

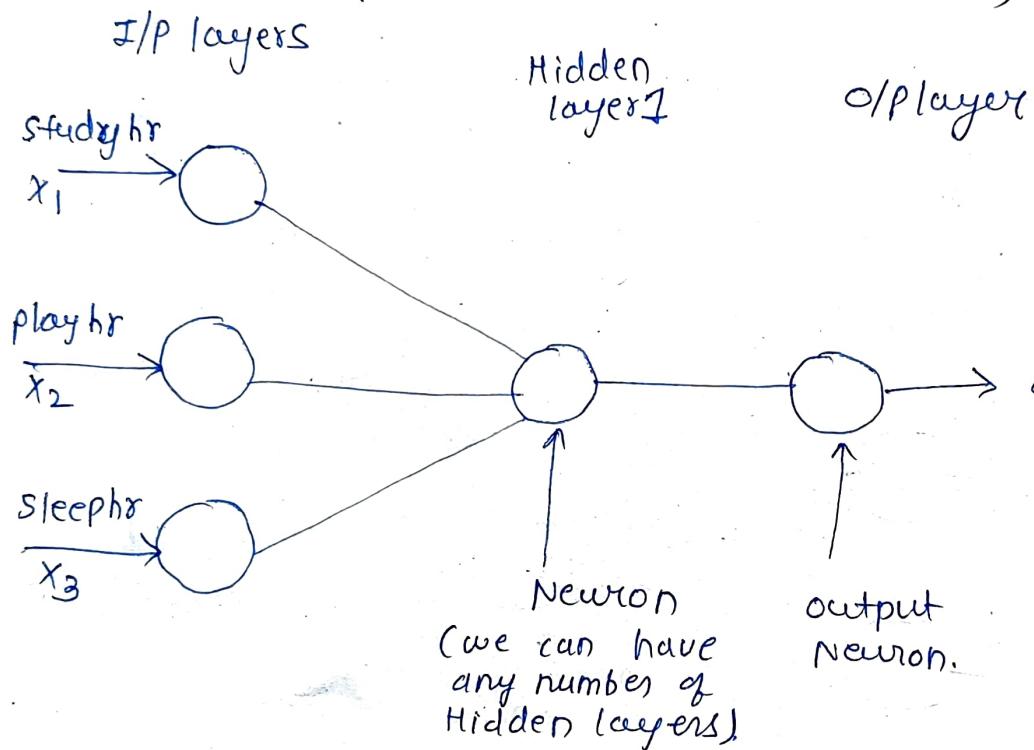
① Deep Learning.

why Deep learning ??

- 1) we have huge amount of data.
- 2) Advancement in computer Hardware (GPU's)

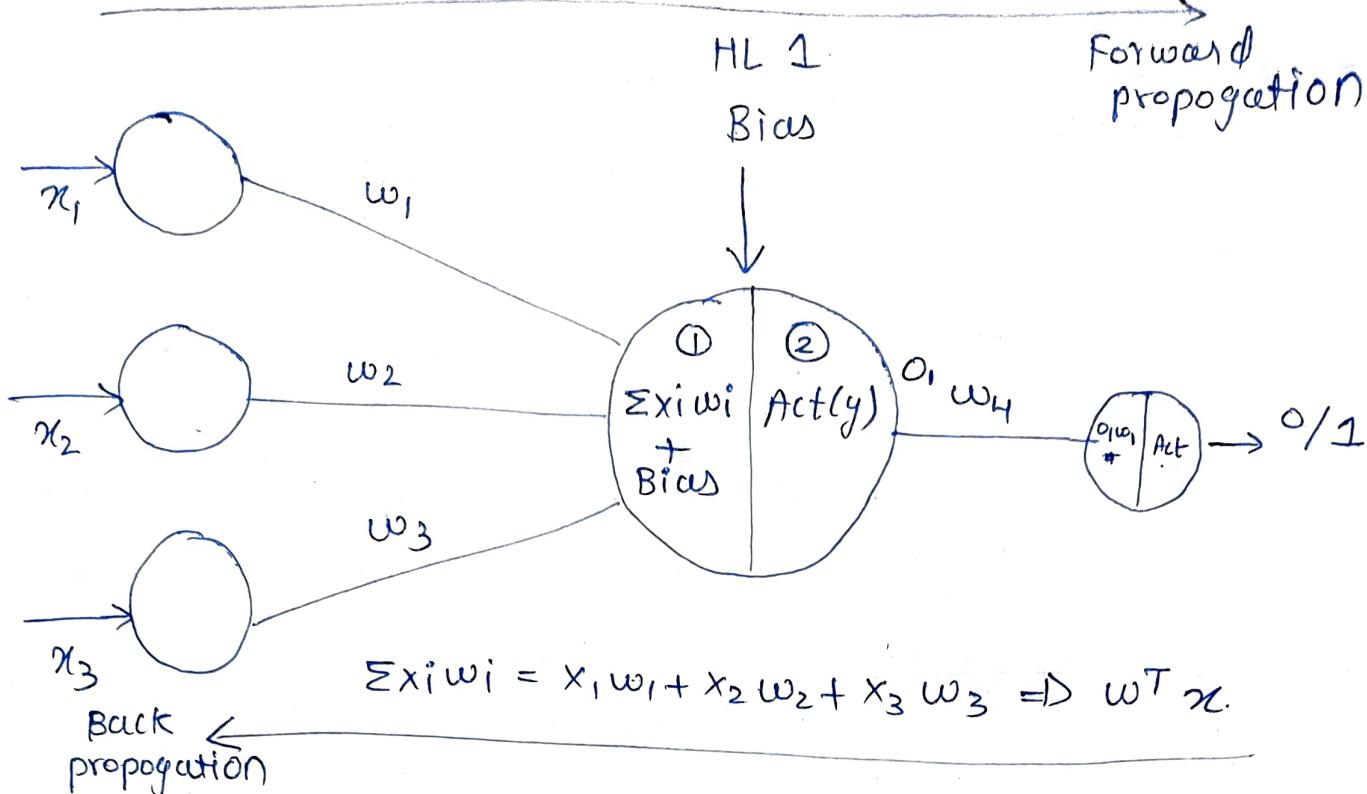
*perceptron

(single layered Neural network)



Dataset (Binary classification)

Study hr	play hr	Sleep hr	Pass/Fail
7	3	7	1
2	5	8	0
4	3	7	1



weights are the real values that are attached to each i/p/feature and they convey the importance of that corresponding feature in predicting the final o/p.

weights convey following things

1) Importance of feature in predicting o/p value.
weight ↑ feature value ↑.

2) Relationship b/w feature and target/ output

Ex: car price, car popularity O/P
Buy car & car popularity Buy or not.

Buy car & /car price

$$w_1(\text{car price}) + w_2(\text{car pop}) + b = 0$$

price of car ↑ $w_1(\text{car price}) \uparrow \Rightarrow$ we buy car.
but this we don't want

so we will assign w_1 a lower value.
or -ve value.

3) wt plays an important role in changing orientation that separates data classes.

