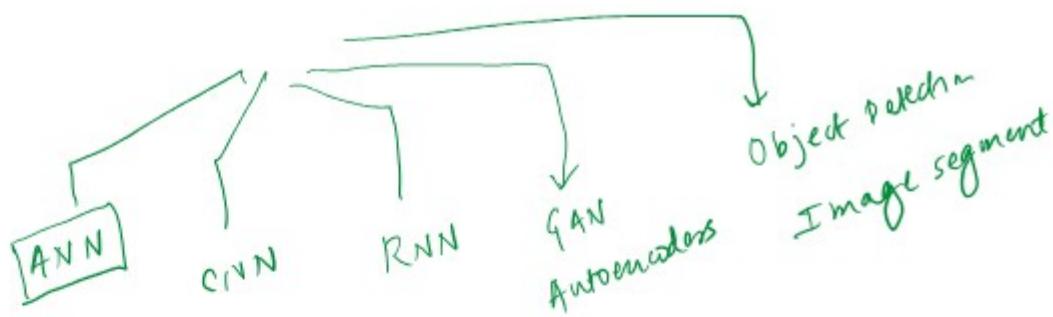


Curriculum

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

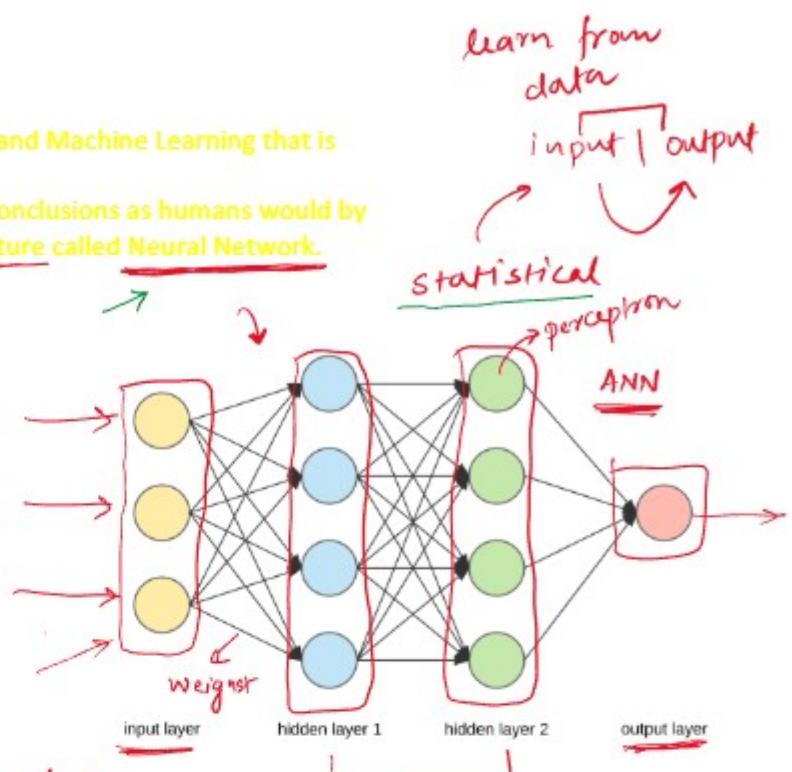
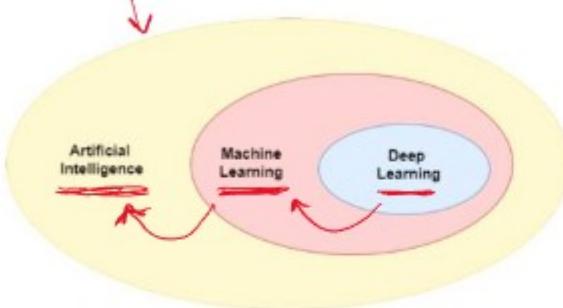
Tuesday, February 15, 2022 3:32 PM

- 1) Deep learning Roadmap
 - 2) Deep learning Project ideas
 - 3) Projects
 - 4) Interview questions
- Next video
→ Three
- 10th AP⁶
→ ANN → CNN
→ 3-4 months

What is Deep Learning?

16 February 2022 06:32

Deep Learning is a subfield of Artificial Intelligence and Machine Learning that is inspired by the structure of a human brain. Deep learning algorithms attempt to draw similar conclusions as humans would by continually analyzing data with a given logical structure called Neural Network.



ANN →

CNN → convolutional
NN ✓ image

RNN → Recurrent NN → speech / text

GAN → generate (text) images

why dl getting so famous?

→ 1) Applicability →

2) Performance → state-of-the-art

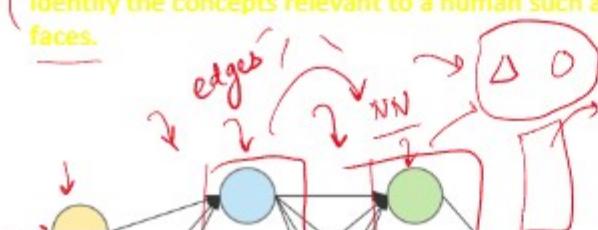
Applications

Go → AlphaGo → AI
5 → ④

self-driving
drugs
→ chemistry

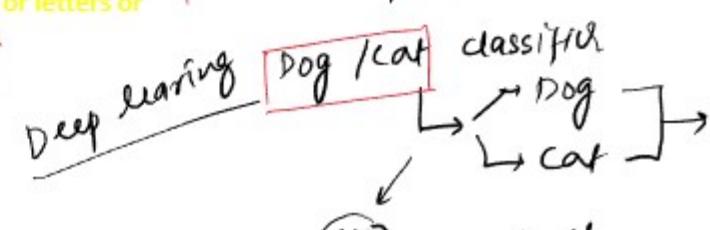
Deep learning is part of a broader family of machine learning methods based on artificial neural networks with representation learning.

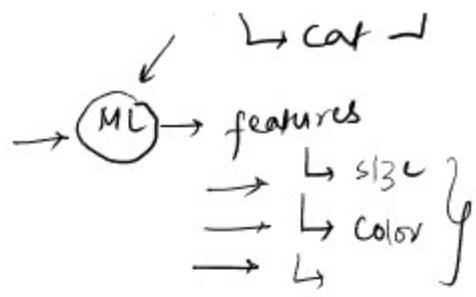
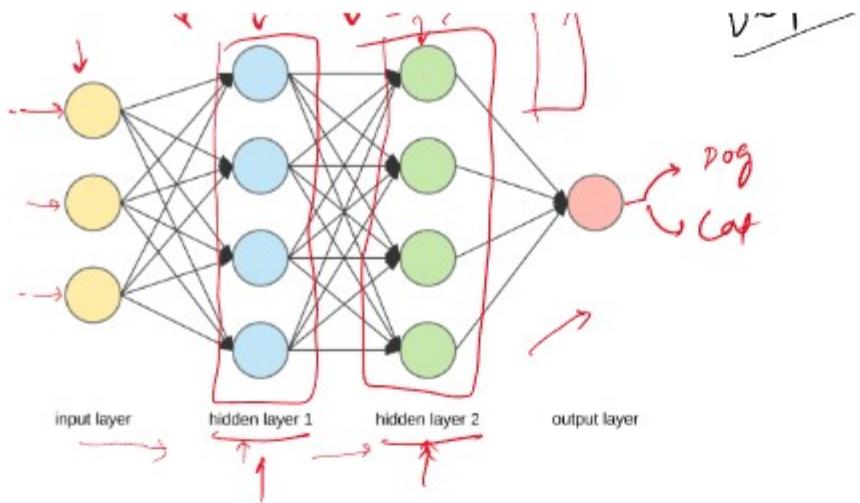
Deep Learning Algorithms uses multiple layers to progressively extract higher-level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits or letters or faces.



→ Representation

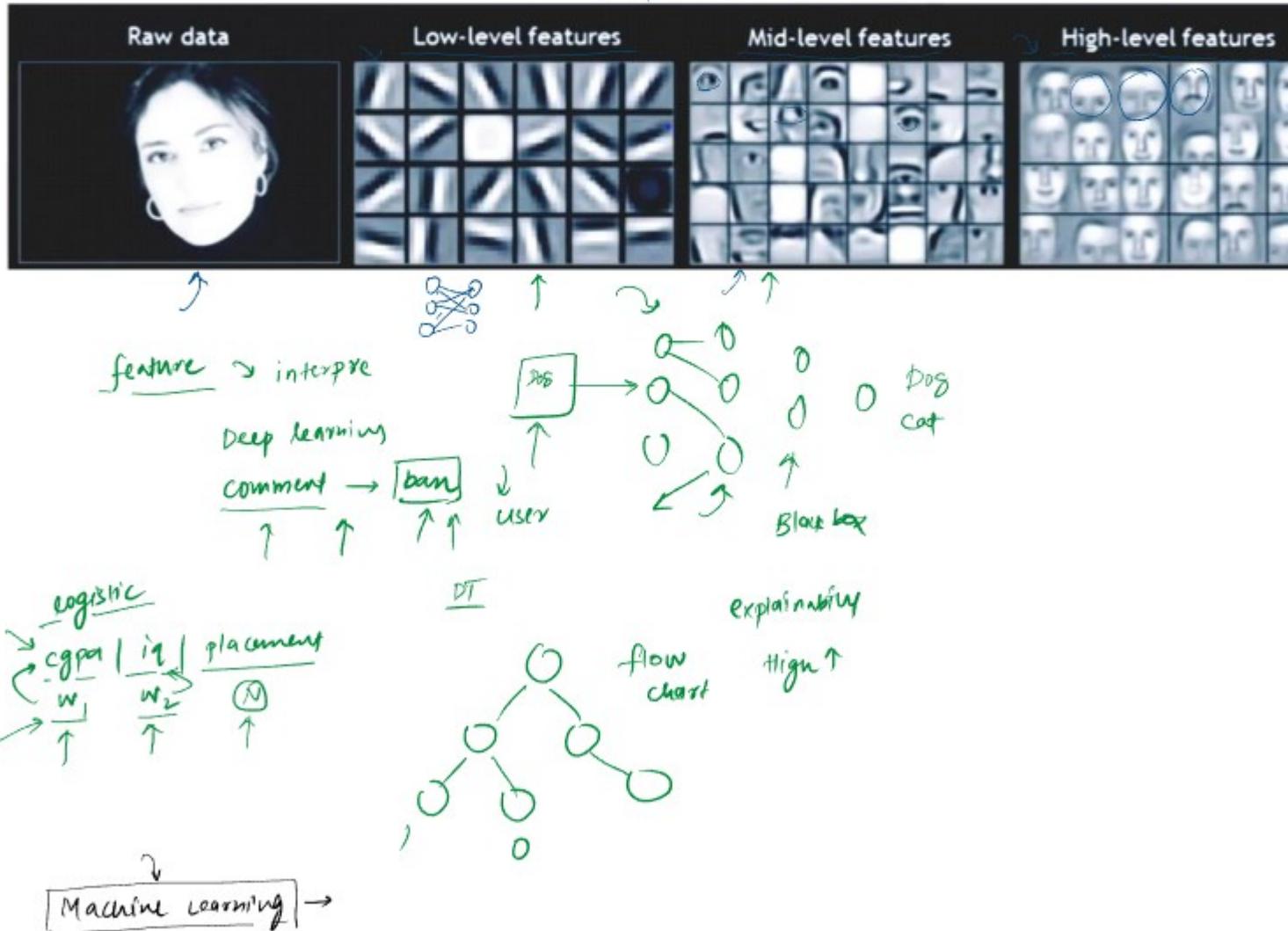
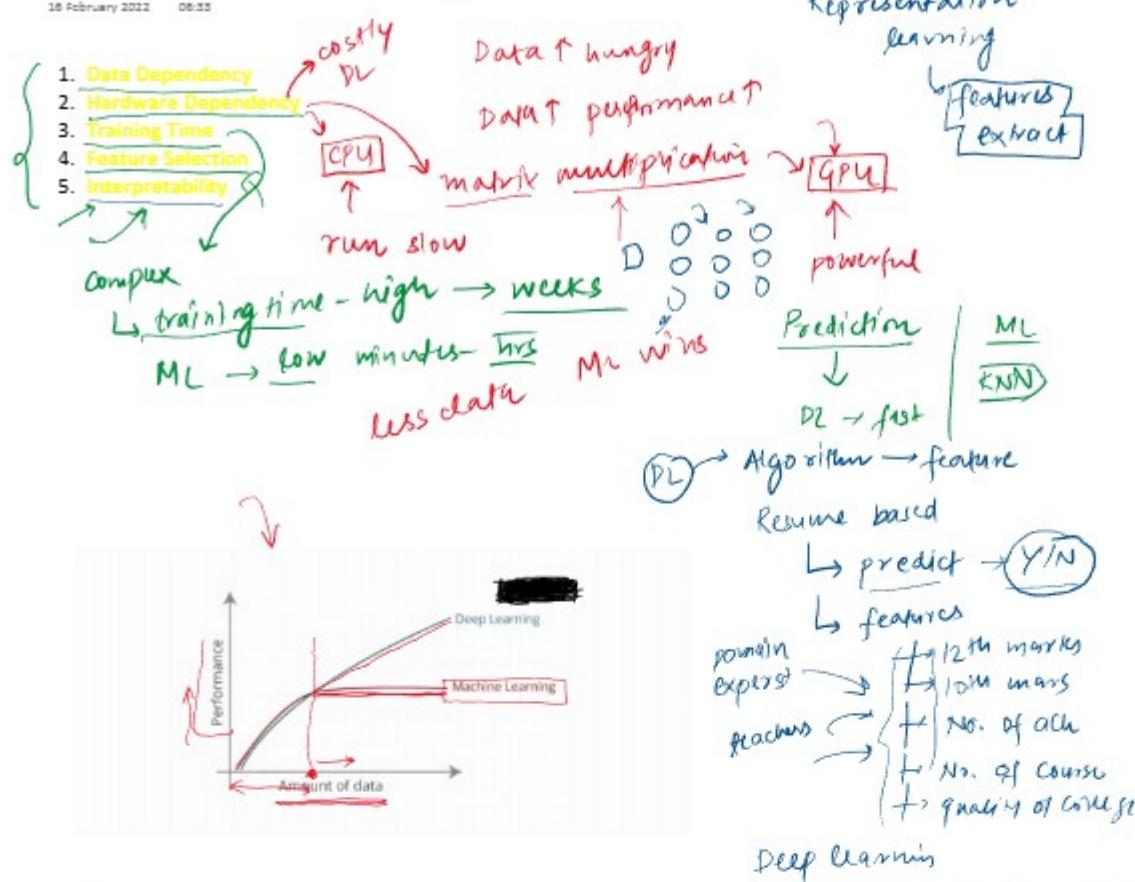
feature extraction
engineering
features





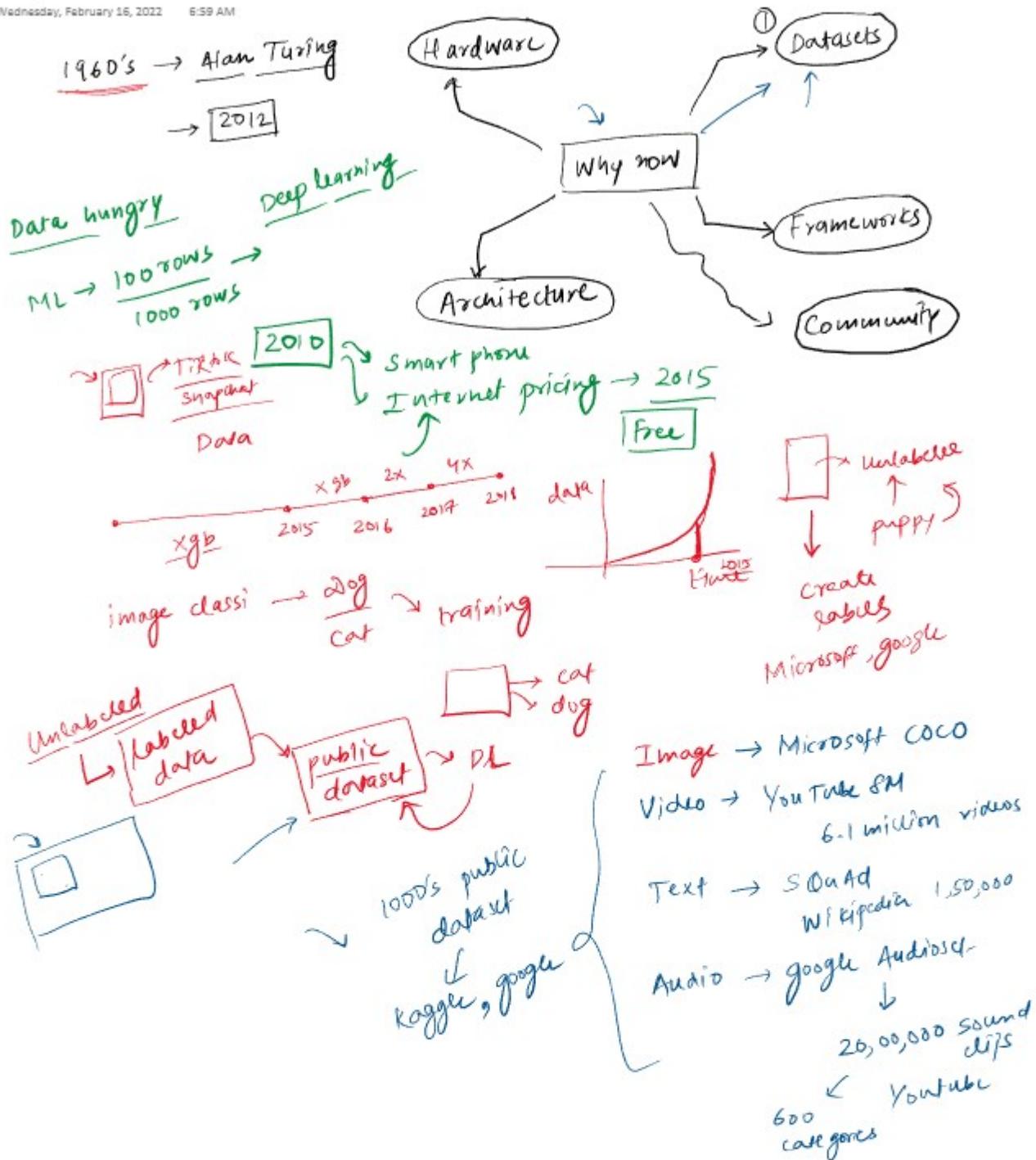
Deep Learning VS Machine Learning

18 February 2022 08:55

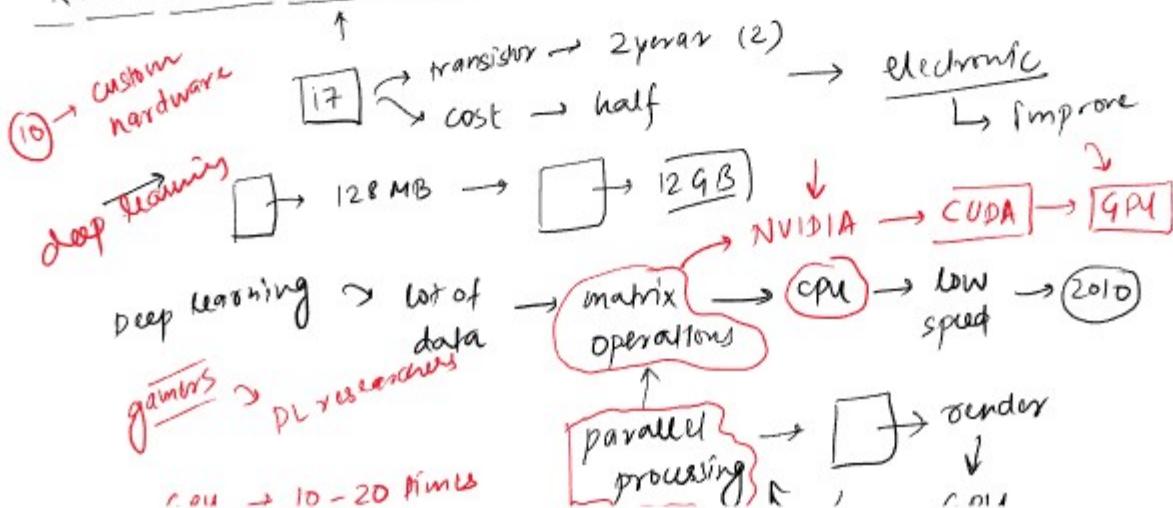


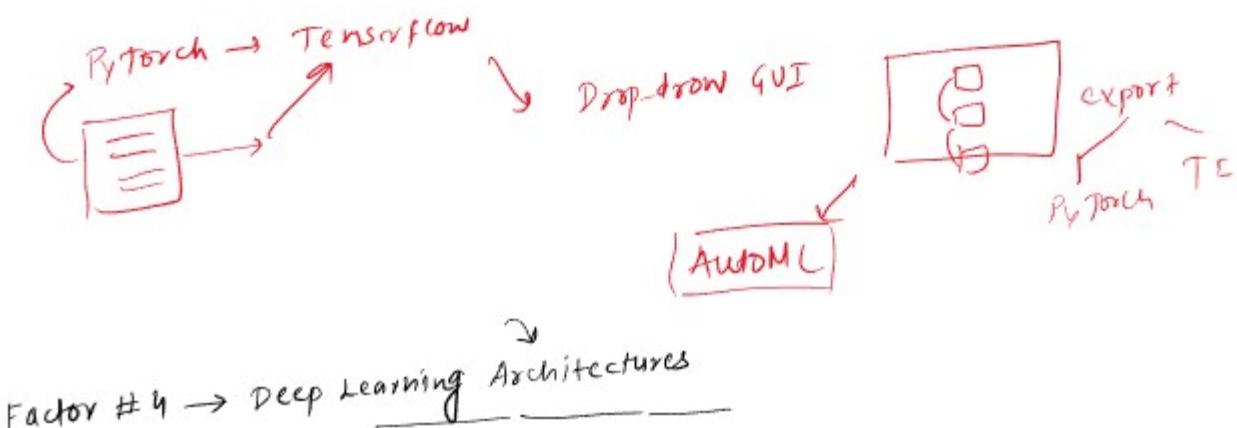
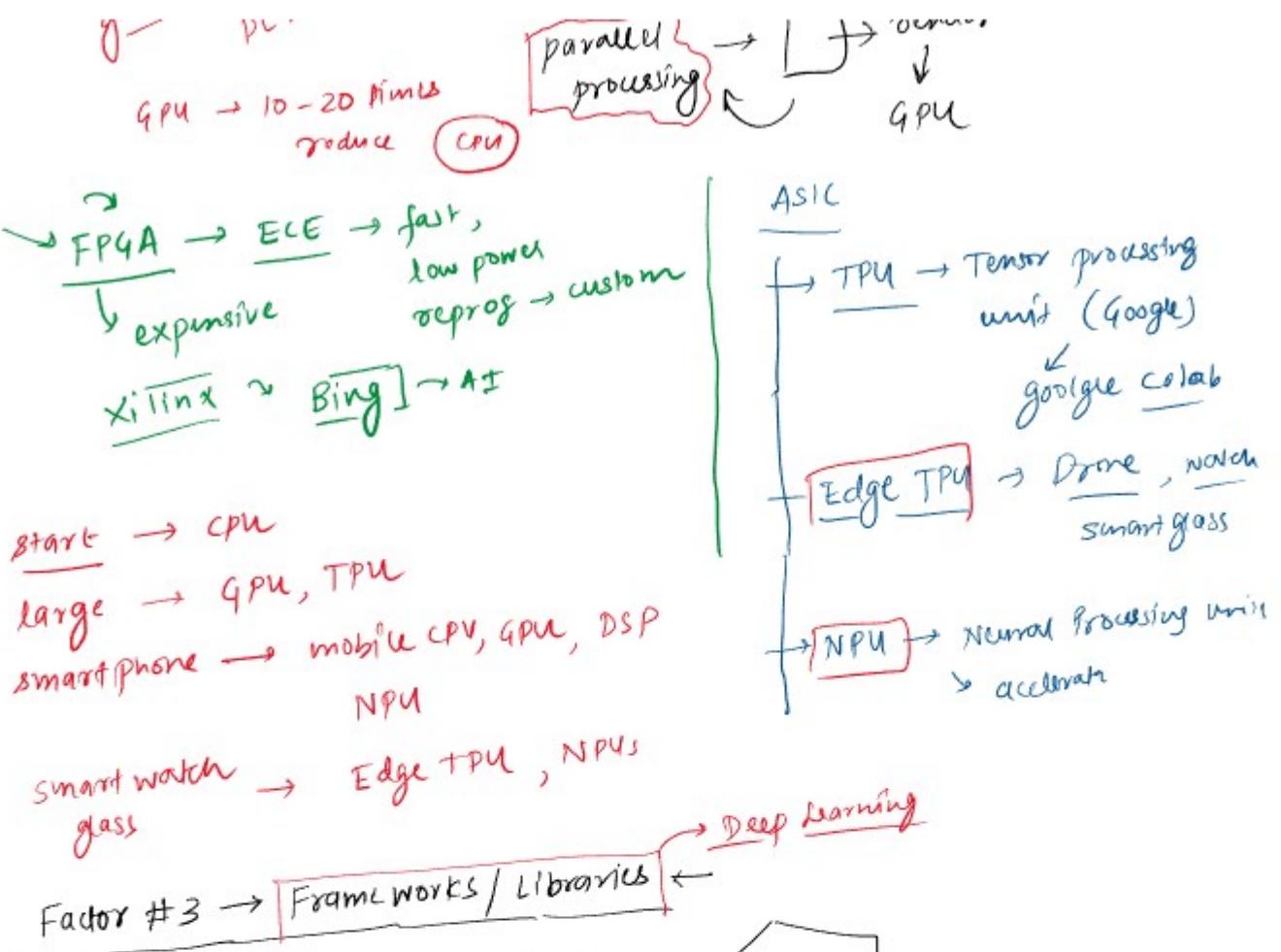
Why now?

Wednesday, February 16, 2022 6:59 AM

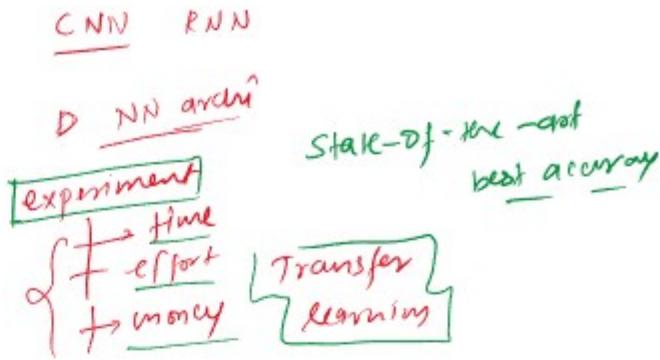
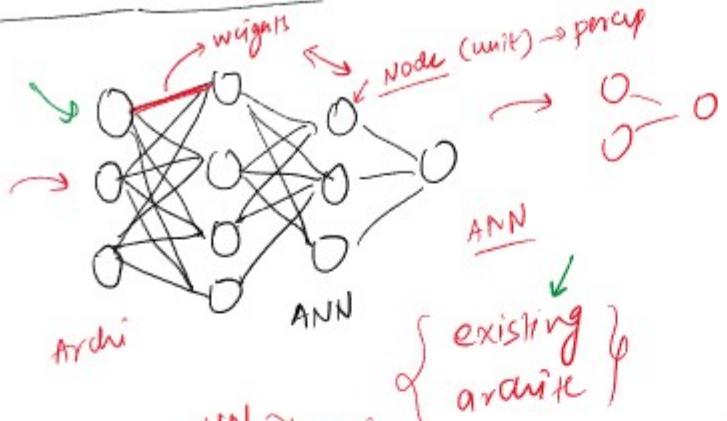


Reason #2 → Hardware





Factor #4 → Deep Learning Architectures



NN archi
Dataset → train → ready to use →
fires down

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 ↗
(PL) →

Research

↳ DL engin

↳ Teacher

→ student → kaggle

Code Example

Saturday, February 19, 2022 7:10 AM

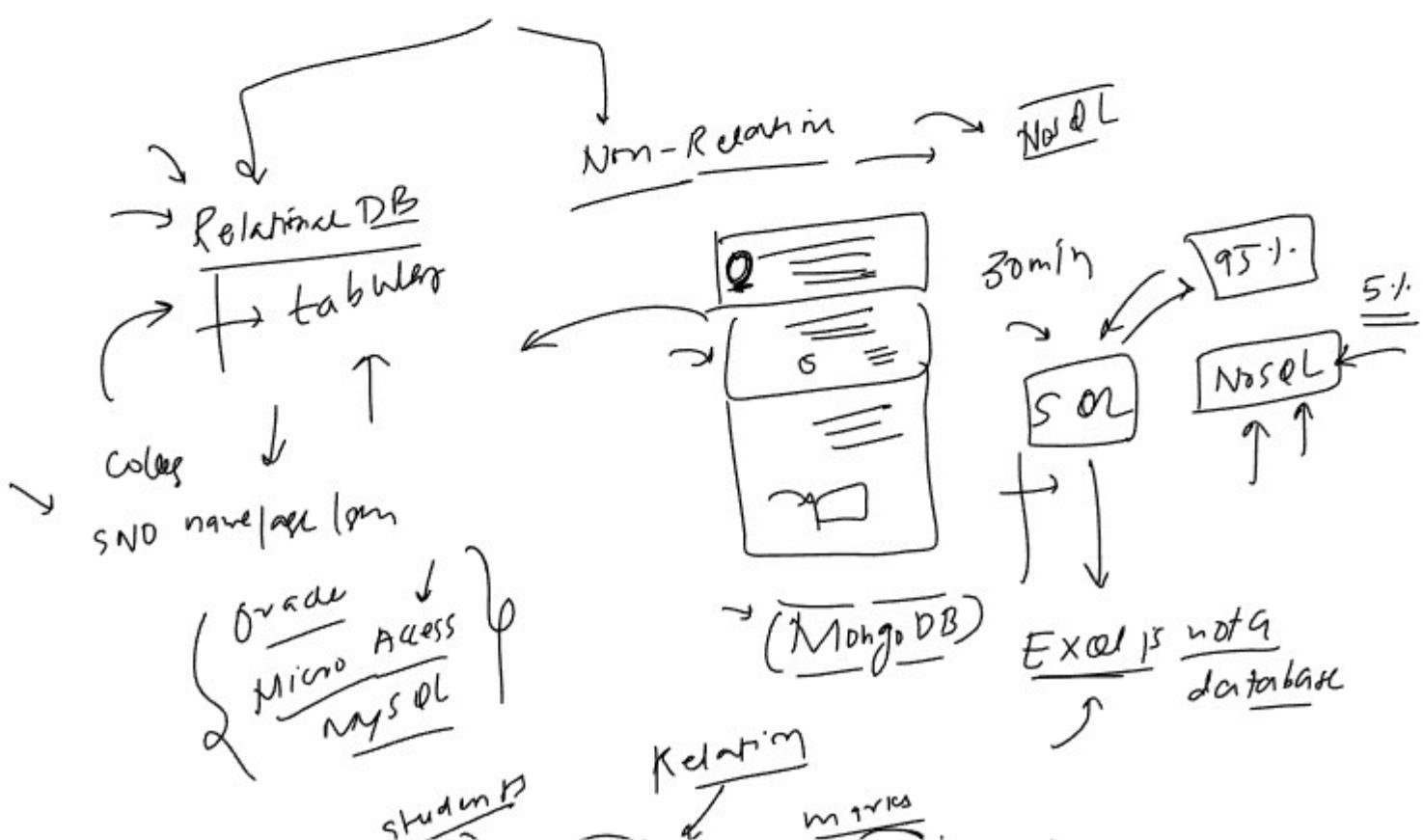
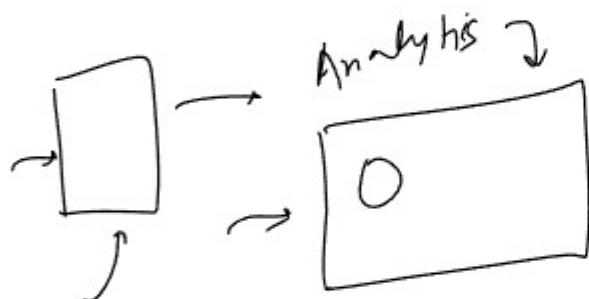
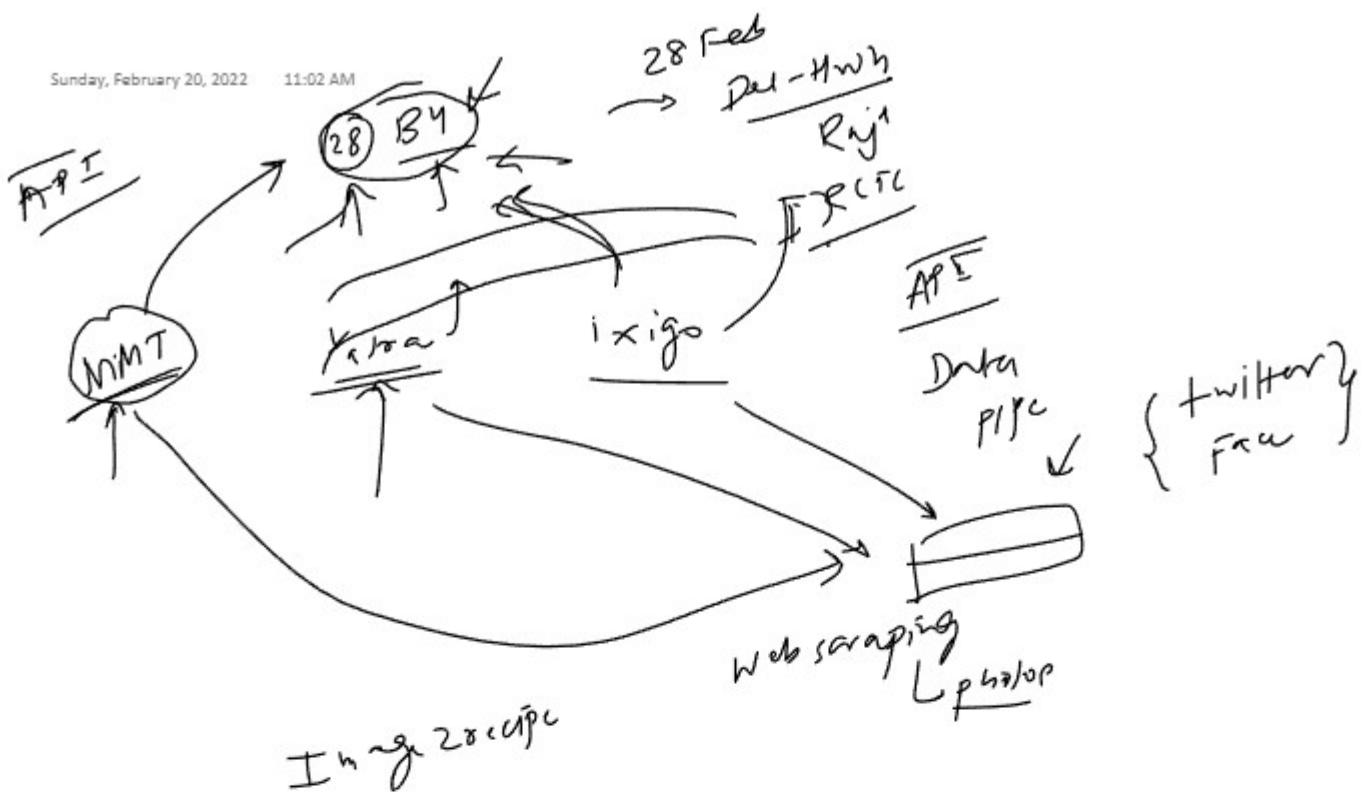
cgpa | resume / placed

