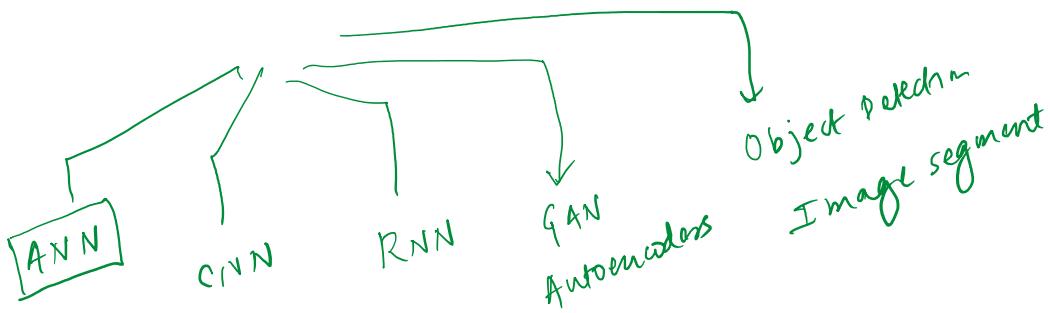


# Curriculum

Tuesday, February 15, 2022 3:30 PM



## Features

Tuesday, February 15, 2022 3:31 PM

- 1) Well researched
  - 2) Easy to consume
  - 3) Well structured
  - 4) Tensorflow + Keras
  - 5) Projects
- 100 days

## Prerequisites

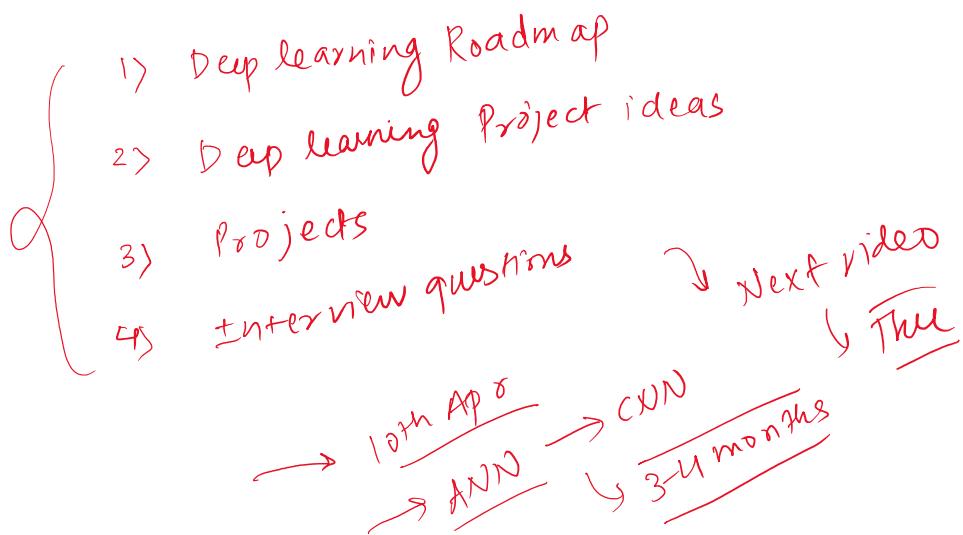
Tuesday, February 15, 2022 3:31 PM

- 1> Python →
- 2> Basics of ML
- 3> Linear Algebra

3 Blue / Brown  
Linear Alg → ⑤

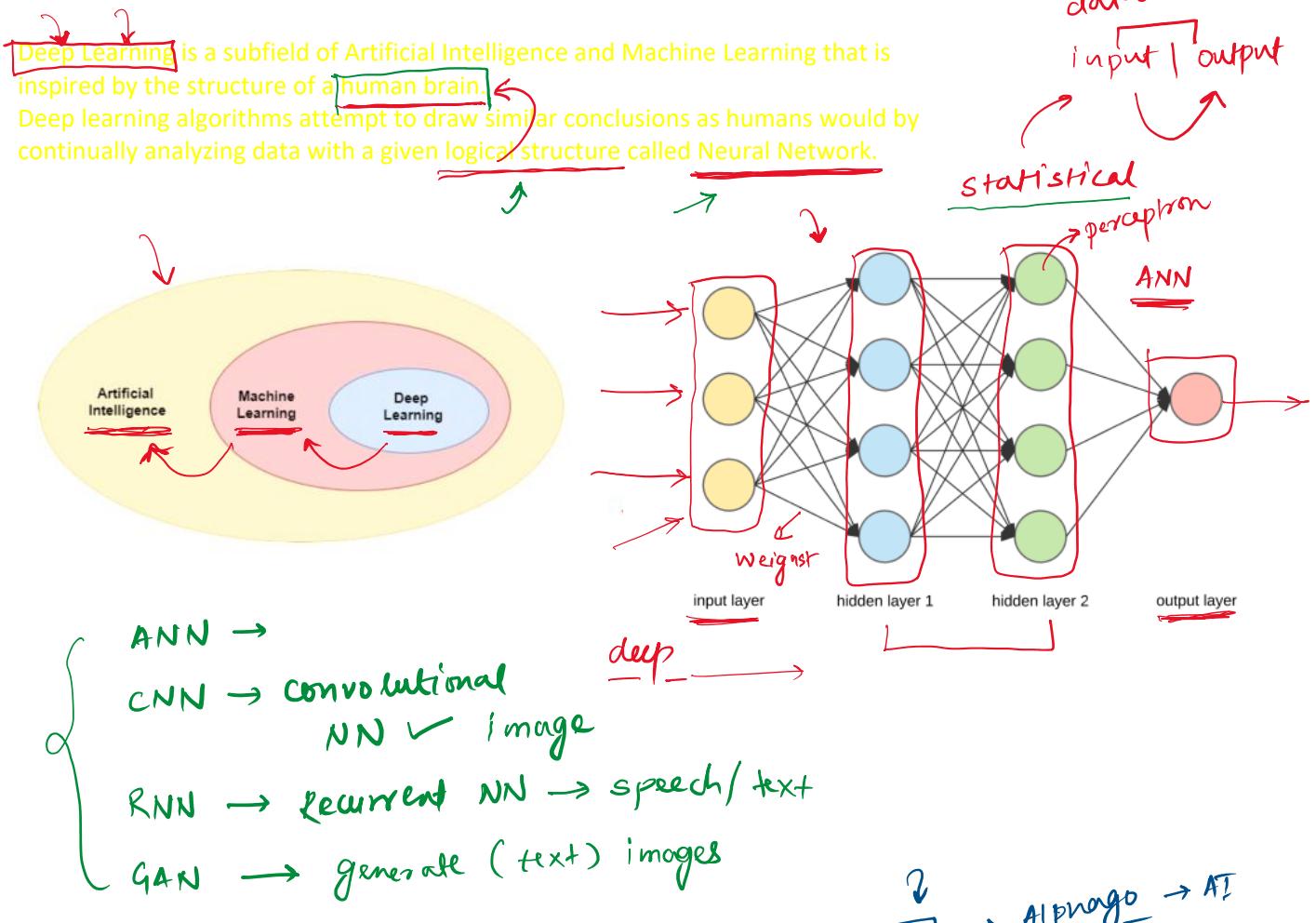
## Extra Content

Tuesday, February 15, 2022 3:32 PM



# What is Deep Learning?

16 February 2022 06:32



## why dl getting so famous?

- 1) Applicability →
- 2) Performance →

state-of-the-art

human experts

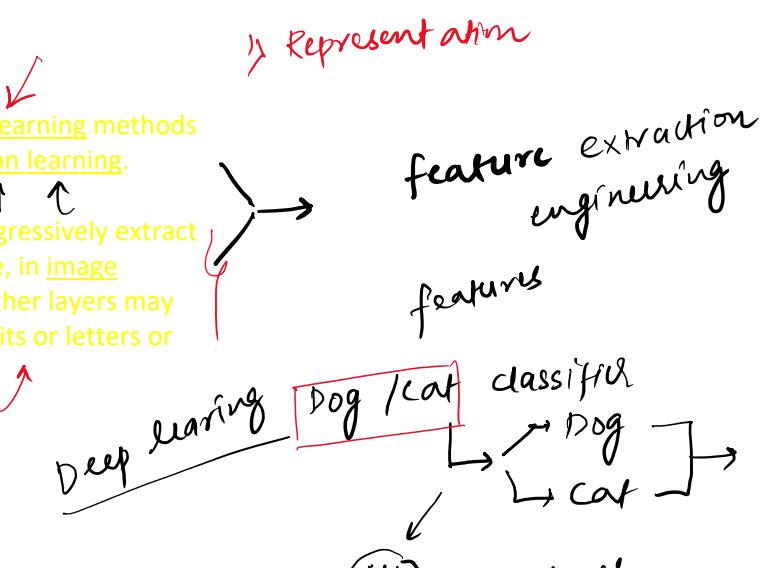
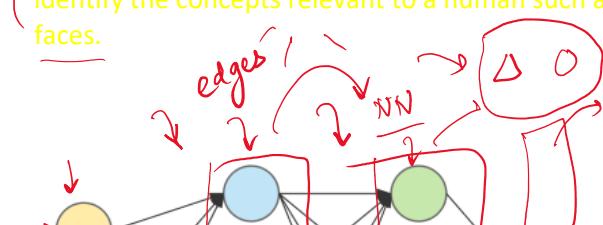
Applications

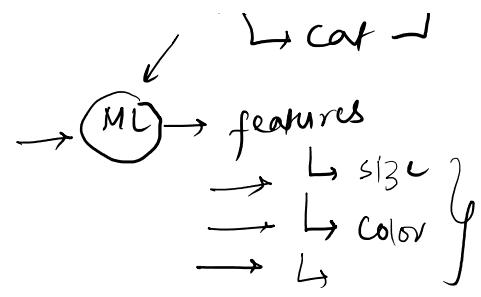
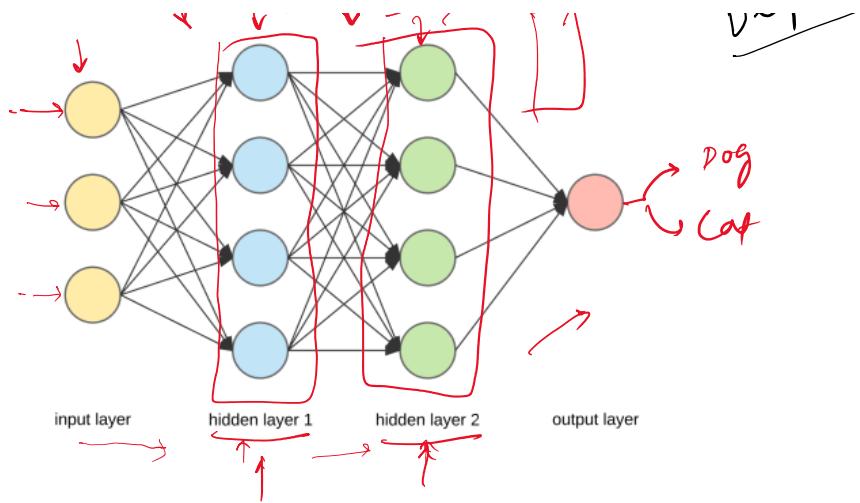
Go → AlphaGo → AI  
5 → 4

self-driving  
drugs  
→ chemistry

Deep learning is part of a broader family of machine learning methods based on artificial neural networks with representation learning.

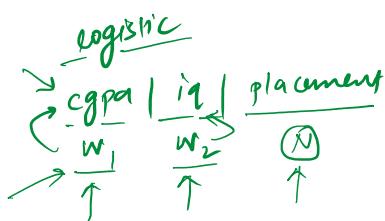
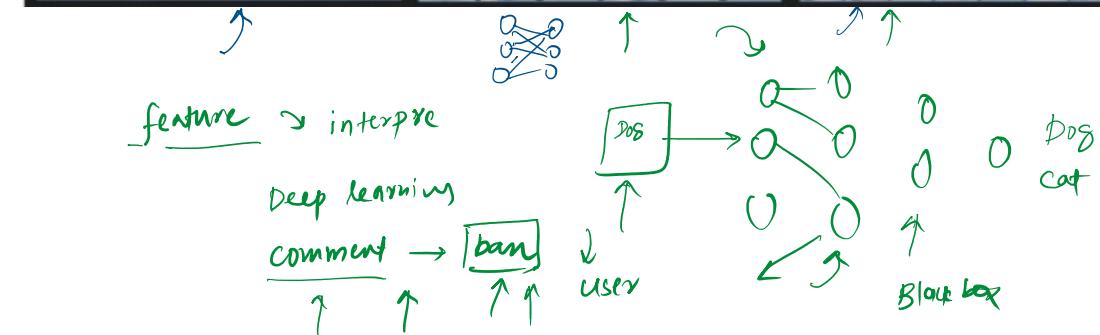
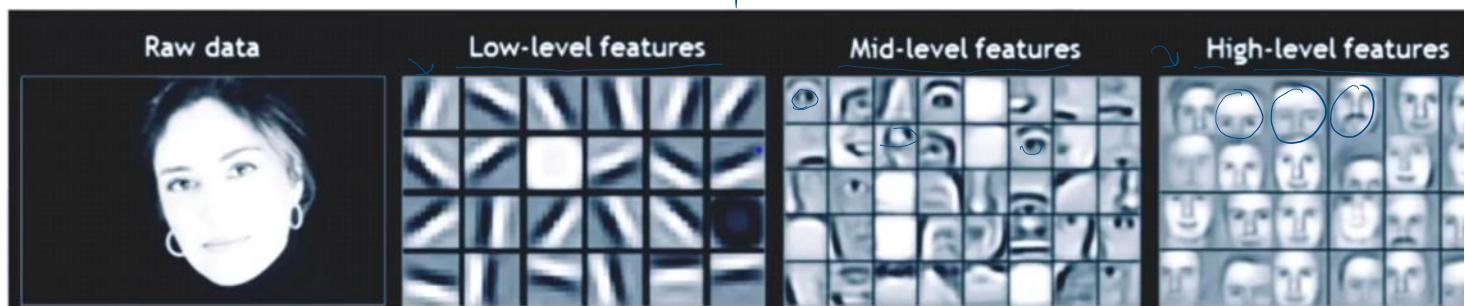
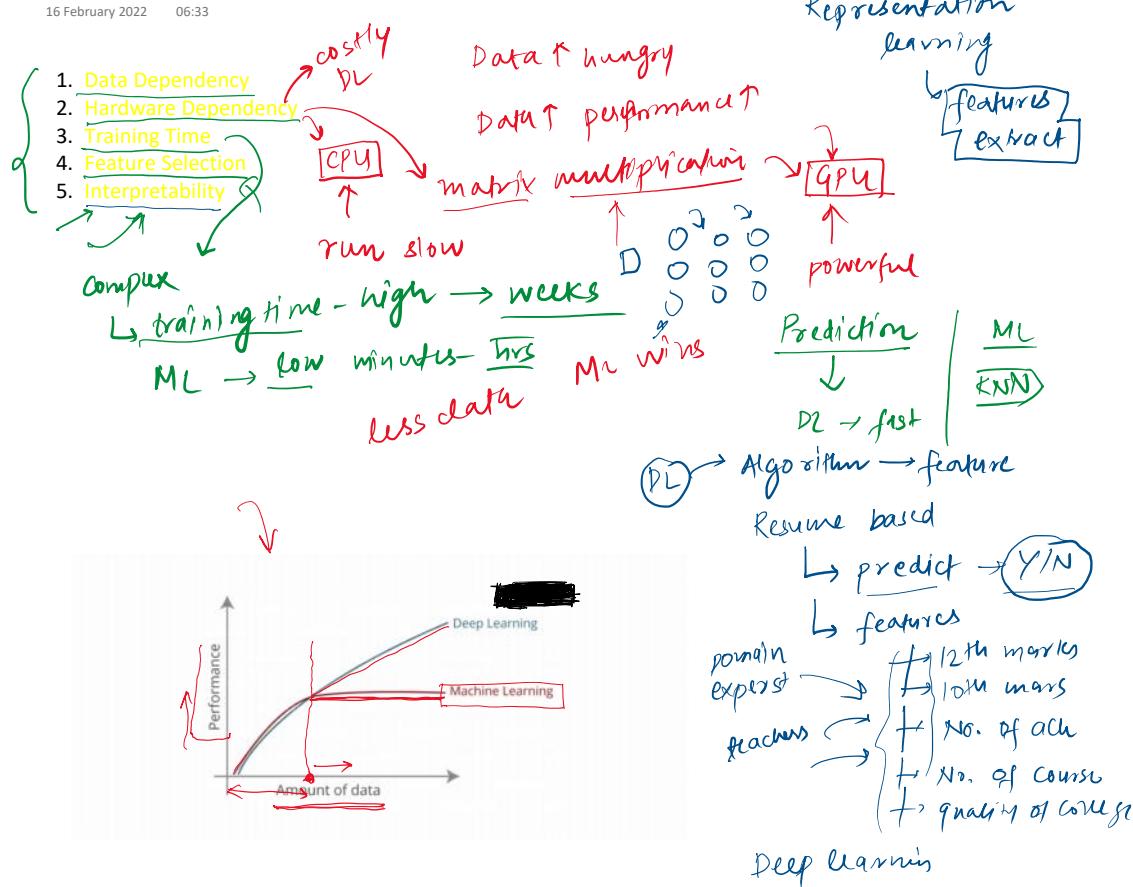
Deep Learning Algorithms uses multiple layers to progressively extract higher-level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits or letters or faces.



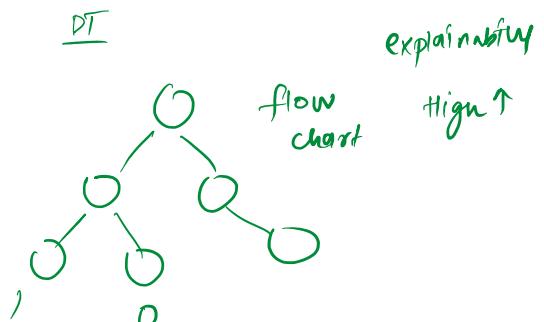


## Deep Learning VS Machine Learning

16 February 2022 06:33

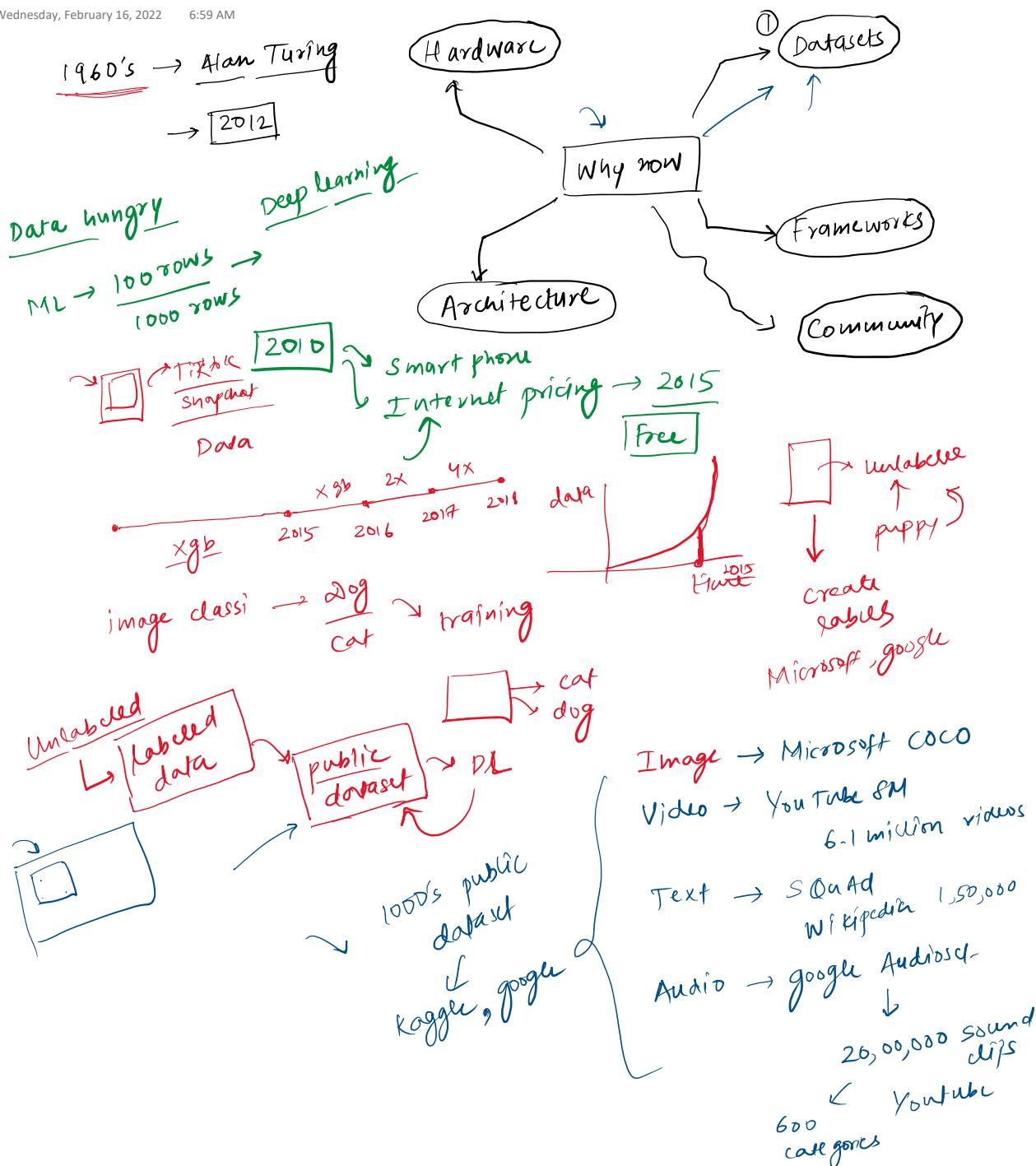


Machine learning →

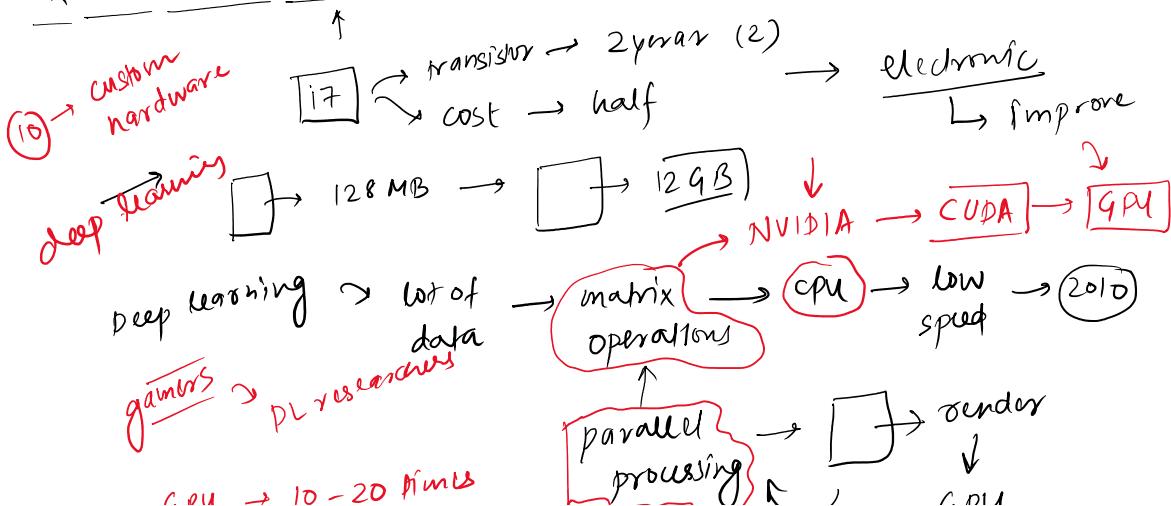


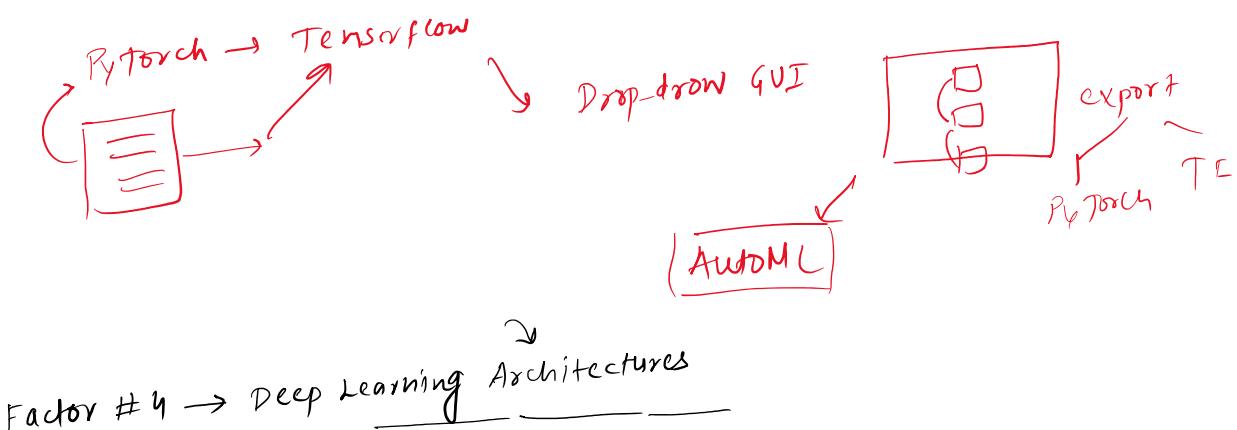
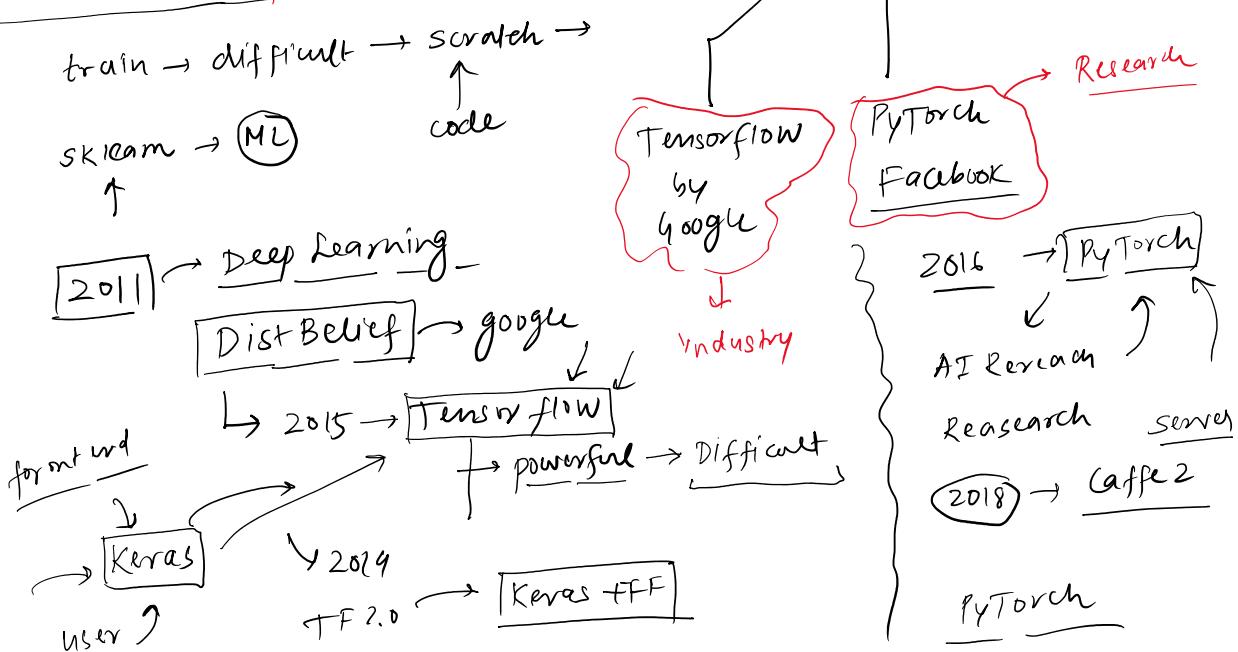
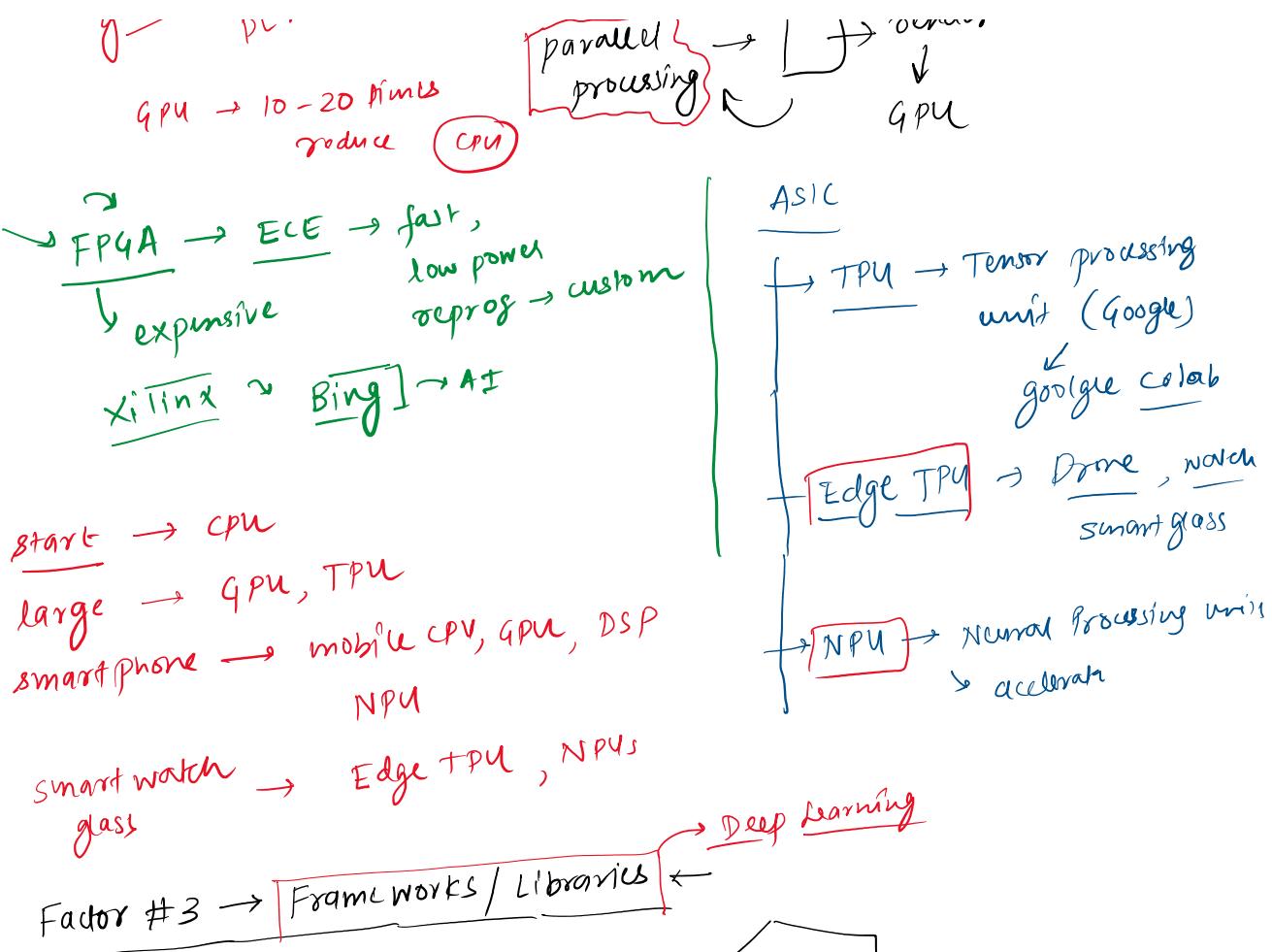
## Why now?

