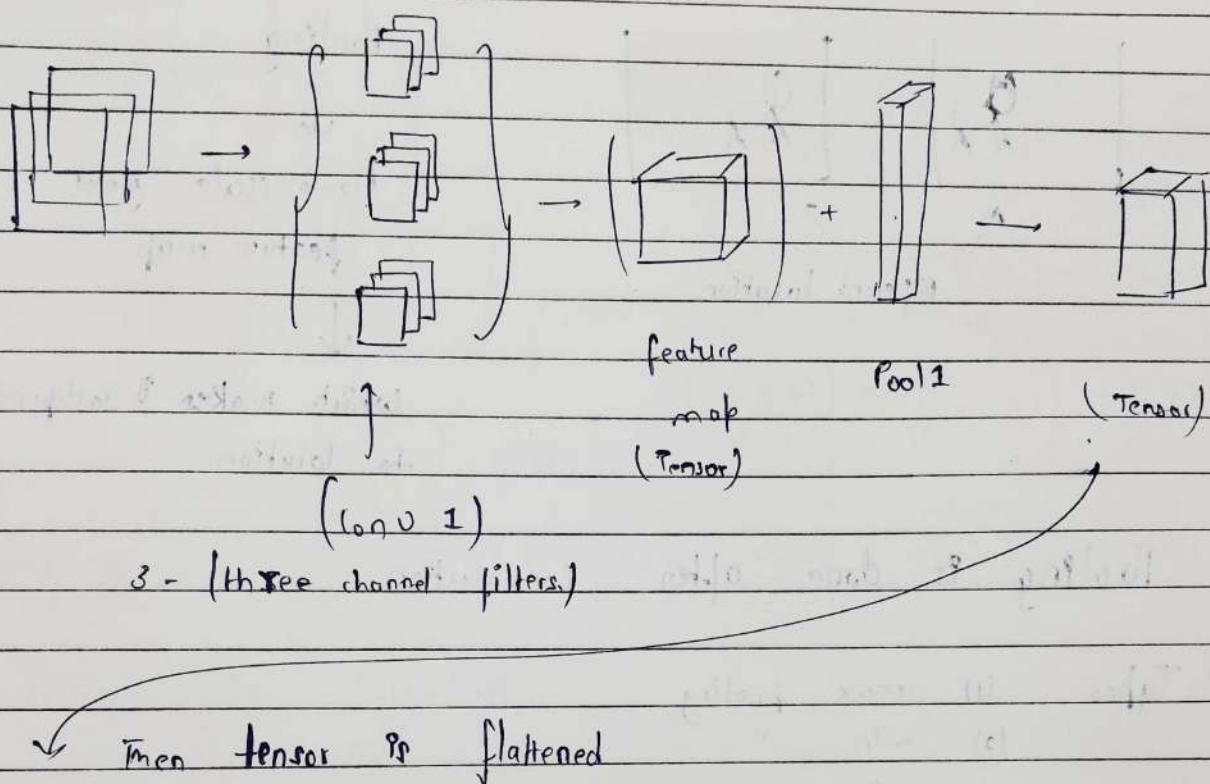
$Padding = 0 = P$

CNN Architecture

PAGE NO.:

DATE: / /

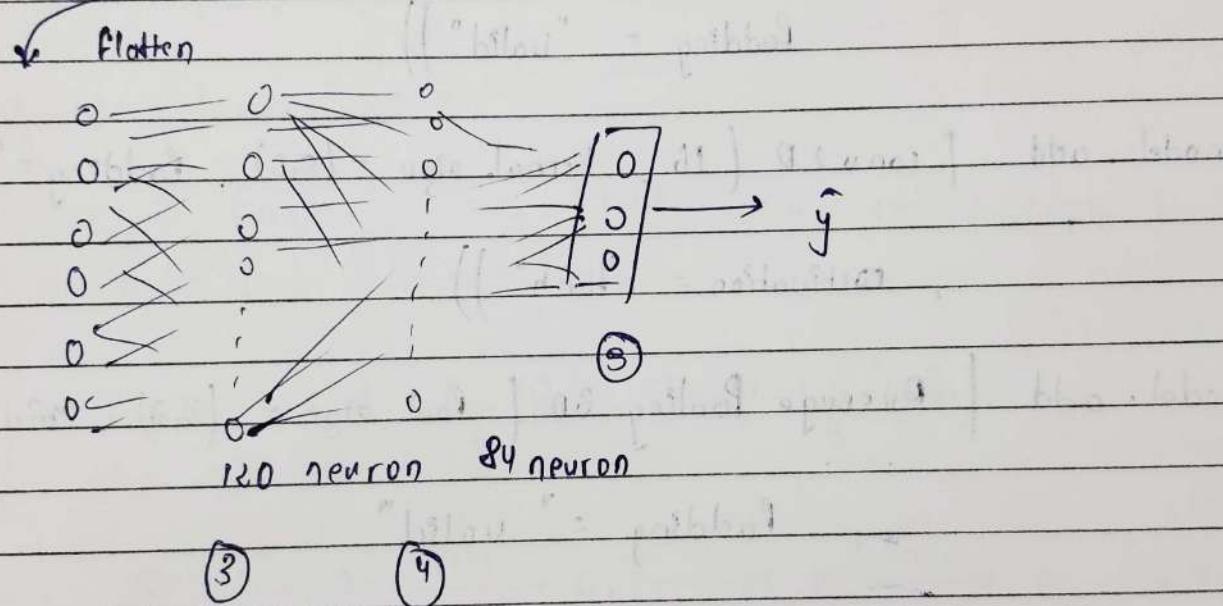
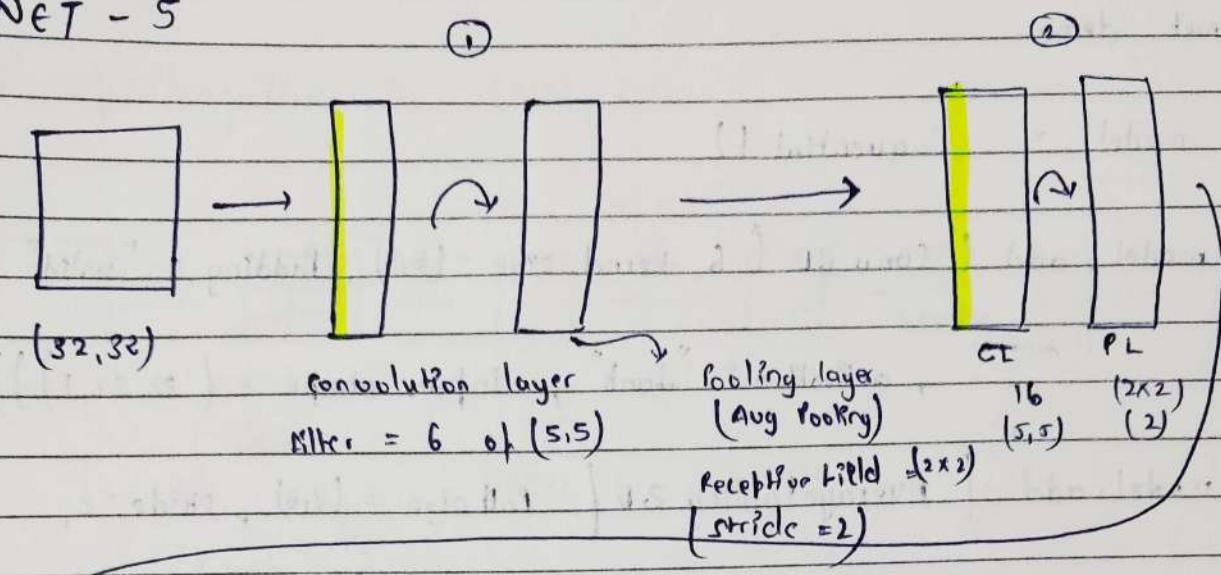
- 1> Convolution
- 2> Padding | stride
- 3> Pooling.



(CNN architectures.)

- | | | |
|----------------|---|---------------|
| (1) Le NET | → | (4) Vgg NET |
| (3) Alex NET | | (5) Res NET |
| (3) Google NET | | (6) Inception |

• LeNET - 5



• Size of Image variation

$$(32, 32) \rightarrow [(28, 28, 6) \rightsquigarrow (14, 14, 16)] \rightarrow [(10, 10, 16) \rightarrow (5, 5, 16)]$$

Flatten, $25 \times 16 = 400 \rightarrow 120 \rightarrow 84 \rightarrow$

features neuron neurons

LeNet demo.

→ model = Sequential()

model.add(Conv2D(6, kernel_size=(5,5), padding="valid",
activation="tanh", input_shape=(32,32,1)))

model.add(AveragePooling2D(pool_size=(2,2), stride=2,
padding="valid"))

model.add(Conv2D(16, kernel_size=(5,5), padding="valid",
activation="tanh"))

model.add(AveragePooling2D(pool_size=(2,2), stride=2,
padding="valid"))

model.add(Flatten())

model.add(Dense(120, activation="tanh"))

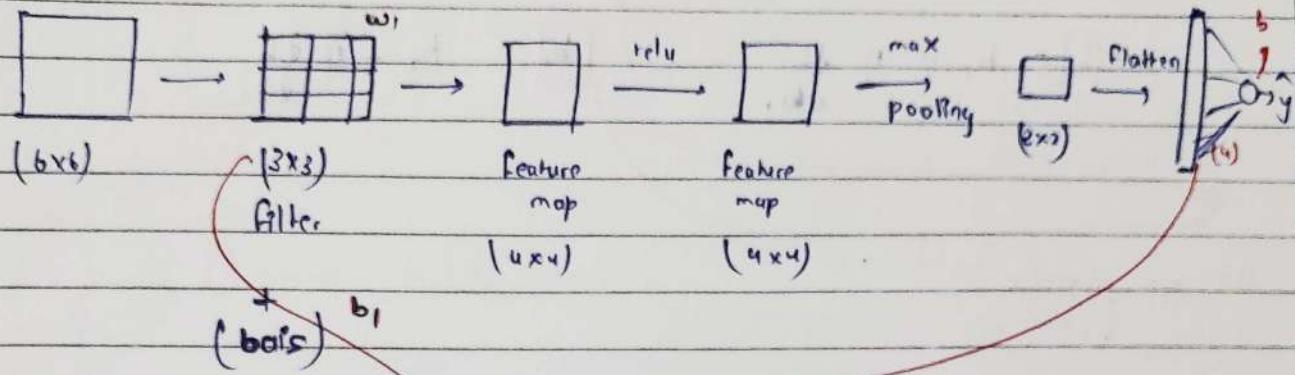
model.add(Dense(84, activation="tanh"))

model.add(Dense(10, activation="softmax"))

