# Deep learning Day-5

## ③ ADAM (Adaptive moment estimation)

Best optimization function for weight and bias
↳ MBSGDA with momentum

Best optimization function for learning rate
↳ Ada Delta (RMS prop)

$$ADAM = MBSGDA \text{ with momentum} + Ada \ Delta.$$

→ ADAM updates weight, bias and learning rate Simultaneously.

✲ MBSGDA with momentum

$$W_{new} = W_{old} - \alpha \ Vdw$$

$$b_{new} = b_{old} - \alpha \ Vdb$$

✲ Ada Delta

$$W_{new} = W_{old} - d_{new} \frac{dc}{dw}$$

$$b_{new} = b_{old} - d_{new} \frac{dc}{db}$$

$$d_{new} = \frac{d_{old}}{\sqrt{Sdwt} + \epsilon}$$

✲ ADAM

$$W_{new} = W_{old} - \frac{d_{old}}{\sqrt{Sdwt} + \epsilon} \times Vdw$$

$$b_{new} = b_{old} - \frac{d_{old}}{\sqrt{sd_{bt}} + \epsilon} \times rdb$$

* Best optimization function to update weight bias and learning rate ⟶ ADAM

1. MBJGDA with momentum
2. Ada Delta (RMS prop) ⎫ 10%
3. ADAM — 90%

* Activation function

It controls the output of any neuron in neural Networks

⟶ Types of Activation function

1. Sigmoid ⟶ Binary classification
2. Tanh ⟶ —ı—
3. Softmax ⟶ multiclass classification (modified version of Sigmoid)
4. Relu (Rectified linear unit) (Base activation function for regression)
5. leaky relu ⎫
6. P-relu ⎬ variant of relu
7. Elu ⎪
8. swish ⎭

(7) Elu

(8) Swish ✓

(1) Sigroid

$$z = \sum_{i=1}^{n} W_i x_i + b$$

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

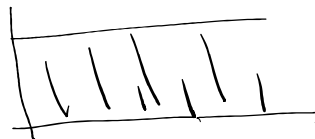↦ Sigmoid always gives probability value of any one class

let's say classes $\langle \begin{smallmatrix} 1 \\ 0 \end{smallmatrix}$

$P(1) = 0.7$

$P(0) = 1 - P(1)$
$= 1 - 0.7$
$= 0.3$

→ Range — 0 to 1

* Advantages of Sigroid (is we use it in activation funn output layer :-

(1) Sigmoid activation funn (range — 0 to 1) gives uniform result.

(2) On the basis of probability value, we can take decision confidently.

classes $\langle \begin{smallmatrix} 1 \\ 0 \end{smallmatrix}$

↦ $P(1) = 0.7$ ✓          threshold — 0.5

$) \longrightarrow P(1) = \underline{0.7}$   threshold — 0.5

$$P(0) = 1 - P(1)$$

$0 \longrightarrow P(0) = 0.3$

* Why we avoid to use sigmoid in our Hidden layer ?.

— First reason

$\quad\quad \hookrightarrow$ Sigmoid is a non zero centric function.

* Zero centricity

$$mean(f(x)) \simeq 0 \longrightarrow f(x) \longrightarrow \text{Zero centric function}$$

$$mean(f(x)) \neq 0 \longrightarrow f(x) \longrightarrow \text{Non Zero centric function}$$

$\longrightarrow$ Zero centric function conversion faster than non zero centric function.