Wednesday, February 16, 2022 6:59 AM

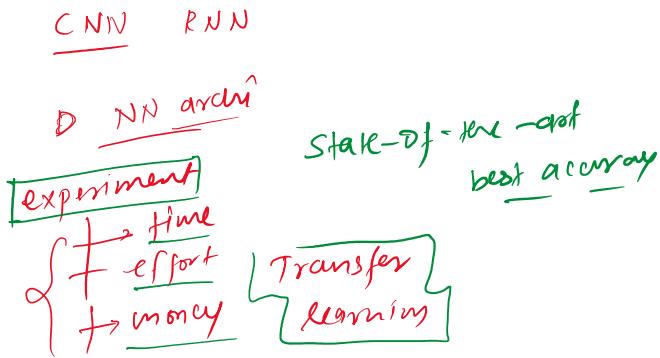
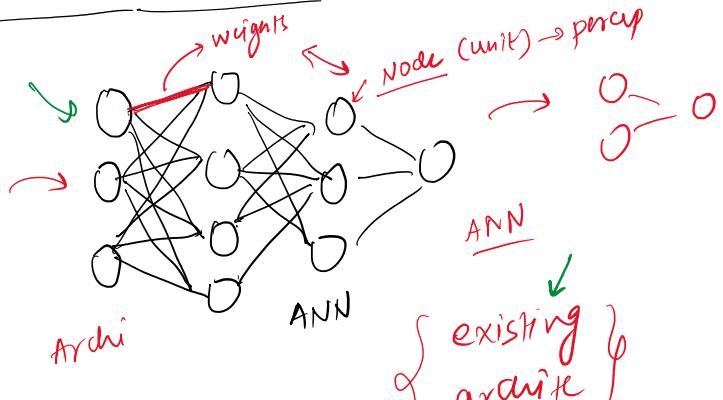


## Reason #2 → Hardware





## Factor #4 → Deep Learning Architectures



NN → Arch  
↓  
fires down  
Dataset → train → ready to use →

image classification → ResNET  
Text classifi → BERT  
Image segment → UNet  
image trans → Pix2Pix

Object det → YOLO  
speech gen → WaveNET

## #5 → People / community

Deep learning → 1960 → 2011 →

Research

↳ DL origin

↳ Teacher

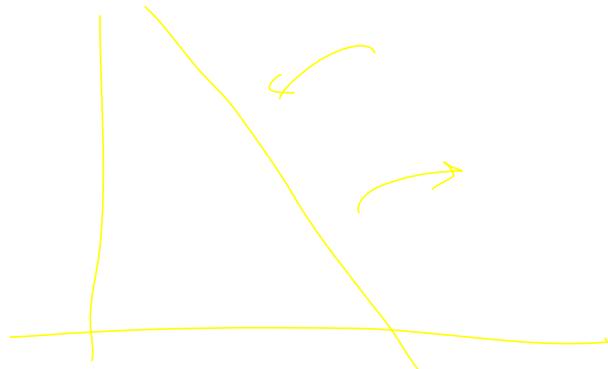
→ student → tagger

(PL) →

## Code Example

Saturday, February 19, 2022 7:10 AM

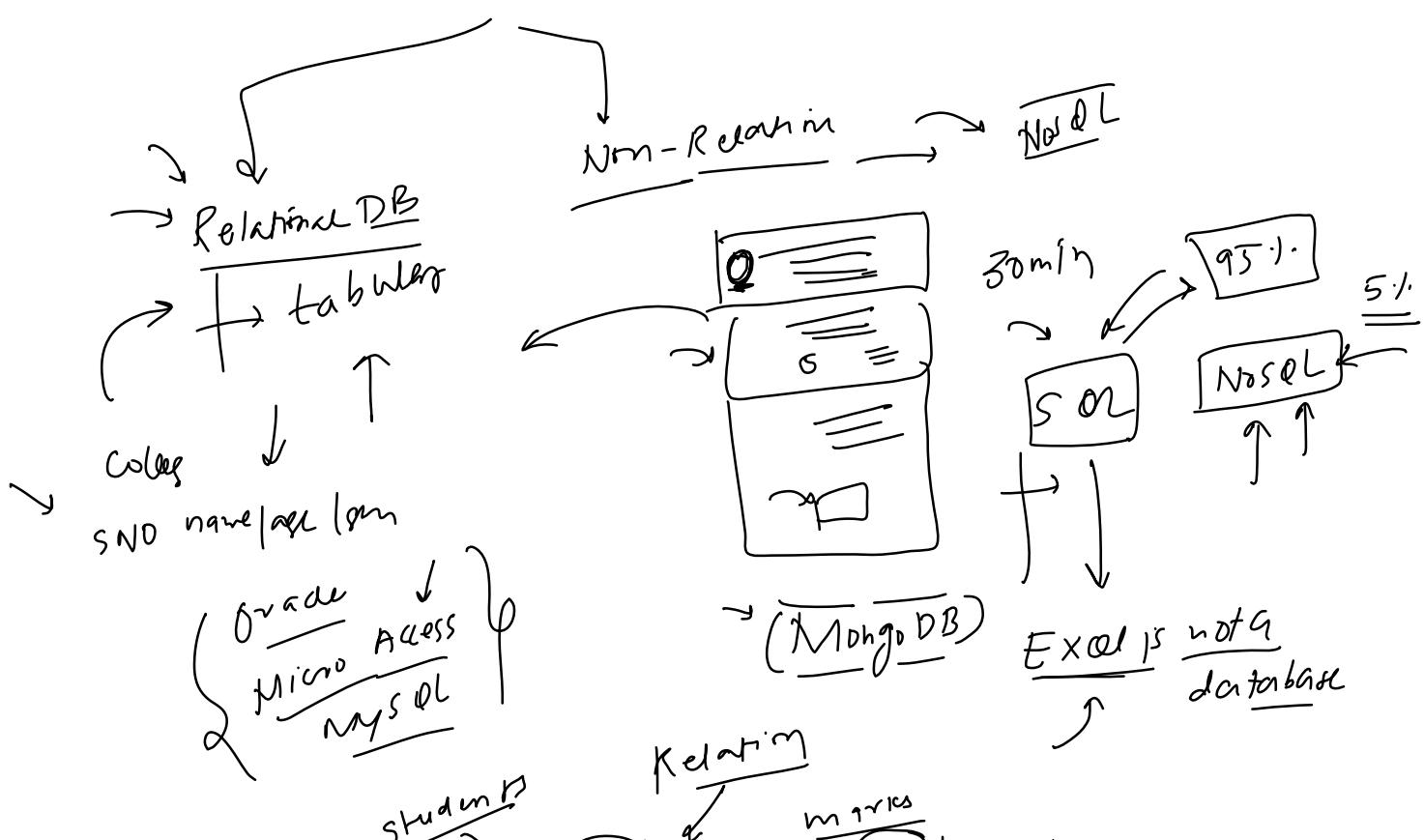
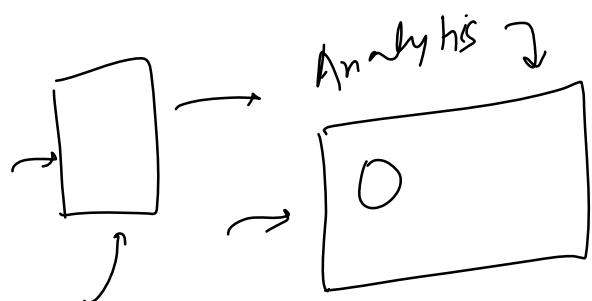
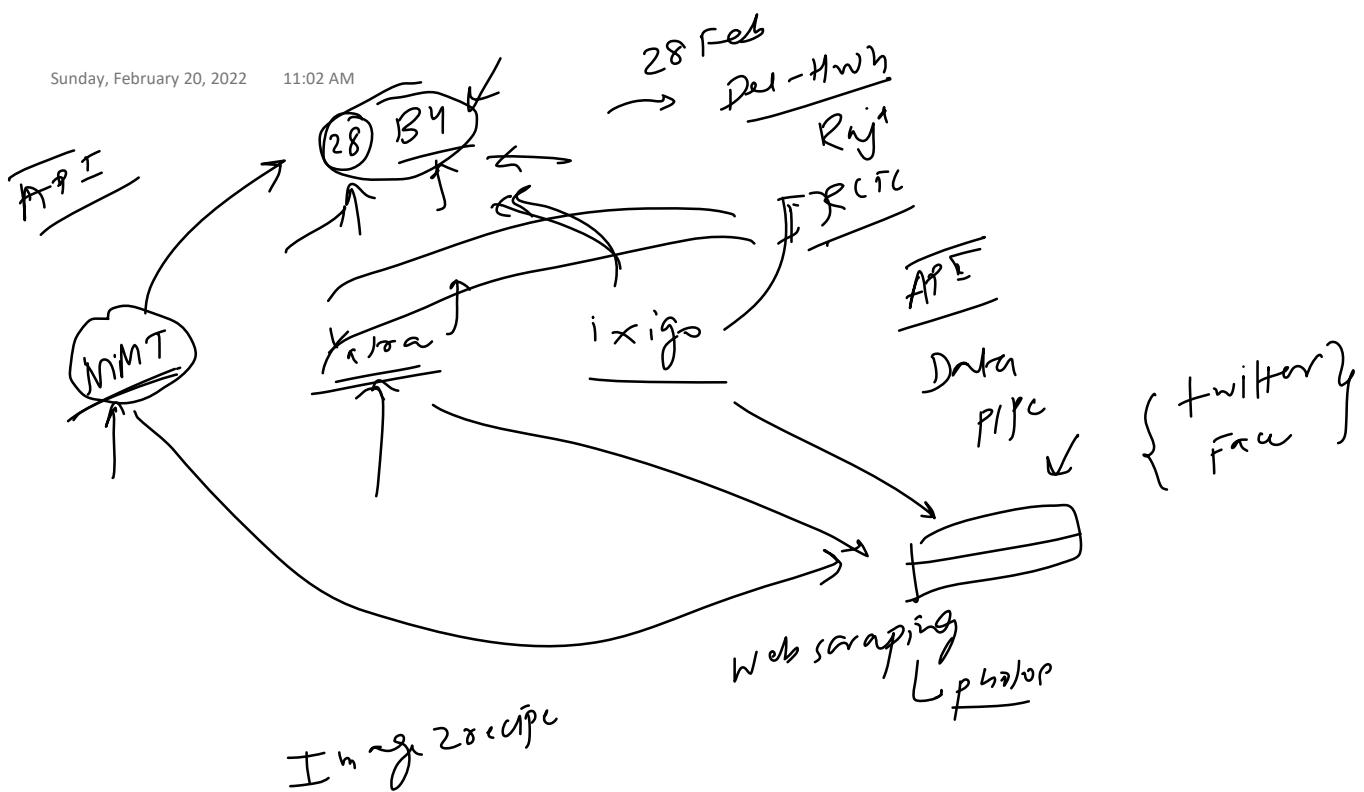
cgpa | resumme / placed

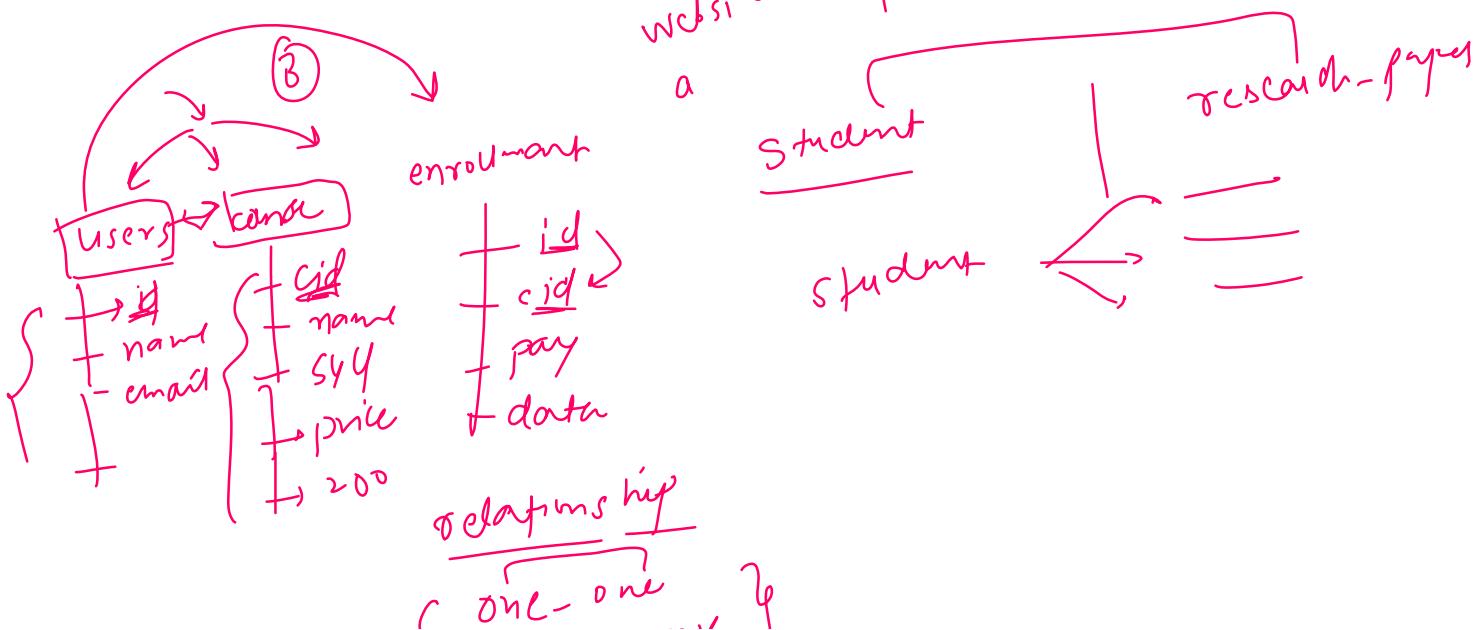
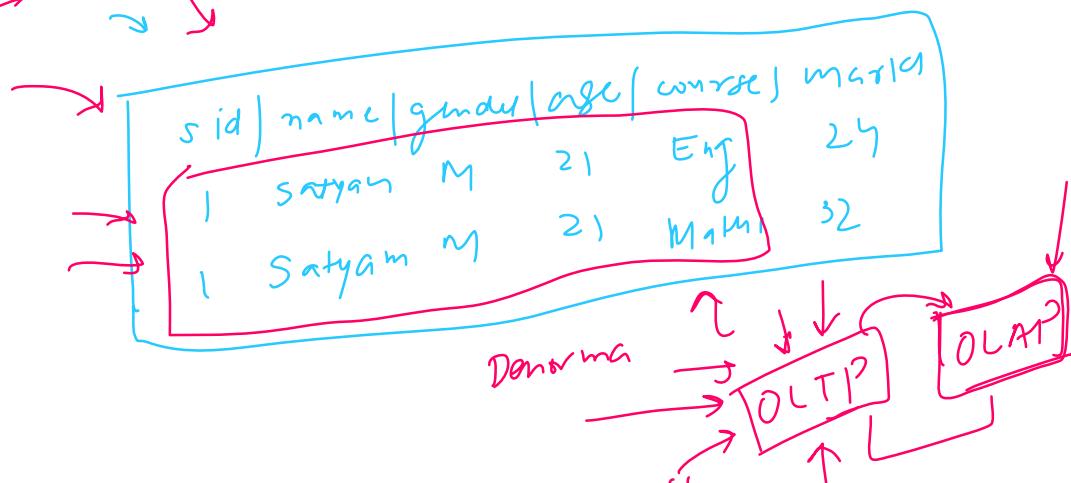
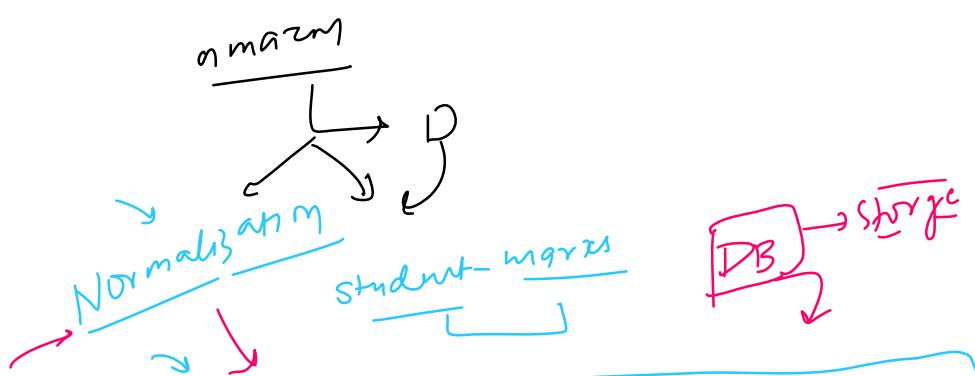
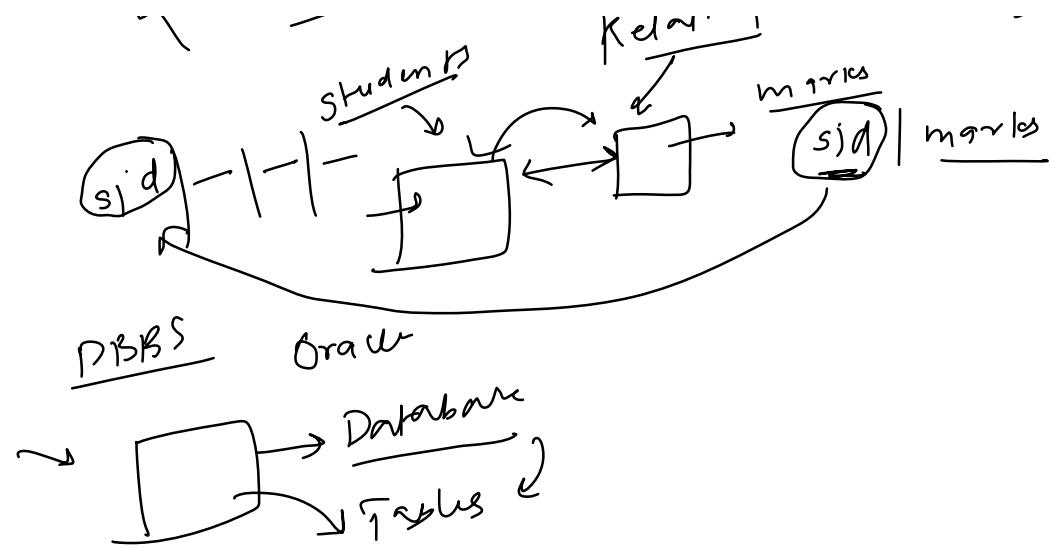


$$w_1 \cdot 31 + w_2 \cdot 1 + b = -25$$



cgpa



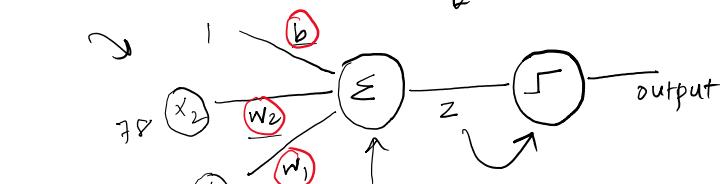


$\rightarrow \{$  one-one  
one-many }  
many-

## Recap

22 February 2022 14:12

$$z = w_1x_1 + w_2x_2 + b$$



cgpa	iq	placed
x1	x2	
7	78	1
6	61	0
9	92	1

↳ mathematical model  $\rightarrow$  neuron

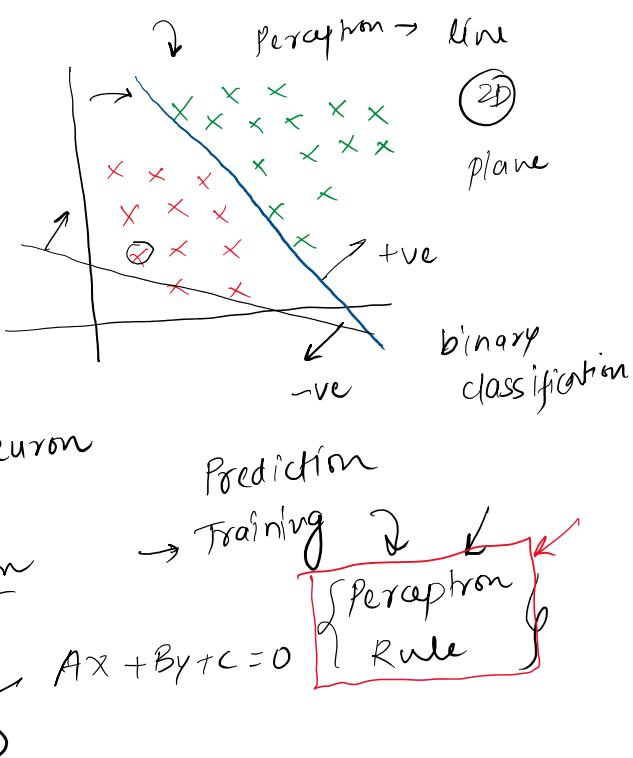
$$z \geq 0 \rightarrow 1$$

$$z < 0 \rightarrow 0$$

100% random

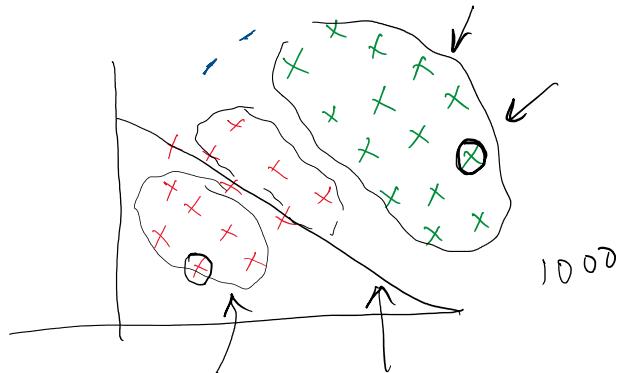
$$A \rightarrow w_1$$

$$B \rightarrow w_2 \quad C \rightarrow b$$



## Problem with Perceptron Trick

22 February 2022 16:10



$$\text{line} \rightarrow \begin{bmatrix} w_1 & w_2 & b \end{bmatrix}$$

randomly per

converge

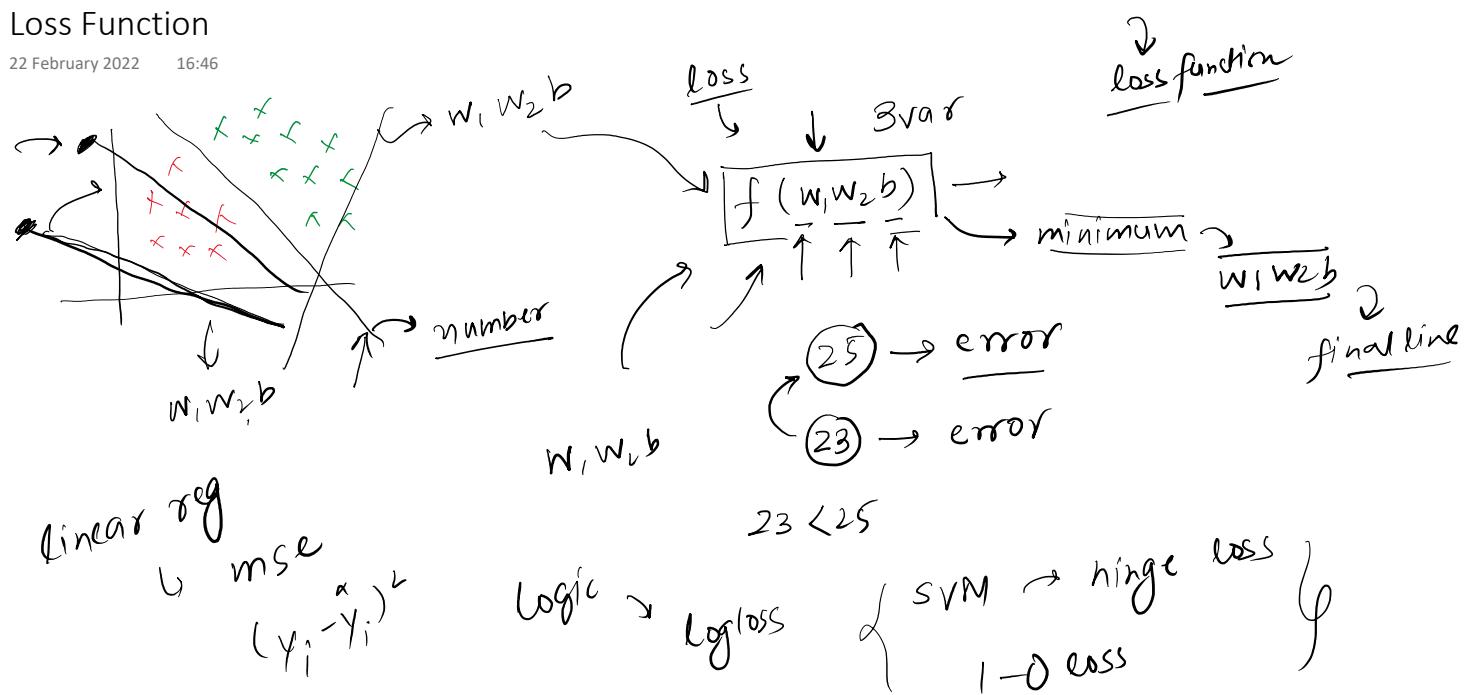
quantify

$\begin{bmatrix} w_1 & w_2 & b \end{bmatrix} \rightarrow \text{loss functions}$

1000

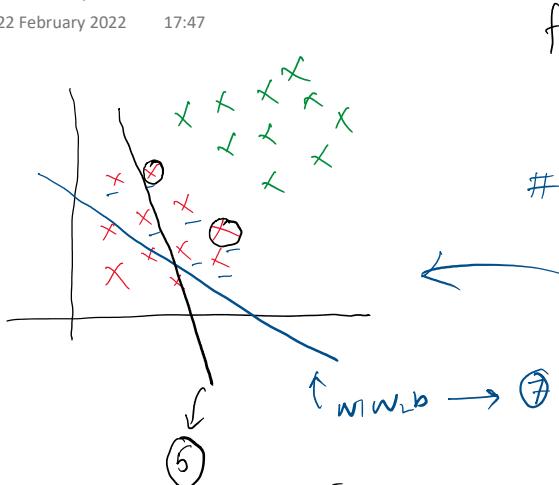
## Loss Function

22 February 2022 16:46



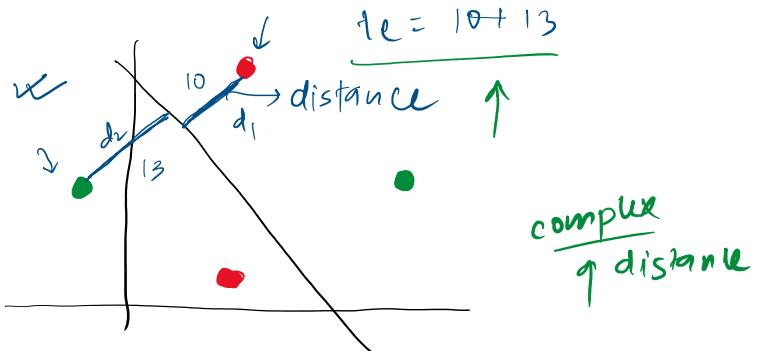
# Perceptron Loss Function

22 February 2022 17:47



$f(\underline{w_1} \underline{w_2} b) \downarrow$   
number  $\rightarrow$  error

# misclassified points



$$\boxed{30+6}$$

→ Skream

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

Regularization

$$L(y_i, f(x_i)) = \boxed{\max(0, -y_i f(x_i))}$$

$$f(x_i) = \frac{w_1 x_1 + w_2 x_2 + b}{z =}$$

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

$L(w_1, w_2, b)$

$n$ : # rows in data

$(w_1, w_2, b)$

$$L = \underset{w_1, w_2, b}{\operatorname{argmin}}$$

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

minimum

## gradient descent

$$L = \underset{w_1, w_2, b}{\text{argmin}} \sum_{j=1}^n$$

gradient descent

## Explanation of Loss Function

23 February 2022 08:07

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

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

$x_1$	$x_2$	$y$
$x_{11}$	$x_{12}$	$y_1$
$x_{21}$	$x_{22}$	$y_2$
$\vdots$	$\vdots$	$\ddots$ points

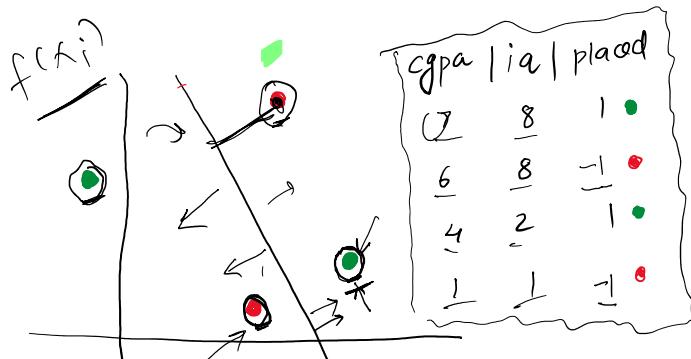
n rows  
n columns  
 $x_{ij}$  → col

$$\max(0, -y_i f(x_i)) = \begin{cases} x &= 0 \\ x > 0 \\ x < 0 \end{cases}$$

$f(x_2) = w_1 x_{21} + w_2 x_{22} + b$

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

2 points n points distance



$$\begin{array}{cc} y_i & \hat{y}_i \\ \rightarrow & \rightarrow \\ \begin{bmatrix} 1 & 1 \\ -1 & -1 \\ 1 & -1 \\ -1 & 1 \end{bmatrix} & \end{array}$$

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

$$\begin{array}{ll} \underline{f(x_i)} = +ve & \\ +ve & +ve = -ve \end{array}$$

$$\max(0, -ve) = 0$$

$$\begin{array}{lll} y_i & f(x_i) & \\ +ve & +ve & \\ -ve & -ve & \\ & -ve & -ve \end{array}$$

$$y_i \quad f(x_i)$$

$$+ve \quad -ve = -ve = +ve$$

$$\max(0, +ve) = +ve$$

$$\begin{array}{ll} y_i & f(x_i) \\ -ve & +ve \rightarrow -ve -1 = +ve \end{array}$$

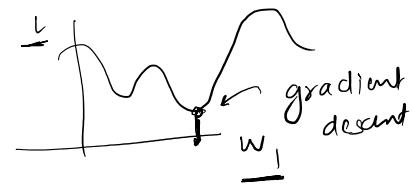
## Gradient Descent

23 February 2022 08:50

$$L = \underset{w_1, w_2, b}{\operatorname{argmin}} \quad \boxed{\frac{1}{n} \sum_{j=1}^n \max(0, -y_j f(x_j))}$$

$$L(w_1, w_2, b)$$

$$L(w_1)$$



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

100 / 1000

$$\boxed{\eta = 0.01}$$

$w_1, w_2, b = 1$

for i in epochs:

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

$$\left[ \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b} \right]$$

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

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

# Loss Function Differentiation

23 February 2022 12:57

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

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

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

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

$\frac{\partial f(x_i)}{\partial w_1} = x_{i1}$

$$\boxed{\frac{\partial L}{\partial w_1} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i1} & \text{if } y_i f(x_i) < 0 \end{cases}}$$

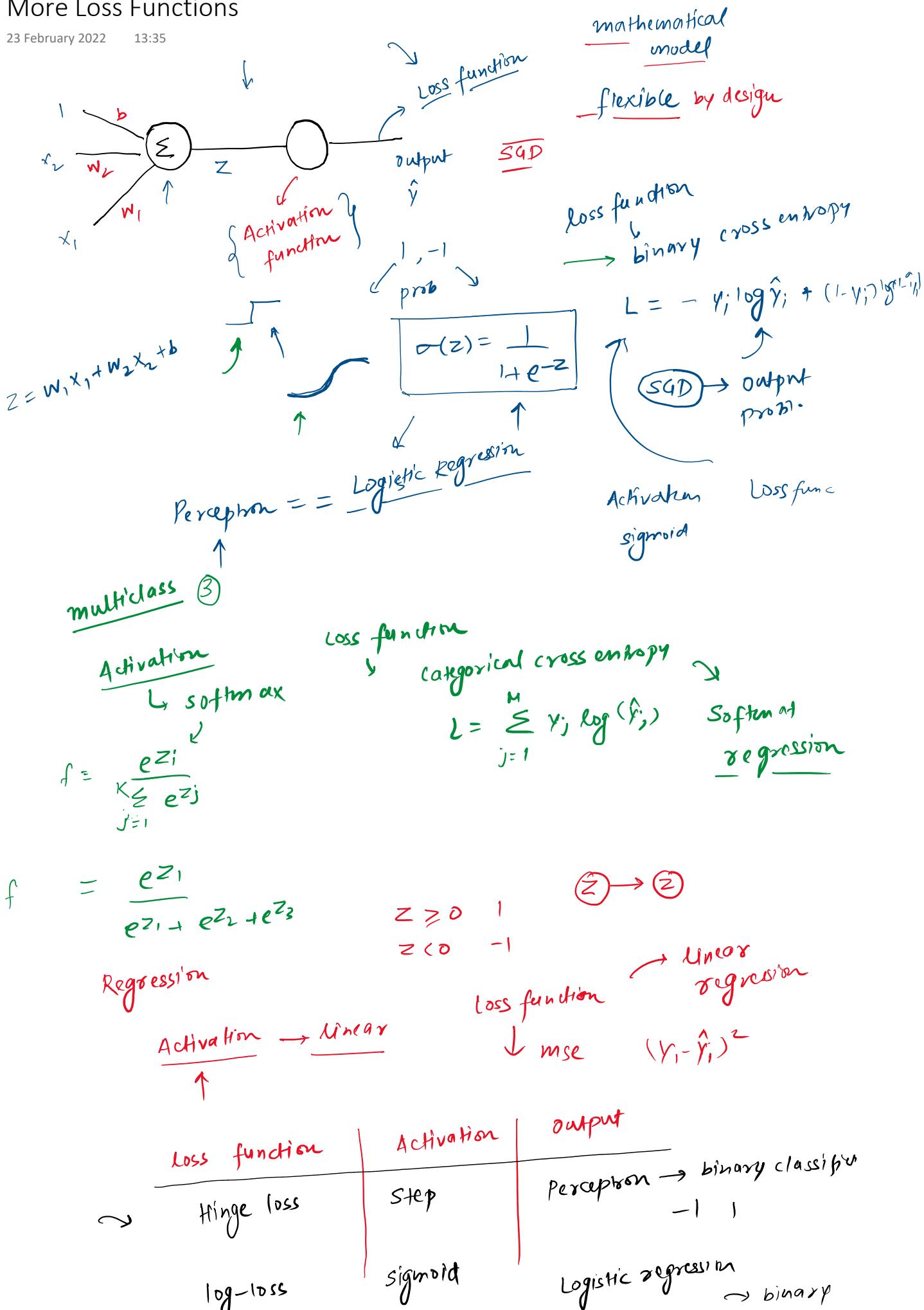
$$\boxed{\frac{\partial L}{\partial w_2} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i2} & \text{if } y_i f(x_i) < 0 \end{cases}}$$