- Bias is used to shift activation function to left or right.
- Bias initialization is automatic by neural network

Sigmoid Activation function (for binary classification)

$$\hookrightarrow \text{sigmoid} = \frac{1}{1+e^{-y}} = \frac{1}{1+e^{-(\sum x_i w_i + b)}} \begin{cases} > 0.5 & 1 \\ < 0.5 & 0 \end{cases}$$

① Forward propagation

Input → Assign weights → Input \times weight + bias.
 ↓
 Actual value y \hat{y} output ← Activation($\sum x_i w_i + b$) ←

$\text{Loss} = y - \hat{y}$ (Reduce Loss to close to zero)
 for this we need to update wt.

This is done using back propagation.
 (using optimizers)

② Back propagation

update wt's using optimizers to reduce loss. ($y - \hat{y}$)

1) I/P layers

2) weights

3) Bias Bias gets added in each hidden layer.

4) Activation function

5) loss function $f(y - \hat{y})$

6) optimizers

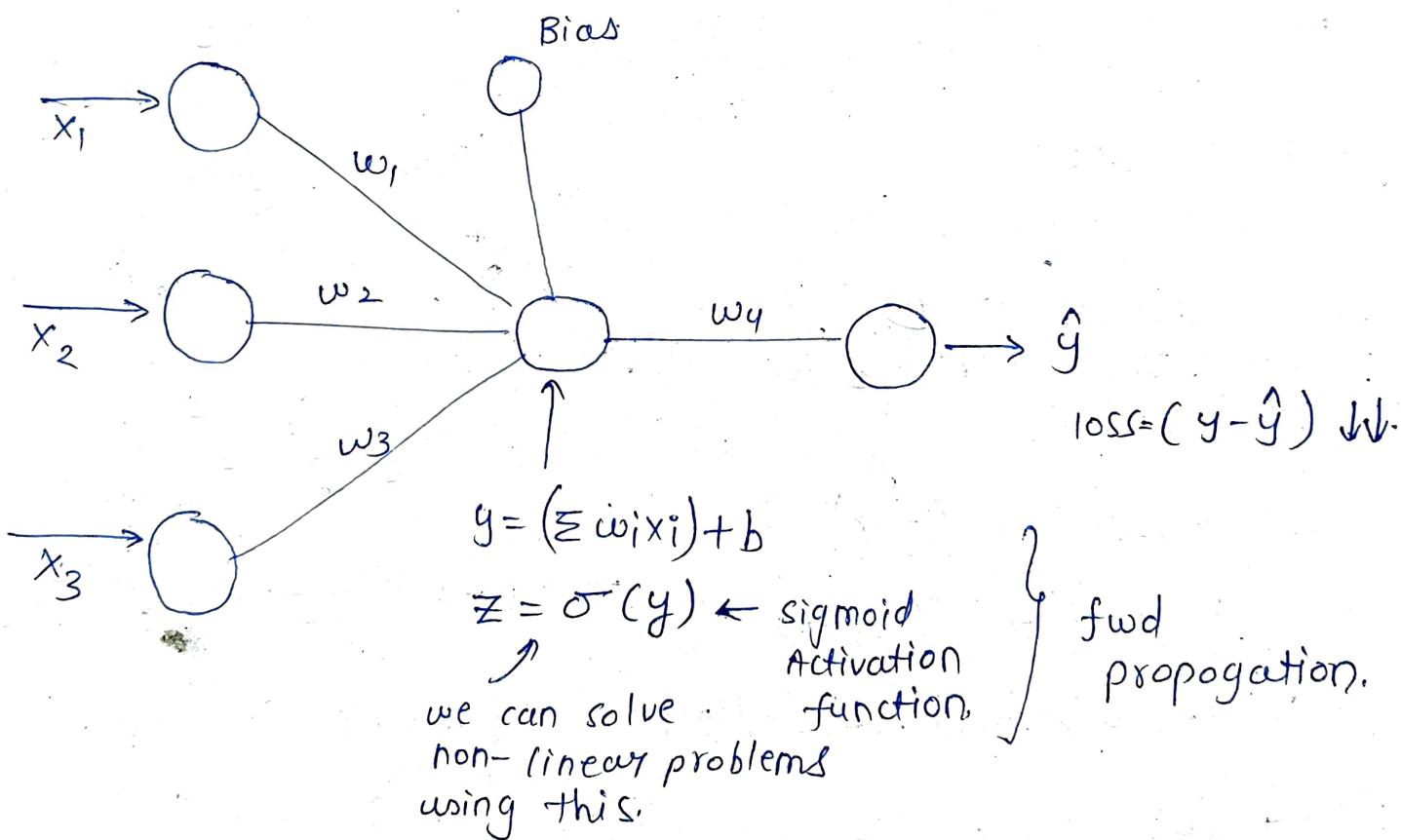
7) update the wt

(2)

forward propagation,

Backward propagation

F

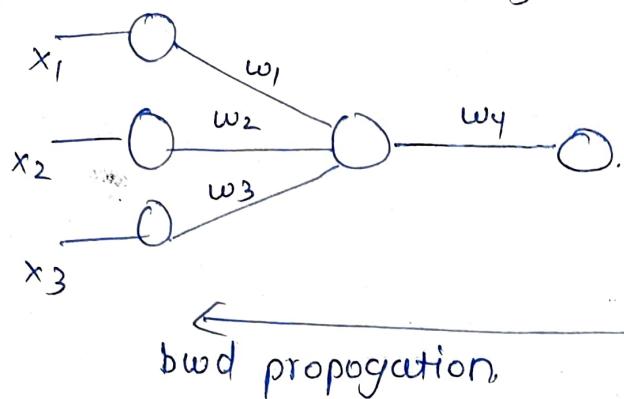


→ we need to update wts using back propogation

① weight updation formula.

② chain rule of differentiation.

→ fwd propagation



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

* weight update formula

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

$$\frac{\partial h}{\partial w_{\text{old}}} = \text{slope}$$

1) when we are at \textcircled{A}

we want to reach w_c for
global minima

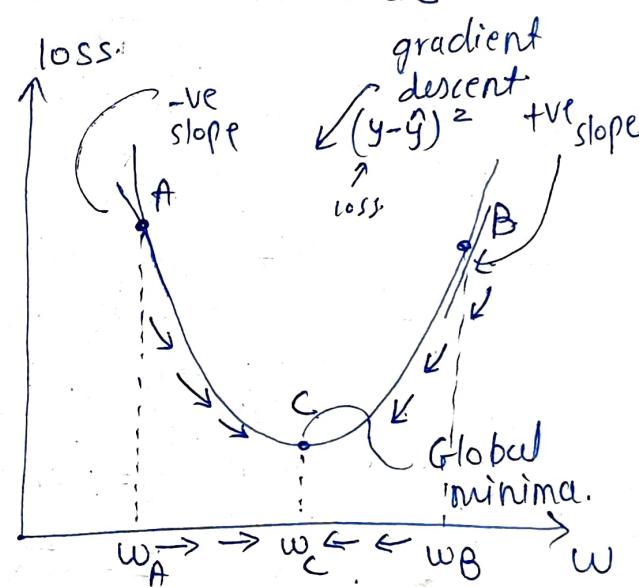
$$\begin{aligned} w_{\text{new}} &= w_{\text{old}} - \eta (-\text{ve}) \\ &= w_{\text{old}} + \eta (\text{const}) \end{aligned}$$

∴ $w_{\text{new}} > w_{\text{old}}$ ⇒ updating this multiple time we will
reach w_c .

