# Different Activation Functions

## ① Binary step fn:

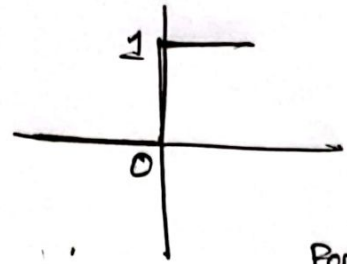In backward Calculation,
For loss function step fn won't be
a good choice.

* multilevel output possible না
* Gradient zero.

$$\frac{\partial L}{\partial \omega} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial \omega}$$
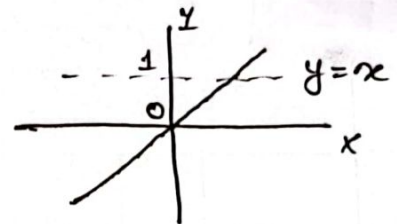
            ↳ zero

* weight update হবে না।

For
$\geqslant 0 : 1 (+)$ positive val
$< 0 : 0 (-)$ neg "

$$\frac{d \mathscr{x} y}{dx} = \frac{d}{dx} 1 = 0$$

$$\frac{dy}{dx} = \frac{d}{dx} 0 = 0.$$

## ② Linear Activation fn:

* Multi level output possible
* Not influential (Gradient র আগে input র অক্ষরর্ত নেই)
* Gradient not zero but ↑

<u>Q:</u> All Layers of the NN collapse into one in linear func

<u>Ans:</u>

Layer-1:   $z^{[1]} = \omega^{[1]} x + b^{[1]}$

       $a^{[1]} = \mathcal{S}(z^{[1]})$   [Sigmoid-ও change হবেনা যা x এবং y র output পাবো]

       $a^{[1]} = z^{[1]}$
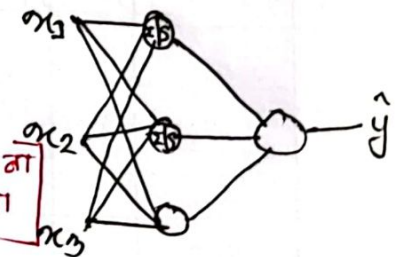
$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$

     $= w^{[2]} z^{[1]} + b^{[2]}$

     $= w^{[2]} (\omega^{[1]} x + b^{[1]}) + b^{[2]}$

     $= \underbrace{w^{[2]} w^{[1]}}_{m} x + \underbrace{w^{[2]} b^{[1]} + b^{[2]}}_{C}$

     $= mx + C$ = linear output

Why do we need nonlinear activation functions?

Ans:

1. NN with linear activation fn result in a linear combination of inputs.

$$z^{[1]} = W^{[1]} x + b^{[1]} = a^{[1]}$$

$$z^{[2]} = \omega^{[2]} a^{[1]} + b^{[2]}$$

$$= \omega^{[2]} \left( W^{[1]} x + b^{[1]} \right) + b^{[2]} = \underbrace{\omega^{[2]} \omega^{[1]}}_{m} x + \underbrace{\omega^{[2]} b^{[1]} + b^{[2]}}_{c}$$

$$= \boxed{m\,x + c}$$

So inputs = outputs.

Collapse in one NN.

2. Linear function are only single grade polynomial $mx+c$ so in gradient decent calculation soon it becomes zero. For single ╲ and sigmoid ╭ so becom. non convex ⤵

3. Multi layer deep NN with nonlinear activation functions can learn hierarchical and abstract representations of features of data.

4. Real world data such as images videos, text often contains nonlinear relationships and high dimentionality. Nonlinear activation fn allow NN to capture and learn these intricate patterns enabling better generalization to unseen data.

## Backpropagation:

$$y = x$$

$$\frac{dy}{dx} = \frac{d}{dx}x = 1$$

$$\frac{d^2}{dx^2} = 0 \quad [\text{Soon becomes zero}]$$

## 3. Sigmoid:

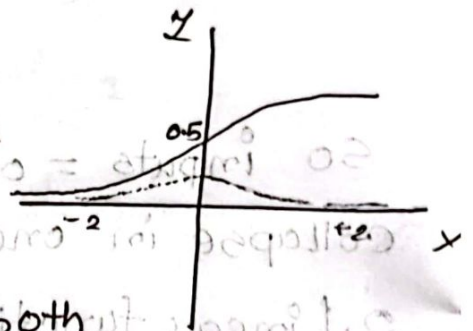→ Multi level output (2 level)

→ Gradient influencial with input

→ 0.5 centered → এটা problem.

