## Dying ReLU Problem:
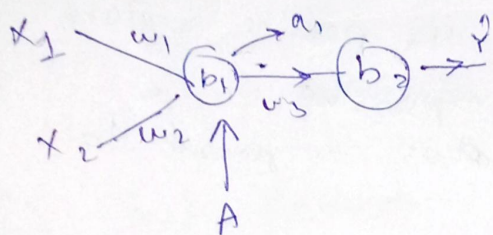
for any input value in the neural network.
some of the. neuron of hidden layer. is
dead (there is no contribution) we con also
say that it is forever dead due to this
(50% neuron are dead) 'the data patterns
are not capture effectively. (not easy easy
easy to find the pattern in the data).

## why dying ReLU problem occur:



$a_1 = max(0, z_1)$

$z_1 = W_1 x_1 + W_2 x_2 + b_1$

if $z_1 < 0$

then $a_1 = 0$

then $\dfrac{\partial a_1}{\partial z_1} = 0$ —— (iii)

due to this $w_1$ and $w_2$ not able to update
because in back propagation

$$w_1 = w_1 - \eta \dfrac{\partial L}{\partial w_1} \quad\text{—— (i)}$$

$$w_2 = w_2 - \eta \dfrac{\partial L}{\partial w_2} \quad\text{—— (ii)}$$

for calculating $\dfrac{\partial L}{\partial w_1}$ and $\dfrac{\partial L}{\partial w_2}$ the (iii) is
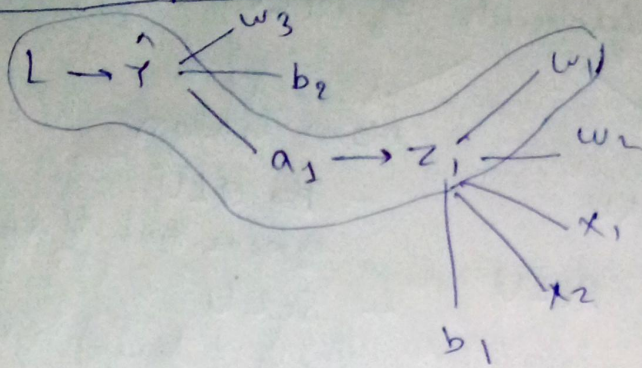used. due to this the equation

$w_1 = w_1$ and $w_2 = w_2$

there is no update hence.
node A is consider as a dead node.

# detail explanation:

$L \rightarrow \hat{Y}$ — $w_3$, $b_2$, $w_1$

$a_1 \rightarrow z_1$ — $w_2$, $x_1$, $x_2$, $b_1$

hence

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{Y}} \cdot \frac{\partial \hat{Y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

similarly,

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{Y}} \cdot \frac{\partial \hat{Y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_2}$$

it is the term (iii)

become -ve?

when $z_1 = w_1 x_1 + w_2 x_2 + b$,

there is two reasons

(I) learning rate is very high

(II) High positive to negative bias

$$b_1 << 0$$

\* dead neuron is not too recoverable.

Solutions to Resolve dying ReLU:

- set low learning rate
- bias set +ve value example 0.01
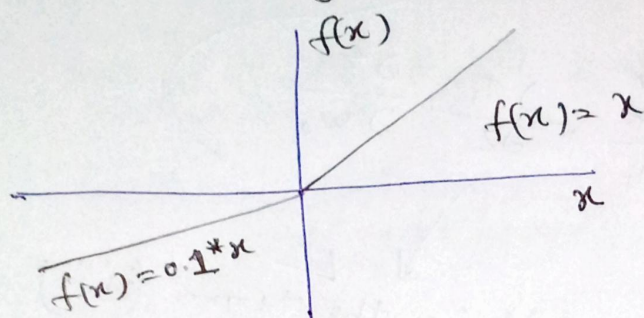- don't use ReLU ReLU, use it variants.

# Variants of ReLU

| Linear | Non-linear |
|--------|-----------|
| — Leaky ReLU | — ELU<br>exponential linear unit |
| — Parametric ReLU | — SeLU<br>Scale linear unit. |

## (i) Leaky ReLU



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

if $z \geqslant 0 \rightarrow z$

if $z < 0 \rightarrow 0.01\,z$

due to this

value of $f'(z)$ for $z \geqslant 0 \rightarrow 1$
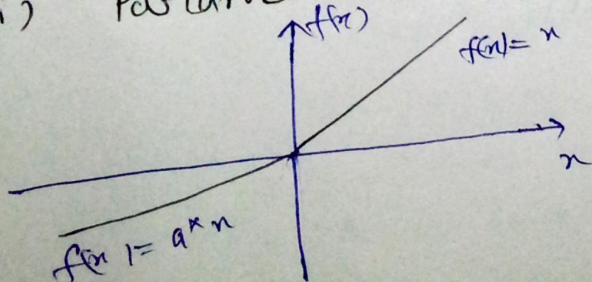
for $z < 0 \rightarrow 0.01$

### advantages:
— Non-saturated (unbounded in both direction)
— easy to compute
— No dying ReLU
— close to zero-centered.

### Disadvantages:
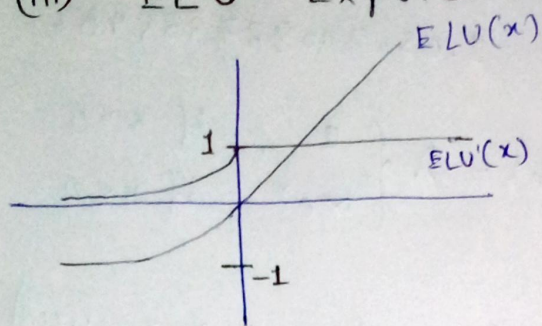— why we use 0.01 value only.

## (ii) Parametric ReLU:



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

Here 'a' is trainable parameter

### advantages:

- all advantage are some as leaky ReLU
- it is flexible and performance better than the leaky ReLU

## (iii) ELU - Exponential Linear Unit:

ELU(x)

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

ELU'(x)

1

-1

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

Here $\alpha$ is constant
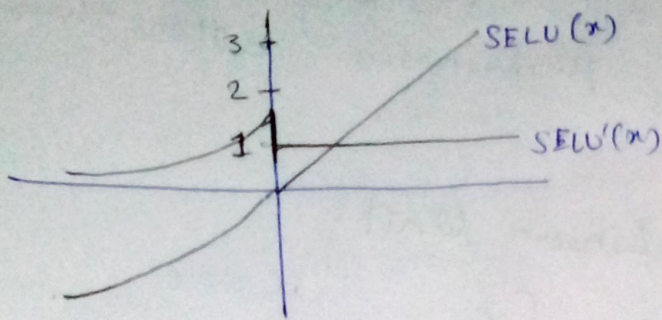$\alpha$ range is 0.1 to 0.3

### advantages:

- performance better than ReLU
- continuous at every point
- always differentiable.
- values are close to zero centered hence convergence is faster.
- generalized results are better (in test data)
- there no dying ReLU problem.

### disadvantages:

- computationally expensive due to $e^x$

(iv) SeLU – Scaled Exponential Linear Unit



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

$\alpha \approx 1.67732632423543$

$\lambda \approx 1.05070098735548$

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

advantages:
- it is self normalizing. (activation is normalized)
  means. mean of actions = 0
         standard deviation = 1
     hence NN converges faster

Disadvantages:
- New in market
- there is less research work on it.