In CNN trainable parameters don't depend on Input Image

PAGE NO.: / /
DATE: / /

* Backpropagation in CNN :-

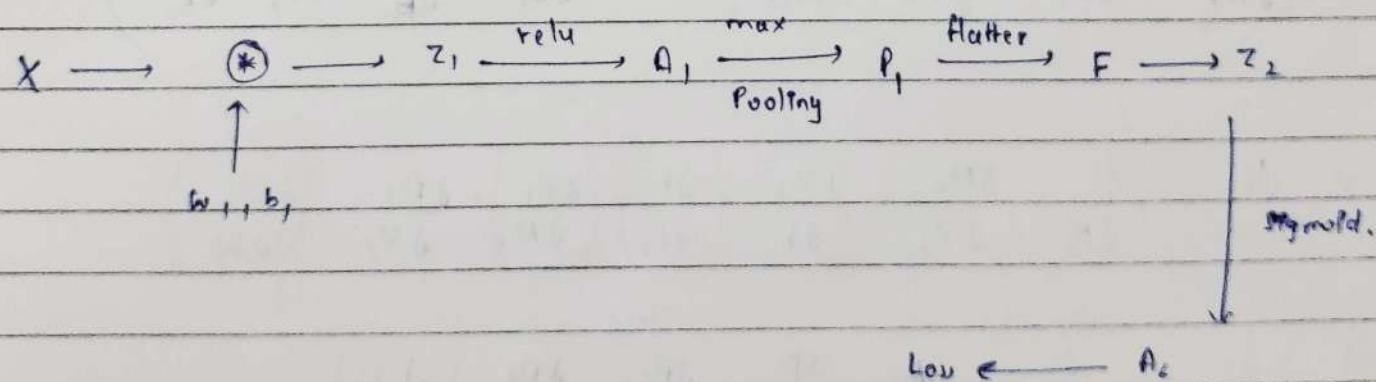


- Trainable parameters

$$w_1 = (3, 3) \quad w_2 = (1, 4) \quad = 15, \text{ trainable parameters}$$

$$b_1 = (1, 1) \quad b_2 = (1, 1)$$

$$\text{Loss} = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$



→ Forward Propagation

$$z_1 = \text{conv}(x, w_1) + b_1$$

$$F = \text{flatten}(P_1)$$

$$A_1 = \text{relu}(z_1)$$

$$z_2 = w_2 F + b_2$$

$$P_1 = \text{max pool}(A_1)$$

$$A_2 = \sigma(z_2)$$

→ Gradient Descent :-

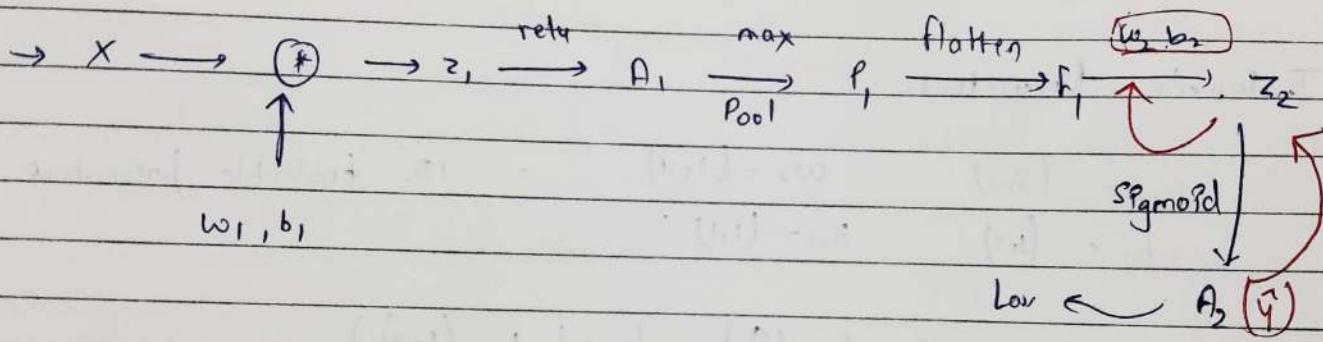
$$w_1 = w_1 - \eta \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

$$b_1 = b_1 - \eta \frac{\partial L}{\partial b_1}$$

$$b_2 = b_2 - \eta \frac{\partial L}{\partial b_2}$$

All are matrix



$$\left(\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}, \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} \right)$$

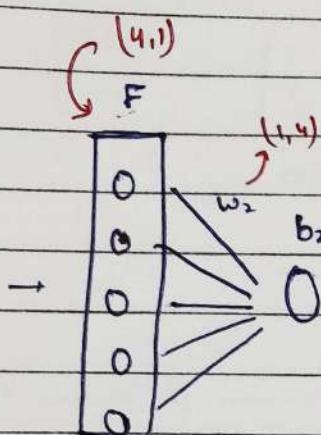
$$\rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial F} \cdot \left(\frac{\partial F}{\partial p_1} \right) \cdot \left(\frac{\partial p_1}{\partial A_1} \right) \cdot \frac{\partial A_1}{\partial z_1} \cdot \left(\frac{\partial z_1}{\partial w_1} \right)$$

$$\rightarrow \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial F}{\partial p} \cdot \frac{\partial p}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \left(\frac{\partial z_1}{\partial b_1} \right)$$

$\text{now } w = ?$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_2} \cdot \frac{\partial A_2}{\partial w_2}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_2} \cdot \frac{\partial A_2}{\partial b_2}$$



→ forward Prop eqⁿ

$$z_2 = w_2 f + b_2$$

$$A_2 = \sigma(z_2)$$

$$\textcircled{1} \rightarrow \frac{\partial L}{\partial A_2} = \frac{\partial}{\partial a_2} [-y_i \log(a_2) - (1-y_i) \log(1-a_2)]$$

$$A_2 \rightarrow \text{matrix} = -\frac{y_i}{a_2} + \frac{(1-y_i)}{1-a_2} = \frac{-y_i(1-a_2) + a_2(1-y_i)}{a_2(1-a_2)}$$

$$= \left\{ \frac{(a_2 - y_i)}{a_2(1-a_2)} \right\}$$

$$\textcircled{2} \rightarrow A_2 = \sigma(z_2)$$

$$\frac{\partial A_2}{\partial z_2} = \sigma(z_2)(1-\sigma(z_2)) = \{a_2(1-a_2)\}$$

$$\textcircled{3} \quad z_2 = w_2 f + b$$

$$\frac{\partial z_2}{\partial w_2} = (F) \quad \frac{\partial z_2}{\partial b_2} = (1)$$

$$\frac{\partial L}{\partial w_2} = \frac{(a_2 - y_i)}{a_2(1-a_2)} \times \cancel{(a_2)(1-a_2)} \times F = (a_2 - y_i) \cdot F$$

PAGE NO.:
DATE: / /

$\downarrow a_2 \rightarrow \text{matrix } (A_2)$

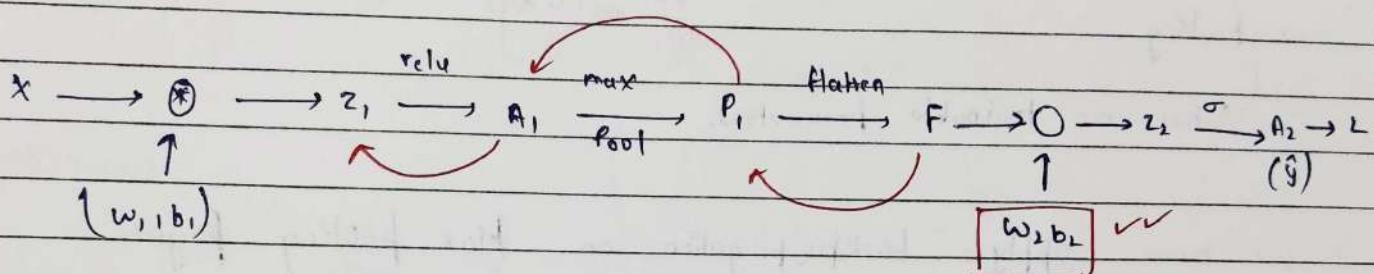
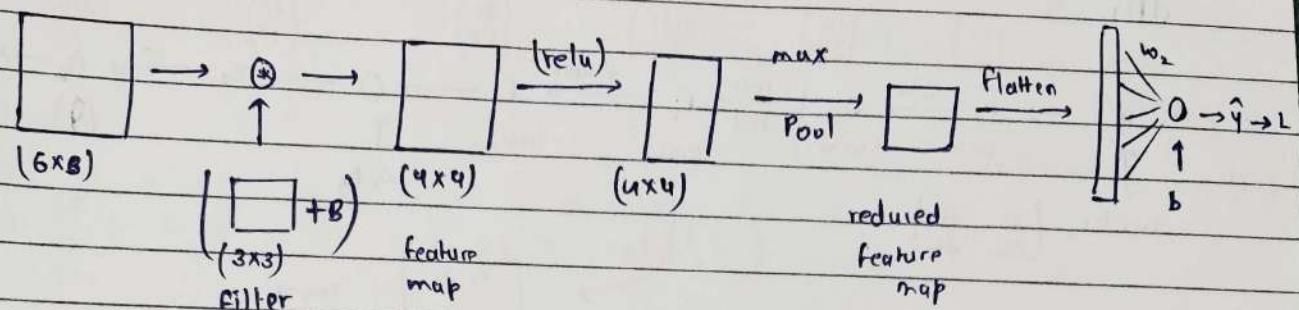
$$\frac{\partial L}{\partial w_2} = (A_2 - Y) F^T$$

$(1,4)$ $(1,1)$ $(4,1)$

$$\frac{\partial L}{\partial b_2} = (A_2 - Y)$$

Back propagation in CNN :-

PAGE NO. :
DATE : / /



→ forward propagation eqn.

$$\begin{aligned} z_1 &= \text{conv}(x, w_1) + b_1 & \left\{ \begin{array}{l} P_1 = \text{maxpool}(A_1) \\ A_1 = \text{relu}(z_1) \end{array} \right. & \left\{ \begin{array}{l} z_2 = w_2 F + b_2 \\ A_2 = \sigma(z_2) \\ F = \text{flatten}(P_1) \end{array} \right. & L = \frac{1}{m} \sum_{i=1}^m \left[-y_i \log(a) - (1-y_i) \log(1-a) \right] \end{aligned}$$

already calculated

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial A_2} \right) \left(\frac{\partial A_2}{\partial z_2} \right) \cdot \frac{\partial z_2}{\partial F} \cdot \frac{\partial F}{\partial P_1} \cdot \frac{\partial P_1}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial F} \cdot \frac{\partial F}{\partial P_1} \cdot \frac{\partial P_1}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$\boxed{\frac{\partial z_2}{\partial F} = w_2}$$

But $\frac{\partial F}{\partial P_1}$ has no trainable parameters.

$$\frac{\partial F}{\partial P_1} = \text{reshape}(P_1, \text{shape})$$

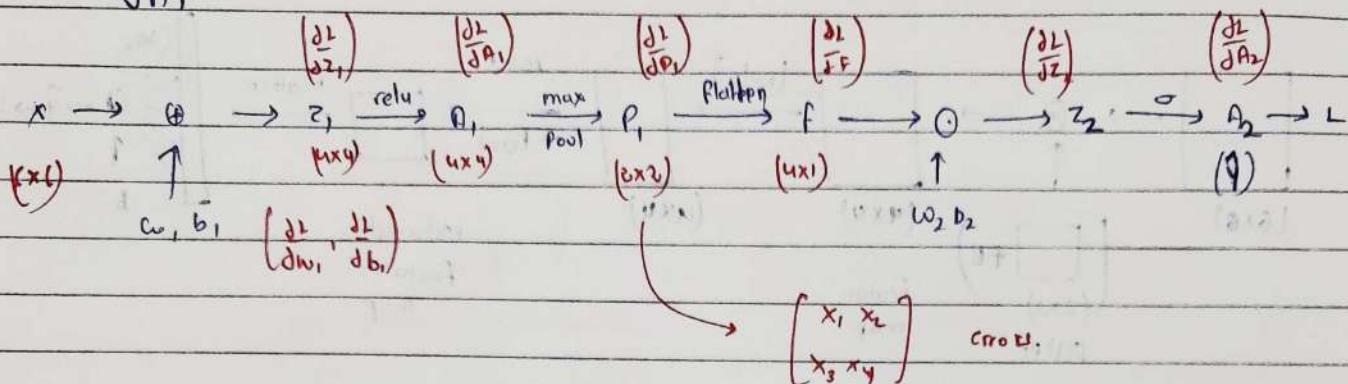
Pooling Operation.

PAGE NO.:

DATE: / /

$\frac{\partial L}{\partial A_1}$

$\frac{\partial A_1}{\partial z_1}$



max pooling

$$\frac{\partial L}{\partial A_1} \Rightarrow (4 \times 4)$$

has no trainable parameters.