2) At \textcircled{B} slope = +ve

$$\begin{aligned} w_{\text{new}} &= w_{\text{old}} - \eta (+\text{ve}) \\ &= w_{\text{old}} - \eta (\text{const}) \end{aligned}$$

∴ $w_{\text{new}} < w_{\text{old}}$ ⇒ updating this multiple
times we will reach
global minima.

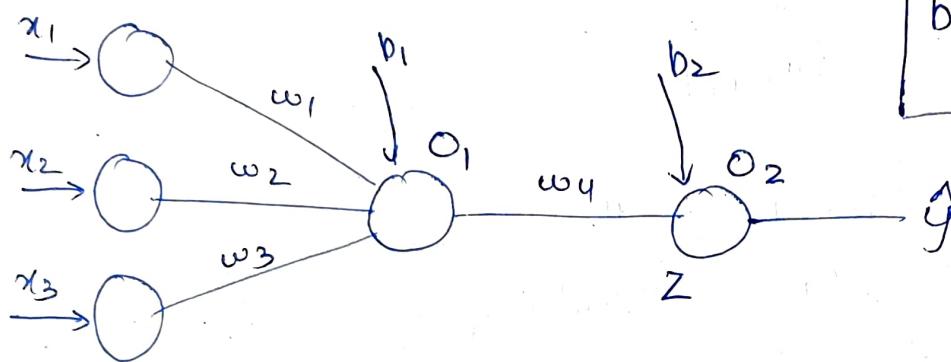


* Importance of learning rate (η)

- 1) It should be a smaller number so that we can slowly reach the global minima.
- 2) For larger learning rate we may not reach global minima.

② chain rule of differentiation

Bias update formula.



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

$$w_4_{\text{new}} = w_4_{\text{old}} - \eta \frac{\delta \text{loss}}{\delta w_4_{\text{old}}}$$

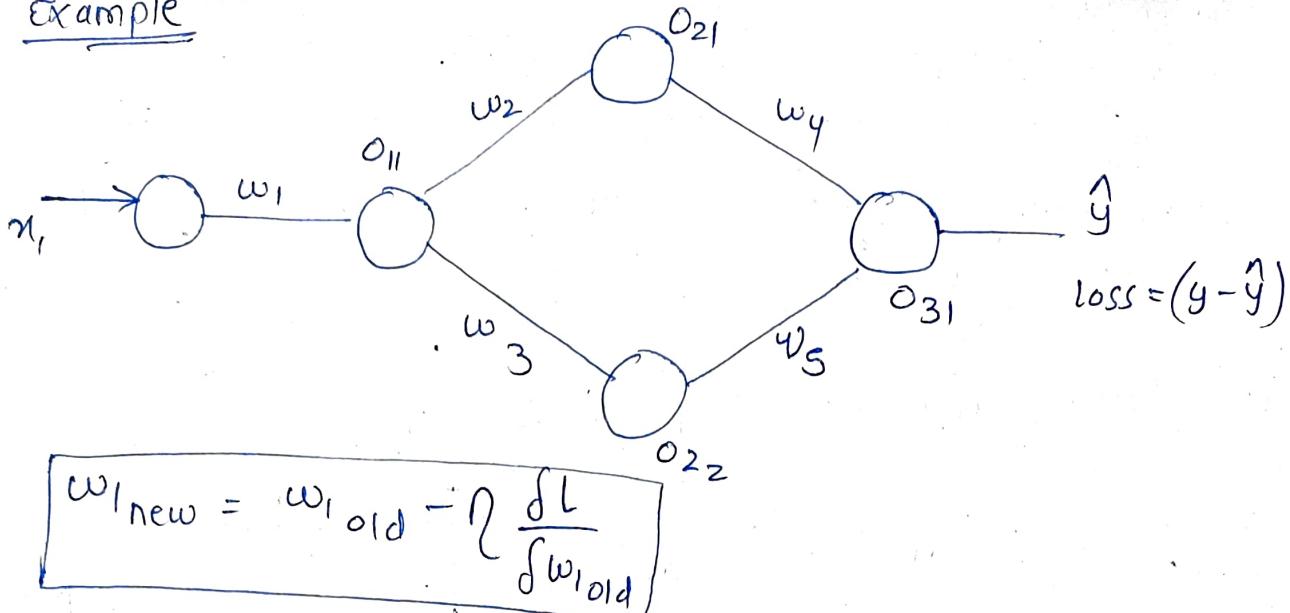
$$\frac{\delta L}{\delta w_4_{\text{old}}} = \frac{\delta L}{\delta O_2} \times \frac{\delta O_2}{\delta w_4} \leftarrow \text{chain rule}$$

$$w_1_{\text{new}} = w_1_{\text{old}} - \eta \frac{\delta L}{\delta w_1_{\text{old}}}$$

$$\frac{\delta L}{\delta w_1_{\text{old}}} = \frac{\delta L}{\delta O_2} \times \frac{\delta O_2}{\delta O_1} \times \frac{\delta O_1}{\delta w_1_{\text{old}}} \leftarrow \text{chain rule}$$

$$\frac{\delta O_2}{\delta w_4} \times \frac{\delta w_4}{\delta O_1}$$

Example



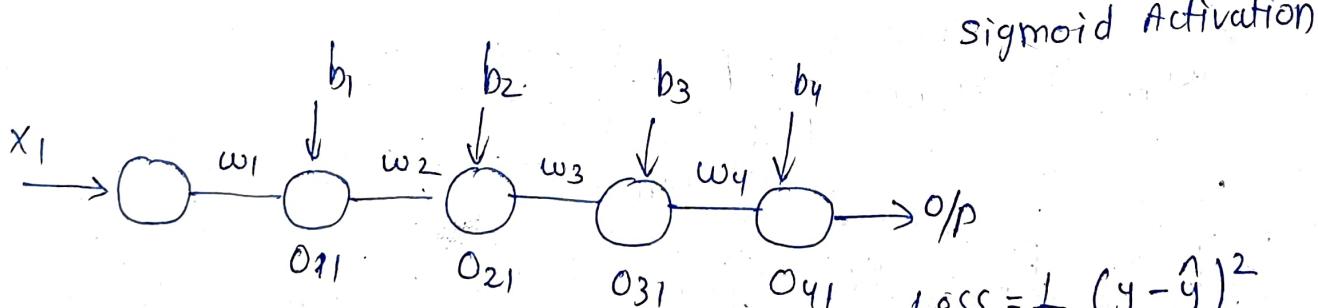
$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\delta L}{\delta w_{1,\text{old}}}$$

$$\begin{aligned} \frac{\delta L}{\delta w_{1,\text{old}}} &= \left[\frac{\delta L}{\delta o_{31}} \times \frac{\delta o_{31}}{\delta o_{21}} \times \frac{\delta o_{21}}{\delta o_{11}} \times \frac{\delta o_{11}}{\delta w_{1,\text{old}}} \right] \\ &\quad + \\ &\quad \left[\frac{\delta L}{\delta o_{31}} \times \frac{\delta o_{31}}{\delta o_{22}} \times \frac{\delta o_{22}}{\delta o_{11}} \times \frac{\delta o_{11}}{\delta w_{1,\text{old}}} \right] \end{aligned}$$

← chain rule of derivatives.