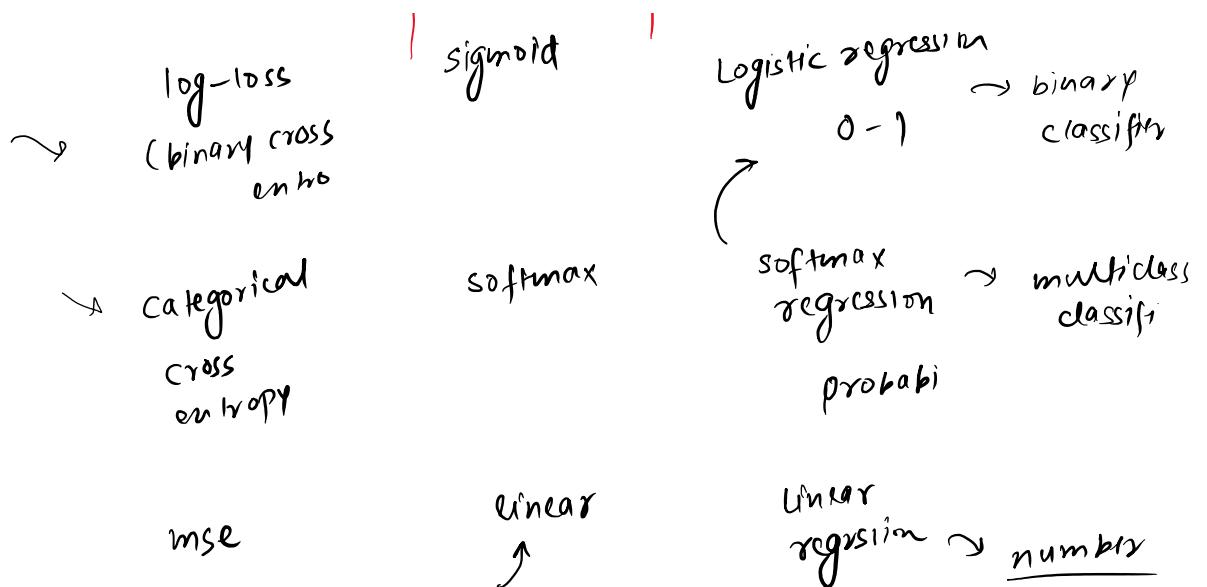
→ code using  
gradient descent

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

# More Loss Functions

23 February 2022 13:35

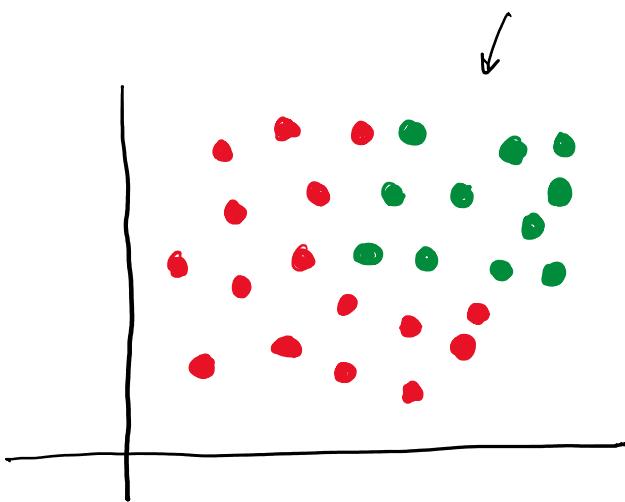




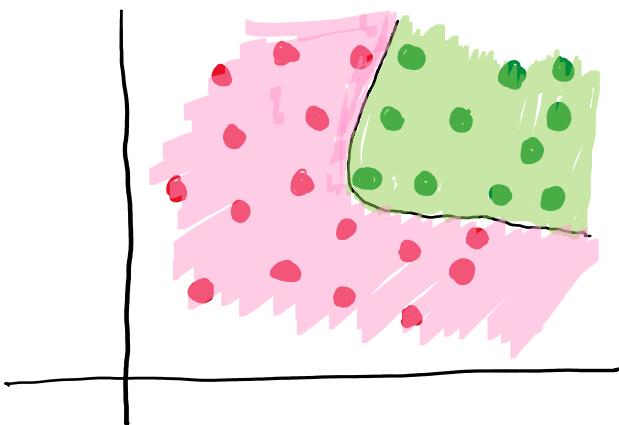
# The Problem

24 February 2022 13:41

perceptron

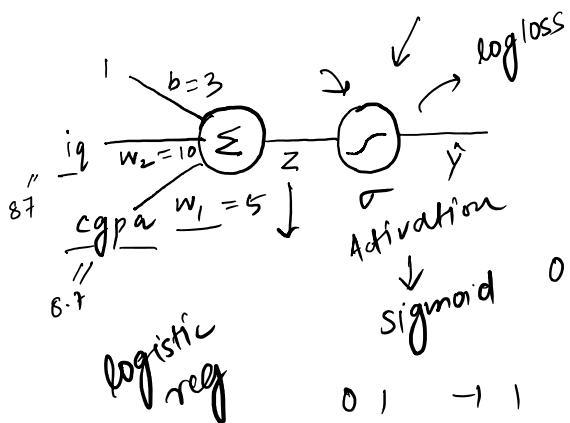


line  
↓  
curve



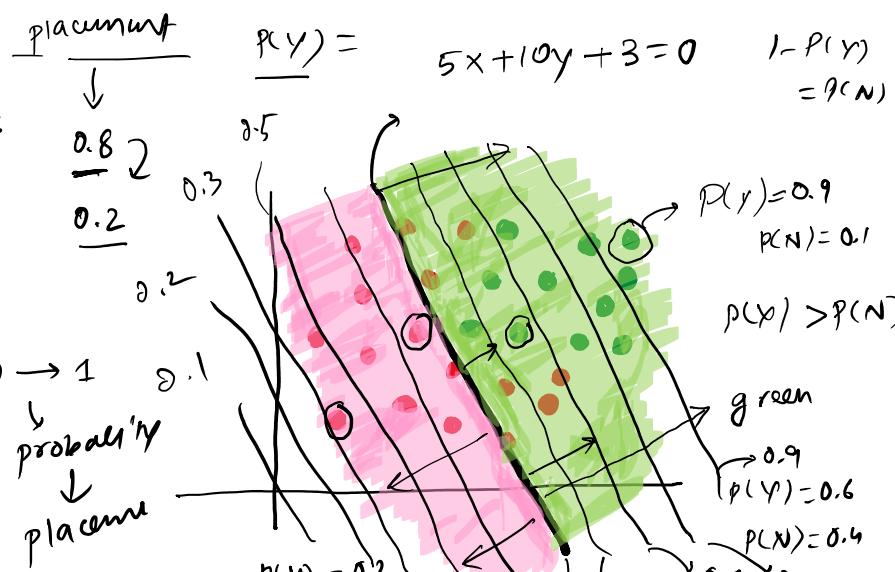
## Perceptron with Sigmoid

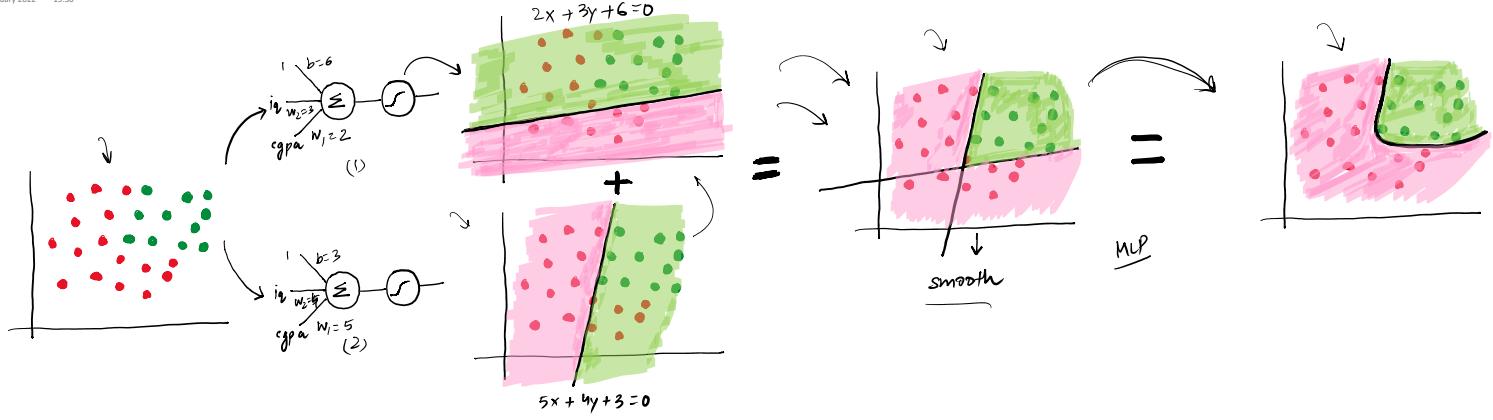
24 February 2022 15:35

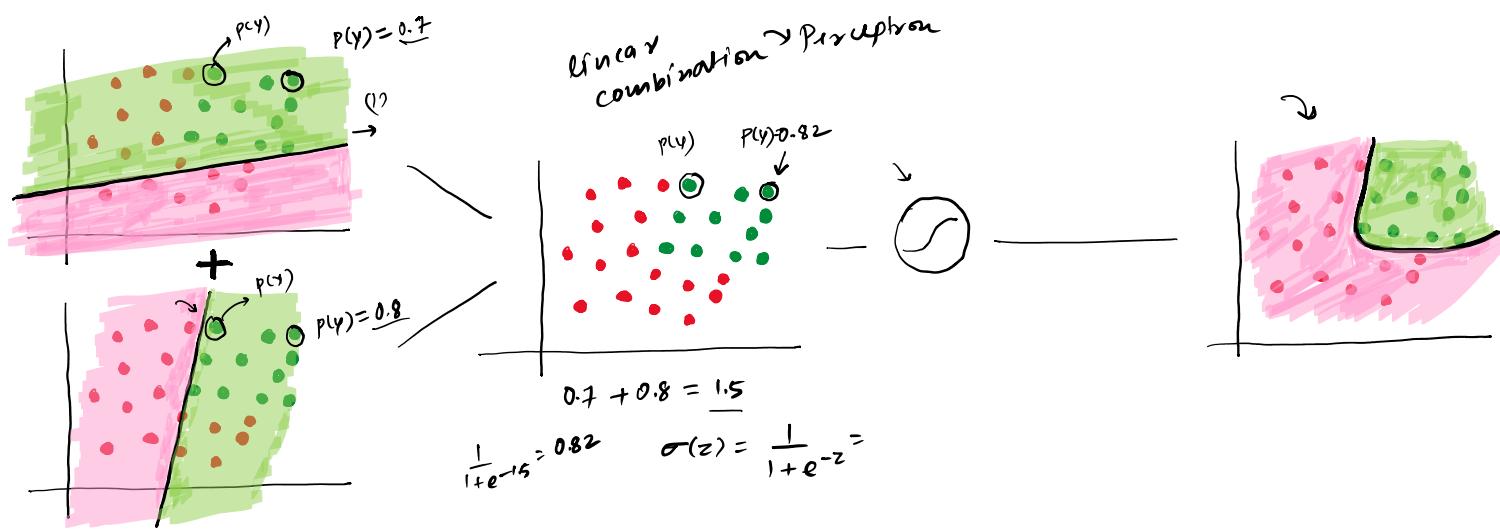


$$\begin{aligned} w_1 \text{ cgp } &+ w_2 \text{ iq } + b \\ \downarrow & \\ 5 \times 8.7 + 10 \times 0.7 + 3 &= z \rightarrow \frac{1}{1+e^{-z}} = 0 \rightarrow 1 \end{aligned}$$

$P(y) = 0.8 \quad P(N) = 0.2$

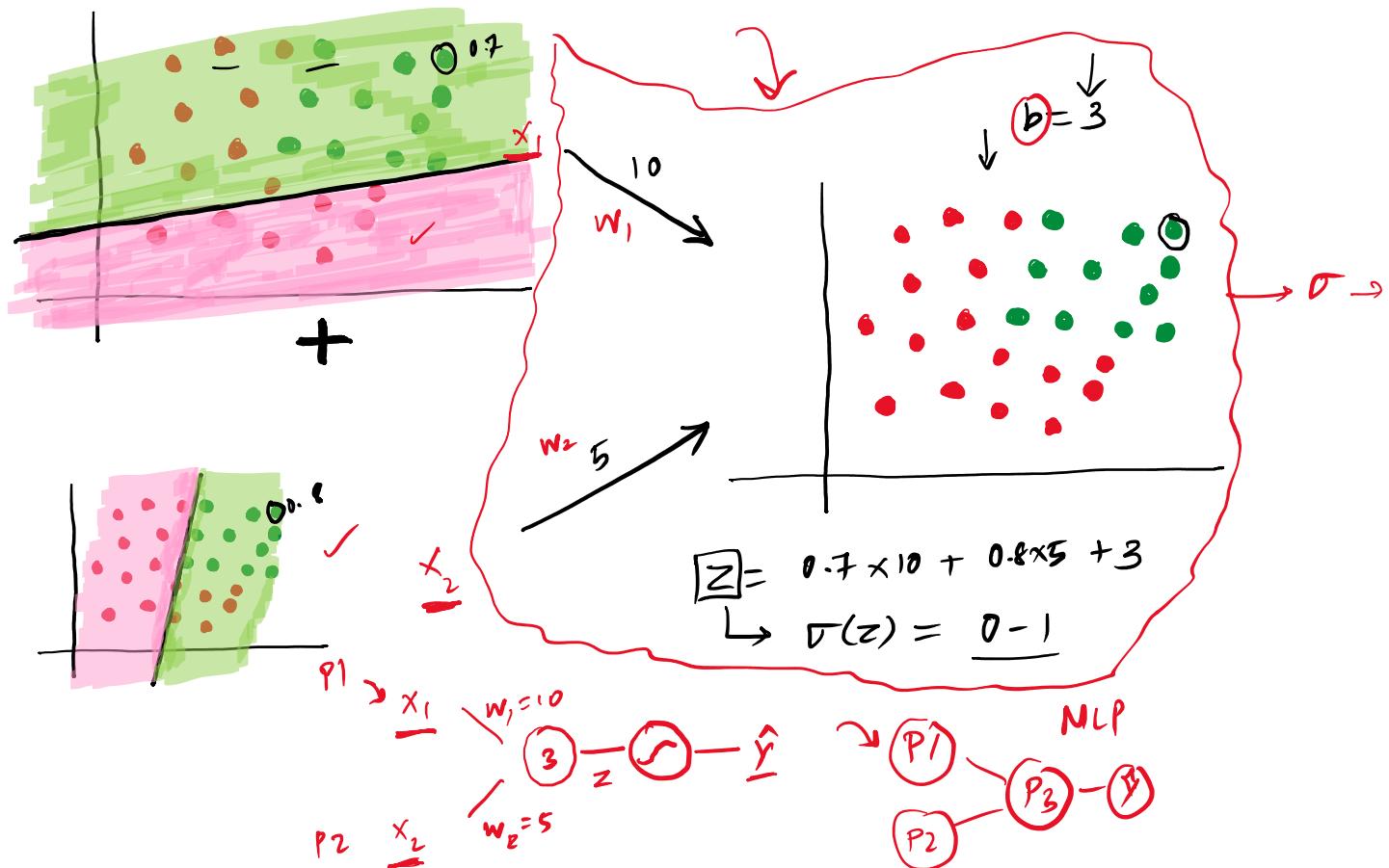
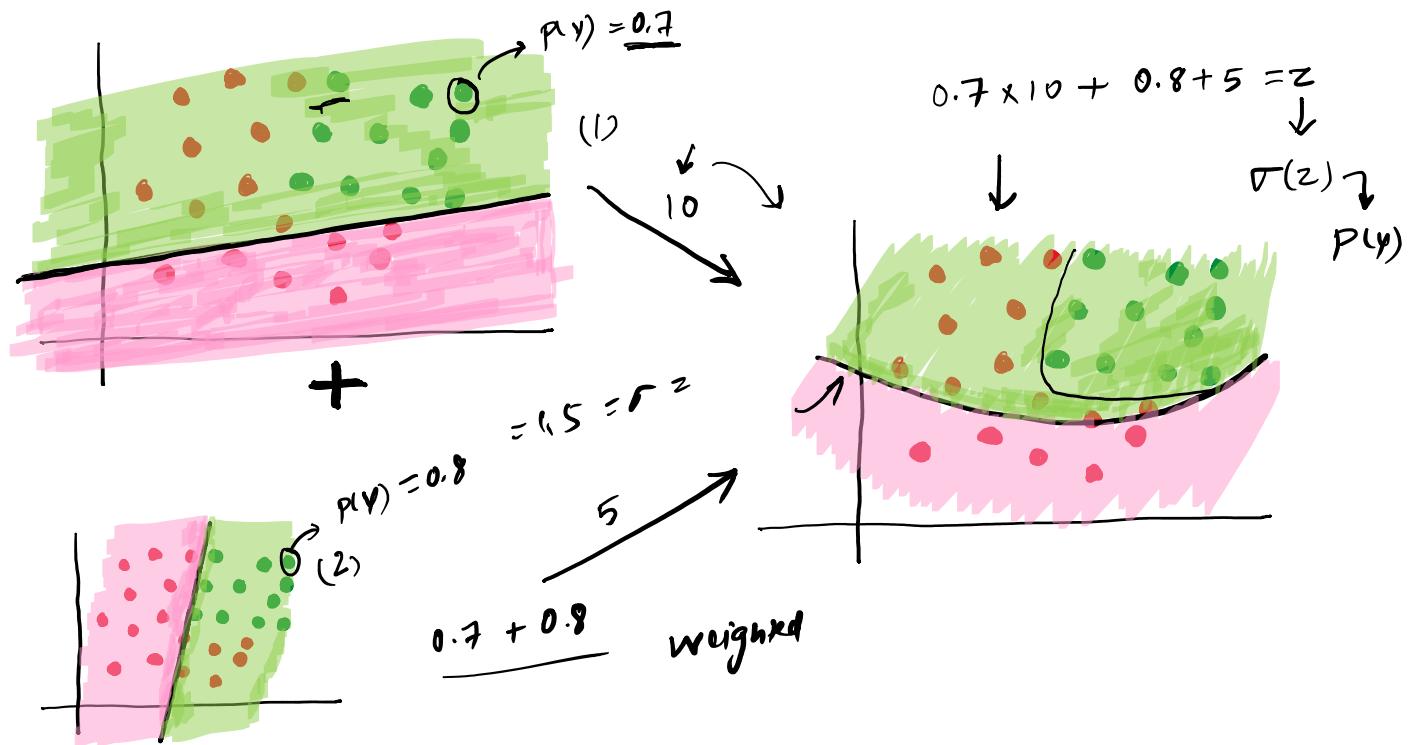






## Adding Weights

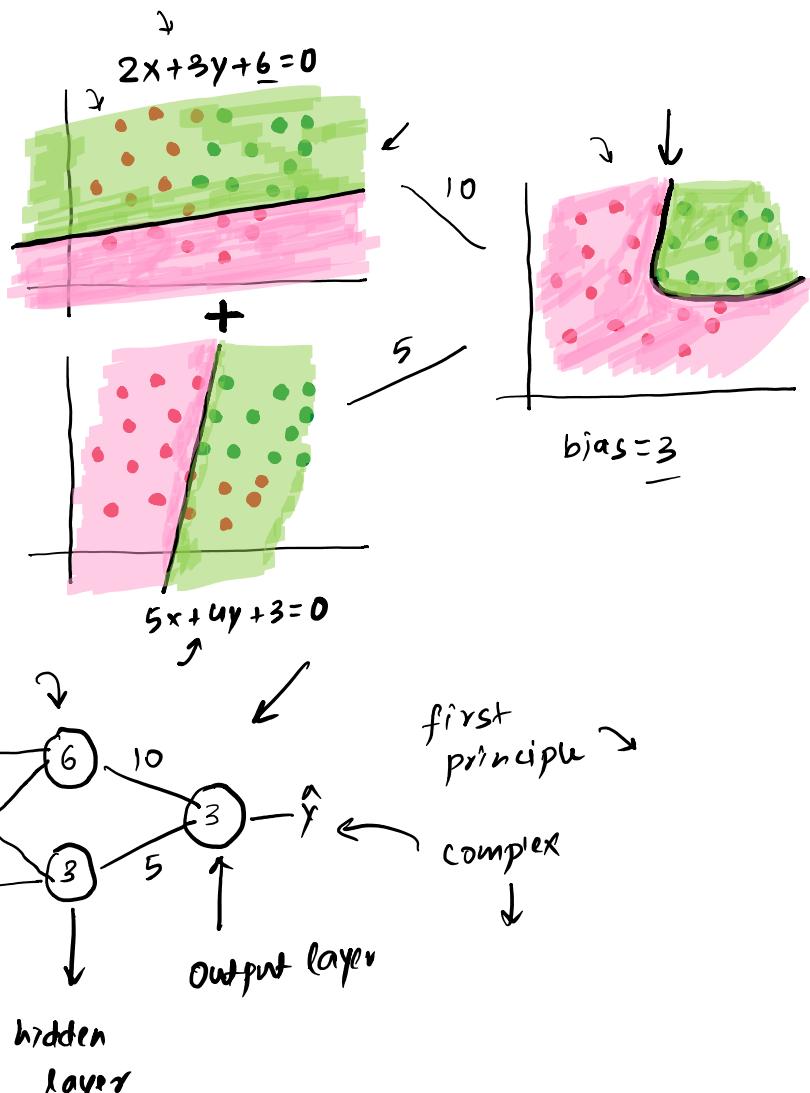
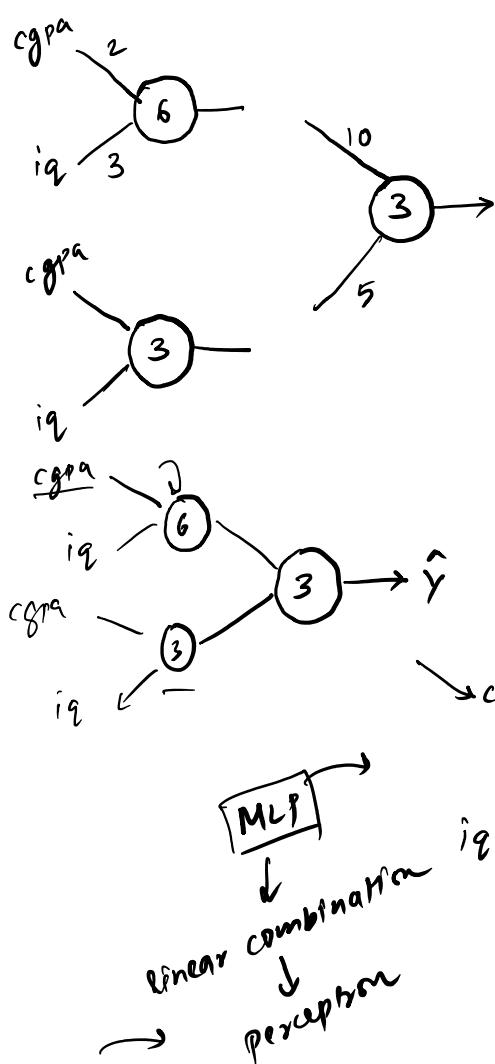
24 February 2022 15:37



# The Idea of MLP

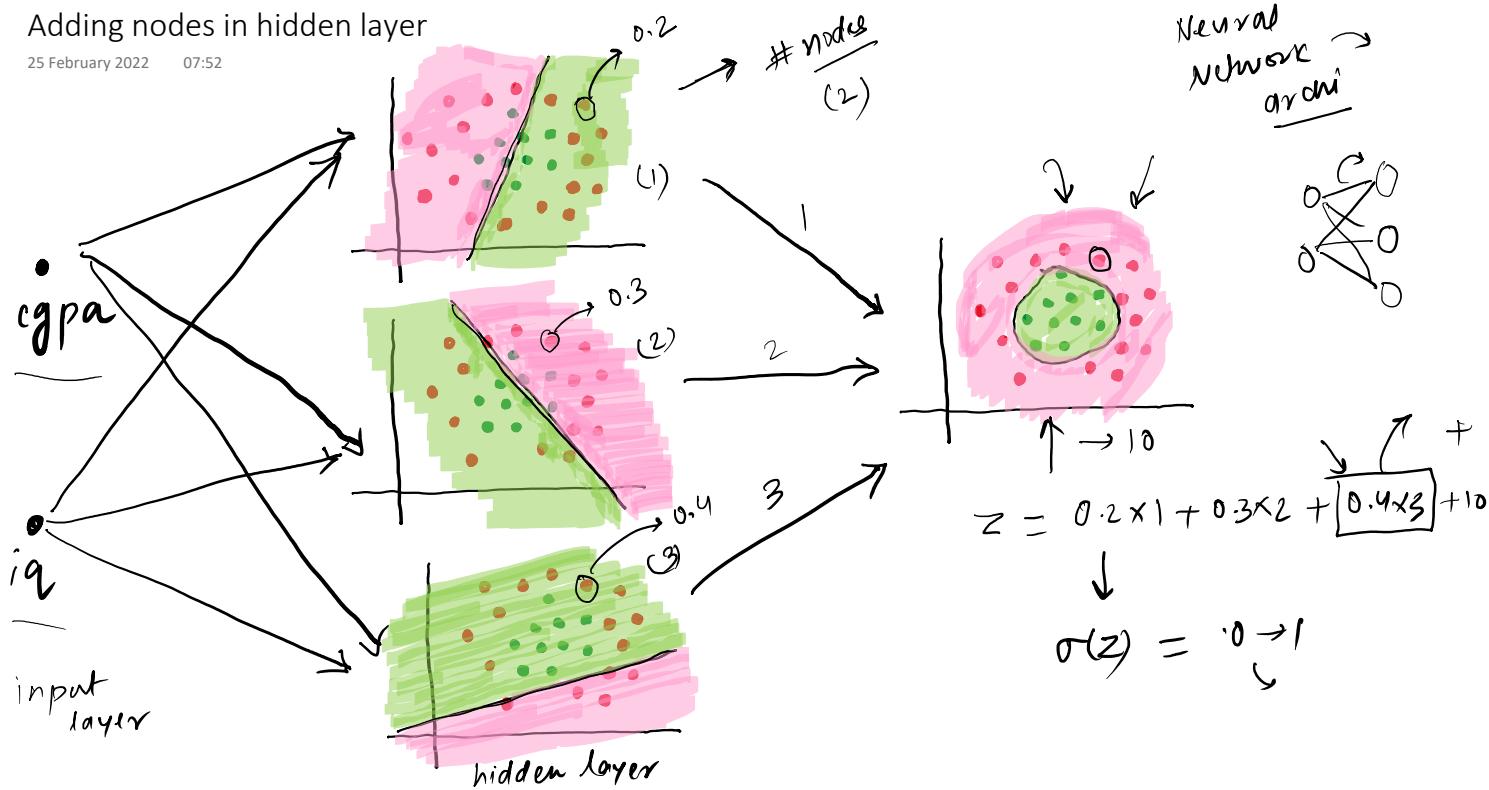
24 February 2022 15:37

c gpa | iq place



## Adding nodes in hidden layer

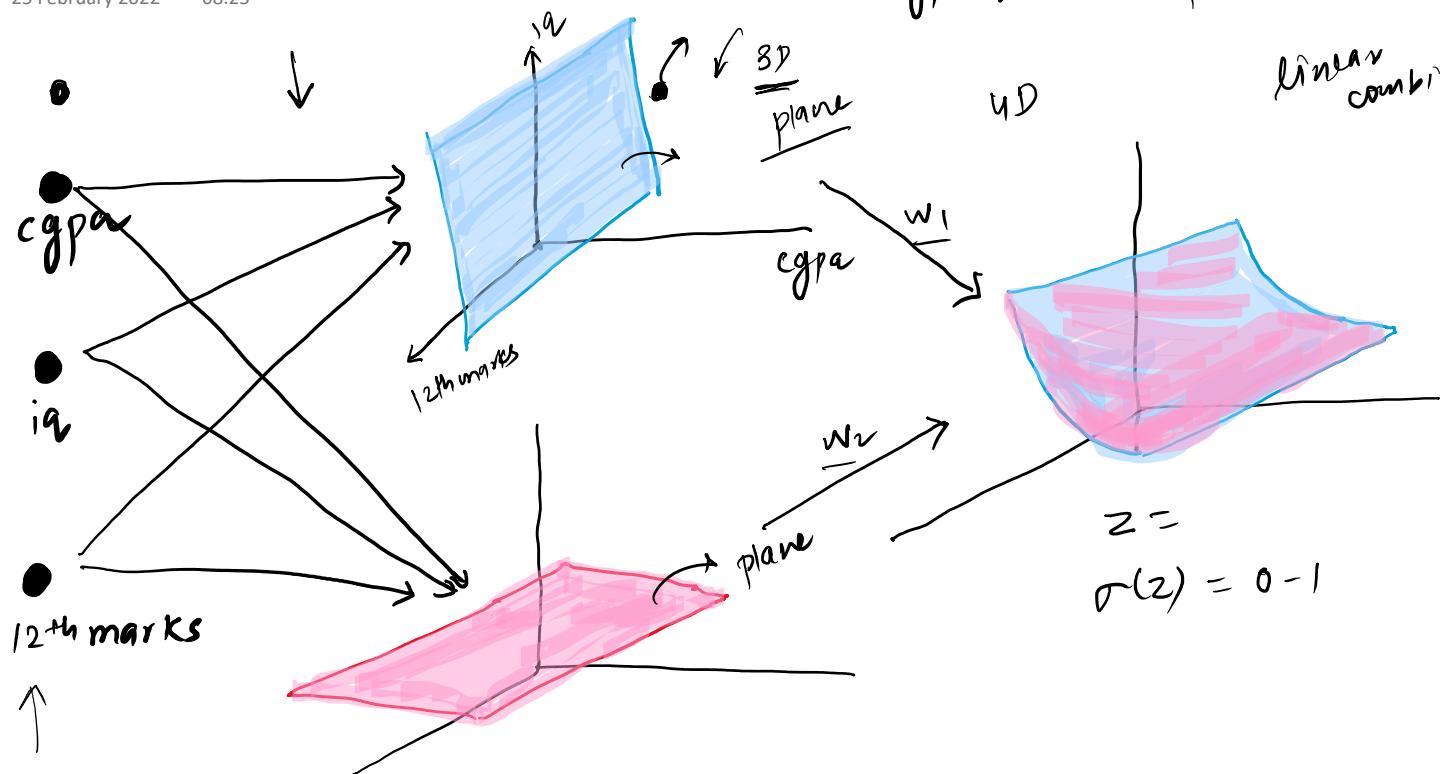
25 February 2022 07:52



## Adding nodes in input

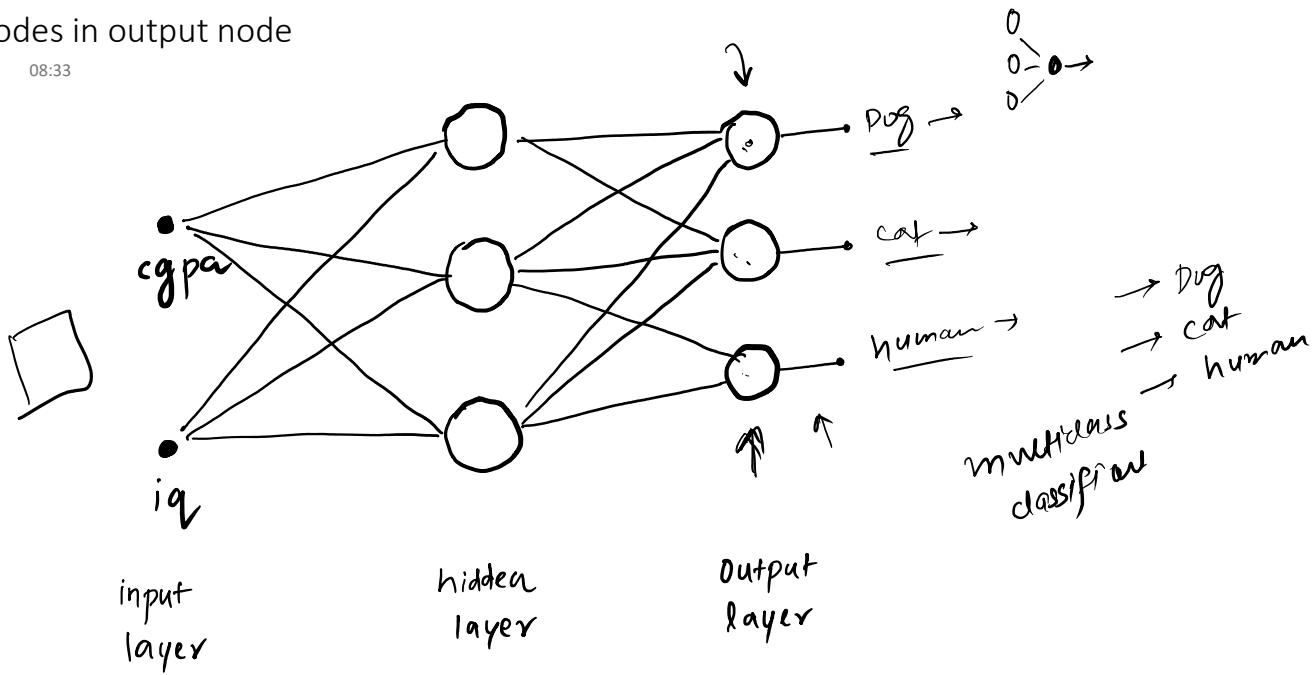
25 February 2022 08:25

cgpa iq 12<sup>th</sup> marks



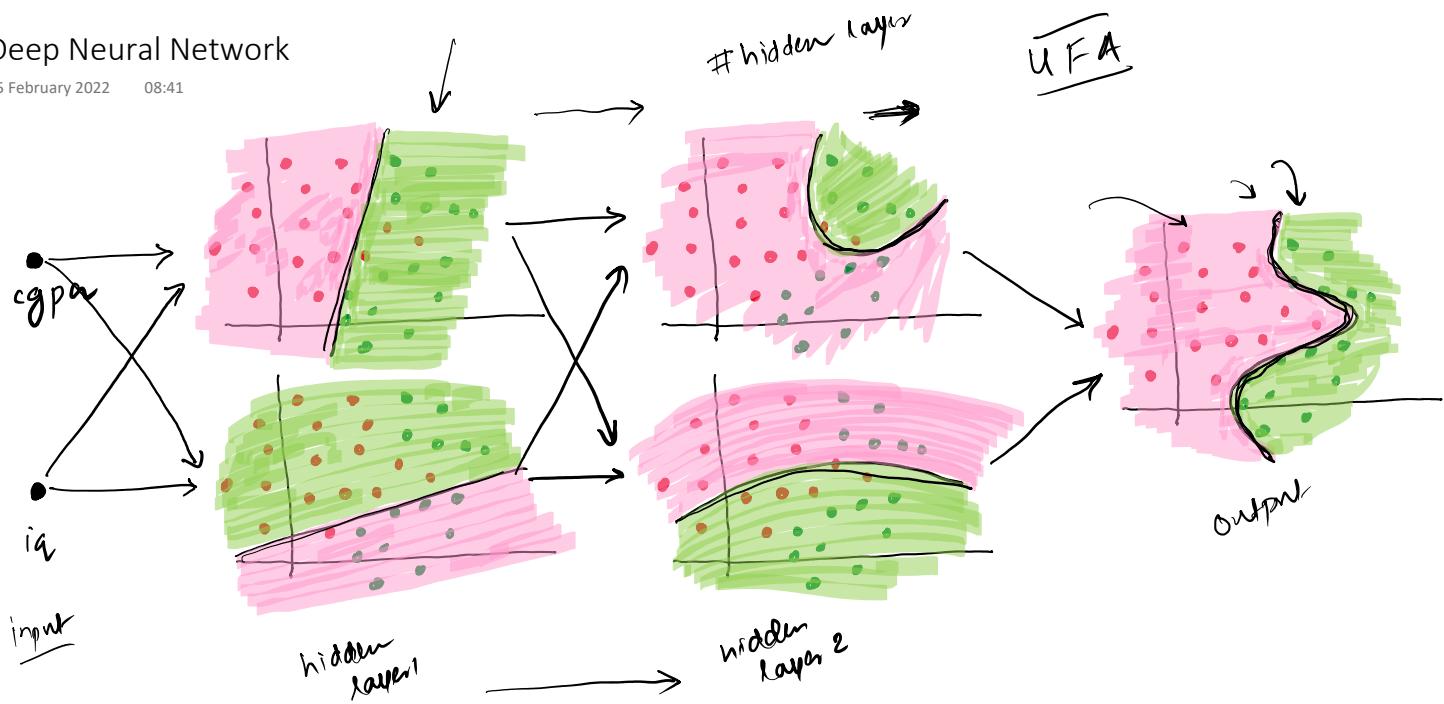
## Adding nodes in output node

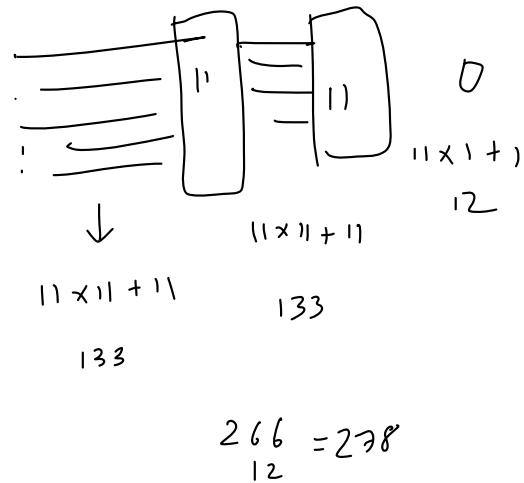
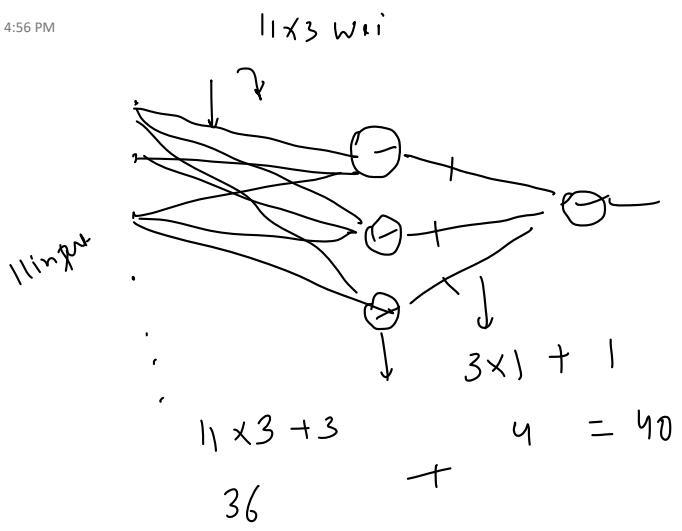
25 February 2022 08:33