করার output positive always (+, -) both
করার center 0 না।

→ vanishing gradient problem [Gradien ক্রমশ zero হবে]

→ Computationally expensive $\quad \sigma = \dfrac{1}{1+e^{-x}}$

## Computing Loss function: (যেখানে total backpropagation calculation দেখানো হয়েছে with sigmoid activation fn)

## To calculate Back Propagation:

$$\frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial \omega_1}$$

$$\frac{\partial L}{\partial a} = \frac{\partial}{\partial a}\left[y \log_2 a + (1-y)\log_2(1-a)\right]$$

$$= \frac{\partial}{\partial a} - \left[y \ln a + (1-y)\ln(1-a)\right]$$

$$= -\left[y\frac{1}{a} + (1-y)\frac{1}{1-a}\cdot\frac{\partial}{\partial a}(1-a)\right]$$

$$= -\left[y/a + \frac{1-y}{1-a}(-1)\right]$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial a}{\partial z} = \frac{\partial}{\partial z}\left(\frac{1}{1+e^{-z}}\right)$$

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

$$= -1(1+e^{-z})^{-2}\frac{\partial}{\partial z}(1+e^{-z})$$

$$= -1(1+e^{-z})^{-2}\; e^{-z}\frac{\partial}{\partial z}(-z)$$

$$= -1(1+e^{-z})^{-2}\; e^{-z}(-1)$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{1+e^{-z}-1}{(1+e^{-z})^2}$$

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

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

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

$$= (1-a)\,a$$

So sigmoid-এ ফিরে আসলে but in different terms.

Now,

$$\frac{\partial z}{\partial \omega_1} = \frac{\partial}{\partial \omega_1}\left(\omega_1 x_1 + \omega_2 x_2 + b\right)$$

$$= x_1$$

$$d\omega_1 = \frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial a}\frac{\partial a}{\partial z}\frac{\partial z}{\partial \omega_1}$$

$$= \frac{-y}{a} + \frac{1-y}{1-a} \times a(1-a) \times x_1$$

$$= \frac{-y + ay + a - ay}{a(1-a)} \times a(1-a) \times x_1$$

$$= (a-y)\, x_1$$

for bias, $b = (a-y)$

$$d\omega_2 = (a-y)x_2$$

$$d\omega_3 = (a-y)x_3$$

$$- - - - - - - -$$

***এখন তুহরা back propagation টাই same, আরকের just activation fn র derivation change হবে এক একটা activation fn র জন্য।

## 4. Tanh Activation fn:

→ Center zero

→ Gradient steep.

→ Vanishing Gradient problem.

→ complex computation $\dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

$$\frac{\partial a}{\partial z} = \frac{\partial}{\partial z}\left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)$$

$$\frac{u}{v} = \frac{vu' - uv'}{v^2}$$

$$= \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2}$$

$$= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$= \left(\frac{e^z + e^{-z}}{e^z + e^{-z}}\right)^2 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)^2$$

$$= 1 - a^2$$

### RELU (Rectified Lere Linear Unit)

$$d\omega_1 = \frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial \omega_1}$$

$$= -\frac{y}{a} + \frac{1-y}{1-a} \times (1 - a^2) \times x_1$$

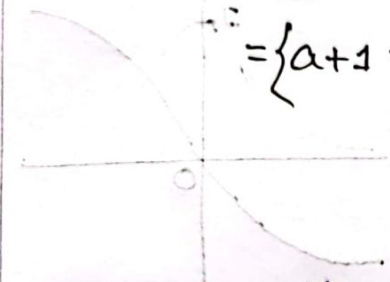$$= \frac{-y + ay + a - ay}{a(1-a)} \times (1 - a^2) \times x_1$$

$$= \frac{a - y}{a(1-a)} \times (1 - a^2) x_1$$

$$= \frac{(a-y)(1+a)}{a(1-a^2)} \times (1 - a^2) x_1$$

$$= \frac{a^2 - ay + a - y}{a} \times x_1$$

$$= \left(a - y + 1 - \tfrac{y}{a}\right) x_1$$

$$\therefore = \left\{a + 1 - y\left(1 + \tfrac{1}{a}\right)\right\} x_1$$

## 4. ReLU: Rectified Linear Unit

$\underline{\text{So As activation fn}:}$

$x \geqslant 0 \quad y = x$

$$\frac{dy}{dx} = \frac{d}{dx} x = 1$$

$x < 0$
$y = 0; \quad \frac{dy}{dx} = 0$

Gradien

$x \geqslant 0 \quad y = x$

$x < 0 \quad y = 0$

**Q:** Proves that ReLU looks like a linear but actions li nonlinear.

linear fn $y = x$

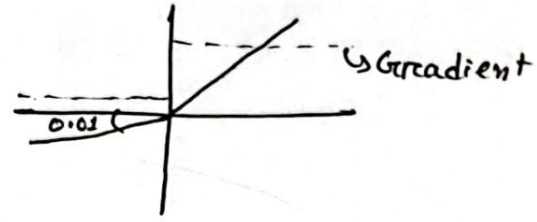Q. if $x < 0$ then $y < 0$

if $x \geqslant 0$ then $y \geqslant 0$

but here for negative value ReLU gives zero value always. So ReLU is not linear.

