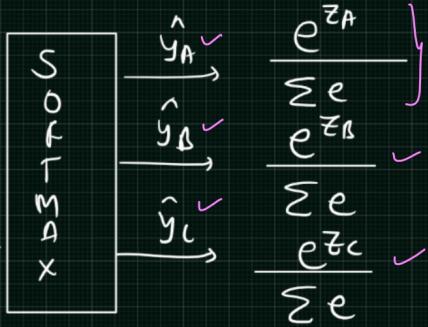
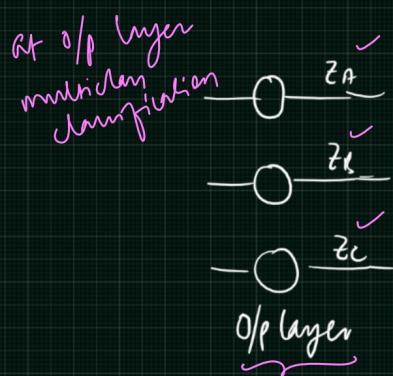


## SOFTMAX ACTIVATION FUNCTION :-

$n.o \text{ of } \text{classes} \geq 2$



$$\sum e = e^{z_A} + e^{z_B} + e^{z_C} \quad \dots \quad (1)$$

$$\hat{y}_A + \hat{y}_B + \hat{y}_C = \frac{e^{z_A} + e^{z_B} + e^{z_C}}{\sum e} = \frac{\sum e}{\sum e} = 1$$

Outcome : Probability distribution

$$\sum y = 1$$

$\hat{y}_A = 0.7 \rightarrow 70\% \text{ of chance is the prediction belongs to class } A$

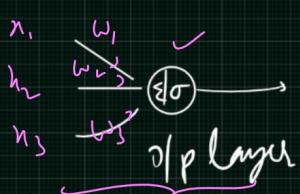
$$\hat{y}_B = 0.1$$

$$\hat{y}_C = 0.0$$

Case Study :-

Case 1: Binary classification { sigmoid }

$$60\% \rightarrow \begin{matrix} 0, 1 \\ 0 \rightarrow 0.1 \\ \downarrow \\ A \\ 1 \rightarrow 0.9 \end{matrix}$$



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

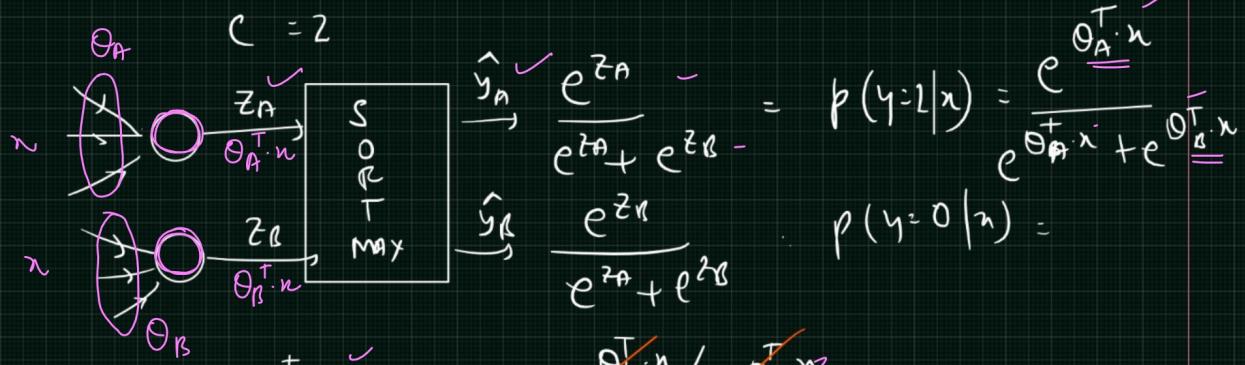
$$\left( \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \right)^T \left[ \begin{matrix} n_1 \\ n_2 \\ n_3 \end{matrix} \right] = \underbrace{\left( \begin{matrix} w_1^T \\ w_2^T \\ w_3^T \end{matrix} \right)}_{(\omega_1, \omega_2, \omega_3)} \cdot \underbrace{\left[ \begin{matrix} n_1 \\ n_2 \\ n_3 \end{matrix} \right]}_x = \Theta^T x$$

$$p(y=1|x) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-\Theta^T x}}$$

$$p(y=0|x) = 1 - p(y=1|x) = 1 - \frac{1}{1 + e^{-\Theta^T x}} = \frac{e^{-\Theta^T x}}{1 + e^{-\Theta^T x}}$$