# Deep Neural Network

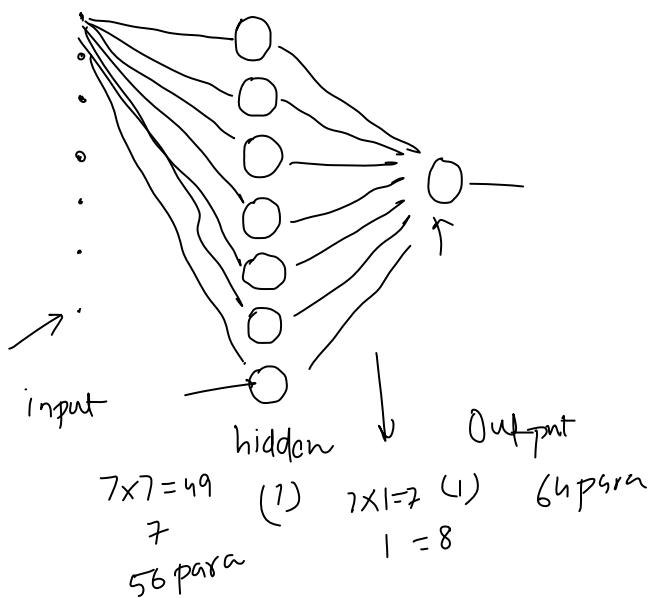
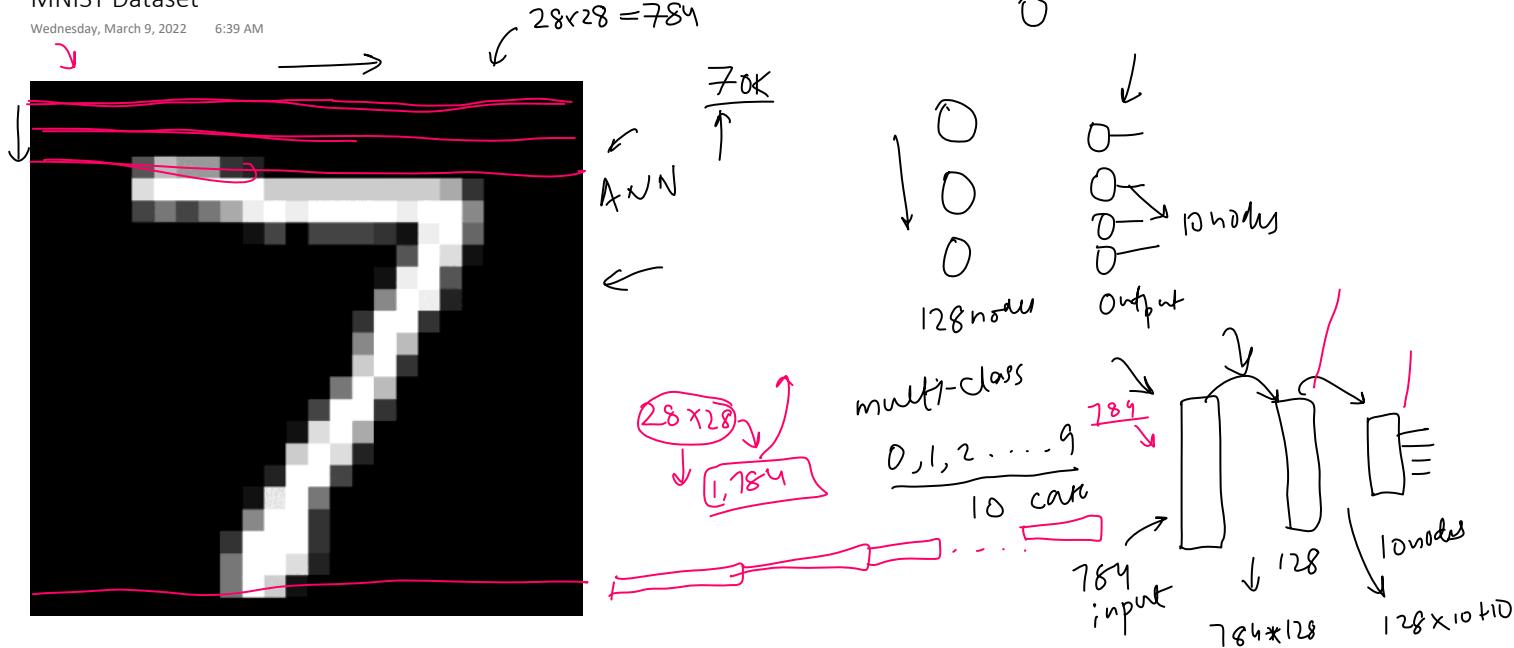
25 February 2022 08:41





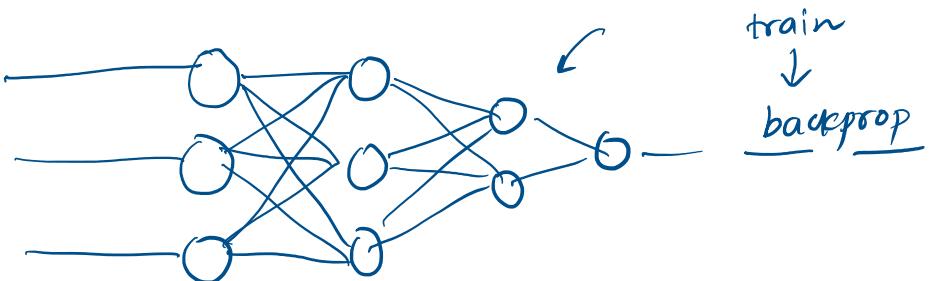
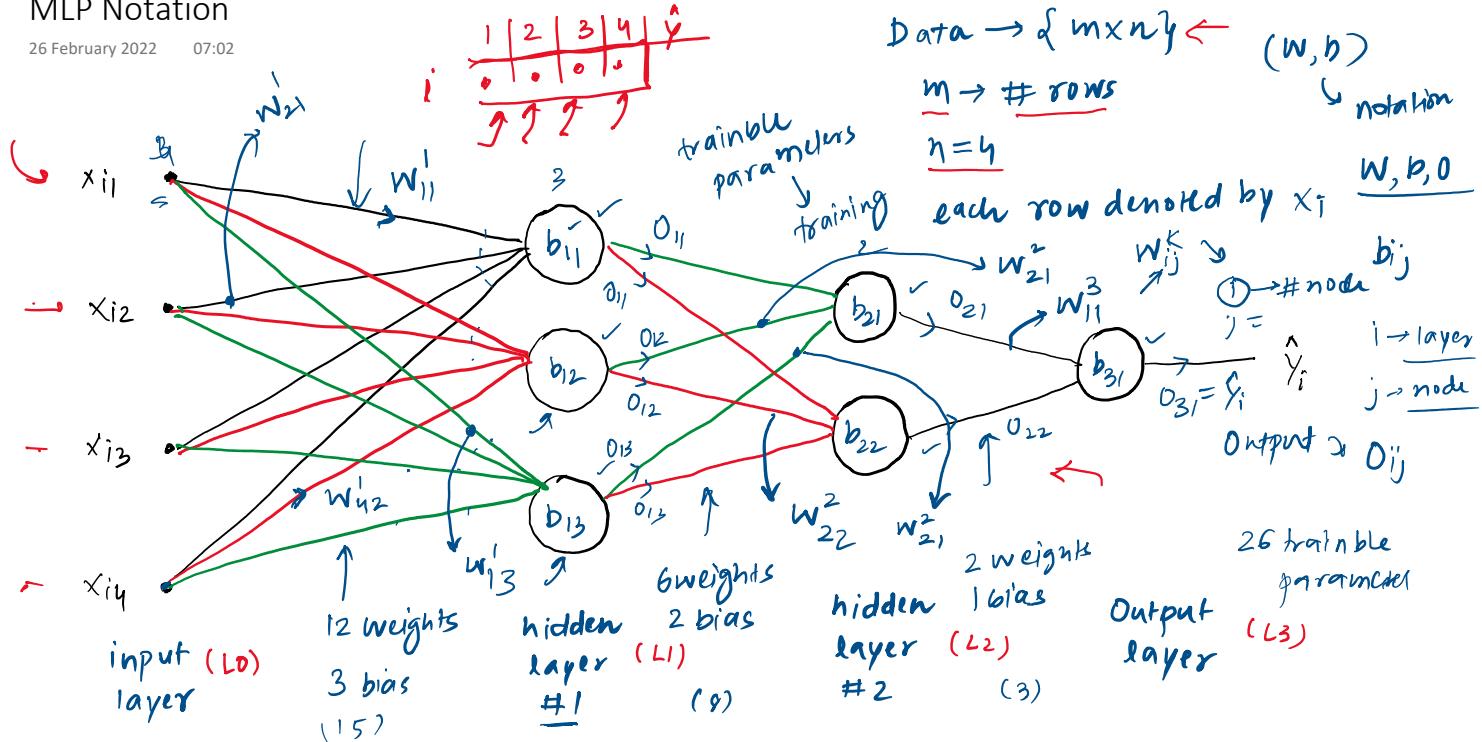
## MNIST Dataset

Wednesday, March 9, 2022 6:39 AM



## MLP Notation

26 February 2022 07:02



## Forward Propagation

03 March 2022

06:07

$\text{c gpa}$  |  $i_1$  |  $10 \text{th m}$  |  $12 \text{th m}$  | placed |  $\hat{y}$  # of trainable parameters

7.2	72	69	81	1
8.1	92	75	76	0

$a^{[1]} = \sigma(a^{[0]} w^{[1]} + b^{[1]})$

$a^{[2]} = \sigma(a^{[1]} w^{[2]} + b^{[2]})$

$a^{[3]} = \sigma(a^{[2]} w^{[3]} + b^{[3]})$

$O_{31} = \hat{y}$

$\text{prediction} \rightarrow \sigma(w^T x + b) \rightarrow \hat{y}$

$\sigma((\sigma((\sigma(a^{[0]} w^{[1]} + b^{[1]})) w^{[2]} + b^{[2]})) w^{[3]} + b^{[3]})$

$a^{[3]} = \hat{y}$

Layer #1

$$= \begin{bmatrix} W_{11}^1 x_{i1} + W_{12}^1 x_{i2} + W_{13}^1 x_{i3} + W_{14}^1 x_{i4} \\ W_{12}^1 x_{i1} + W_{22}^1 x_{i2} + W_{23}^1 x_{i3} + W_{24}^1 x_{i4} \\ W_{13}^1 x_{i1} + W_{23}^1 x_{i2} + W_{33}^1 x_{i3} + W_{34}^1 x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$= \sigma \left( \begin{bmatrix} W_{11}^1 x_{i1} + W_{12}^1 x_{i2} + W_{13}^1 x_{i3} + W_{14}^1 x_{i4} + b_{11} \\ W_{12}^1 x_{i1} + W_{22}^1 x_{i2} + W_{23}^1 x_{i3} + W_{24}^1 x_{i4} + b_{12} \\ W_{13}^1 x_{i1} + W_{23}^1 x_{i2} + W_{33}^1 x_{i3} + W_{34}^1 x_{i4} + b_{13} \end{bmatrix} \right)$$

$$= \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} \rightarrow a^{[1]}$$

Layer #2

$$= \sigma \left( \begin{bmatrix} W_{11}^2 O_{11} + W_{12}^2 O_{12} + W_{13}^2 O_{13} + b_{21} \\ W_{12}^2 O_{11} + W_{22}^2 O_{12} + W_{23}^2 O_{13} + b_{22} \end{bmatrix} \right) = \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} \rightarrow a^{[2]}$$

Layer #3

$$= \sigma \left( \begin{bmatrix} W_{11}^3 O_{21} + W_{21}^3 O_{22} + b_{31} \end{bmatrix} \right) = \hat{y} = [O_{31}] \rightarrow a^{[3]}$$

## What is Loss Function?

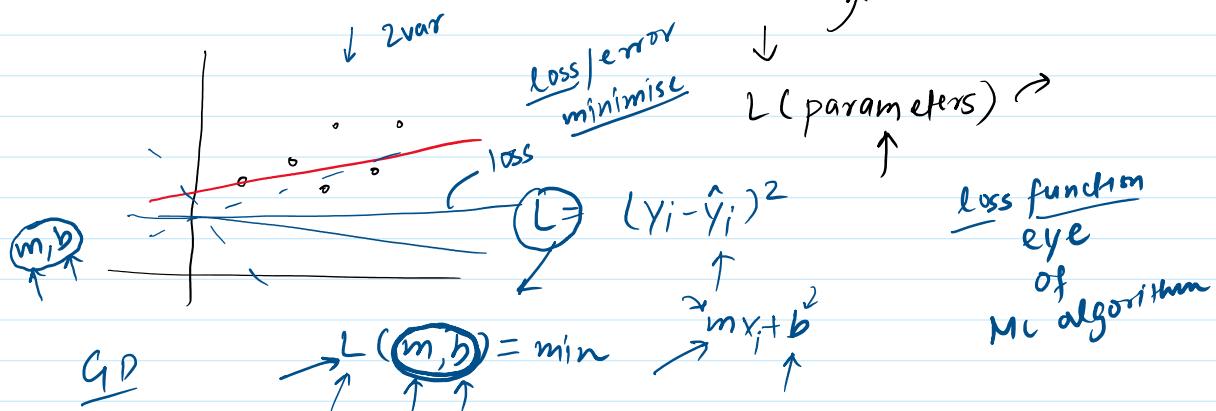
23 March 2022 07:15

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

high → poor

small → great

$$f(x) = x^2 + 2$$



Why is Loss function important?

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

Peter Drucker

multiplic  
epoche

Loss Function in Deep Learning?

Backprop

random w, b

gradient descent

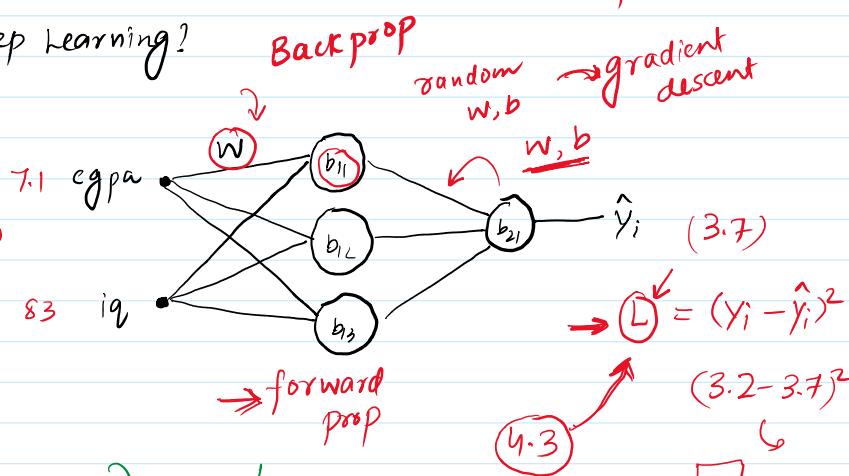
✓ cgpa | iq | package(epa)

	7.1	83
	8.5	91
	6.3	102
	5.1	87
..	..	..

7.1 cgpa

3.2  
4.5  
6.1  
2.7

83 iq



Keras

- object detection
- classification
- regression
- loss

{ mse  
mae  
huber loss

Loss functions in DL

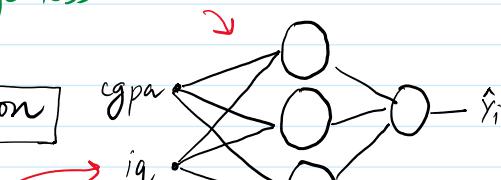
{ classification  
binary crossentropy  
categorical cross entro  
hinge loss

Autoencoders  
+ KL divergence

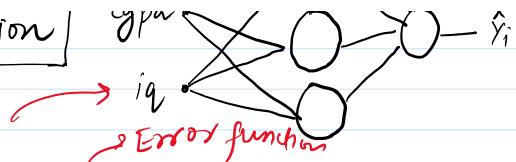
Embedding  
Tripled loss

GAN  
+ discriminator  
loss  
minmax  
gan loss

LOSS Function vs Cost Function



## Loss Function vs Cost Function



Loss function → single training eg

$$y_i = 6.3 \quad \hat{y}_i = 6.1$$

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

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

$$\frac{1}{4} [(6.1 - 6.3)^2 + (4.1 - 4)^2 + (3.5 - 3.7)^2 + (7.2 - 7)^2] = CF$$

batch Cost function

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

$$\textcircled{5} \text{ lpa } \sqrt{\sum (y_i - \hat{y}_i)^2}$$

$$(lpa)^2$$

### 1. Mean Squared Error (MSE)

squared loss L2 loss

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

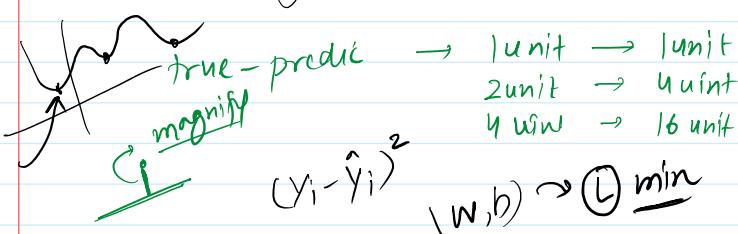
(true - predict)<sup>2</sup>

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

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

quadratic punish

magnitude



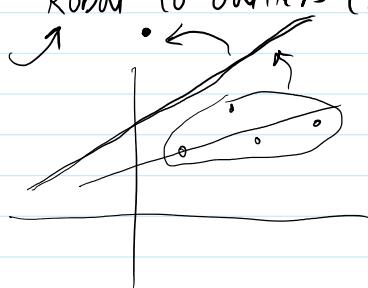
### Advantages

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

### Disadvantages

- 1) Error unit (squared) → diff
- 2) Robust to Outliers (Not)

$$CF = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



### 2. Mean Absolute Error (MAE) → L1 loss

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

$$2 |abs|$$

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|^2$$

punish

$$|true - predict| \rightarrow \frac{1}{n} \sum |y_i - \hat{y}_i|$$

$$csp |iq| pack lpa |(y_i - \hat{y}_i)|$$

$$lpa$$

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

true-predict  
punish

Advantages

- 1) Intuitive and easy
- 2) Unit  $\rightarrow$  same  $-y$
- 3) Robust to outliers

Disadvantage

- 1) Not differentiable

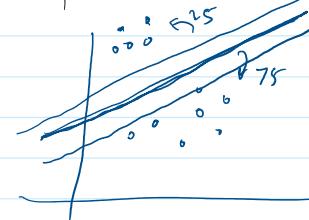
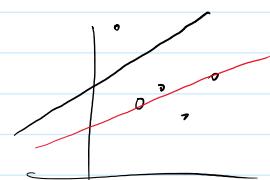
$\downarrow$  GD  $\approx$  differentiate  
Subgradient

mse - outliers ✓  
mae - normal point ✓

3. Huber Loss ✓

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

mae



4. Binary Cross Entropy

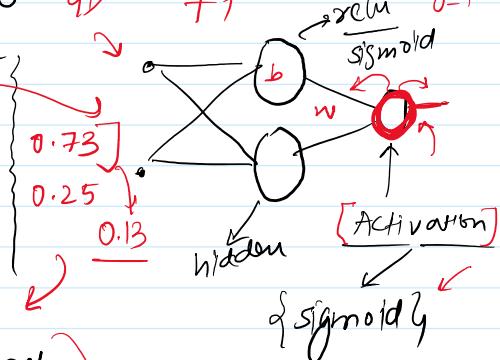
- classification ✓
- Two classes ✓

100 days  
cgpa | iq | placement

8	80
7	70
6	60
0	0

1 0

GD FP



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

Keras

$$0.12 = -(1-0) \log(1-0.25) - 1 \log(0.75)$$

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

$$-1 \times 0.12 = 0.12$$

maximum LLL

Logistic Reg

Advantage

- \* Differentiable

D/advan

- multi local minima
- Intuitive

5. Categorical Cross Entropy [used in Softmax Regression]

- multi-class classifications ✓

OHE

cgpa	iq	placed?	Yes	No	Maybe
------	----	---------	-----	----	-------

5. Categorical Cross Entropy [Used in softmax regression]

→ Multi-class classification

1 point

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

Where  $K$  is # classes in the data

cgpa	iq	placed?	Yes	No	Maybe
8	80	Yes 1	1	0	0
6	60	No 2	0	1	0
7	70	Maybe 3	0	0	1

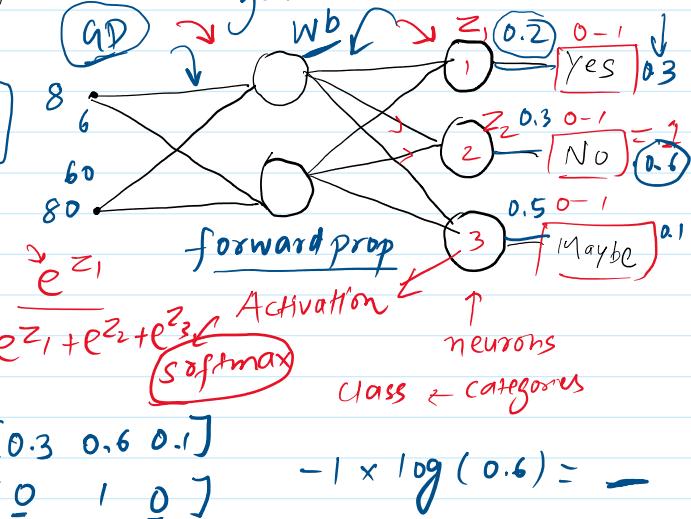
1 point

$$L = - (y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2) + y_3 \log(\hat{y}_3))$$

$$[0.2 \ 0.3 \ 0.5] \xrightarrow{e^{z_1} + e^{z_2} + e^{z_3}} f(z) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$[1 \ 0 \ 0] \xrightarrow{e^{z_2}} \text{Activation}$$

$$L = -1 \times \log(0.2) \xrightarrow{e^{z_1} + e^{z_2} + e^{z_3}} -1 \times \log(0.6) = -$$



### Sparse Categorical Cross Entropy

(L) ~

$$-1 \times \log(0.1)$$

SCE

Loss function  
fast

cgpa | iq | placed OHE

①	1/2/3	7	70	Yes	1
②		8	80	No	2
③		6	60	Maybe	3

$$[0.1 \ 0.4 \ 0.6]$$

$$[0.6 \ 0.2 \ 0.2]$$

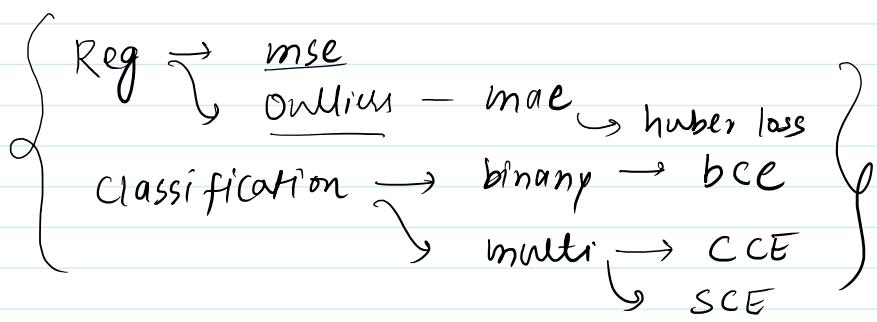
OHE — fast

SCE

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

Cost function

$$C = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$



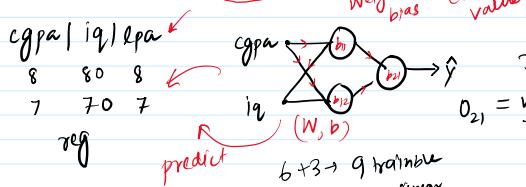


Pre requisite

- Gradient Descent (✓) →
- Forward Propagation (✓) →

Backprop  
↓  
Algo → train nn

(How)

trainingWeights  
bias → correct value $\sigma(z)$ 

$$O_{21} = W_{21}^T O_{11} + W_{21}^T O_{12} + b_{21}$$

predict

6 → 9 trainable

Backprop

Students		
iq	cgpa	lpa
80	8	3
60	9	5
70	5	8
120	7	11

- Steps  
0) Init  $w, b$  ✓  
1) You select a point (row) ↴  
↳ student

$$L \min \rightarrow \hat{y} \rightarrow O_{21}$$

activation = linear

linear

mse

Random val

 $w \rightarrow \frac{1}{2}, b \rightarrow 0$ 

error

- 2) Predict ( $lpa$ ) → forward prop [dot product]  
↳

$$(3-18)^2 = 225$$

- 3) Choose a loss function (mse)

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

- 4) Weights and bias update ✓  
↳ Gradient descent

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

$$(3-18)^2 = 225$$

$$\text{error}$$

$\hat{y} = O_{21}$

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

$W_{21}^2_{\text{new}} = W_{21}^2_{\text{old}} - \eta \frac{\partial L}{\partial W_{21}^2}$

$b_{21 \text{ new}} = b_{21 \text{ old}} - \eta \frac{\partial L}{\partial b_{21}}$

$W_{21}^2_{\text{new}} = W_{21}^2_{\text{old}} - \eta \frac{\partial L}{\partial W_{21}^2}$

$b_{21 \text{ new}} = b_{21 \text{ old}} - \eta \frac{\partial L}{\partial b_{21}}$

derivative of loss wrt weight

$$\frac{\partial L}{\partial W_{21}^2}$$

↳ partial derivative

derivative

$$\frac{\partial L}{\partial W_{21}^2}$$

chain rule of differen  
(L min)  
 $\frac{dL}{dx}$

$\frac{\partial L}{\partial W_{21}^2}, \frac{\partial L}{\partial W_{21}^1}, \frac{\partial L}{\partial b_{21}}$

$\frac{\partial L}{\partial W_{21}^1}, \frac{\partial L}{\partial W_{21}^1}, \frac{\partial L}{\partial b_{21}}$

$\frac{\partial L}{\partial W_{21}^2}, \frac{\partial L}{\partial W_{21}^2}, \frac{\partial L}{\partial b_{21}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

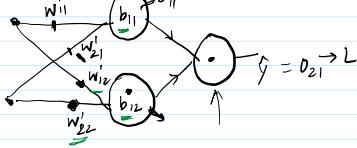
$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial \hat{y}}$

$$\frac{\partial L}{\partial w_{21}^2} = -2(y - \hat{y}) \delta_{12} \quad | \quad 2 \quad \begin{array}{c} \hat{y} \rightarrow b \\ \uparrow \\ L \end{array} \quad \frac{\partial \hat{y}}{\partial b_{21}} = 1$$

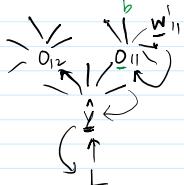
$$\frac{\partial L}{\partial b_{21}} = -2(y - \hat{y}) \quad | \quad 3 \quad \frac{\partial L}{\partial b_{21}} = \left[ \frac{\partial L}{\partial \hat{y}} \right] \times \frac{\partial \hat{y}}{\partial b_{21}}$$



$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{11}} \frac{\partial \delta_{11}}{\partial w_{11}^2} \rightarrow [-2(y - \hat{y}) w_{11}^2 x_{11}] \quad 4$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{12}} \frac{\partial \delta_{12}}{\partial w_{21}} \rightarrow [-2(y - \hat{y}) w_{21}^2 x_{12}] \quad 5$$

$$\frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{11}} \frac{\partial \delta_{11}}{\partial b_{11}} \rightarrow [-2(y - \hat{y}) w_{11}^2] \quad 6$$



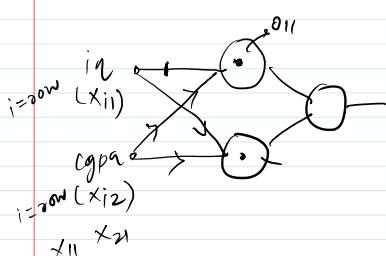
$$\frac{\partial L}{\partial w_{12}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{12}} \frac{\partial \delta_{12}}{\partial w_{12}^2} \rightarrow [-2(y - \hat{y}) w_{12}^2 x_{12}] \quad 7$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{12}} \frac{\partial \delta_{12}}{\partial w_{22}^2} \rightarrow [-2(y - \hat{y}) w_{22}^2 x_{12}] \quad 8$$

$w_{11}^2$

$$\frac{\partial \hat{y}}{\partial \delta_{12}} = w_{21}^2$$

$$\frac{\partial L}{\partial b_{12}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{12}} \frac{\partial \delta_{12}}{\partial b_{12}} \rightarrow [-2(y - \hat{y}) w_{21}^2] \quad 9$$



$$\frac{\partial \delta_{11}}{\partial w_{11}^1} = \underbrace{\partial \delta_{11}}_{\text{forward}} \frac{\partial \delta_{11}}{\partial w_{11}^1} + \underbrace{cgpa w_{21}^1}_{\text{backward}} + \underbrace{b_{11}}_{\text{bias}} = (iq)$$

$$\frac{\partial \delta_{11}}{\partial w_{11}^1} = x_{11}$$

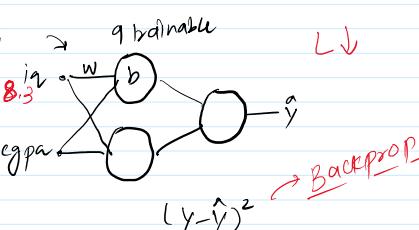
$$\frac{\partial \delta_{11}}{\partial b_{11}} = 1$$

$$\frac{\partial \delta_{12}}{\partial w_{12}^1} = \frac{\partial}{\partial w_{12}^1} [ \underbrace{iq w_{12}^1}_{\text{forward}} + \underbrace{cgpa w_{22}^1}_{\text{backward}} + \underbrace{b_{12}}_{\text{bias}} ] = iq(x_{12})$$

$$\frac{\partial \delta_{12}}{\partial w_{22}^1} = x_{12}$$

$$\frac{\partial \delta_{12}}{\partial b_{12}} = 1$$

		multiple choice	
		cgpa	iq
		lpa	y
1	epochs	8	8
2	steps (once again)	6	6
3	loop - 100/100	7	7
4	for i in range(4):	9	9



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

12 min

1a) 1 student  $\rightarrow$  forward prop  $\rightarrow$

lpa  $\rightarrow$

1b) Loss calculate (mse)  $\rightarrow$

1c) Adjust all weights and bias  $\rightarrow$

optimus

The How

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

### Backpropagation Algorithm

epochs=5

for i in range(epochs):

for j in range(x.shape[0]):

$$\frac{\partial L}{\partial w_{11}^2} = \underbrace{-2(y - \hat{y}) \delta_{11}}_{\text{Ready}}$$

$$\frac{\partial L}{\partial w_{21}^2} = \underbrace{-2(y - \hat{y}) \delta_{12}}_{\text{Ready}}$$

for i in range(epochs):

    for j in range(x.shape[0]):

        → Select 1 row (random)

        → Predict (using Forward prop)

        → Calculate loss (using Loss function → mse)

        → Update weights and bias using GD

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

    → Calculate avg loss for the epoch

$L_{avg}$

$$\frac{\partial L}{\partial w_{21}} = -2(y - \hat{y}) o_{21} \quad \checkmark$$

$$\frac{\partial L}{\partial b_{21}} = -2(\hat{y} - y) \quad \checkmark$$

$$\frac{\partial L}{\partial w_{11}} = -2(y - \hat{y}) w_{11}^2 x_{11} \quad \checkmark$$

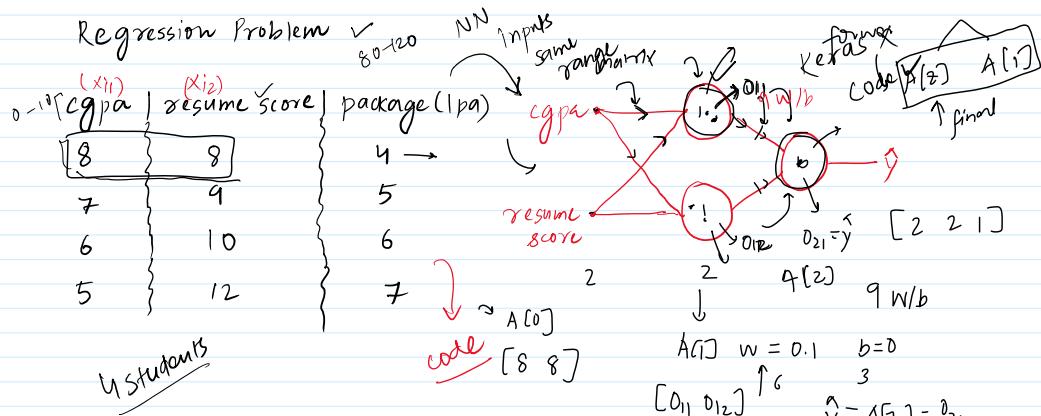
$$\frac{\partial L}{\partial w_{12}} = -2(y - \hat{y}) w_{12}^2 x_{12} \quad \checkmark$$

$$\frac{\partial L}{\partial b_{11}} = -2(y - \hat{y}) w_{11}^2 \quad \checkmark$$

$$\frac{\partial L}{\partial w_{12}} = -2(y - \hat{y}) w_{12}^2 x_{12} \quad \checkmark$$

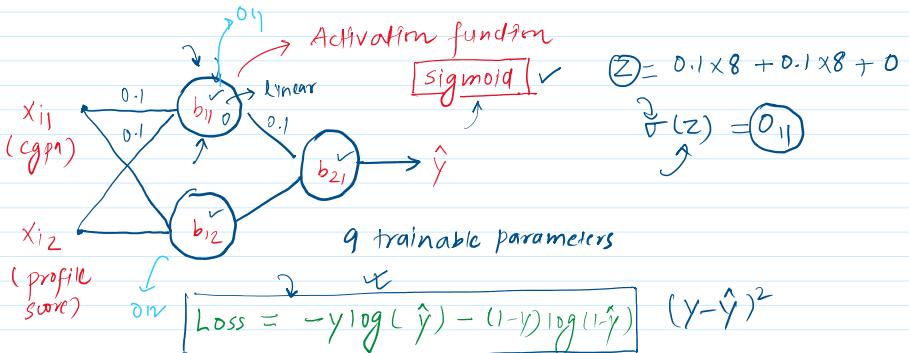
$$\frac{\partial L}{\partial w_{22}} = -2(y - \hat{y}) w_{22}^2 \quad \checkmark$$

$$\frac{\partial L}{\partial b_{12}} = -2(y - \hat{y}) w_{22}^2 \quad \checkmark$$



Classification Example

cgpa	profile score	placement
8	8	1
7	9	1
6	10	0
5	5	0



Backpropagation Algorithm

epochs = 5

for i in range(epochs): (rows)

    for j in range(x.shape[0]):

        → Select 1 row (random)  $\rightarrow 1by1$

        → Predict (using Forward prop)  $\rightarrow$  predict

        → Calculate loss (using Loss function  $\rightarrow$  bce)  $\rightarrow$  code convert

        → Update weights and bias using GD

## Backpropagation Algorithm

epochs = 5 ✓

for i in range(epochs): (rows)

for j in range(x.shape[0]):

→ Select 1 row (random) → by

→ Predict (using Forward prop) → predict

→ Calculate loss (using Loss function → bce) → code convert

→ Update weights and bias using GD

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

q w, b update

→ Calculate avg loss for the epoch

Lavg

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

$$z = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_2$$

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

der wrt z

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial w_{11}} = -(y-\hat{y}) o_{11} \quad (1)$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial w_{21}} = -(y-\hat{y}) o_{12} \quad (2)$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial b_2} = -(y-\hat{y})$$

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

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

$$\hat{y} = \sigma(z) \quad \frac{\partial \sigma(z)}{\partial z} = \sigma(z)[1 - \sigma(z)] \quad \sigma(z) = \hat{y}$$

$$\boxed{\frac{\partial \hat{y}}{\partial z}} = \frac{\partial}{\partial z} (\sigma(z)) = \sigma(z)[1 - \sigma(z)] = \hat{y}(1-\hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{(y-\hat{y})}{\hat{y}(1-\hat{y})} \quad \frac{\partial \hat{y}}{\partial z} = \hat{y}(1-\hat{y}) \Rightarrow \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} = \frac{-(y-\hat{y}) \times \hat{y}(1-\hat{y})}{\hat{y}(1-\hat{y})} = -(y-\hat{y})$$

$$\frac{\partial z_p}{\partial w_{11}} = w_{11}' x_{11} + w_{21}' x_{12} + b_{11} \quad \rightarrow -(y-\hat{y})$$

$$z_f = (w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_2) o_{11} (1-o_{11})$$

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

cross entropy

$$z_f = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21} o_{11}(1-o_{11})$$

$$o_{11} = \sigma(z_{\text{prev}}) = \sigma(z_{\text{prev}})[1-\sigma(z_{\text{prev}})]$$

$$\frac{\partial L}{\partial w_{11}} = \left[ \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_f} \right] \times \frac{\partial z_f}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial z_{\text{prev}}} \times \frac{\partial z_{\text{prev}}}{\partial w_{11}} \rightarrow x_{i1}$$

$\downarrow$   
 $w_{11}^2$   
 $\downarrow$   
 $o_{11}(1-o_{11})$

$$\boxed{\frac{\partial L}{\partial w_{11}} = -(y - \hat{y}) w_{11}^2 o_{11}(1-o_{11}) x_{i1}} \quad (3)$$

$$\boxed{\frac{\partial L}{\partial w_{21}} = -(y - \hat{y}) w_{21}^2 o_{12}(1-o_{12}) x_{i2}}$$

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

$$\frac{\partial z_f}{\partial o_{12}} = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21} o_{12}(1-o_{12})$$

$$\frac{\partial L}{\partial w_{12}} = \left[ \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_f} \right] \times \frac{\partial z_f}{\partial o_{12}} \times \frac{\partial o_{12}}{\partial z_p} \times \frac{\partial z_p}{\partial w_{12}} \rightarrow x_{i1}$$

$\downarrow$   
 $w_{21}^2$   
 $\downarrow$   
 $o_{12}(1-o_{12})$

$$\boxed{\frac{\partial L}{\partial w_{12}} = -(y - \hat{y}) w_{21}^2 o_{12}(1-o_{12}) x_{i1}} \quad (3)$$

$$\boxed{\frac{\partial L}{\partial w_{22}} = -(y - \hat{y}) w_{21}^2 o_{12}(1-o_{12}) x_{i2}}$$

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

Backpropagation → The Why?