* Vanishing gradient problem.

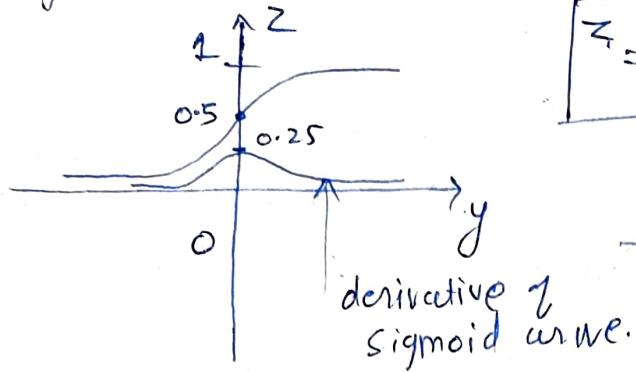
consider a very deep neural network



$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\delta L}{\delta w_{1,\text{old}}}$$

$$\begin{aligned} \frac{\delta L}{\delta w_{1,\text{old}}} &= \frac{\delta L}{\delta o_{41}} \times \frac{\delta o_{41}}{\delta o_{31}} \times \frac{\delta o_{31}}{\delta o_{21}} \times \frac{\delta o_{21}}{\delta o_{11}} \times \frac{\delta o_{11}}{\delta w_{1,\text{old}}} \end{aligned}$$

Sigmoid Activation.



$$z = \frac{1}{1+e^{-x}} \quad \begin{cases} > 0.5 = 1 \\ < 0.5 = 0 \end{cases}$$

(3)

$$0 < \delta(\text{sigmoid}) \leq 0.25$$

↑ derivative condition.

$$\frac{\delta L}{\delta w_{\text{old}}} = 0.25 \times 0.15 \times 0.10 \times 0.05 \times 0.02$$

derivative value decreases
= very small value.

$$w_{\text{new}} = w_{\text{old}} - \eta \times (\text{small number}),$$

Also a small No.

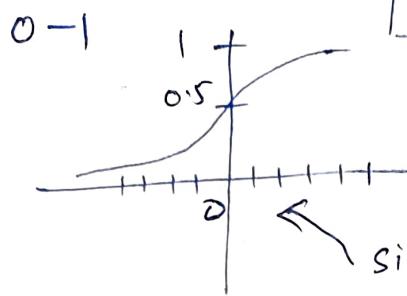
$$\therefore [w_{\text{new}} \approx w_{\text{old}}] \quad \text{very small}$$

the wt will not get updated or very very slowly updated.

This is called vanishing gradient problem
solution is change Activation function

- ① sigmoid activation function
- ② Tanh
- ③ ReLU
- ④ Leaky ReLU
- ⑤ PReLU

① Sigmoid

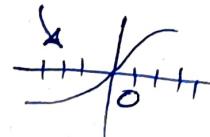


sigmoid is
not zero
centered curve

Adv

- 1) smooth gradient, prevents jumps in o/p values
- 2) since o/p b/w 0 and 1, it normalises o/p of each neuron.
- 3) clear prediction very close to 0 or 1.

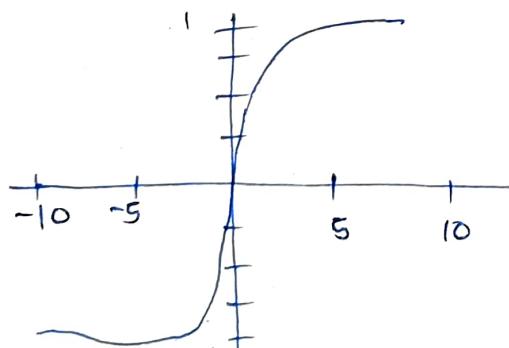
* wt updation is
very easy in
case of zero
centered curves



Disadv

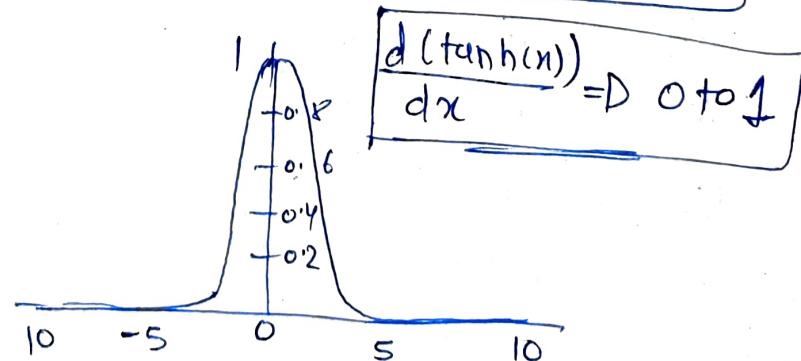
- 1) prone to gradient vanishing / we cannot create deep neural network.
- 2) Function o/p is not zero centred, wt updation
- 3) Sigmoid has exponential operation, issue slower for computers to calculate which is

② tanh function (hyperbolic tangent function)



$$\tanh(x) \Rightarrow 0 \rightarrow 1$$

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



$$\frac{d(\tanh(x))}{dx} = 0 \rightarrow 1$$

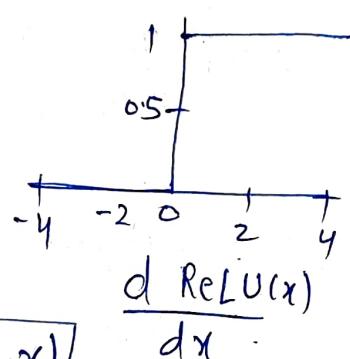
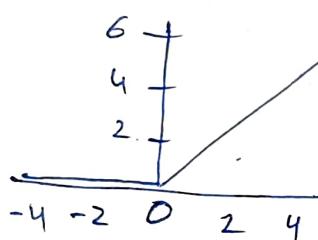
Adv 1) zero centric output.

Disadv 1) still can create vanishing gradient problem.

for binary classification $\rightarrow \tanh \leftarrow$ hidden layer
 $\rightarrow \text{sigmoid} \leftarrow \text{o/p layer}$

(3) ReLU function.

↑ Rectified Linear unit.



→ not zero centric

not differentiable at zero.

* when $\frac{d \text{ReLU}}{dx} = 0$
↑ neuron is dead.
this is a problem.

Adv 1) no vanishing gradient

2) faster calculation

3) when input is +ve, no gradient

saturation problem.

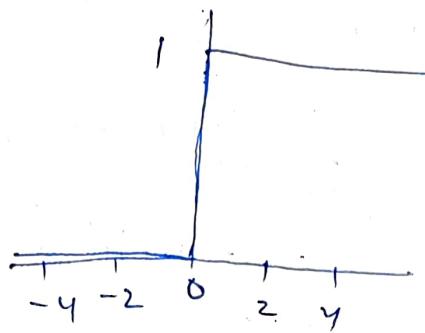
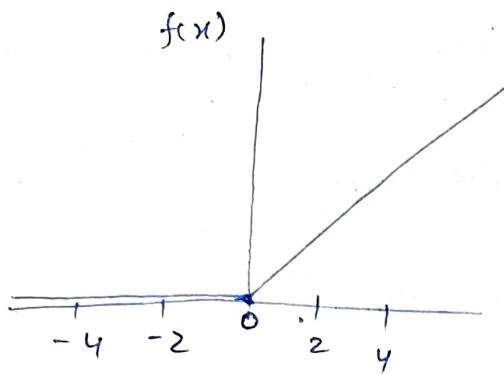
Disadv: 1) not a zero centric function since o/p is either 0 or 1.