$$p(y=0|x) = \frac{e^{-\Theta^T x}}{1 + e^{-\Theta^T x}}$$

Case 2: and fn: Saffman Binary Classification



$$p(y=1|x) = \frac{e^{\theta_A^T \cdot n}}{e^{\theta_A^T \cdot n} + e^{\theta_B^T \cdot n}} = \frac{e^{\theta_A^T \cdot n}}{e^{\theta_A^T \cdot n} + e^{\theta_B^T \cdot n}} + \frac{e^{\theta_B^T \cdot n}}{e^{\theta_A^T \cdot n} + e^{\theta_B^T \cdot n}}$$

$$= \frac{1}{1 + e^{\theta_B^T \cdot n - \theta_A^T \cdot n}}$$

$$= \frac{1}{1 + e^{-(\theta_A^T - \theta_B^T) \cdot n}}$$

Assume  $\underbrace{\theta_A^T - \theta_B^T}_{=} = \underbrace{\theta^T}$

$$p(y=1|x) = \frac{1}{1 + e^{-\theta^T \cdot n}}$$

$$p(y=0|x) = \frac{e^{\theta_B^T \cdot n}}{e^{\theta_A^T \cdot n} + e^{\theta_B^T \cdot n}}$$

$$= \frac{e^{\theta_B^T \cdot n - \theta_A^T \cdot n}}{1 + e^{\theta_A^T \cdot n - \theta_B^T \cdot n}}$$

$$= \frac{e^{-(\theta_A^T - \theta_B^T) \cdot n}}{1 + e^{-(\theta_A^T - \theta_B^T) \cdot n}}$$

$$p(y=0|x) = \frac{e^{\theta^T \cdot n}}{1 + e^{\theta^T \cdot n}}$$

Conclusion: For binary classification Softmax is equivalent to sigmoid

Drawback :-

- i) It contains exponential terms ✓
- ↓  
Computationally intensive ✓

Softmax is always used in output layer

2) ELU (Exponential Linear Units)

$$\text{elu}(z, \alpha) = \begin{cases} z & z \geq 0 \\ \alpha(e^z - 1) & z < 0 \end{cases}$$

forward pass

domain  $z \in (-\infty, \infty)$

range  $\text{elu}(z, \alpha) \in (-\alpha, \infty)$

$\alpha$  = scale value | its a slope of -ve section

at  $z=0$ , let see continuity of elu ✓

$$\left\{ \begin{array}{l} \text{LHL} = \lim_{z \rightarrow 0^-} \text{elu}(z, \alpha) = \lim_{z \rightarrow 0^-} \alpha(e^z - 1) \approx \alpha(1^0 - 1) \\ \quad = \alpha(1 - 1) = 0 \quad \checkmark \\ \text{RHL} = \lim_{z \rightarrow 0^+} \text{elu}(z, \alpha) = \lim_{z \rightarrow 0^+} z = 0 \quad \checkmark \end{array} \right.$$

$$\text{elu}(z=0, \alpha) = 0$$

$$\therefore \text{LHL} = \text{RHL} = \text{elu}(z=0, \alpha)$$

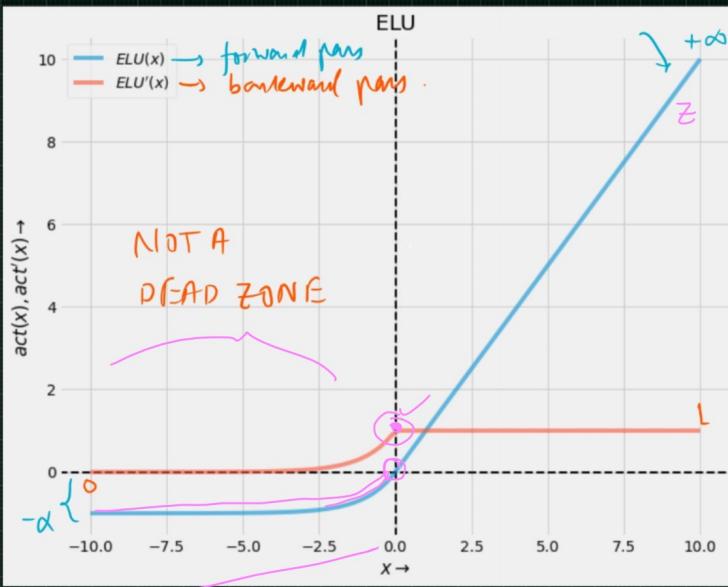
$\Rightarrow$  elu is continuous. ✓

backward pass

domain  $z \in (-\infty, \infty)$

range  $\text{elu}'(z, \alpha) \in (0, 1]$

$$\text{elu}'(z, \alpha) = \begin{cases} 1 & z > 0 \\ \alpha e^z & z \leq 0 \end{cases}$$