The intuition behind the algorithm ✓

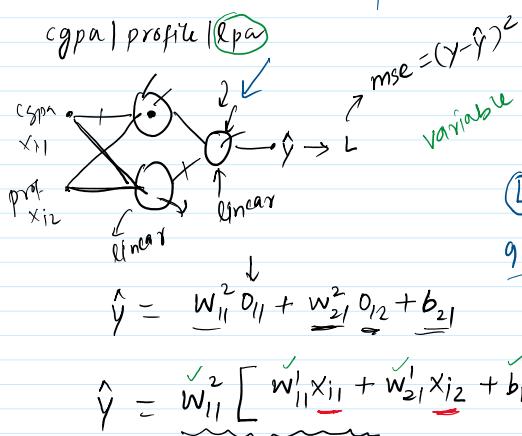
## Backpropagation Algorithm

```

epochs = 5
for i in range(epochs):
    for j in range(x.shape[0]):
        → Select 1 row (random)
        → Predict (using Forward prop) →  $\hat{y}$ 
        → Calculate loss (using loss function → mse)
        → Update weights and bias using GD
             $W_n = W_0 - \eta \frac{\partial L}{\partial W}$  rule
            magic
        → Calculate avg loss for the epoch
             $L_{avg}$ 
    → bce

```

→ Loss function is a function of all trainable parameters



$$L = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

variable

$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$

mathematical function

$f(x) = y$

→ Concept of Gradient

Gradient descent

fancy word for derivative

$$y = f(x) = x^2 + x$$

$$\frac{dy}{dx} = \frac{d}{dx}(fx) = \frac{d}{dx}(x^2 + x) = 2x + 1$$

derivative

$y \rightarrow x \rightarrow$  derivative

$$\frac{d}{dx}$$

$$z = f(x, y) = x^2 + y^2$$

$$\frac{\partial z}{\partial x} = 2x \quad \frac{\partial z}{\partial y} = 2y$$

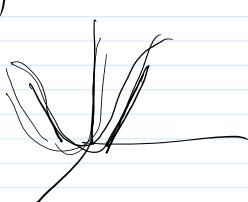
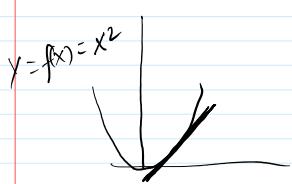
gradient  $\left[ \begin{array}{c} \frac{\partial z}{\partial x} \\ \frac{\partial z}{\partial y} \end{array} \right]$

gradient complex  $L(w_{11}, w_{12}, \dots, w_{21}, b_{11}, b_{12})$

$$q \text{ param} \quad \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$$

$qD$  function  $\rightarrow$  qdiff steps wrt each dim

$$3D \rightarrow z = x^2 + y^2 \quad z = f(x, y)$$



→ Concept of Derivative → Derivative  
at or  
point  
intuition

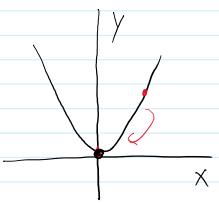
$$\frac{\partial L}{\partial w_i}$$

$w_{11}^1 = 1 \text{ unit}$

$$\begin{aligned} \text{derivative} & \quad \boxed{x=5} \rightarrow \text{deriv} \\ y = x^2 + 2x & \\ \frac{dy}{dx} = (2x+1)_{x=5} & \\ \left( \frac{dy}{dx} \right)_{x=5} = 11 & \rightarrow \text{slope} \end{aligned}$$

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

→ the concept of minima



$$Y = X^2$$

$$\frac{dy}{dx} = 2x = 0$$

$\downarrow$

$x=0$

$$\left. \begin{array}{l} \boxed{\begin{aligned} z &= x^2 + y^2 \\ (x,y) &\rightarrow z \min \end{aligned}} \quad \begin{array}{c} \text{graph of } z = x^2 + y^2 \\ \text{at } (0,0) \end{array} \\ \begin{array}{l} \frac{\partial z}{\partial x} = 0 \quad x \rightarrow \\ \frac{\partial z}{\partial y} = 0 \quad y \rightarrow \end{array} \quad \begin{array}{l} z \min \\ (x,y) \\ (0,0) \\ 2x = 0 \Rightarrow x = 0 \\ 2y = 0 \Rightarrow y = 0 \end{array} \end{array} \right\}$$

$L$  (9 param)  
 $w, b$   $\min$

L ↴ J

$$\frac{\partial L}{\partial w_{11}} = \dots = \frac{\partial L}{\partial b_{12}} = 0$$

9dim 2) minima

## → Backprop Intuition ↗ ↘ $\eta = 1$

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W} \quad \xleftarrow{\text{qstep}}$$

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

Diagram illustrating the backpropagation of error through a neural network. A bottom layer node labeled "b" receives input from two nodes above it. The top-left node has a weight "w" and is labeled "b". The top-right node is labeled "d". An arrow points from the bottom layer node to the right, labeled "L=".

9  
8  
constant

$$\underline{L} \left( \begin{matrix} b_{21} \\ \vdash \end{matrix} \right)$$

$$b_{21} = \boxed{b_2} \quad \boxed{\frac{\partial L}{\partial b_{21}}} \quad b_{21} = 5 \text{ v/e}$$

$$\frac{\partial L}{\partial b_{21}} \xrightarrow{\text{Smart}} \begin{matrix} \text{L derivative w.r.t } b_{21} \\ \uparrow \\ ? \end{matrix} \xrightarrow{\text{uni}}$$

LJ

$$\frac{\partial L}{\partial b_{21}} = -ve$$

$b_{21} \uparrow \quad L \downarrow$

game . . . nine

-ve of the gradient

$$b_2 +$$



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

$$W_n = W_0 - \boxed{\frac{\partial L}{\partial x}}$$

2

$$W_1 = W_0 - \eta \frac{\partial L}{\partial w}$$

## → Effect of Learning Rate ( $\eta$ )

$$W_n = W_0 - \eta \frac{\partial L}{\partial w}$$

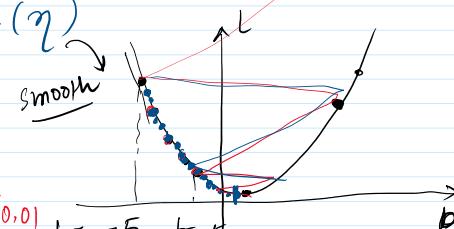
$$\frac{\partial L}{\partial b} = -50$$

$$10 = -5 - (-50) = 45$$

$$b = -5 \quad \frac{\partial L}{\partial b} = -50 \quad q = 0,01$$

$$b = -5 + (0.01 \times 50) = 0.5$$

$$= -4.5$$



train

$q$  = parameter  
↓  
0.01

$$\eta = 0.00001$$

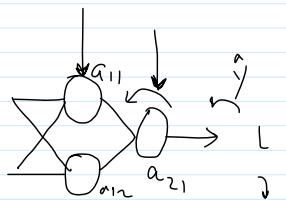
$$\frac{\partial L}{\partial b} = 50$$

$$b = 45 - 50$$

→ What is convergence?

The diagram illustrates the backpropagation update rule and its execution flow:

- Update Rule:**  $w_n = w_0 - \eta \frac{\partial L}{\partial w} \approx 0$
- Initial Condition:**  $w_n = w_0$
- Execution Flow:**
  - while  $\rightarrow$  com
  - for 1000
- Graph:** A U-shaped curve representing a function, likely the loss function.



$$\frac{\partial L}{\partial w_{ii}^2} = \left( \frac{\partial L}{\partial a_{ii}} \right) \times \frac{\partial a_{ii}}{\partial w_{ii}^2}$$

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

$$-(y - \hat{y})$$

$$\hat{y}(y - \hat{y}) \mathbf{0}_1$$

$$-\gamma \log(\alpha_{z_1}) - (1-\gamma) \log(1-\alpha_{z_1}) - (\gamma - \hat{\gamma}) \alpha_{z_1}$$

$$-\frac{y}{\gamma} + \frac{(1-y)}{(1-\gamma)}$$

activation = ②

activation = ②

$$\frac{\partial L}{\partial w_{21}^2} = \text{curr\_activa} = a[1]$$

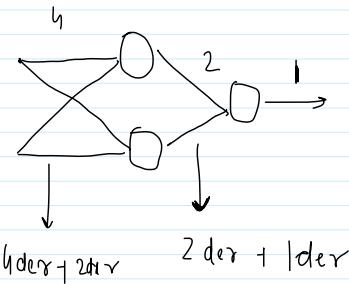
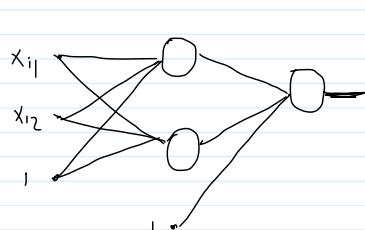
$$\frac{\partial L}{\partial b_{21}} =$$

$$\delta_{1+a} = (y - \hat{y}) a^{[0]} (1 - a^{[2]})$$

$$W[1] \quad [1 \times 2]$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{-\frac{1}{\gamma} + \frac{(1-\gamma)}{(1-\hat{y})}}{\left( \begin{array}{c} -(y - \hat{y}) a_{21} (1 - a_{21}) a_{11} \\ -(y - \hat{y}) a_{21} (1 - a_{21}) a_{12} \\ -(y - \hat{y}) a_{21} (1 - a_{21}) \end{array} \right)}$$

$$\frac{\partial L}{\partial b_{21}} = \frac{-(y - \hat{y}) a_{21} (1 - a_{21})}{\begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix}}$$



$1 \times 2 \quad 1 \times 2$

$1 \times 2 \quad x \quad 2 \times 2$

$$\frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial w}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial w}$$

$$\frac{\partial L}{\partial w_{11}} \quad \frac{\partial L}{\partial w_{12}}$$

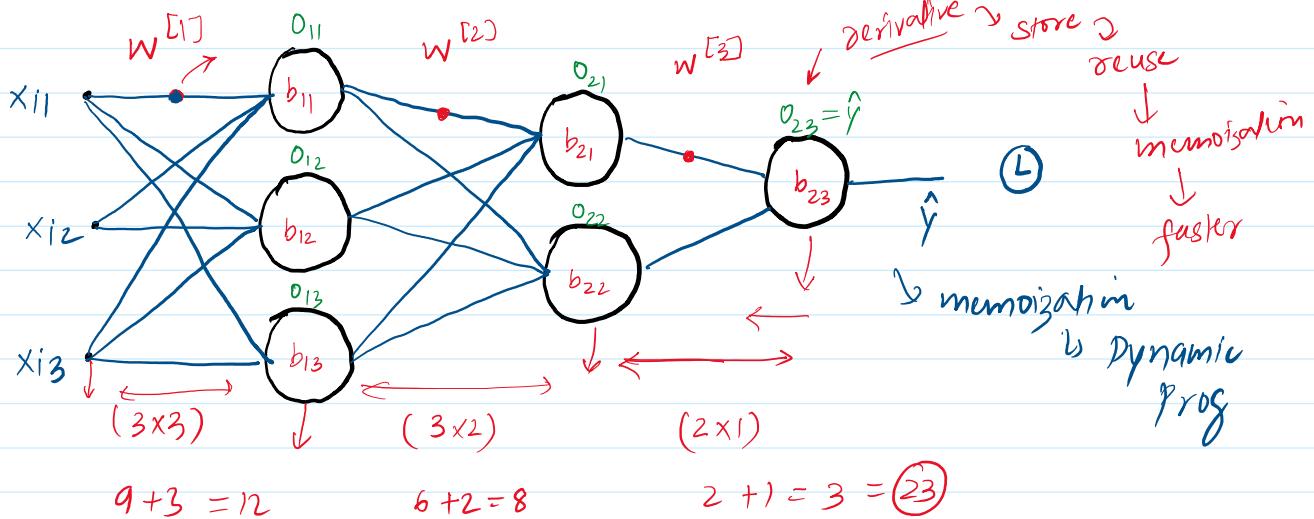
$$\frac{\partial L}{\partial w_{21}} \quad \frac{\partial L}{\partial w_{22}}$$

$$\frac{\partial L}{\partial w_{11}} \quad \frac{\partial L}{\partial w_{12}}$$

## MLP Memoization

Thursday, April 7, 2022 7:45 AM

backprop update  $\rightarrow$  w/b deriv



Chain rule

$$L \rightarrow \hat{y} \rightarrow o_{21} \rightarrow o_{11} \rightarrow w_{11}^1$$

$$\frac{\partial L}{\partial w_{11}^3} = \boxed{\frac{\partial L}{\partial \hat{y}}} \times \boxed{\frac{\partial \hat{y}}{\partial w_{11}^3}}$$

$$\boxed{\frac{\partial L}{\partial \hat{y}}} = \boxed{\frac{\partial L}{\partial \hat{y}}} \times \boxed{\frac{\partial \hat{y}}{\partial o_{21}}} \times \boxed{\frac{\partial o_{21}}{\partial w_{11}^1}}$$

$$x \rightarrow f(x), g(x) \rightarrow h(f(x), g(x)) = \frac{\partial h}{\partial x} = \left[ \frac{\partial h}{\partial f(x)} \times \frac{\partial f(x)}{\partial x} \right] + \left[ \frac{\partial h}{\partial g(x)} \times \frac{\partial g(x)}{\partial x} \right]$$

$$\frac{\partial \hat{y}}{\partial w_{11}^1} \rightarrow w_{11}^1 \rightarrow o_{11} \rightarrow o_{21} \rightarrow \hat{y} \rightarrow L$$

$$\frac{\partial L}{\partial w_{11}^1}$$

$$\frac{\partial L}{\partial w_{11}^1} = \boxed{\frac{\partial L}{\partial \hat{y}}} \left[ \boxed{\frac{\partial \hat{y}}{\partial o_{21}}} \times \boxed{\frac{\partial o_{21}}{\partial o_{11}}} \times \boxed{\frac{\partial o_{11}}{\partial w_{11}^1}} + \boxed{\frac{\partial \hat{y}}{\partial o_{22}}} \times \boxed{\frac{\partial o_{22}}{\partial o_{11}}} \times \boxed{\frac{\partial o_{11}}{\partial w_{11}^1}} \right]$$

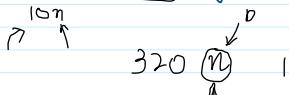
Backprop  $\rightarrow$  Chain diff (rule) + Memoization  
maths

Keras  
TF



50 values → smart replacement → warp  
 $y = 50$  values → vectorization → faster loop  
 $\hat{y}, y \rightarrow$  loss  
 $w, b$  update  $w_n = w_0 - \eta \frac{\partial L}{\partial w}$  Optimized  
→ loss

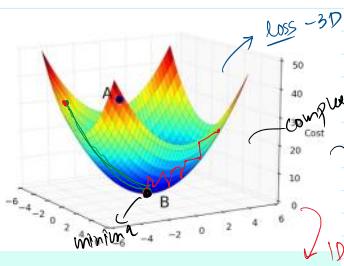
Which is faster (given same no. of epochs)



Which is the faster to converge (given same # epochs)



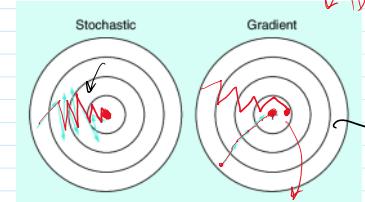
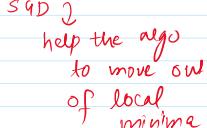
Loss function → 2 algo



Stochastic → random → point updates

Batch → dataset → all point update

Spirky SGD useful



Exact solution  
approximate diff  
faster → works

Vectorization → big dataset  
 $\text{np.dot}(x, w) + b$  RAM  
 $x \rightarrow$  dataset of 10 more

Mini Batch Gradient Descent → SGD → BGD

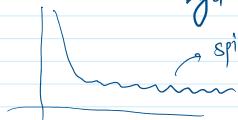
Best of both worlds

```
for i in range(nb_epochs):
    np.random.shuffle(data)
    for batch in get_batches(data, batch_size=50):
        params_grad = evaluate_gradient(loss_function, batch, params)
        params = params - learning_rate * params_grad
```

$\frac{n}{x} = \# \text{ of batches}$   
 $\# \text{ updates/epoch}$

Keras batch\_size =  $x$   
 $bgd > mbgd > sgd$   
 $mbgd < sgd$

MBGD

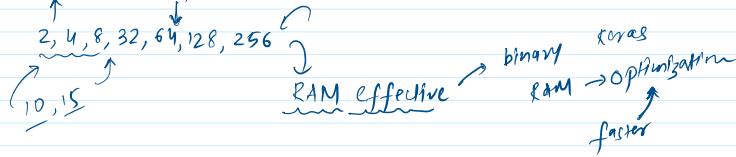


for i in epochs →  
for j in num of batch

1 batch  
 $\text{np.dot}$  →  $y_{pred}$  (vector)  
loss  
update

Vectorization → smaller batch

→ Why batch\_size is provided in multiple of 2?



→ What if batch\_size doesn't divide # rows properly

e.g. # of rows  $n = 400$   
batch\_size = 150

$$\# \text{ of batch} = \frac{400}{150} = 2.66$$

3  
150, 150, left 100  
↑ ↑ ↑  
1 batch 2 batch 3rd batch

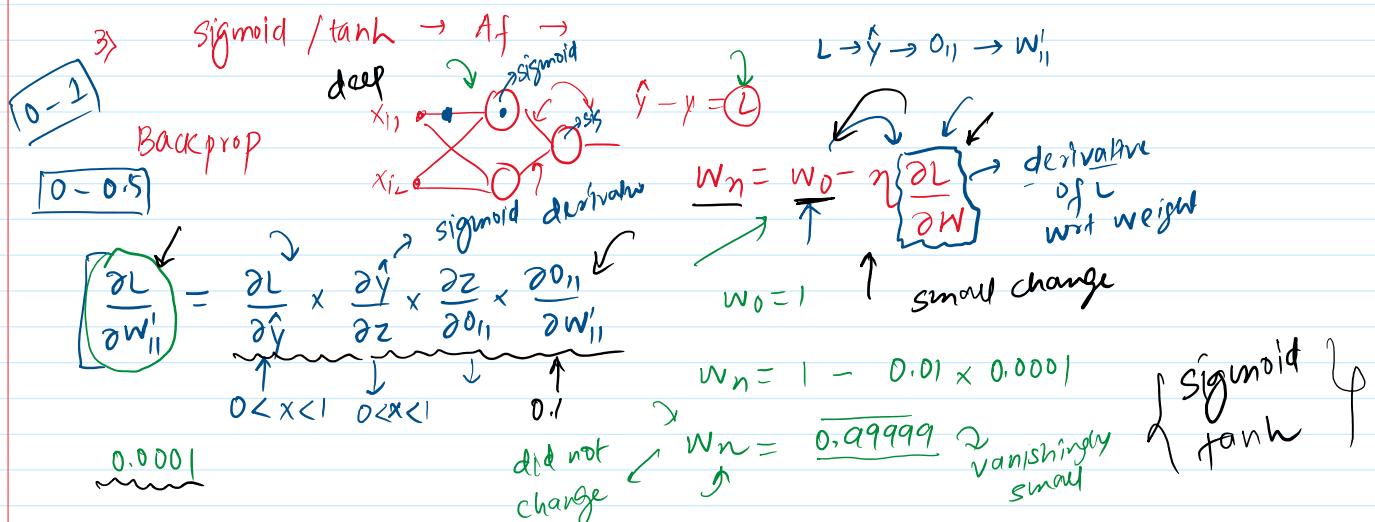
## Vanishing Gradient Problem

Thursday, April 7, 2022 7:45 AM

In machine learning, the **vanishing gradient problem** is encountered when training artificial neural networks with gradient-based learning methods and backpropagation. In such methods, during each iteration of training each of the neural network's weights receives an update proportional to the partial derivative of the error function with respect to the current weight. The problem is that in some cases, the gradient will be vanishingly small, effectively preventing the weight from changing its value. In the worst case, this may completely stop the neural network from further training. An example of the problem can be found in the diagram below.

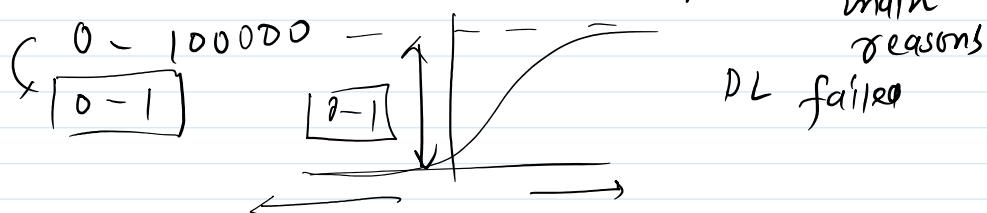
$$1) 0.1 \times 0.1 \times 0.1 \times 0.1 = [0.0001] \rightarrow VGP$$

2) Deep NN  $\rightarrow \square \square \square \square \square$



Backprop 2  
model training  $\rightarrow x$

$$\frac{\partial L}{\partial w_{11}} = 0.000001$$



How to recognize?

- 1) Loss focus  $\rightarrow$  epoch  $\rightarrow$  no changes  $\rightarrow$  VGP
  - 2) weights  $\rightarrow$  graph value
- Keras  
 $w_{11}$   
loss after epoch
- epoch

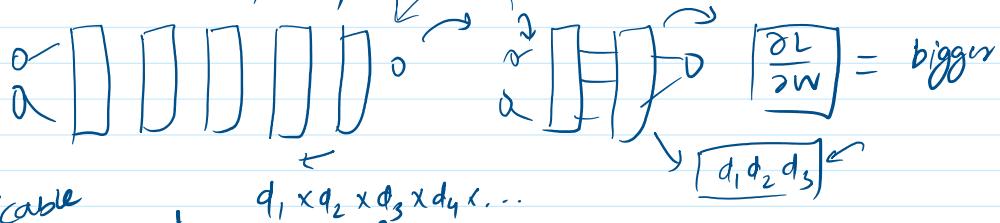
$$w_n = w_0 - \eta \left[ \frac{\partial L}{\partial w} \right]$$

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

$$w_n = w_0 - \eta \left[ \frac{\partial L}{\partial w} \right] \rightarrow \sqrt{\frac{w_0 - w_n}{\eta}} =$$

## How to handle Vanishing Gradient Problem →

- 1) Reduce model complexity



Applicable  
→ find out  
complex  
patterns

- 2) Using ReLU Activation functions

- 3) Proper weight init  $\xrightarrow{\text{Xavier}} \text{Glorot}$

- 4) Batch norm → layer →

- 5) Residual Network  $\xrightarrow{\text{building block}}$  ANN

## Exploding Gradient Problem

$$10, 10, 10, 10 \rightarrow 10000$$

random

~~Loss↓~~

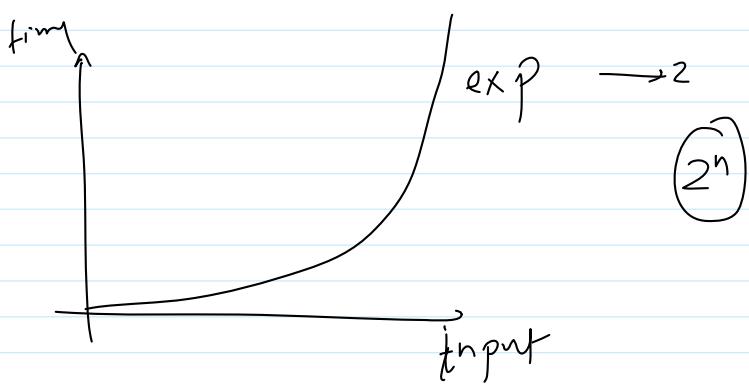
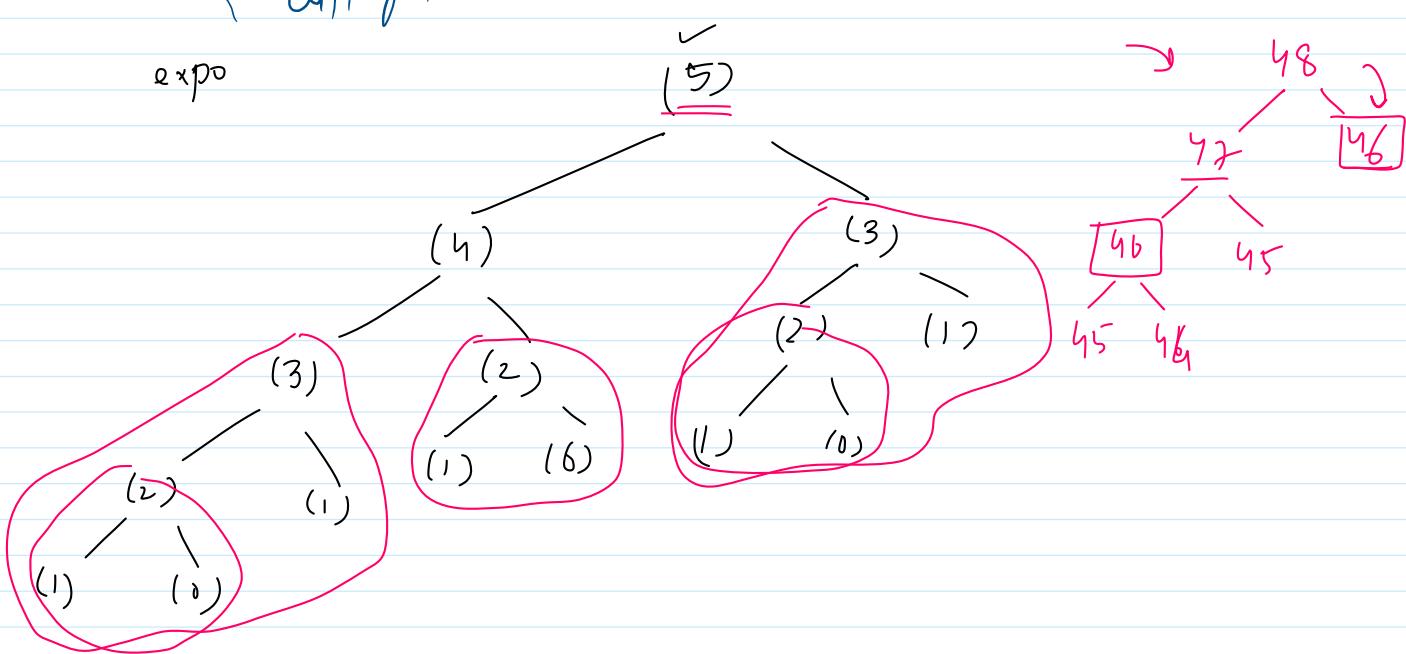
## Gradient Clipping

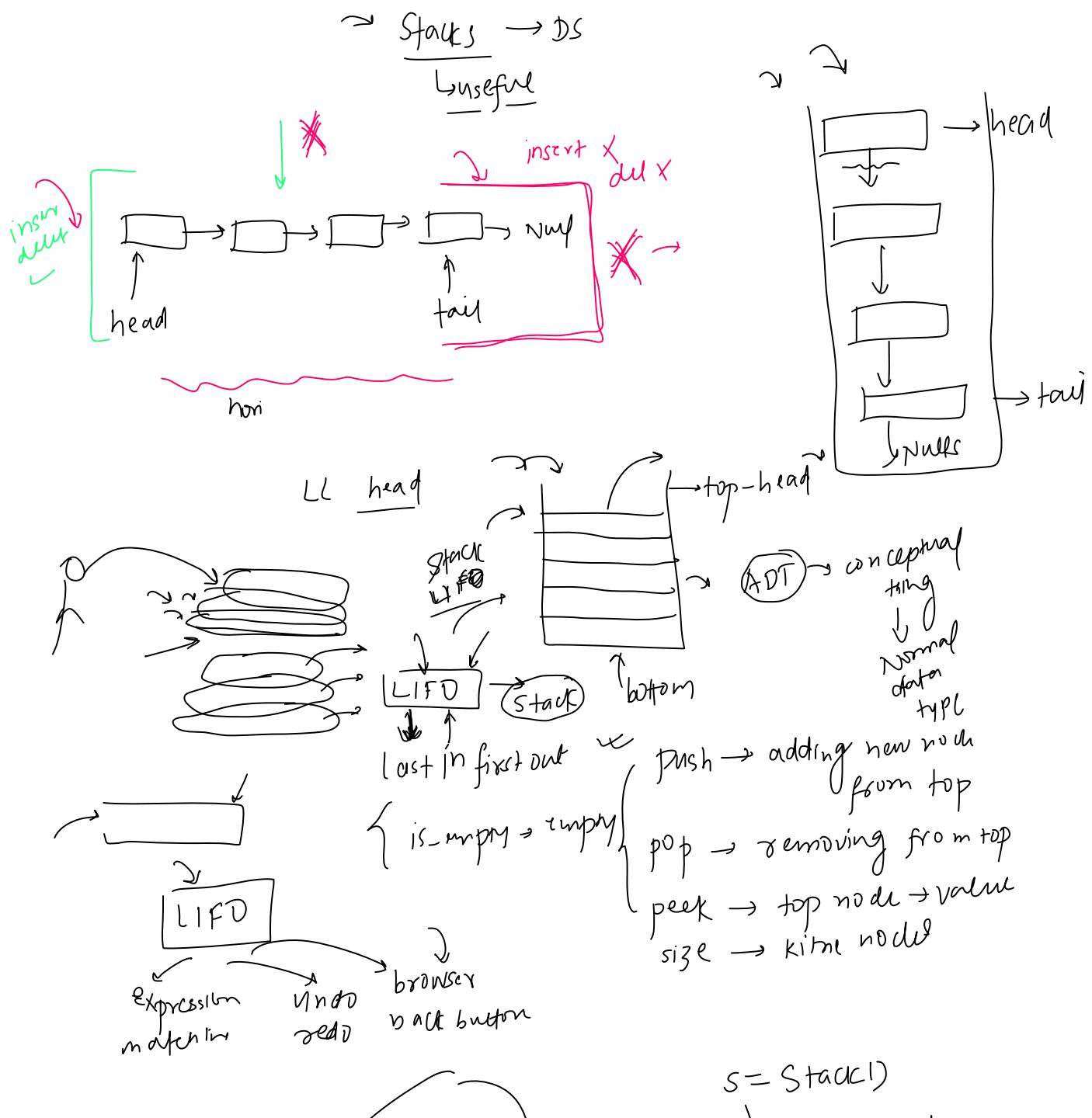
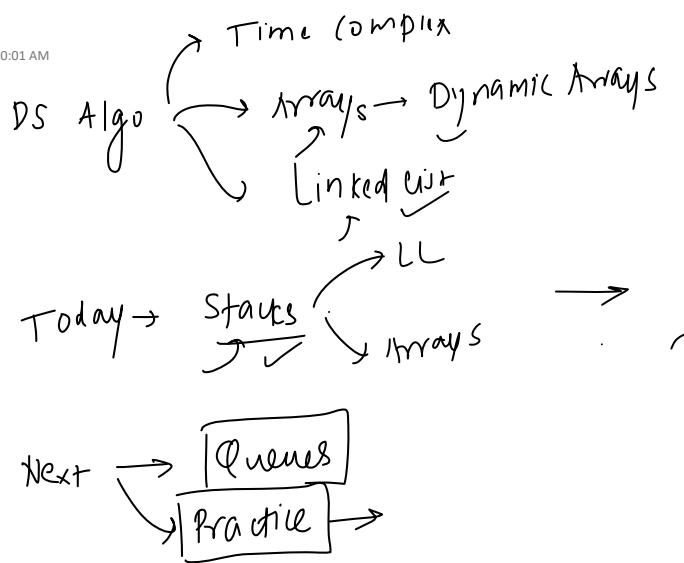
$$L_{NN} = \frac{\partial L}{\partial w_{11}} = d_1 \times d_1 \times \dots \times d_3$$