2) when input is -ve, ReLU is inactive completely.
⇒ ReLU is dead for input < 0 .

↳ gradient = zero } same problem as sigmoid and tanh.

4) Leaky ReLU function

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



Adv: 1) solves dead neuron problem in ReLU.

5) Parametric ReLU function

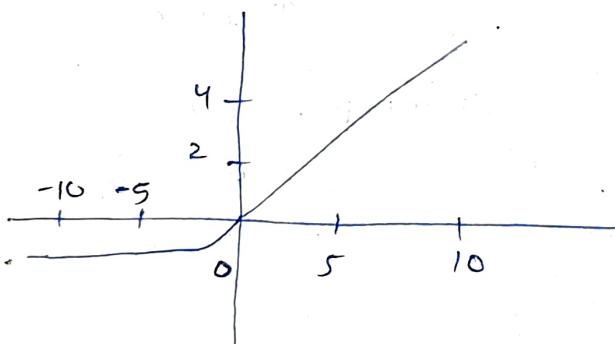
$$a=0 \quad f(x) \rightarrow \text{ReLU}$$

$$a>0 \quad f(x) \rightarrow \text{Leaky ReLU}$$

a is learnable \rightarrow PReLU

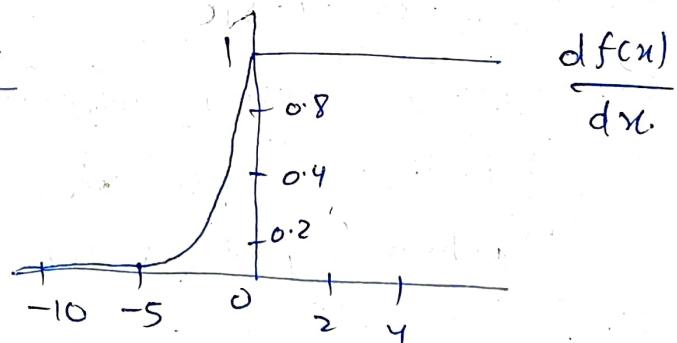
6) ELU (Exponential Linear units) function

$$f(x)$$



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

$$f(u) = \begin{cases} u & u > 0 \\ \alpha(e^u - 1), & \text{otherwise} \end{cases}$$



It solves problems of ReLU.

Adv: 1) No dead ReLU issue

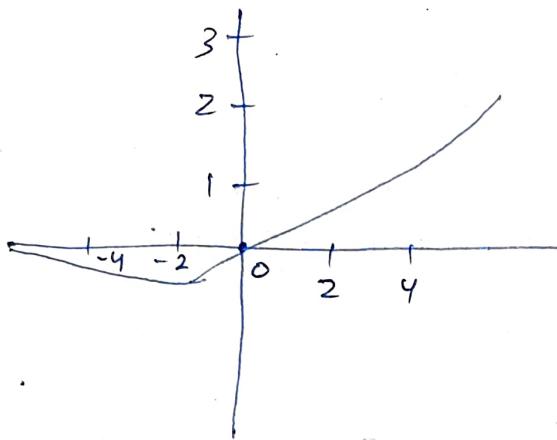
2) zero centred.

3) Better than ReLU (proved)

Disadv 1) computation heavy (due to exponential calcn)

8). SWISH (self gated) function

(4)



$$f(x) = x \times \text{sigmoid}(cx)$$

↳ gating in LSTM

9) max out

$$f(x) = \max(w_1^T x + b_1, w_2^T x + b_2).$$

It is generalised version of ReLU and Leaky ReLU.

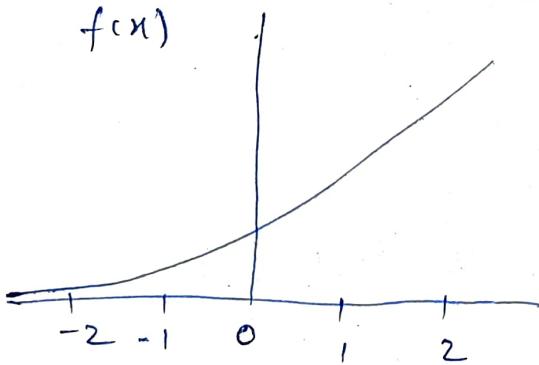
for ReLU $w_1, b_1 = 0$

for Leaky ReLU $w_1, b_1 = 0, w_2 = 0.01$

It is a learnable function.

10) soft plus

$$f(x) = \ln(1+e^x)$$



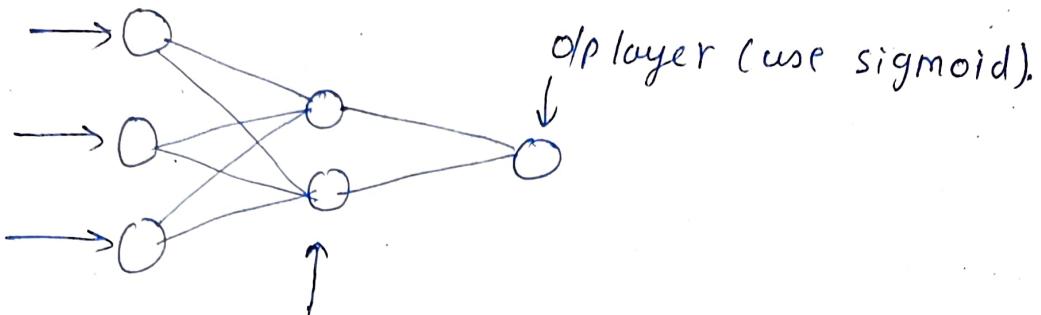
→ differentiable throughout

→ similar to ReLU, but relatively smooth.

→ $0 \rightarrow +\infty$ } Accepted range

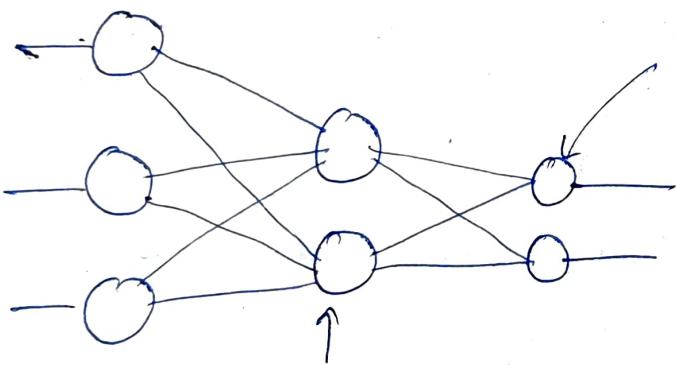
Technique which activation function we should use?