**Advantage:**

- i) No dead ReLU issues ✓
- ii) Approx zero centered

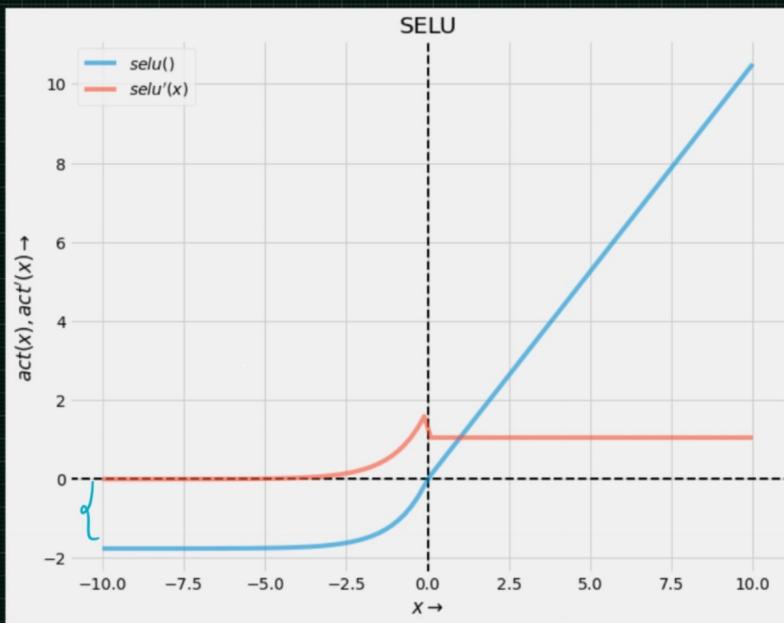
**Disadvantage:**

- i) exponential term
- ii) compute intensive

3) SELU → Scaled ELU

$$\text{selu}(z, \alpha) = \underbrace{\text{scale}}_{\gamma} \times \underbrace{\text{elu}(z, \alpha)}_{\text{elu}}$$

$$z = \underbrace{k}_{\gamma} \times \underbrace{\text{elu}(z, \alpha)}_{\text{elu}}$$



## 4) Swish

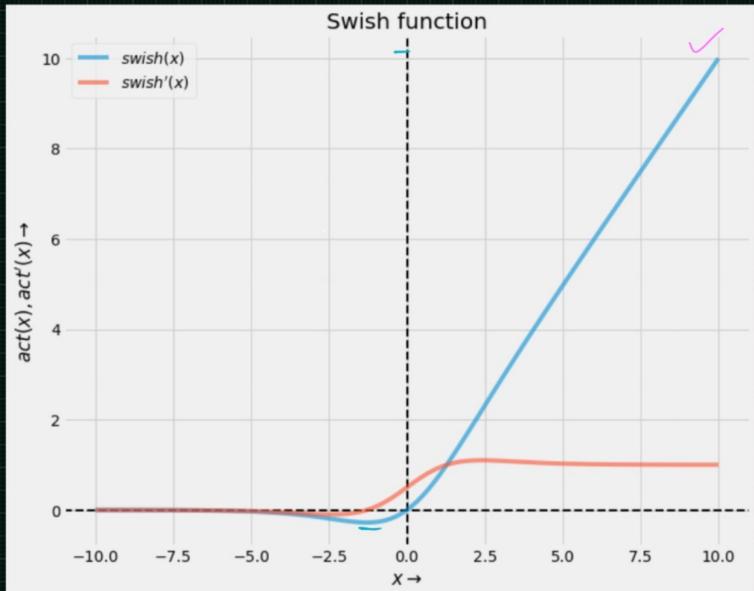
$$\text{swish}(z) = z \cdot \sigma(\beta \cdot z) = \frac{z}{1 + e^{-\beta z}}$$