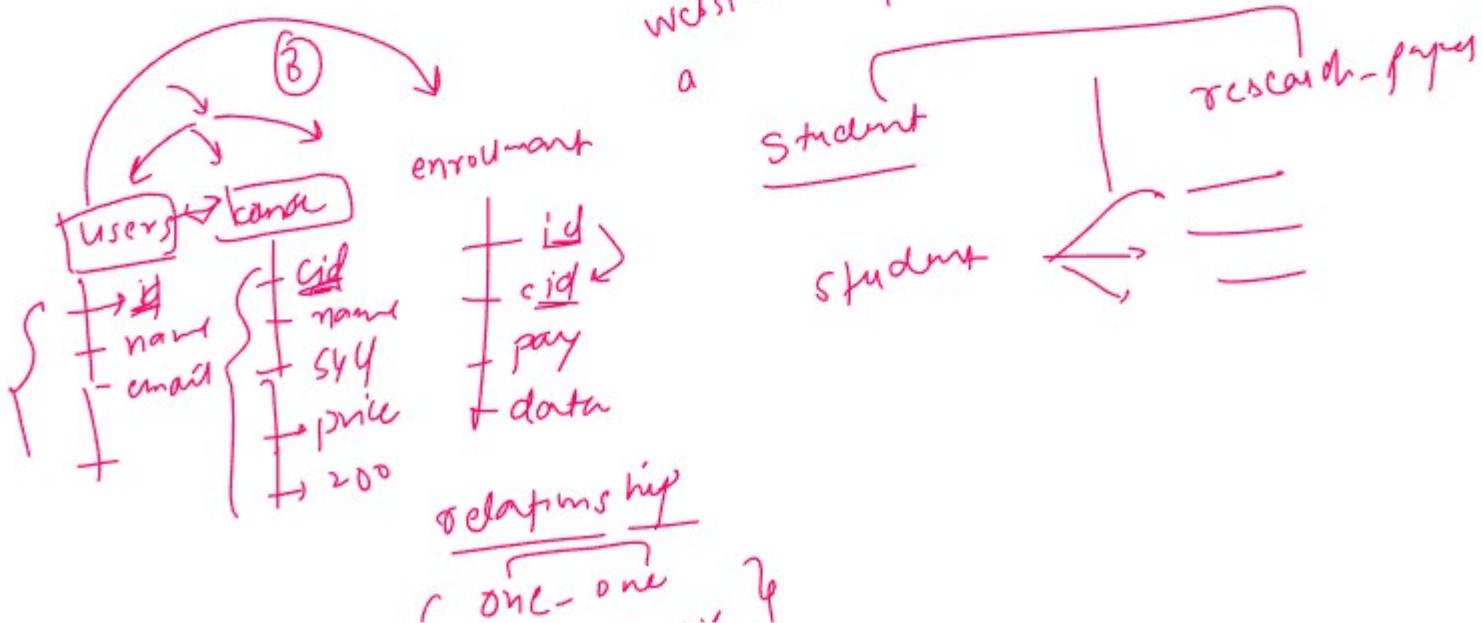
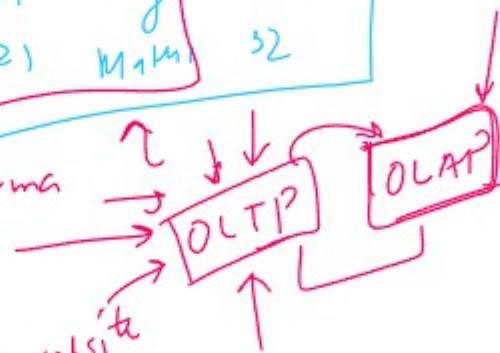
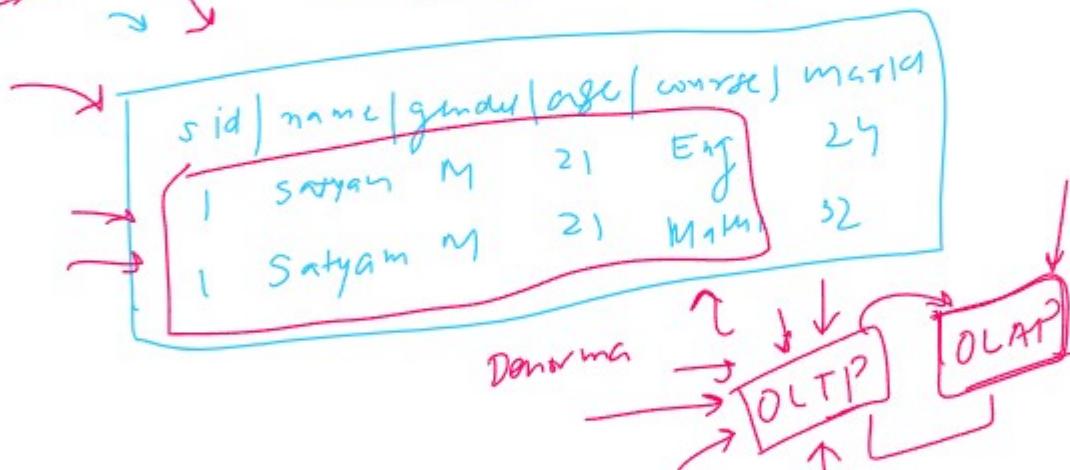
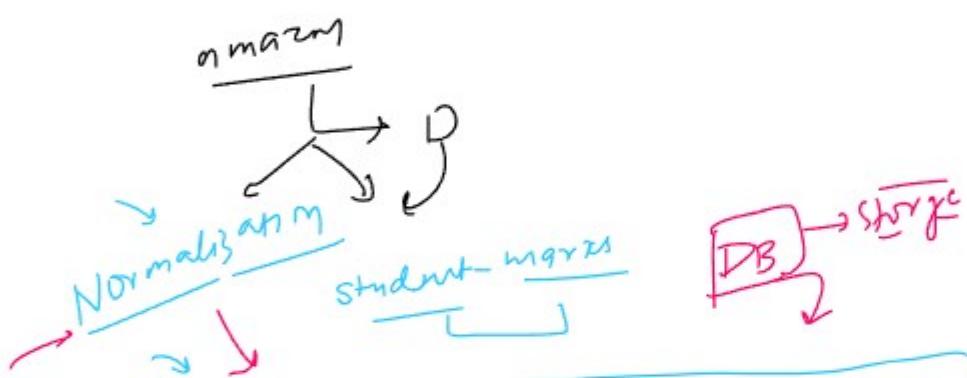
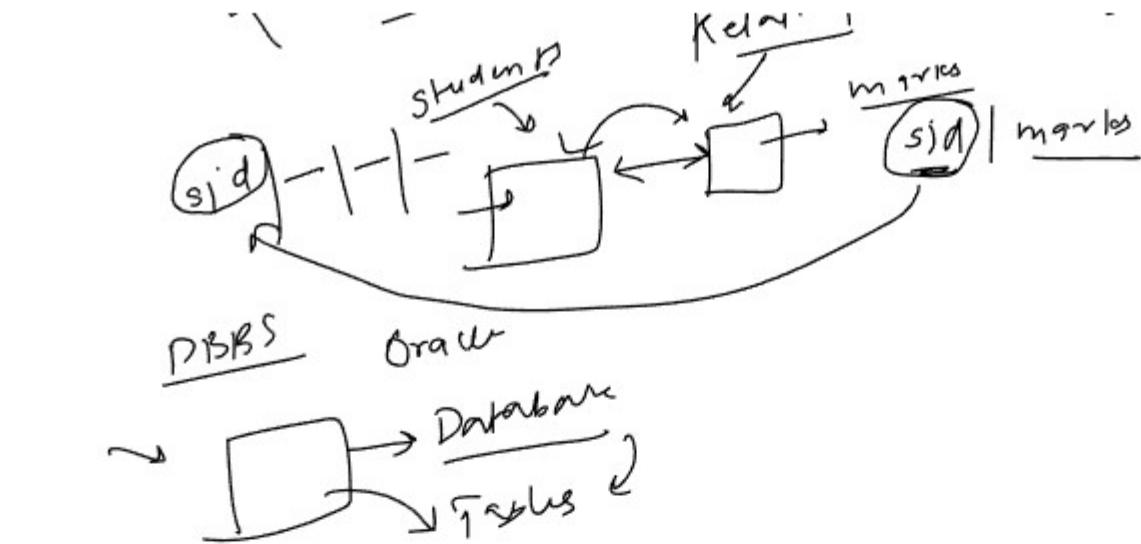


$$w_1 \quad w_2 \quad b = -25$$



cgpa

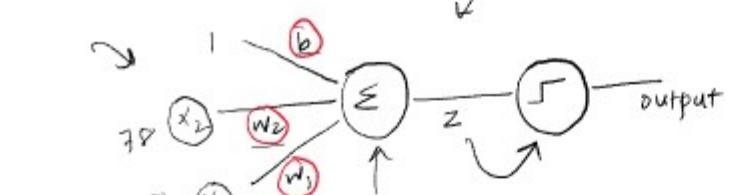




$\rightarrow \{$ one-one
one-many }
many-

Recap

22 February 2022 14:12



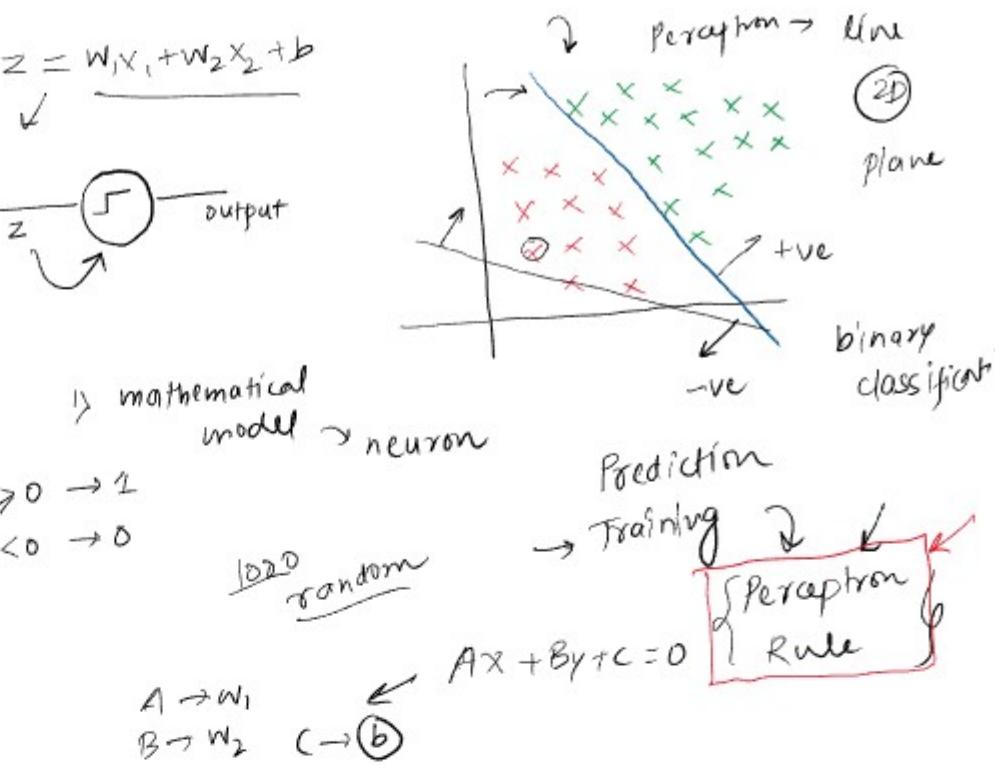
x_1	x_2	placed
cgpa	iq	
7	78	1
6	61	0
9	92	1

↳ mathematical model \rightarrow neuron

$$\begin{aligned} z \geq 0 &\rightarrow 1 \\ z < 0 &\rightarrow 0 \end{aligned}$$

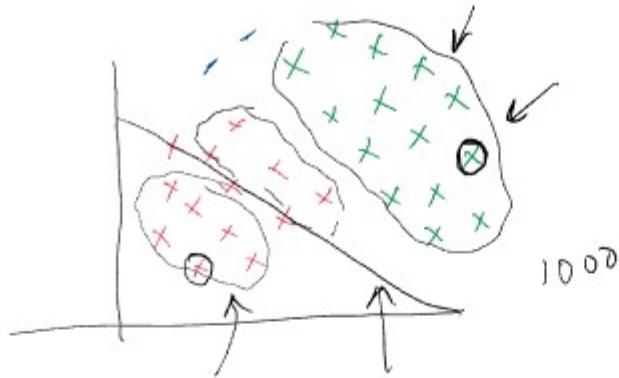
$\xrightarrow{1020 \text{ random}}$

$$\begin{aligned} A &\rightarrow w_1 \\ B &\rightarrow w_2 \\ C &\rightarrow b \end{aligned}$$



Problem with Perceptron Trick

22 February 2022 16:10



$$\text{line} \rightarrow [w_1, w_2, b]$$

randomly per

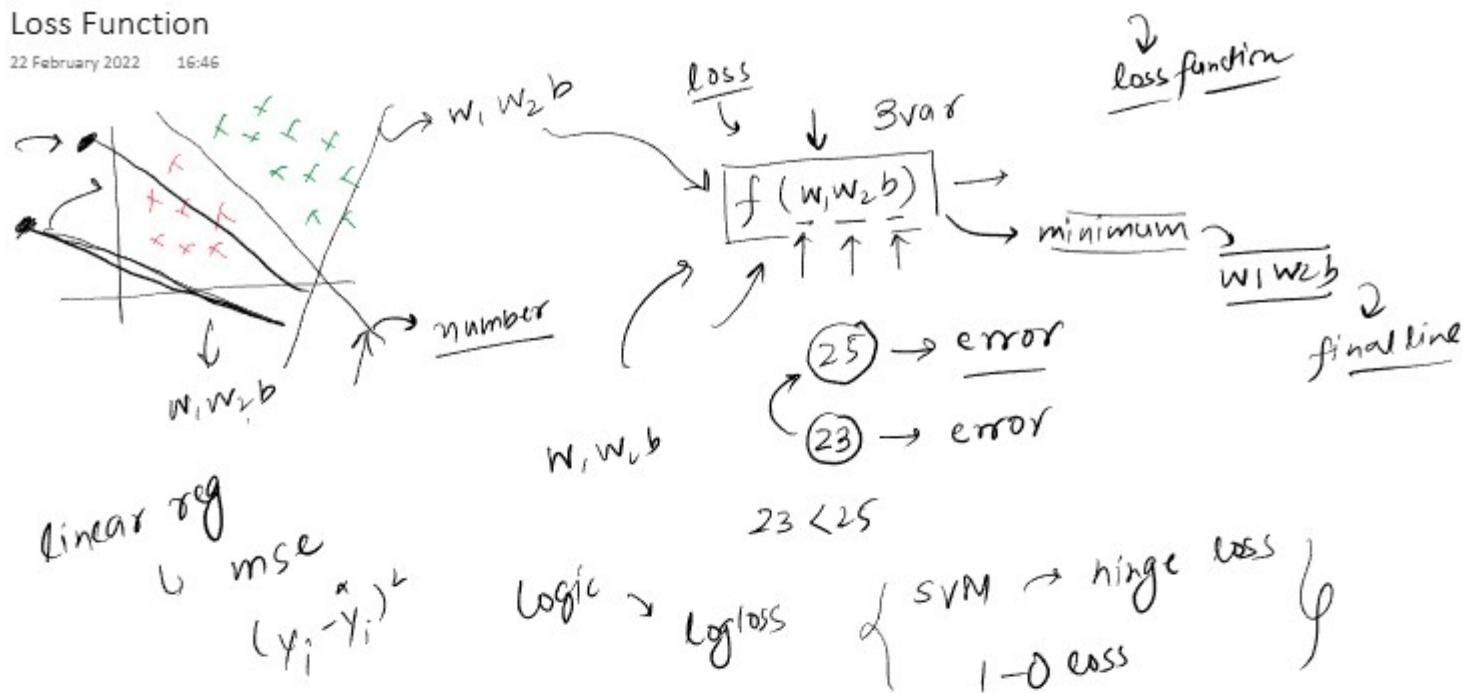
converge

quantify

$[w_1, w_2, b] \rightarrow$ loss functions

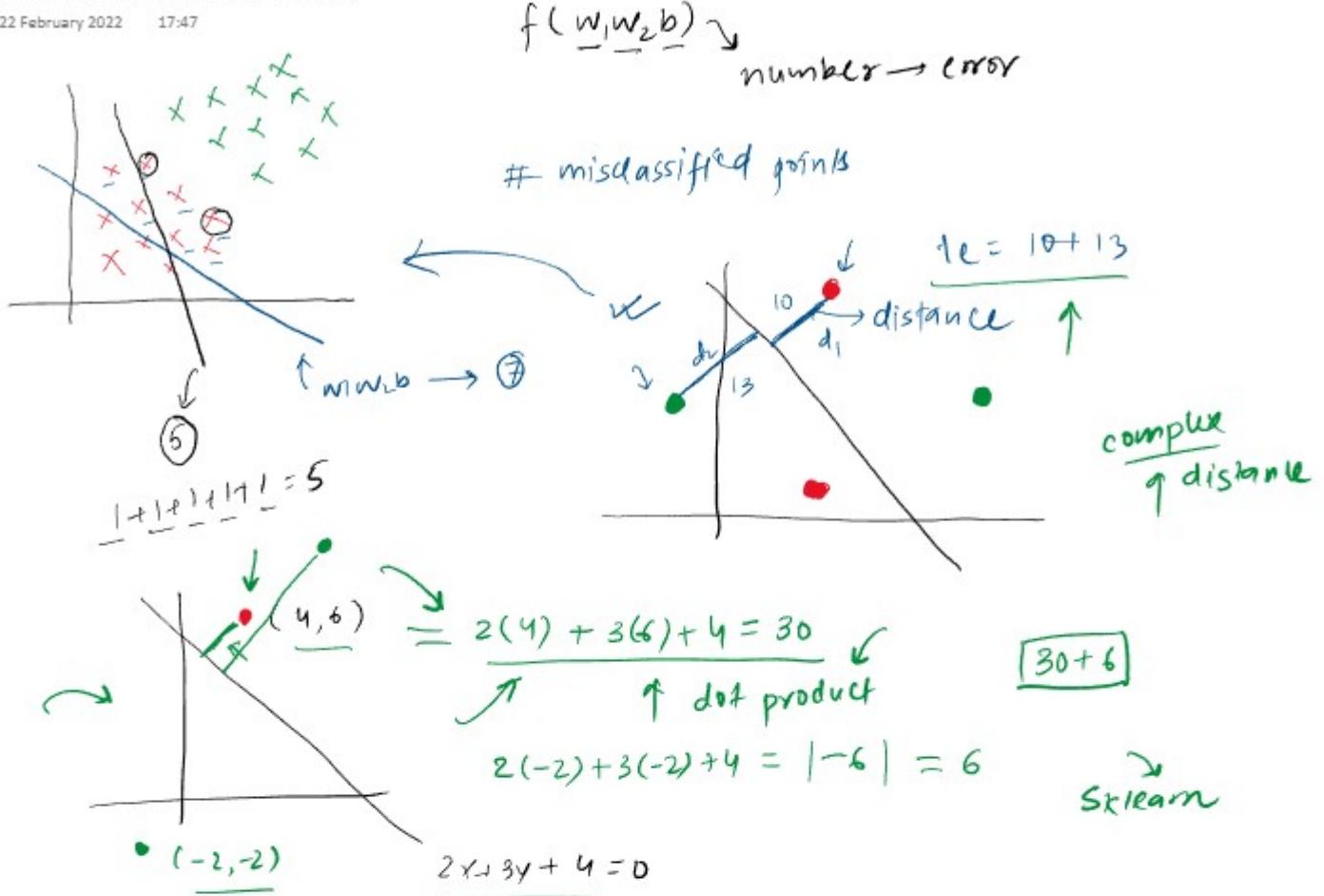
Loss Function

22 February 2022 16:46



Perceptron Loss Function

22 February 2022 17:47



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

Regularization

categorical plausibility

y_i

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

$f(x_i) = \frac{w_1 x_1 + w_2 x_2 + b}{z = 1}$

hinge

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

$L(w, w_2 b)$

$w_1 w_2 b$

n: # rows in data

$$L = \underset{w, w_2 b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

minimum

gradient descent

$$L = \frac{1}{n} \sum_{i=1}^n \text{max}(0, w_1 x_1 + w_2 x_2 + b)$$

gradient descent

Explanation of Loss Function

23 February 2022 08:07

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

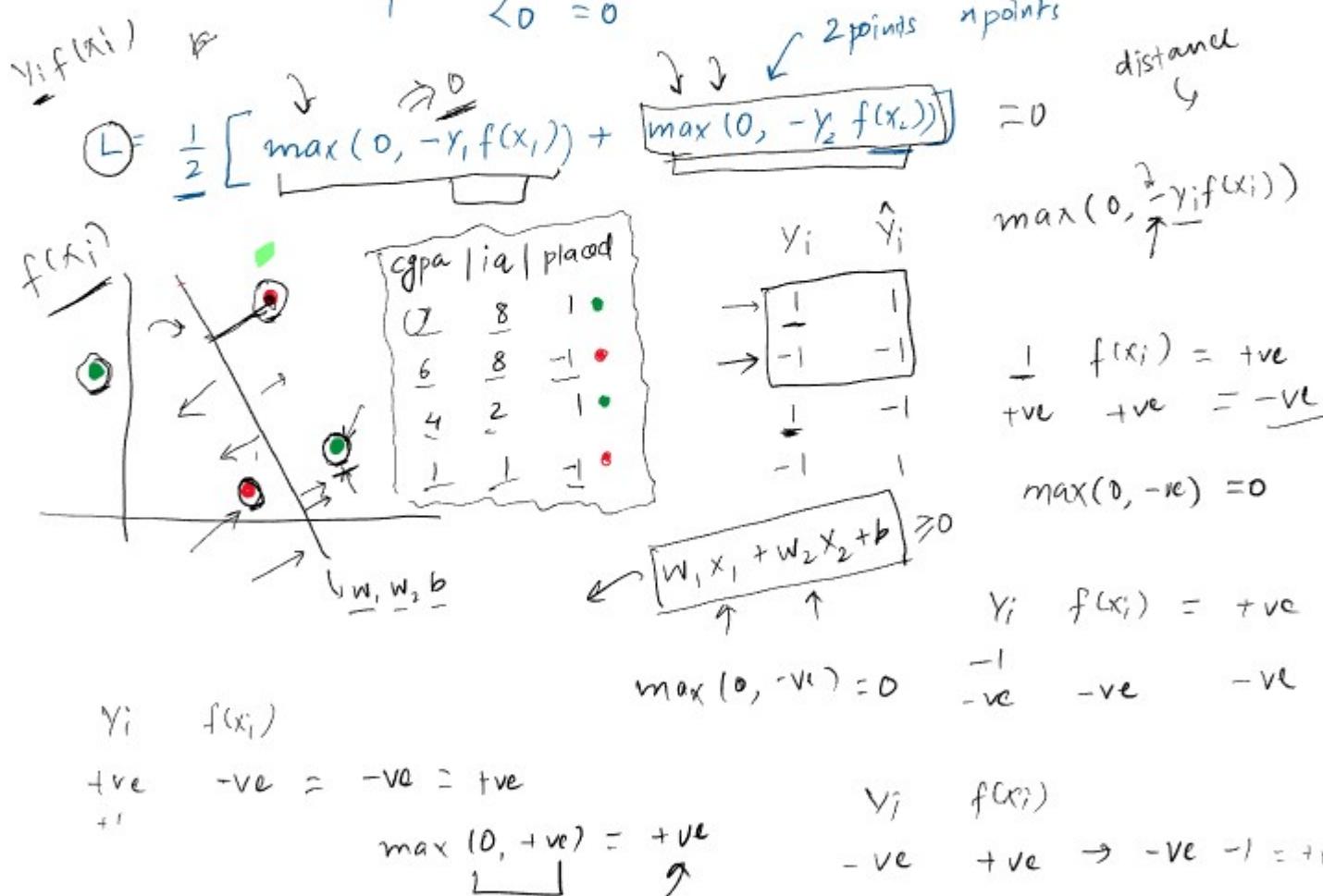
where $f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$

x_1	x_2	y
x_{11}	x_{12}	y_1
x_{21}	x_{22}	y_2
\vdots	\vdots	\ddots points
x_{i1}	x_{i2}	y_i

n rows
i points
row col

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

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



Gradient Descent

23 February 2022 08:50

$$L = \underset{w_1, w_2, b}{\operatorname{argmin}} \left[\frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i)) \right]$$

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

$$L(w_1)$$



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

100 / 1000

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

$\overbrace{w_1, w_2, b}^{\rightarrow} = 1$

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

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

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

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

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

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

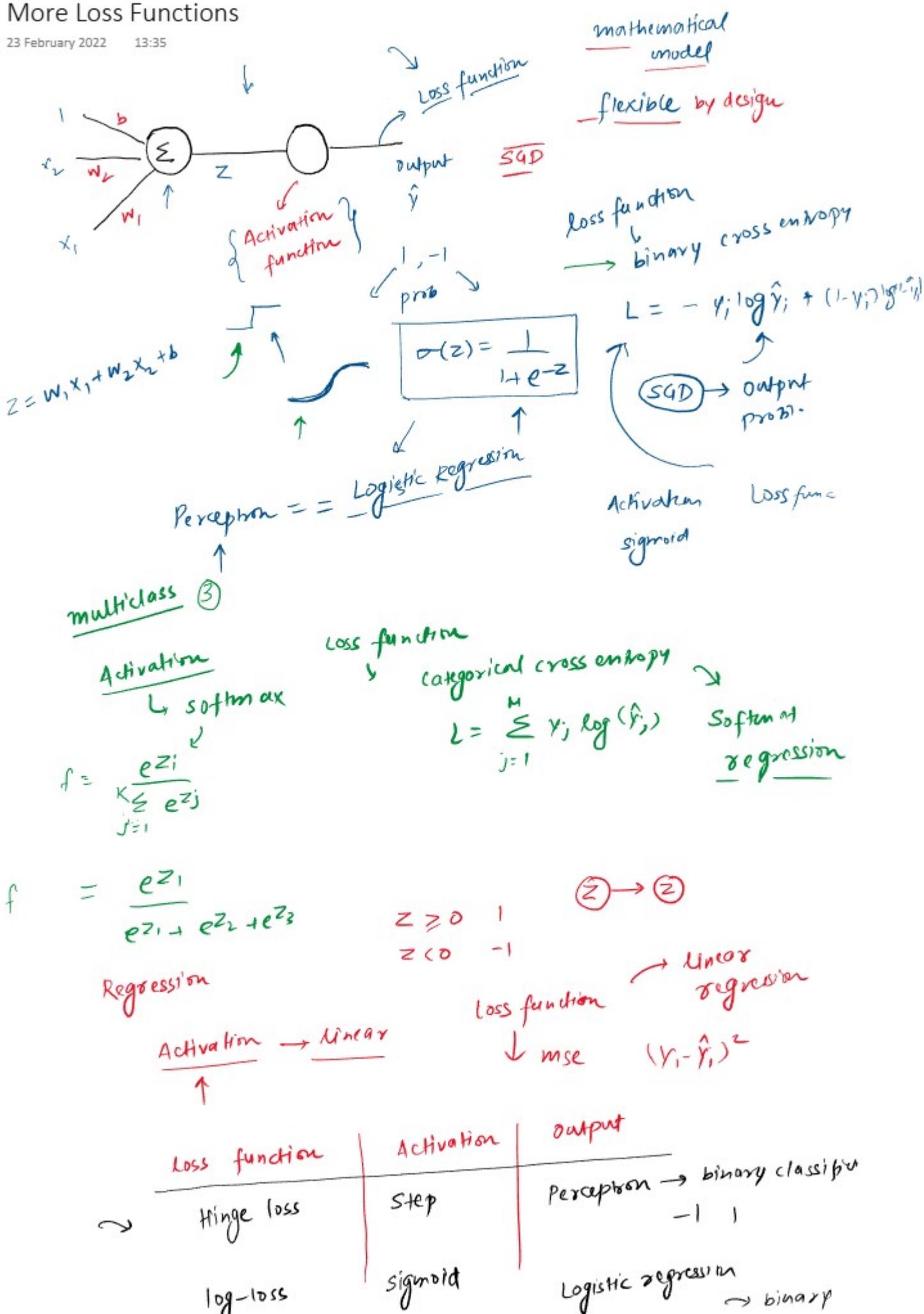
$$\boxed{\frac{\partial L}{\partial b} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}}$$

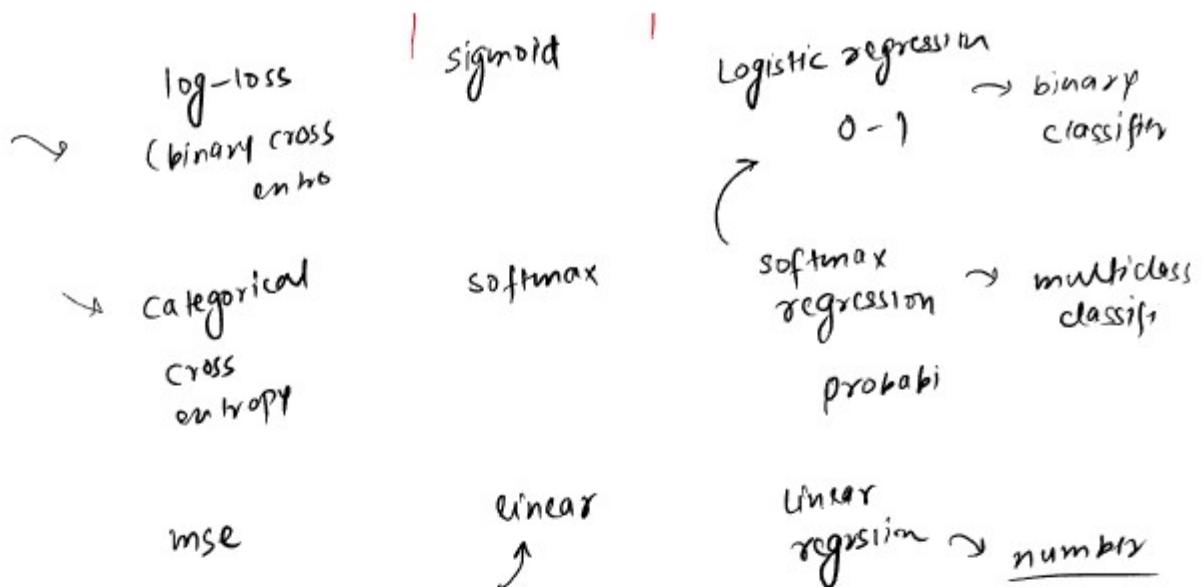
(cgpai) is | played
→ x_{i1}

→ code using
gradient descent

More Loss Functions

23 February 2022 13:35

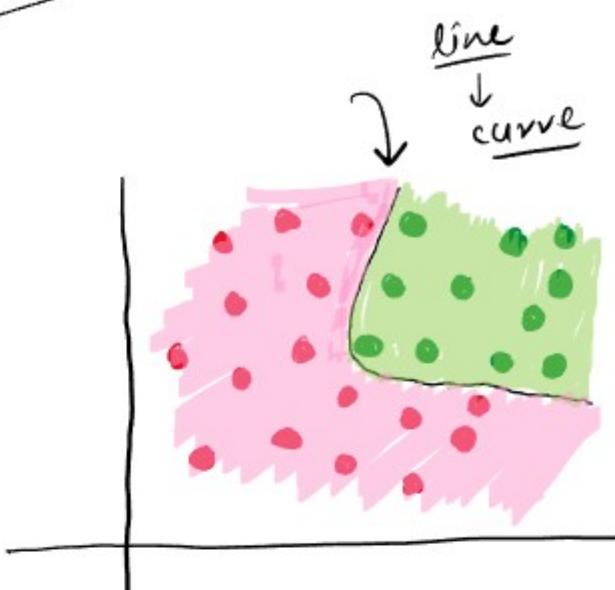
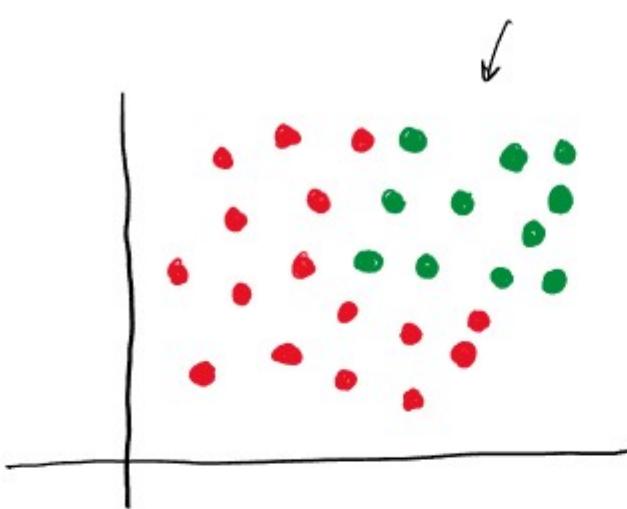




The Problem

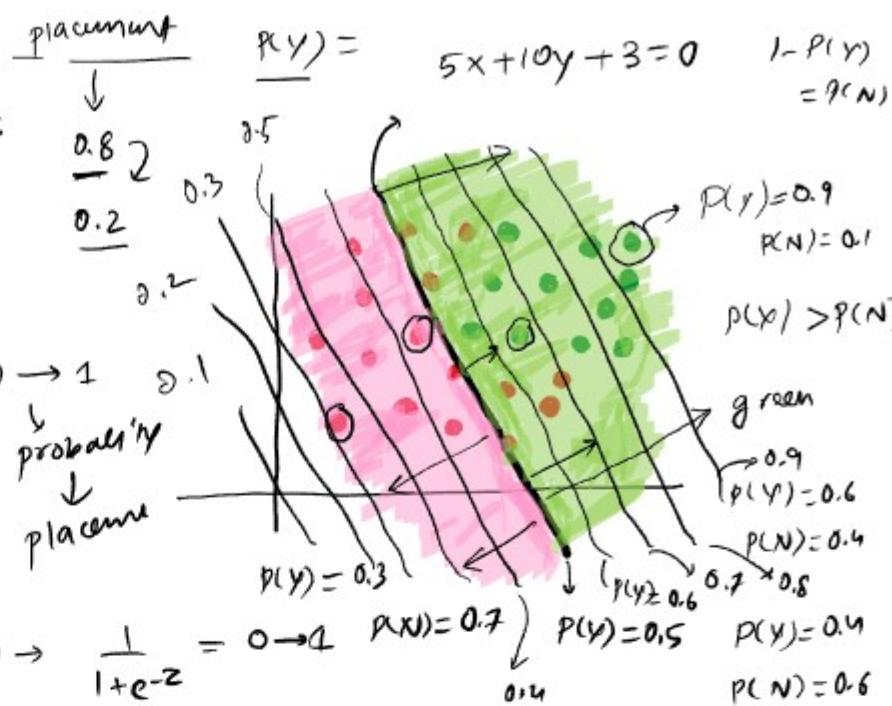
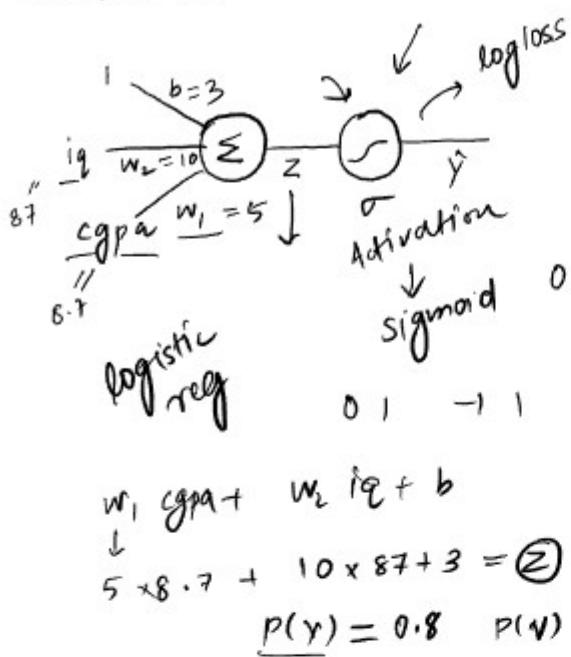
24 February 2022 13:41

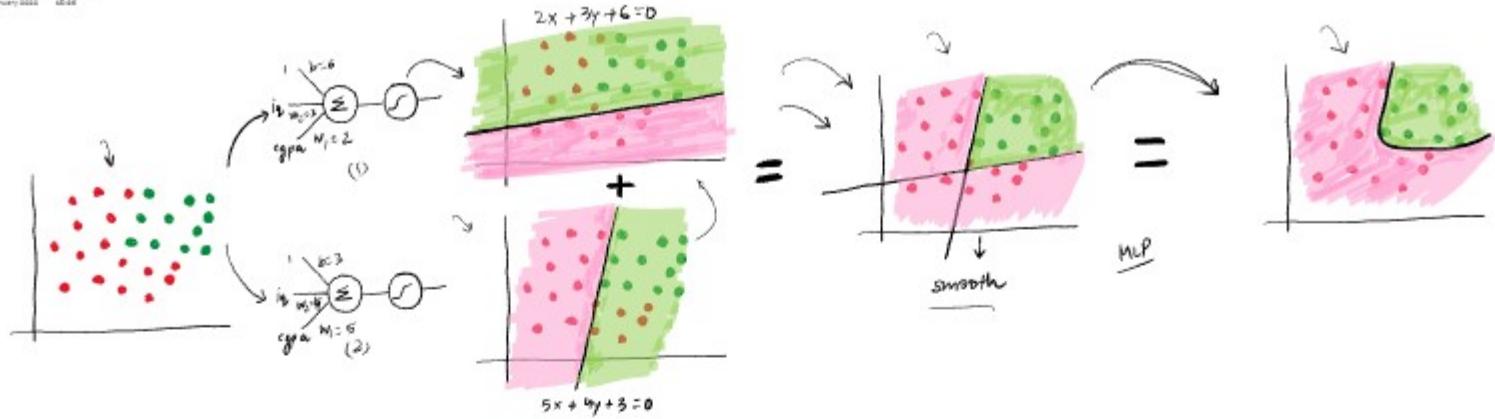
perceptron

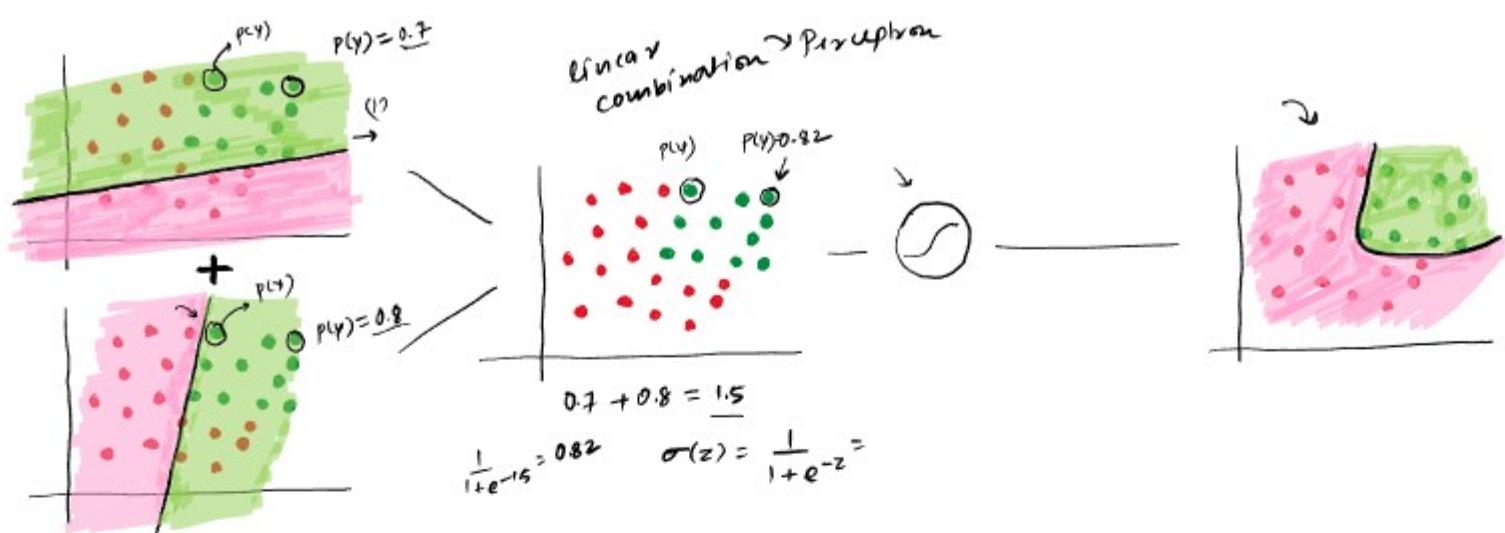


Perceptron with Sigmoid

24 February 2022 15:35

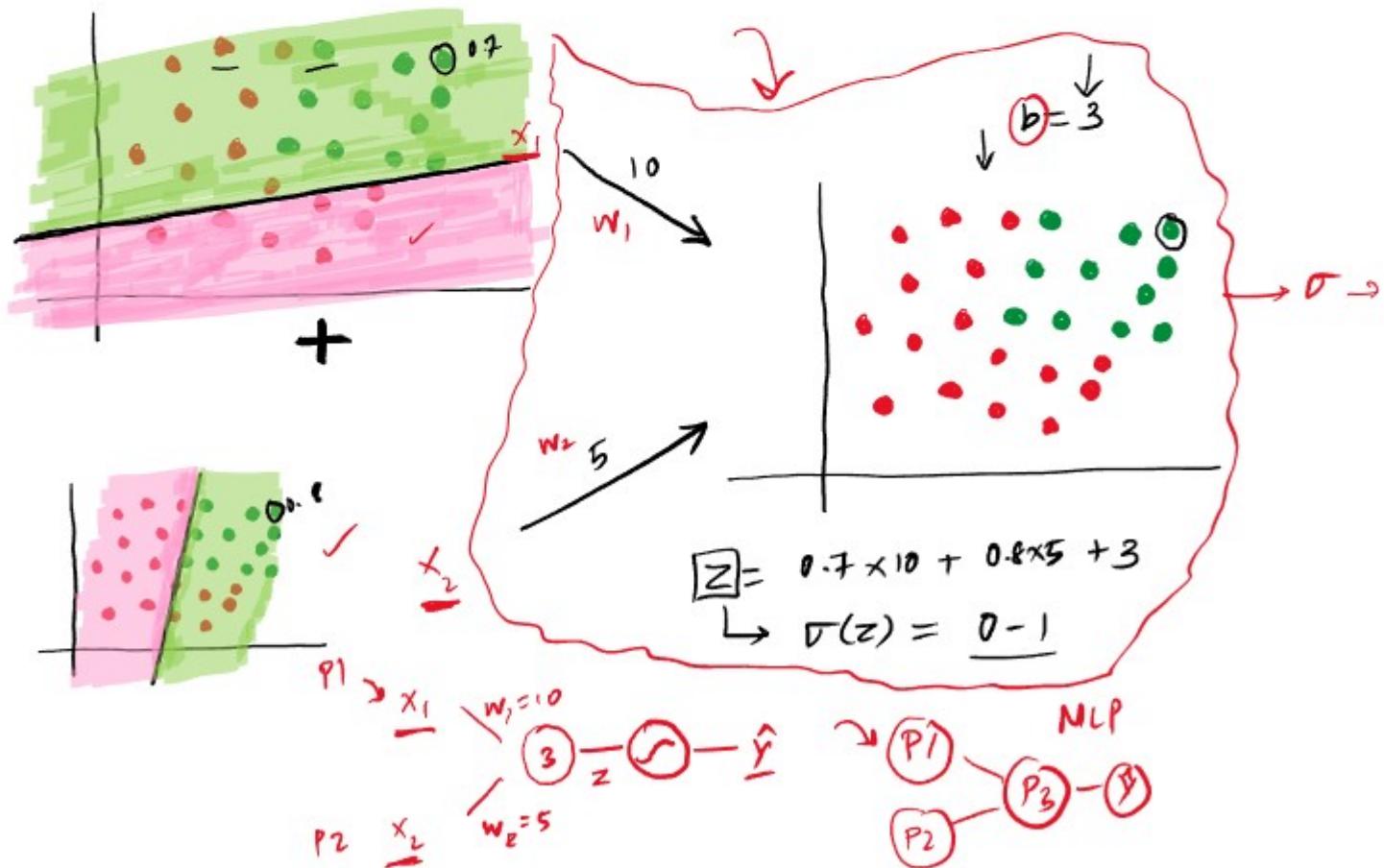
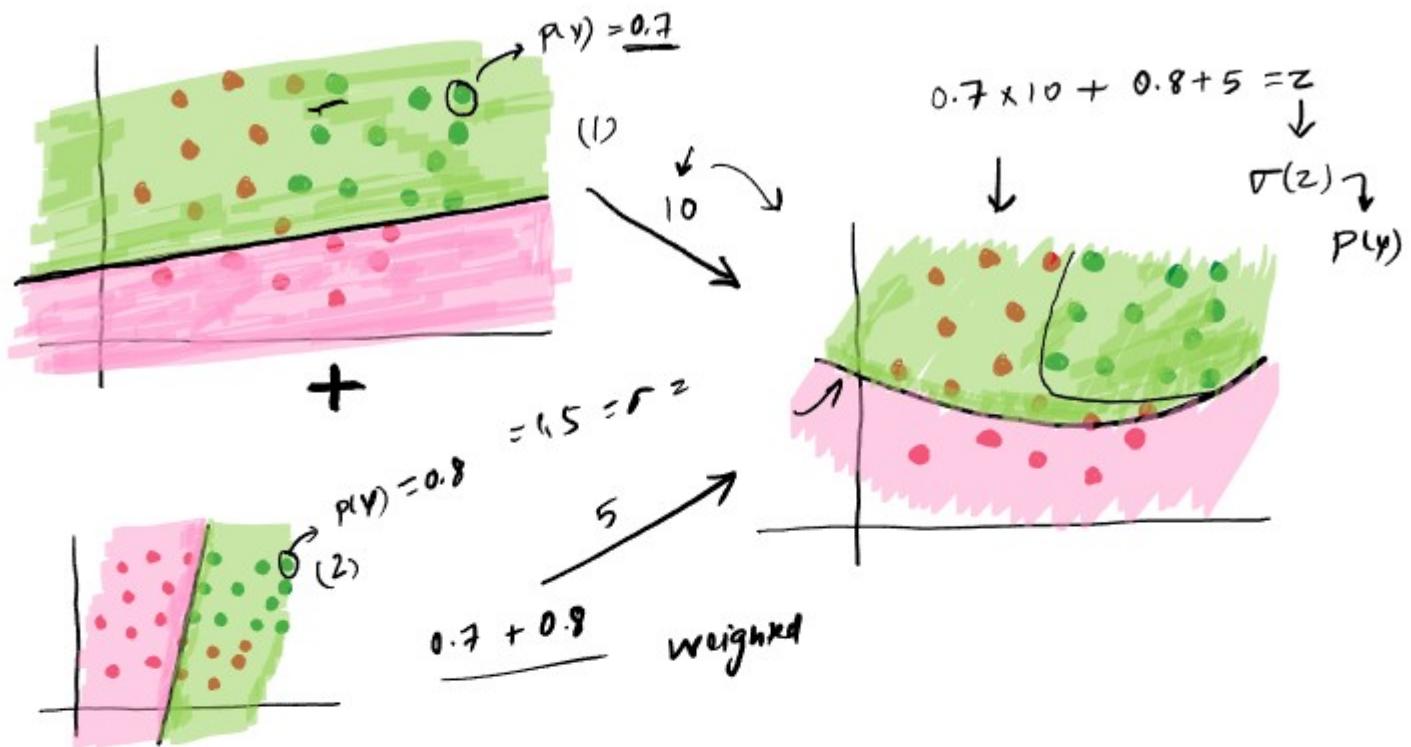