Say Binary classification



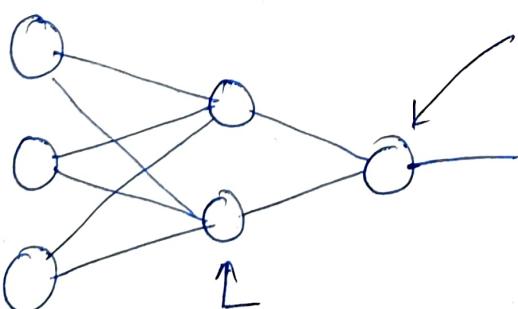
Hidden Layer } if convergence is not
always use ReLU } happening use PReLU,
ELU.

multi class classification



use ReLU } if convergence not
happening use PReLU, ELU.

Regression

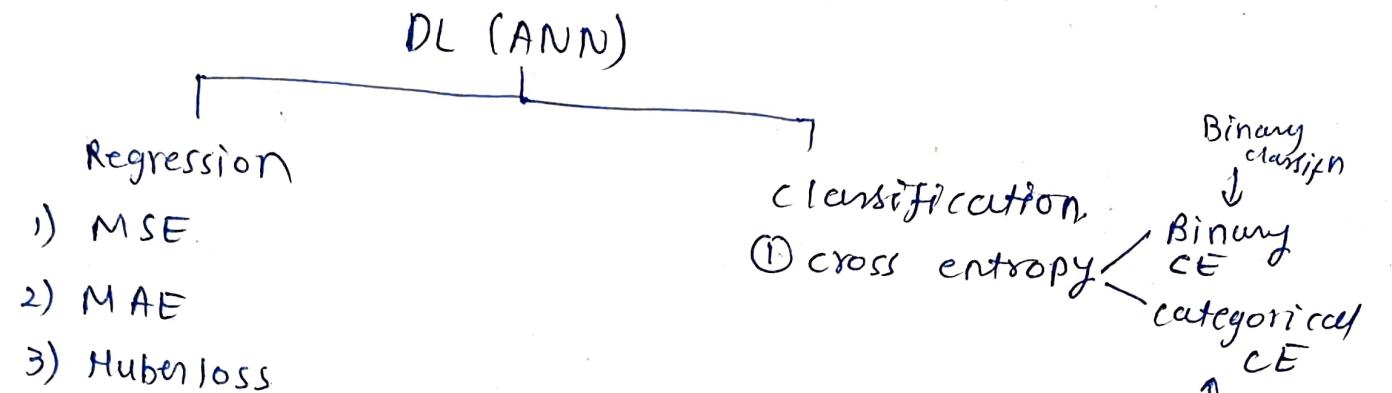


use ReLU or
any variation of
ReLU

use linear activation fn.

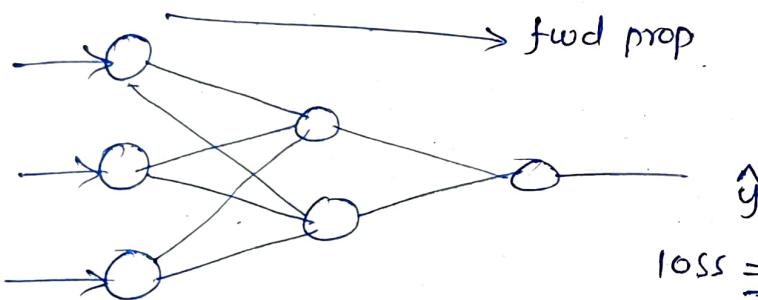
↳ Also there is a
separate loss function.

* Loss functions



* Loss v/s cost function

* 100 recs in dataset.



$$\text{loss} = \frac{1}{2} (y - \hat{y})^2 \leftarrow \text{single value.}$$

∇ cost function for batch input

$$\text{cost fn} = \frac{1}{2} \sum_{i=1}^n (y - \hat{y})^2.$$

(A) Regression

① MEAN squared error

$$\text{LOSS fn} = \frac{1}{2} (y - \hat{y})^2$$

$$\boxed{\text{cost fn} = \frac{1}{2n} \sum_{i=1}^n (y - \hat{y})^2}$$

quadratic eqn.

gradient descent

Adv

1) It is differentiable

2) It has only one local/global minimizer.

3) It converges faster.

Disadv

1) not robust to outliers. (we are squaring the error)

error penalising.

② Mean Absolute Error

$$LF = \frac{1}{2} |y - \hat{y}|$$

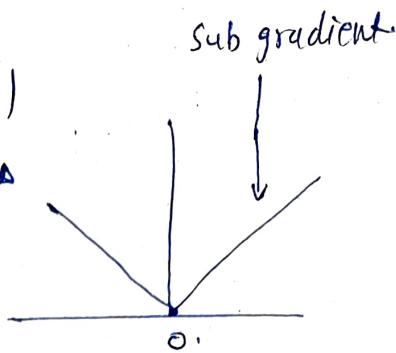
$$CF = \frac{1}{2n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Adv:

- 1) Robust to outliers

disadv

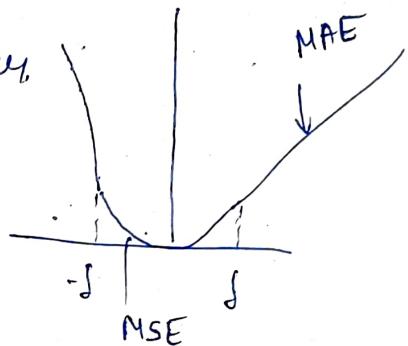
- 1) not diff at $x=0$
- 2) Time consuming due to sub gradient



③ Huber loss. ($MSE + MAE$)

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

Hyper parameters



④ classification

① Binary cross entropy

$$\boxed{\text{loss} = -y \times \log(\hat{y}) - (1-y) \times \log(1-\hat{y})}$$

$$\text{loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

\downarrow log loss. \uparrow log reg.

we use
sigmoid in
last layer for
 \hat{y} calculation

② Categorical cross entropy

⑤

			O/P	OHE		
				j=1	j=2	j=3
i=1	f ₁	f ₂	f ₃	Good	Bad	Neutral
i=1	2	3	4	good	1	0
i=2	5	6	7	bad	0	1
i=3	8	9	10	neutral	0	0

$$L(x_i, y_i) = - \sum_{j=1}^c y_{ij} \times \ln(\hat{y}_{ij})$$

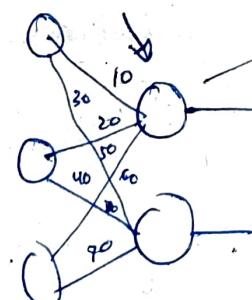
c → no. of categories.
 i = Row
 j = column of o/p

$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

$$y_{ij} = \begin{cases} 1 & \text{if element in class.} \\ 0 & \text{otherwise} \end{cases}$$

$\hat{y}_{ij} \Rightarrow$ use softmax Activation function → O/P layer

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}}$$



$$\frac{e^{10}}{e^{10+20+\dots+70}} = 0.4$$

Prob
+ always one

0.6

Conclusion

↓
HL J/o/p Loss.

multiclass → ReLU, softmax — cat C.E

Binary → ReLU, sigmoid — Binary CE

Regression → ReLU, Linear — Activation MSE, MAE, huber loss.