domain  $z \in (-\infty, \infty)$   
range  $\in (0, \infty)$

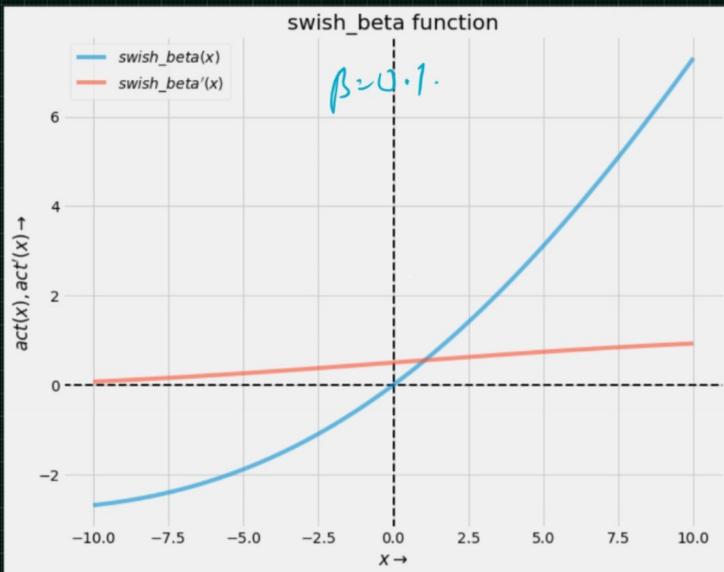
$\beta$  can be a learnable parameter

for  $\beta = 1$ .

$$\text{swish-1} \Rightarrow \text{swish}(z) = z \cdot \sigma(z) = \frac{z}{1 + e^{-z}}$$



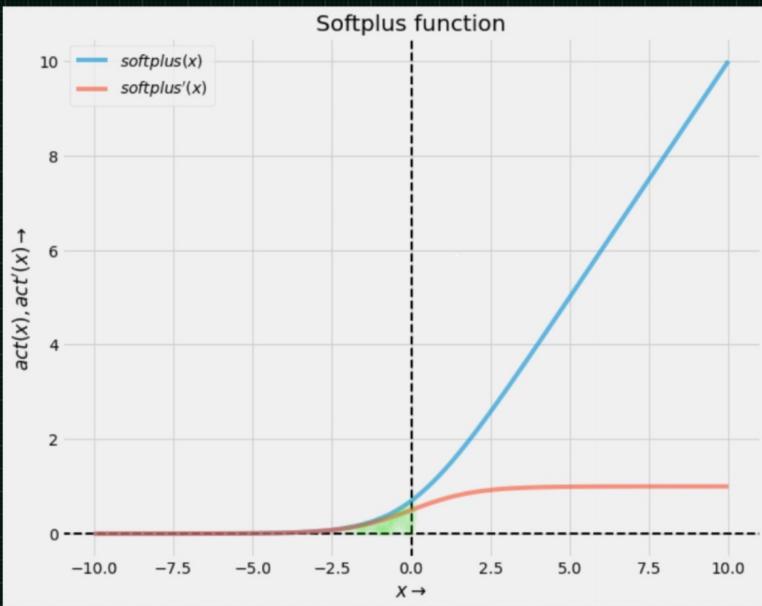
derivative of swish  
domain  $z \in (-\infty, \infty)$   
range  $\in (0, 1)$



- Disadvantages :-
- (1) presence of  $e^x$  term
  - ↓  
computation intensive
- Advantages :-
- (1) vs. moving the origin