**\*\*\*** যখন negative value র জন্য Gradient হবে না হয়ে যায় ওটা die ReLU problem. যেহেতু neg =0 তাই কখনো কিছু layer-এ কিছু node/neuron এখন output পারই না। ReLU automatic dropout করু হয়েছে neg=0 যা

→ computationally efficient.

## 5. Leaky ReLU

→ Parameterized ReLU

কার্ভর neg (→) ঢাল manually চেঞ্জ
হয় (0.02 or 0.01)

→ Exponential ReLU ⟋

neg এ slop কিছু curved হয়
তাই।

$x \geqslant 0$  $y = x$  $\dfrac{dy}{dx} = 1$

$x < 0$  $y = 0.01x$  $\dfrac{dy}{dx} = 0.01$

[যাদের derivate হবে Loss বের করার সময় দুইটা কথা
লাগবে।

so the derivation turns into $(1 - \tanh^2)$ where except $\tanh$ no other variable is existed. Hence it is visiable that derivation of $\tanh$ only depends on itself.

Q: "The derivative of the hyperbolic tangent function is more steep than the sigmoid function" - Justify the statement with proper evidance.

Ans: We know, for sigmoid activation function if large value is assigned then the gradient becomes zero same goes for very small value.

Now the sigmoid function, $g(z) = \dfrac{1}{1+e^{-z}}$

And the derivation of sigmoid is,

$$g'(z) = g(z) \cdot (1 - g(z))$$

Let's, if $z = 10$; $g(z) = \dfrac{1}{1+e^{-10}} = 0.999 \approx 1$

So $g'(z) = 1(1-1) = 0$ ✓

Again, if $z = -10$; $g(z) = \dfrac{1}{1+e^{-(-10)}} = 0.0000045 \approx 0$

So, $g'(z) = 0(1 - 0) = 0$ ✓

for $z = 0$; $g(z) = \dfrac{1}{1+e^{-0}} = \dfrac{1}{2}$

So, $g'(z) = \dfrac{1}{2}(1 - \dfrac{1}{2}) = \dfrac{1}{4} = 0.25$ ✓

So it is proved that, for large and low value sigmoid shows vanishing gradient problem and

It's steepness is about 0·25v long if it is drawn graphically.

Now for tanh,

tanh, $g(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-x}}$

The derivation of $tanh(z)$ is ⓪·25

$g'(z) = 1 - \big(g(z)\big)^2$

Let's if $z = 10$; $g(z) = \dfrac{e^{10} - e^{-10}}{e^{10} + e^{-10}} \approx 0.9999 \approx 1$.

So $g'(z) = 1 - (1)^2 = 0$ ✓

Again if $z = -10$; $g(z) = \dfrac{e^{-10} - e^{-(-10)}}{e^{-10} + e^{-(-10)}} \approx -1$

So $g'(z) = 1 - (1)^2 = 0$ ✓

Now, $z = 0$; $g(z) = \dfrac{e^0 - e^{-0}}{e^0 + e^{-0}} = 0$

So $g'(z) = 1 - 0^2 = 1$ ✓

Tanh also suffers from vanishing G. problem but it's steepness much higher than sigmoid about 1·1
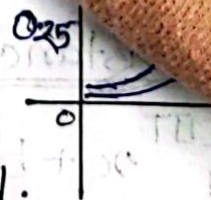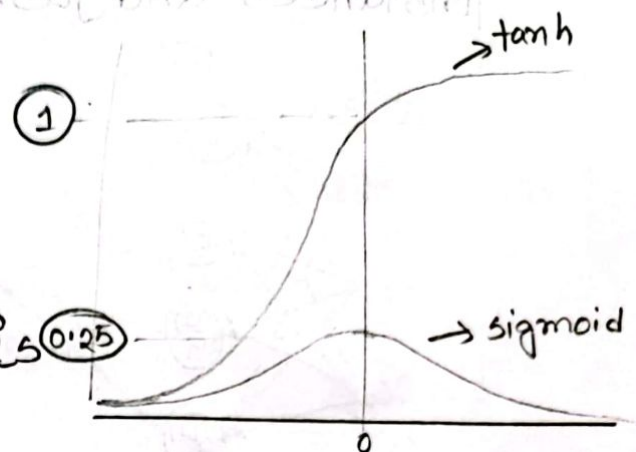
Therefore the statement is justified.



Fig 1: Representation of steepness of tanh & sigmoid