* Optimizers

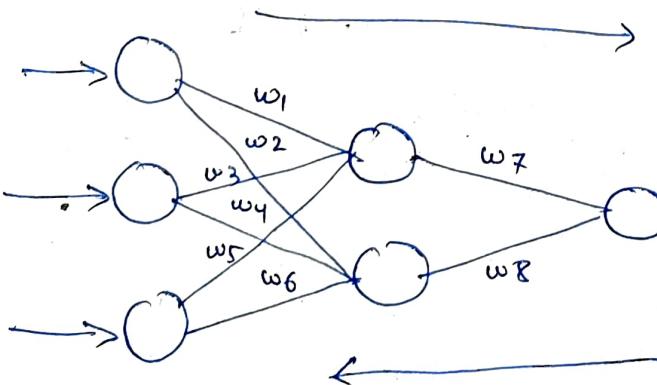
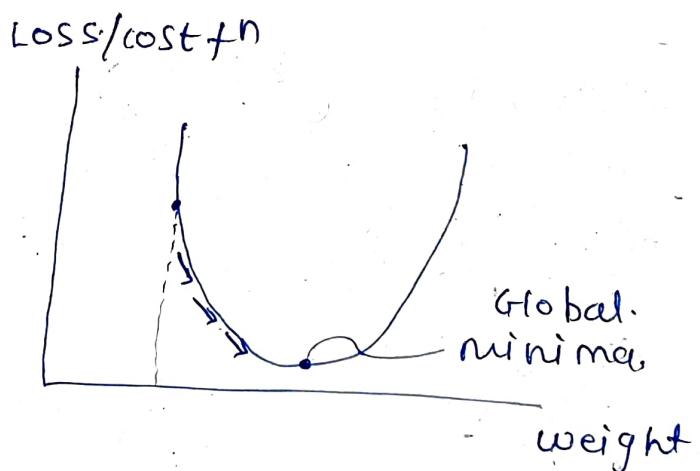
- 1) Gradient descent
- 2) SGD (stochastic GD)
- 3) Mini Batch SGD.
- 4) SGD with momentum
- 5) ADgrad
- 6) RMS-PROP
- 7) Adam optimizer.

① Gradient Descent

wt updation formula

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

η = learning rate



$$\text{Cost funcn} = \frac{1}{2n} \sum_{i=1}^n (y_i - g_i)^2$$

MSE

we use optimizes in back propagation to update wt so that we can minimise loss/error.

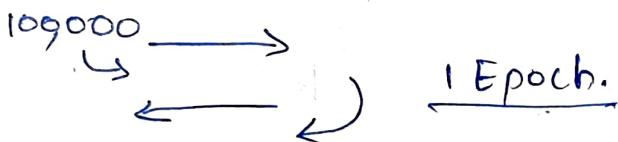
Disadv

- 1) It is resource extensive technique
(more team required)

Epoch

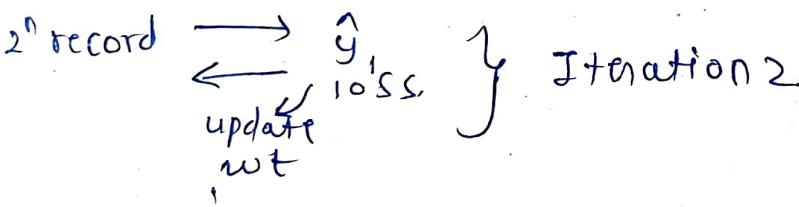
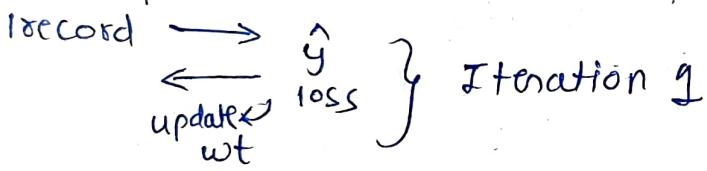
It means training neural network with full dataset for one cycle (fwd and bwd propagation)

Ex! 100,000 data points



② Stochastic Gradient Descent (SGD)

Epoch 1 1 million records in dataset



complete dataset 1 million Iteration.

Say we need 10 Epochs to train model

$$\Rightarrow 10 \text{ Epoch} \times 1 \text{ million Iterations} \Rightarrow \text{Total } 10 \text{ million runs.}$$

Disadv

- 1) converge is very slow.
- 2) Time complexity high.

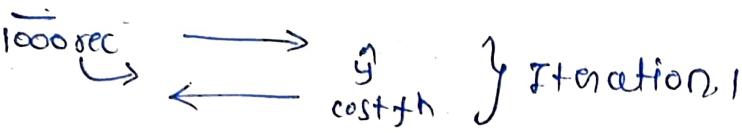
Adv:

- 1) low memory required.

3) Mini - Batch SGD

1 million records \Rightarrow

Epoch 1



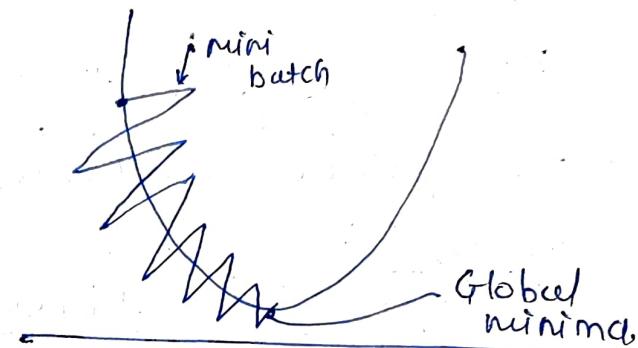
:

1000 Iterations

$$\text{Batch size} = 1000 \text{ recs}$$

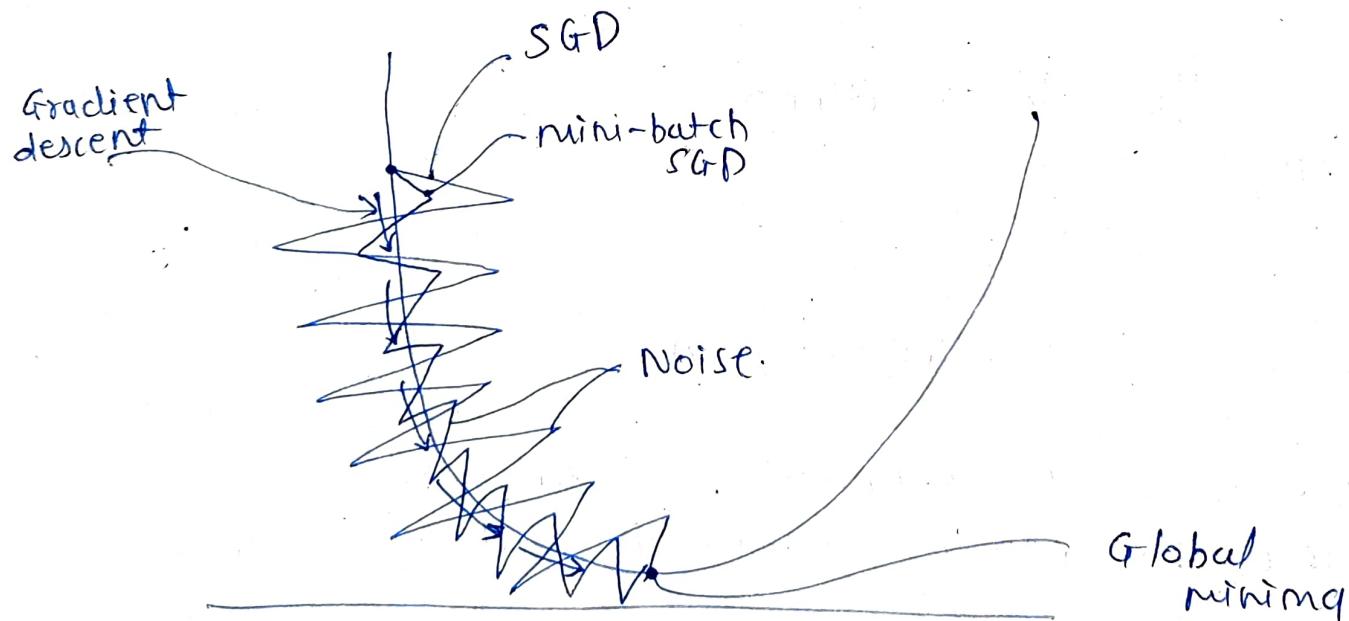
$$\text{No. of Iterations} = \frac{1 \text{ million}}{1000}$$