Adding Weights

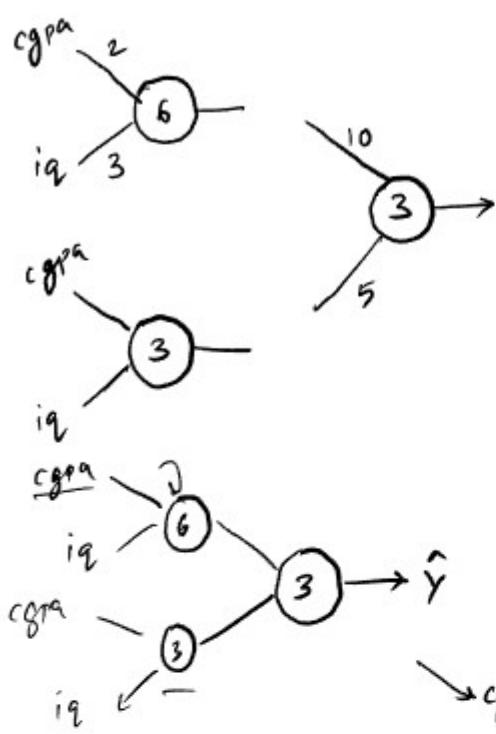
24 February 2022 15:37



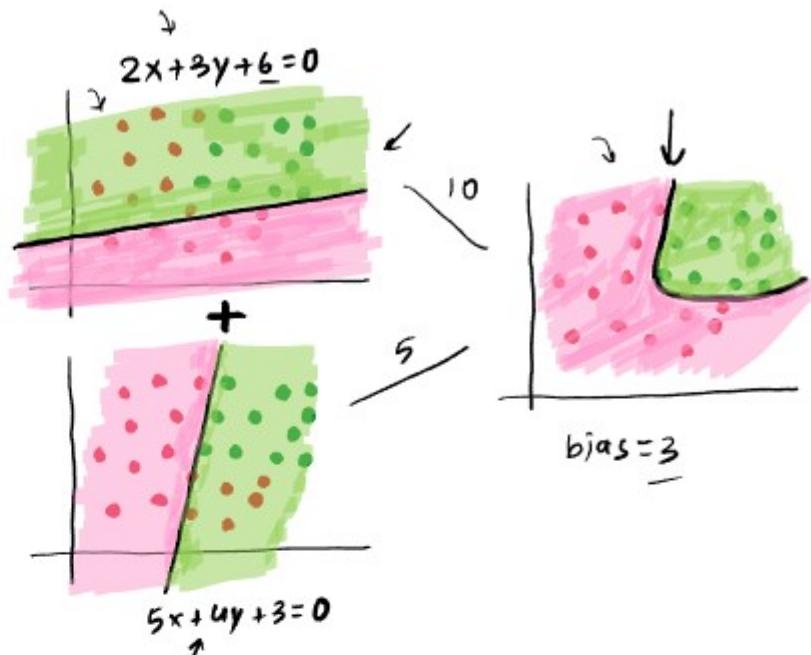
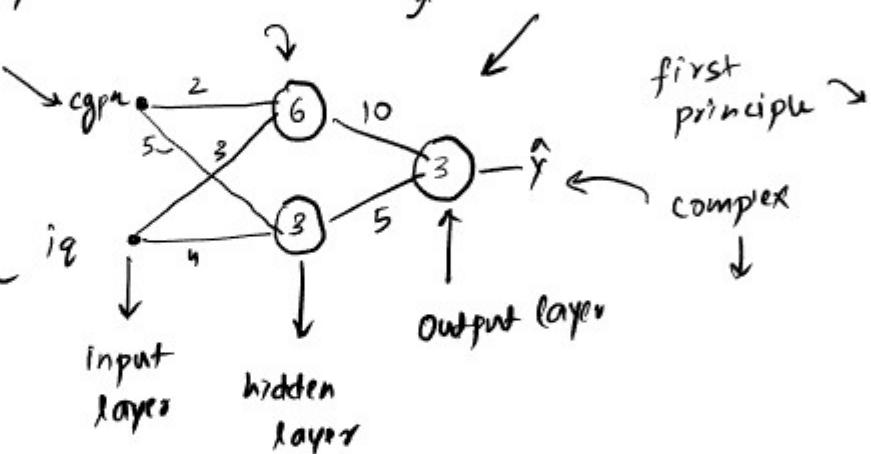
The Idea of MLP

24 February 2022 19:37

cgpa | iq | place

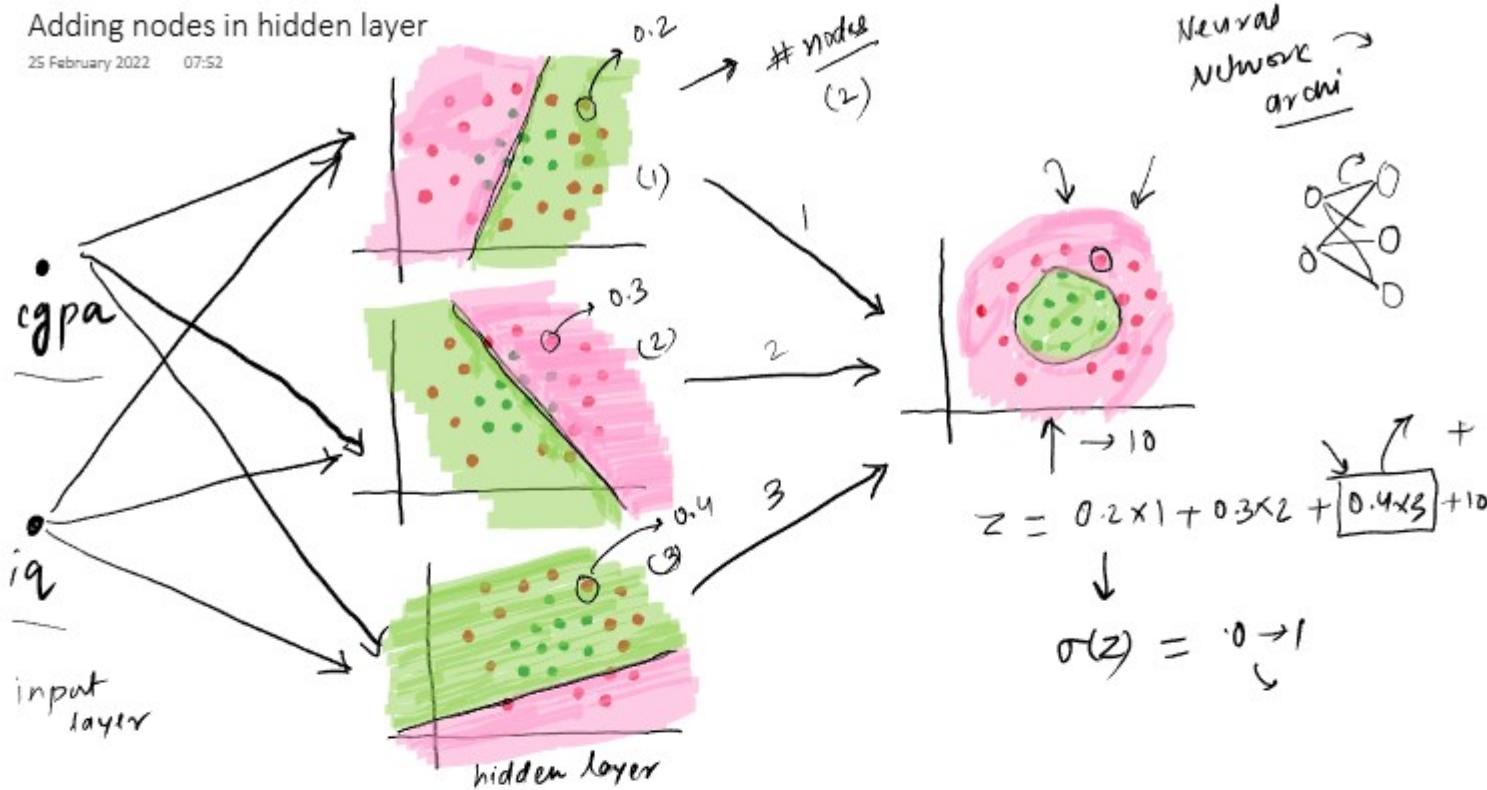


MLP
↓
linear combination
↓
perception



Adding nodes in hidden layer

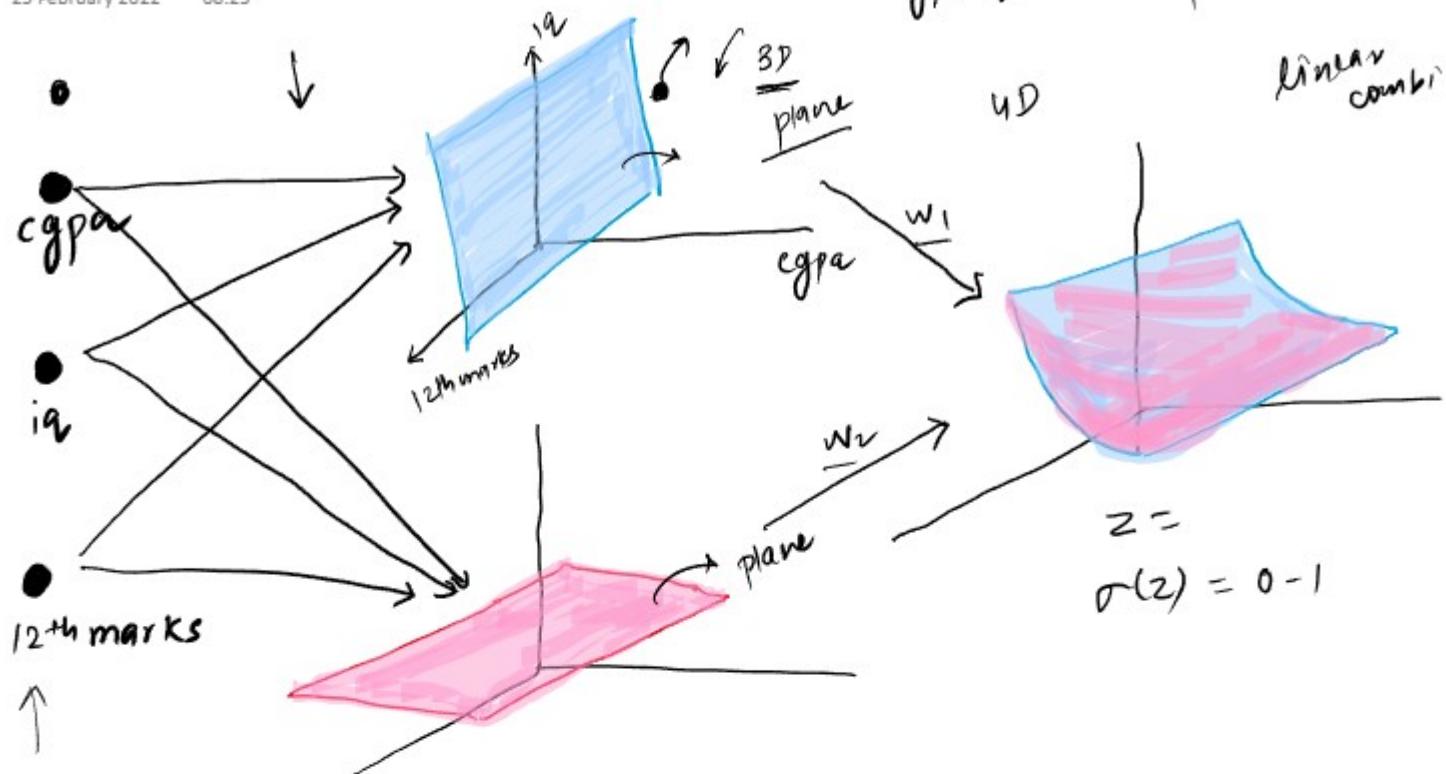
25 February 2022 07:52



Adding nodes in input

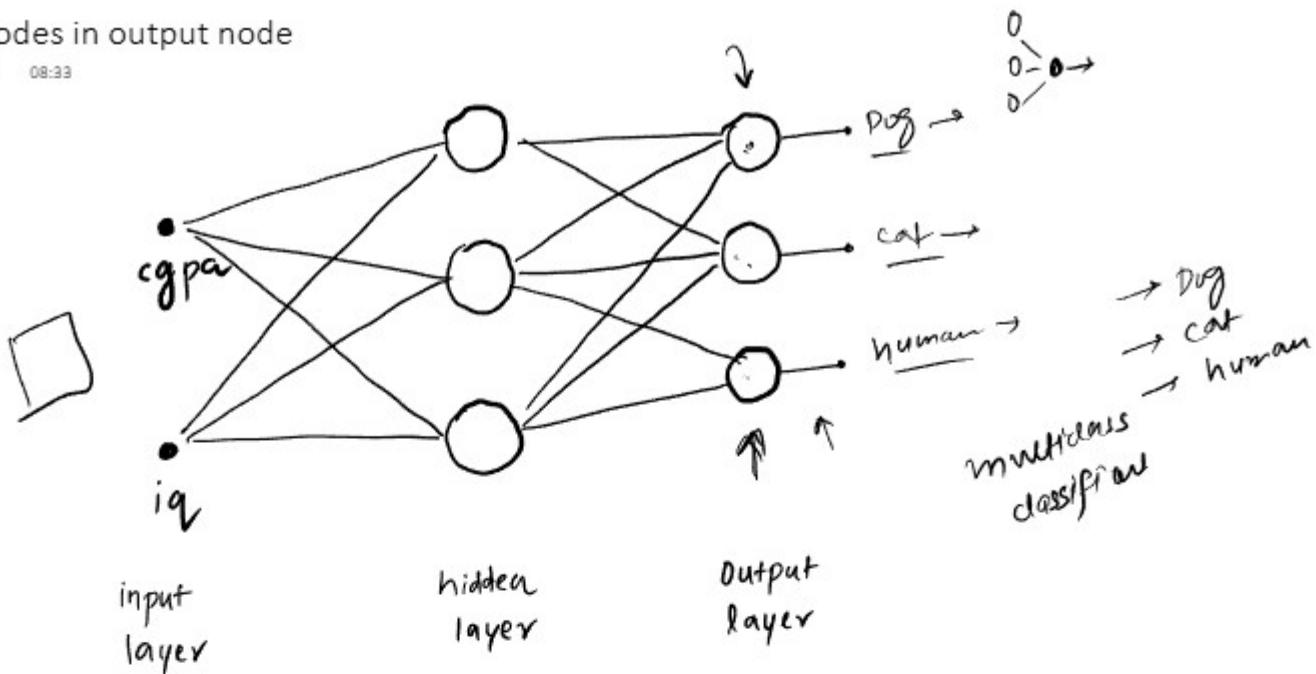
25 February 2022 08:25

cgpa iq 12th marks



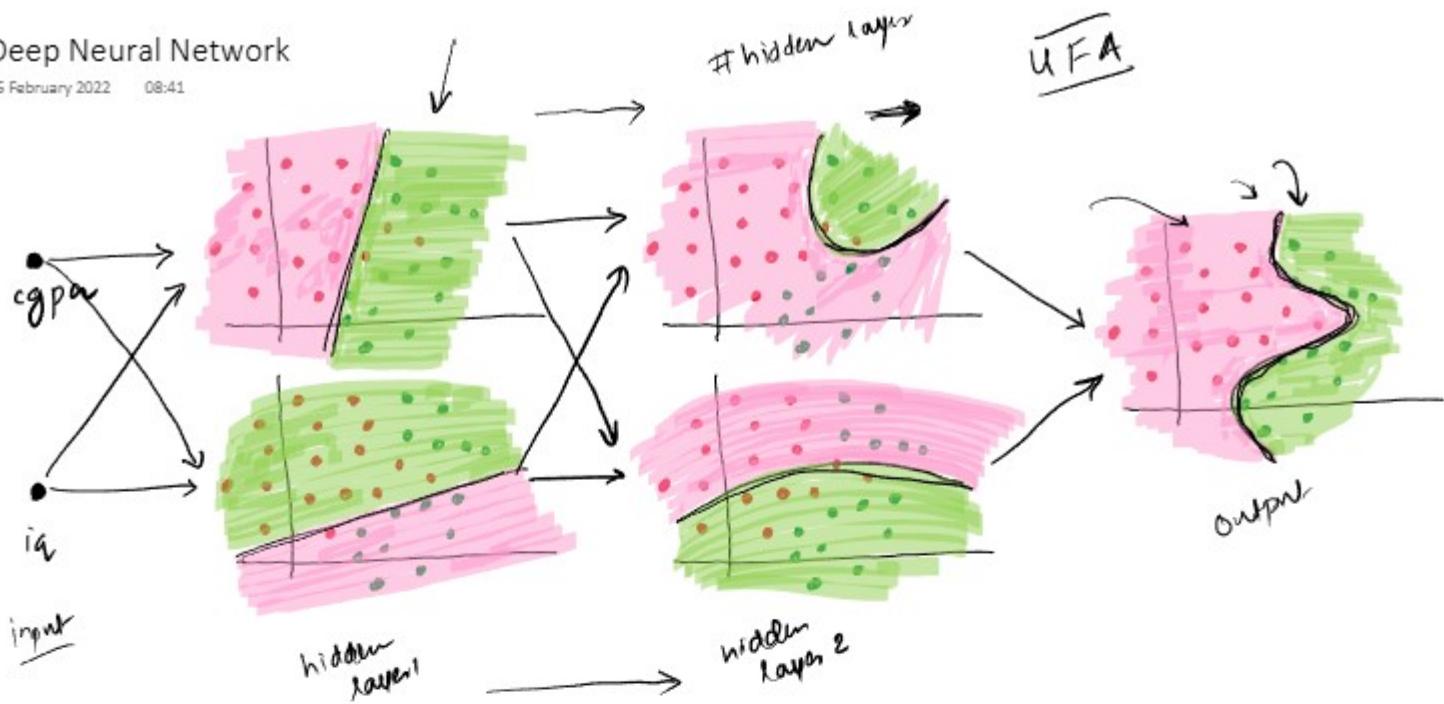
Adding nodes in output node

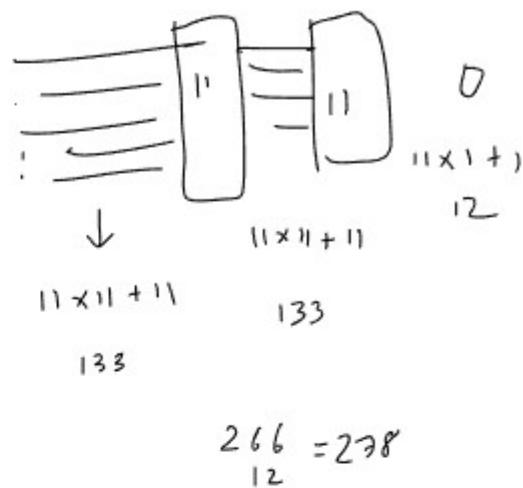
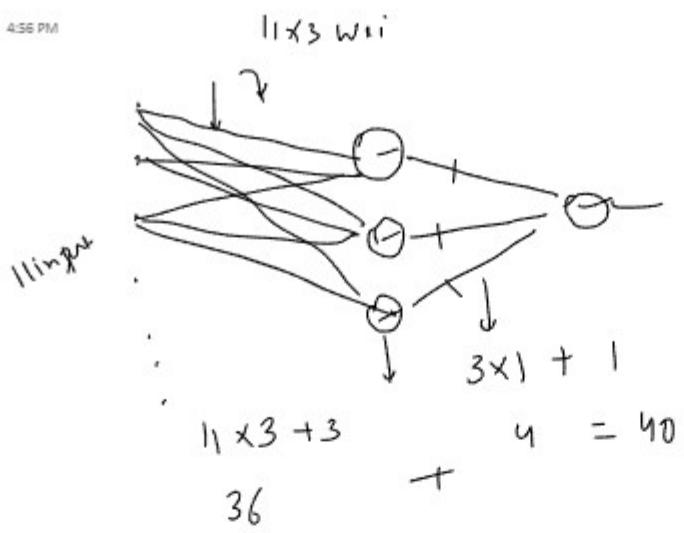
25 February 2022 08:33



Deep Neural Network

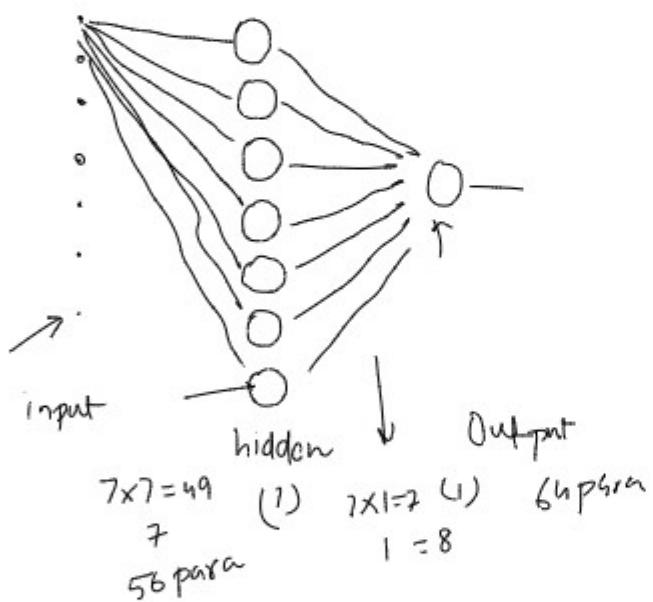
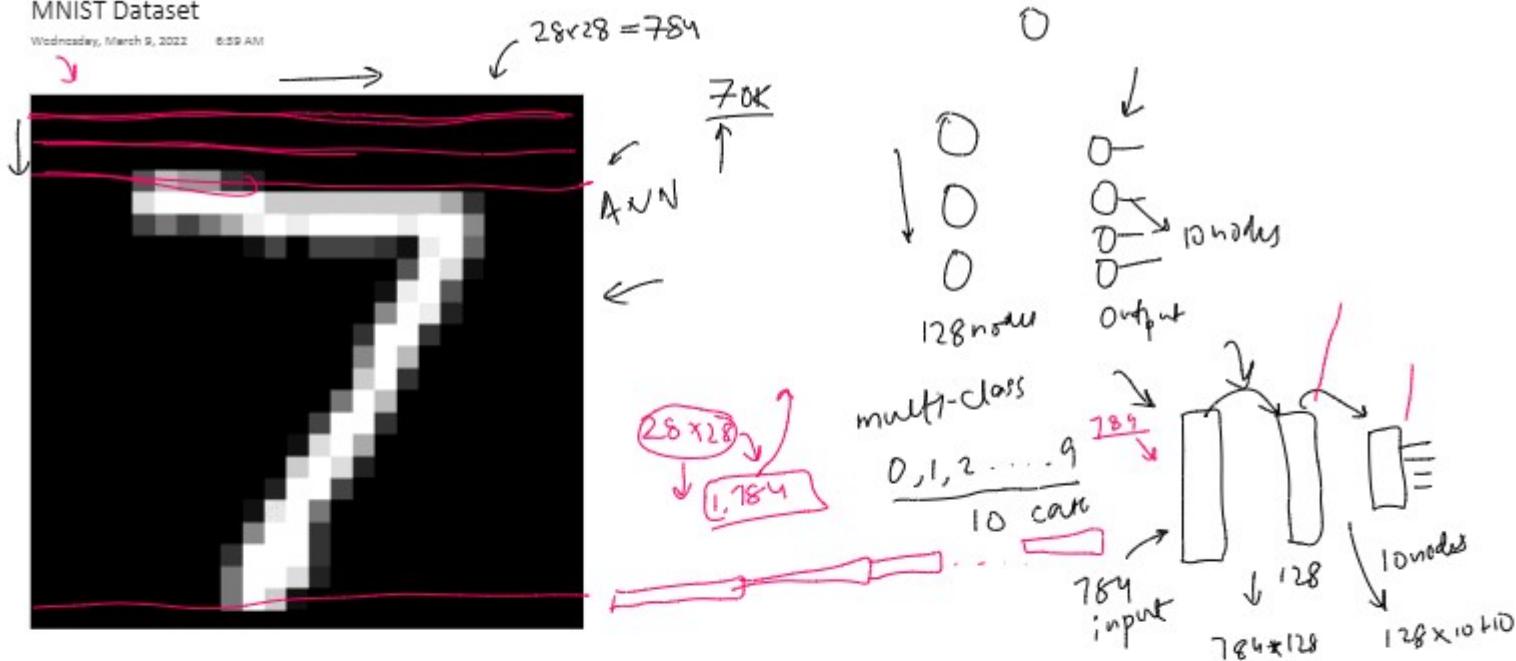
25 February 2022 08:41





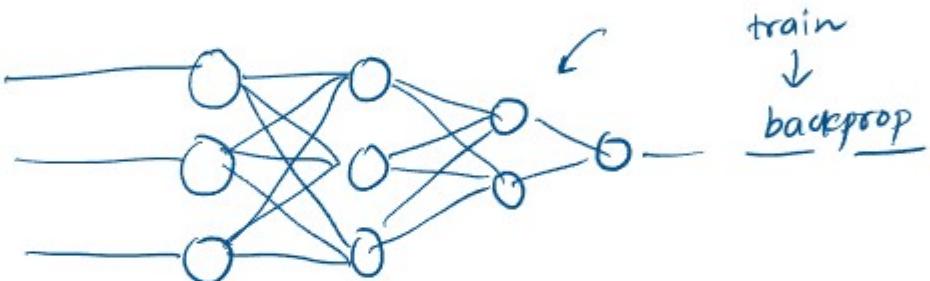
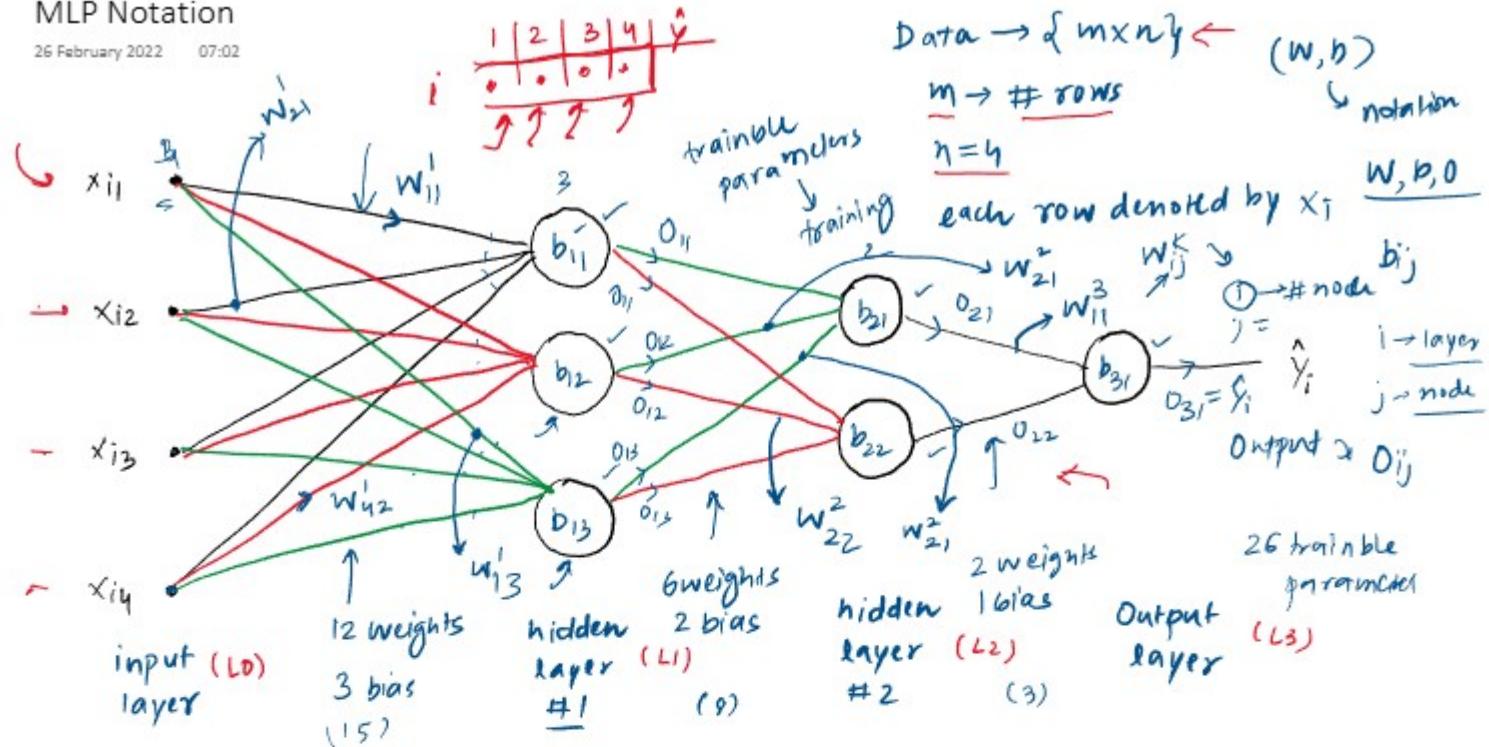
MNIST Dataset

Wednesday, March 9, 2022 8:59 AM



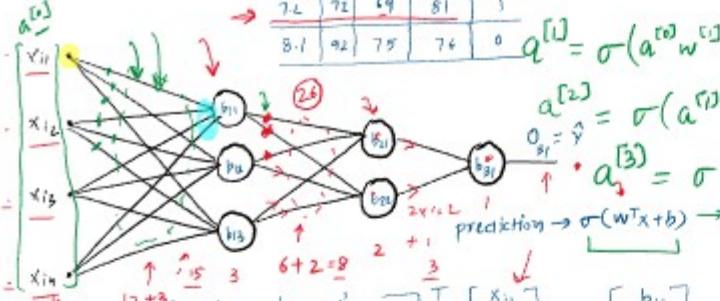
MLP Notation

26 February 2022 07:02



Forward Propagation

28 March 2022 09:07



Forward Propagation placed \hat{y} # of trainable parameters

28 Mar 2022	09:07	19	104m	116m	1
		7.2	72	69	81
		8.1	92	75	76
					0

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

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

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

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

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

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

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

$$= \begin{bmatrix} W_{11}^1 x_{11} + W_{21}^1 x_{12} + W_{31}^1 x_{13} + W_{41}^1 x_{14} \\ W_{12}^1 x_{11} + W_{22}^1 x_{12} + W_{32}^1 x_{13} + W_{42}^1 x_{14} \\ W_{13}^1 x_{11} + W_{23}^1 x_{12} + W_{33}^1 x_{13} + W_{43}^1 x_{14} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$= \sigma \left(\begin{bmatrix} W_{11}^1 x_{11} + W_{21}^1 x_{12} + W_{31}^1 x_{13} + W_{41}^1 x_{14} + b_{11} \\ W_{12}^1 x_{11} + W_{22}^1 x_{12} + W_{32}^1 x_{13} + W_{42}^1 x_{14} + b_{12} \\ W_{13}^1 x_{11} + W_{23}^1 x_{12} + W_{33}^1 x_{13} + W_{43}^1 x_{14} + b_{13} \end{bmatrix} \right)$$

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

Layer #2

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

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

$$= \begin{bmatrix} 3 \times 2 & 2 \times 3 & 3 \times 1 \\ & & 2 \times 1 \end{bmatrix}$$

Layer #3

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

$$= \begin{bmatrix} W_{11}^3 \\ W_{21}^3 \end{bmatrix}^T \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \end{bmatrix}$$

$$= \begin{bmatrix} 2 \times 1 & 2 \times 2 & 2 \times 1 \\ & & 1 \times 1 \end{bmatrix}$$

What is Loss Function?

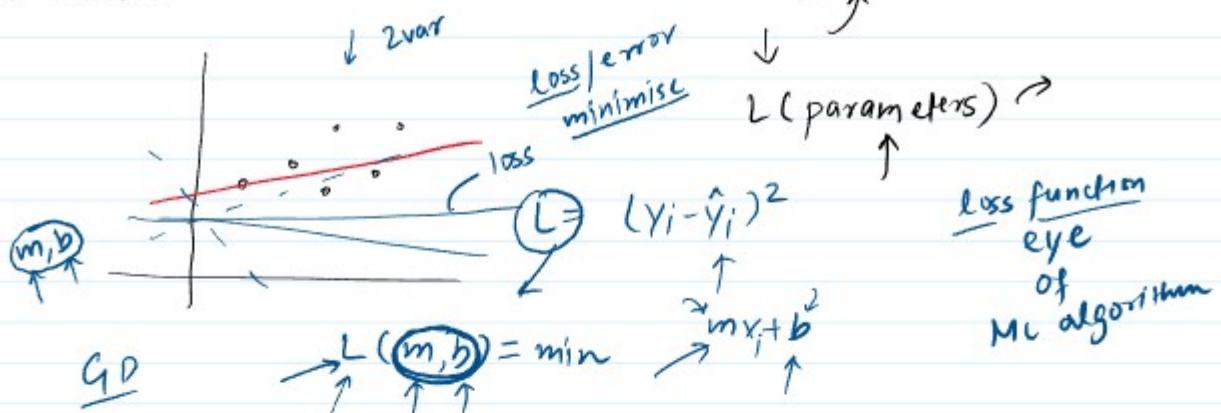
23 March 2022 07:15

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

high \rightarrow poor

small \rightarrow great

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



Why is Loss function important?

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

Peter Drucker

multiple epochs

Loss Function in Deep Learning?

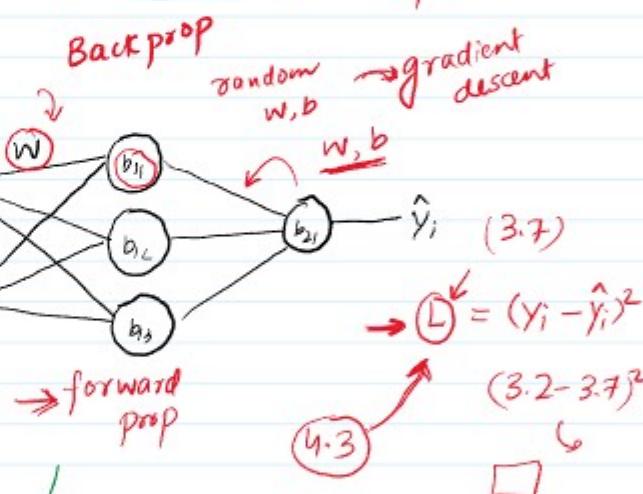
✓ cgpa | iq | package(epa)

7.1	83
8.5	91
6.3	102
5.1	87

3.2
4.5
6.1
2.7

7.1 cgpa

83 iq



Keras

- object detection
- face loss
- Regression

{ mse
mae
huber loss

Loss functions in DL

Classification

- + binary crossentropy
- + categorical cross entropy
- + hinge loss

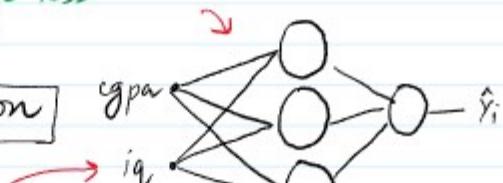
Autoencoders
+ KL divergence

Embedding
Triplet loss

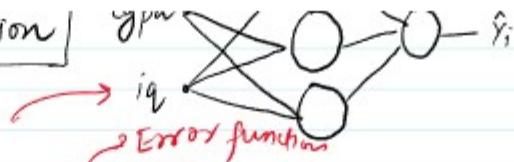
GAN

- + discriminator loss
- + minmax gan loss

Loss Function vs Cost Function



Loss Function vs Cost Function



Loss function → single training eg

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

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

batch Cost function

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

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

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

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

$$(6.3 - 6.1)^2$$

1. Mean Squared Error (MSE)

squared loss L2 loss

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

$$(true - predict)^2$$

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

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

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

Advan
DBAdv
mae
 $(y_i - \hat{y}_i)^2$
quadratic
punish

6.3	100
7.1	91
8.5	83
9.2	102
8.1	100

6.3
4.1
3.5
7.2
7

6.1
4
3.7
7

$$\left\{ \begin{array}{l} 0.2 \\ \rightarrow 0.1 \\ \downarrow -0.2 \\ 0.2 \end{array} \right\}$$

$$\left\{ \begin{array}{l} \text{overall} \\ \text{error} \end{array} \right\}$$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

cgpa
1g
package
Prediction
$\{(y_i - \hat{y}_i)^2\}$

Advantages

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

Disadvantage

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

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



$$[mse]$$

2. Mean Absolute Error (MAE) → L1 loss

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

$$2|abs|$$

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

$$|true - predict|$$

$$CGPA \quad 1g \quad package \quad Prediction \quad |(y_i - \hat{y}_i)|$$

$$L1 \quad l1$$

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

true - predict
punish

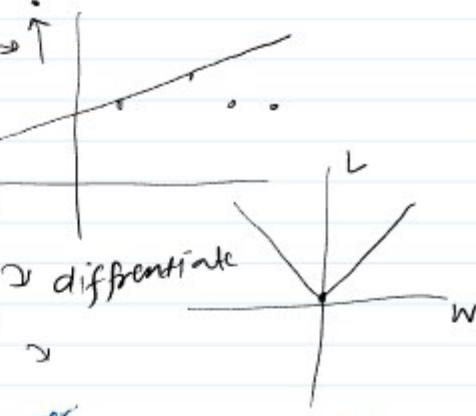
Advantages

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

Disadvantage

- 1) Not differentiable

\downarrow GD \approx differentiable
Subgradient

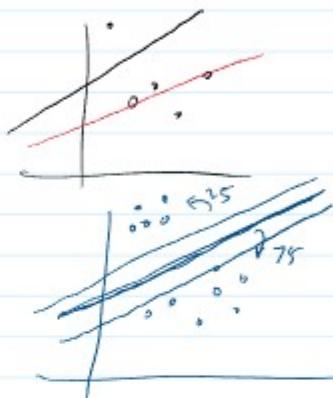


3. Huber Loss ✓

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

mae mse

mse - outliers ✓
mae - normal points ✓



4. Binary Cross Entropy \rightarrow log loss

- classification ✓
- Two classes ✓

100 days
cgpa | iq | placement

8	80
7	70
6	60
5	50

1
0
0
0

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

Keras

$$0.12 = -0.73 \log(0.73) - 0.25 \log(0.25)$$

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

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

$$-1 \times -0.13 = 0.13$$

maximum likelihood

Logistic Reg

Advantage

- * Differentiable

Disadvantage

- multi local minimum
- Intuitive

5. Categorical Cross Entropy [used in Softmax Regression]

- multi-class classification ✓

OHE

cgpa	{ iq }	{ placed? }
?	?	Yes
?	?	No
?	?	Maybe

5. Categorical Cross Entropy [Used in softmax regression]

→ Multi-class classification

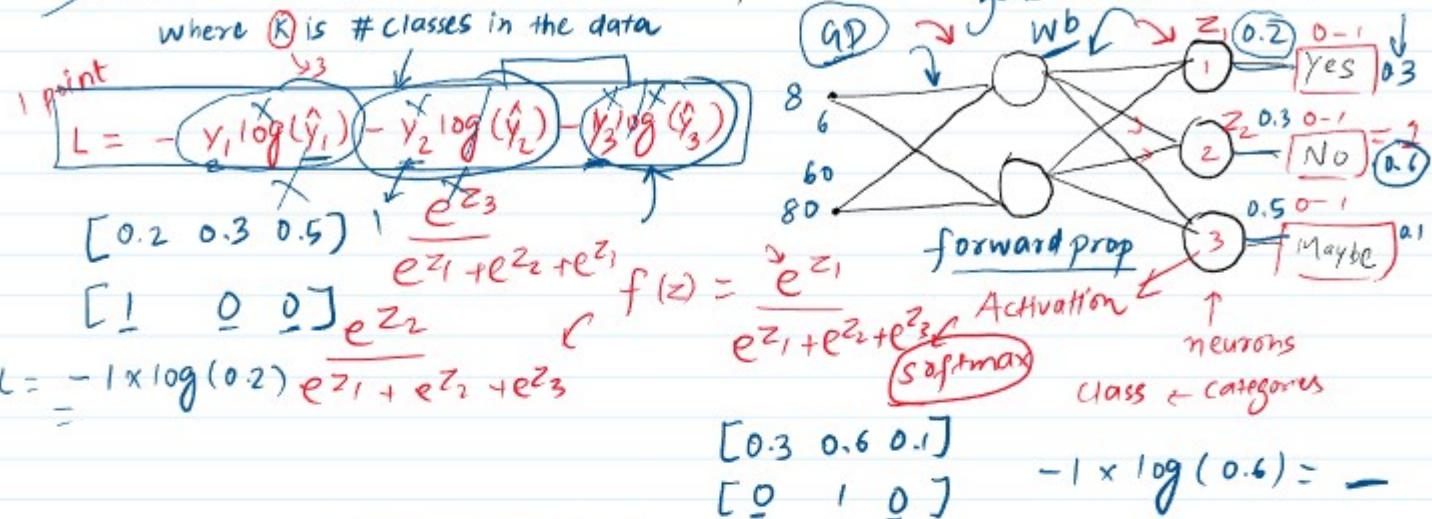
1 point

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

Where K is # classes in the data

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

Yes	No	Maybe
1	0	0
0	1	0
0	0	1



Sparse Categorical Cross Entropy

(L) ~

SCE

Loss function
fast

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

Cost function

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

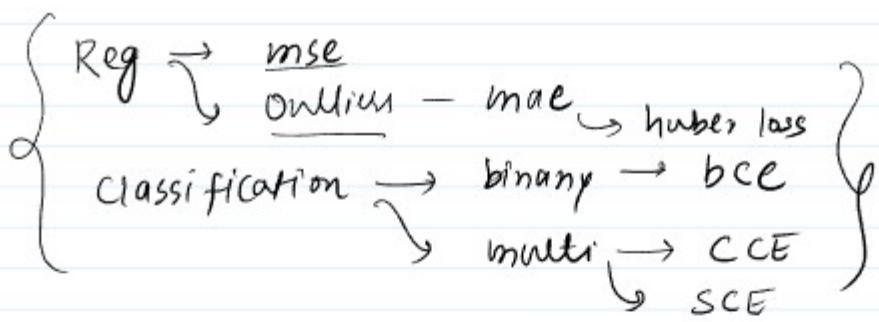
DHE

cgpa iq	placed	1
7	70	Yes
8	80	No
6	60	Maybe

$$[0.1 \ 0.4 \ 0.6]$$

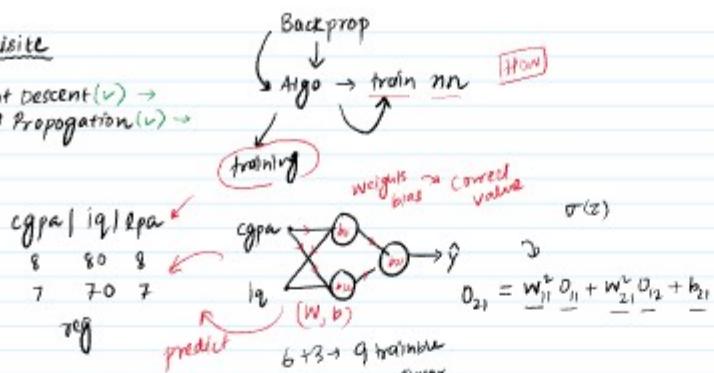
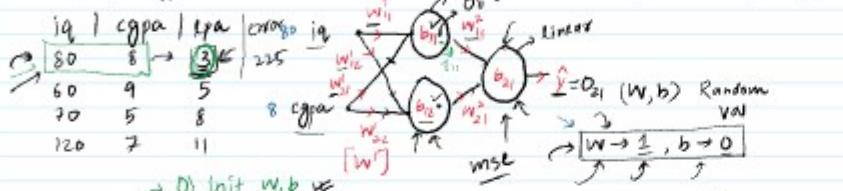
$$[0.6 \ 0.2 \ 0.2]$$

OHE - fast
SCE



Pre requisite

- Gradient descent (v) →
- Forward Propagation (v) →

**Backprop** 4 students

- Steps**
- 1) Init w, b v
 - 2) You select a point (row) ↳ student
 - 3) Predict (\hat{y}_{pa}) → forward prop [out predict]
 - 4) Choose a loss function (mse)
 - 5) Weights and bias update v ↳ gradient descent
- $L = (y - \hat{y})^2$
- $(3 - 18)^2 = \frac{225}{e^{-100}}$

Backprop

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

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

$w_{11}^T \text{ new} = w_{11}^T \text{ old} - \eta \frac{\partial L}{\partial w_{11}^T}$ ↳ **q derivative**

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

$w_{21}^T \text{ new} = w_{21}^T \text{ old} - \eta \frac{\partial L}{\partial w_{21}^T}$

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

derivative of loss wrt weight

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

↳ partial derivative

derivative

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

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

↳ chain rule of diff.