We have apply backpropagation on Max pooling layer.

max pool convert $(4 \times 4) \rightarrow (2 \times 2)$

e.g.

$$A_1 = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} \left[\begin{array}{c|c} 5 & 6 \\ \hline 7 & 8 \end{array} \right] \rightarrow \begin{bmatrix} 4 & 8 \\ 12 & 16 \end{bmatrix}$$

The numbers other than the max no

In their window has no

role in predicting 4 / 1

So in back propagation when $\begin{bmatrix} 4 & 8 \\ 12 & 16 \end{bmatrix}$ will be converted

$$\text{to } \begin{bmatrix} 0 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 0 \\ 0 & n & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 \\ 0 & 8 \\ 0 & 0 \\ 0 & 16 \end{bmatrix}$$

either of these
will be asked from
 A_1

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial F} \cdot \frac{\partial F}{\partial p_i} \cdot \frac{\partial p_i}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_i}$$

PAGE NO.:

DATE: / /

$$\frac{\partial L}{\partial p_i}$$

$$\frac{\partial L}{\partial p_i} = \begin{cases} \frac{\partial L}{\partial z_{xy}}, & \text{if } A_{m,n} \text{ is} \\ & \text{the max} \\ & \text{element} \\ 0, & \text{otherwise} \end{cases}$$

→ now

$$\frac{\partial A_1}{\partial z_1} = ?$$

$$A_1 = \text{relu}(z_1)$$

$$\frac{\partial A_1}{\partial z_1} = \begin{cases} 1 & \text{if } z_{1,xy} > 0 \\ 0 & \text{if } z_{1,xy} \leq 0 \end{cases}$$

→ now

$$\frac{\partial z_1}{\partial w_1} = ?$$

$$\frac{\partial z_1}{\partial b_1} = ?$$

Backprop on Convolution

$$x \rightarrow \textcircled{*} \rightarrow z_1$$

$(3,3)$
 \uparrow
 (w, b)
 $(2,2)$

$$\left(\frac{\partial L}{\partial z_1} \right) \rightarrow (2,2)$$

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$$

$$w_1 = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \quad z_1 = \begin{pmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \end{pmatrix}$$

$$z_{11} = x_{11} w_{11} + x_{12} w_{12} + x_{21} w_{21} + x_{22} w_{22} + b_1$$

$$z_{12} = x_{12} w_{11} + x_{13} w_{12} + x_{22} w_{21} + x_{23} w_{22} + b_1$$

$$z_{21} = x_{21} w_{11} + x_{22} w_{12} + x_{31} w_{21} + x_{32} w_{22} + b_1$$

$$z_{22} = x_{22} w_{11} + x_{23} w_{12} + x_{32} w_{21} + x_{33} w_{22} + b_1$$

after convolution
operation

$$\frac{dL}{db_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{db_1}$$

PAGE NO.: _____
DATE: / /

$$\frac{dL}{dz_1} = \begin{bmatrix} \frac{dL}{dz_{11}} & \frac{dL}{dz_{12}} \\ \frac{dL}{dz_{21}} & \frac{dL}{dz_{22}} \end{bmatrix}$$

$$\frac{dL}{db_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{db_1} = \frac{dL}{dz_{11}} \cdot \frac{dz_{11}}{db_1} + \frac{dL}{dz_{12}} \cdot \frac{dz_{12}}{db_1} + \frac{dL}{dz_{21}} \cdot \frac{dz_{21}}{db_1} + \frac{dL}{dz_{22}} \cdot \frac{dz_{22}}{db_1}$$

↑ ↑ ↑ ↑

$$\Rightarrow \left(\frac{dL}{dz_{11}} + \frac{dL}{dz_{12}} + \frac{dL}{dz_{21}} + \frac{dL}{dz_{22}} \right)$$

$$= \text{sum} \left(\frac{\frac{dL}{dz_i}}{\text{of all terms}} \right)$$

$$\rightarrow \frac{dL}{dw_1} = \{$$

$$\frac{dL}{dw_1} = \begin{bmatrix} \frac{dL}{dw_{11}} & \frac{dL}{dw_{12}} \\ \frac{dL}{dw_{21}} & \frac{dL}{dw_{22}} \end{bmatrix}$$

$$\frac{dL}{dw_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{dw_1}$$

$$\frac{dL}{dz_1} = \begin{bmatrix} \left(\frac{dL}{dz_{11}} \right) & \frac{dL}{dz_{12}} \\ \frac{dL}{dz_{21}} & \frac{dL}{dz_{22}} \end{bmatrix}$$

$$\frac{dL}{dw_{11}} = \frac{dL}{dz_{11}} \cdot \frac{dz_{11}}{dw_{11}} + \frac{dL}{dz_{12}} \cdot \frac{dz_{12}}{dw_{11}} + \frac{dL}{dz_{21}} \cdot \frac{dz_{21}}{dw_{11}} + \frac{dL}{dz_{22}} \cdot \frac{dz_{22}}{dw_{11}}$$

$$\frac{dL}{dw_{12}} =$$

$$\frac{dL}{dw_{21}} =$$

$$\frac{dL}{dw_{22}} =$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial b_{11}} \cdot X_{11} + \frac{\partial L}{\partial z_{12}} \cdot X_{12} + \frac{\partial L}{\partial z_{21}} \cdot X_{21} + \frac{\partial L}{\partial z_{22}} \cdot X_{22}$$

PAGE NO.: / /
DATE: / /

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial z_{11}} \cdot X_{12} + \frac{\partial L}{\partial z_{12}} \cdot X_{13} + \frac{\partial L}{\partial z_{21}} \cdot X_{22} + \frac{\partial L}{\partial z_{22}} \cdot X_{23}$$

$$\frac{\partial L}{\partial w_{21}} = - - - - -$$

$$\frac{\partial L}{\partial w_{22}} = - - - - -$$

we can say.

$$\frac{\partial L}{\partial w_1} = \text{sum } V \left(X, \frac{\partial L}{\partial z_1} \right)$$

↑ ↑
matrixs matrix

$$X = \begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix} \quad \frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

* Conclusion :-

$$\frac{\partial L}{\partial w_1} = \text{sum } \left(X, \frac{\partial L}{\partial z_1} \right)$$

$$\frac{\partial L}{\partial b_1} = \text{sum } \left(\frac{\partial L}{\partial z_1} \right)$$

• Pretrained model.

→ image NET dataset.
 ↳ images 1.4 cores

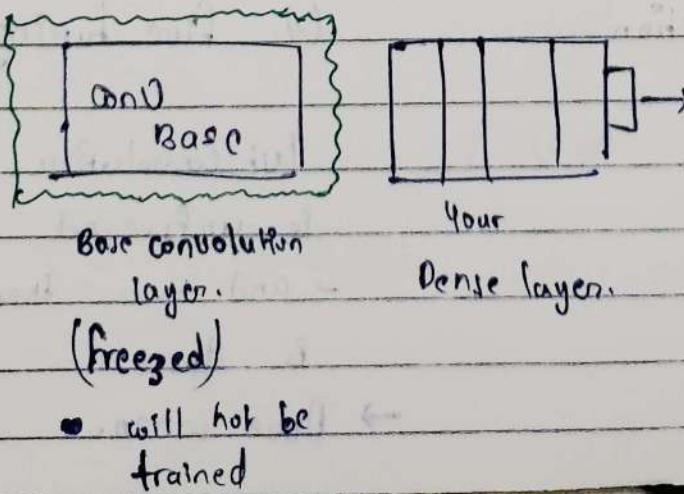
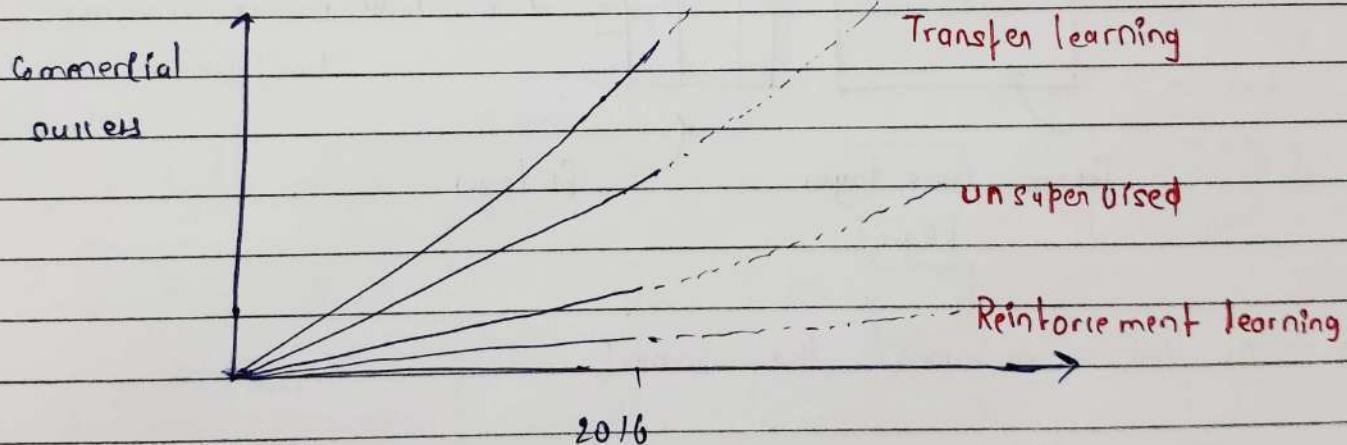
2010	→	ML model	→	28%	Revolution
2011	→	ML model	→	26%	
2012	→	Alex Net	→	16.4%	
2013	→	ZF Net	→	11.7%	
2014	→	VGG	→	7.3%	
2015	→	Google Net	→	6.7%	
2016	→	RESNET	→	3.5%	

Humans → 5% error rate.

→ code is implemented in jupyter notebook.

Transfer learning.

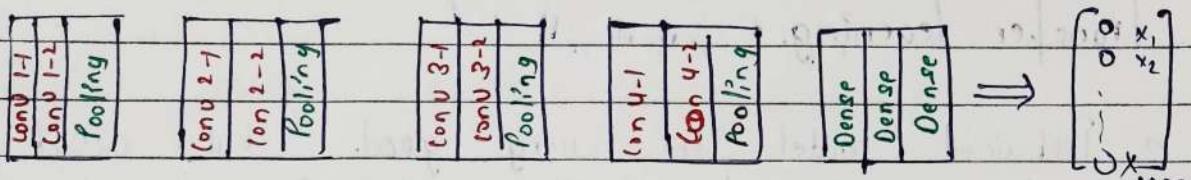
- Pretrained model are very good. but what to do when pretrained model do not have the categories that we want to classify.
- That's where transfer learning comes in.
- TF learning is a research problem in machine learning that focuses on ~~sharing~~ knowledge gained while solving one problem and applying it to a different but ~~similar~~ related problem



trained on ImageNet
data set

PAGE NO.:
DATE: / /

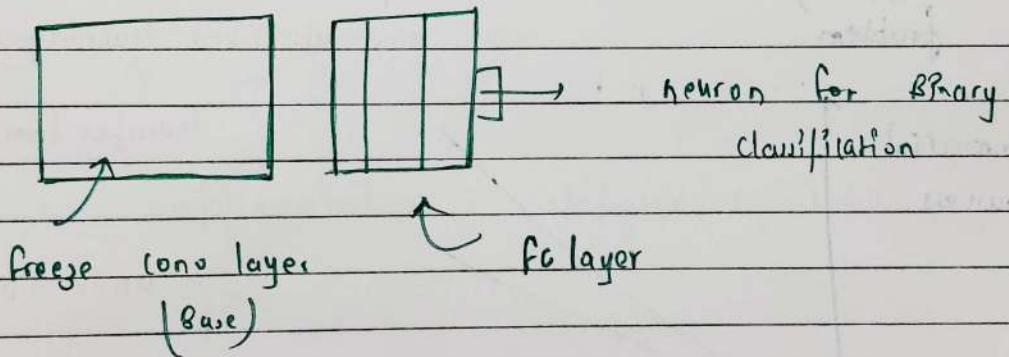
Ugg 16
model.



Convolutional Base
layers.

Fc layer.

→ now, convolutional layer is freezed and then used to perform another problem by taking convolutional layer as it is



Do not reinvent the wheel.

• Two Type of T.F.

(1) Feature extraction

(Discussed above)

(2) Fine tuning,

last convolution layer
is unfreezed

→ And then then training
is done.

→ Done when,

Fine Tuning

→ Done when we have to do classification which is not in our training dataset classification labels.

* Non linear neural networks.

↳ Why needed.

Example →

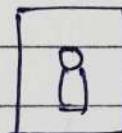


image
dataset

humans → faces

25000

images

(image)

To solve there can be two ways.