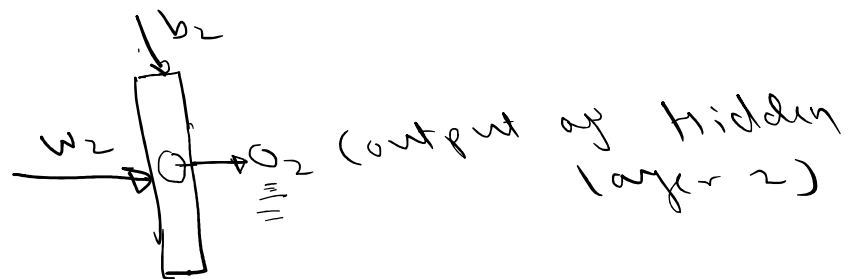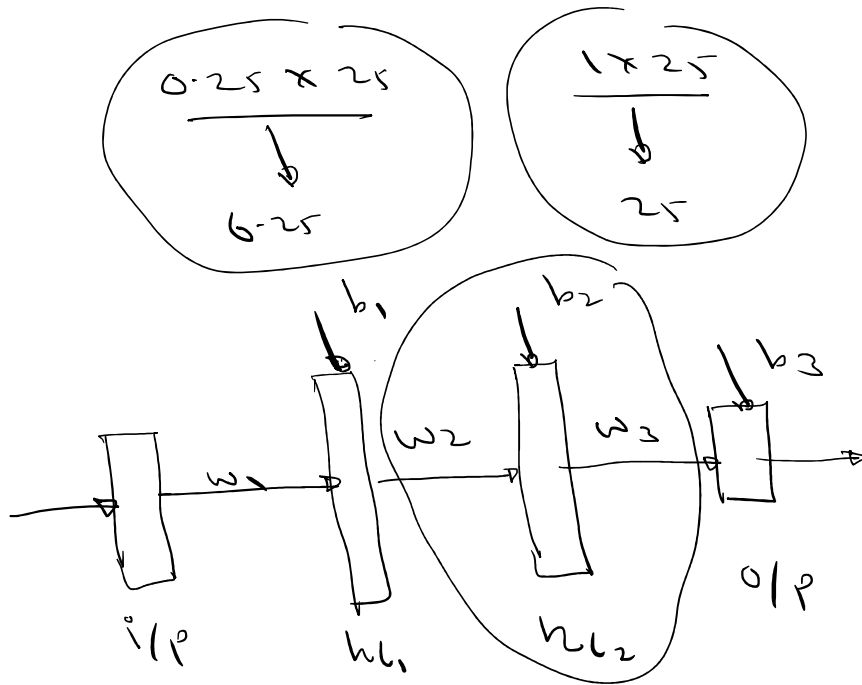Sigmoid $\longrightarrow$ Range — 0 to 1

$\quad\quad\quad\quad\quad\quad\quad$ mean

$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ 0.5

Sigmoid $\longrightarrow$ Non zero centric function

Second Reason

## Second Reason

derivative of sigmoid = Small Number

$$\frac{0.25 \times 25}{6.25}$$

$$\frac{1 \times 25}{25}$$



$b_1$   $b_2$   $b_3$

$w_1$   $w_2$   $w_3$

i/p   $hl_1$   $hl_2$   o/p



$b_2$

$w_2$   $O_2$   $O_2$ (output of hidden layer 2)

$O_2$

activation funn $\left( \sigma(z) = \frac{1}{1+e^{-z}} \right)$

Summation function $\left( z = \sum_{i=1}^{n} w_i \times_i + b \right)$

$w, b$

log

$\frac{\partial C}{\partial w}, \frac{\partial C}{\partial b}$

$$\frac{d}{dz} \sigma(z)$$

$$\frac{\partial C}{\partial w} = \frac{\partial O_2}{\partial z}$$

$$\frac{\partial C}{\partial w} = \frac{\partial O_2}{\partial z} \times \frac{\partial z}{\partial w}$$
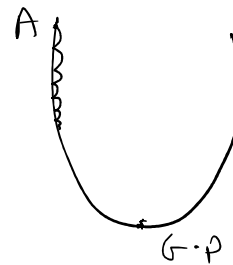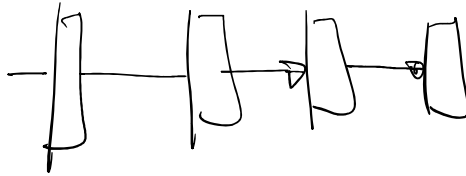
$0.25 \times 25$   $1 \times 25$

$$\frac{dc}{dw} = \frac{d\sigma_2}{dz} \times \frac{dz}{dw}$$

$$0.25 \times 25 \qquad .1 \times 25$$
$$\downarrow \qquad\qquad |$$
$$6.25 \qquad\qquad 25$$

$$O_2 = \sigma(z)$$

$$\downarrow \frac{dc}{dw} = \boxed{\frac{d\sigma(z)}{dz}} \times \boxed{\frac{dz}{dw}}$$

derivative of sigmoid
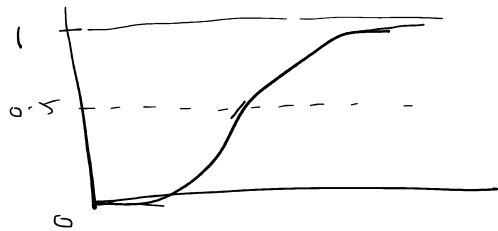
vanishing gradient