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

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

$\frac{\partial \hat{y}}{\partial w_{11}^T} = \frac{\partial}{\partial w_{11}^T} [O_{11} w_{11}^T + O_{12} w_{21}^T + b_{21}] = O_{11}$

$\frac{\partial L}{\partial w_{11}^T} = -2(y - \hat{y}) O_{11}$

$$\frac{\partial \hat{y}}{\partial w_{21}^T} = \frac{\partial}{\partial w_{21}^T} [O_{11} w_{11}^T + O_{12} w_{21}^T + b_{21}] = O_{12}$$

$\frac{\partial L}{\partial w_{21}^T} = -2(y - \hat{y}) O_{12}$

$\frac{\partial L}{\partial b_{21}} = \frac{\partial}{\partial b_{21}} \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial b_{21}} = 1$

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

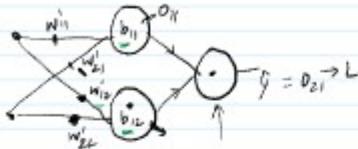
$y \rightarrow b$

$\hat{y} \rightarrow b$

$$\frac{\partial \hat{y}}{\partial b_{21}} = 1$$

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

$$\frac{\partial L}{\partial b_{21}} = \left[\frac{\partial L}{\partial \hat{y}} \right] \times \frac{\partial \hat{y}}{\partial b_{21}}$$



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

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

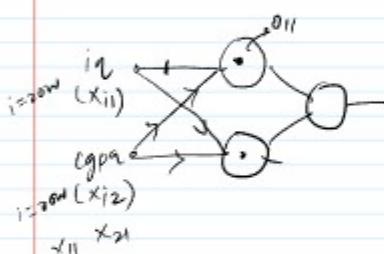
$$\frac{\partial L}{\partial b_{11}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \delta_{11}} \frac{\partial \delta_{11}}{\partial b_{11}} \rightarrow -2(y - \hat{y}) w_{11}^2$$



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

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

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



$$\frac{\partial \delta_{11}}{\partial w_{11}} = \underbrace{\frac{\partial \hat{y}}{\partial w_{11}}}_{w_{11}^2} + \underbrace{\frac{\partial \hat{y}}{\partial w_{21}}}_{w_{21}^2} + \underbrace{\frac{\partial \hat{y}}{\partial b_{11}}}_{w_{11}^2} = (iq)$$

$$\frac{\partial \delta_{11}}{\partial w_{21}} = x_{11}$$

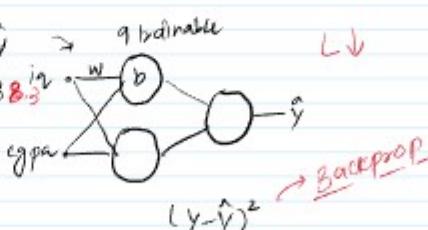
$$\frac{\partial \delta_{11}}{\partial b_{11}} = 1$$

$$\frac{\partial \delta_{12}}{\partial w_{12}} = \frac{\partial \delta_{11}}{\partial w_{12}} \left[iq w_{12}^2 + cgpa w_{22}^2 + b_{12} \right] = iq(x_{12})$$

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

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

multiple choice	
w_{11}	iq
w_{12}	$cgpa$
b_{11}	\hat{y}
b_{12}	L



steps (once again)

0) weights/bias → init $w=0$

loop - 100/100

1) for i in range(100):

1a) 1 student → forward prop → $l_{pq} \rightarrow$ 1/2 min

1b) Loss calculate (mse) → .

1c) Adjust all weights and bias $\frac{\partial L}{\partial w}$.

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

The How

Backpropagation Algorithm

epochs = 5

for i in range(epochs):

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

$$\frac{\partial L}{\partial w_{11}} = -2(y - \hat{y}) \delta_{11} \times \text{Row}$$

$$\frac{\partial L}{\partial w_{21}} = -2(y - \hat{y}) \delta_{12}$$

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

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

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

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

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

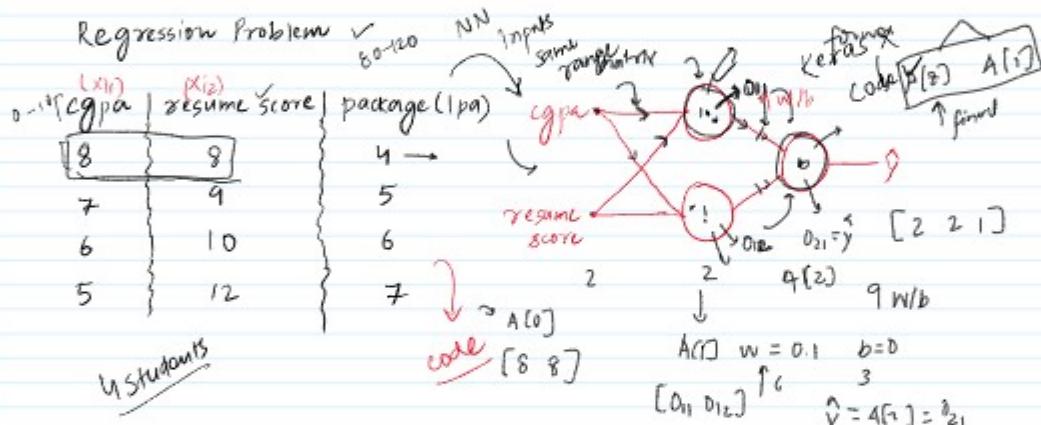
$$\frac{\partial L}{\partial w_{12}} = -2(y - \hat{y}) w_{12}^T x_{12} \quad \checkmark$$

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

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

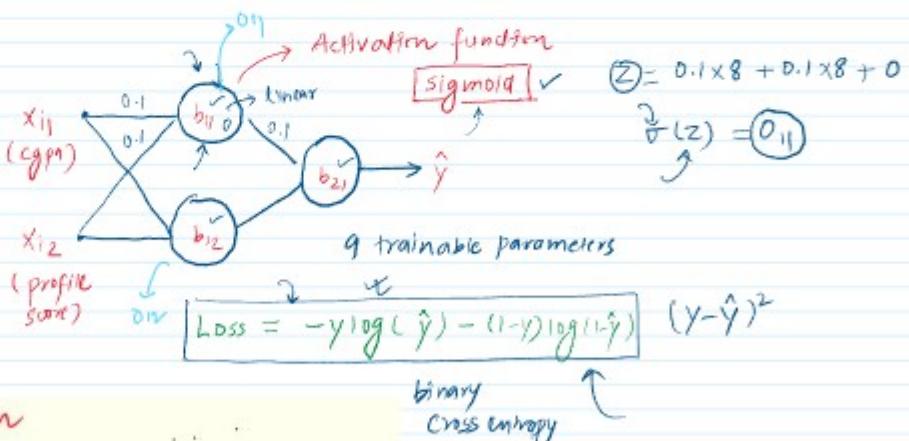
$$\frac{\partial L}{\partial w_{22}} = -2(y - \hat{y}) w_{22}^T \quad \checkmark$$

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



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) $\rightarrow lby$

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

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

 → 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 1

→ 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

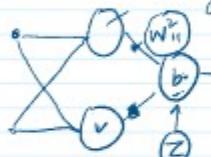
$$z = W_{11}^2 \delta_{11} + W_{21}^2 \delta_{21} + b_{21}$$

→ Calculate avg loss for the epoch

Lavg

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

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



der w11 L

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

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

$$z_{pred} = W_{11}^2 O_{11} + W_{21}^2 O_{12} + b_{21}$$

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

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

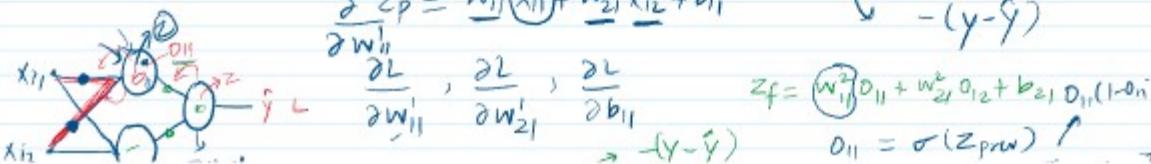
$$= -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} = \frac{-y(1-\hat{y}) + \hat{y}(1-y)}{\hat{y}(1-\hat{y})} = \frac{-y + y\hat{y} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})}$$

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

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

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

$$\frac{\partial z_p}{\partial W_{11}} = W_{11}^2 x_{11} + W_{21}^2 x_{12} + b_{21} \quad -(y - \hat{y})$$



$$z_f = (W_{11}^2 O_{11} + W_{21}^2 O_{12} + b_{21} O_{11})(1-\hat{y})$$

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

cross entropy

$$z_f = w_{11}^2 z_{f1} + w_{21}^2 z_{f2} + b_{21}$$

$$z_{f1} = w_{11}^1 x_{i1} + b_{11}$$

$$z_{f2} = w_{21}^1 x_{i2} + b_{21}$$

$$\frac{\partial L}{\partial w_{11}^1} = \left[\frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_f} \right] \times \frac{\partial z_f}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial z_{\text{prev}}} \times \frac{\partial z_{\text{prev}}}{\partial w_{11}^1} \rightarrow x_{i1}$$

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

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

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

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

$$z_f = w_{12}^2 z_{f1} + w_{22}^2 z_{f2} + b_{22}$$

$$z_{f1} = w_{12}^1 x_{i1} + b_{12}$$

$$\frac{\partial L}{\partial w_{12}^1} = \left[\frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_f} \right] \times \frac{\partial z_f}{\partial o_{12}} \times \frac{\partial o_{12}}{\partial z_p} \times \frac{\partial z_p}{\partial w_{12}^1} \rightarrow x_{i1}$$

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

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

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

$$\frac{\partial L}{\partial b_{12}} = -(y - \hat{y}) w_{12}^2 o_{12} (1 - o_{12}) \rightarrow x_{i1}$$

Backpropagation → The Why?

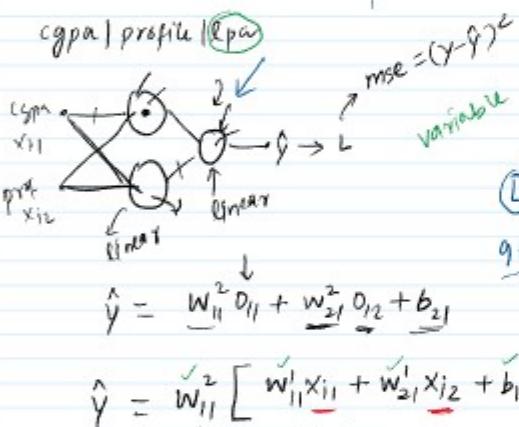
The intuition behind the algorithm ✓

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
             $W_n = W_0 - \frac{\partial L}{\partial W}$  rule
            magic
        → Calculate avg loss for the epoch
             $L_{avg}$ 
    
```

→ Loss function is a function of all trainable parameters



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

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

$L(\cdot)$ → $L(y) = \frac{1}{2} y^2$

$L(w, b) \downarrow$

9 things → w, b

$\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$

→ Concept of Gradient \rightarrow 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 $\frac{\partial}{\partial x}$

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

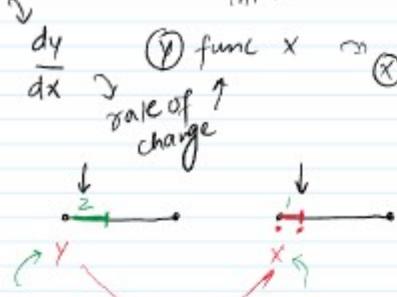
9 params $\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$

9D function \rightarrow 9 diff steps wrt each dim

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



→ Concept of Derivative



↑ intuition
at a point

$$\frac{dy}{dx} \quad \textcircled{⑦} \text{ func } x \quad \textcircled{⑧} \text{ rate of change}$$

$$\frac{dy}{dx} = \underline{\underline{\underline{z}} \text{ (+ve)}} \quad \text{mag}$$

$$\frac{dy}{dx} = \underline{\underline{\underline{-z}} \text{ (-ve)}} \quad \text{sign}$$

$$x \rightarrow \quad y \rightarrow$$

derivative

$$y = x^2 + 2x$$

$$\frac{dy}{dx} = (2x+1)_{x=5}$$

$$\left(\frac{dy}{dx} \right)_{x=5} = \textcircled{⑨} \text{ slope}$$

$$\frac{\partial L}{\partial w_{ii}}$$

$$w_{ii} \rightarrow 1 \text{ unit}$$

$$L \rightarrow ? \quad \begin{matrix} \text{(mag)} \\ \text{(sign)} \end{matrix}$$

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

→ the concept of minima



$$y = x^2$$

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

$$\boxed{x=0}$$

$$\begin{cases} z = x^2 + y^2 \\ (x,y) \rightarrow z \text{ min} \end{cases}$$

$$\frac{\partial z}{\partial x} = 0 \rightarrow x \rightarrow (x,y) \text{ (0,0)}$$

$$\frac{\partial z}{\partial y} = 0 \rightarrow y \rightarrow 2y = 0 \Rightarrow y = 0$$

$$L(\text{param})$$

$$(w, b) \min$$

$$L \downarrow$$

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

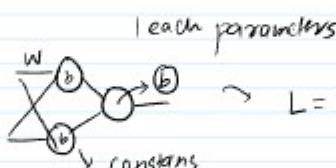
$$\begin{matrix} \text{qdim} \\ \text{minima} \end{matrix}$$

→ Backprop Intuition

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

qstep

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



$$L = L(w, b)$$

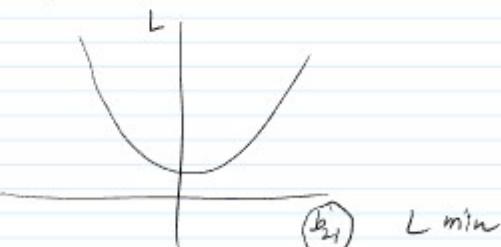
$$\downarrow L = (y - \hat{y})^2$$

$$L(b_{21}) \quad \hat{y} = (b_{21})$$

$$\frac{L(b_{21})}{T} \quad L = b_{21}^2$$

$$\begin{matrix} b_{21} = \boxed{b_{21}} \\ \text{smart} \\ \text{L derivative w.r.t } b_{21} \\ \uparrow ? \end{matrix}$$

$$\frac{\partial L}{\partial b_{21}} \quad \begin{matrix} b_{21} = \textcircled{⑤} w \\ +ve \\ -ve \end{matrix}$$



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

$$L \downarrow$$

$$\frac{\partial L}{\partial b_{21}} = +ve \quad b_{21} \downarrow \quad L \uparrow$$

$$b_{21} \downarrow$$

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

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

$$b_{21} \uparrow$$

$$b_{21} +$$

game

intuit

-ve of the gradient

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

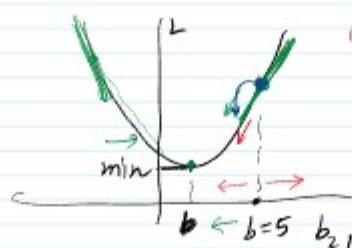
game derivative
(b_2)

-ve of the gradient

$$b_2 = b_1 \boxed{\frac{\partial L}{\partial b_2}}$$

$$\frac{\partial L}{\partial b} = \text{slope}$$

+ve



algo

$$b = -5$$

$$\frac{\partial L}{\partial b} = -\text{ve}$$

-ve

slope(gradient)
-ve

Slope

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

$$\boxed{\eta}$$

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

→ Effect of Learning Rate (η)

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

smooth

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

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

$$b = 45$$

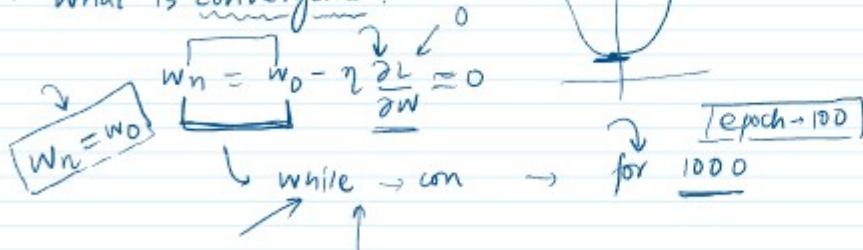
$$\frac{\partial L}{\partial b} = 50$$

$$b = 45 - 50$$

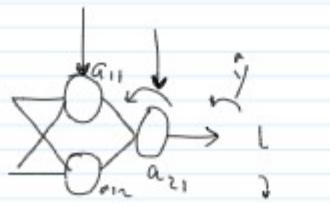
$$b = -5 \quad \frac{\partial L}{\partial b} = -50 \quad \eta = 0.01$$

$$b = -5 + (0.01 \times 50) = 0.5 \\ = -4.5$$

→ What is convergence?



$$i = 1, 0$$



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

activation = ②

$$\hat{y}(1 - \hat{y}) O_{11}$$

$$-y \log(a_{11}) - (1-y) \log(1-a_{11})$$

$$-(y - \hat{y}) O_{11}$$

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

$$\frac{\partial L}{\partial w_{11}} = \boxed{\frac{\partial L}{\partial a_{11}}} \times \boxed{\frac{\partial a_{11}}{\partial w_{11}}}$$

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

$$\sim \eta \sim \frac{\partial L}{\partial w}$$

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

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

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

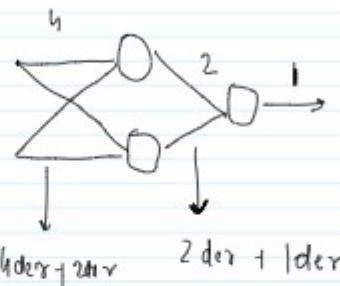
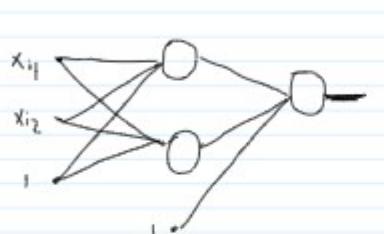
$$W[1] \quad 1 \times 1 \quad 1 \times 2$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{- (y - \hat{y}) a_{21} (1-a_{21}) a_{11}}{-(y - \hat{y}) a_{21} (1-a_{21}) a_{12}}$$

$$\frac{\partial L}{\partial b_{21}} = - (y - \hat{y}) a_{21} (1-a_{21})$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{-(y - \hat{y}) a_{21} (1-a_{21}) a_{11}}{-(y - \hat{y}) a_{21} (1-a_{21}) a_{12}}$$

$$\frac{\partial L}{\partial b_{21}} = - (y - \hat{y}) a_{21} (1-a_{21})$$



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

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

$$\begin{bmatrix} a_{01} \\ a_{02} \end{bmatrix} \quad \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} \quad \begin{bmatrix} a_{21} \end{bmatrix}$$

$$\begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} \quad \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} \quad \begin{bmatrix} W_{11}^2 \\ W_{12}^2 \end{bmatrix} \quad \begin{bmatrix} b_{21} \end{bmatrix}$$

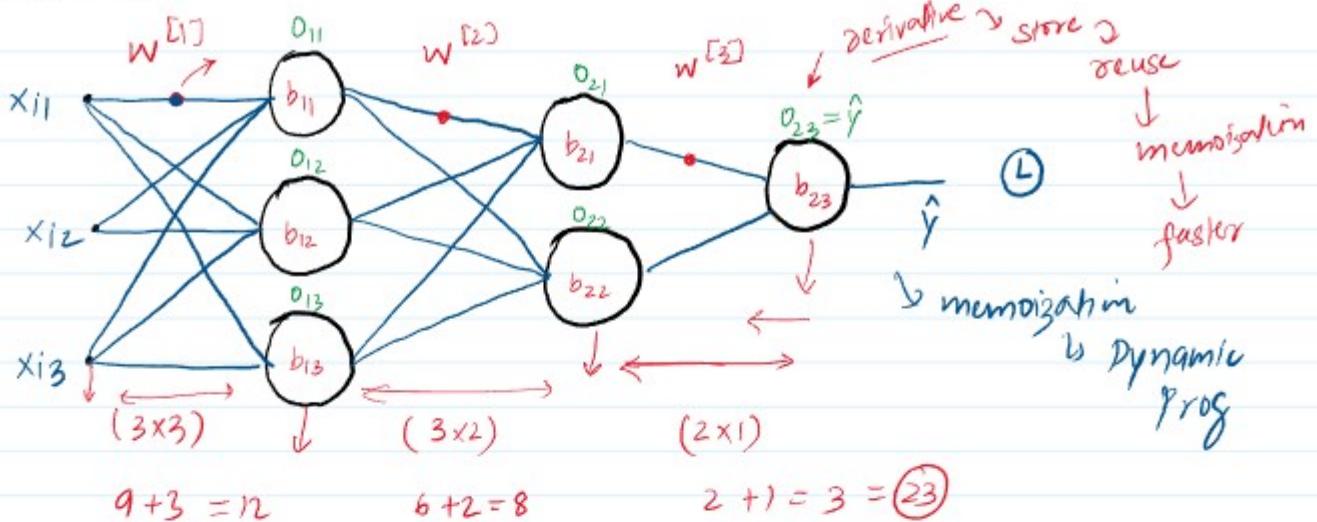
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial a_1} \times \frac{\partial a_1}{\partial w}$$

$$\begin{bmatrix} \frac{\partial L}{\partial w_{11}^1} & \frac{\partial L}{\partial w_{12}^1} \\ \frac{\partial L}{\partial w_{21}^1} & \frac{\partial L}{\partial w_{22}^1} \end{bmatrix} \quad \begin{bmatrix} \frac{\partial L}{\partial w_{11}^2} \\ \frac{\partial L}{\partial w_{12}^2} \end{bmatrix}$$

MLP Memoization

Thursday, April 7, 2022 7:45 AM

backprop update \rightarrow w/b driv



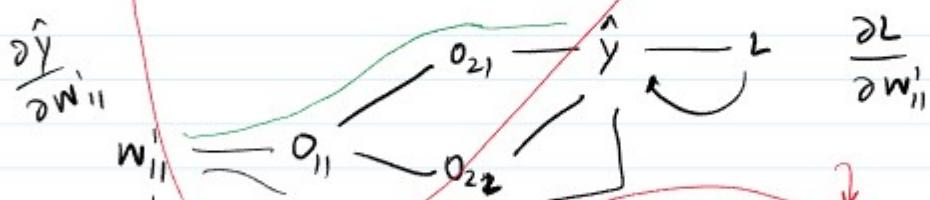
Chain rule

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

$$\frac{\partial L}{\partial w_{11}^3} = \left[\frac{\partial L}{\partial \hat{y}} \right] \times \left[\frac{\partial \hat{y}}{\partial w_{11}^3} \right]$$

$$\frac{\partial L}{\partial w_{11}^2} = \left[\frac{\partial L}{\partial \hat{y}} \right] \times \left[\frac{\partial \hat{y}}{\partial o_{21}} \right] \times \left[\frac{\partial o_{21}}{\partial w_{11}^2} \right]$$

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



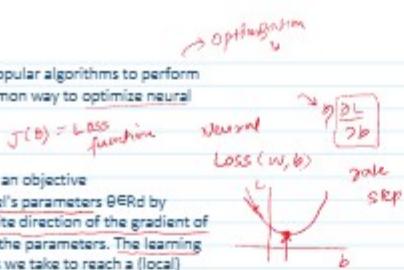
$$\frac{\partial L}{\partial w_{11}^1} = \left[\frac{\partial L}{\partial \hat{y}} \right] \left[\frac{\partial \hat{y}}{\partial o_{21}} \right] \left[\frac{\partial o_{21}}{\partial o_{11}} \right] \times \frac{\partial o_{11}}{\partial w_{11}^1} + \left[\frac{\partial \hat{y}}{\partial o_{22}} \right]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

Backprop \rightarrow Chain diff (rule) + Memoization
maths

Keras
TF

(13)

Gradient descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks.



Gradient descent is a way to minimize an objective function $J(B)$ parameterized by a model's parameters B and by updating the parameters in the opposite direction of the gradient of the objective function $\nabla J(B)$ w.r.t. to the parameters. The learning rate η determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley.

Backpropagation Algorithm

```

epoch = 5
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
    Avg
  
```

{Batch / Stochastic / mini batch}

There are three variants of gradient descent, which 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 between the accuracy of the parameter update and the time it takes to perform an update.

Batch GD (Vanilla GD)

```

for i in range(nb_epochs):
    param_grad = evaluate_gradient(loss_function, data, param)
    param = param - learning_rate * param_grad
  
```

entire dataset \mathcal{D}
update

epoch = # of update

$\frac{\partial L}{\partial W}$ derivative
3 types
accuracy → time → budget.

Stochastic GD vs

```

for i in range(nb_epochs):
    np.random.shuffle(data)
    for example in data:
        param_grad = evaluate_gradient(loss_function, example, param)
        param = param - learning_rate * param_grad
  
```

epoch → 10 (50 steps) 50 times

for i in range(10):
for i in range(x.shape[0]):
shuffle
 $10 \times 50 = 500$ update

↳ 1 random point
↳ $y_{\text{hat}} \rightarrow \text{forward}$
↳ loss
↳ w, b update → $W_n = W_0 - \eta \frac{\partial L}{\partial W}$

avg loss print → for the epoch

frequency of weight update higher

Batch GD → 1 loop

epoch = 10 epoch
for i in range(10):

$y_{\text{hat}} = \text{np.dot}(x, w) + b$
50 values
dot → smart replacement → loops
→ total → 10 times
w, b update

↳ vectorization ~ faster loop

50 points

1 epoch

↳ 50 times

w, b update

Batch GD

epoch = 5

↳ 50 points → predict \mathcal{J}

dot product

$y_{\text{hat}} = \text{np.dot}(x, w) + b$

↳ 50 predict

single $y = 50$ actual out.

y y_{hat} \mathcal{J} loss

↳ $f = \sum_{i=1}^{50}$

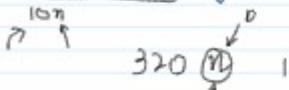
5 epochs

↳ 5 times

w, b update

50 values dot → smart replacement → w/o
 $y = 50 \text{ values}$ ↳ vectorization ~ faster loop
 $y\text{-hat}, y \rightarrow \text{loss}$
 $\hookrightarrow w, b \text{ update } W_n = W_0 - \eta \frac{\partial L}{\partial W}$ Optimized!
 → loss

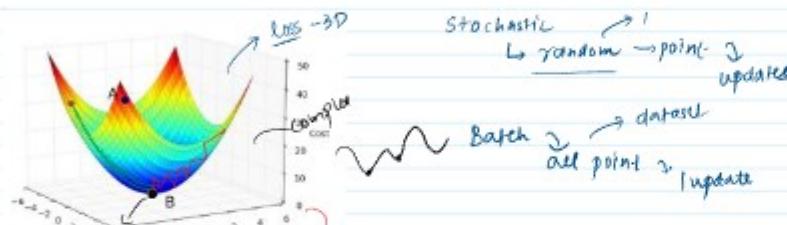
Which is faster (given same no. of epochs)



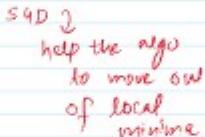
Which is the faster to converge (given same # epochs)



Loss function → 2 algo



Stochastic SGD useful



SGD ↳ help the algo to move out of local minima

Exact solution
approximate diff
lower → loops

Vectorization ↳ big dataset
 $\hookrightarrow np.\text{dot}(x, w) + b$ RMS
 $x \rightarrow \text{dataset of 10 rows}$

Mini Batch Gradient Descent ↳ us
 SGD → MBGD

MBGD ↔ SGD Best of both worlds
 ↳ 320 rows ↳ batches
 in every epoch n 10 batches ↳ 10 update