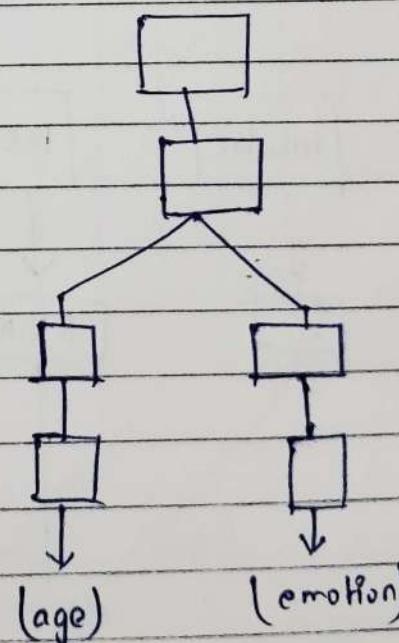
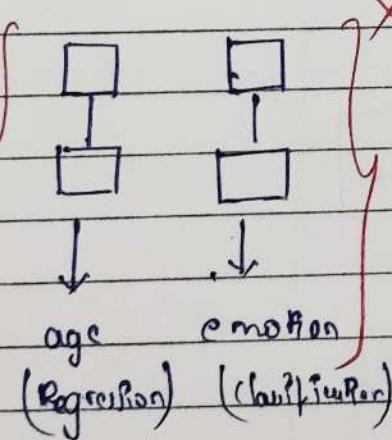
age ?

emotion

→ happy

→ sad

→ angry.

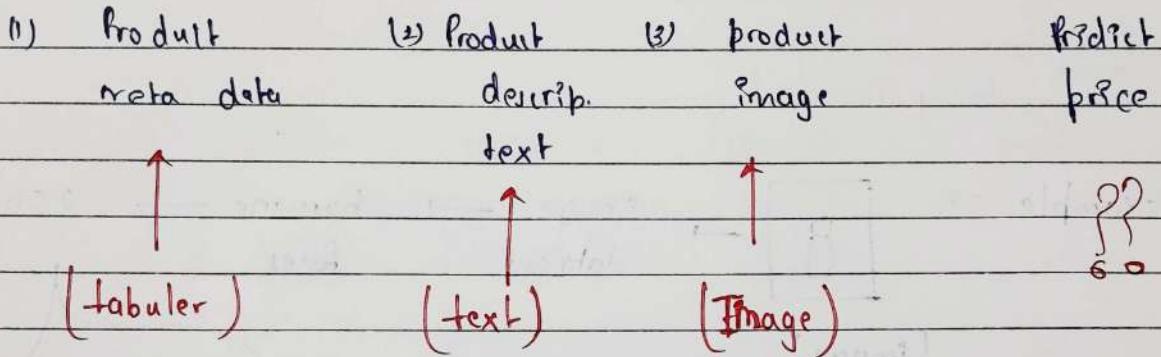


ex :-

let say you are working in commerce company.

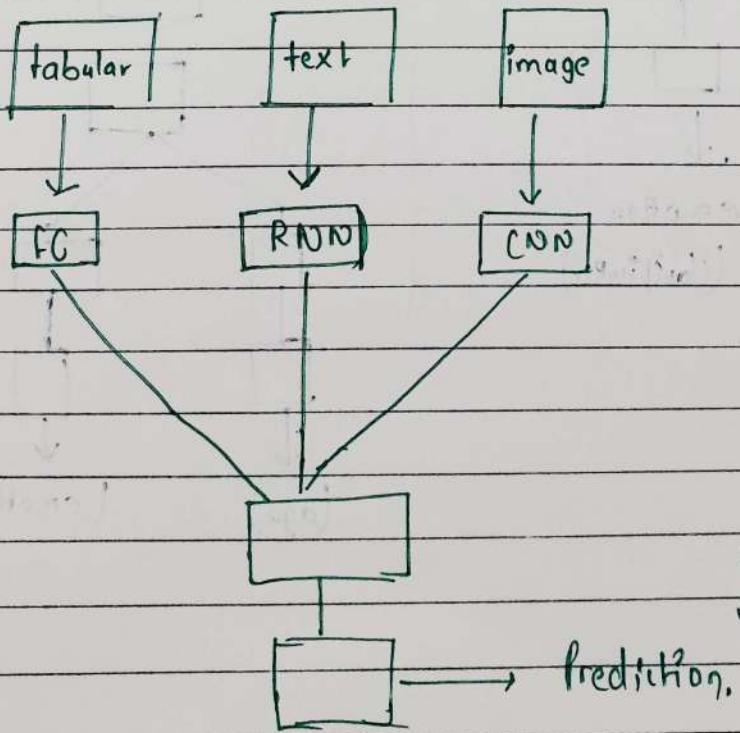
↪ certain information is given to you.

↪ Data given of product.



→ Good approach to solve this problem
can be :-(non-linear neural network)

* *

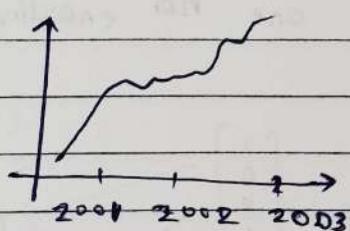


Recurrent neural network (RNN)

Depend on sequential data. (sequence matters).

eg. ① → text { Hi my name is Akash }

③ → Time series.



③ speech

④ DNA sequence.

→ why are ANN?

let's say we need to do text classification using ANN.

(Input)	Output
Hi my name is Akash	0
I love cambur(x)	0
India won the match	1

ANN will not understand words
we need to vectorize them to
feed them in neural network.

methods.

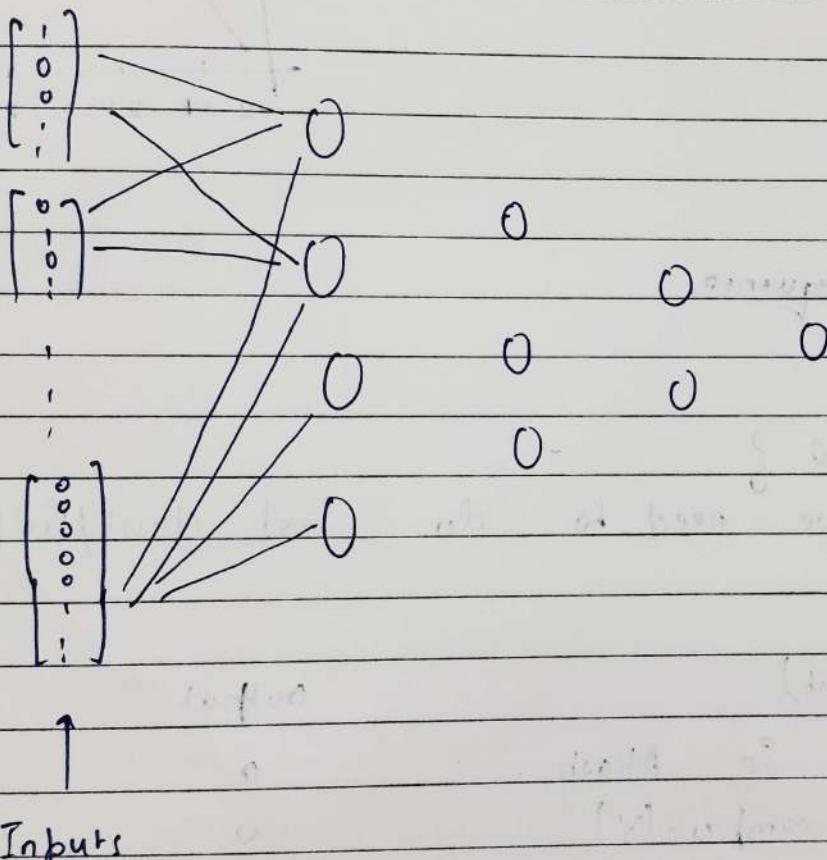
- one hot encoding
- bag of words.

Hi my name is Akash.

$[1, 0, 0, \dots]$ $[0, 1, 0, \dots]$ $[0, 0, 1, 0, \dots]$... $[0, 0, 0, 0, 1, 0, \dots]$

12 items

one hot encoding vectorization.



Problem → Input size can be different with every sentence.

Solution → zero padding.

eg. I love Jaipur.

PAGE NO.:

DATE: / /

But input needs to be of 5 vectors. so.

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

0

0

0

0

0 0

0

0

0

← This is zero padding.

But it has problem like

(1) unnecessary computation

(2) text input → varying size.

(3) prediction problem

(4) totally disregarding the sequence info.

Projects

- Sentiment analysis demo
- Sentence prediction
- Generate image caption

PAGE NO.:

DATE: / /

* Advanced Model → Sent. (Question Answer Demo)

How data is given to RNN.

e.g. → movie was good.

Vocab → 5 words. (word with max digits)

movie was good.
 $(1, 0, 0, 0, 0)$ $[0, 1, 0, 0, 0]$ $[0, 0, 1, 0, 0]$

→ $\left[[1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0] \right]$

$\underbrace{[[1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0]]}_{(3, 5)}$

no. of time steps

of features.

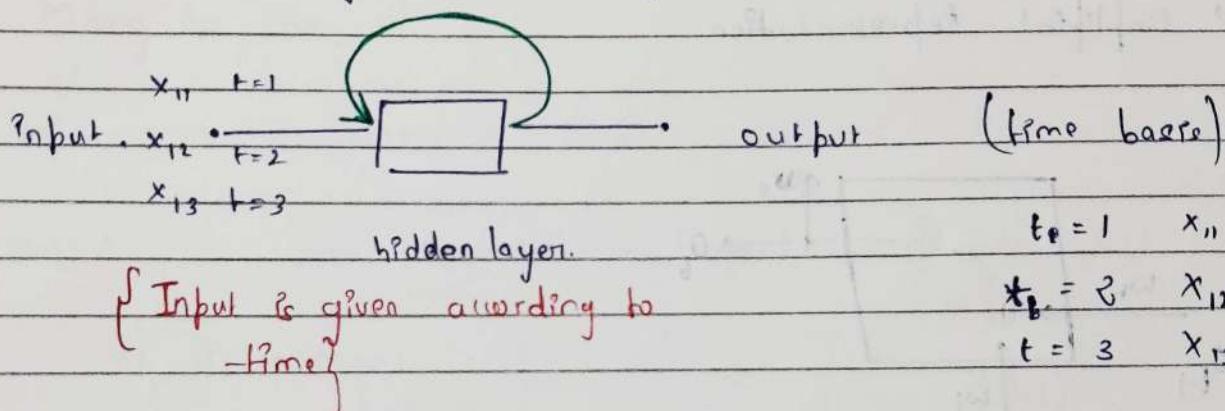
(Keras) → simple RNN → batch size, time step, input features
 $(3, 4, 5)$

How RNN works.