A

G.P

**\* Drawback of Sigmoid**

① It is a Non Zero centric function.
hence it will take longer time to
conversion.

② Derivative of sigmoid is very small
number and because of that we
can face issue of vanishing gradient.
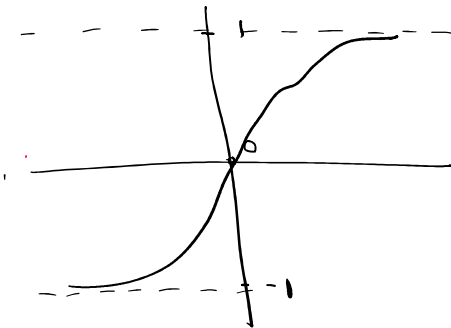
② Tanh

Sigmoid

Range → 0 to 1

mean → 0.5

Tanh

Range → -1 to 1

mean → 0

$\rightarrow$ it is a non zero
centric funn

$\rightarrow$ it is zero centric funn

$\rightarrow$ $\sigma(z) = \dfrac{1}{1+e^{-z}}$
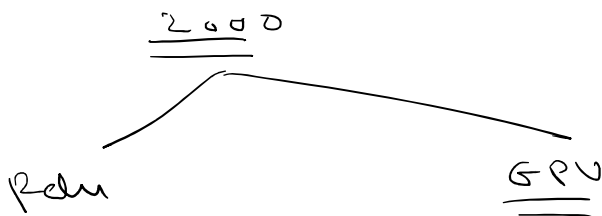
$\rightarrow$ $Tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

* **Advantage of Tanh :-**

(1) It is a zero centric function hence it will conversion faster than sigroid.

* **Disadvantage of Tanh :-**

(1) It is mathematically heavey as compare to Sigroid

(2) Derivative of Tanh is still small number.

$\underline{\underline{2000}}$

Relu                    GPU
                        $\underline{\underline{\phantom{GP}}}$

(3) Relu ( Rectified linear unit )

→ Relu is a non zero centric function.

→ it is works on maximize theory

$$relu(z) = max(0, z)$$
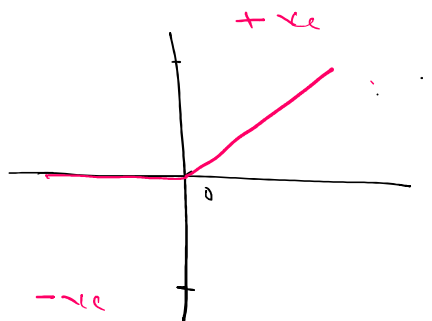
**Case I** → $z$ → +ve ($z = 55$)

$$relu(z) = max(0, 55)$$

$$relu(z) = 55$$

**Case II** → $z = -ve$ ($z = -55$)

$$relu(z) = max(0, -55)$$

$$relu(z) = 0$$

Rectified linear unit



+ve

-ve

\* How Relu solved vanishing gradient problem

**Case I** → $z$ → +ve

$$\frac{d}{dz} relu(z) = \frac{d}{dz} max(0, z)$$

$$\frac{d}{dz} relu(z) = \frac{d}{dz}(z)$$

$0.25 \times 25$   $1 \times 25$

$$\frac{\partial}{\partial z} relu(z) = \frac{\partial}{\partial z}(z)$$

$$\boxed{\frac{\partial L}{\partial w} = \frac{\partial}{\partial z} relu(z) \times \frac{\partial z}{\partial w} = 1}$$