```
for i in range(n_epochs):
    np.random.shuffle(data)
    for batch in get_batches(data, batch_size=10):
        parson_grad = evaluate_gradient(loss_function, batch, parsons)
        parsons = parsons - learning_rate * parson_grad
```

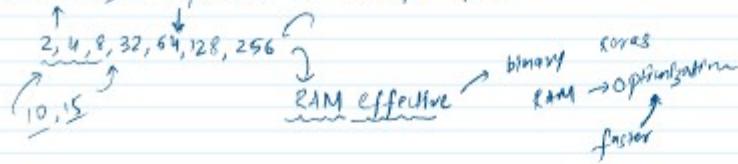
$\frac{n}{x} = \# \text{ of batches}$
 Keras batch_size = x
 $\# \text{ updates/epoch}$

MBGD
 $bgd > mbgd > sgd$
 con
 $bgd < mbgd < sgd$

for i in epochs →
 for j in num of batch
 1 batch
 np.dot ↳ y-pred / vector
 ↳ loss
 ↳ update

spiky SGD
 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
↑ ↑ ↑
1st batch 2nd 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. As one example of the problem, here, traditional activation

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

2) Deep NN $\rightarrow \boxed{\square} \boxed{\square} \boxed{\square} \boxed{\square}$

3) Sigmoid / tanh $\rightarrow A_f \rightarrow$
 deep $x_1, x_2 \rightarrow$ sigmoid $f - p = L$
 $L \rightarrow \hat{y} \rightarrow 0_{11} \rightarrow W_{11}$

0-2

Backprop

0-05

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}}$$

$0 < x < 1$ $0 < x < 1$ 0.1

0.0001

did not change

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

$w_0 = 1$ small change

derivative of L wrt weight

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

$$w_n = 0.9999 \approx \text{vanishingly small}$$

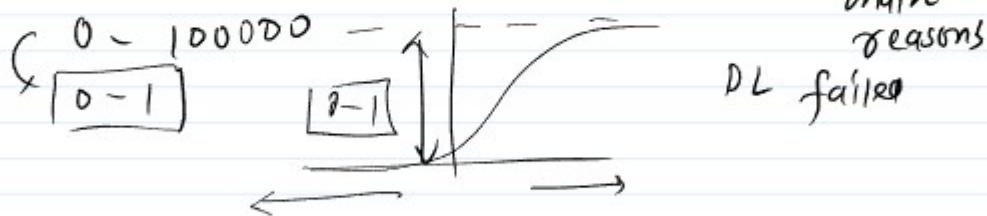
Sigmoid
tanh

$$\begin{array}{ccccccc} & 0 & 0 & 0 & 0 & \frac{\partial L}{\partial w_{11}} = 0.000001 \\ \begin{array}{c} \nearrow \\ \searrow \end{array} & 0 & 0 & 0 & 0 & \nearrow \\ & 0 & 0 & 0 & 0 & \end{array}$$

Backprop 2

model training

$\rightarrow x$



How to recognize?

1) Loss focus \rightarrow epoch \rightarrow no changes \rightarrow VGP

2) weights \rightarrow graph value \rightarrow Tensorboard \rightarrow callbacks

w_{11}

Keras

loss after epoch

epoch

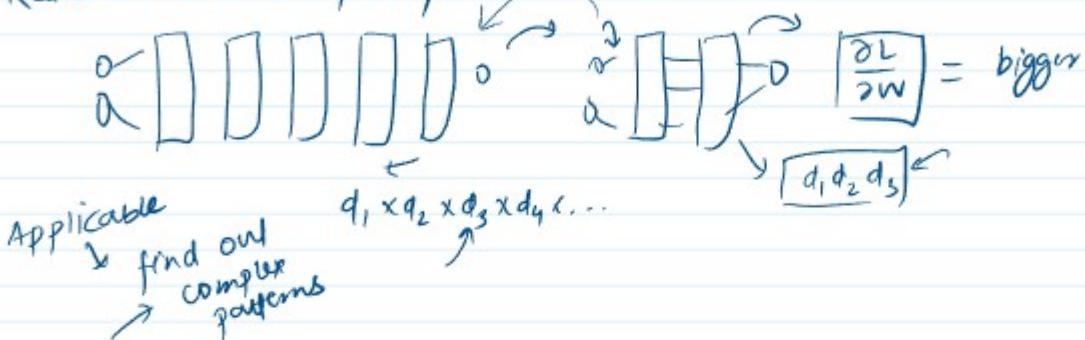
$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

$$\left| \frac{w_0 - w_n}{\eta} \right| =$$

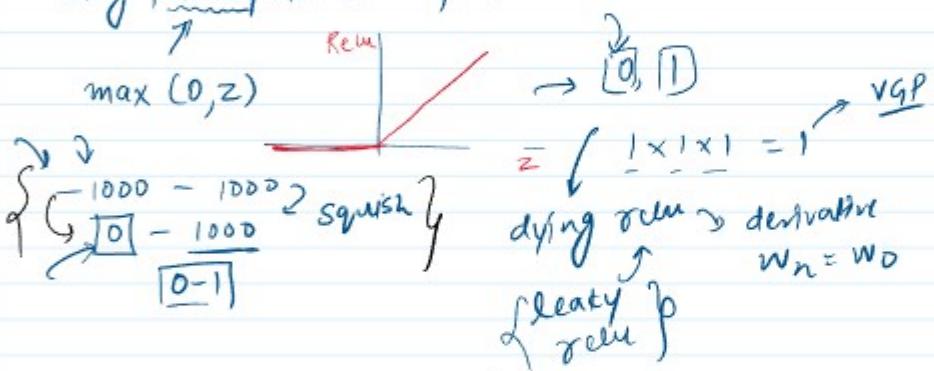
$$w_n = w_0 - \eta \left[\frac{\partial L}{\partial w} \right] \rightarrow \boxed{\frac{w_0 - w_n}{\eta}} =$$

How to handle Vanishing Gradient Problem →

1) Reduce model complexity ↗



2) Using ReLU Activation functions



3) Proper weight init ↗ Glorot ↗ Xavier

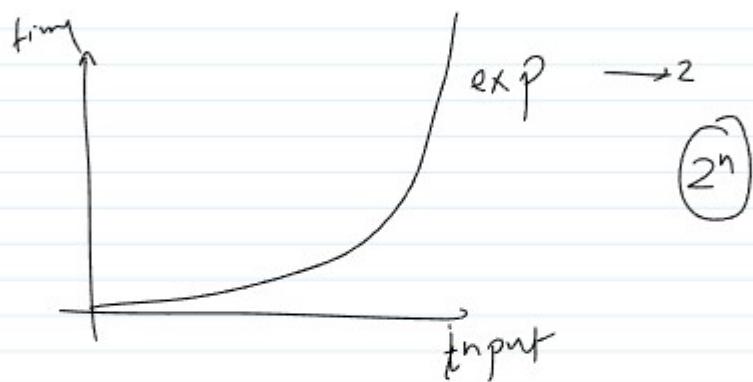
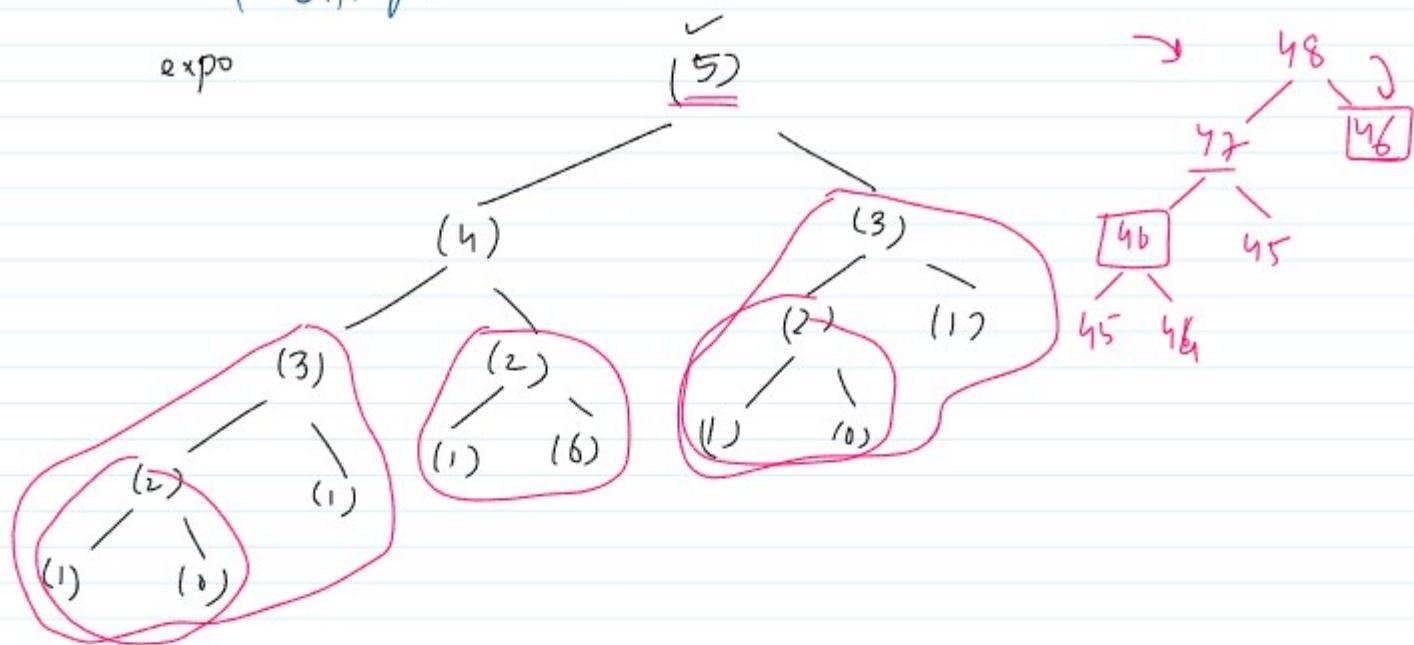
4) Batch norm → layer →

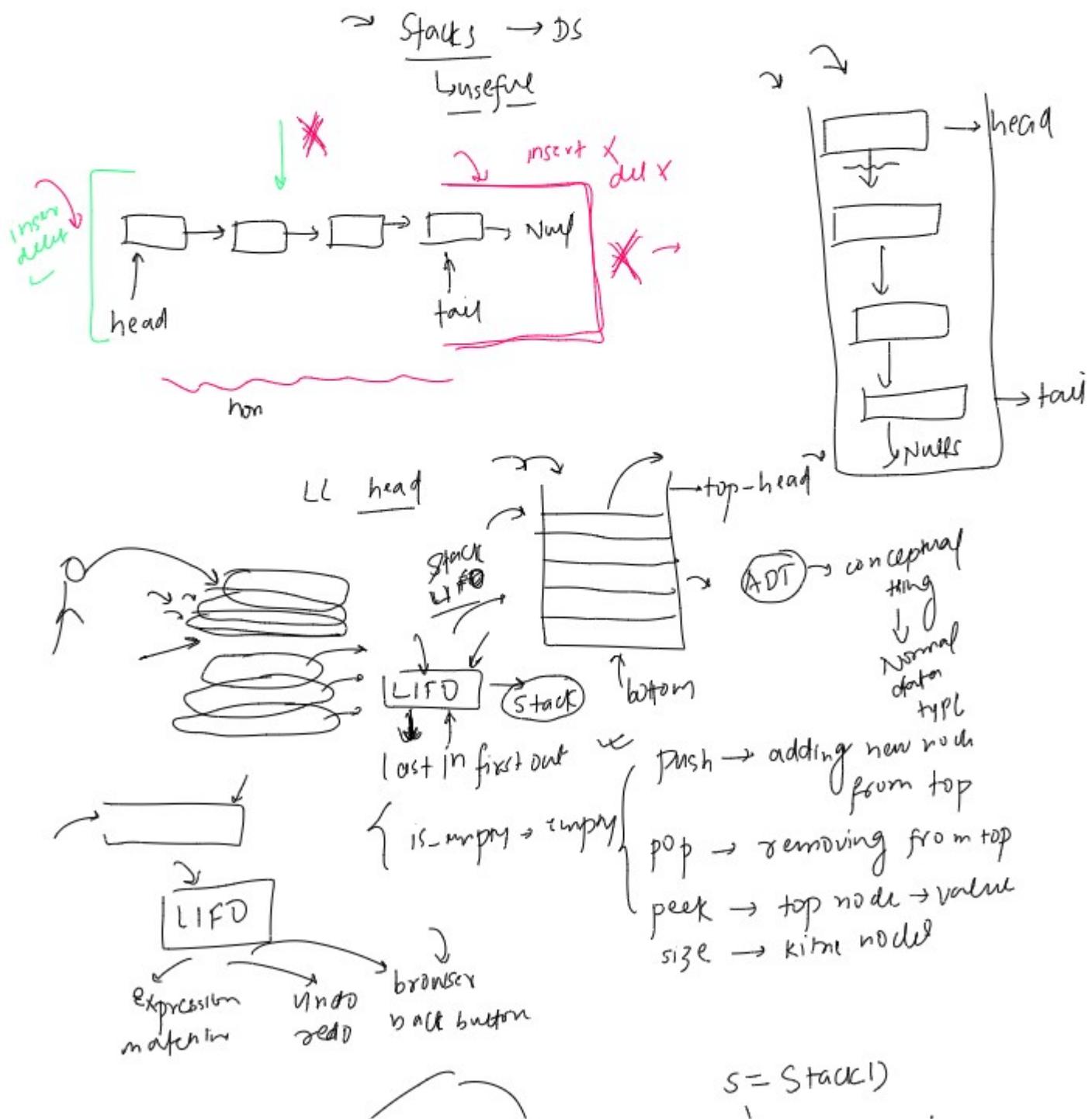
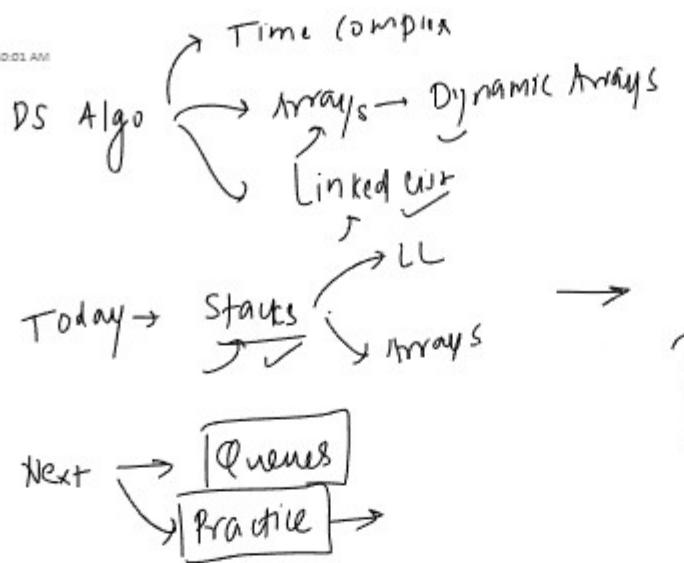
5) Residual Network ↗ CNN → ResNET
↳ building block → ANN

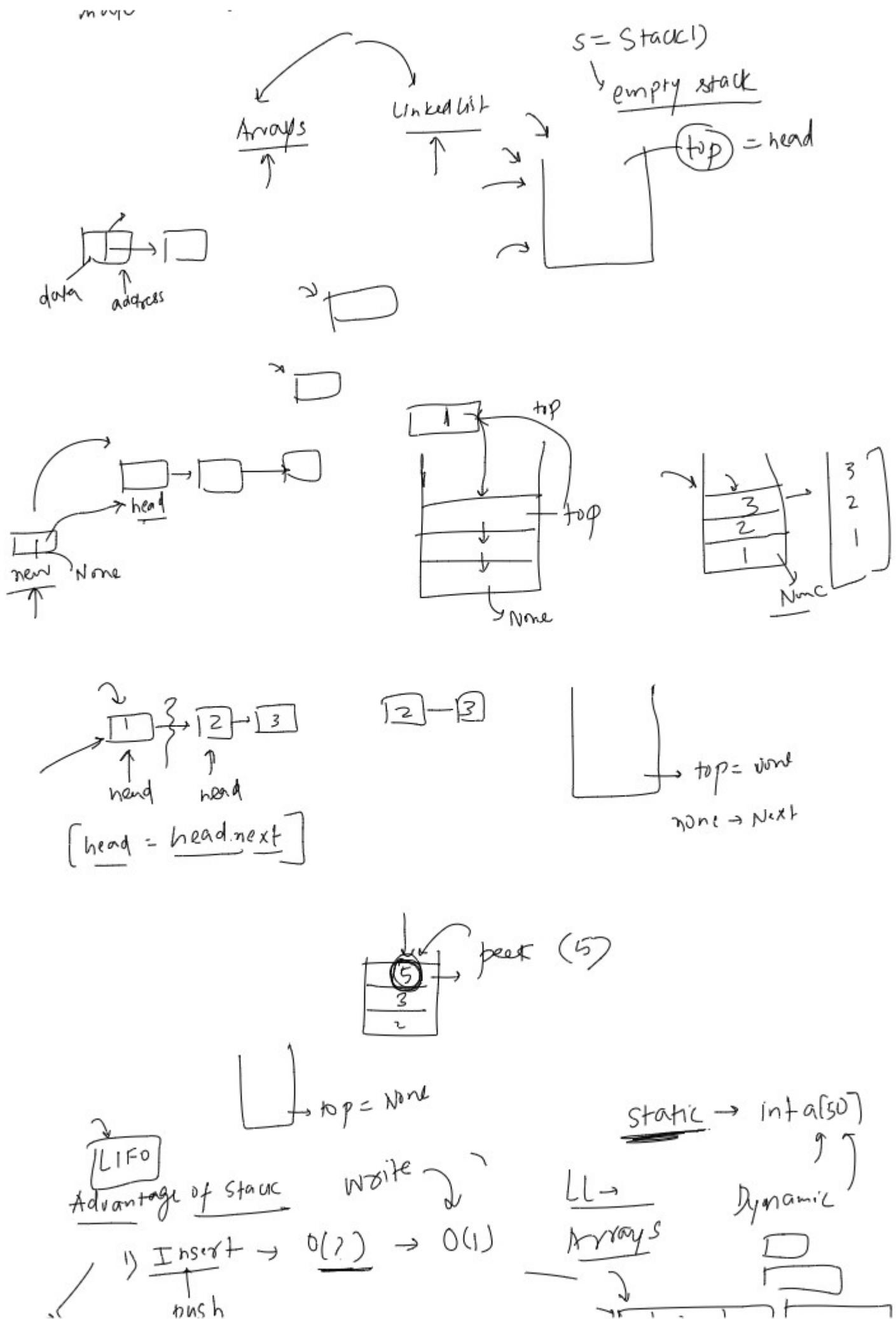
Exploding Gradient Problem → RNN
10, 10, 10, 10 → 10000
random Loss ↓ × ↗ Gradient clipping ↗ 1 - 100 ↗ 1000 ↗ 1000 ↗ big number
 $\frac{\partial L}{\partial w_{ii}} = d_1 \times d_2 \times d_3 \times \dots$
 $w_n = w_0 - \eta \frac{\partial L}{\partial w}$

Vanishing Gradients

exp⁰





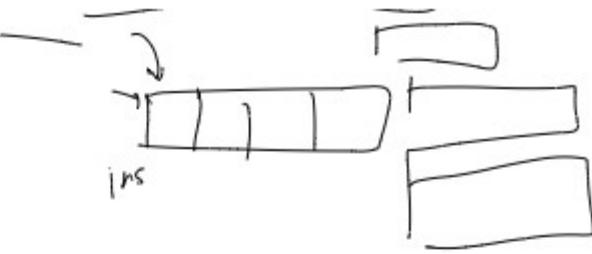


\rightarrow \downarrow \leftarrow

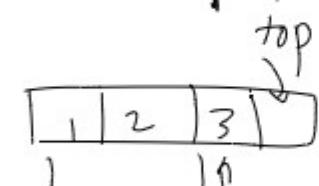
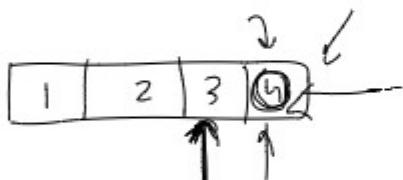
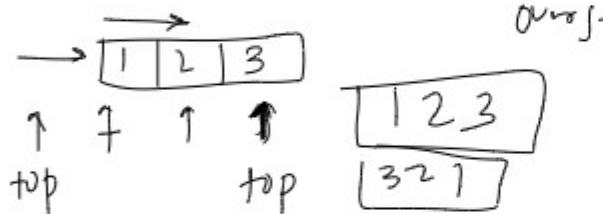
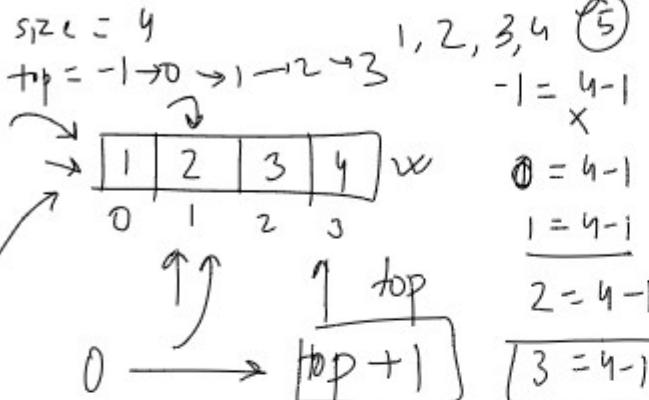
push

$\hookrightarrow \underline{\text{delete}} \rightarrow O(?) \rightarrow O(1)$

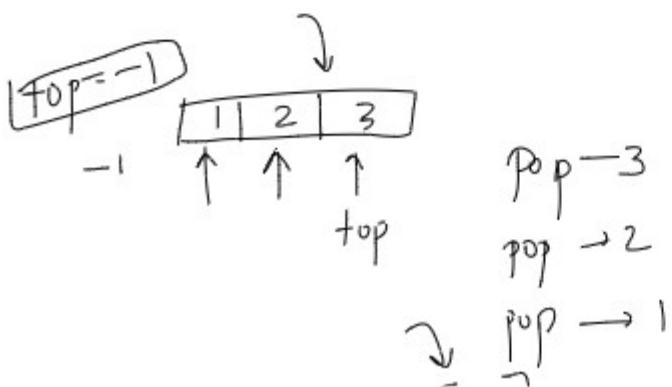
int $a[5]$

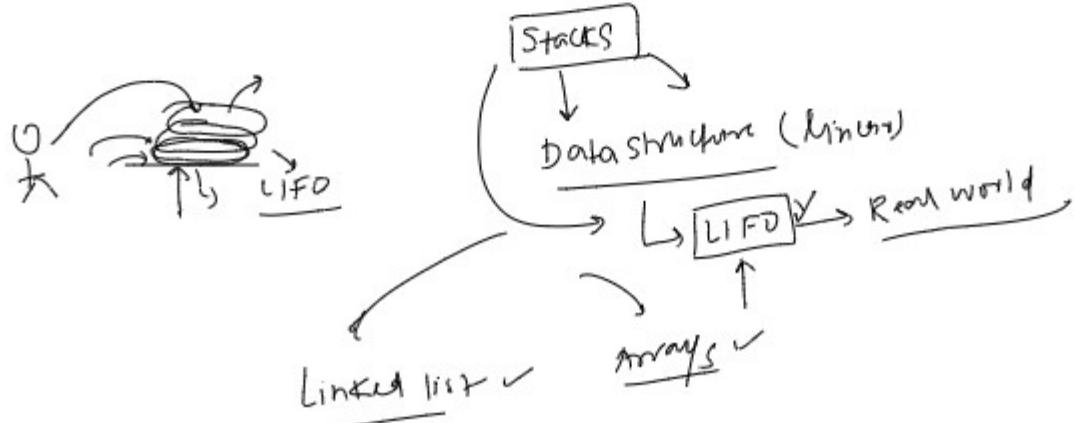
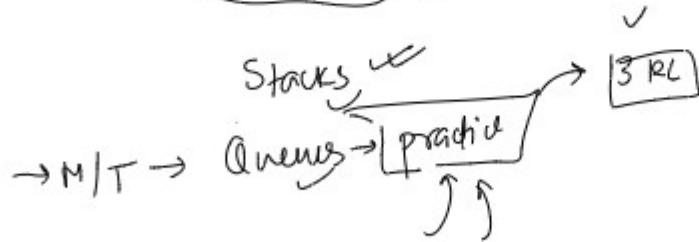
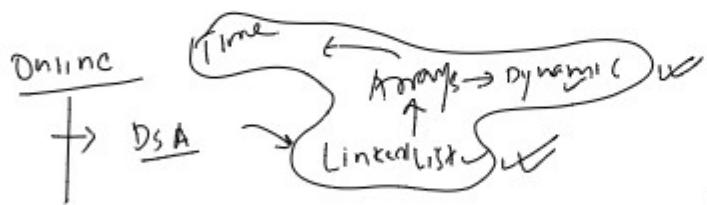
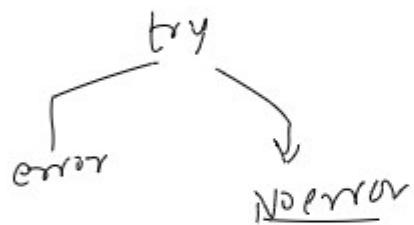
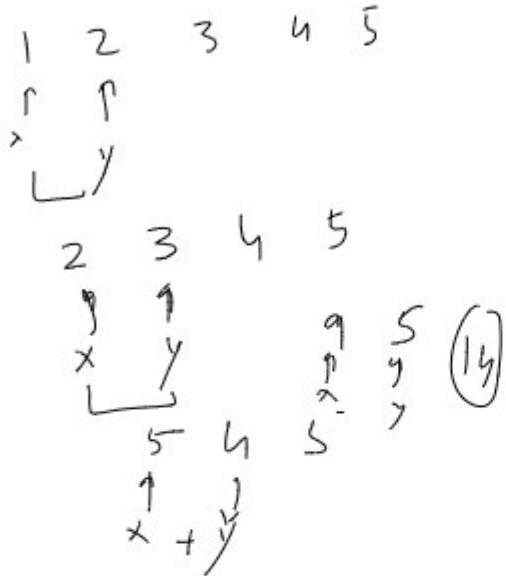
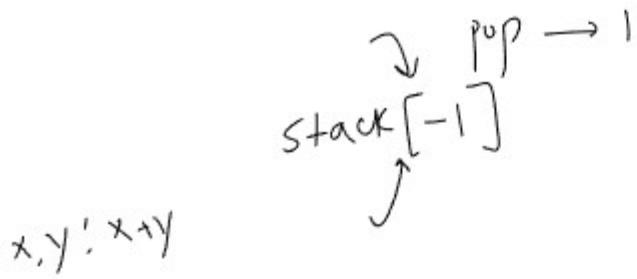


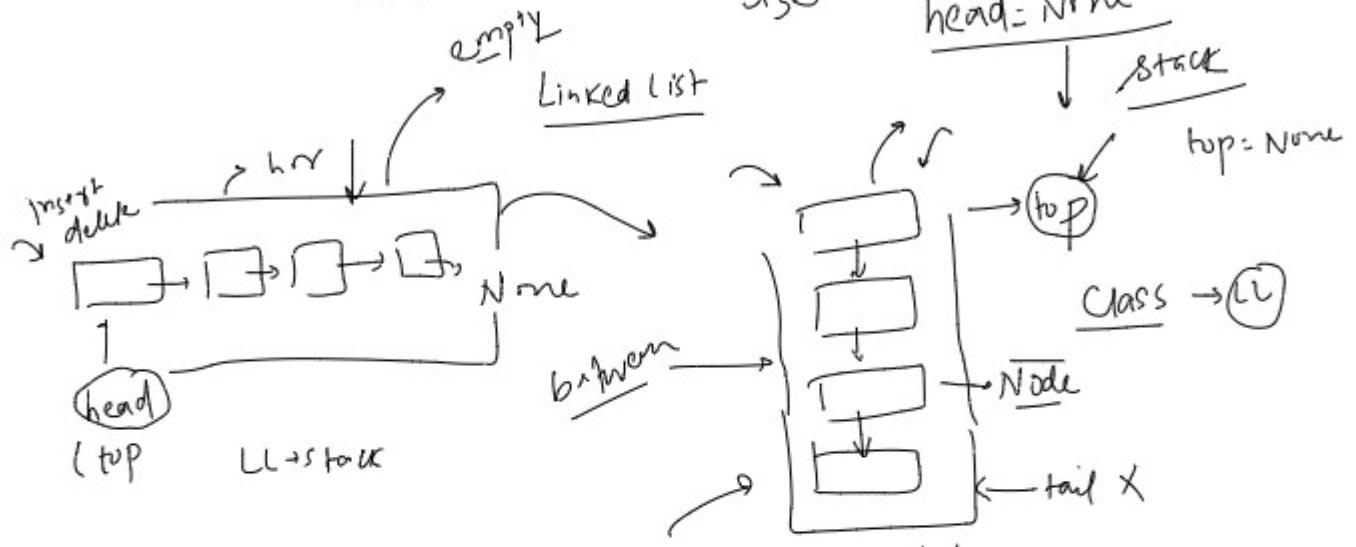
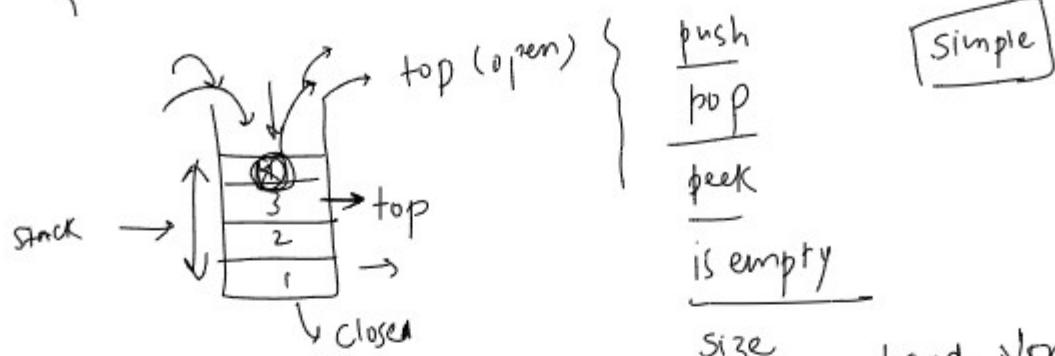
$\text{size} = 4$



\uparrow
 top





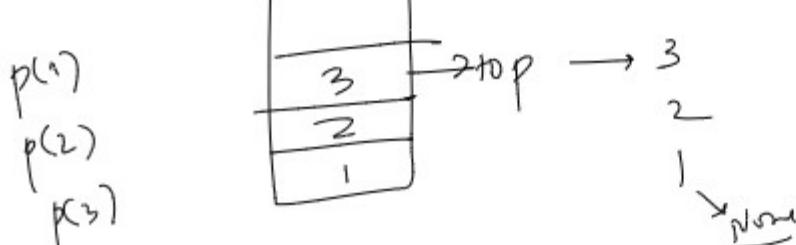
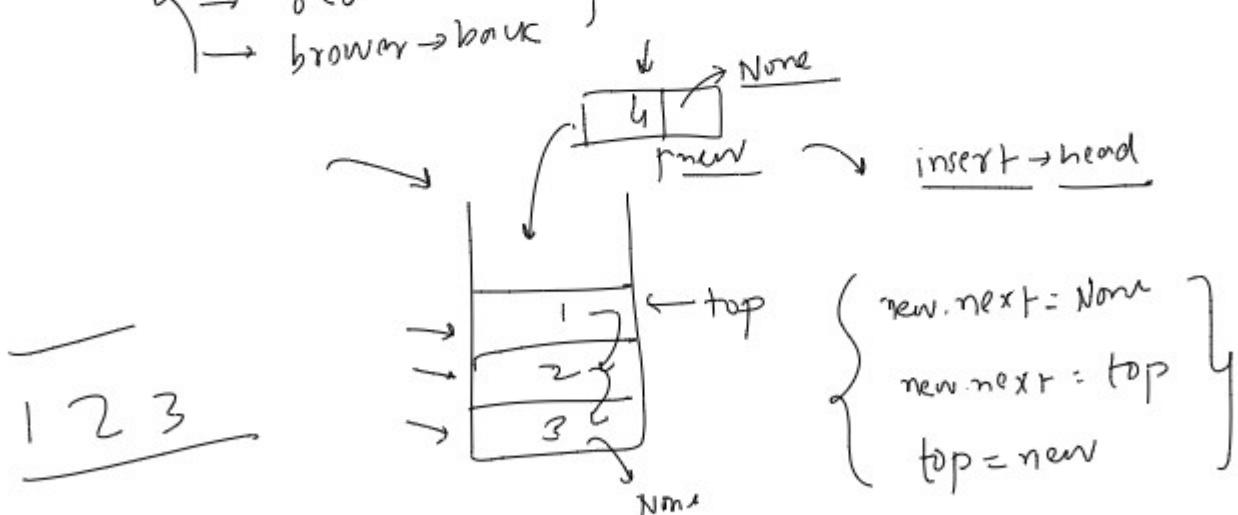


Application

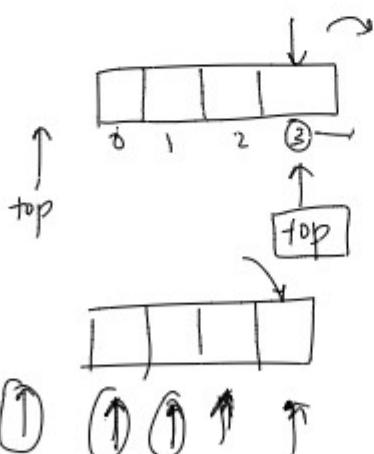
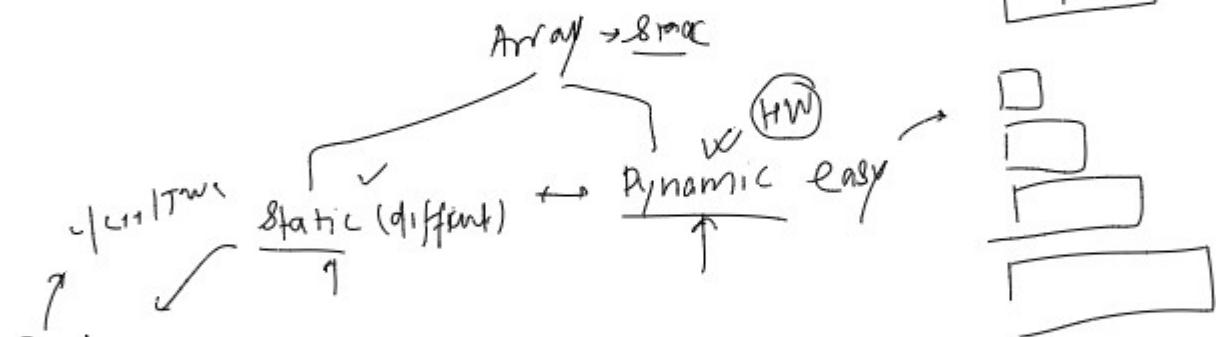
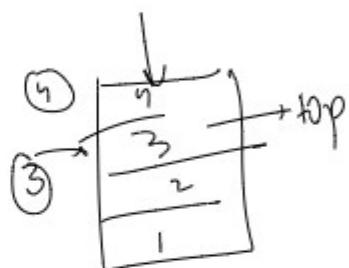
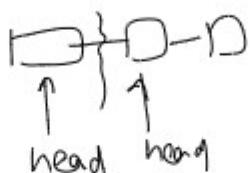
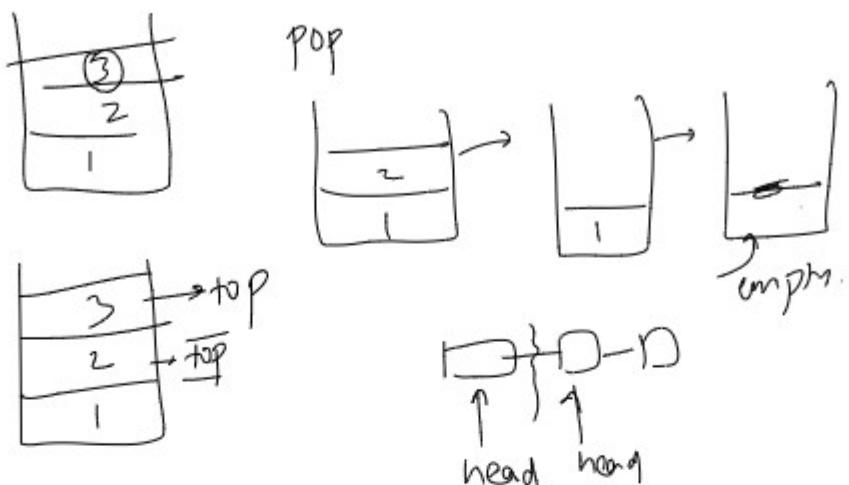
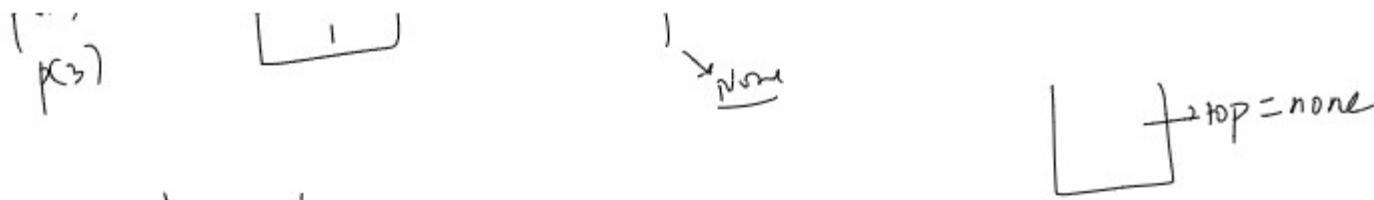
- expression
- undo/redo
- returning
- browser → back

LL

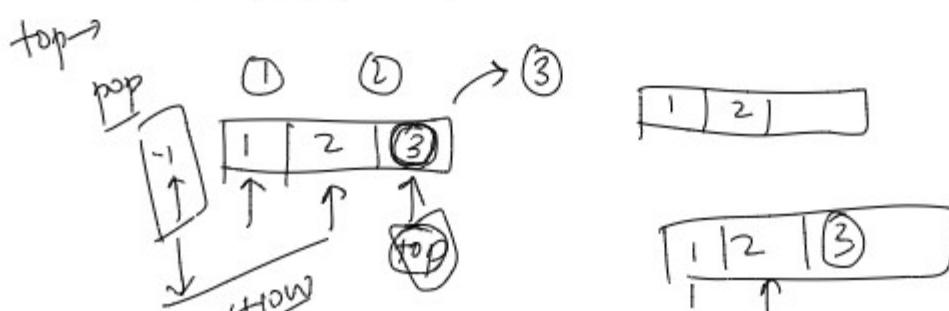
LIFO

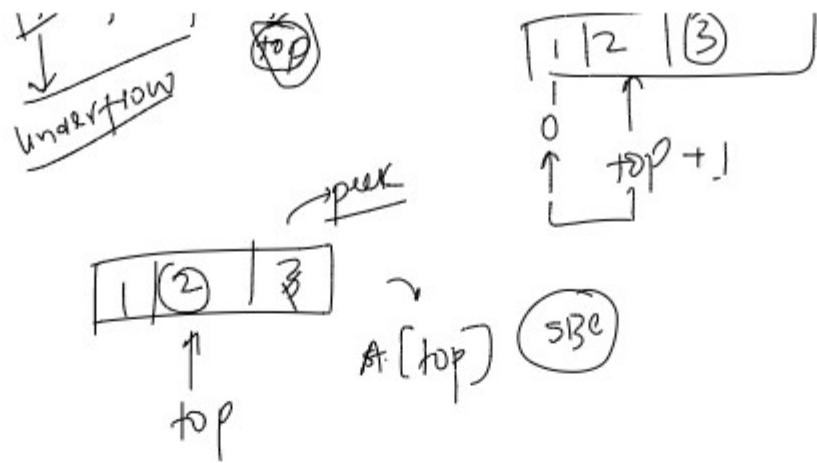


1 ... n - nodd



check
 $\boxed{\text{top} == \text{size}-1}$





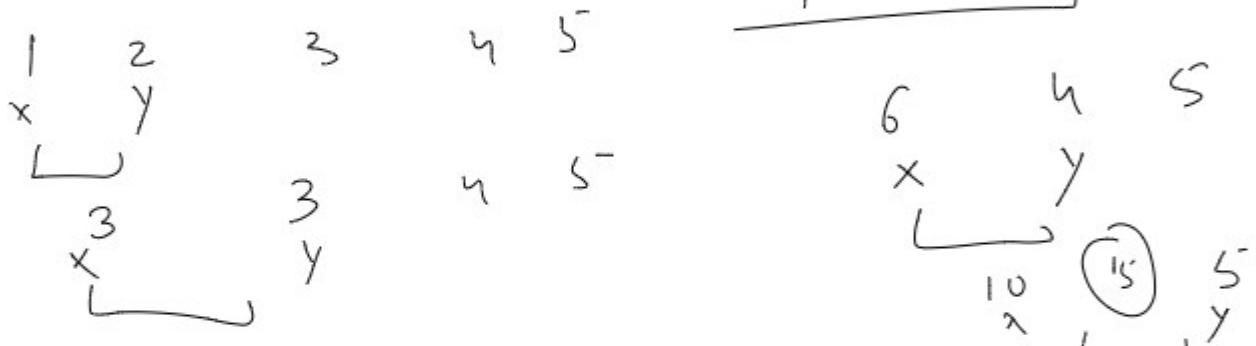
Benefit - LIFO

write

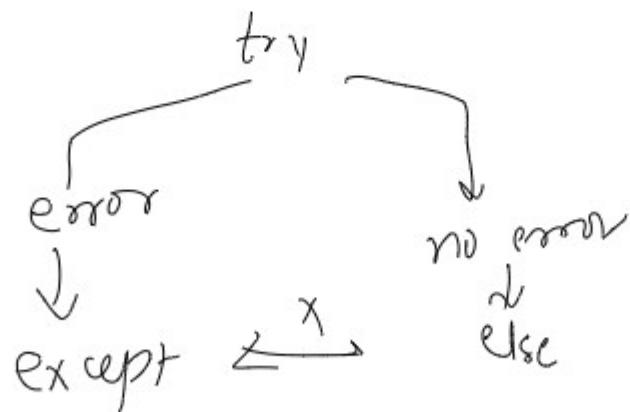
delete
insert

$O(1)$

stack SBC



try
 error
 except
 no error
 else



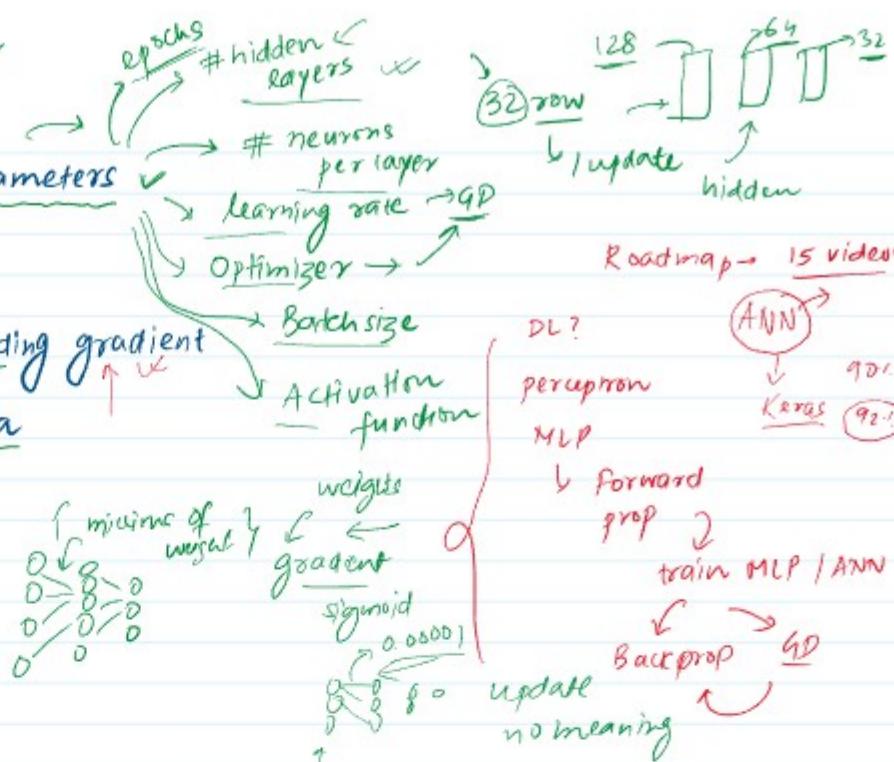
How to improve a neural network ✓

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

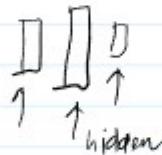
Learning rate

Optimizer

Activation function

Batch size

1. No. of hidden layers



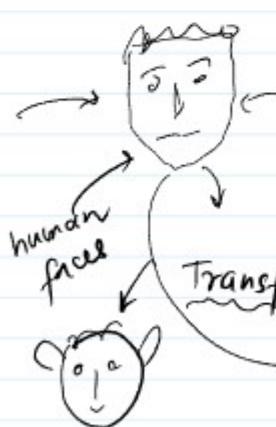
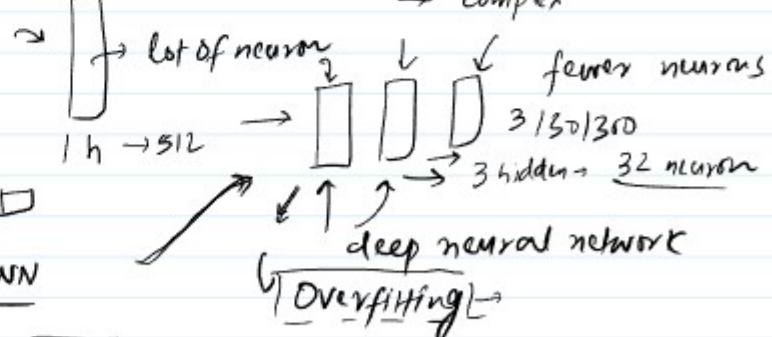
1 hidden layer

↳ 512 neurons

↳ complex

Representation

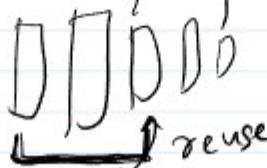
1 ↗ 0 ↗



Transfer learning → CNN

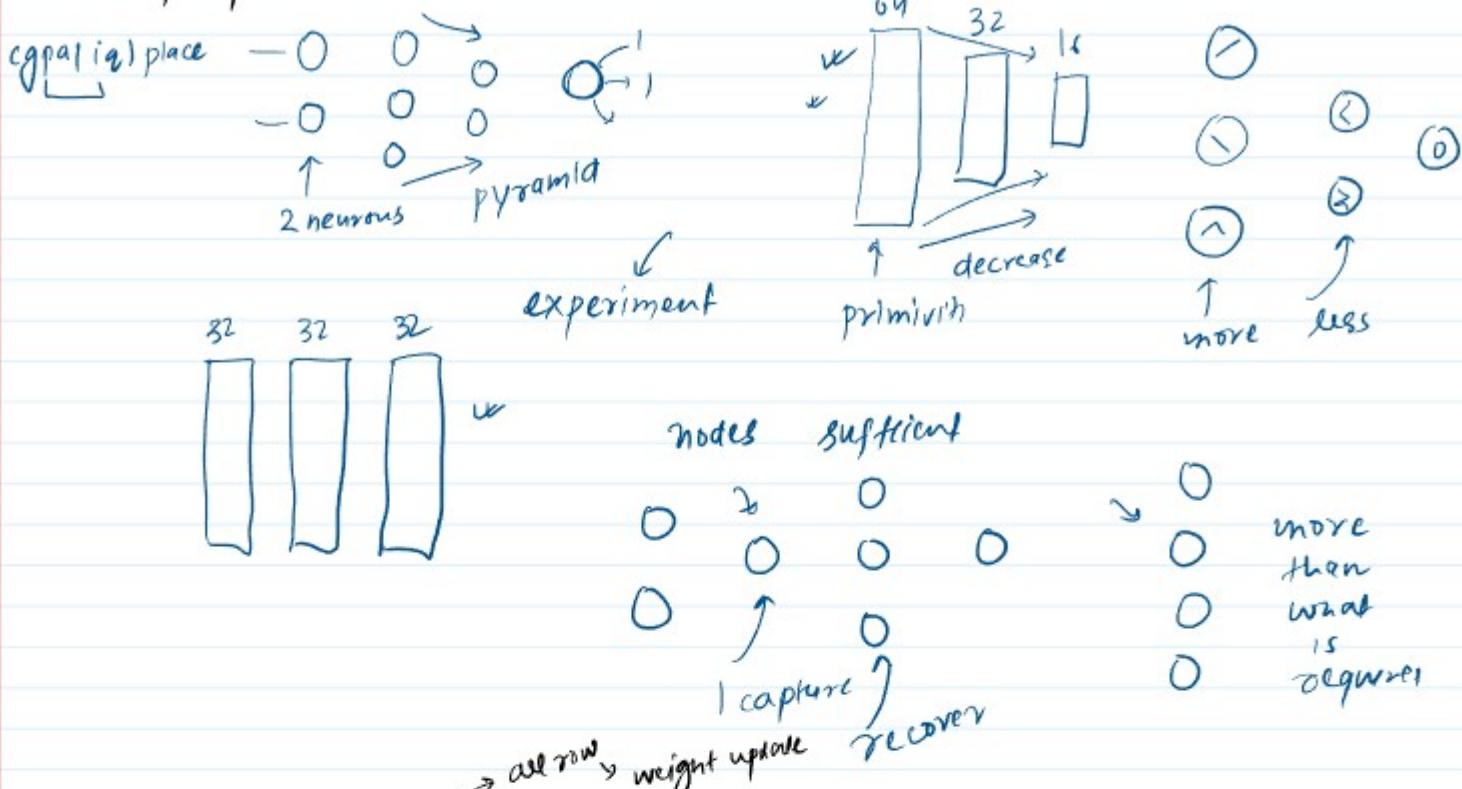
NN

reuse



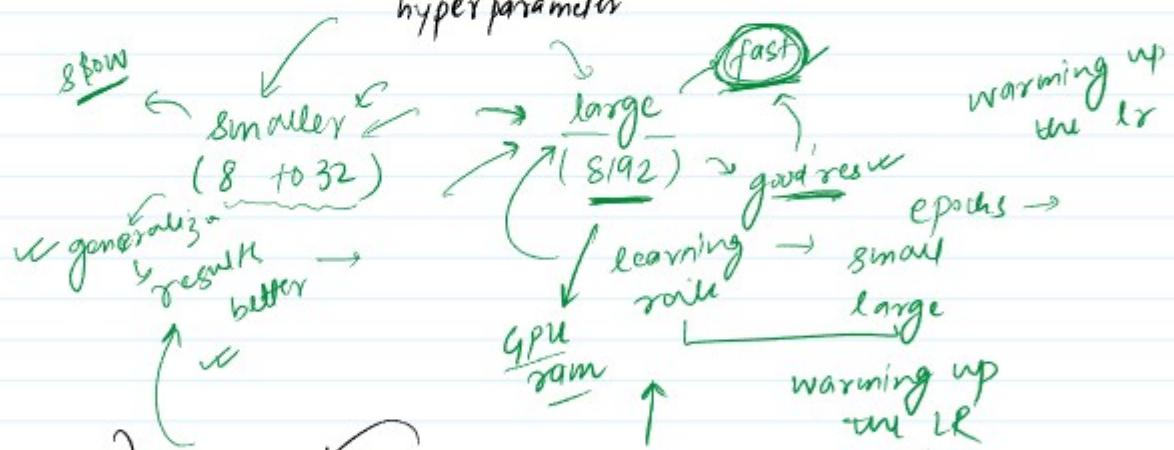
fine, shapes

2 Neuron / layer



3. Batch size

Batch → Stochastic → 1 row → weight
 MiniBatch → batch → 32 → row → update
 hyperparameter



4. Epochs: 100, 500, 1000 → early stopping → Keras → stable

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

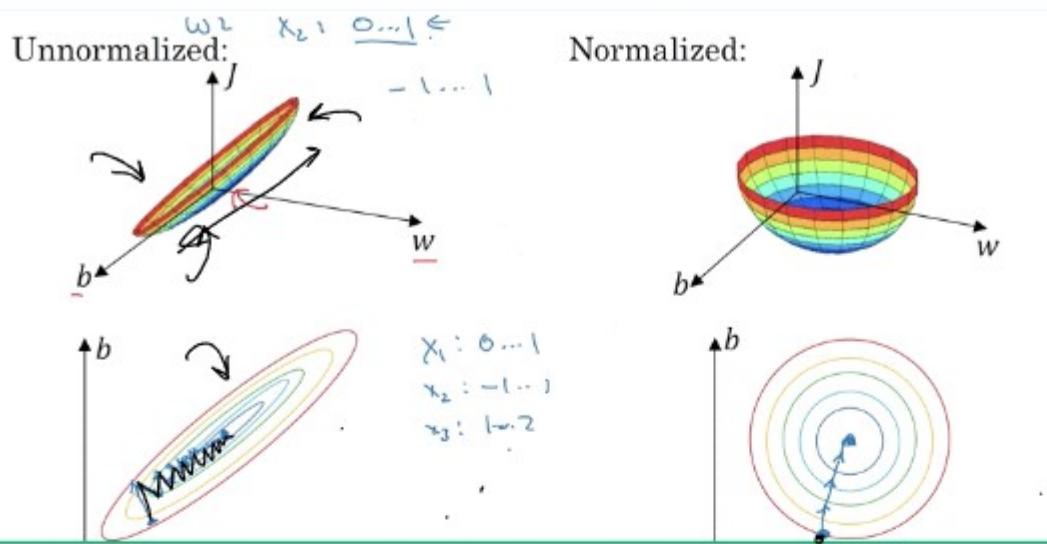
- Optimizers
- Learning rate scheduler

Overfitting

- l_1 and l_2 reg
- Dropouts

Feature Scaling

05 May 2022 17:37



(age)

$$x_1 \curvearrowright x_2 \text{ (Salary)}$$

$$w_1 \quad w_2$$

$$\nabla \rightarrow \quad \nabla \rightarrow \quad w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

Standard

$$\frac{x_i - \mu}{\sigma}$$



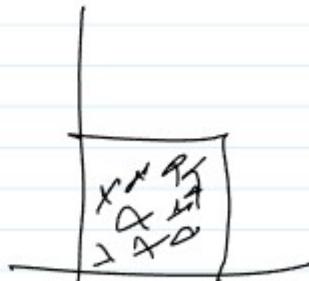
Normalize

$$\frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

$$\min = 0$$

$$\max = 1$$

$$0 \text{ to } 1$$



$$[-1 \text{ to } 1]$$