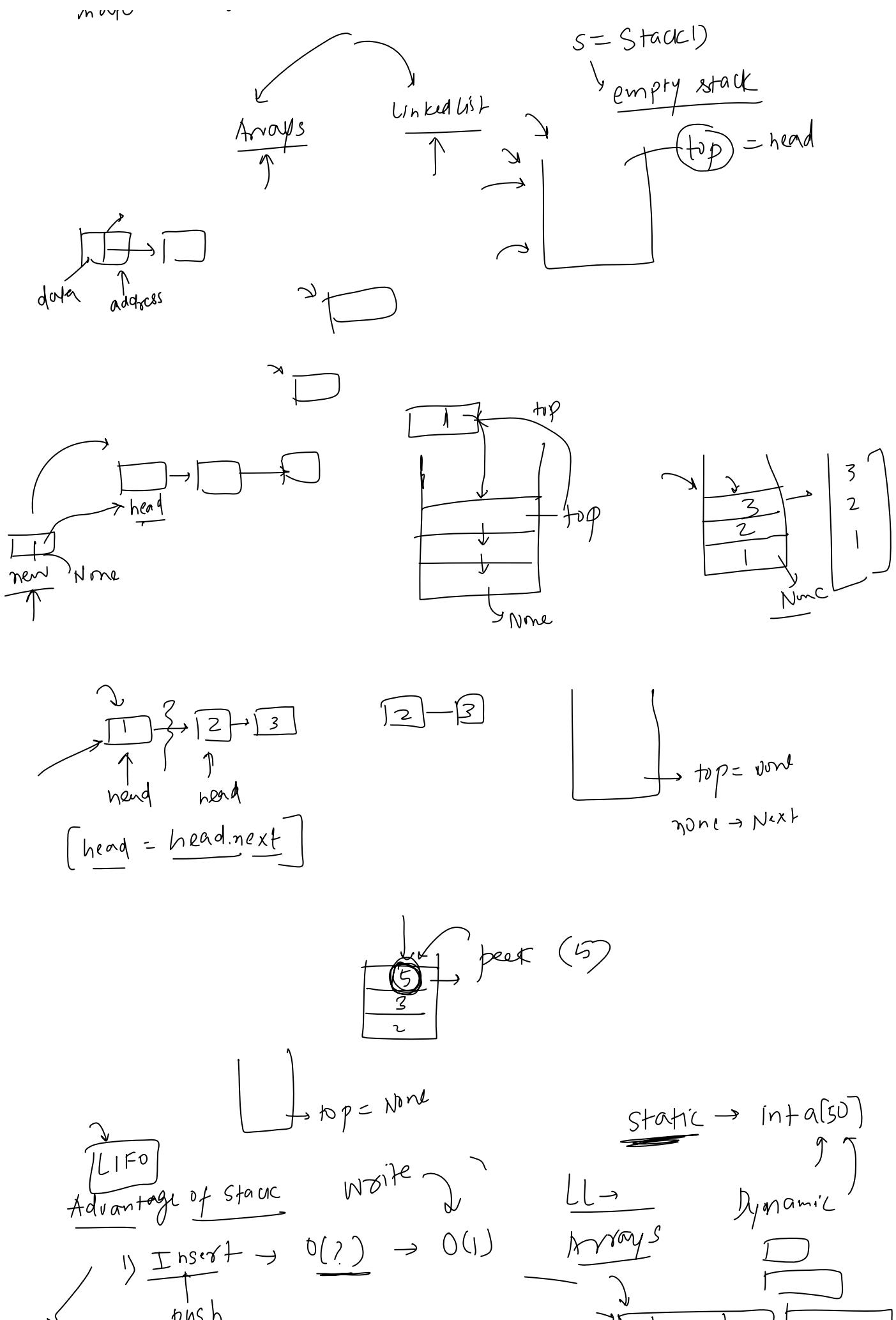
big number

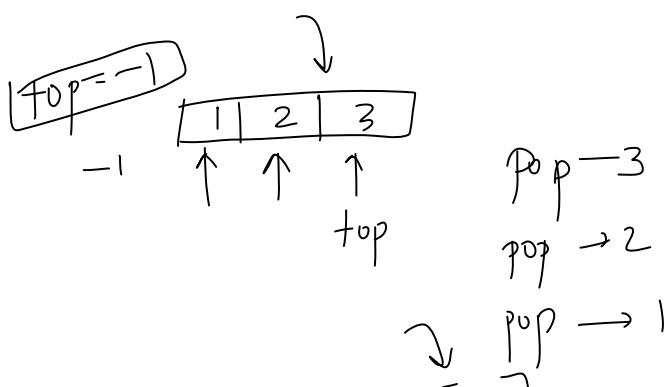
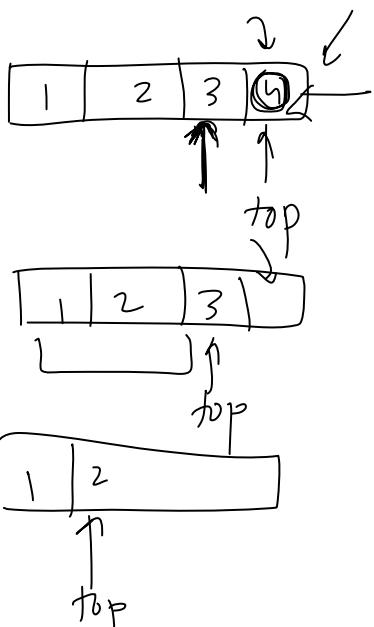
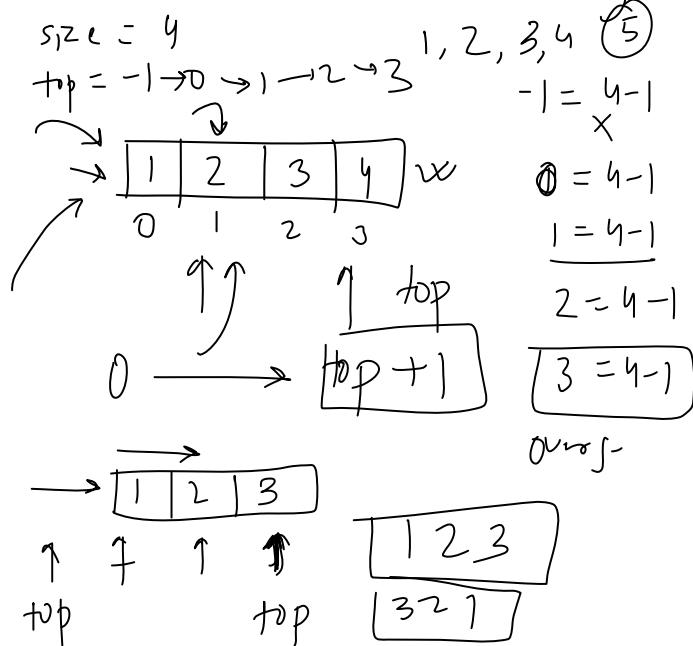
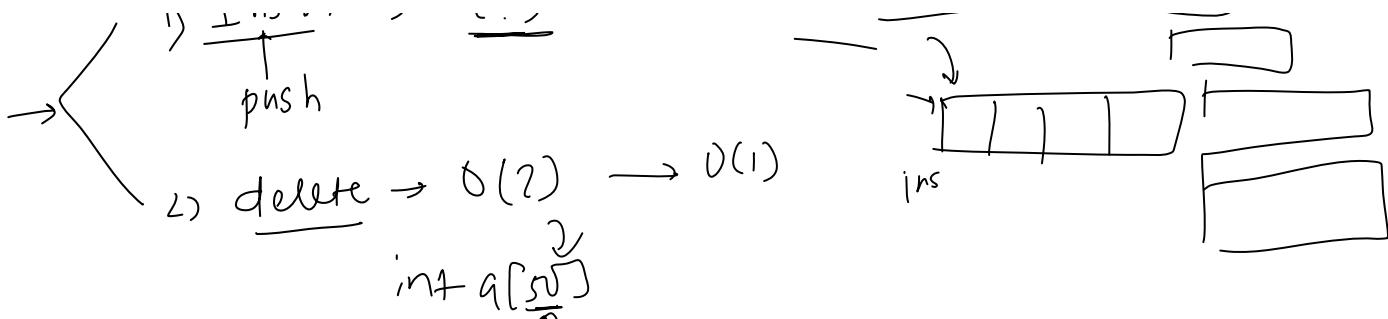
$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

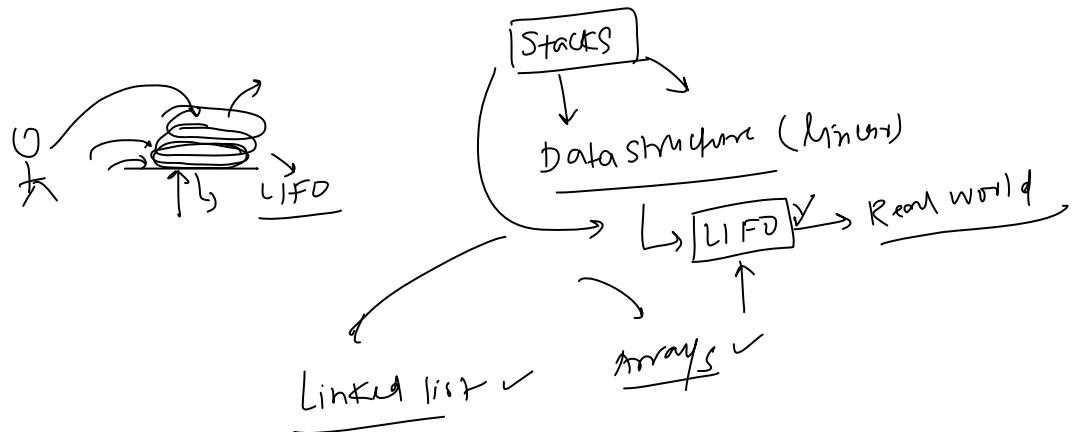
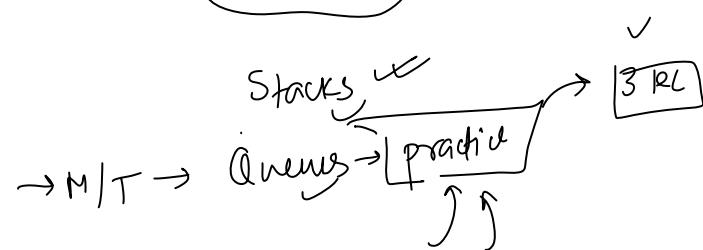
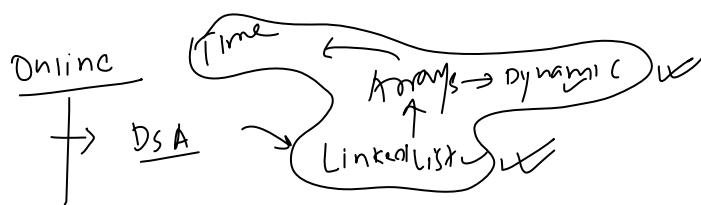
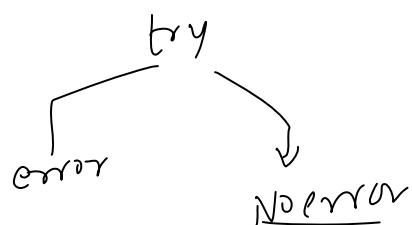
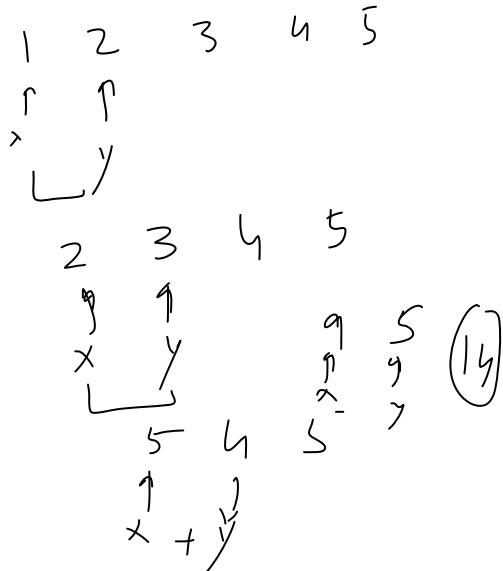
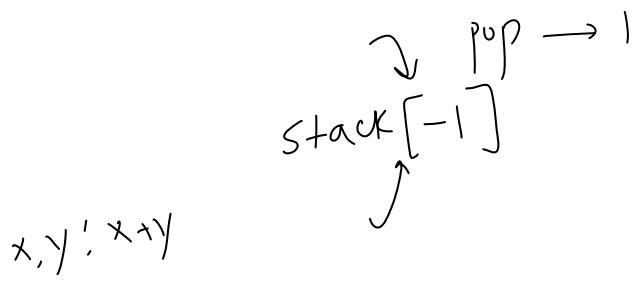
✓  
15)

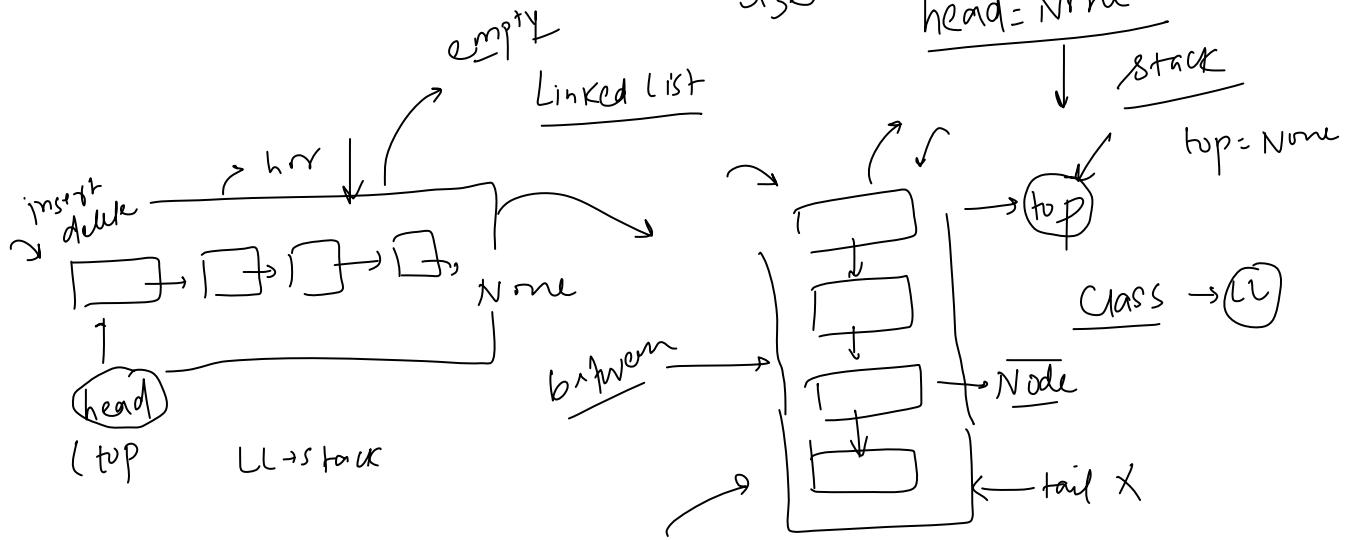
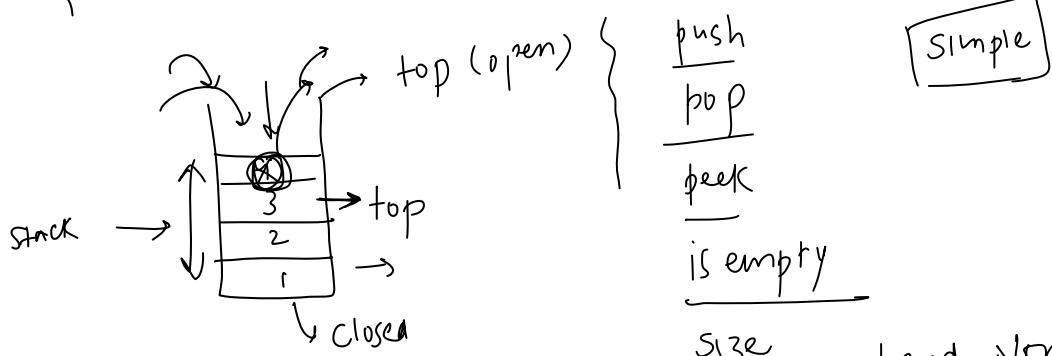












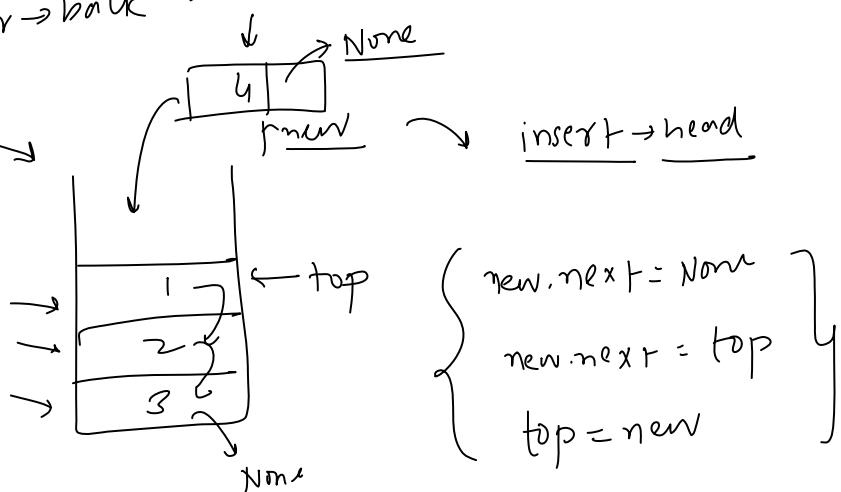
Application

- expression
- undo/redo
- returning
- browser → back

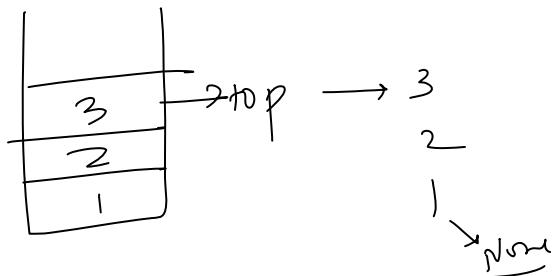
↓ LIFO

LL

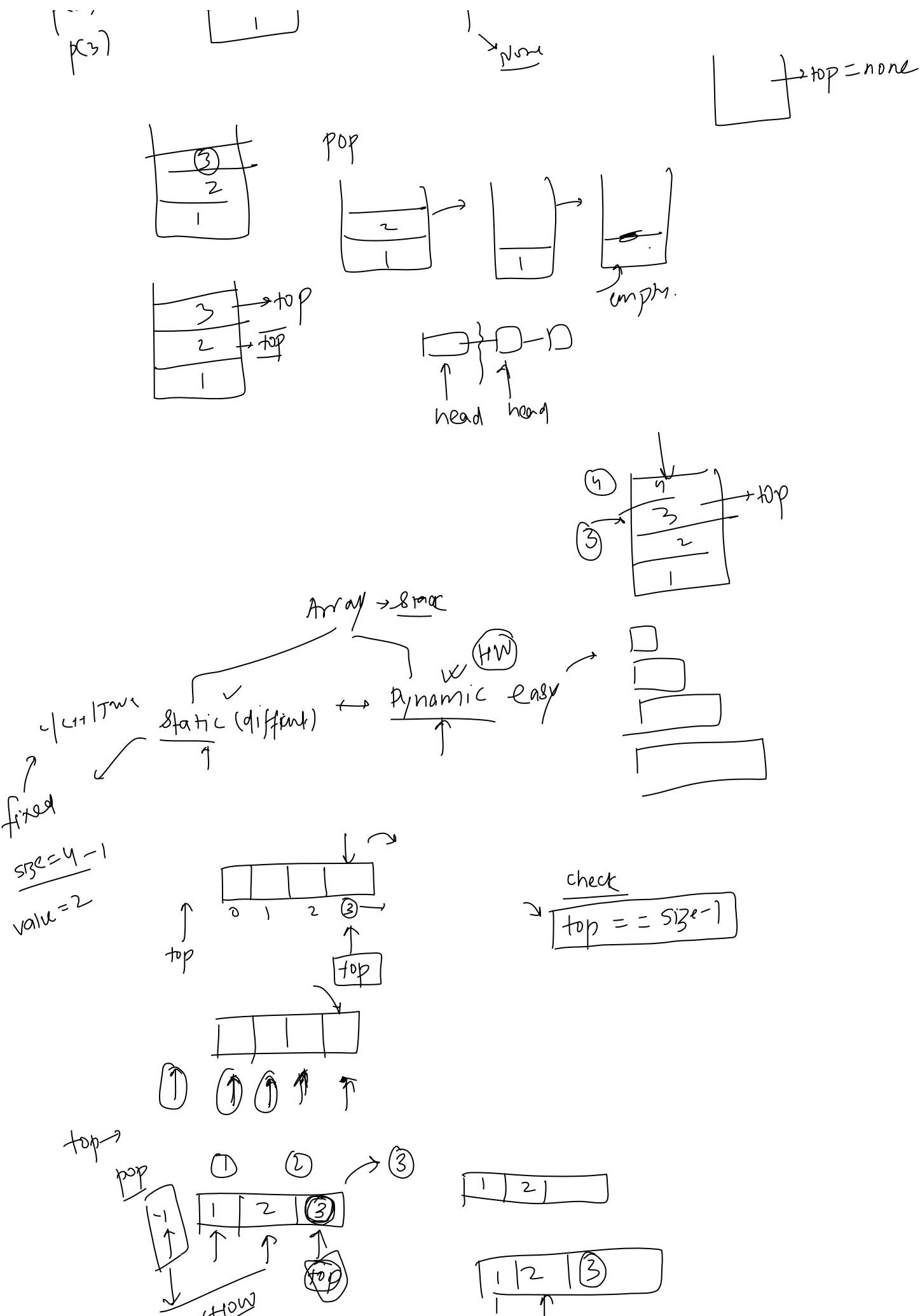
1 2 3

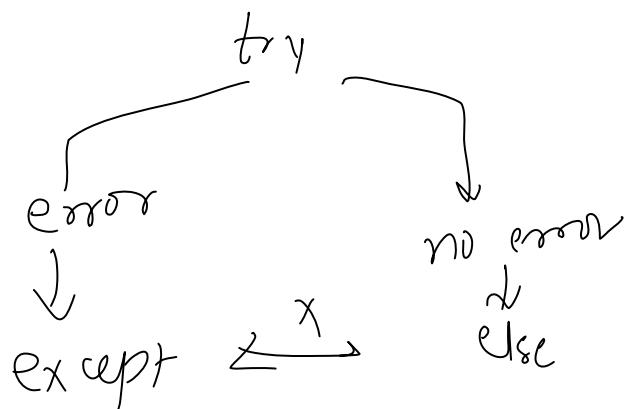
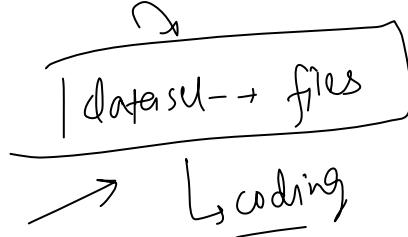
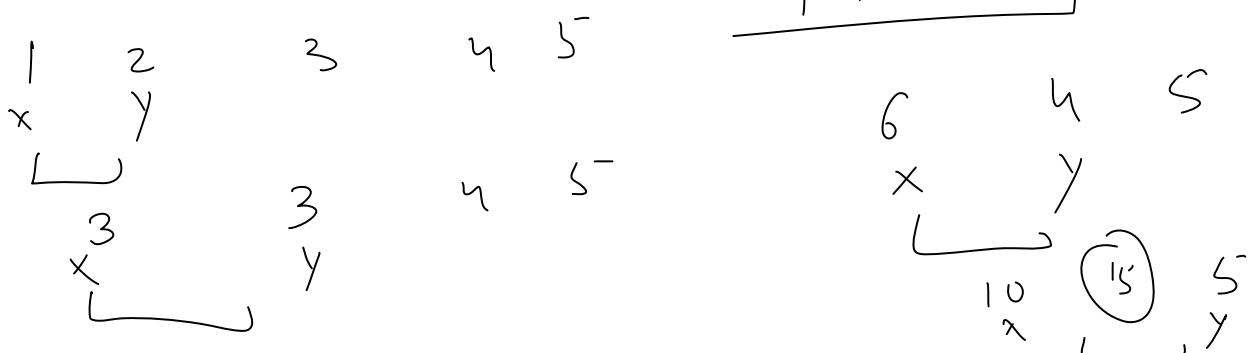
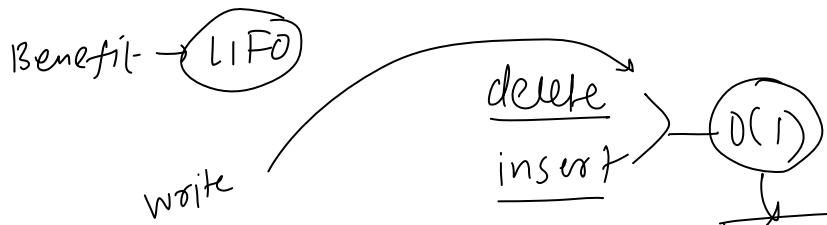
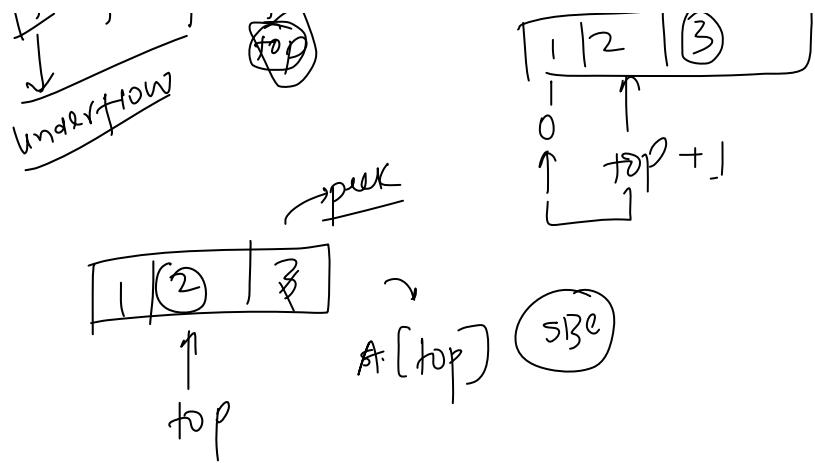


p(1)  
p(2)  
p(3)



1 ... n - nodd





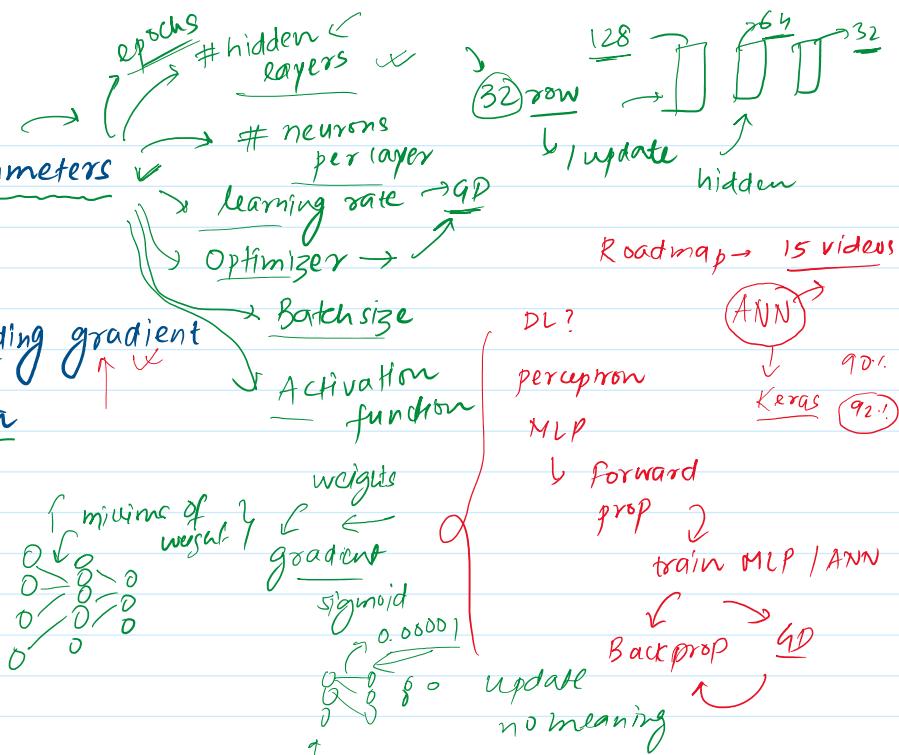
## How to improve a neural network ✓

29 April 2022 13:51

### 1. Fine tuning NN hyperparameters

### 2. By solving problems:

- ✓ → Vanishing / Exploding gradient
- ✓ → Not enough data
- Slow training
- Overfitting



### Fine tuning Hyperparameters

✓ No. of hidden layers

✓ No. of neurons per layer

Learning rate

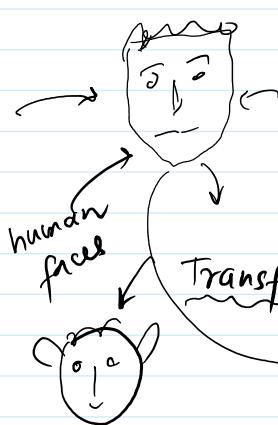
Optimizer

Epochs

Batchsize

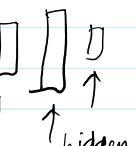
Activation function

### 1. No. of hidden layers



Representation

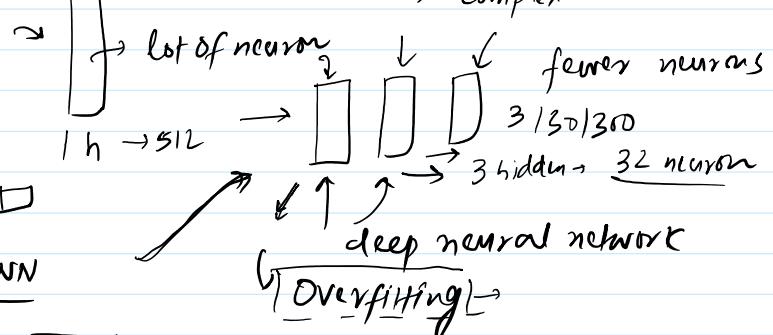
1 ↗ 0 ↘



1 hidden layer

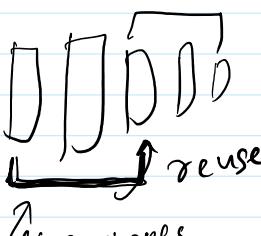
512 neurons

complex

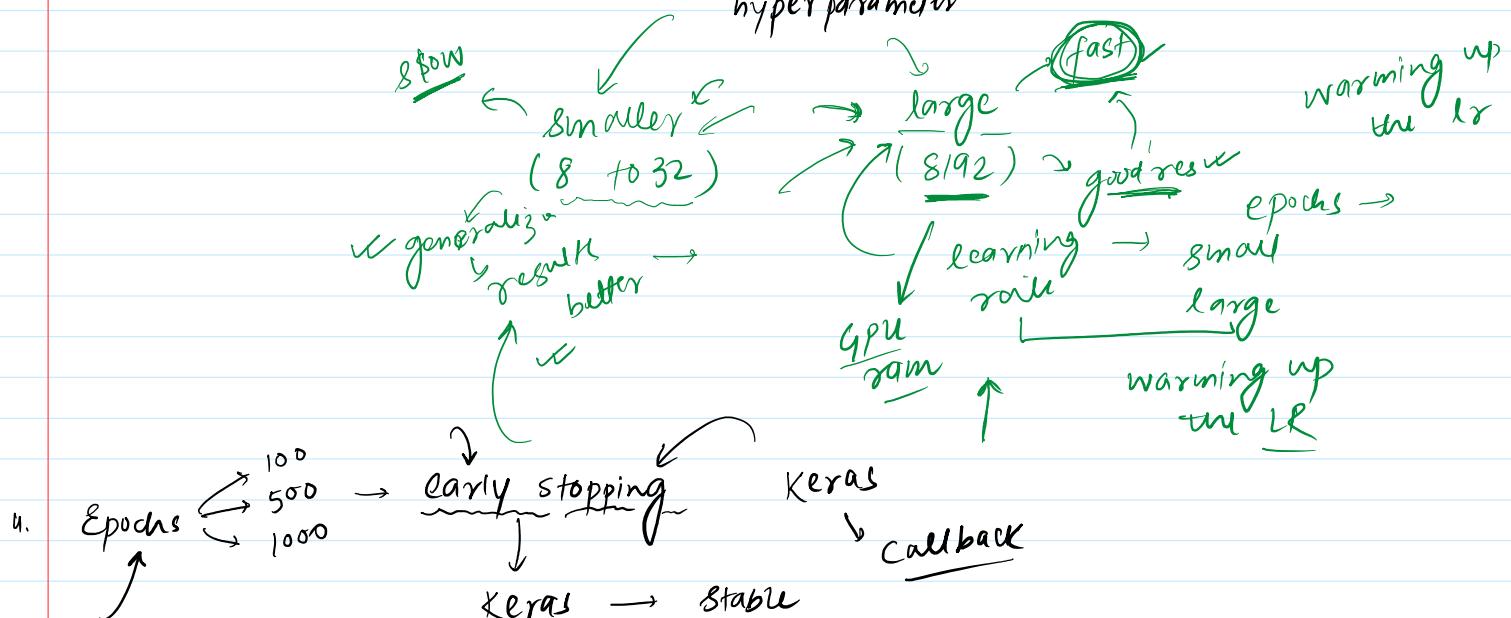
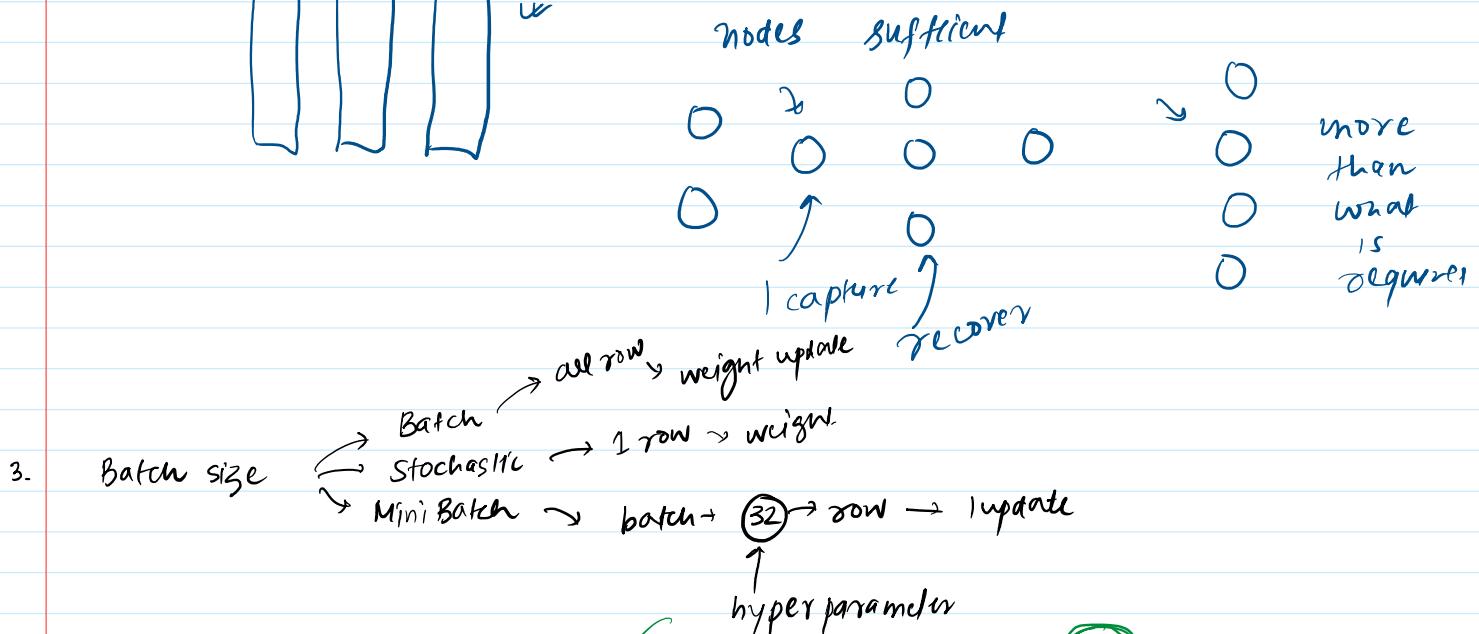
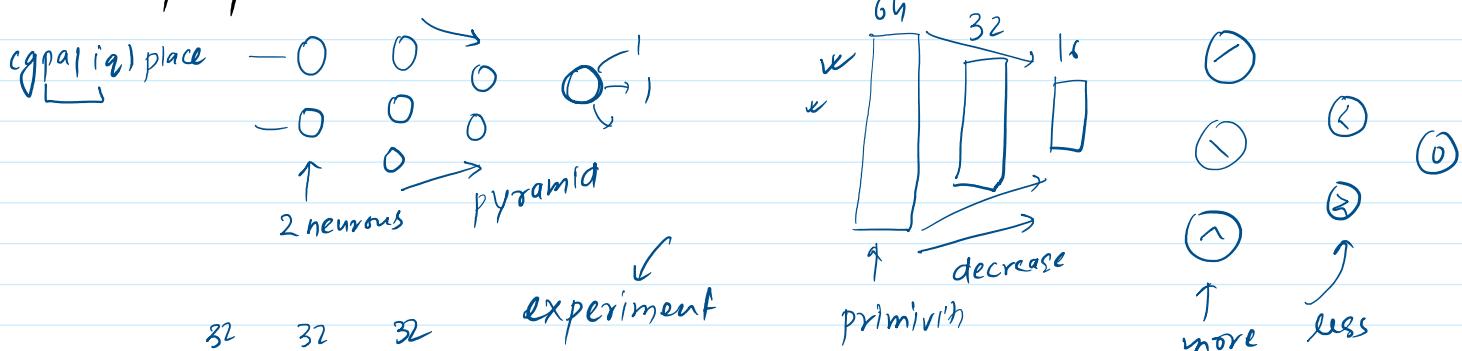


NN

reuse



## 2 Neuron / layer



## Problems with Neural Networks

Vanishing and exploding gradients

- Weight init →
- Activation function →
- Batch Norm →
- Gradient Clipping →

Not enough data

- Transfer learning
- Unsupervised pretraining

Slow training  
→ Adam

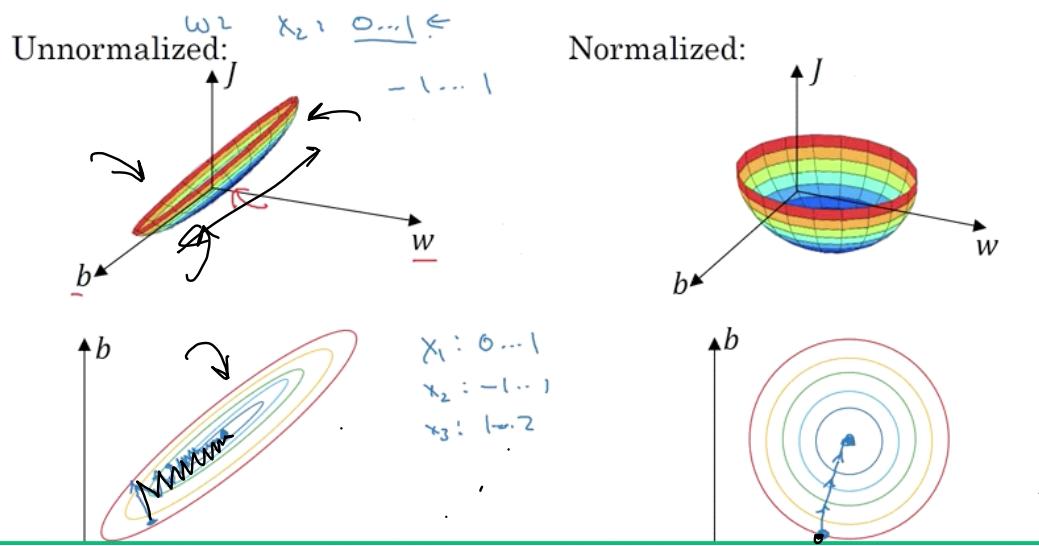
- Optimizers
- learning rate
- scheduler