## SOFTPLUS

$$\text{softplus}(z) = \log(1 + e^z)$$



**Disadvantage :-**

- i) presence of  $e^z$  &  
 $\log$   $\Rightarrow$  computa  
intensive

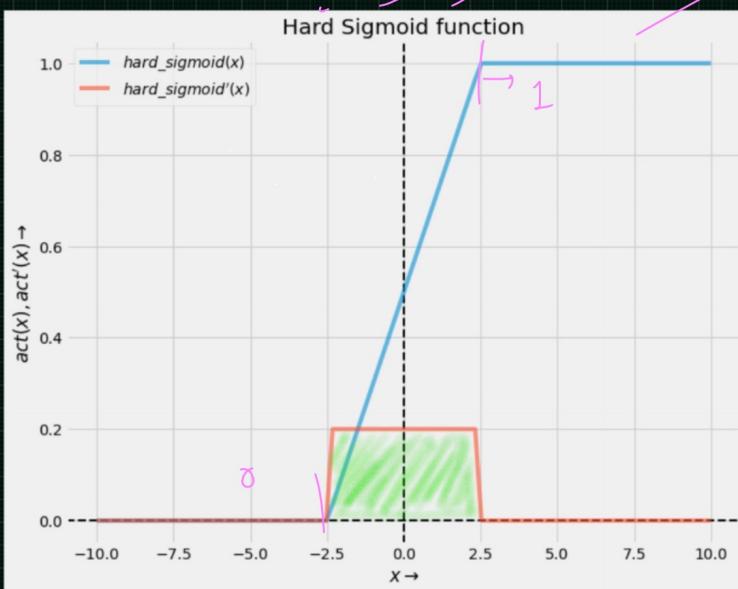
**Advantages :-**

- i) it's differentiable  
everywhere

## Hard Sigmoid

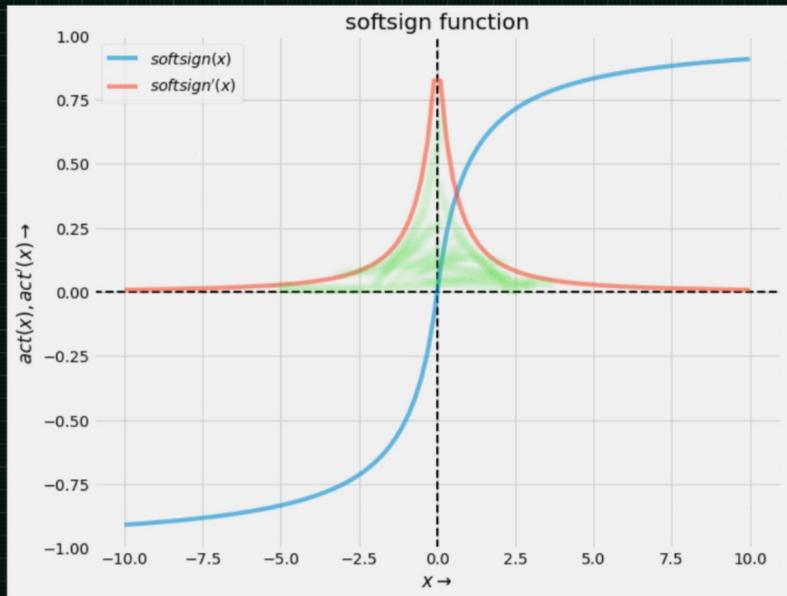
$$\phi(z) = \begin{cases} 0 & z < -2.5 \\ 1 & z > 2.5 \\ 0.2z + 0.5 & -2.5 \leq z \leq 2.5 \end{cases}$$

} if else

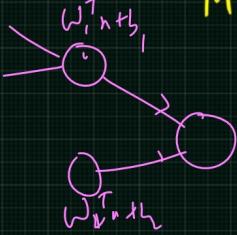


## SOFTSIGN

$$\text{SS}(z) = \frac{z}{|z|+1} = \begin{cases} \frac{z}{z+1} & z > 0 \\ \frac{z}{-z+1} & z < 0 \end{cases}$$

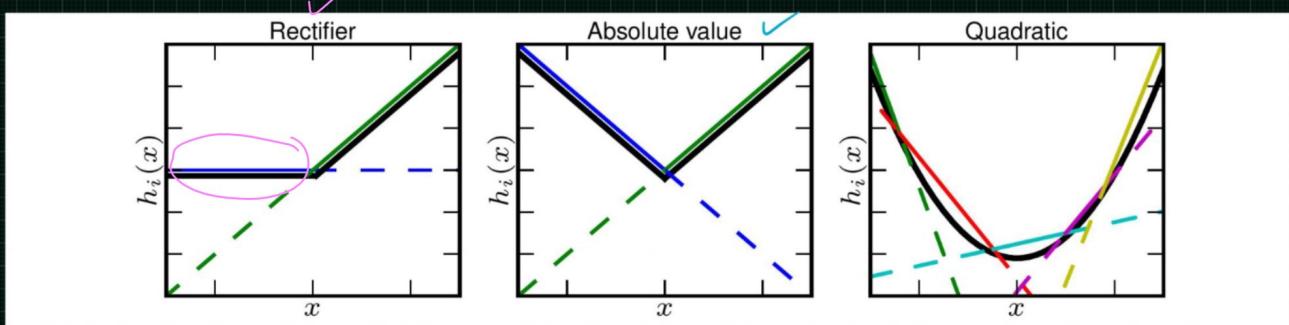


## MAXOUT



Generalization of ReLU & Leaky ReLU

$$f(x) = \max(\underbrace{w_1^T x + b_1}_{(0)}, \underbrace{w_2^T x + b_2}_{(+)})$$



Drawbacks :-

- i) It's computation intensive as it doubles the number of parameters.

→ How to start / begin selection of activation function?