$$\begin{bmatrix} 2 \times 2 & 2 \times 3 \\ \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix}$$

$$\begin{array}{c} 2 \times 2 \quad 2 \times 3 \\ \boxed{\begin{array}{cc} 1 & 2 \\ \cancel{3} & \cancel{4} \end{array}} \quad \boxed{\begin{array}{c} 1 \\ \cancel{4} \\ - \\ 2 \\ 5 \\ 3 \\ \cancel{6} \end{array}} \end{array}$$

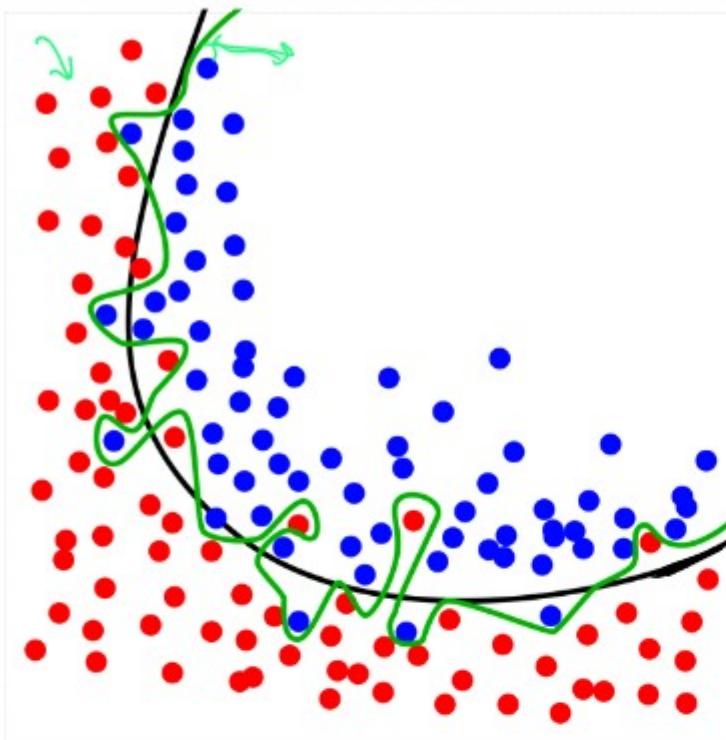
$$\begin{bmatrix} 1+8 & 2+10 & 3+12 \\ 3+16 & 6+20 & 9+24 \end{bmatrix}$$

$$\begin{bmatrix} 9 & 12 & 15 \\ 19 & 26 & 33 \end{bmatrix}$$

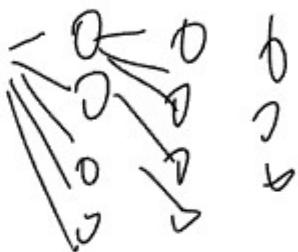
↑

The Problem of Overfitting

12 May 2022 09:32



ANN
overfitting
Complex
→
reduces



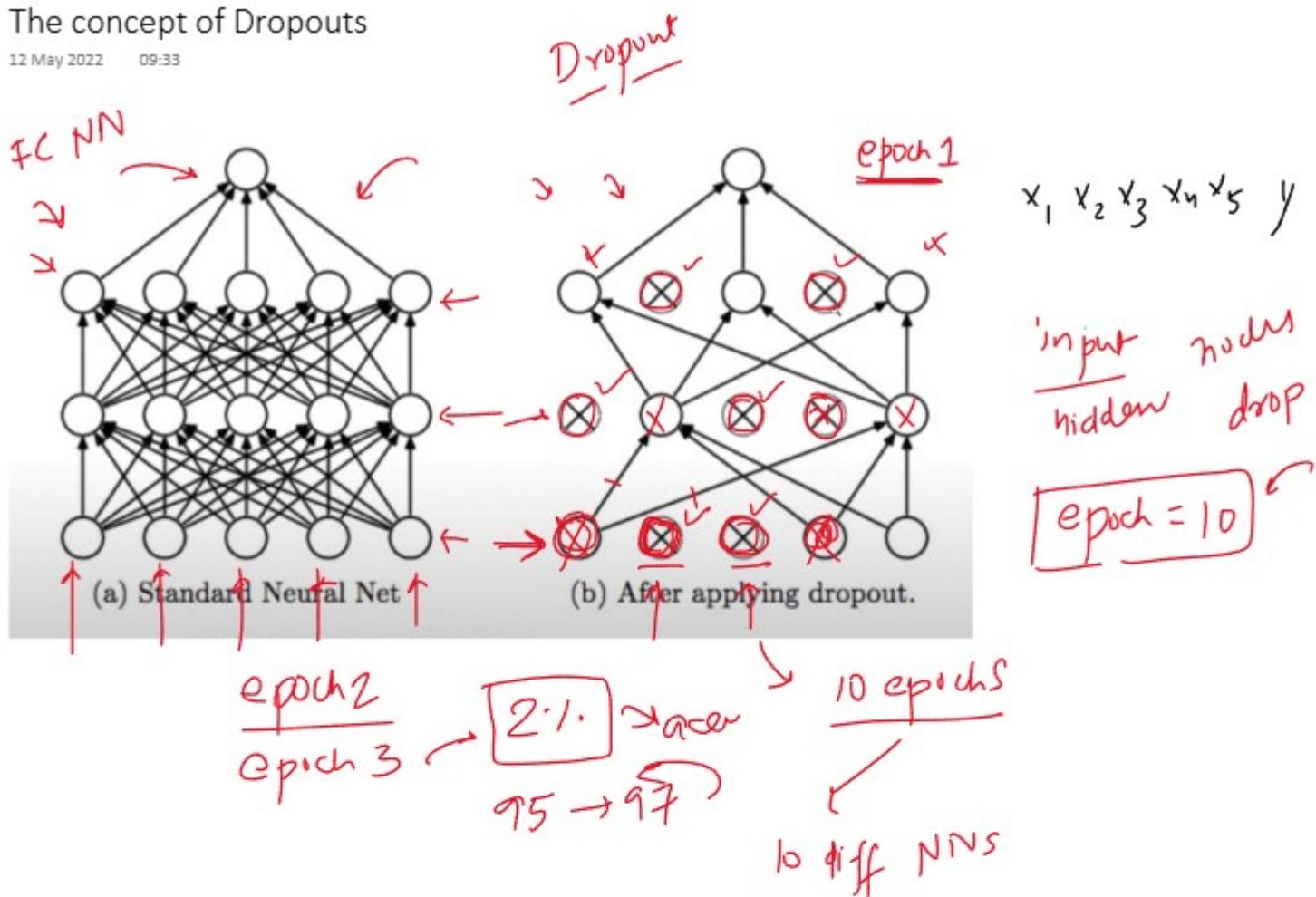
Possible Solutions

12 May 2022 09:33

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

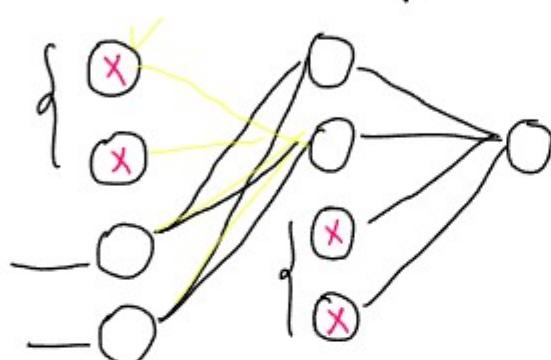
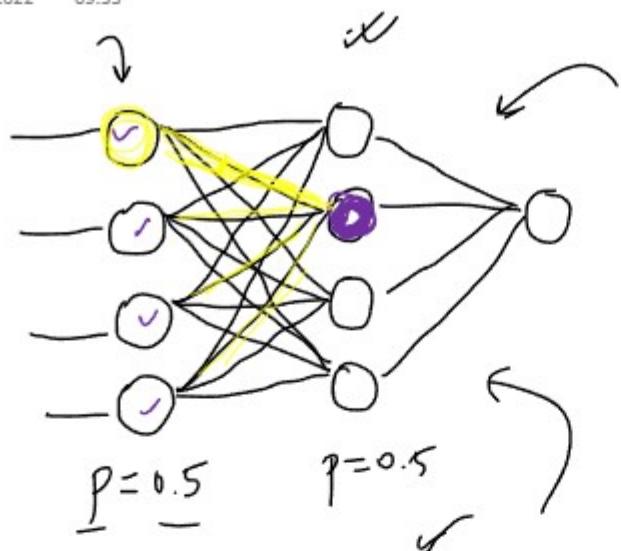
The concept of Dropouts

12 May 2022 09:33

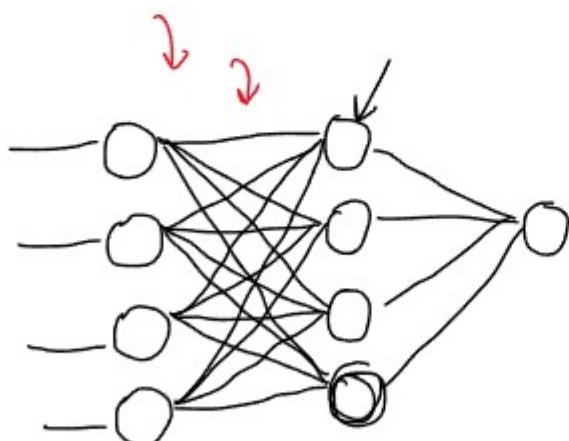
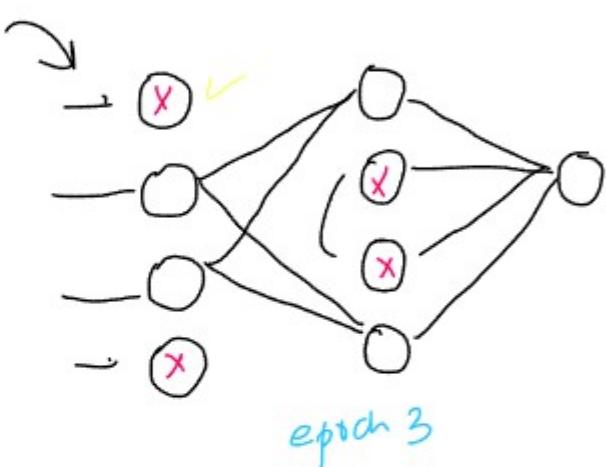
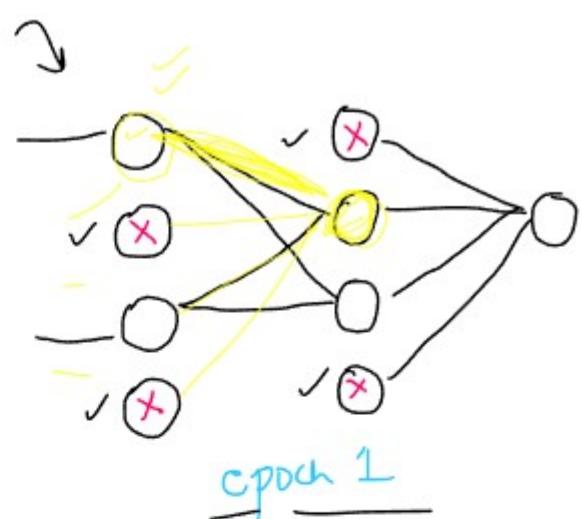


Why this works?

12 May 2022 09:33

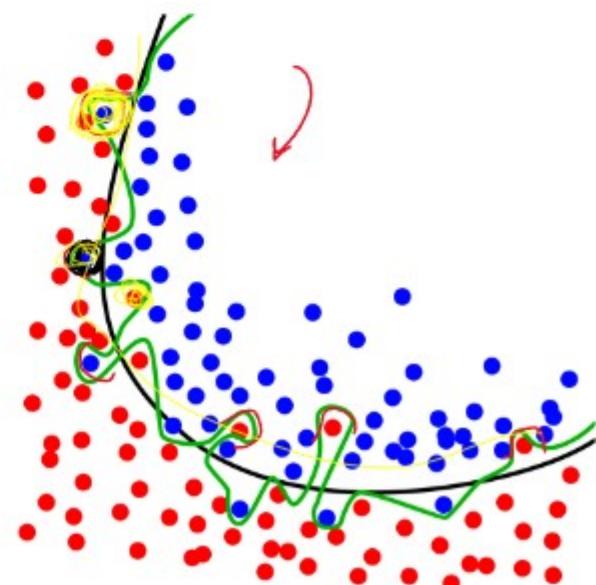


epoch 2



→ Reduce the # nodes

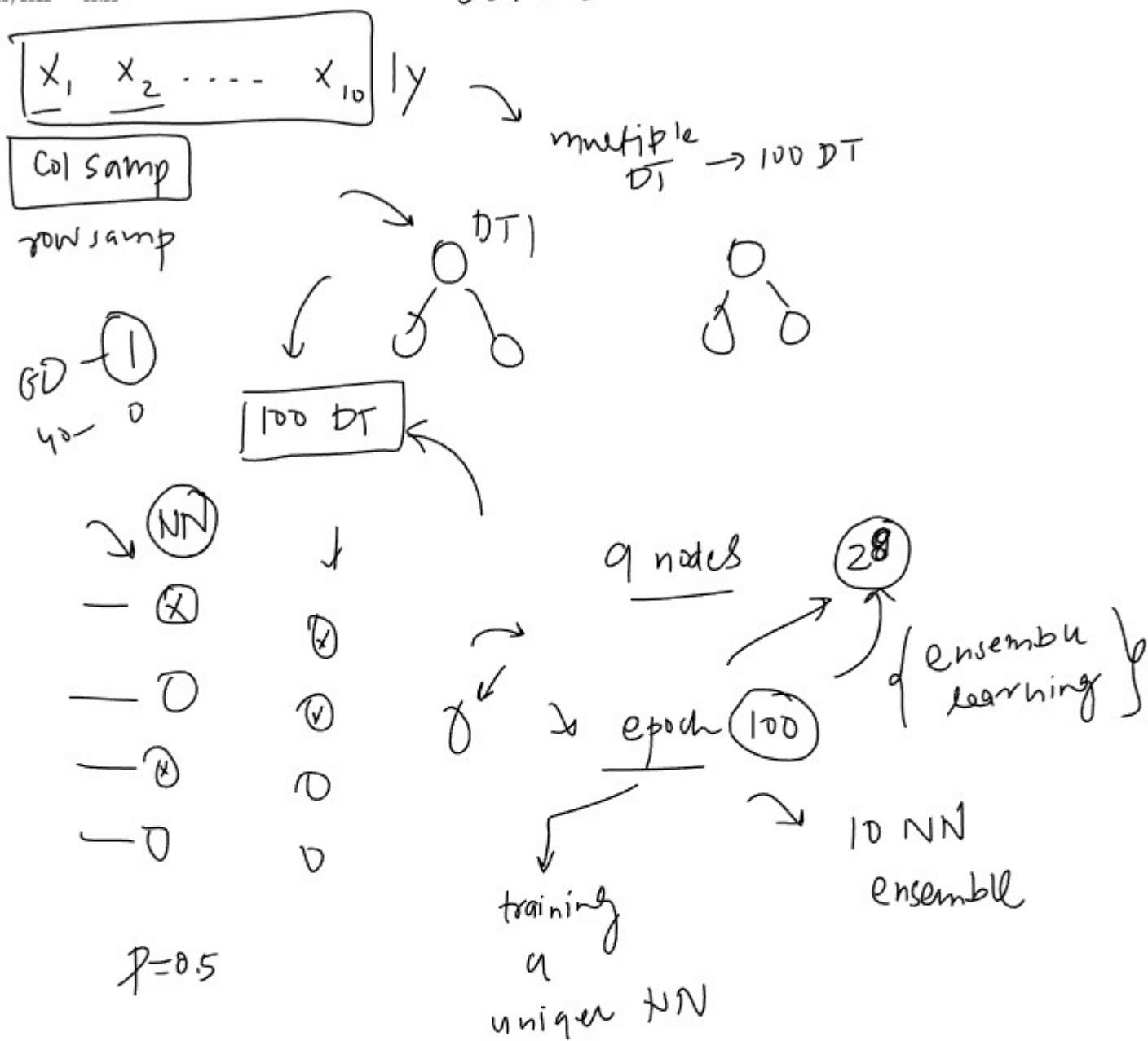
→



Random Forest Analogy

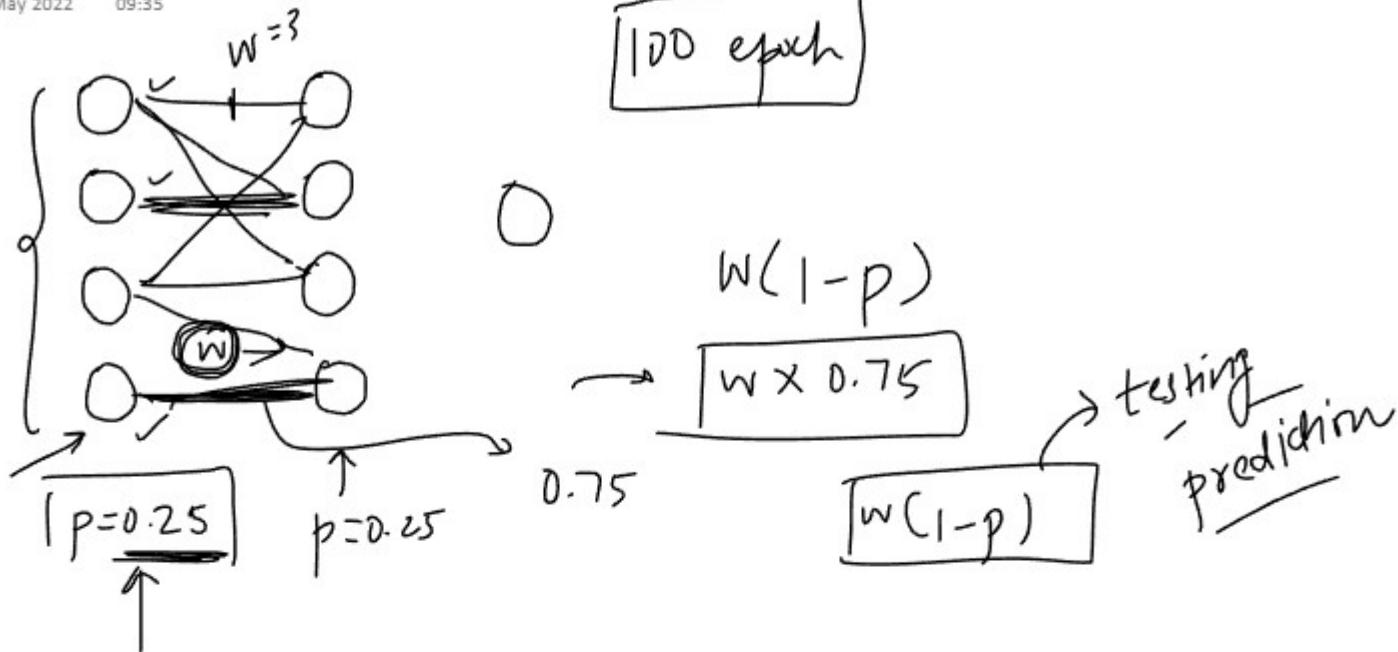
12 May 2022 09:33

50% w/o



How prediction works?

12 May 2022 09:35



Regression Code Example

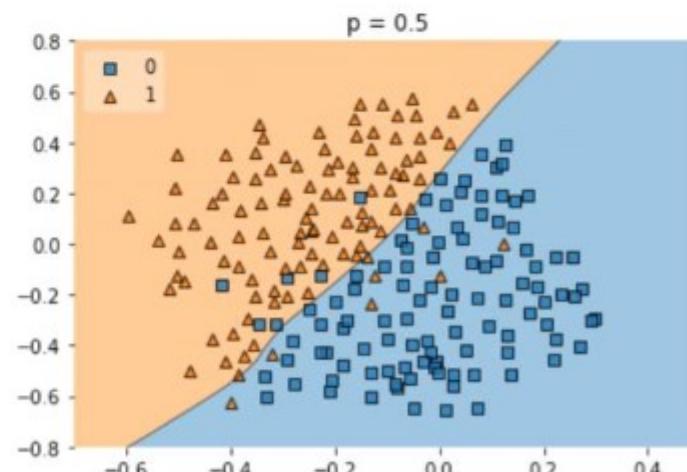
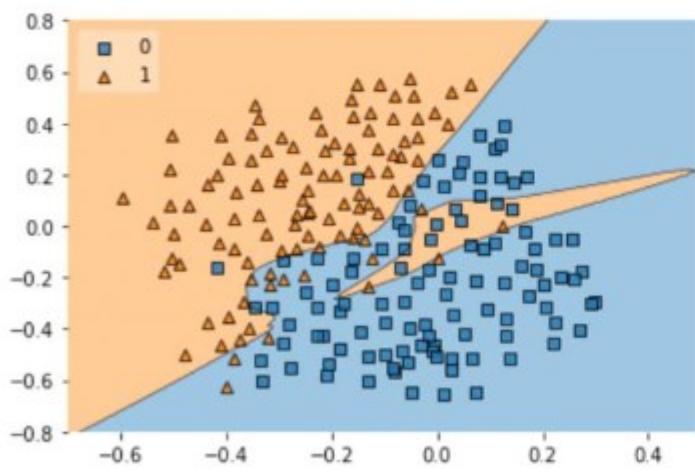
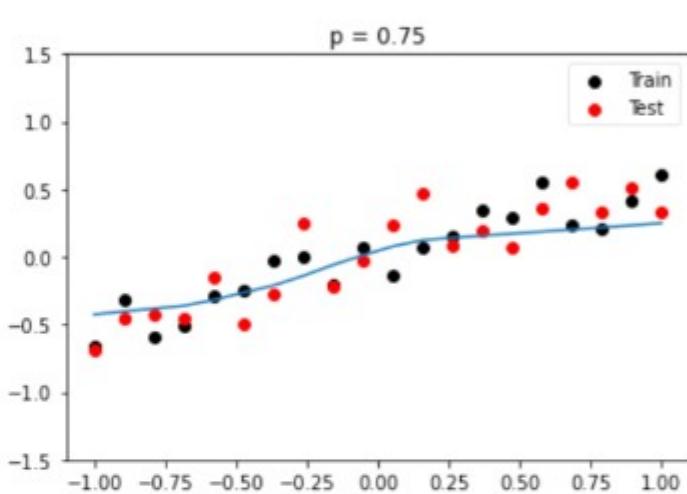
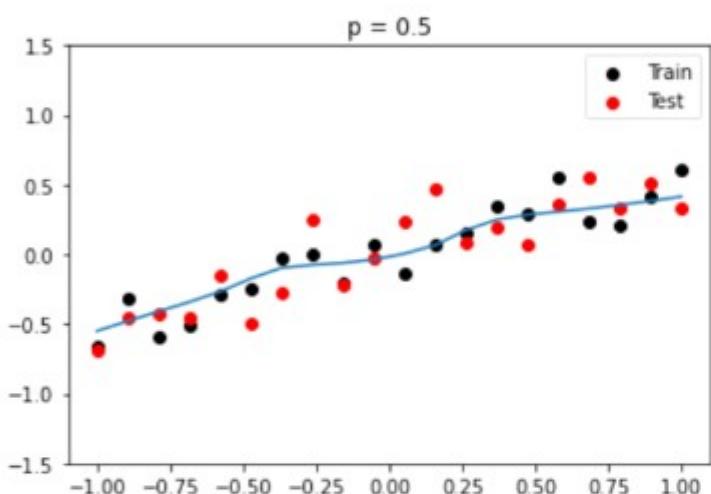
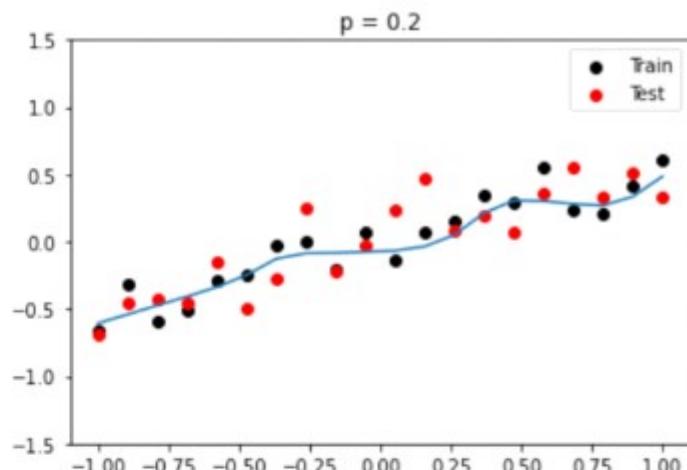
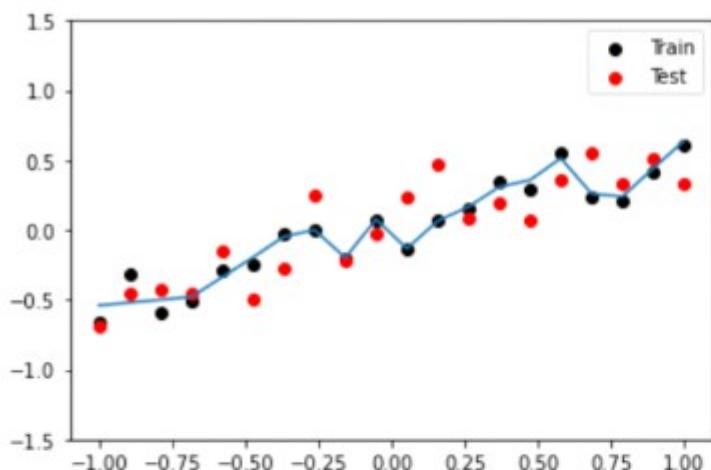
12 May 2022 09:34

Classification Code Example

12 May 2022 09:34

Effect of p

12 May 2022 10:06

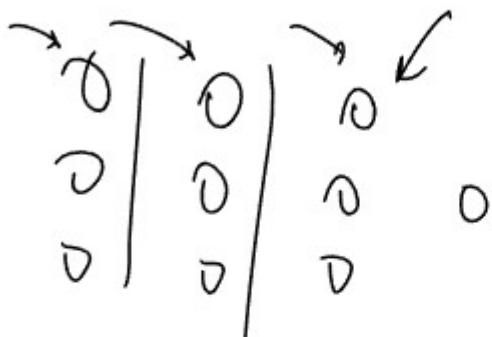


Practical Tips and Tricks

12 May 2022 09:34

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

2) Last layer \rightarrow dropout



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

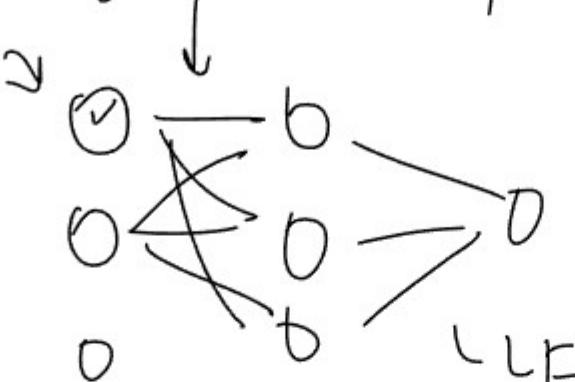
20 - 30 RNN

50 >

ANN \rightarrow 10 - 50

Drawbacks

12 May 2022 09:34

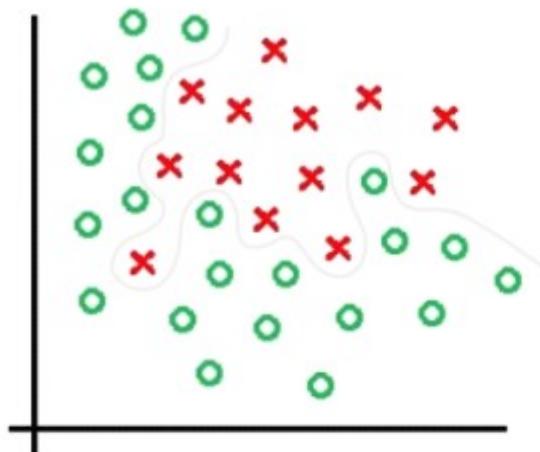
- 1) convergence \rightarrow delay
- 2) 
loss function
value changes
- LF - has to change

Resources

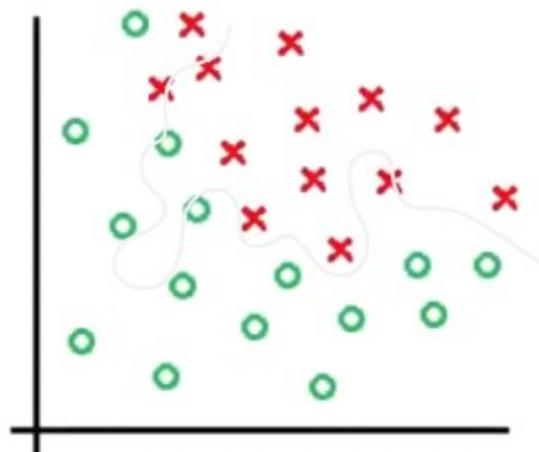
12 May 2022 12:38

Overfitting

20 May 2022 08:14



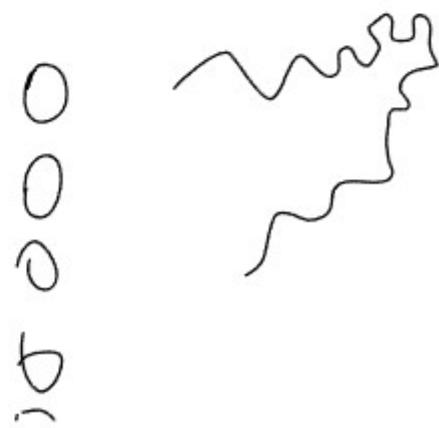
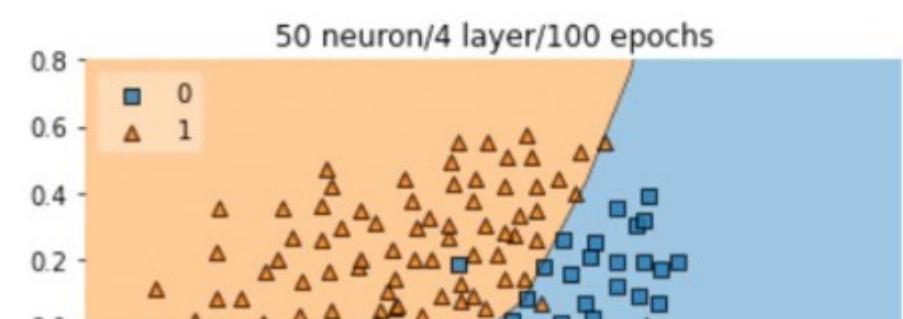
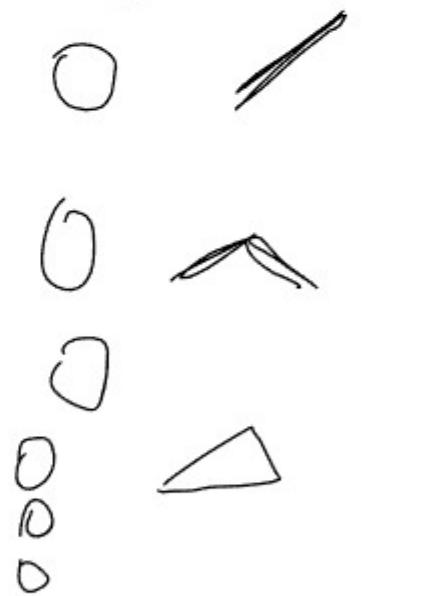
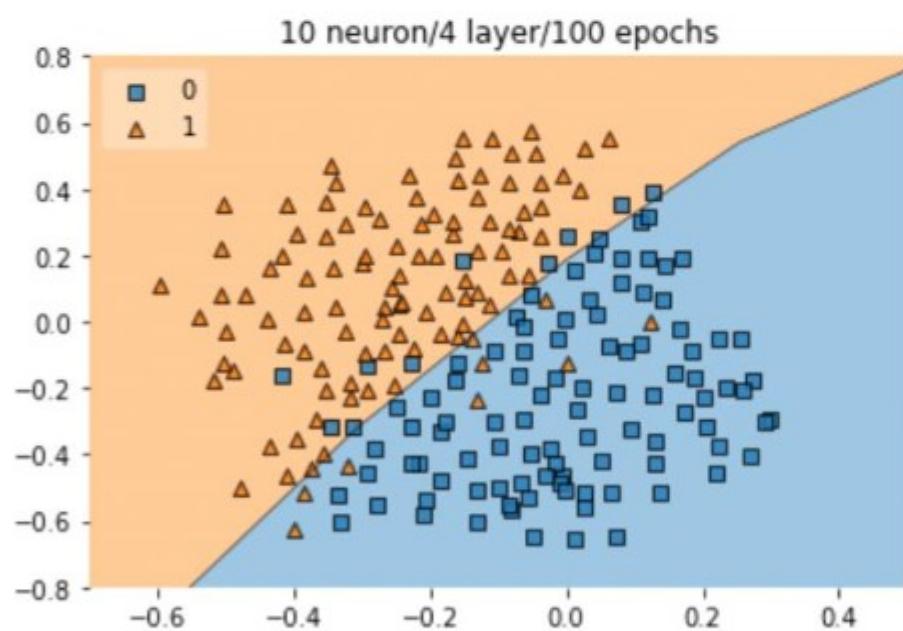
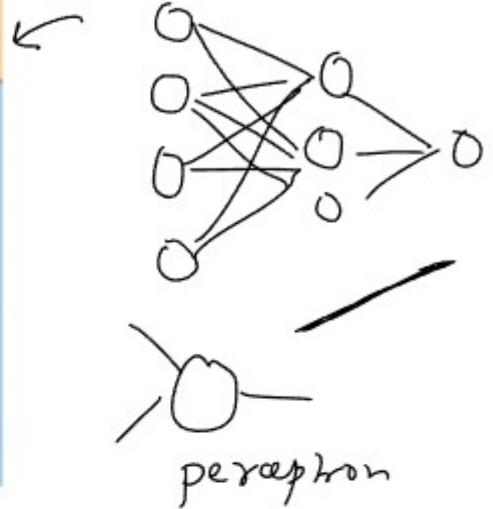
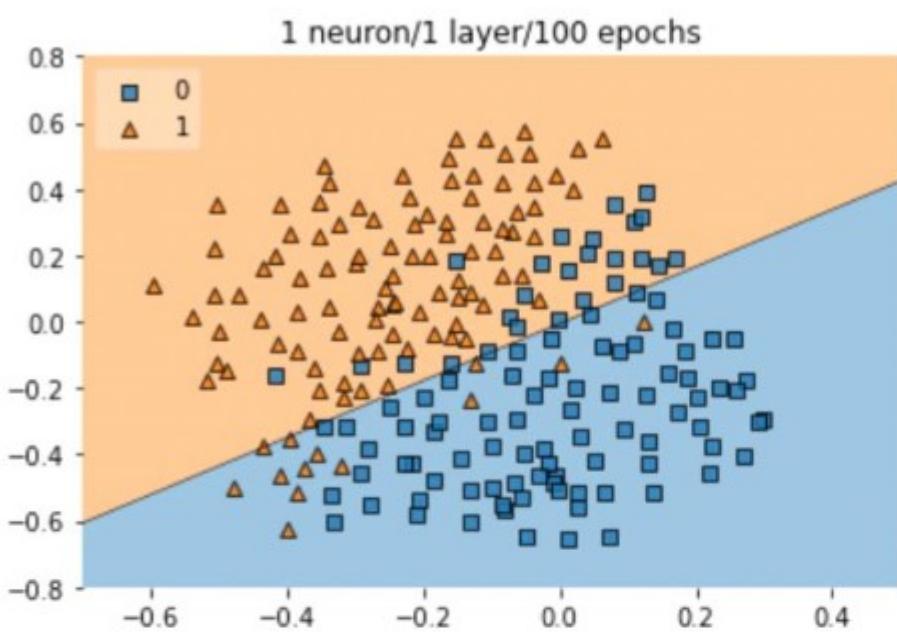
Train Dataset

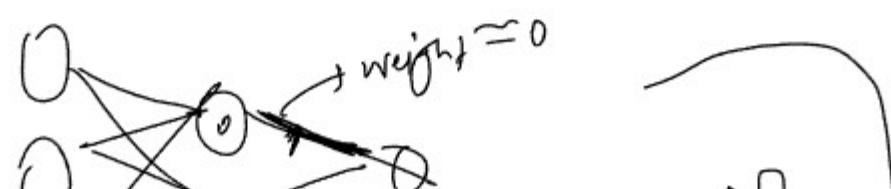
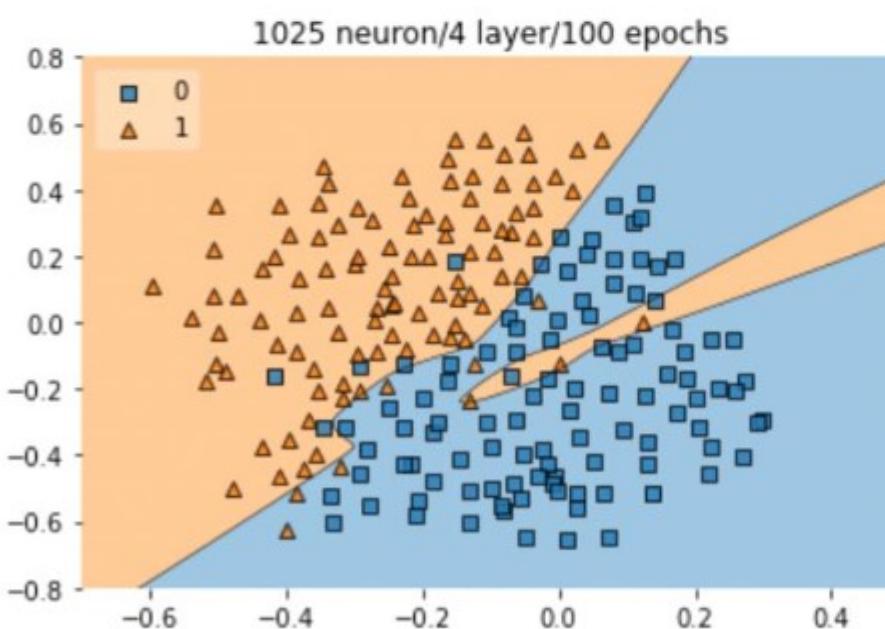
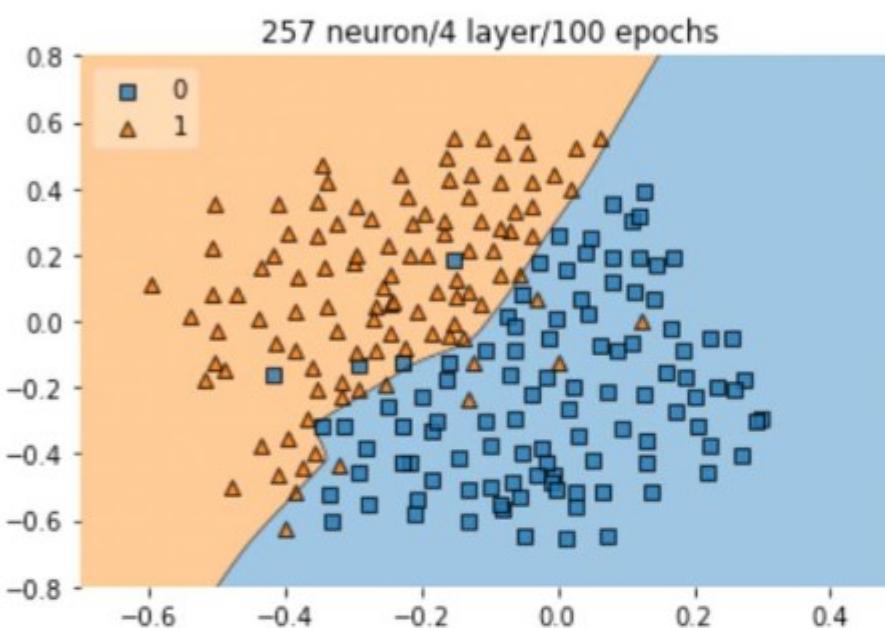
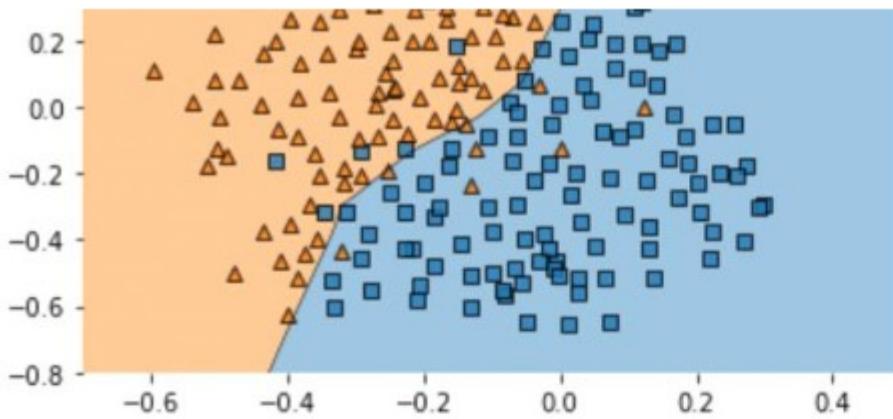


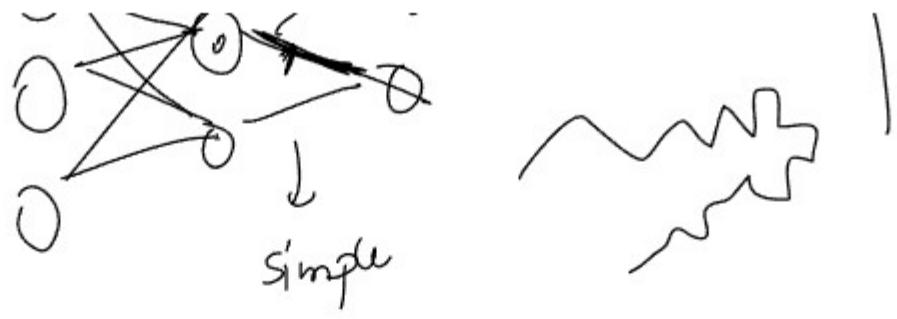
Test Dataset

Why Neural Networks Overfit?

19 May 2022 07:34

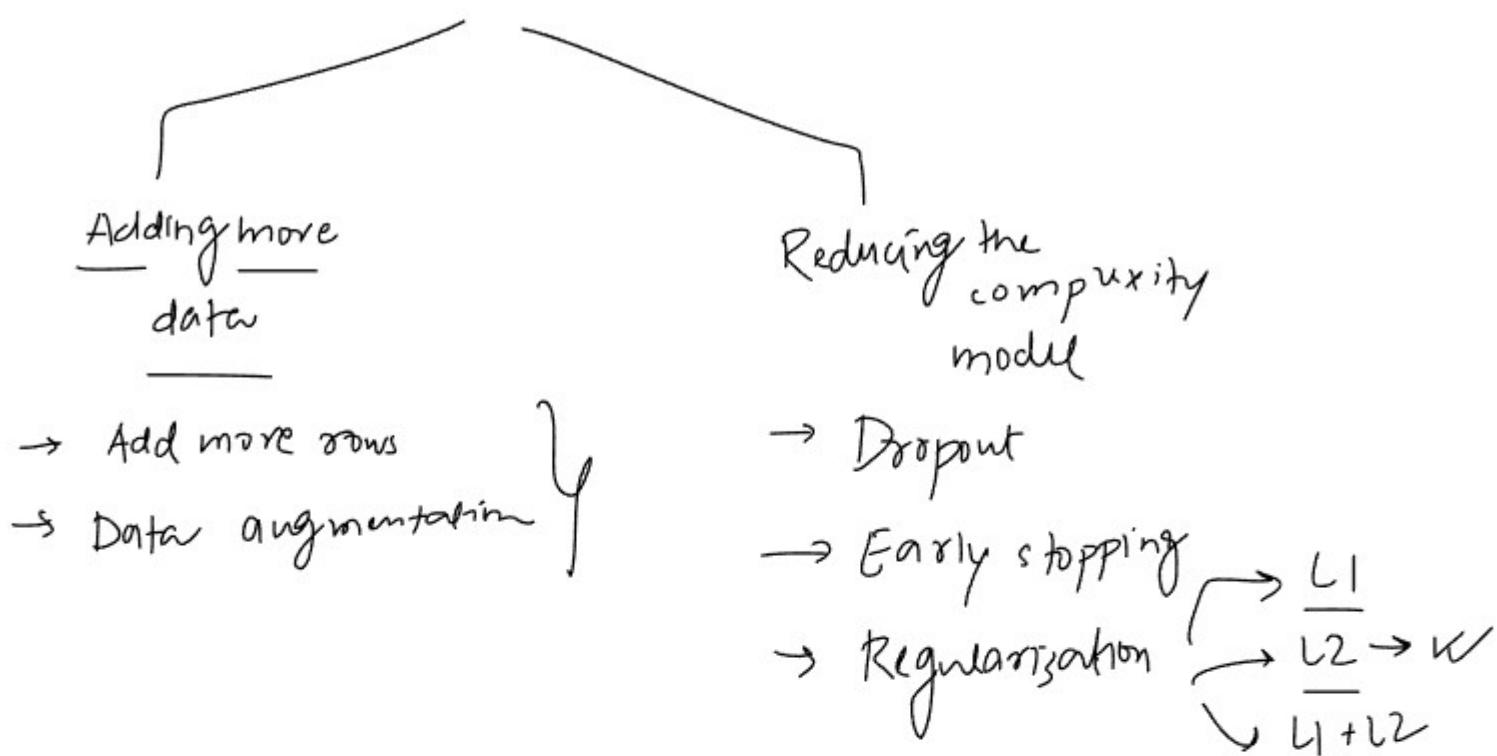






Ways to solve overfitting

20 May 2022 08:13

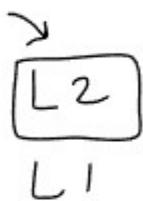


$A_{NN} \rightarrow \text{weights} | \text{bias}$

$\hookrightarrow \min \text{ Lossfunction}$

$$L = \text{mse}$$

↓
binary



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

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

$$w_1 \rightarrow w_{10}$$

$$\left(\frac{\lambda}{2n} \right) [w_1^2 + w_2^2 + \dots + w_{10}^2]$$

↑ $\lambda = \text{hyperparameter} \uparrow$

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

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

$$\downarrow$$

$w \approx 0$

$\rightarrow \underline{\text{l1 norm}}$

$$C = L + \frac{\lambda}{2n} \sum \|w_i\|_1$$

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

↗ 0 ↗ 15 ↗ n

$$d, j = d$$

$$\begin{array}{c}
 \text{Diagram of a neural network layer: } \\
 \text{Input layer } L=1: \quad \text{Three nodes} \\
 \text{Hidden layer } l=2: \quad \text{Three nodes} \\
 \text{Output layer } l=3: \quad \text{One node} \\
 \text{Bias: } \boxed{\text{bias } X}
 \end{array}
 \quad
 \begin{aligned}
 \psi, J &= \alpha \\
 w_1^2 + w_2^2 + w_3^2 & \\
 \sum_{i=1}^k \|b_i\|^2 &
 \end{aligned}$$

Intuition behind Regularization

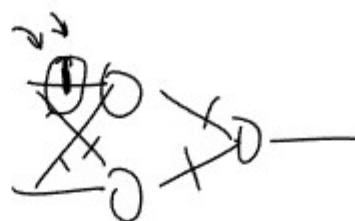
20 May 2022 08:15

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

$1 - \eta \lambda$ positive

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

L



$$w_1^T [w_2^T + w_3^T + \dots]$$

zw

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

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

$$w_0 \ll w_0$$

L2 reg \rightarrow weight decay

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

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

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

weight decay $\uparrow w_0$

Code Demo

20 May 2022 08:15

Further Resources

20 May 2022 08:15

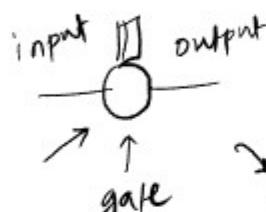
What are Activation Functions?

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

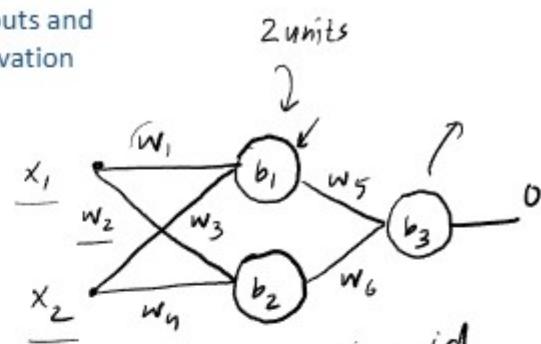
$$a = g(w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b)$$

This function g is referred to as the activation function.



activated
(Y/N)
how much

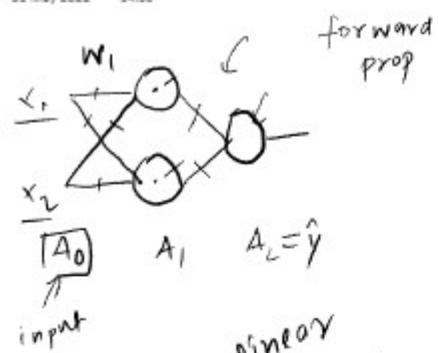
$$\underbrace{g}_{\text{activation}}(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$



2 units
sigmoid
relu
softmax

Why Activation Functions are needed?

31 May 2022 14:50



Non-linear
output = input

$$A_2 = \hat{y} = \underline{\underline{w' A_0 + b'}}$$

Linear

$$z_1 = \underline{\underline{w_1 A_0 + b_1}}$$

$$A_1 = \underline{\underline{g(z_1)}} = z_1$$

~~linear~~ \rightarrow sigmoid

$$A_2 = \underline{\underline{g(w_2 A_1 + b_2)}}$$

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

ANN



Non-linear

$$\begin{aligned} \text{input} &= w_2 A_1 + b_2 = w_2 (w_1 A_0 + b_1) + b_2 \\ &= \underline{\underline{w_2 w_1 A_0}} + \underline{\underline{w_2 b_1 + b_2}} \\ &= \underline{\underline{w' A_0 + b'}} \end{aligned}$$