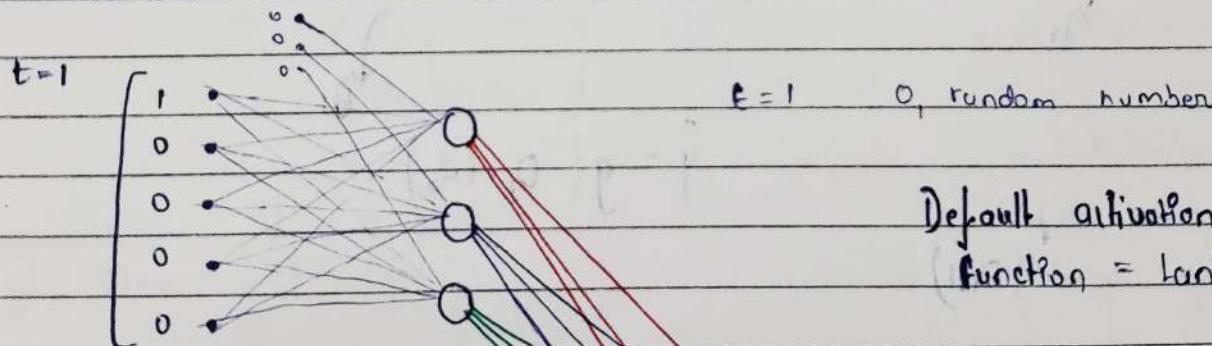
PAGE NO:

DATE: / /

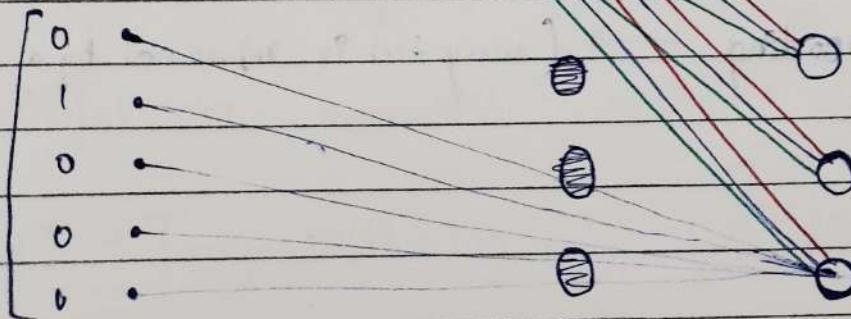
Review	{ sentiments }			
$x_{11} \quad x_{12} \quad x_{13}$ x_1 movie was good		1		RNN \rightarrow ANN \rightarrow 2 big difference
$x_{21} \quad x_{22} \quad x_{23}$ x_2 movie was bad		0		
$x_{31} \quad x_{32} \quad x_{33} \quad x_{34}$ x_3 movie was not good		0		

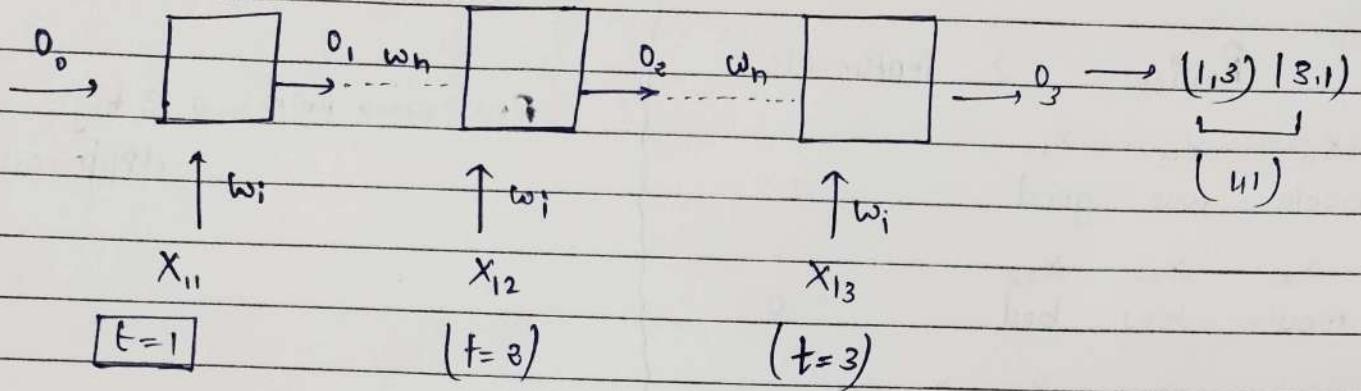


→ Feedback mechanism :-

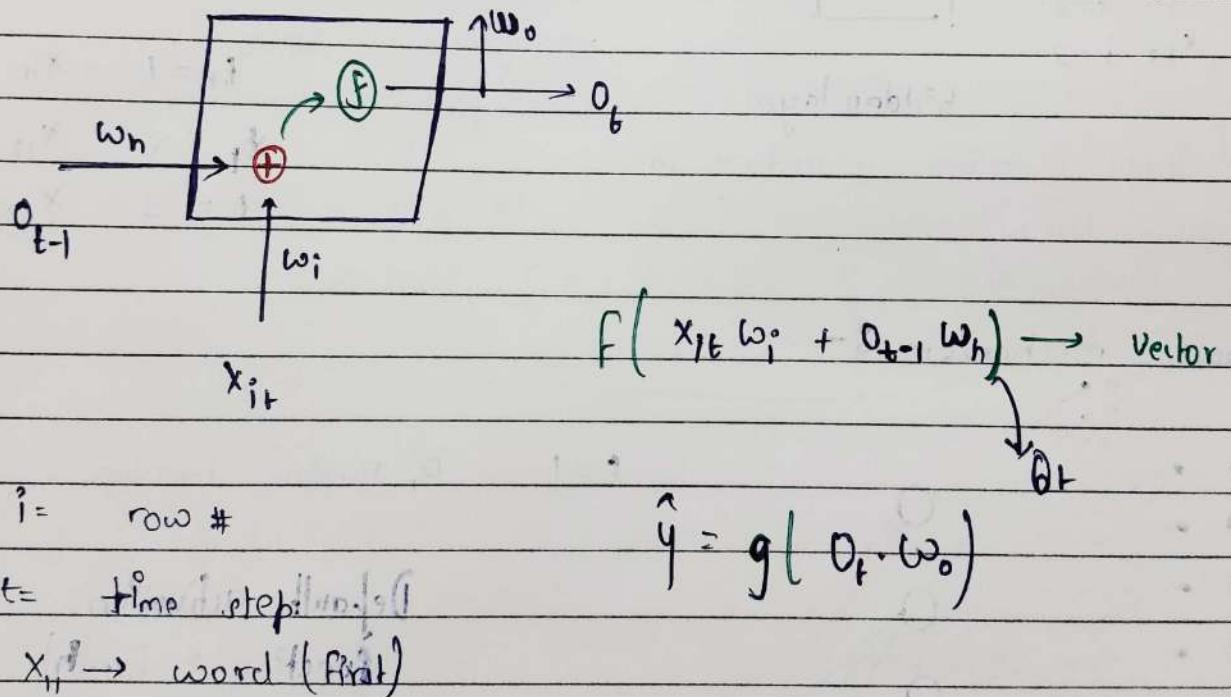


$t=2$





→ Simplified Representation



* Text to Vector conversion techniques.

- ① integer encoding → (every text is represented by a number)
- ② embedding

* Types of RNN

① Many to one.

③ Many to many

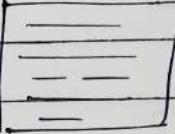
② One to many

④ One to One.

→ Many to one.

Input → Sequence → Sentence, characters, time series.

Output → non-sequential → integer/num → scalar (0,1)

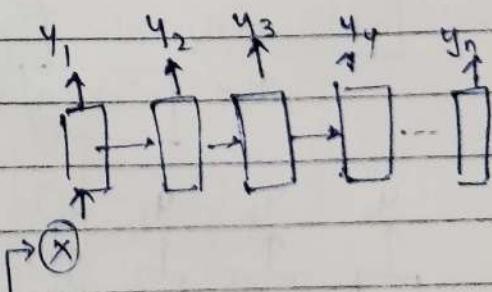
e.g. ① sentiment analysis] →  ⇒ (0,1)

② Rating Prediction] → 

multiclass classification

② One to Many :-

• Input → normal non sequential
e.g. (Image)



• Output → sequential

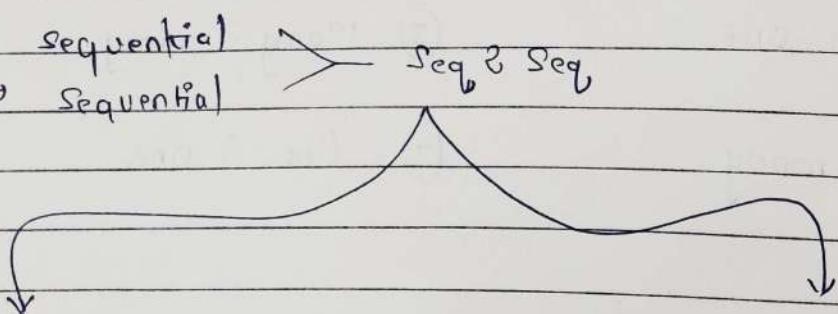
e.g. → Image captioning.

 → [NN] → A man is playing cricket.

→ music generation

③ Many to many →

input → sequential
output → sequential



Same length

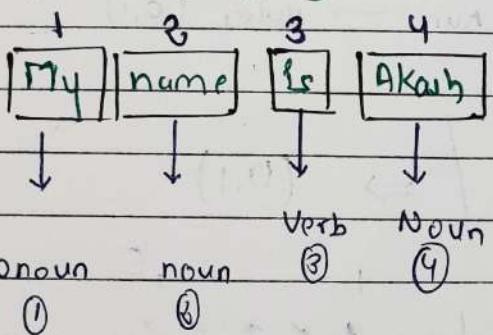
Variable length

input Seq = output seq.

e.g.: machine translation.

e.g. → POS Tagging → NLP

1 language → other language



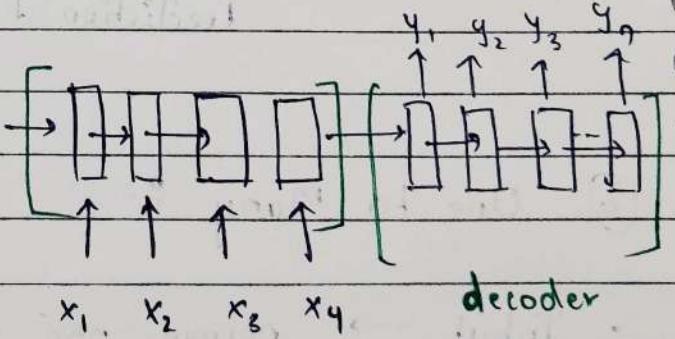
→ Hi my name is Akash

→ मेरा नाम आकाश है

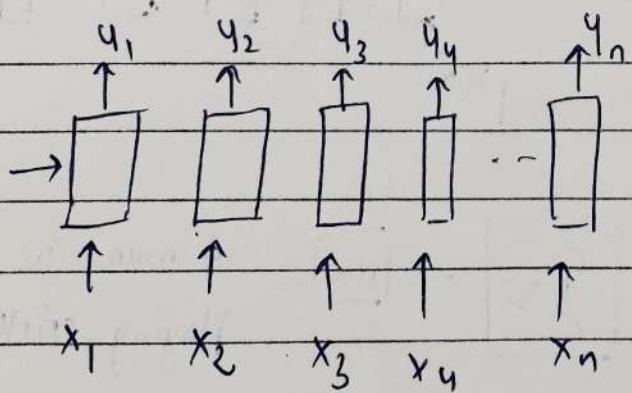
(google translate)

② Name entity Recognition (NER)

→ lets meets at [7pm] at
the [airbork]



encoder



④ One-to-One [Technically not a neural network]

input → non-sequential

→ image classification.

output → non-sequential

Back Propagation in RNN

Through an example of : (Many to one RNN)

text → I/O

→ let's take a toy dataset.

text

			X	y
Cat	mat	rat	$x_1 [1\ 0\ 0] [0\ 1\ 0] [0\ 0\ 1]$	1
rat	rat	mat	$x_2 [0\ 0\ 1] [0\ 0\ 1] [0\ 1\ 0]$	1
mat	mat	Cat.	$x_3 [0\ 1\ 0] [0\ 1\ 0] [1\ 0\ 0]$	0

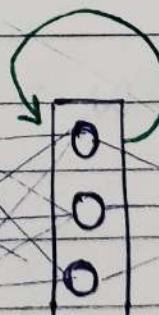
Vocab → cat mat rat
 $[1\ 0\ 0] [0\ 1\ 0] [0\ 0\ 1]$