Case II $\longrightarrow$ $z \longrightarrow -ve$

$$\frac{\partial}{\partial z} relu(z) = \frac{\partial}{\partial z} max(0, z)$$

$$= \frac{\partial}{\partial z}(0) \longrightarrow constant$$

$$\frac{\partial}{\partial z} relu(z) = 0$$

$$\frac{\partial L}{\partial w} = \underset{0}{\frac{\partial}{\partial z} relu(z)} \times \frac{\partial z}{\partial w}$$

$$w_{new} = w_{old} - \alpha \boxed{\underset{0}{\frac{\partial L}{\partial w}}}$$

$$w_{new} = w_{old}$$

↳ Dead Neuron
or
Dead Relue

$$z \longrightarrow -ve$$

$$z = \underset{i=1}{\overset{N}{\sum}} w_i x_i + b$$

10 lakh

$x_1 \quad x_2 \quad x_3 \qquad y$

$\uparrow$ Humidity $\longrightarrow$ Temp $\uparrow$

$$\boxed{\text{Humidity} \longrightarrow \text{Temp T}}$$

$$z \xrightarrow{\quad} -ve$$

$5\%$

$95\% \longrightarrow z \longrightarrow +ve$

* **Advantage of Relu**

1. mathematically light weight.

2. with the help of Relu we can avoid vanishing gradient problem.

* **Disadvantage of Relu**

1. It is a Non zero centric function

2. we might face Dead Neuron or Dead Relu issue.

③ **leaky Relu**

$$relu(z) = max(0, z)$$

$$lr(z) = max(0.001z, z)$$

Case I $\longrightarrow z = +ve$

$$\frac{\partial}{\partial z} lr(z) = \frac{\partial}{\partial z} max(0.001z, z)$$

$$= \frac{\partial}{\partial z}(z)$$

$$lr(z) = 1$$

$$\frac{\partial}{\partial z} r(z) = 1$$

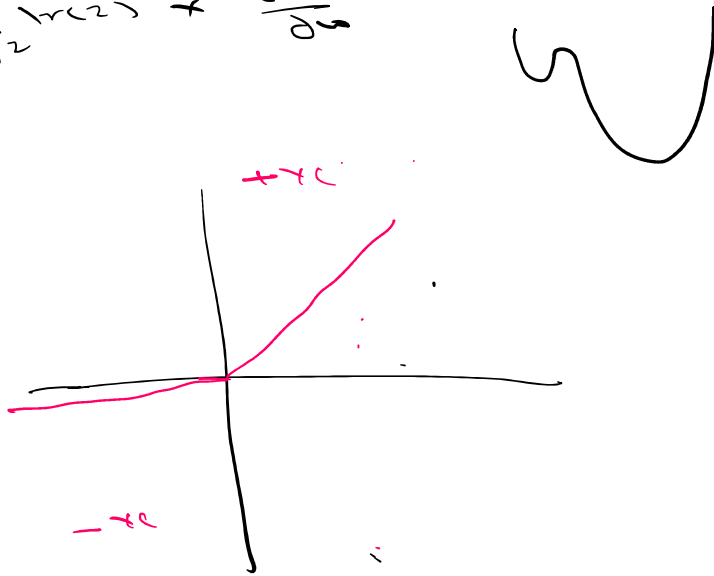Case II $\longrightarrow$ $z = -ve$

$$\frac{\partial}{\partial z} r(z) = \frac{\partial}{\partial z} max(0.001z, z)$$

$$= -0.001 \frac{\partial}{\partial z}(z)$$

$$\frac{\partial}{\partial z} r(z) = -0.001 \longrightarrow small$$

$$\frac{\partial c}{\partial w} = \frac{\partial}{\partial z} r(z) \times \frac{\partial z}{\partial w}$$

$= -ve$



$+ve$

$-ve$

(5) Parametric Relu (P-relu)

$$relu(z) = max(0, z)$$

$$lr(z) = max(0.001z, z)$$

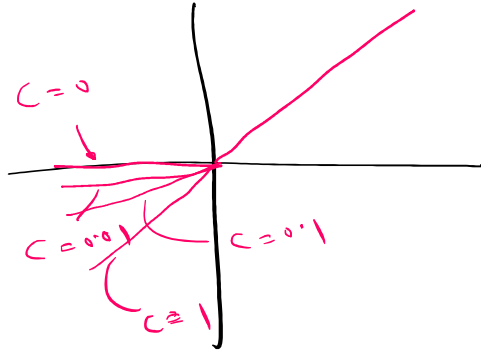$$Prelu(z) = max(cz, z)$$

learnable parameter

$$[0, 0.001, 0.1, 1]$$

$\longrightarrow$ when $c = 0$

$$Prelu(z) = max(0, z) \longrightarrow Relu$$

Prelu (z) - ...

when $c = 0.001$

Prelu (z) = max $(0.001z, z)$ → lr.