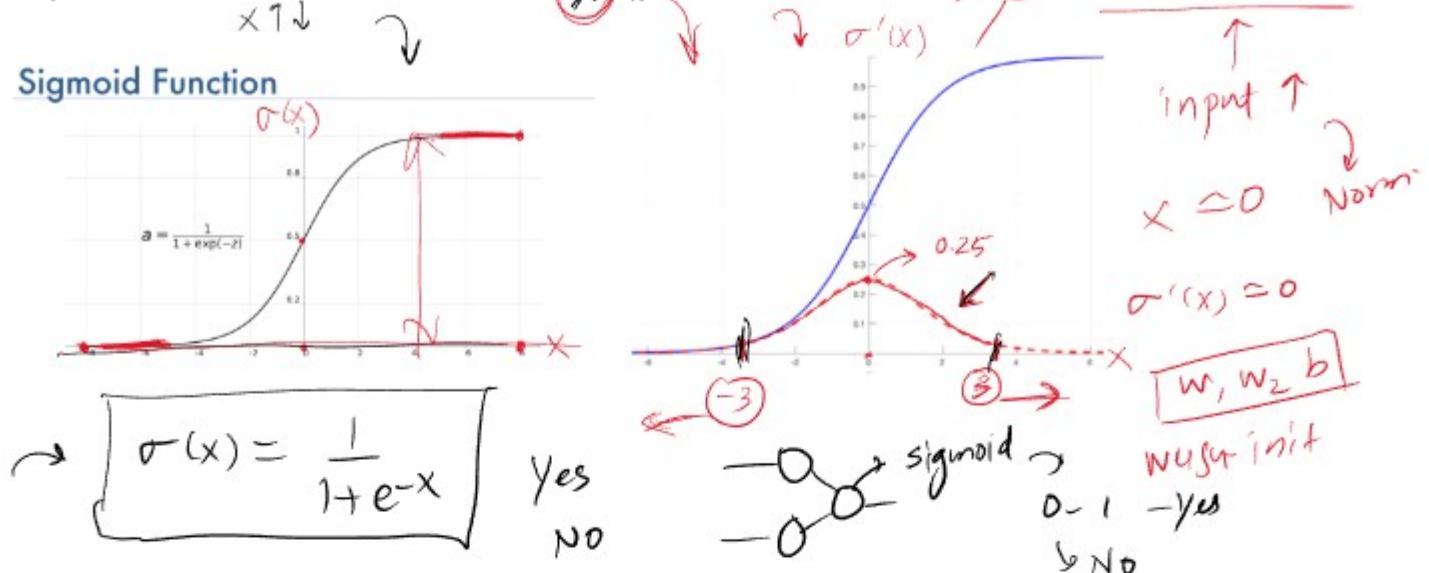
Ideal Activation Function

31 May 2022 14:50

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

Sigmoid Activation Function

31 May 2022 14:50

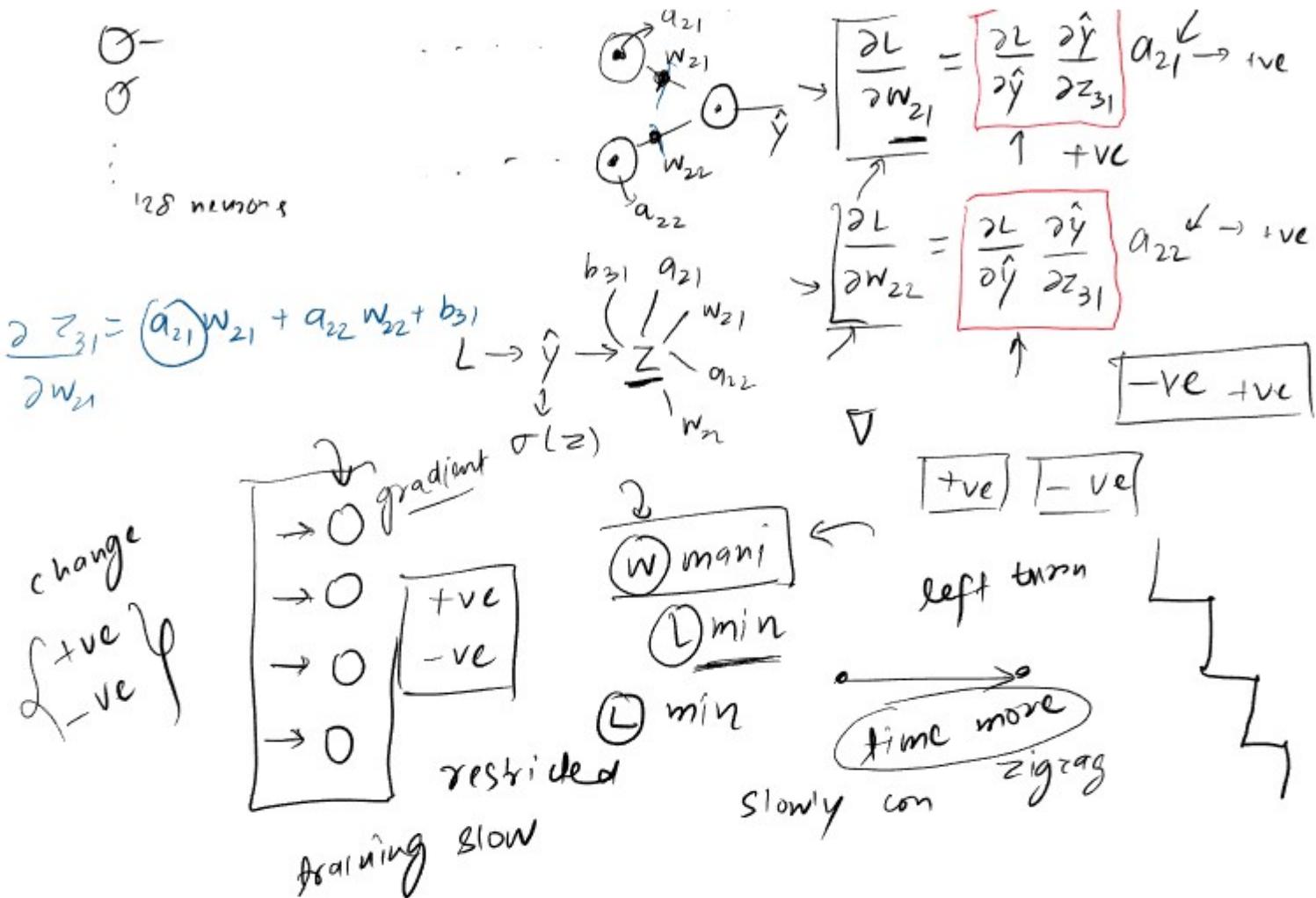


Advantages

- 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 \rightarrow update $w_n = w_o - \eta \frac{\partial L}{\partial w}$
 - \downarrow No update will take place training x
- 2) $\boxed{\text{Non Zero centered}}$
 - $\circlearrowleft \text{Normalized}$ \downarrow +ve $\rightarrow \boxed{\text{Training slow}}$
 - $\circlearrowleft \text{mean} = 0$ \downarrow -ve
 - \downarrow +ve - ve convergence problem $\boxed{0 \rightarrow 1}$
 - \downarrow a_{21} w_{21} $\rightarrow \frac{\partial L}{\partial w_{21}} = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}} & \frac{\partial \hat{y}}{\partial z_2} \end{bmatrix} a_{21} \downarrow$ +ve



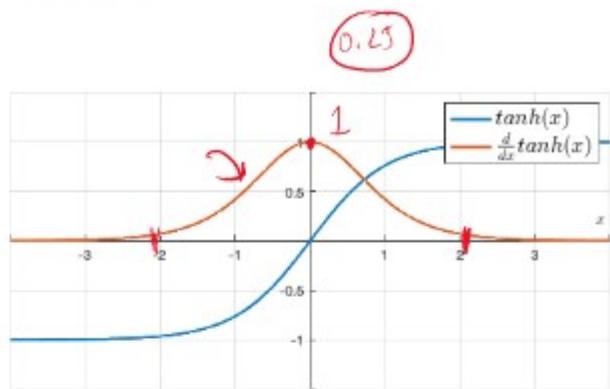
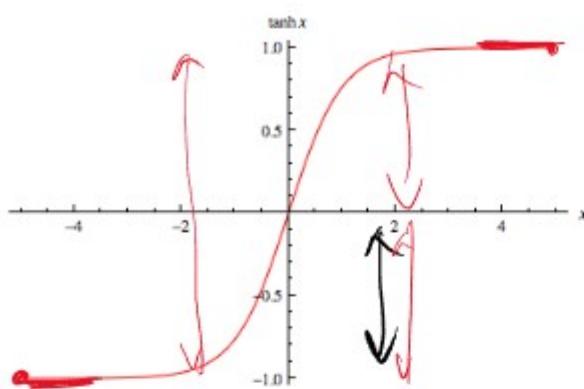
3) $\sigma(x) = \frac{1}{1+e^{-x}}$ \rightarrow computation time
~~exp~~

hidden layer
 \downarrow
 sigmoid
 \uparrow
 $\overline{\text{outpt}} \rightarrow b \text{ op}$

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

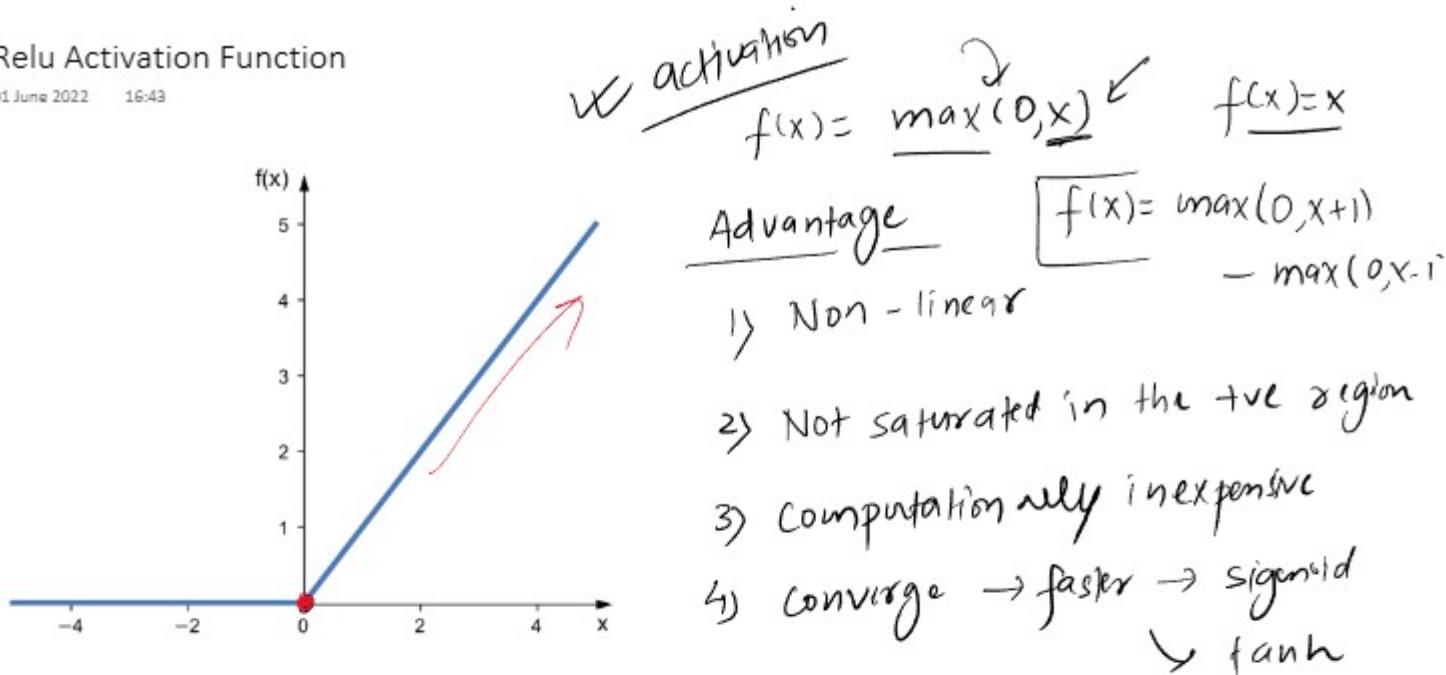
- 1) Non-linear
- 2) Differentiable
- 3) zero centered (+ve) training faster sigmoid (-ve)

Disadvantage

- saturating function → vanishing gradient prob
- computationally exp [exp] slow

Relu Activation Function

01 June 2022 16:43



Disadvantage

→ Differentiability

$$\begin{aligned} x < 0 &\rightarrow 0 \\ x \geq 0 &\rightarrow 1 \end{aligned}$$

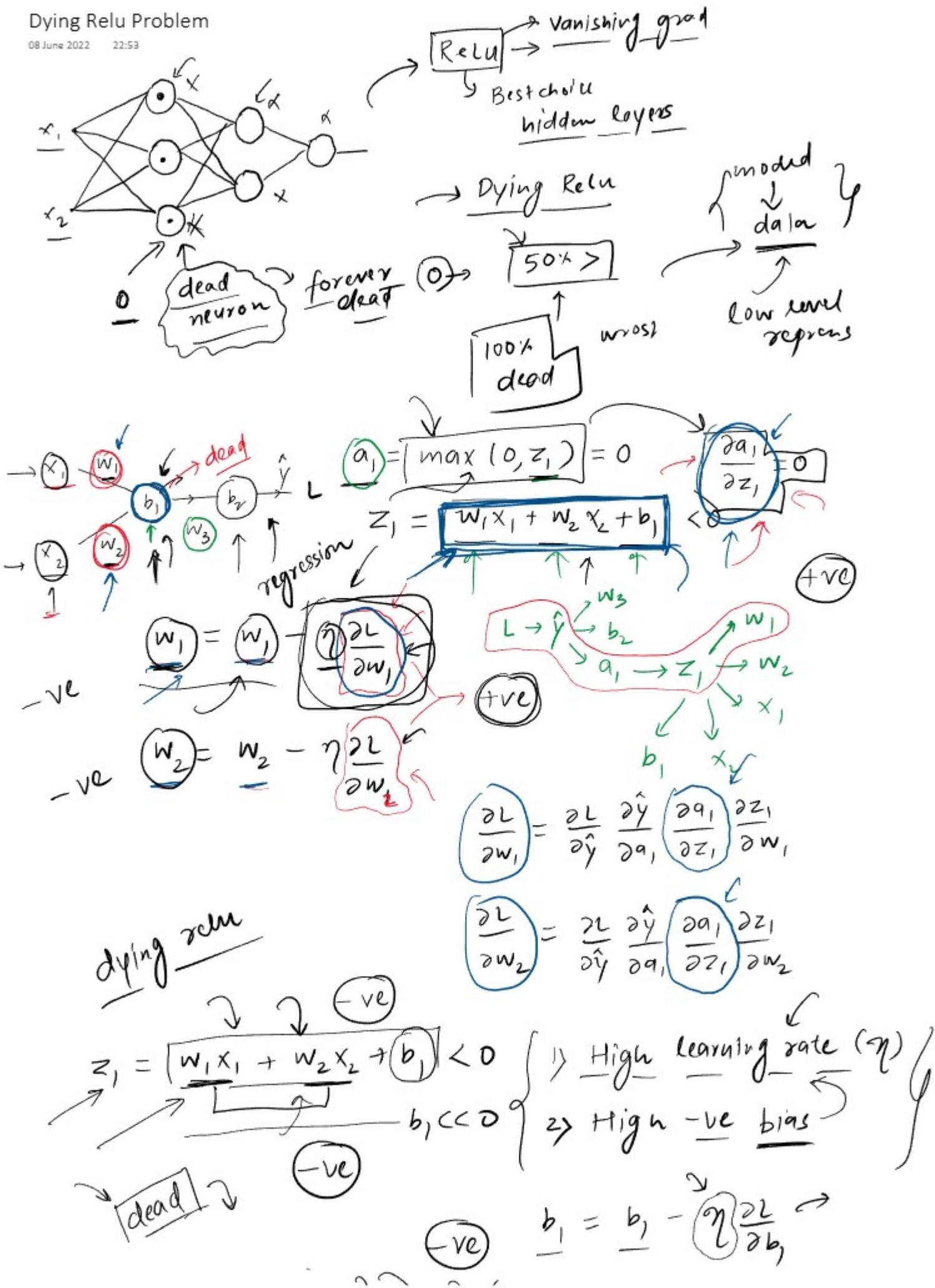
→ NonZero centered → sigmoid → Batch Normalization

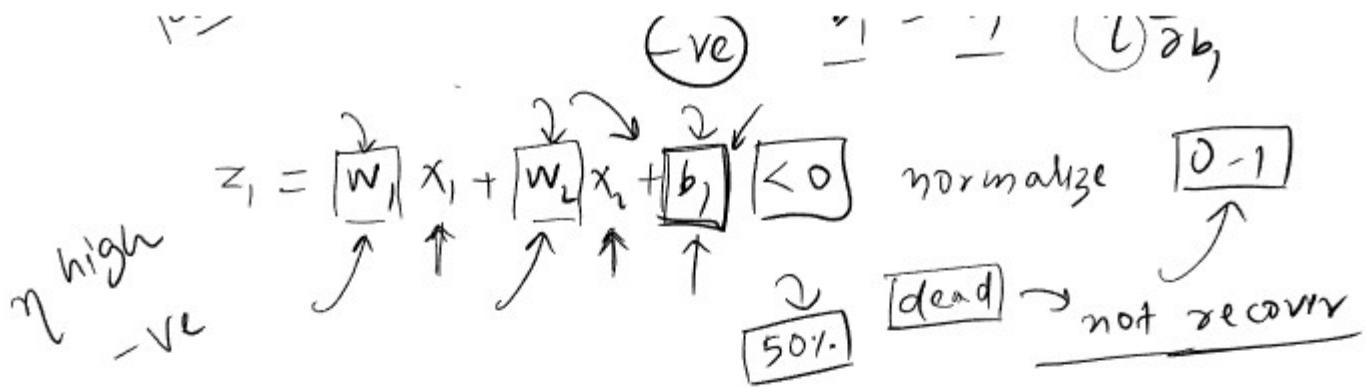
 ○ normalize →

↳ Dying ReLU problem

Dying Relu Problem

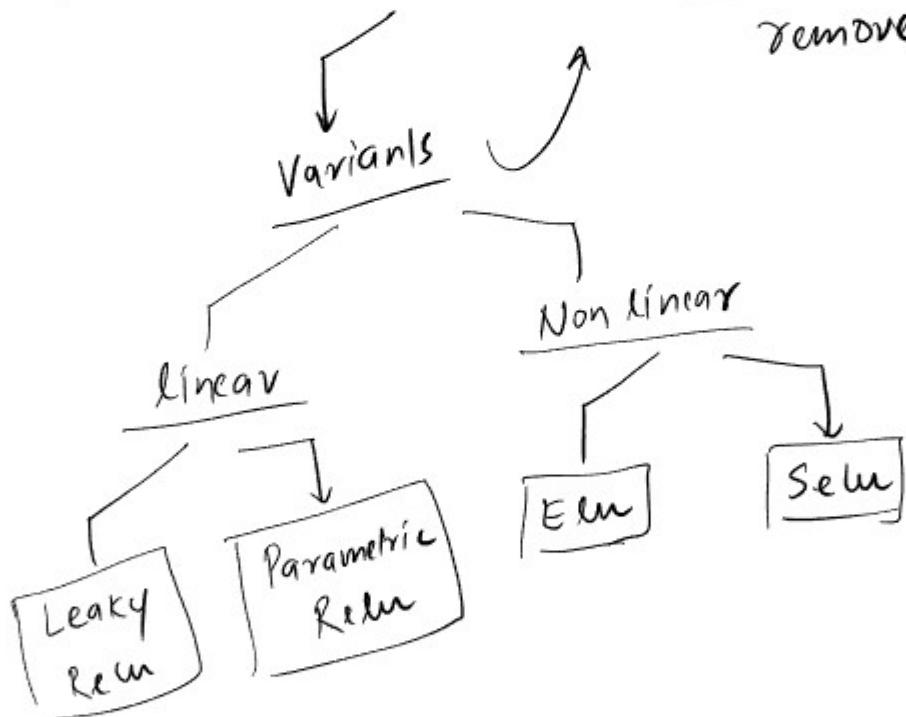
08 June 2022 22:53





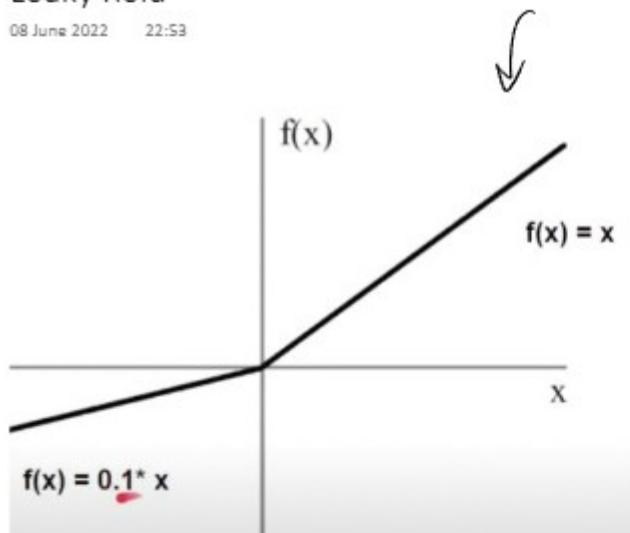
Solutions

- Set low learning rate
- bias \rightarrow +ve value \rightarrow [0.01]
- Don't use $\text{relu} \rightarrow$ variants $\frac{(+\text{ve keep})}{\text{remove}}$



Leaky Relu

08 June 2022 22:53



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

\downarrow

$$z \geq 0 \rightarrow z$$

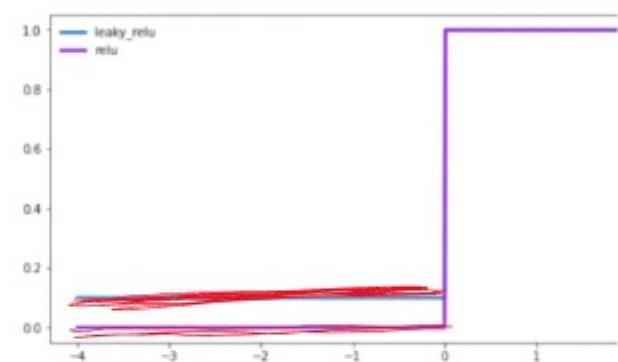
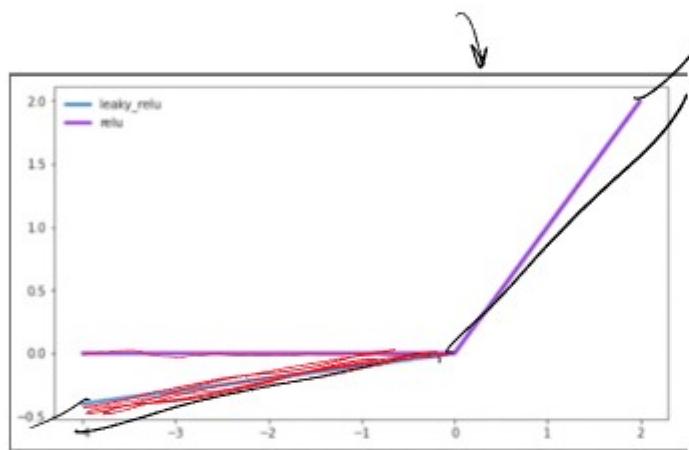
$$z < 0 \rightarrow \frac{1}{100}z \text{ (fraction of } z\text{)}$$

$$-ve(z) \rightarrow \neq 0 \quad \boxed{0.01}z$$

\hookrightarrow derivative

$$f'(z) \quad z \geq 0 \rightarrow \boxed{1}$$

$$z < 0 \rightarrow \boxed{0.01}$$



Advantages

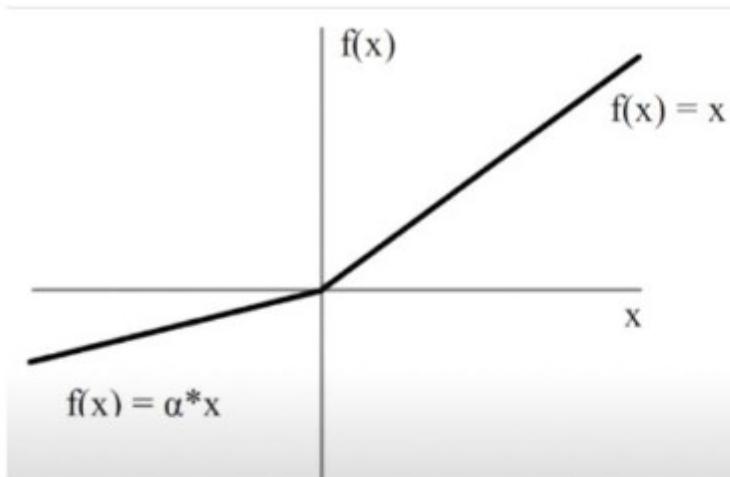
- Non-saturated \rightarrow Unbounded
- Easily computed
- No dying relu problem
- Close to 0 centered

Disadv

-ve/+ve

Parametric Relu

08 June 2022 22:53



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \\ 0.01x & \end{cases}$$

$\alpha \rightarrow$ trainable parameter

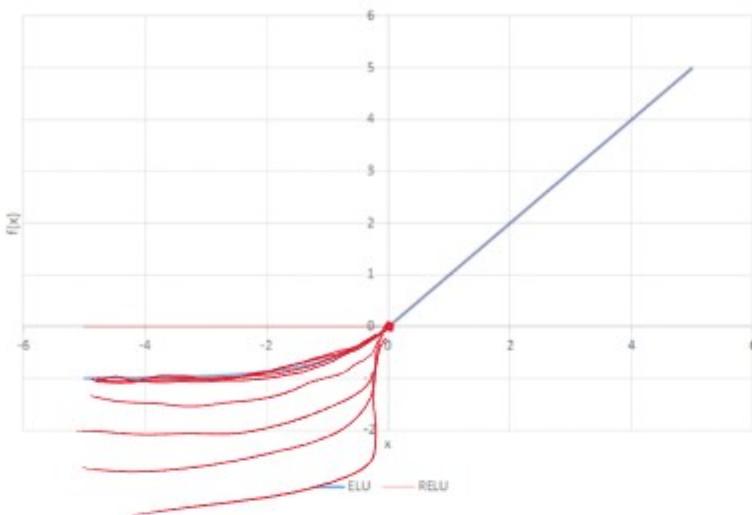
flexibility

depend on data

Elu - Exponential Linear Unit

09 June 2022 00:29

performance better ReLU



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

0.1 to 0.3

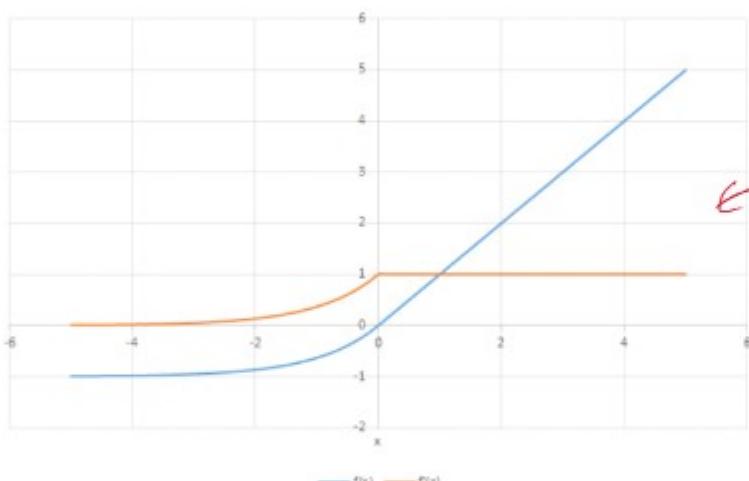
$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \text{ELU}(x) + \alpha & \text{if } x \leq 0 \end{cases}$$



Advantages

→ Close to zero centered
↳ convergence faster

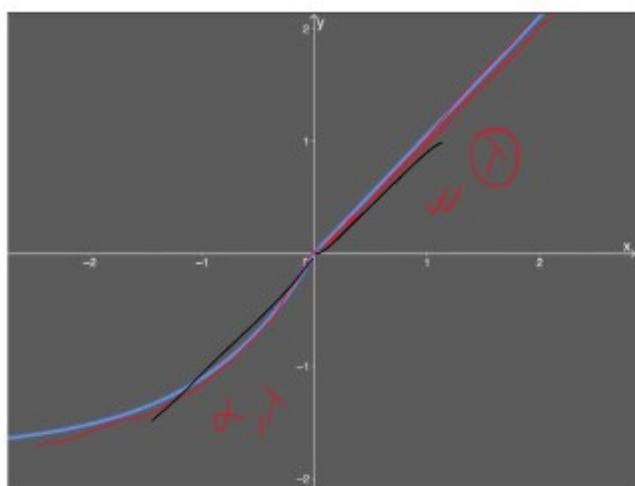
- Better generalized
- Dying ReLU x
- Always continuous as well as differentiable



Disadvantages

→ computation expensive $\textcircled{ex} 2$

Recent \rightarrow new \rightarrow 90+ pages



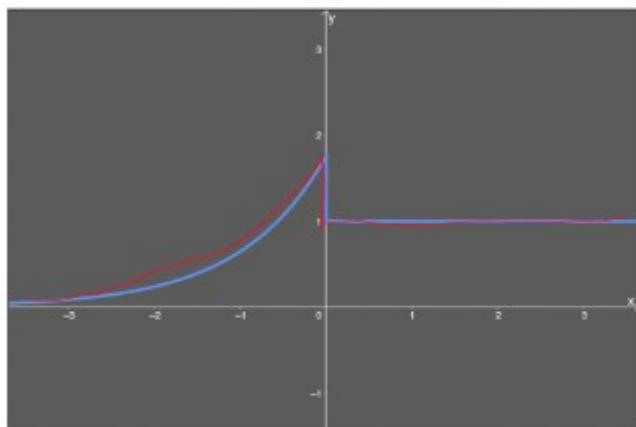
$$\text{SELU}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$

$$\alpha \approx 1.6732632423543772848170429916717$$

$$\lambda \approx 1.0507009873554804934193349852946$$

\rightarrow fixed
train but \propto

$$\text{SELU}'(x) = \lambda \begin{cases} 1 & \text{if } x > 0 \\ \alpha e^x & \text{if } x \leq 0 \end{cases}$$



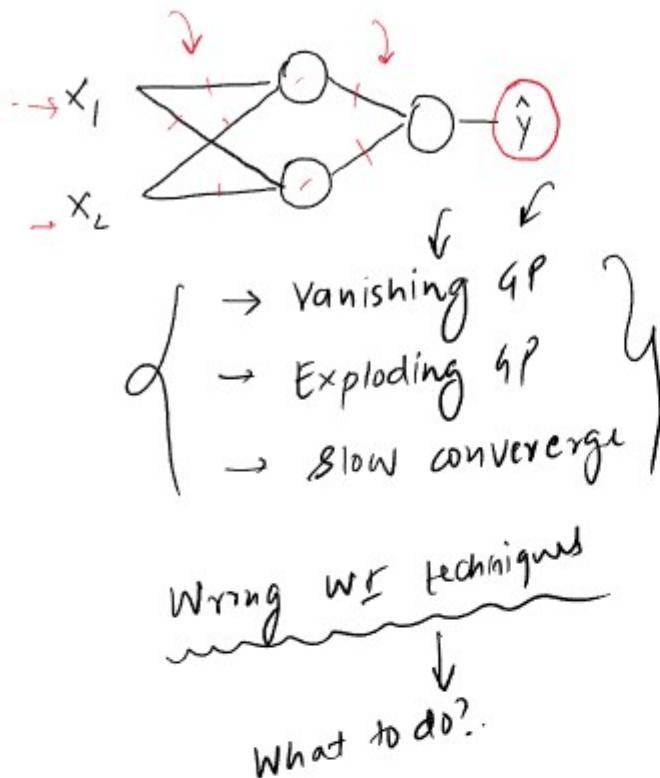
Self-normalizing \rightarrow activation
 \downarrow
normalized

$$m = 0 \quad \sigma = 1$$

converge faster

Why Weight Initialization is Important?

22 June 2022 12:48



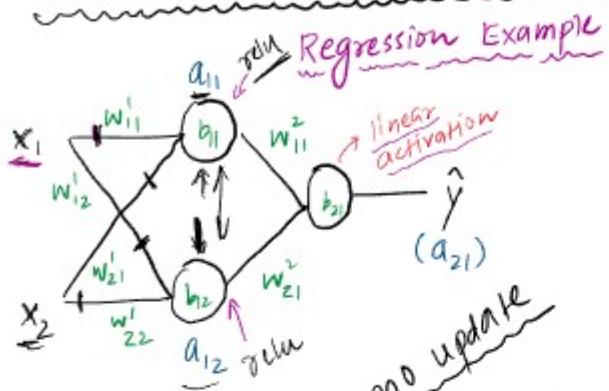
- 2000's Deep Learning
1. Initialize the parameters
 2. Choose an *optimization algorithm* ✓
 3. Repeat these steps:
 1. Forward propagate an input
 2. Compute the cost function
 3. Compute the gradients of the cost with respect to parameters using backpropagation
 4. Update each parameter using the gradients, according to the optimization algorithm

Sigmoid Activ of wrong weigh't init

What not to do?

22 June 2022 12:49

Case 1 → zero Initialization



$$w_{11}^1 = \underline{w_{11}^1} - \eta \frac{\partial L}{\partial w_{11}^1}$$

relu
tanh
sigmoid

$$w=0, b=0$$

$$a_{11} = \frac{e^{z_{11}} - e^{-z_{11}}}{e^{z_{11}} + e^{-z_{11}}} = 0$$

$$a_{11} = \max(0, z_{11}) = \boxed{0}$$

$$z_{11} = \underline{w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}} = 0$$

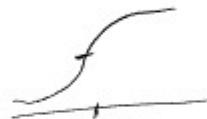
$$a_{12} = \max(0, z_{12}) = \boxed{0}$$

$$z_{12} = \underline{w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}} = 0$$

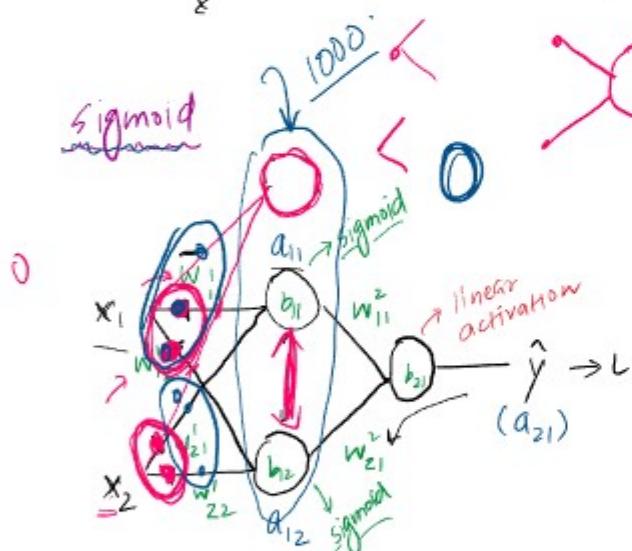
$$\boxed{a_{11} = a_{12}} \rightarrow \text{equal and } 0$$

$w=0$

No training will place



Non-linear part



$$a_{11} = \sigma(z_{11}) = 0.5$$

$$z_{11} =$$

$$a_{12} = 0.5$$

$$a_{11} = a_{12}$$

$$w = \underline{w} - \eta \frac{\partial L}{\partial w}$$



$$a_{11} = a_{12}$$

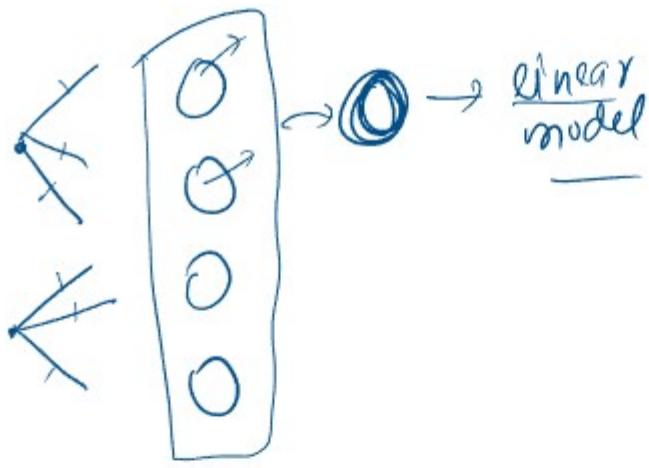
$$z_{11} = z_{12}$$

$$\frac{\partial L}{\partial w_{11}^1} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}}} X_1$$

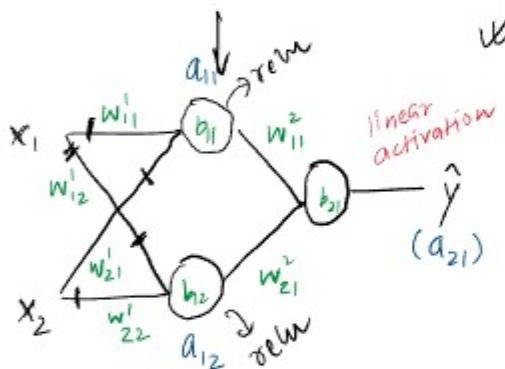
$$\frac{\partial L}{\partial w_{12}^1} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}}} X_1$$

$$\frac{\partial L}{\partial w_{21}^1} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}}} X_2$$