$t=3 \quad b=2 \quad g=1$

0 0 1

0 1 → 0

1 0 0



$w_h (3 \times 3)$

$a, b, c, d, e, f, g, h, i$

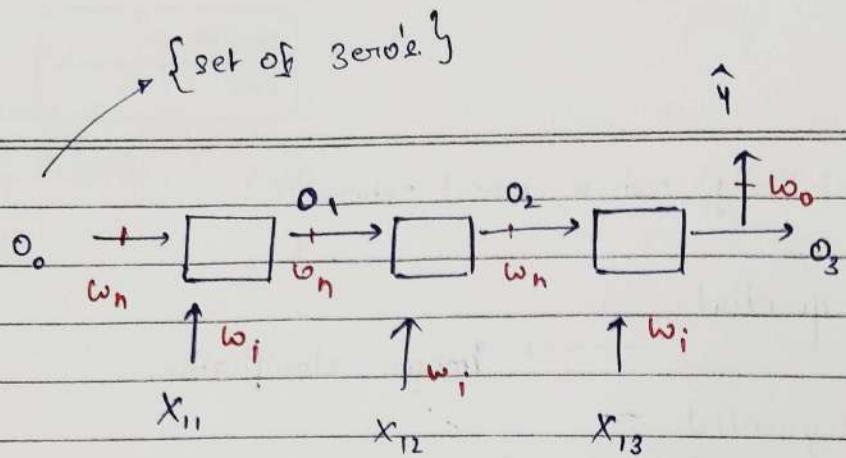
$o - \hat{y}$

w_o

w_i
 $(3, 3)$

$(3, 1) \quad (1 \text{ bias})$

(3 bias)



$$o_1 = f(u_1 w_i + o_0 w_n) \quad y \rightarrow \hat{y} = \sigma(o_3 w_0)$$

$$o_2 = f(x_{12} w_i + o_1 w_n) \quad L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

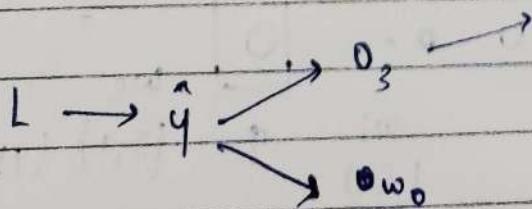
$$o_3 = f(x_{13} w_i + o_2 w_n)$$

→ gradient descent

$$w_i = w_i - \eta \boxed{\frac{\partial L}{\partial w_i}}$$

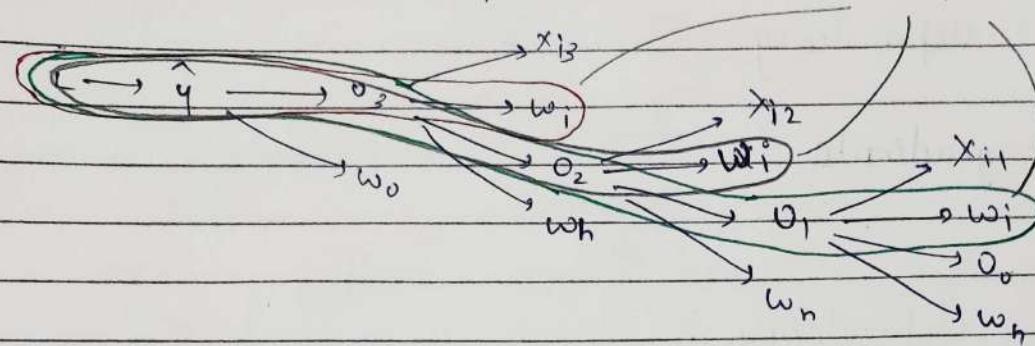
$$w_n = w_n - \eta \boxed{\frac{\partial L}{\partial w_n}}$$

$$w_0 = w_0 - \eta \boxed{\frac{\partial L}{\partial w_0}}$$



$$\frac{dL}{dw_0} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_0}$$

3 paths.



$$\frac{dL}{dw_i} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_i} +$$

$$\frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_2} \cdot \frac{d\omega_2}{d\omega_i} +$$

$$\frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_2} \cdot \frac{d\omega_2}{d\omega_1} \cdot \frac{d\omega_1}{d\omega_i}$$

$$\frac{dL}{d\omega_i} = \sum_{j=1}^3 \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_j} \cdot \frac{d\omega_j}{d\omega_i}$$

$$\frac{dL}{dw_n} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_n} +$$

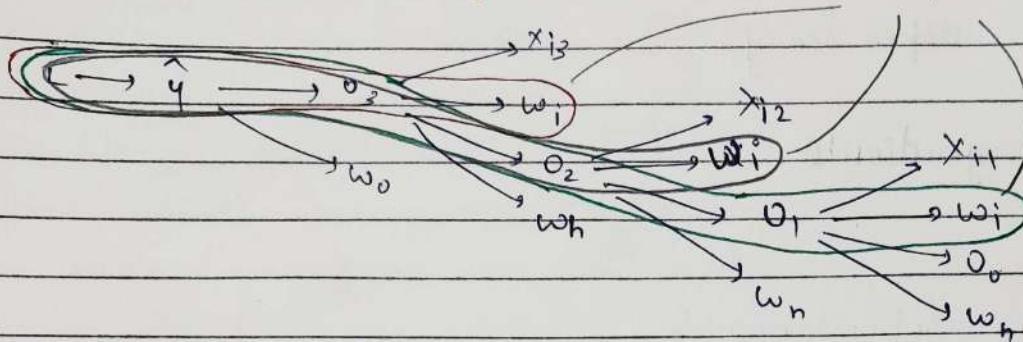
$$\frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_2} \cdot \frac{d\omega_2}{d\omega_n} +$$

$$\frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_3} \cdot \frac{d\omega_3}{d\omega_2} \cdot \frac{d\omega_2}{d\omega_1} \cdot \frac{d\omega_1}{d\omega_i}$$

$$\frac{dL}{d\omega_n} = \sum_{j=1}^3 \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\omega_j} \cdot \frac{d\omega_j}{d\omega_n}$$

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial w_0}$$

3 paths.



$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_i} +$$

$$\frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_2} \cdot \frac{\partial o_2}{\partial w_i} +$$

$$\frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_2} \cdot \frac{\partial o_2}{\partial w_1} \cdot \frac{\partial o_1}{\partial w_i}$$

$$\left[\frac{\partial L}{\partial w_i} = \sum_{j=1}^3 \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_i} \right]$$

$$\frac{\partial L}{\partial w_n} = \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial w_n} +$$

$$\frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial w_n} +$$

$$\frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_i}$$

$$\rightarrow \frac{\partial L}{\partial w_n} = \sum_{j=1}^3 \frac{\partial L}{\partial \hat{q}} \cdot \frac{\partial \hat{q}}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_n}$$

* Problem of RNN

① long term dependency.

② Unstable gradients

① long term dependency :-

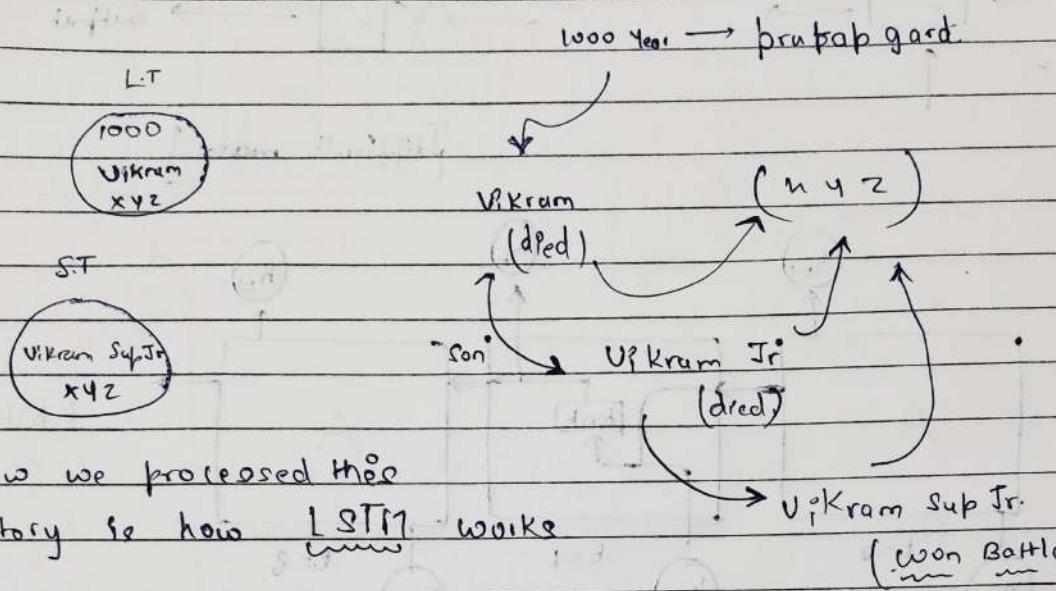
eg :- Marathi is spoken in Maharashtra.] short term dependence

→ Maharashtra :- a beautiful place

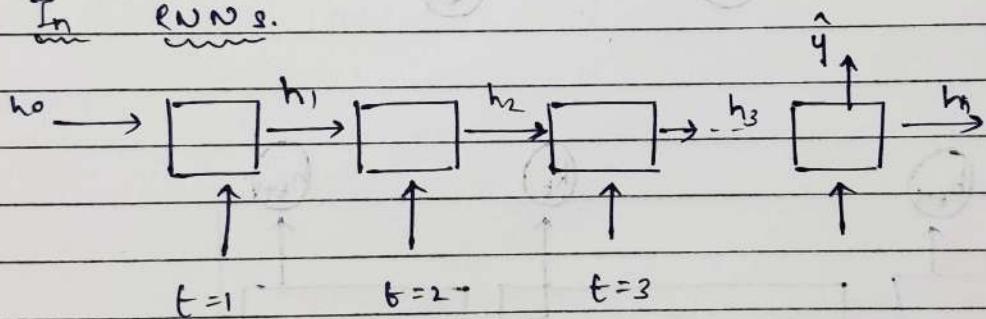
I went there last year

But I could not enjoy there → long term
B/c I could not speak Marathi dependence

LSTM (long short Term memory)

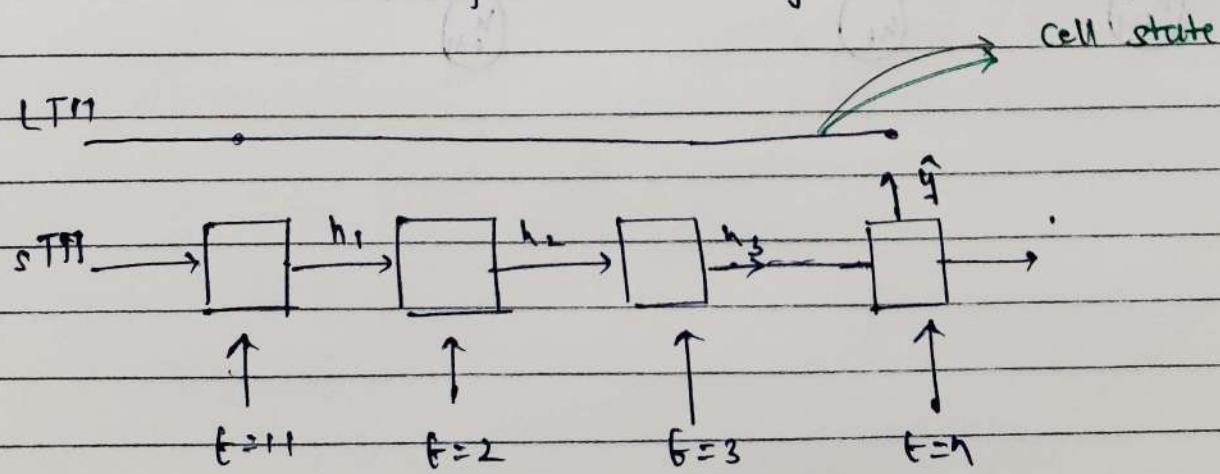


In RNNs.