$c = 0$

$c = 0.01$    $c = 0.1$

$c = 1$

(6)   Elu (Exponential linear unit)

Exponential function
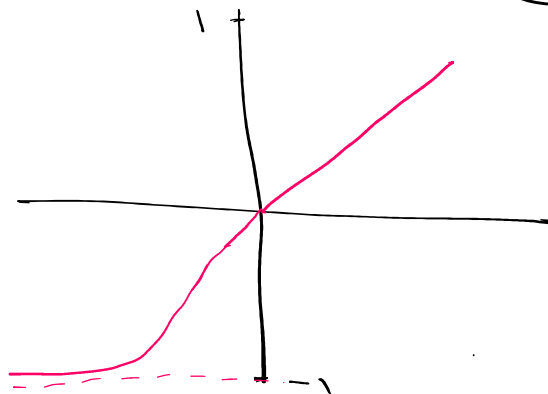
↳ learning curve

↳ Smooth

↳ Elu is a zero centric function.

→ Elu(z) = max $(\alpha(e^z - 1), z)$

learnable parameter.



$$elu(x) = f(x) : \begin{cases} x : & \text{if } x > 0 \\ \alpha(e^x - 1) : & \text{otherwise} \end{cases}$$

$$elu(x) = f(x) : \left\{ \alpha \underline{(e^x-1)} : otherwise \right\}$$
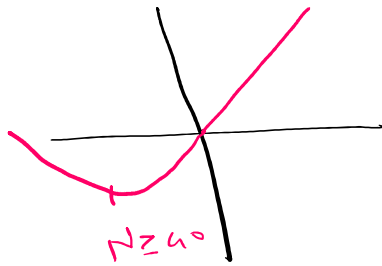
$$elu(z) = max(\alpha(e^z-1), z) \checkmark$$

⑦ Swish

$$Swish(z) = z \times \sigma(z)$$

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

→ we can not use swish in case of shallow Neural Networks

→ we can use swish when $N \geq 40$ where $N \to$ No of hidden layer.



N≥40

⑧ Softmax

→ It is a modified version of sigmoid and it is specially design for multi class classification.

let's say we are dealing with multiclass classification and we have 5 classes.

Class

A

B

C

D

E

Output of Softmax

$P(x \in y = A) = P(A)$

$P(x \in y = B) = P(B)$

$P(x \in y = C) = P(C)$

$P(x \in y = D) = P(D)$

$P(x \in y = E) = P(E)$

$P(A) + P(B) + P(C) + P(D) + P(E) = 1$

$max[\ P(A), P(B), P(C), P(D), P(E)\ ]$

Classification

$$S(x_j) = \frac{e^{x_j}}{\sum_{k=1}^{K} e^{x_k}}$$

$e^A$ — Real Number.

$$S(A) = P(A) = \frac{e^A}{e^A + e^B + e^C + e^D + e^E}$$

$-10$

$-20$

$-40$

$-60$

$-70$

$$\sigma(10) = \frac{1}{1 + e^{-10}} = —$$

A — 10

B — 20

$$S(A) = P(A) = \frac{e^A}{e^A + e^B + e^C + e^D + e^E}$$

(A) — 10
B — 20
C — 40
D — 60
E — 70

$$S(A) = P(A) = \overline{\frac{}{e^A + e^B + e^C + e^D + e^E}}$$

$$= \frac{e^{10}}{e^{10} + e^{20} + e^{40} + e^{60} + e^{70}}$$

$$= \overline{\quad}$$

→ Softmax convert a vector of $k$ real number into a probability distribution of $k$ possible outcome.