$$\frac{\partial L}{\partial w_{22}^1} = \boxed{\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}}} X_2$$



Case 2 → Non-0 constant value



$$w = 0.5 \quad b = 0.5 \quad \text{sigmoid} \rightarrow 0$$

$$a_{11} = \max(0, z_{11}) \rightarrow \text{non-zero}$$

$$z_{11} = w_{11}^T x_1 + w_{21}^T x_2 + b_{11} \neq 0$$

$$a_{12} = \max(0, z_{12}) \neq 0$$

$$z_{12} = w_{12}^T x_1 + w_{22}^T x_2 + b_{12} \neq 0$$

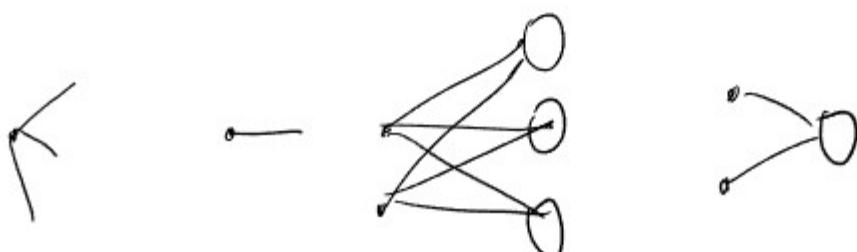
$$z_{12} = z_{11} \quad \boxed{a_{11} = a_{12}}$$

$$\boxed{\frac{\partial L}{\partial w_{11}^T}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} x_1$$

$$\boxed{\frac{\partial L}{\partial w_{12}^T}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} x_1$$

$$\frac{\partial L}{\partial w_{21}^T} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} x_2$$

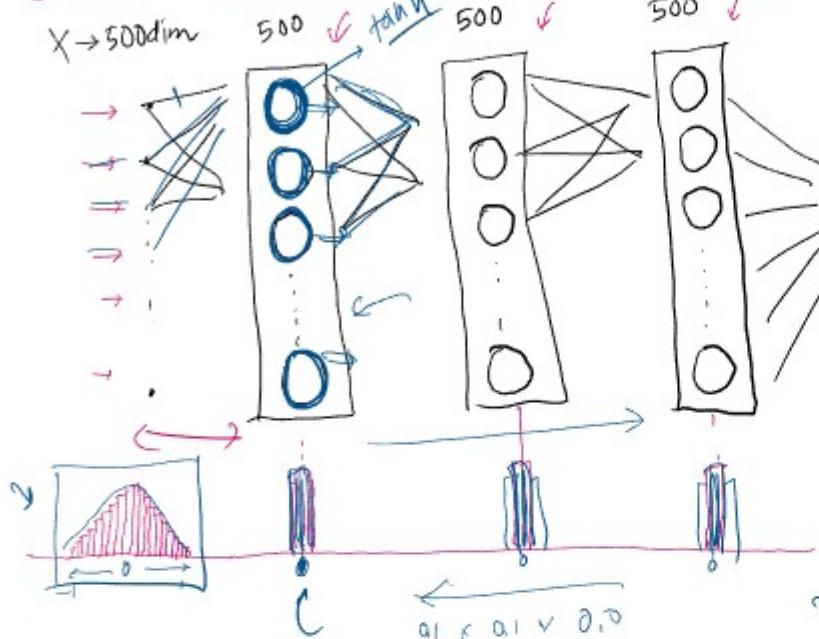
$$\frac{\partial L}{\partial w_{22}^T} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} x_2$$



of random init

VGP sigmoid

Case 3 - Random Initialization

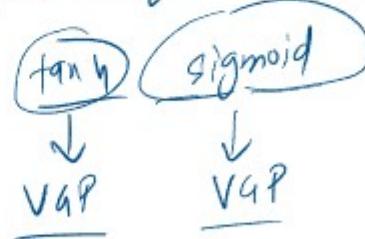


$$w = w - \eta \frac{\partial L}{\partial w} \approx 0$$

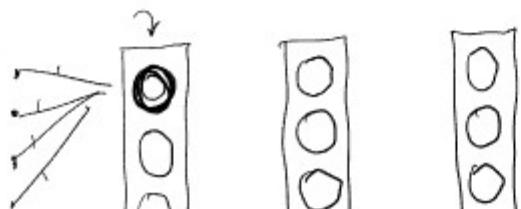
② vanishing gradients problem

ReLU \rightarrow small values
 \rightarrow slow convergence

small weights $\rightarrow 0.0007$



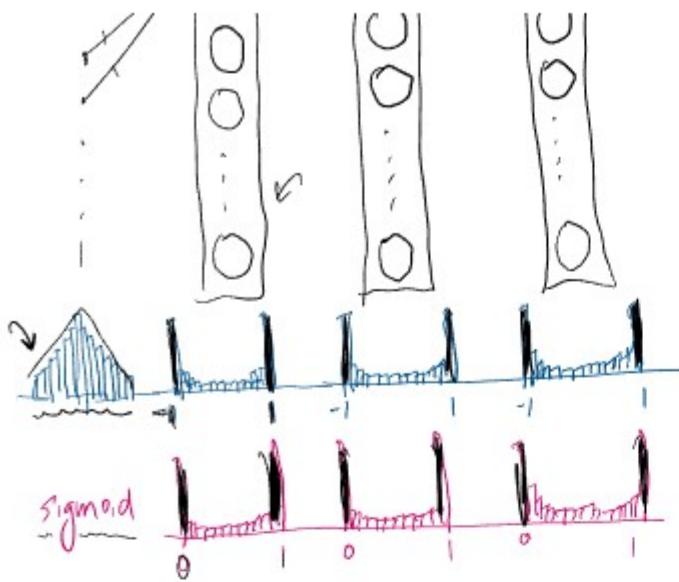
Random Initialization (Large values)



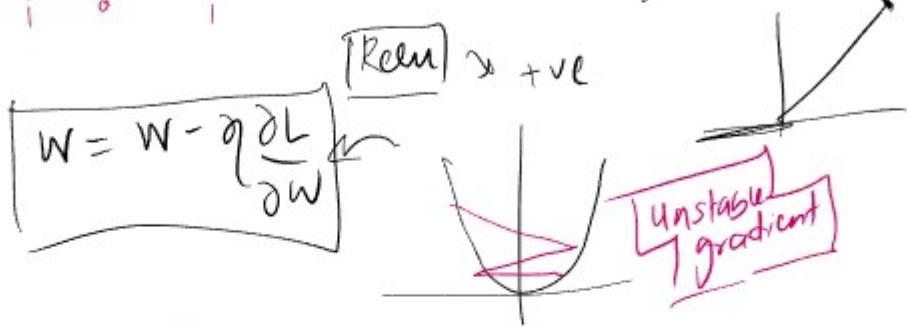
$\text{np.random.randn}(500, 500) \rightarrow [0-1]$

$\sum w_i x_i \rightarrow 500 \rightarrow [100 \rightarrow 500]$

0-1 tanh / sigmoid



0-1 tanh / sigmoid
 Saturates $\rightarrow g(z)$
 $\sum w_i x_i \rightarrow [100] \rightarrow$ saturate
 $\sum w_i x_i \rightarrow$ slow training
 Vanishing gradient problem
 gradients

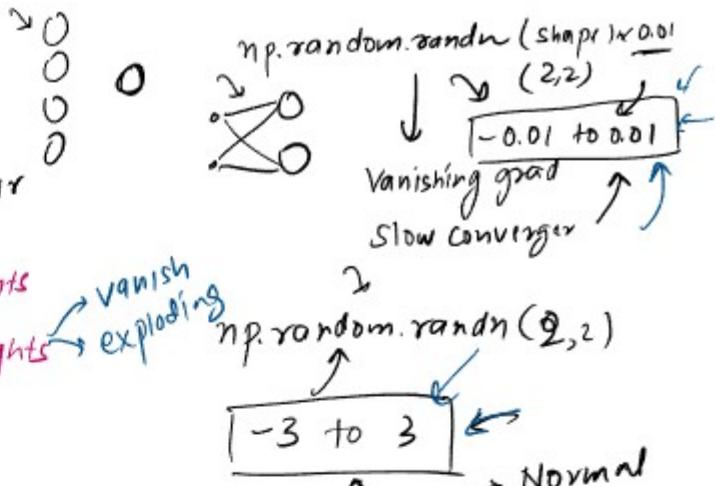


- zero x
- non-zero same x
- small random!
- large random x

Key Insights

22 June 2022 12:49

- Things not to do -
- 1) Zero initialization → no training
- 2) Non-zero constant initialization → linear
- 3) Random initialization with small weights → vanishing grad
- a) Random initialization with large weights → exploding



What can be done

→ Heuristics (Jugad) → practical sol^④

-1 w¹ 0-1
Intuition

0 to -0.01 to 0.01

Xavier/Glorot Init
(tanh, sigmoid)
He Init
(relu)
Keras

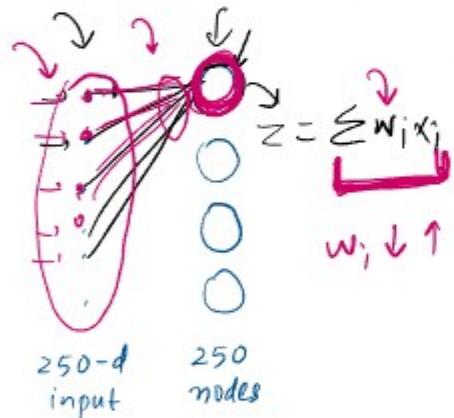
Normal

Uniform

Normal

uniform

formulas



small
np.random.randn(250, 250) * 0.01

large
np.random.randn(250, 250) * 1

250
-250 to 250

$$\text{Variance} = \frac{1}{n} = \frac{1}{250}$$

$$* \sqrt{\frac{1}{250}} \rightarrow [sd]$$

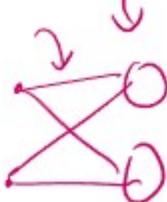
:

$$n \rightarrow \uparrow \downarrow \quad \frac{1}{n} \downarrow \rightarrow \text{var}$$

Xavier init
(Normal)

$$\sqrt{\frac{1}{\text{fan-in}}}$$

of inputs
coming to
the node



$$\text{np.random.randn}(2,2) * \sqrt{\frac{2}{2}}$$

$$\sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$$

He Normal
 $\text{np.random.normal}(\mu, \sigma^2)$

$$\text{fan-in} \times \sqrt{\frac{2}{\text{fan-in}}}$$

Uniform Distribution

Xavier Uniform

$[-\text{limit}, \text{limit}]$



$$\text{limit} = \sqrt{\frac{6}{(\text{fan-in} + \text{fan-out})}}$$

He Uniform

$[-\text{limit}, \text{limit}]$

$$\text{limit} = \sqrt{\frac{6}{\text{fan-in}}}$$

Keras

Weight Initialization Techniques

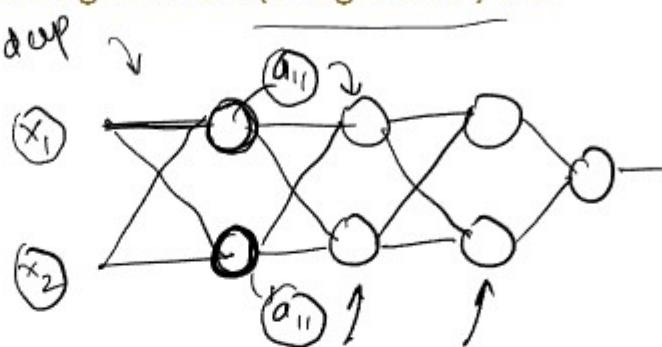
22 June 2022 12:49

What is Batch Norm?

27 June 2022 11:00

Batch-Normalization (BN) is an algorithmic method which makes the training of Deep Neural Networks (DNN) faster and more stable.

- It consists of normalizing activation vectors from hidden layers using the mean and variance of the current batch. This normalization step is applied right before (or right after) the nonlinear function.



normalized

$$m=0 \quad sd=1$$

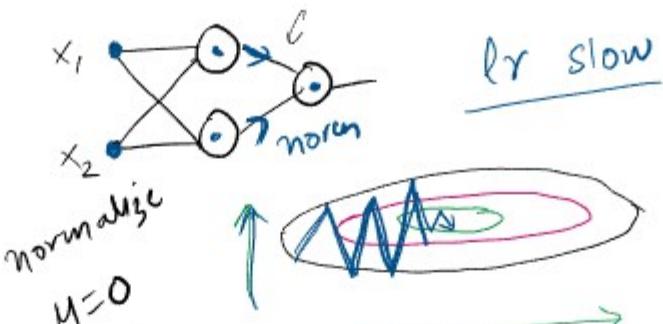
normalize

$$m=0 \quad sd=1$$

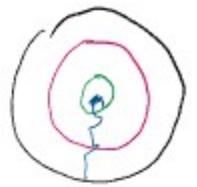
Why use Batch Norm?

30 June 2022 16:08

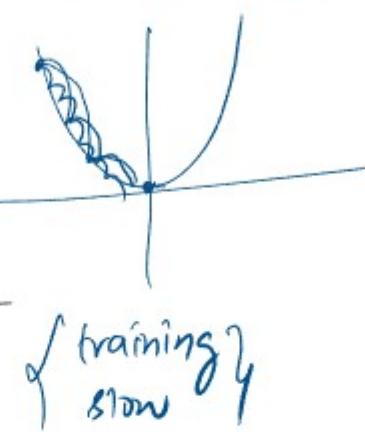
<u>cgpas</u>	<u>iq</u>	<u>placed</u>
7	70	1
8	80	0
9	90	1
6	60	0



lr slow



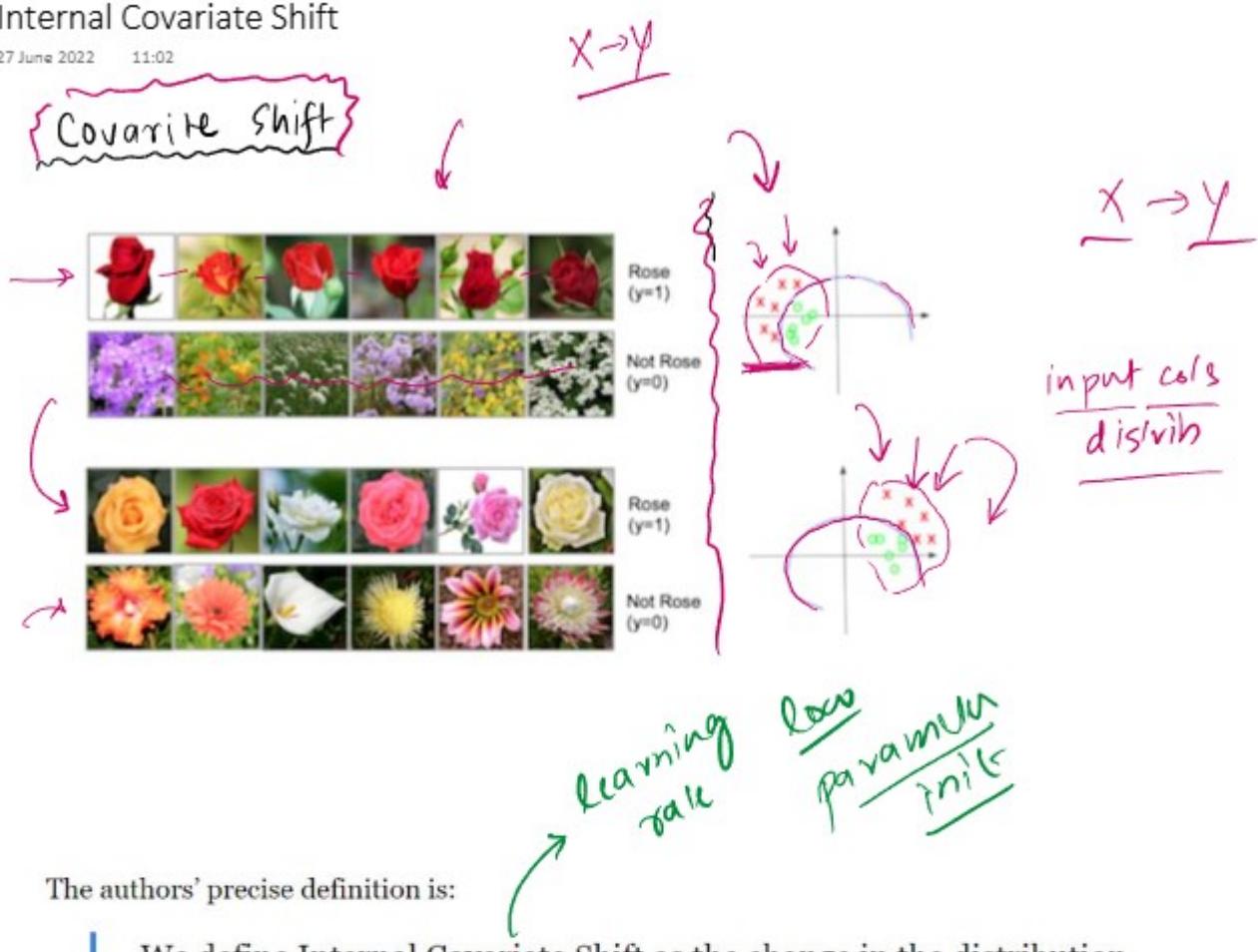
training faster stable



{ training }
slow

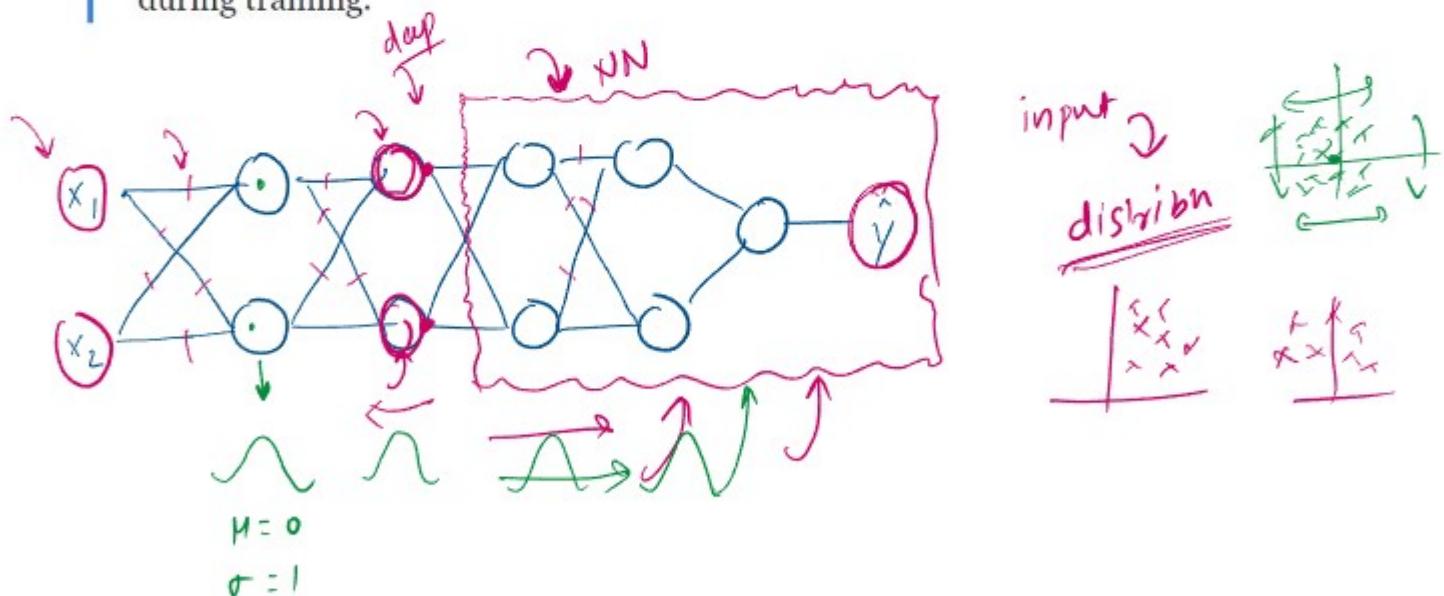
Internal Covariate Shift

27 June 2022 11:02



The authors' precise definition is:

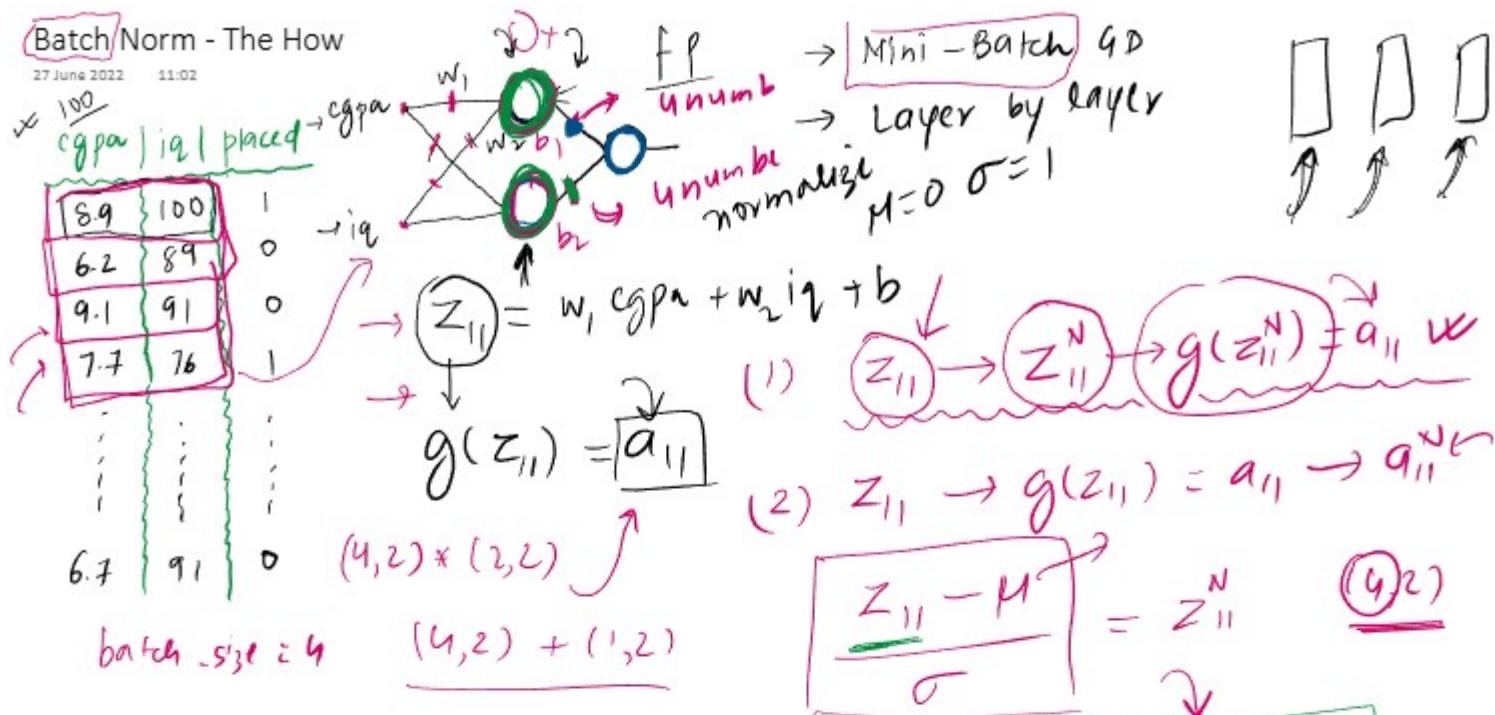
We define Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.





Batch Norm - The How

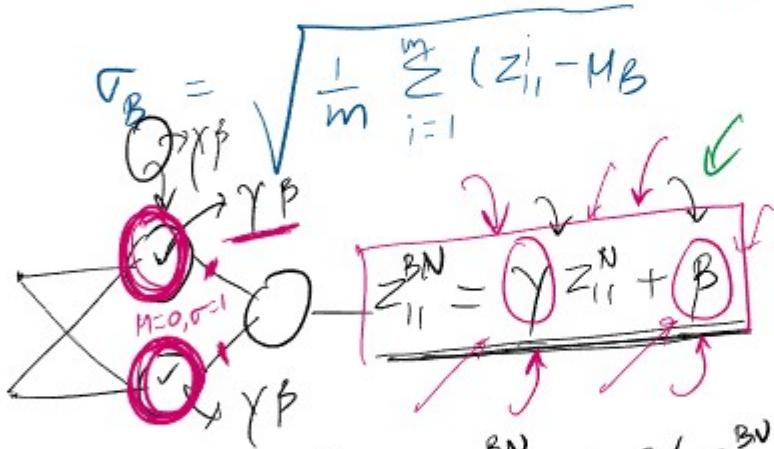
27 June 2022 11:02



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

$$m = 4$$

$$z_{11}^i = \frac{z_{11}^i - M_B}{\sigma_B + \epsilon}$$



$$\begin{aligned} & \text{learnable parameters } \gamma, \beta \\ & \gamma = 1, \beta = 0 \end{aligned}$$

$$z_{11} \rightarrow z_{11}^N \rightarrow z_{11}^{BN} \rightarrow g(z_{11}^{BN}) = a_{11}$$

$z_{11} \rightarrow \text{normall} \quad \mu=0 \quad \sigma=1$

$$\gamma = \sigma + \epsilon \quad \beta = M$$

flexibility

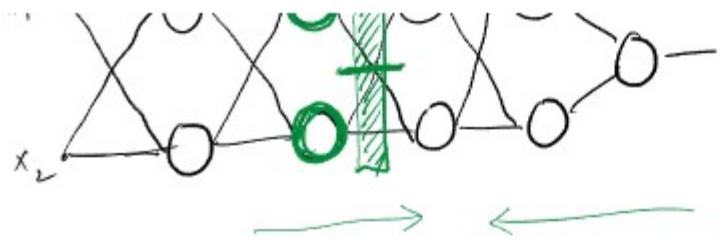
$$z_{11} \rightarrow z_{11}^N \rightarrow z_{11}^{BN}$$

Batch Norm → Layer



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

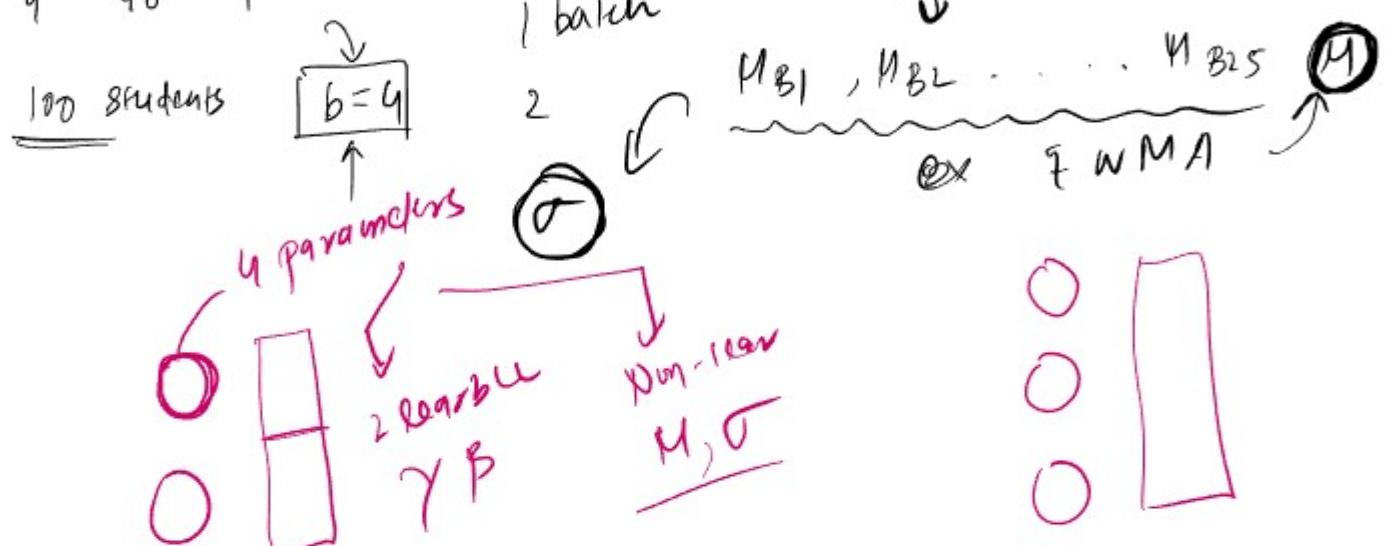
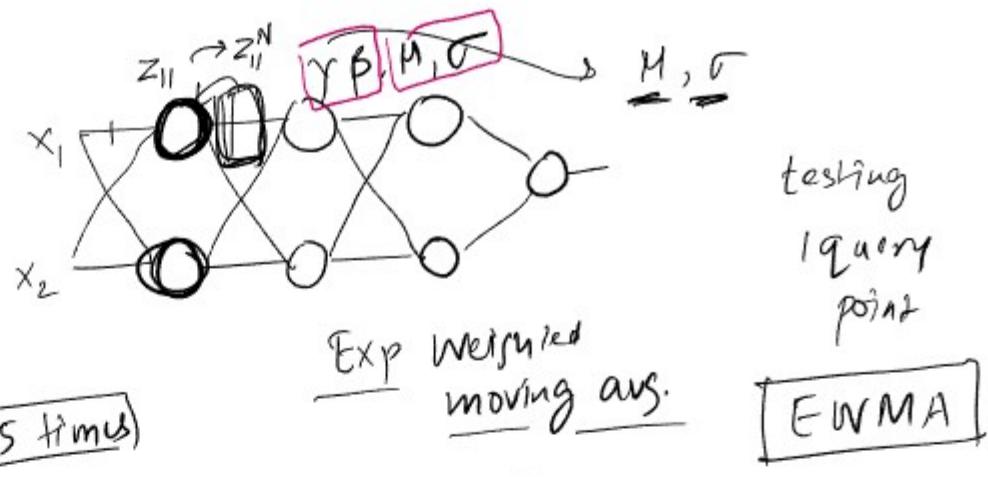
Keras



Batch Norm during test

27 June 2022 11:03

cgpa	ia	placed
8	80	1
7	70	0
6	60	1
:	:	:
9	90	1

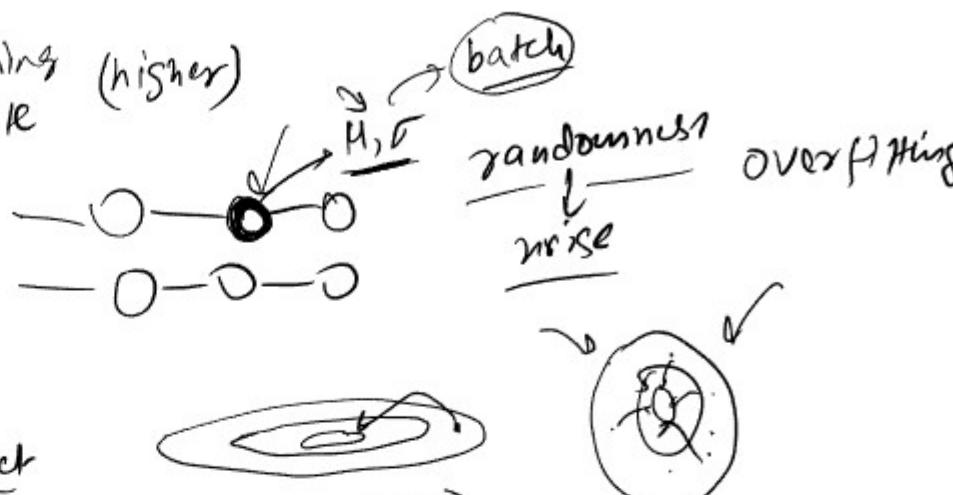


$$3 \times 4 = 12$$

(6) (6)

Advantages

27 June 2022 11:02

- 1) stable → hyper ^{para} → wider range of values
- 2) faster → learning rate (higher)
- 3) Regularizer → 
dropout
- 4) weight init impact reduce

Keras Implementation

27 June 2022 11:03

```
model = Sequential()  
  
model.add(Dense(3,activation='relu',input_dim=2))  
model.add(BatchNormalization()) ←  
model.add(Dense(2,activation='relu')) ←  
model.add(BatchNormalization()) ←  
model.add(Dense(1,activation='sigmoid'))
```