Overfitting

- $l_1$  and  $l_2$  reg
- Dropouts

## Feature Scaling

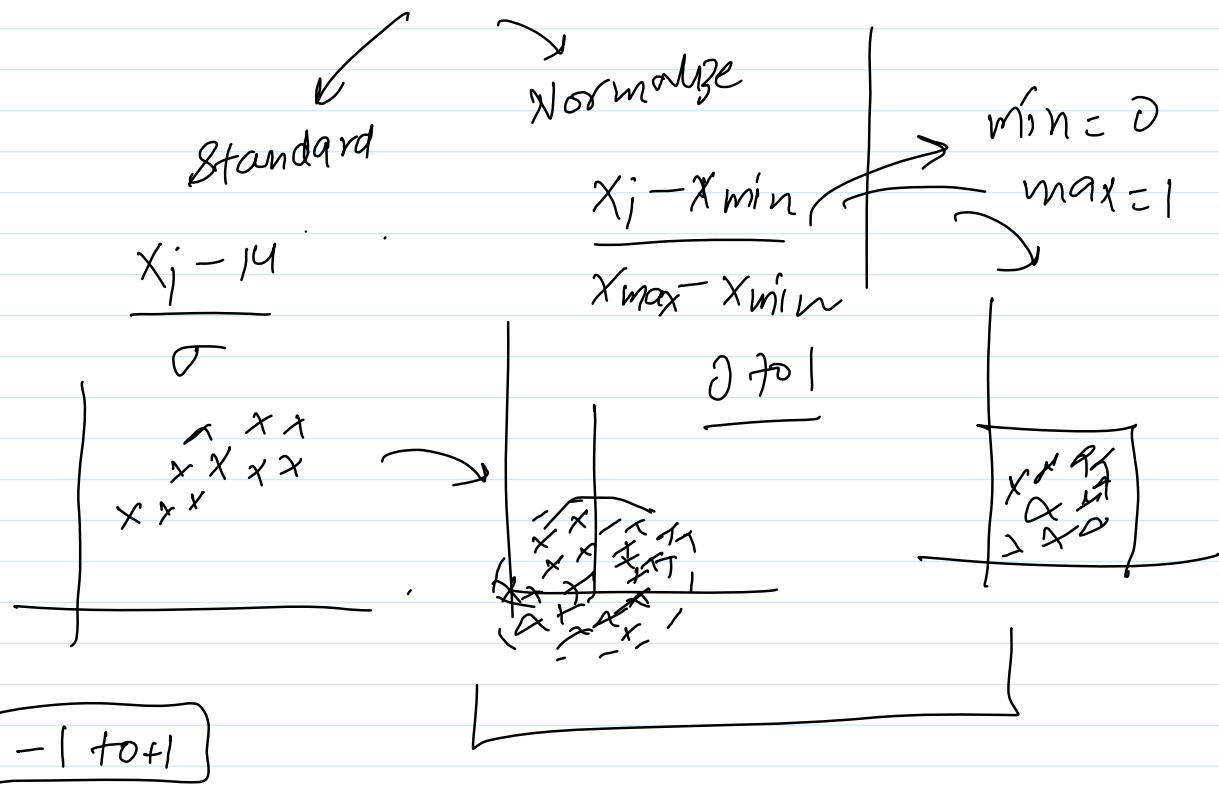
05 May 2022 17:37



(age)  $x_1 \rightarrow x_2$  (Salary)

$$w_1 \quad w_2$$

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



$$\begin{matrix} 2 \times 2 & 2 \times 3 \end{matrix}$$

$$\begin{array}{c} \boxed{2 \times 2} \quad \boxed{2 \times 3} \\ \left[ \begin{array}{cc} 1 & 2 \\ \cancel{3} & \cancel{4} \\ \hline \cancel{3} & \cancel{4} \end{array} \right] \quad \left[ \begin{array}{c} 1 \\ \cancel{4} \\ \hline \end{array} \quad \begin{array}{c} 2 \\ 5 \\ \hline \end{array} \quad \begin{array}{c} 3 \\ 6 \\ \hline \end{array} \right] \end{array}$$

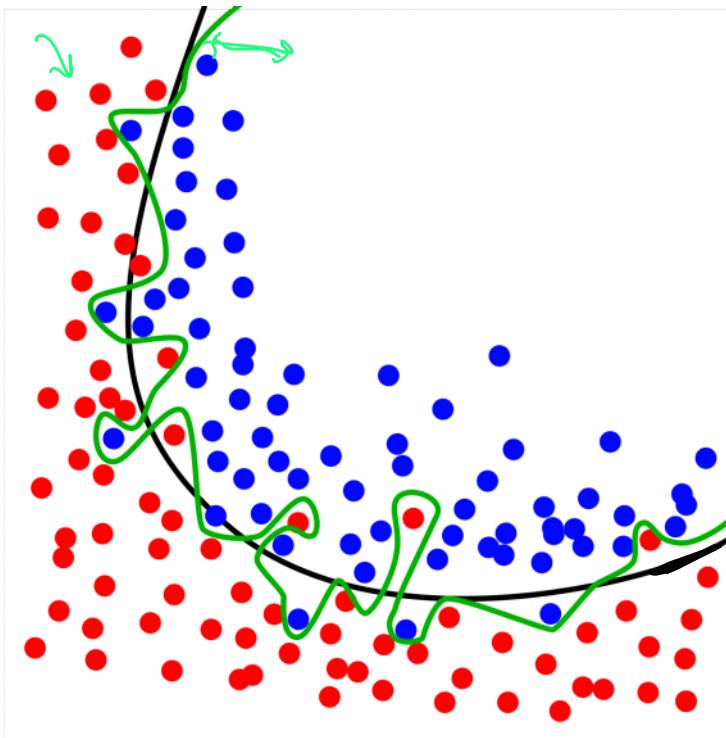
$$\left[ \begin{array}{ccc} 1+8 & 2+10 & 3+12 \\ 3+16 & 6+20 & 9+24 \end{array} \right]$$

$$\left[ \begin{array}{ccc} 9 & 12 & 15 \\ 19 & 26 & 33 \end{array} \right]$$

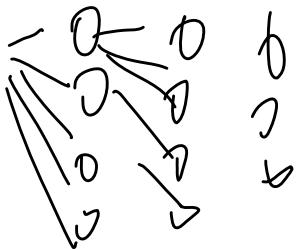
↑

# The Problem of Overfitting

12 May 2022 09:32



ANN overfitting  
Complex  
→ epochs



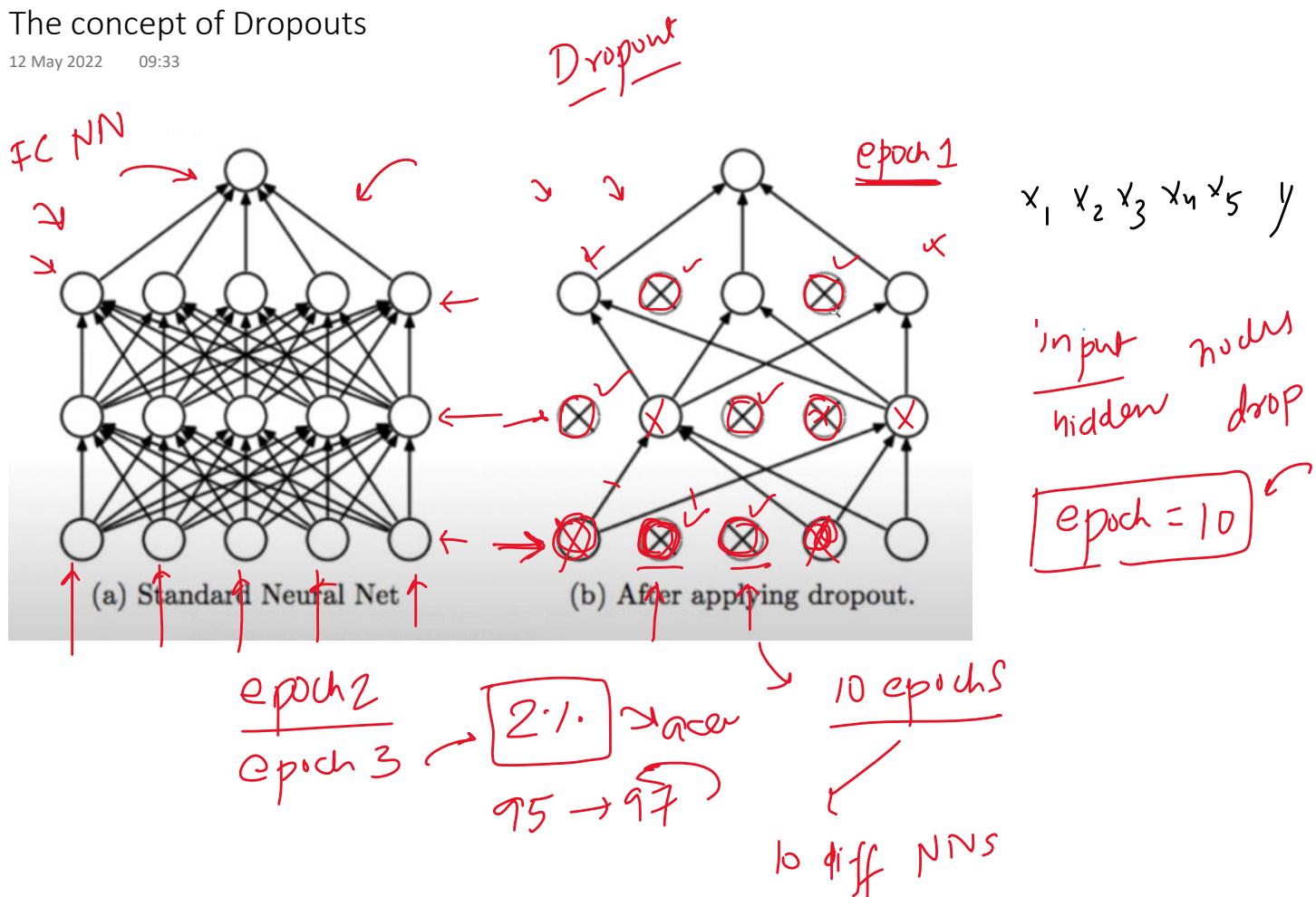
# Possible Solutions

12 May 2022 09:33

- 1) Add more data
- 2) Reduce complexity
- 3) Early stopping
- 4) Regularization  $\rightarrow L_1 \cup L_2$
- 5) Dropout

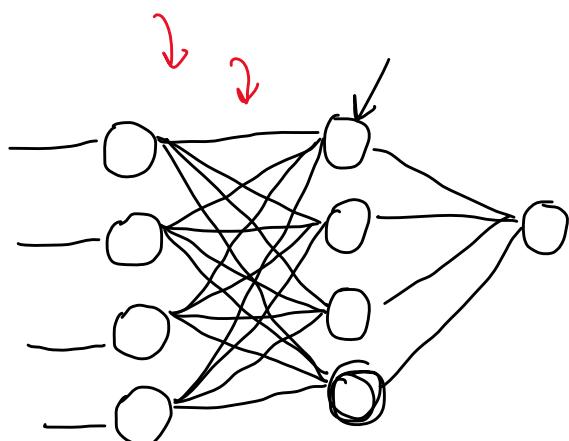
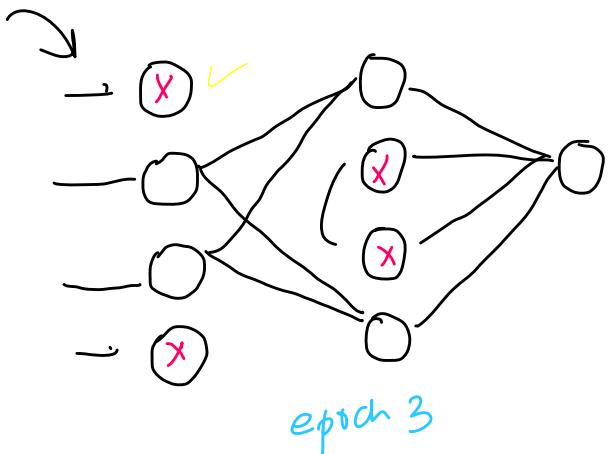
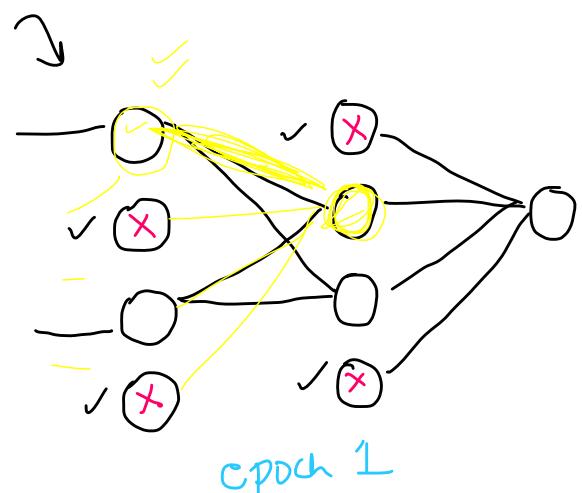
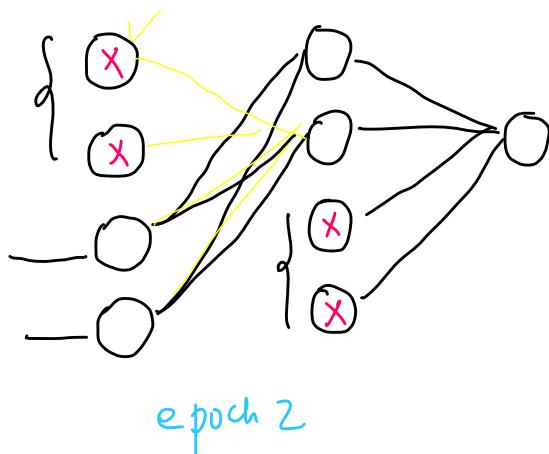
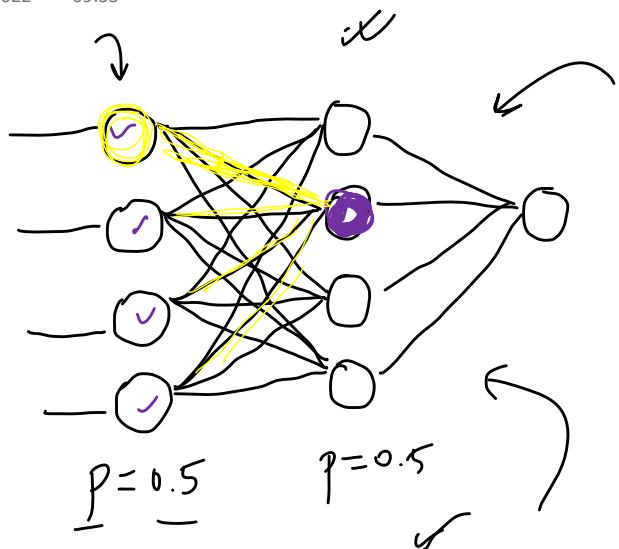
# The concept of Dropouts

12 May 2022 09:33



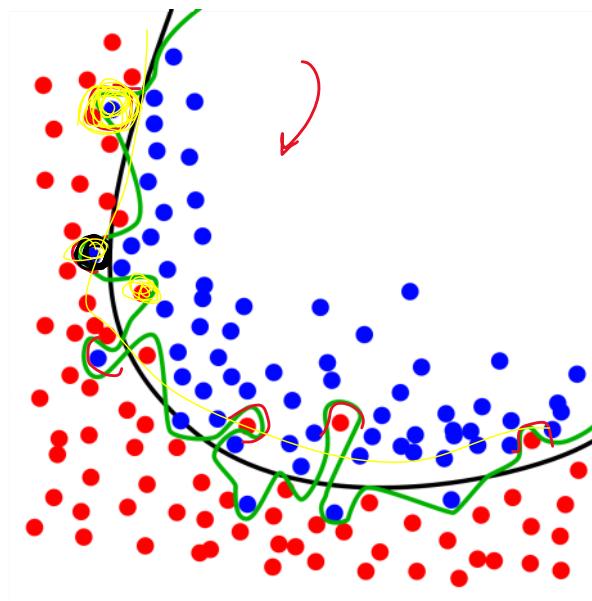
## Why this works?

12 May 2022 09:33



→ Reduce the # nodes

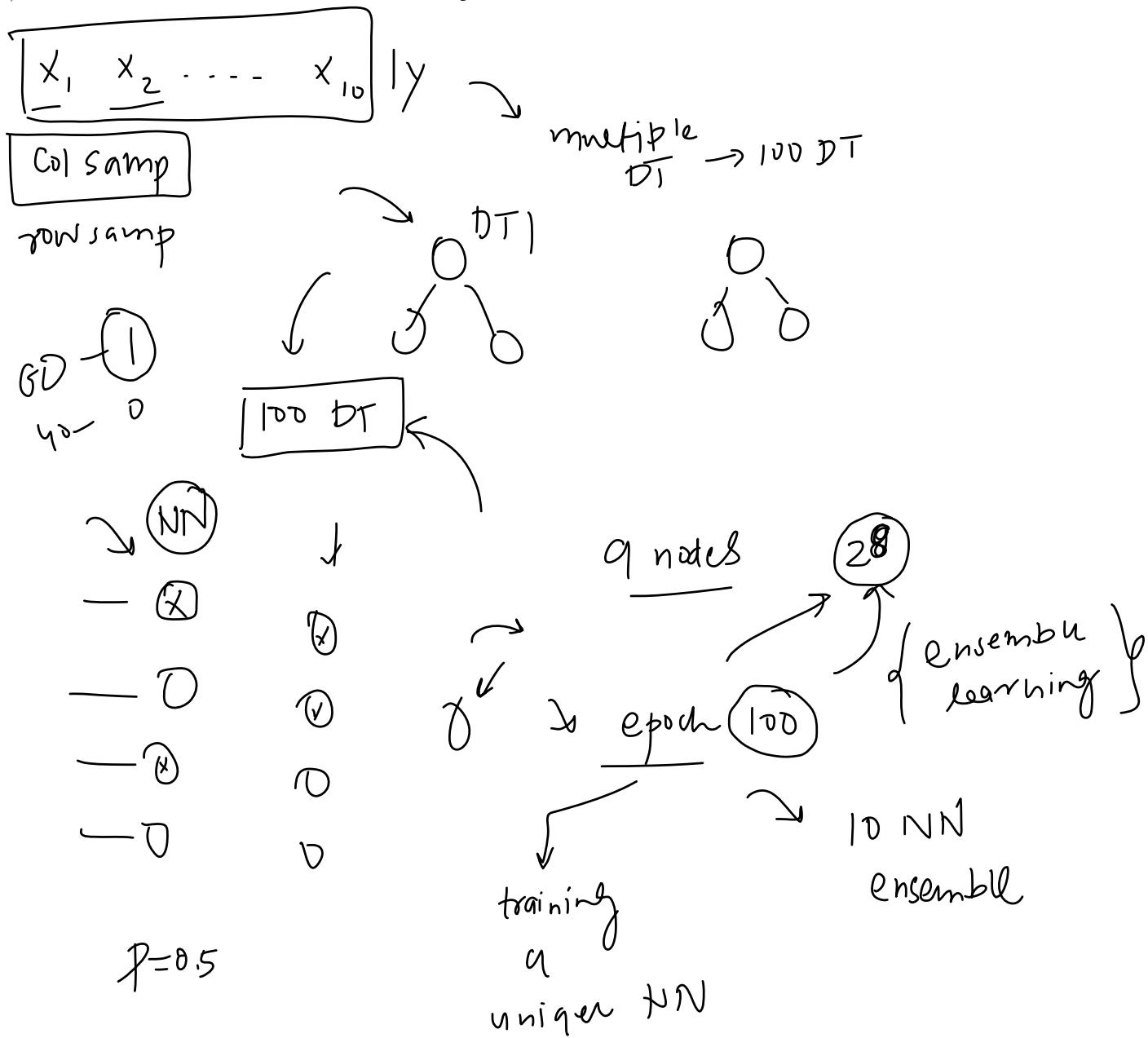
→



# Random Forest Analogy

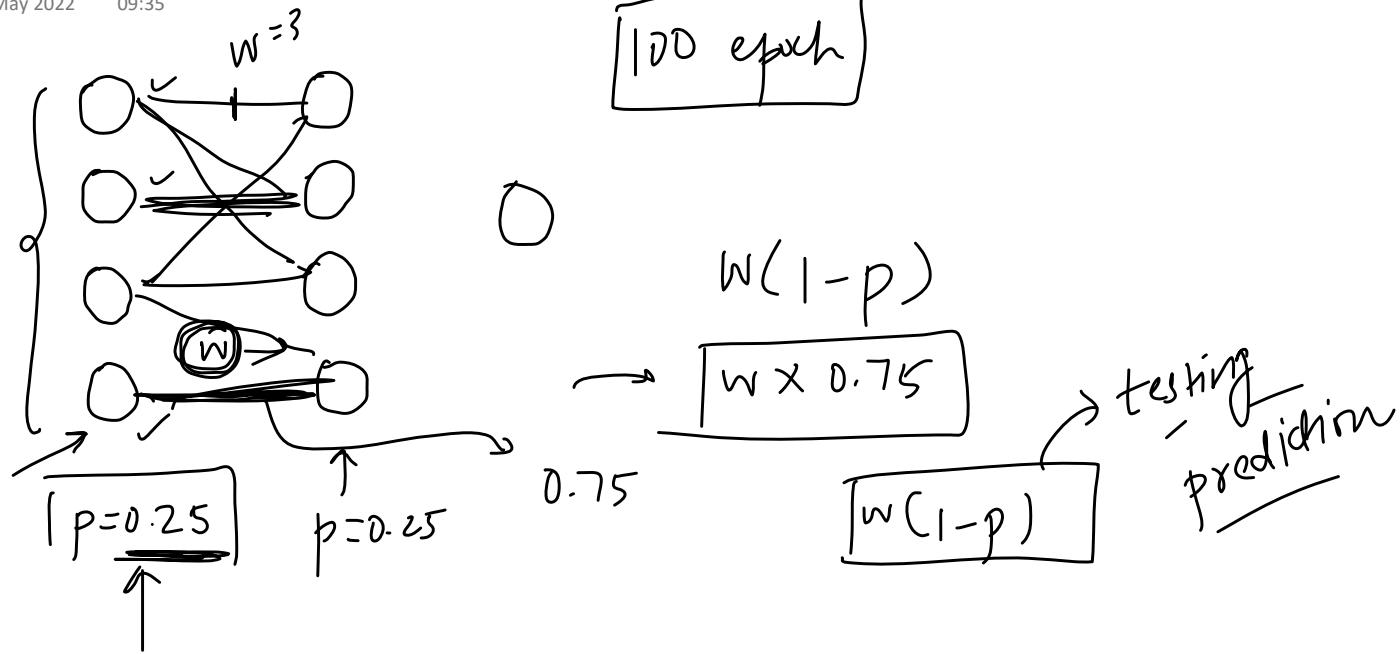
12 May 2022 09:33

50% w/w



## How prediction works?

12 May 2022 09:35



# Regression Code Example

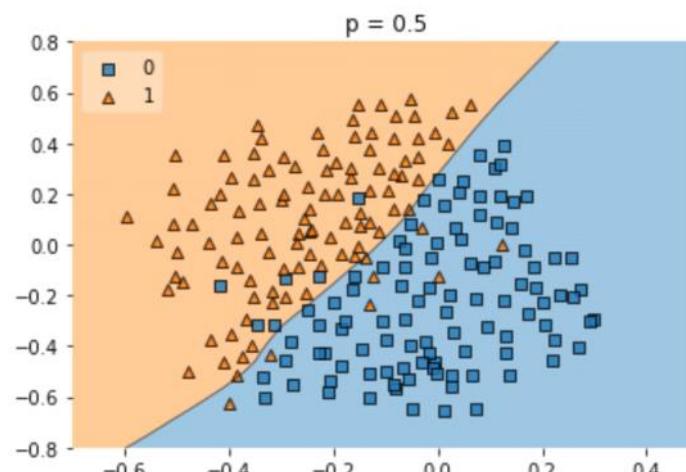
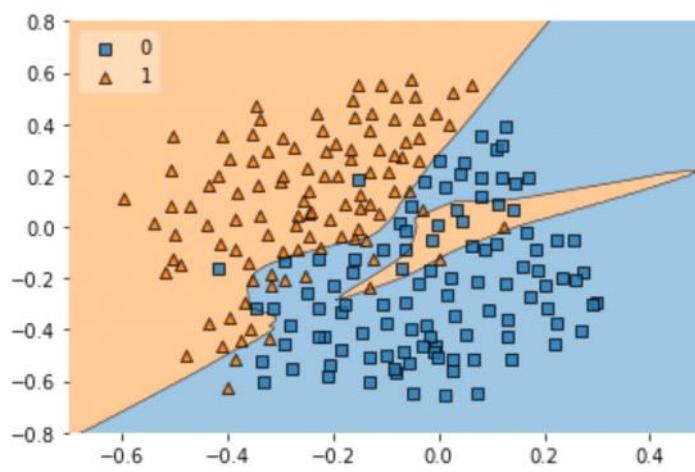
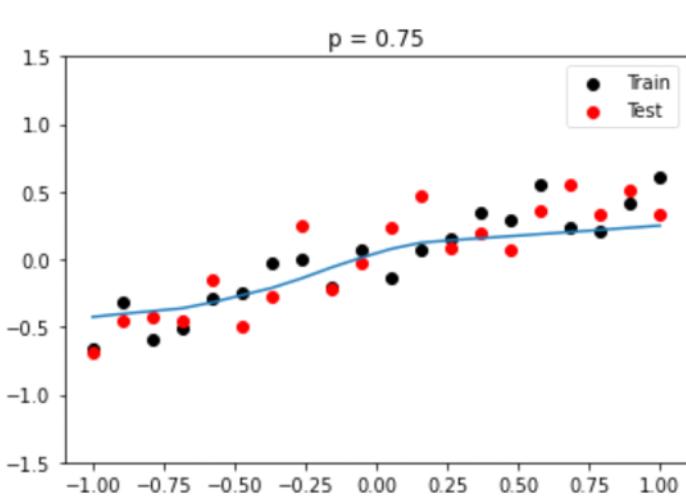
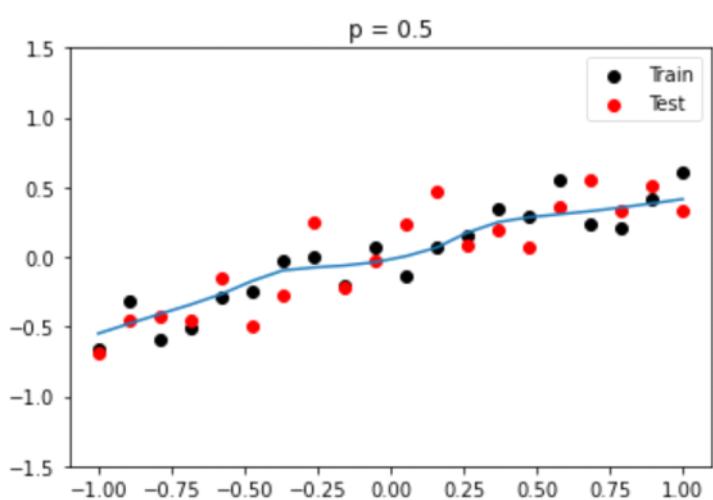
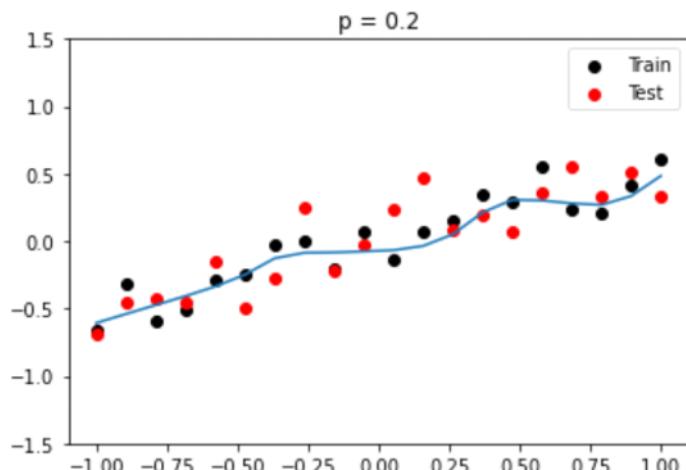
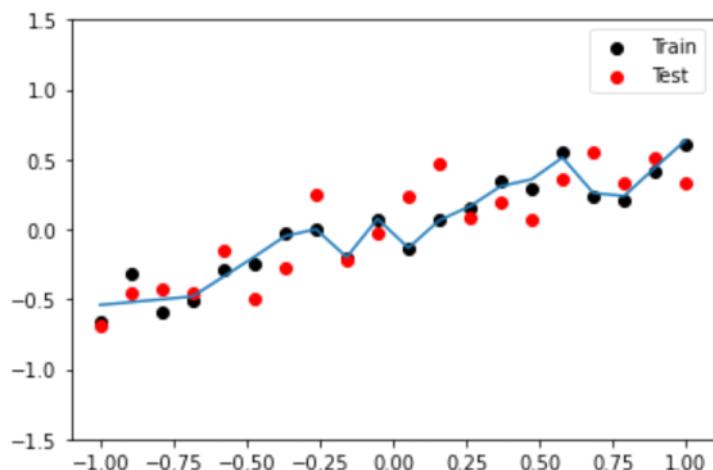
12 May 2022 09:34

# Classification Code Example

12 May 2022 09:34

# Effect of p

12 May 2022 10:06

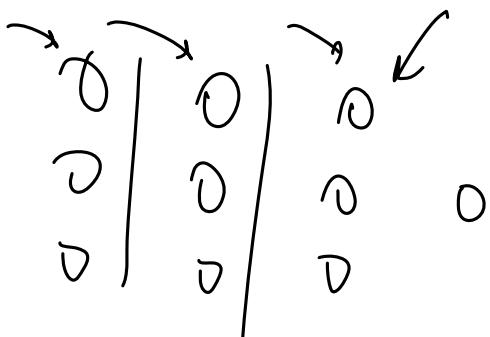


## Practical Tips and Tricks

12 May 2022 09:34

1) Overfitting  $P \uparrow$ , underfitting  $P \downarrow$

2) Last layer  $\rightarrow$  dropout



3) CNN  $\rightarrow$  40 - 50 - 1. (P)  $\rightarrow \checkmark$

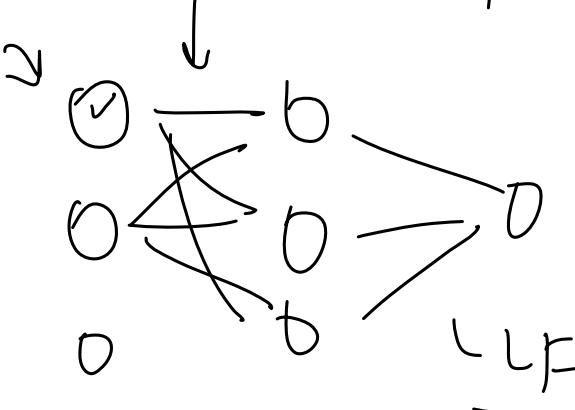
20 - 30 RNN

| 50 >

ANN  $\rightarrow$  [10 - 50]

## Drawbacks

12 May 2022 09:34

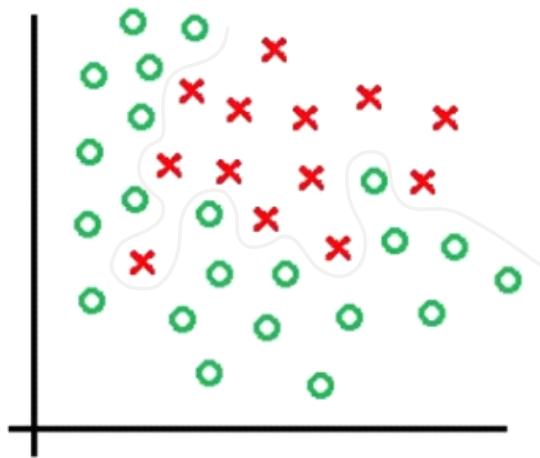
- 1) convergence  $\rightarrow$  delay
- 2) 
- { loss function }
- value changes
- LF - has to change