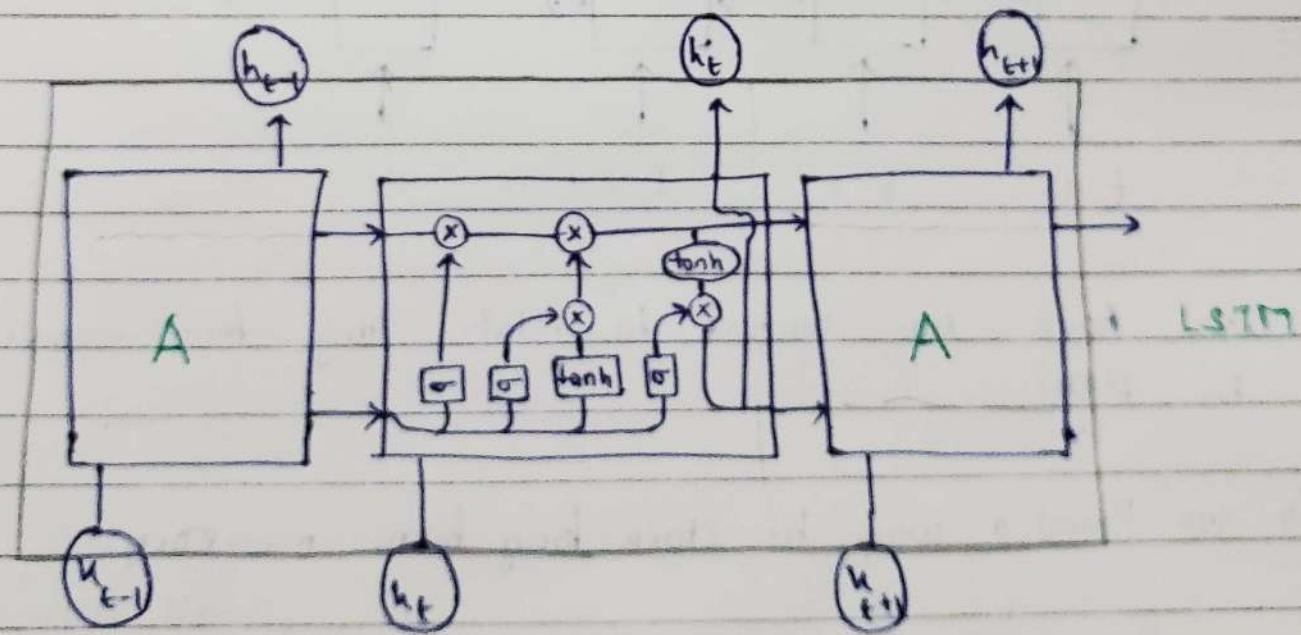
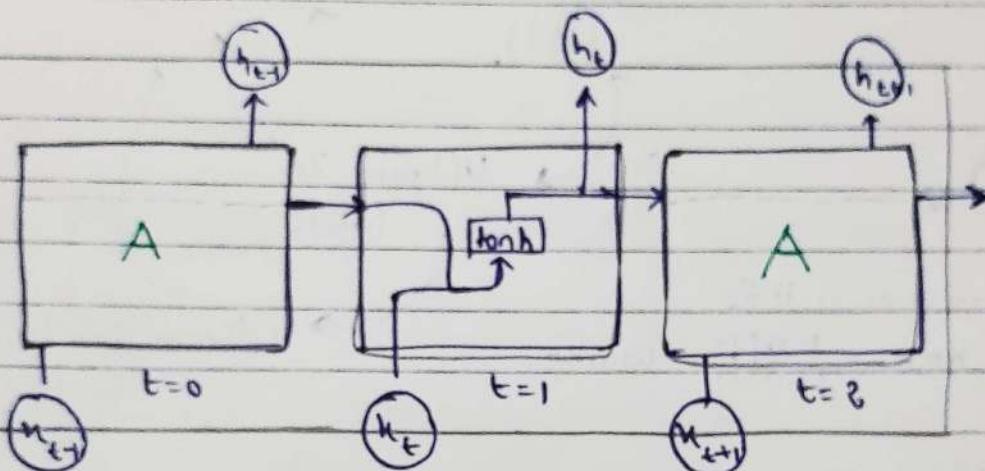
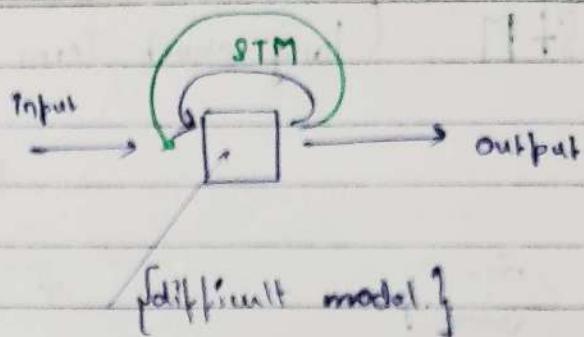
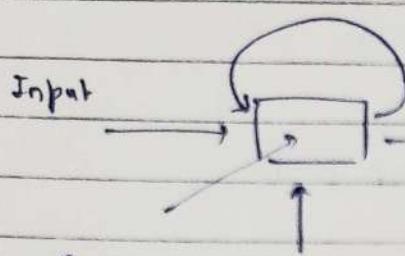


- we have no context to retain long term memory in RNN

→ so we Found a way to store long term memory.



LTM



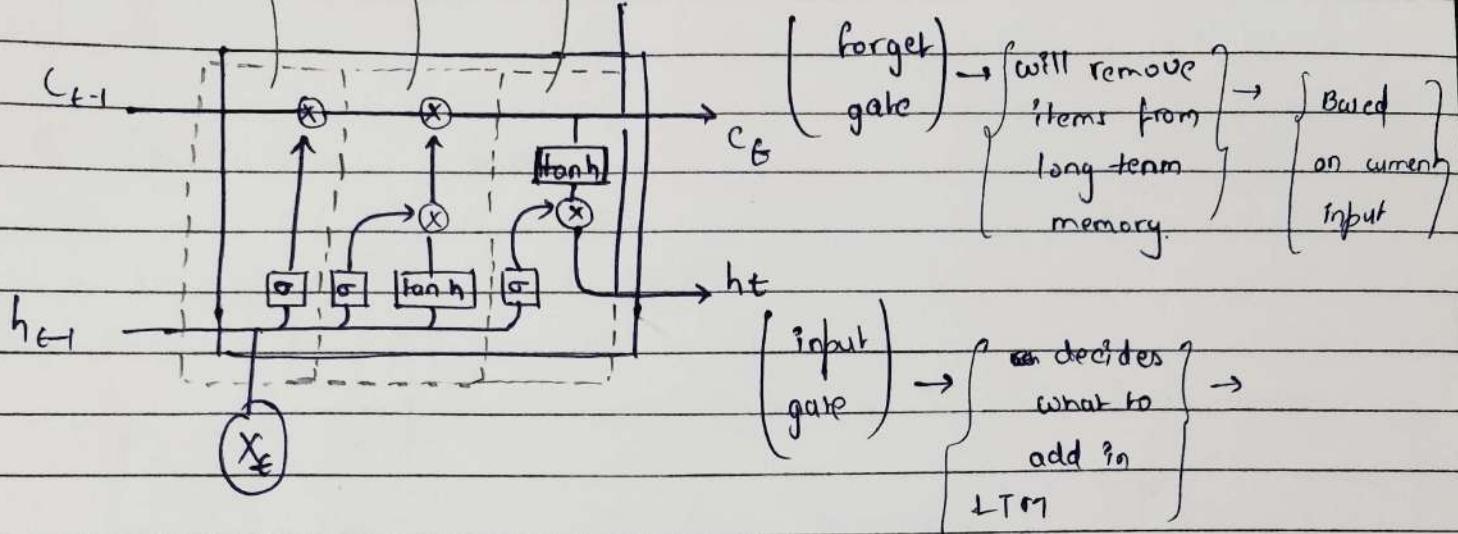
forget
gate

input
gate

output
gate

PAGE NO.:

DATE: / /



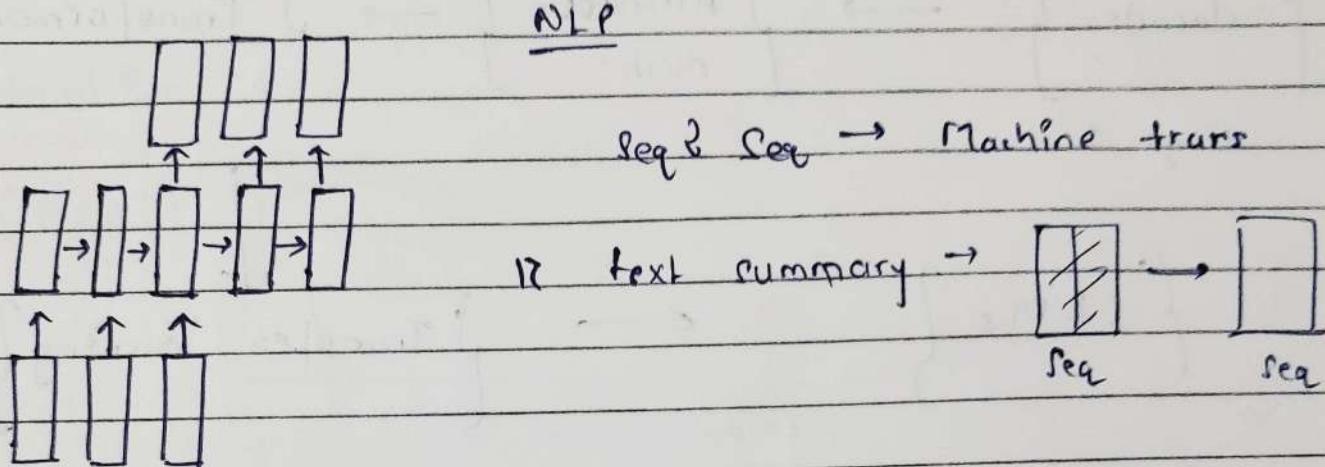
(Output) \rightarrow { what should be shown } \rightarrow { Based on input }

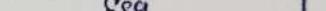
gate

as output from LTM

Large language Models.

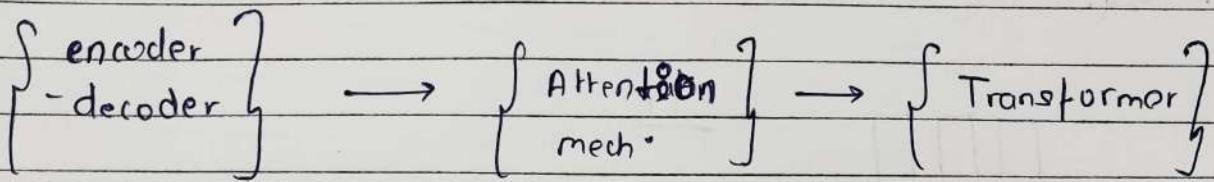
* Seq. & Seq tasks



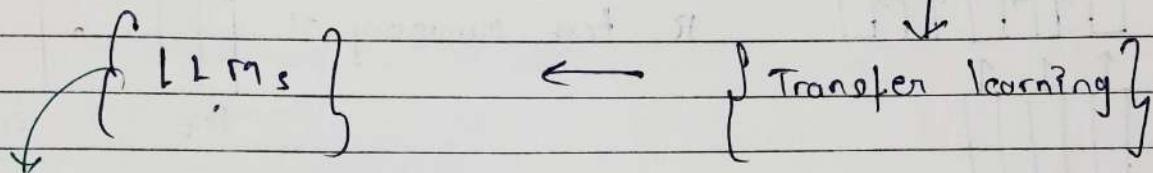
3) Question Answering → 

4) speech → text →
↓ ↓
Seq. Seq.