$$= 1000 \text{ Iterations}$$



Adv

- 1) Less resource intensive
- 2) Better convergence
- 3) Better/improved Time complexity.



noise \Rightarrow SGD > mini-batch SGD.

more noise less noise

To remove noise we use momentum.

(6)

4) SGD with momentum
 ↑ mini batch.

{exponential weighted Average?

momentum is used to eliminate noise using moving avg.
 exponential.

→ This is also used in Time series

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

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

$$w_t = w_{t-1} - \eta \frac{\delta L}{\delta w_{t-1}}$$

$t \in \text{time}$

Exponential weighted average → For smoothening the noise.

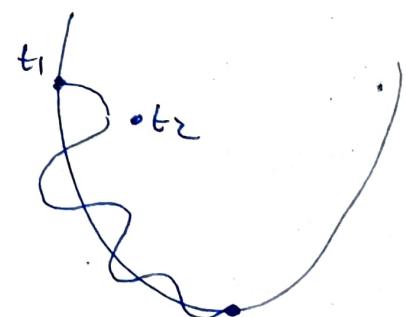
Time	t_1	t_2	t_3	t_n
Value(v.)	a_1	a_2	a_3	a_n

Value at t_1 , $v_{t_1} = a_1$

$\beta \equiv$ hyper parameter
 $\rightarrow 0$ to 1

$$v_{t_2} = \beta \times v_{t_1} + (1-\beta) \times a_2$$

$$v_{t_3} = \beta v_{t_2} + (1-\beta) a_3$$



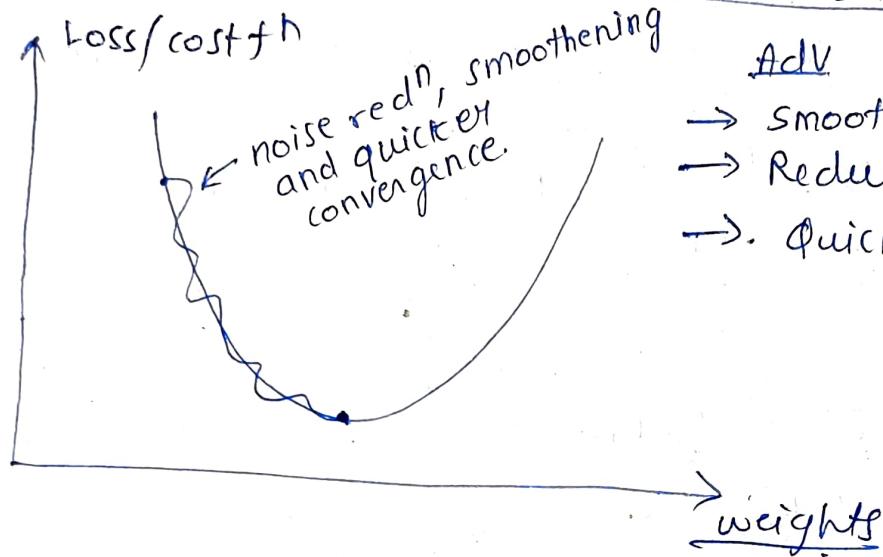
$$v_{t_n} = \beta v_{t_{n-1}} + (1-\beta) a_n$$

→ This is applied in $\frac{\delta L}{\delta w}$

Now

$$w_t = w_{t-1} - \eta V_{dw_t}$$

where $V_{dw_t} = \beta \times V_{dw_{t-1}} + (1 - \beta) \times \frac{\delta L}{\delta w_{t-1}}$



⑤ Adagrad \Rightarrow Adaptive Gradient descent

$$w_t = w_{t-1} - \eta \frac{\delta L}{\delta w_{t-1}}$$

η = fixed in all above optimizers.

Replace with η'

$$w_t = w_{t-1} - \eta' \frac{\delta L}{\delta w_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

Requirement:

Initially large learning rate will make quicker convergence initially, but as we approach global minima learning rate value should reduce so that we do not miss global minimizer.

$$\alpha_t = \sum_{i=1}^t \left(\frac{\delta L}{\delta w_t} \right)^2$$

increases with time ie with every epoch α_t increases.

small no. to avoid zero division error

as the time increase $\sum \frac{\delta L}{\delta w_t}$ will increase

$\therefore \alpha_t$ will increase

$\therefore \eta'$ will decrease as we will reach near to global minima.

\therefore we have made learning rate adaptive.

Disadv

for deep neural network

α_t value will become large

$\therefore \eta'$ will almost become negligible

\therefore our model will be very slow or stop converging to global minima. i.e $w_t \approx w_{t-n}$

6) Ada delta and RMSprop

$$\eta' = \frac{\eta}{\sqrt{sdw + \epsilon}}$$

initially $sdw = 0$

$$sdw_t = \beta sdw_{t-1} + (1-\beta) \left(\frac{\delta L}{\delta w_{t-1}} \right)^2$$

\therefore learning rate will decrease slowly by implementing sdw . i.e exponential weighted average

→ Here smoothening is missing i.e Vdw .

⑦ Adam optimizer

(momentum + RMS prop intuition)
 v_{dw} Adaptive learning
 s_{dw} rate

Initially $v_{dw}, v_{db}, s_{dw}, s_{db} = 0$

$$w_t = w_{t-1} - \eta^1 v_{dw}$$

$$b_t = b_{t-1} - \eta^1 v_{db}$$

$$\eta^1 = \frac{\eta}{\sqrt{s_{dw} + \epsilon}}$$

$$v_{dw_t} = \beta \times v_{dw_{t-1}} + (1-\beta) \frac{\delta L}{\delta w_{t-1}}$$

$$s_{dw_t} = \beta s_{dw_{t-1}} + (1-\beta) \left(\frac{\delta L}{\delta w_{t-1}} \right)^2$$

Adv

- 1) Smoothening the Gradient descent / noise
- 2) Learning rate is adaptive
- 3) Quicker convergence
- 4) Less resource intensive

* for which algorithms feature scaling is req?

Req	Not Req
ANN	Dt
LR	RF
Linear Reg	XGBoost
KNN	AdaBoost
k-means	

Distance based,
optimizers involved } for quicker
convergence