# Resources

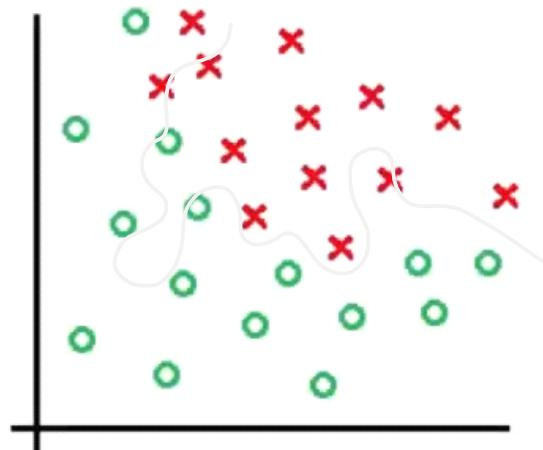
12 May 2022 12:38

# Overfitting

20 May 2022 08:14



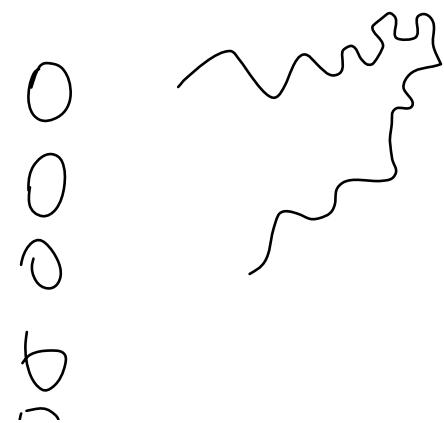
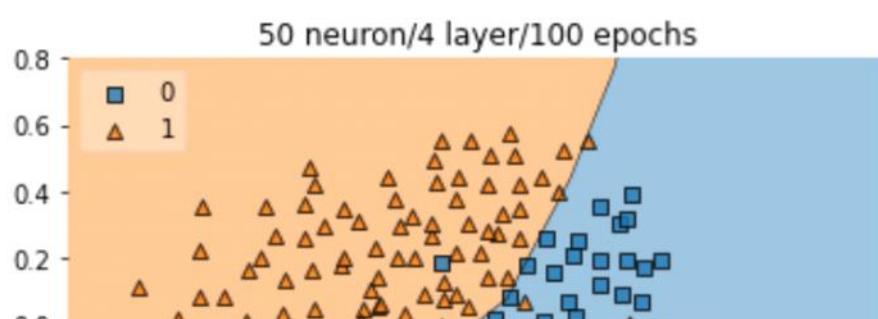
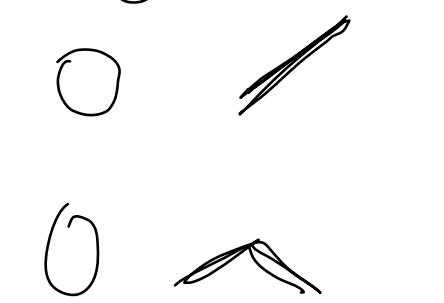
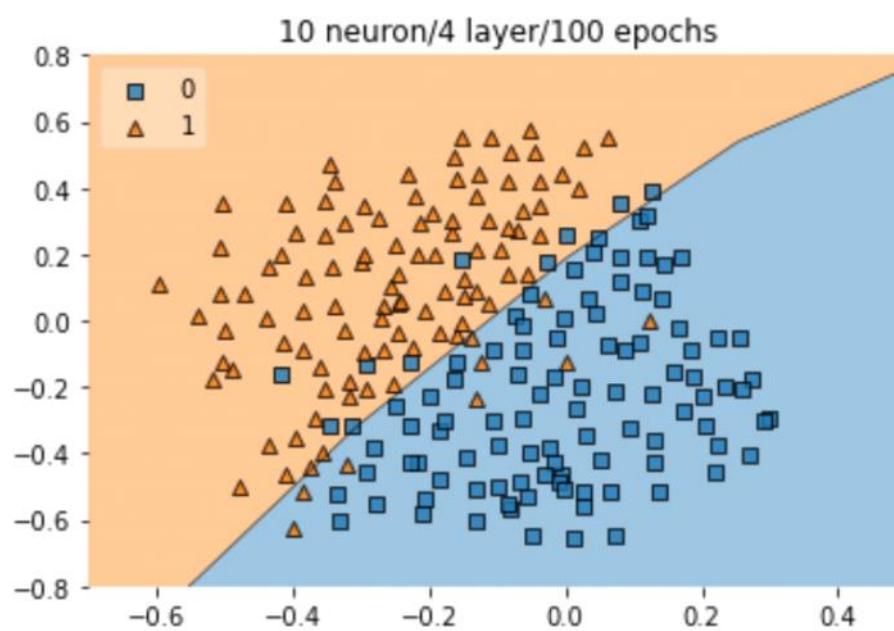
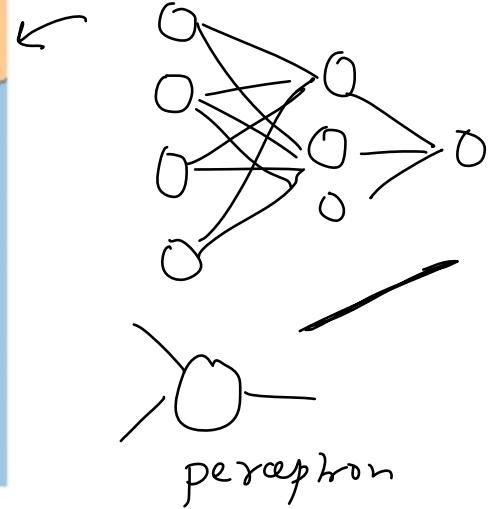
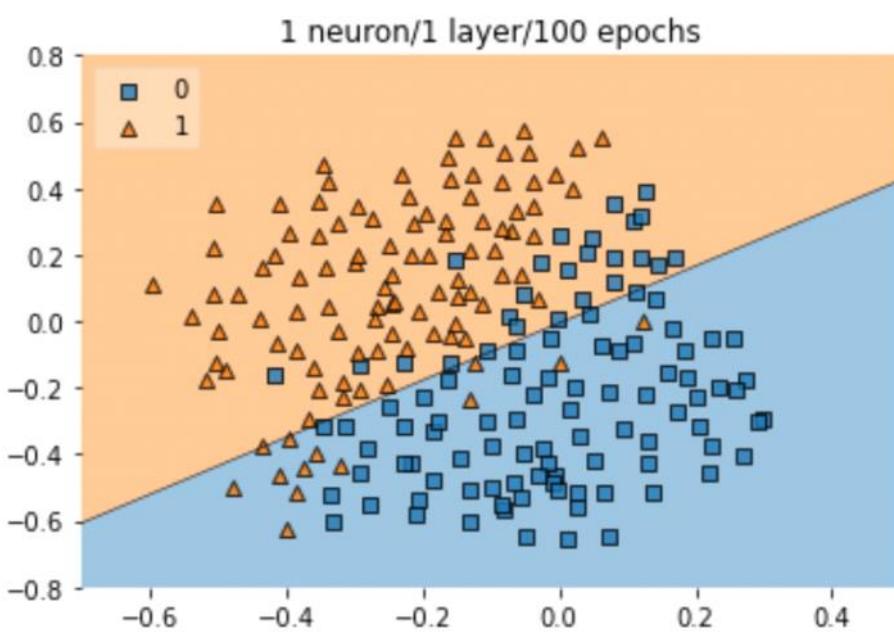
Train Dataset

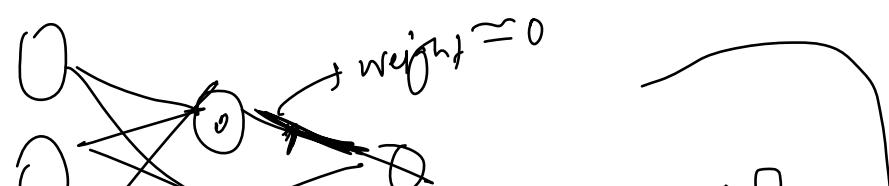
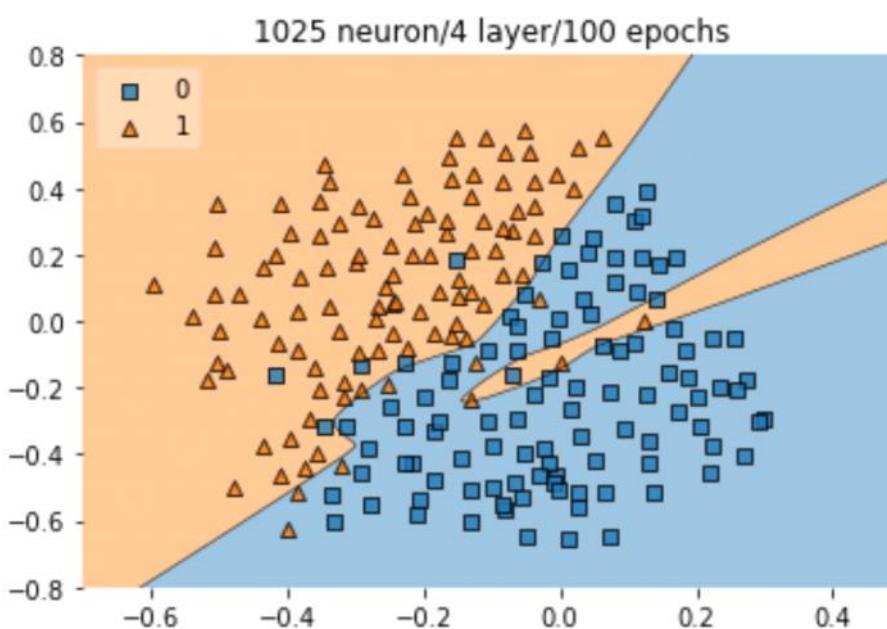
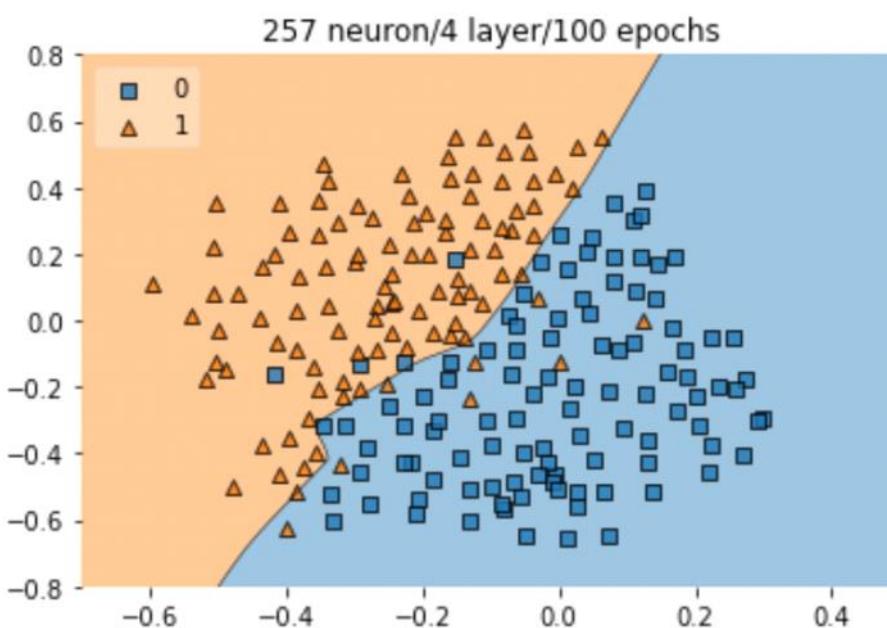
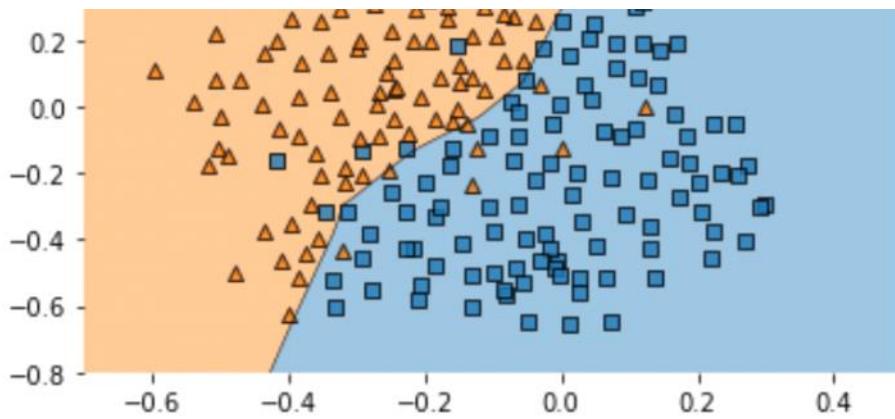


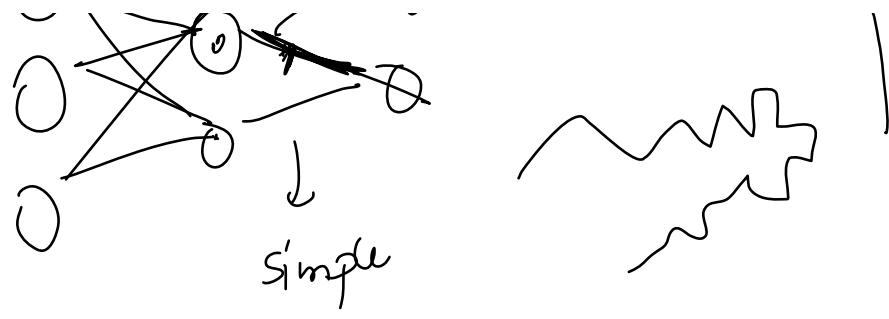
Test Dataset

# Why Neural Networks Overfit?

19 May 2022 07:34

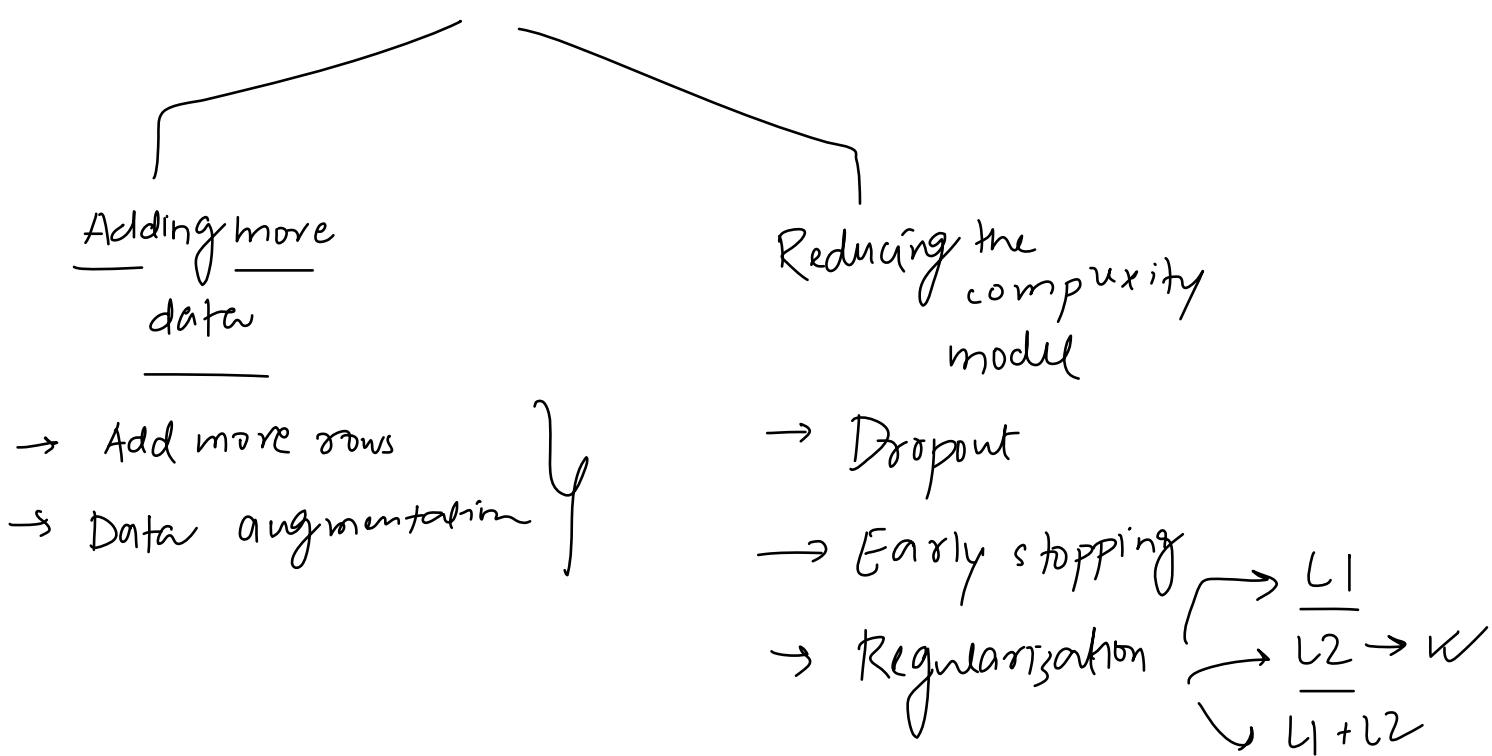






## Ways to solve overfitting

20 May 2022 08:13



## Regularization

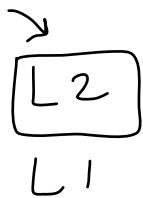
20 May 2022 08:14

A NN → weights | bias

↳ min Lossfunktion

$$L = \text{mse}$$

↓  
binary



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

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

$$w_1 \rightarrow w_{1,0}$$

$$\frac{\lambda}{2n} [w_1^2 + w_2^2 + \dots + w_{1,0}^2]$$

↑  $\lambda$  = hyperparameter ↑

$$\boxed{\lambda=0}$$

$$C = \sum L(y_i, \hat{y}_i) + P$$

$$L_2 \rightarrow L_1$$

$$w \approx 0$$

→  $L_1$  norm

$$C = L + \frac{\lambda}{2n} \sum \|w_i\|$$

$$C = \sum_{i=1}^n L(y_i, \hat{y}_i) + \left[ \sum_{l=1}^L \sum_{i=1}^n \sum_{j=1}^d \|w_{ij}^l\|^2 \right] \quad \boxed{w \approx 0}$$

$$\overrightarrow{0} \quad \textcircled{15} \quad \overleftarrow{n}$$

$$d, j = d$$

$$\begin{array}{c}
 \text{Diagram showing layers of neurons: } L=1, L=2, L=3 \\
 \text{Number of neurons in each layer: } 3, 3, 1 \\
 \text{Bias term: } \boxed{\text{bias } X}
 \end{array}$$

$$\psi, J = w_1^2 + w_2^2 + w$$

$$\sum_{i=1}^k \|b_i\|^2$$

# Intuition behind Regularization

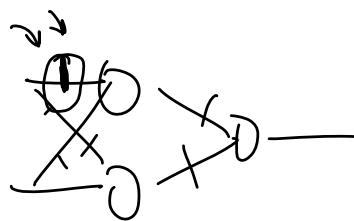
20 May 2022 08:15

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

$|1-\eta\lambda|$  positive

$$L' = L + \frac{\lambda}{2} \sum \|w_i\|^2$$

L



$$w_2^T [w_2^2 + w_3^2 + \dots]$$

zw

$$\frac{\partial L'}{\partial w_0} = \frac{\partial L}{\partial w_0} + \frac{\lambda}{2} 2w_0$$

$$w' = \frac{\partial L}{\partial w_0} + \lambda w_0$$

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

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

$$w_n = (1 - \eta \lambda) w_0 - \eta \frac{\partial L}{\partial w_0}$$

$$w_0 \ll w_n$$

L2 reg  $\rightarrow$  weight decay

weight decay  $\uparrow$   $w_0$

# Code Demo

20 May 2022 08:15

# Further Resources

20 May 2022 08:15

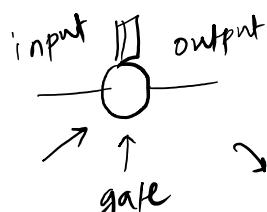
## What are Activation Functions?

31 May 2022 14:49

In artificial neural networks, each neuron forms a weighted sum of its inputs and passes the resulting scalar value through a function referred to as an activation function or transfer function. If a neuron has  $n$  inputs then the output or activation of a neuron is

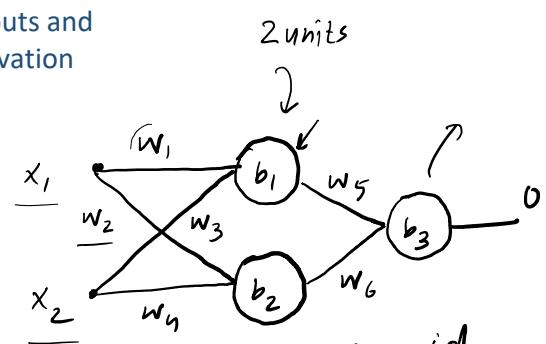
$$a = g(w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b)$$

This function  $g$  is referred to as the activation function.



activated  
(Y/N)  
how much

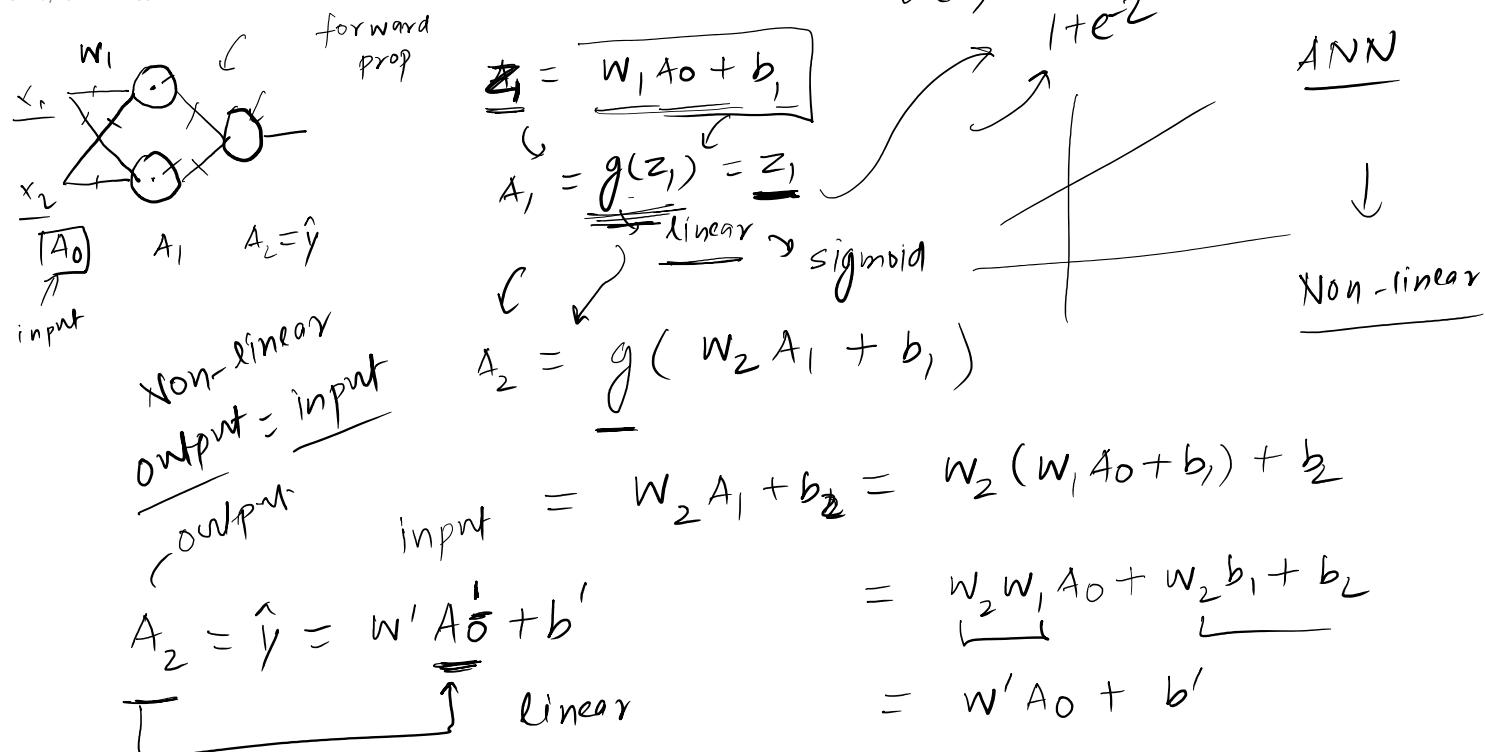
$$\underbrace{g}_{\text{activation}}(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$



sigmoid  
relu  
softmax

## Why Activation Functions are needed?

31 May 2022 14:50



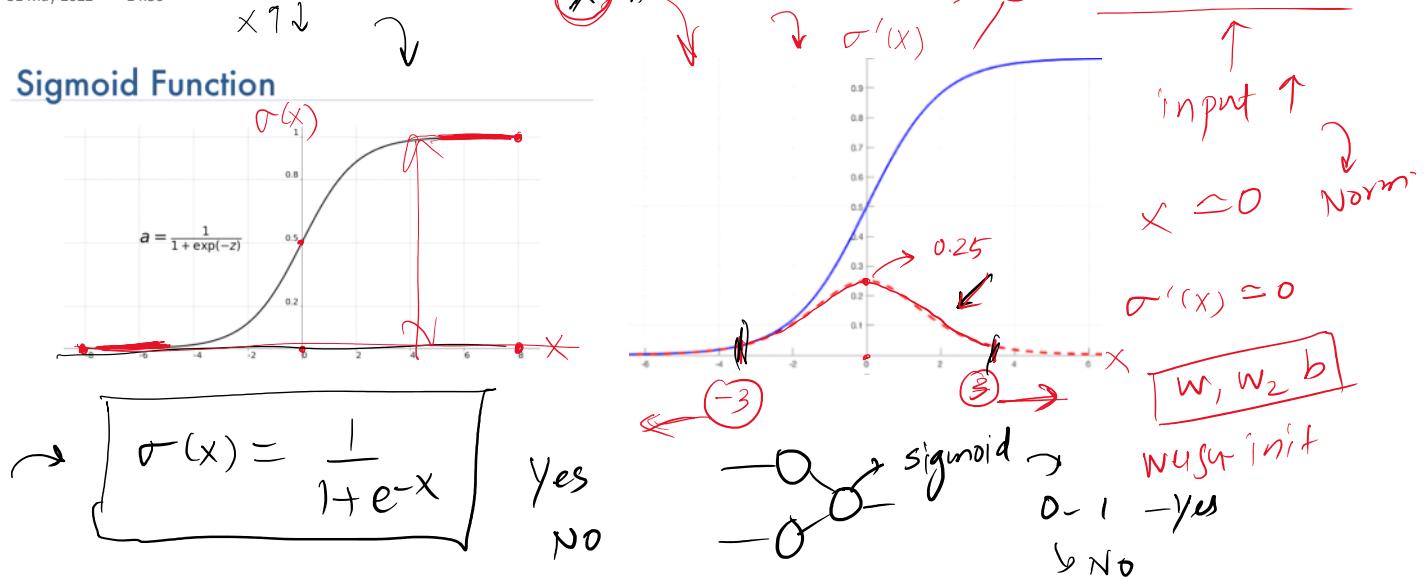
## Ideal Activation Function

31 May 2022 14:50

- 1) Non-linear ✓  $y = f(x)$  ✓ linear non-linear  $\sigma(z) = \frac{1}{1+e^{-z}}$
  - 2) Differentiable ✓ Universal Approx Theorem gradient des  $\rightarrow$  derivative  $\rightarrow$  ReLU (0)
  - 3) Computationally inexpensive ✓ derivative  $\rightarrow$  simple easy fast training slow
  - 4) Zero centered ✓ zero cent normalized mean=0 normalized data input later layers faster tanh
  - 5) Non-saturating ✓ Sigmoid  $\rightarrow [0, 1]$   $f(x) = \max(0, x)$   
tanh  $\rightarrow (-1, 1)$  ReLU
- update  $x$   $[0, 1]$  backprop saturating AF  $\rightarrow$  vanishing gradient problem  
 $\sim 0.0001$   $\leftarrow 0.0001$   $0.01$

## Sigmoid Activation Function

31 May 2022 14:50

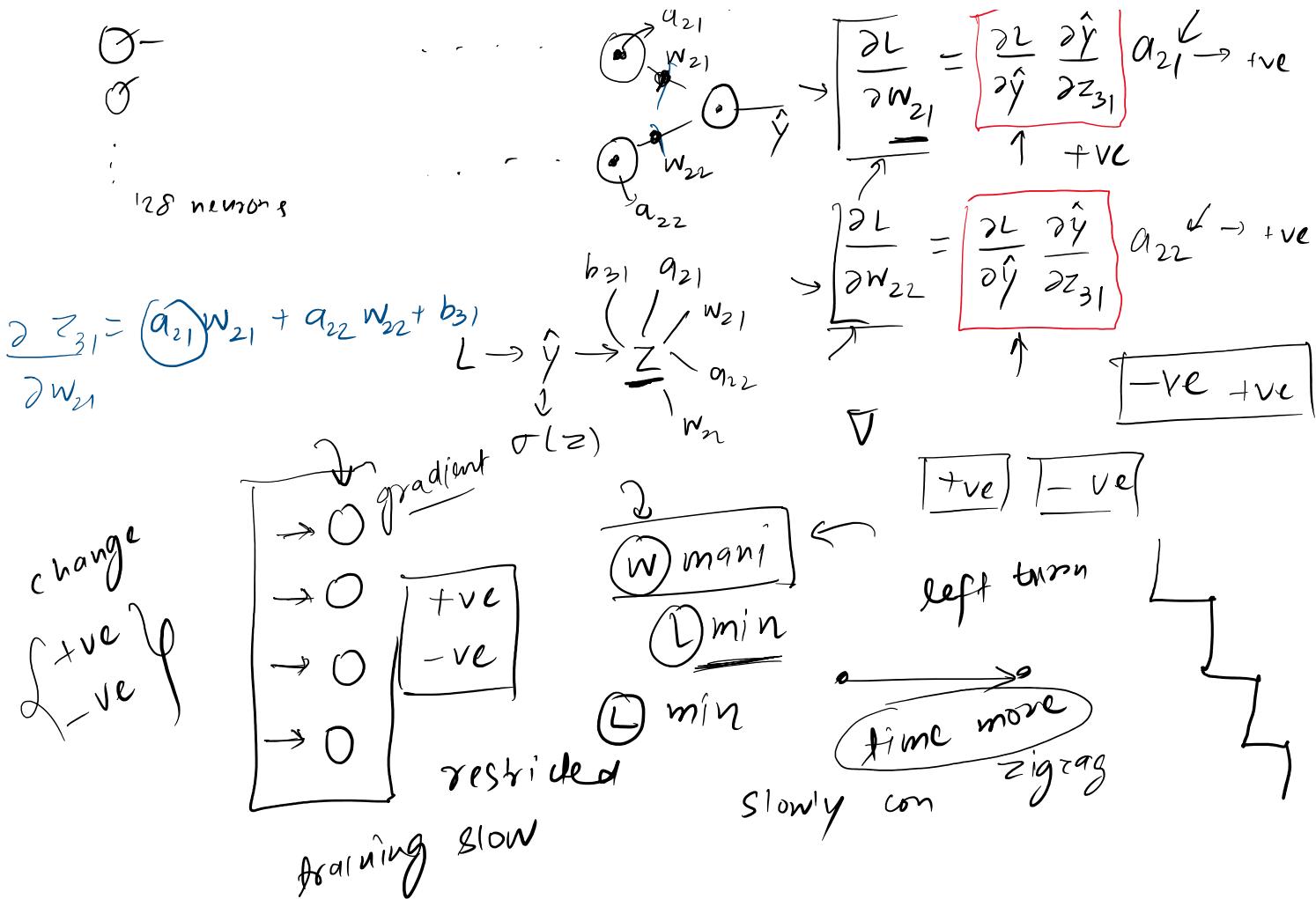


### Advantages

- 1)  $[0, 1] \rightarrow \text{probability} \rightarrow \text{output layer} \rightarrow \boxed{\text{Binary classification}}$
- 2)  $\text{Non-linear} \rightarrow \boxed{\text{Non-linear data}} \rightarrow \text{good option}$
- 3)  $\text{Differentiable} \rightarrow \text{Backprop} \rightarrow \boxed{\frac{\partial L}{\partial w}}$

### Disadvantages

- 1)  $\text{Saturating function} \rightarrow [-\infty, \infty] \rightarrow [0, 1]$ 
  - $\hookrightarrow \boxed{\text{Vanishing gradient problem}}$
  - $\downarrow$  Backprop → update
  - $w_n = w_0 - \eta \frac{\partial L}{\partial w}$
  - $w_n = w_0$
  - $\downarrow$  No update will take place
  - $\downarrow$  training X
- 2)  $\boxed{\text{Non Zero centered}}$ 
  - $\circlearrowleft$  Normalized
  - $\circlearrowleft$  +ve
  - $\circlearrowleft$  -ve
  - $\rightarrow \boxed{\text{Training slow}}$
  - $\downarrow$   $a_{21}$
  - $\downarrow$   $w_{21}$
  - $\downarrow$   $\frac{\partial L}{\partial w_{21}} = \begin{bmatrix} \frac{\partial L}{\partial v} & \frac{\partial L}{\partial z_1} \end{bmatrix} a_{21} \rightarrow \text{+ve}$



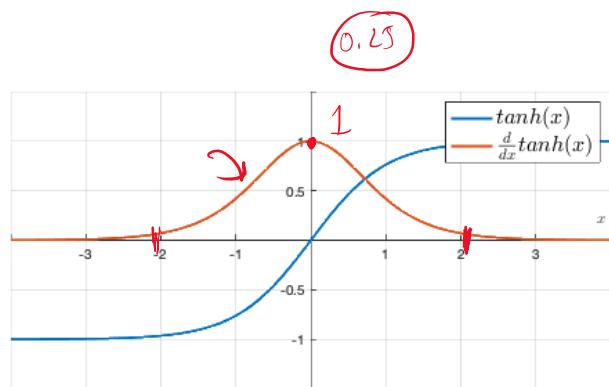
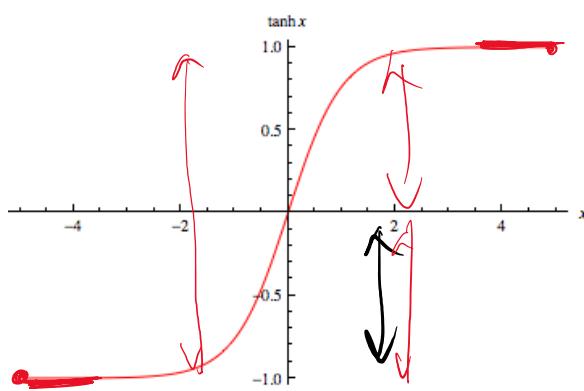
3)  $\sigma(x) = \frac{1}{1+e^{-x}}$   $\rightarrow$  computation time  
~~exp~~

hidden layer  
 $\downarrow$   
 sigmoid  
 $\uparrow$   
 output  $\rightarrow$  b c p

# Tanh Activation Function

31 May 2022 14:50

$(0,1)$   $(-1,1)$



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

$$f'(x) = (1 - \tanh^2(x))$$

## Advantages

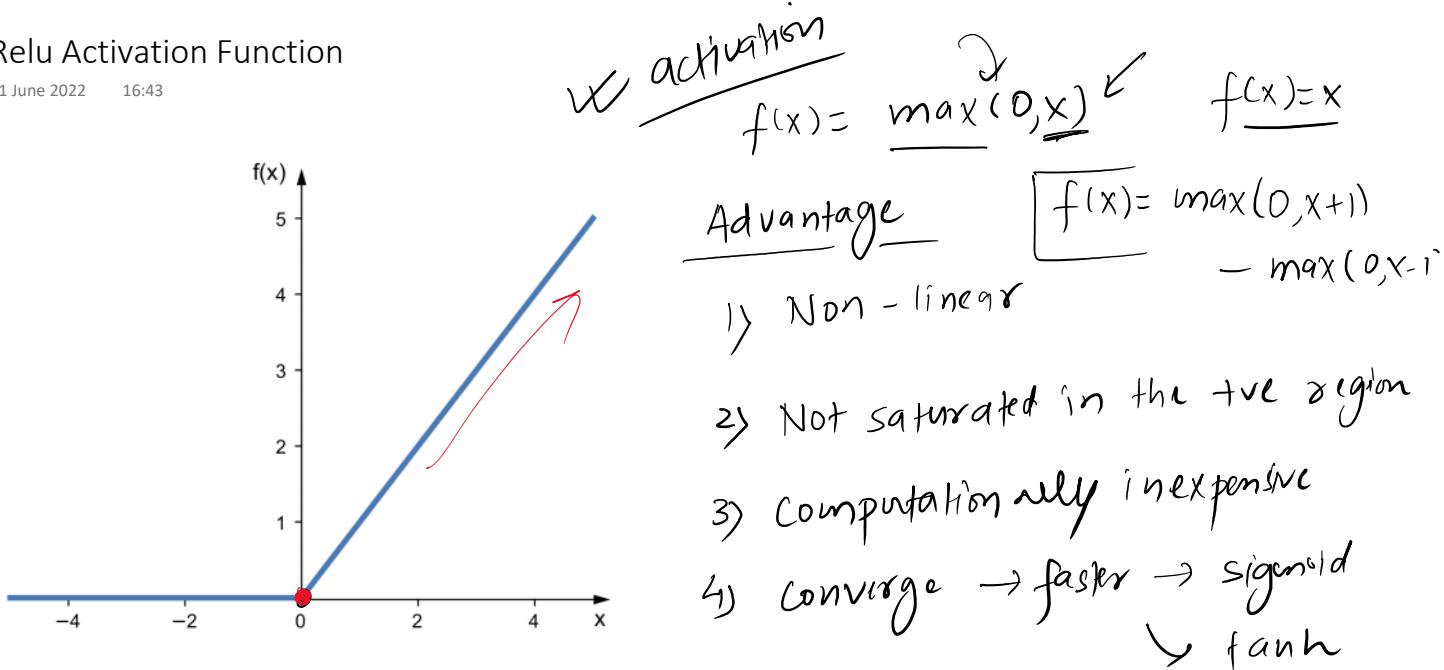
- 1) Non-linear
- 2) Differentiable
- 3) zero centered  $\rightarrow$  +ve  $\downarrow$  -ve  $\downarrow$  training faster sigmoid

## Disadvantage

- saturating function  $\rightarrow$  vanishing gradient prob
- computationally exp slow

## Relu Activation Function

01 June 2022 16:43



Disadvantage

→ Differentiability

$$x < 0 \rightarrow 0$$

$$x \geq 0 \rightarrow 1$$

→ NonZero centered → sigmoid → Batch Normalization

~~→~~ normalize →

↳ Dying ReLU problem

## Dying Relu Problem

08 June 2022 22:53