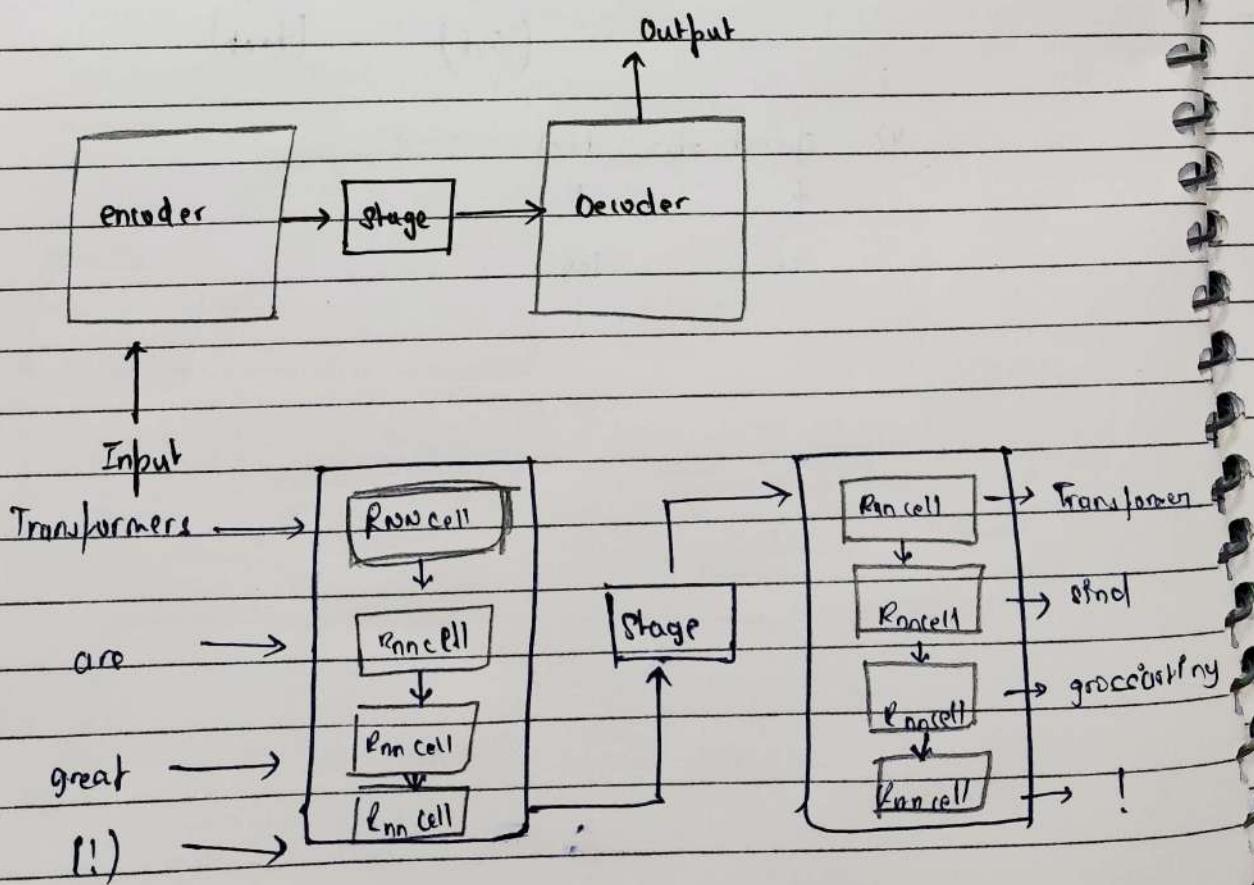
History of seq. & seq. models.



stage 1

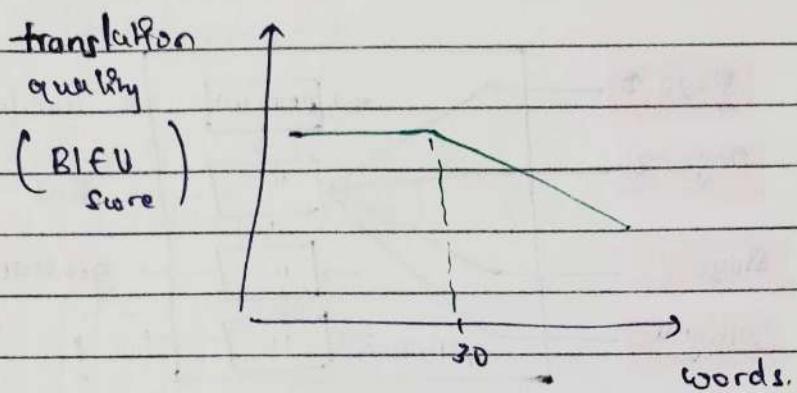


(1) Encoder-decoder



(stage) → consist of the summary vector of the input

Problem → ?- very long input seq. then output
? not good.

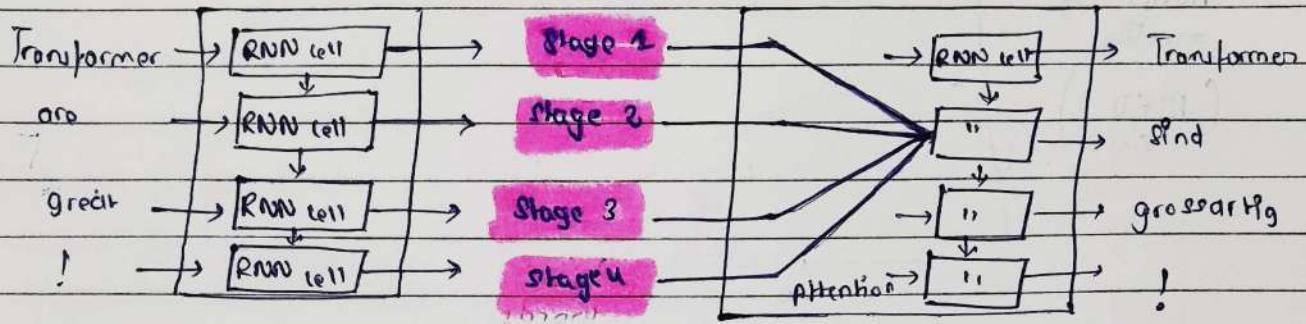


→ A potential issue with this encoder-decoder approach is that a neural network has to be able to, compress all the necessary information of a source sentence into a fixed length vector. This may make it difficult for the neural network to cope with long sentences.

* Attention mechanism.

Paper \Rightarrow Neural Machine translation

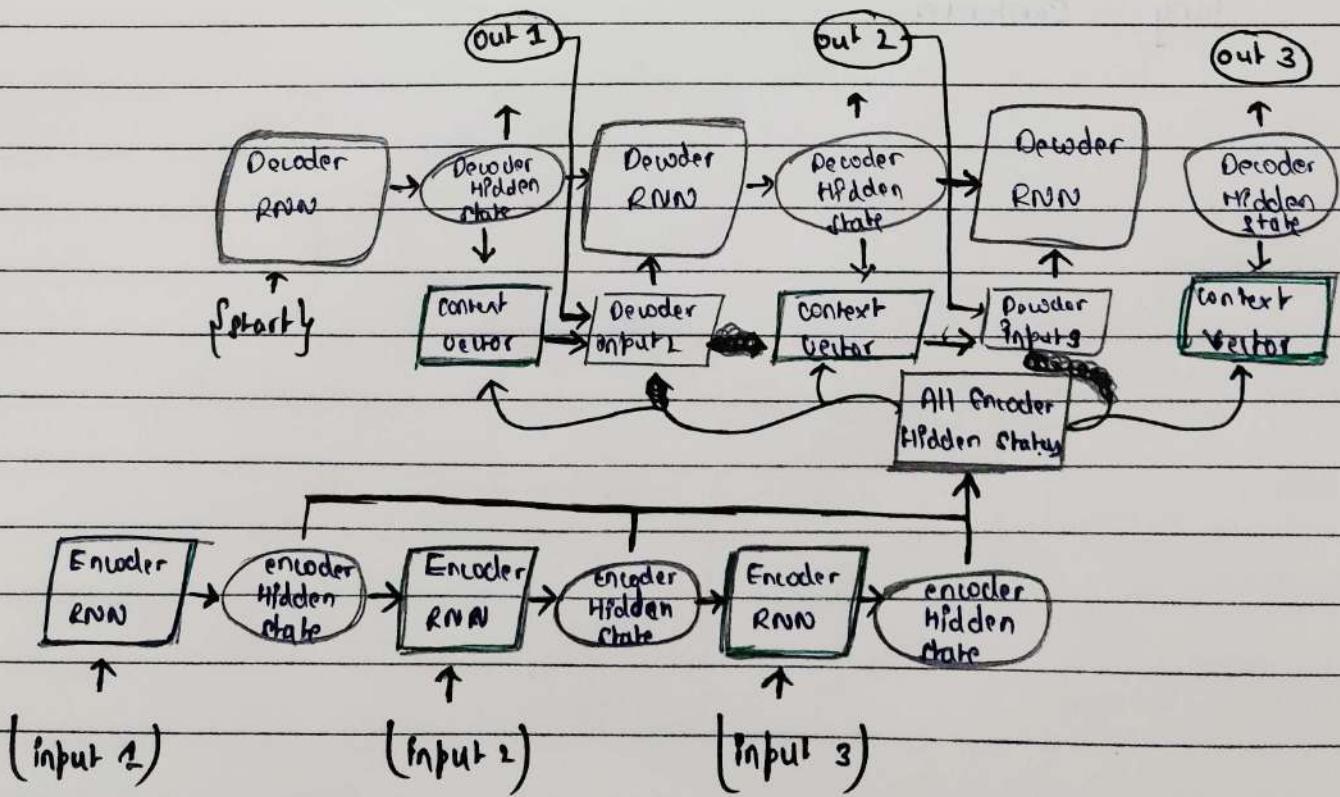
By Jointly learning to Align and translate.



Attention mech^o

In transformer there is no one single state vector

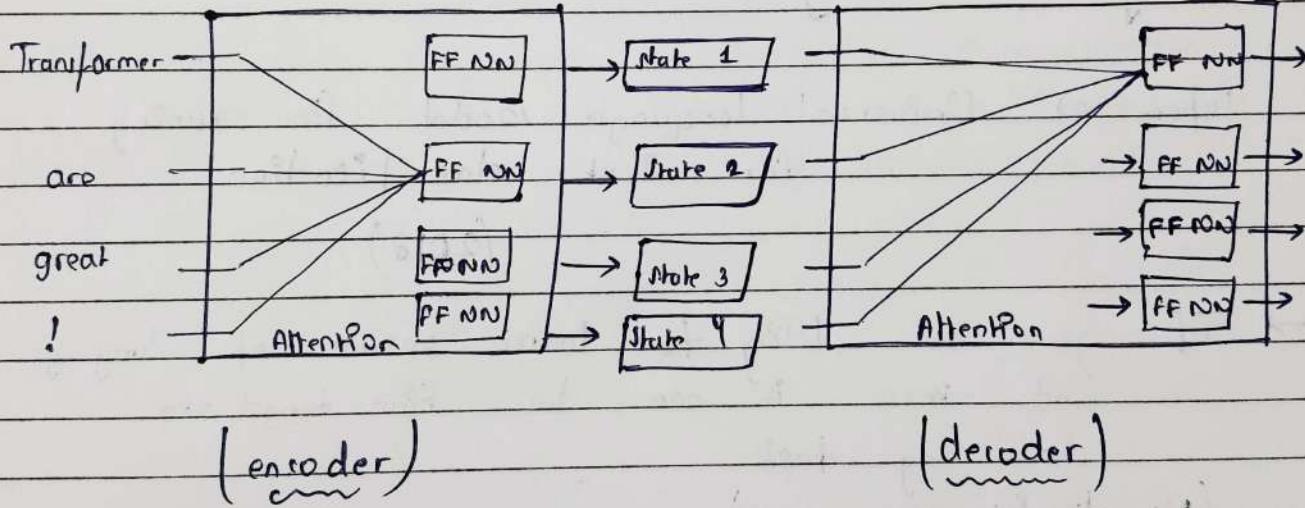
at every step of predicting output it has
 notion of every step of input to perform
 better than encoder decoder



Problem :- Computational complexity is very high.
(using LSTM's are the problem)

* Transformers.

Google released a paper. (2017) : Attention is All You need.



There are no (LSTM's / RNN cells.) X

↓
(Attention) ✓
(self Attention) ✓

→ It has ability to see all the words of input at every ~~at~~ ^{every} state of predicting of words.
- which is called parallel processing.

- Reduce computational complexity.

Problem →

- 1) Hardware → GPU → costly.
- 2) Time → still not that less
for training from scratch
- 3) Data → needs a lot of data

Transfer learning :-

Paper ⇒ Universal language Model fine tuning
for Text classification
(2018)

→ T.F gave the ability to learn basic of language
and men if can be fine tuned to
do any task
(No transformers used)

LLM's (2018)

Google and open AI released language model
Based On transformer architecture.

Google (BERT)

[Encoder only]

open AI (chat gpt)

[decoder only]

→ Both trained on huge dataset.

→ Tasks that they were
able to do. ⇒ (1) text summarization
(2) question / Answering
(3) sentiment Analysis