

Agenda

- 1) Problems in Training NN -
- 2) Weight initialization & activation func "brain"
- 3) Act function.
- 4) Transfer learning
- 5) Batch Normalization

Remaining topic

- (1) Optimizer
- (2) Regularization
- (3) Loss function

1)

i) Vanishing & exploding Gradients

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \text{Gradients } \downarrow \downarrow & & \text{Grad } \uparrow \uparrow \\ \nabla C \downarrow \downarrow & & \nabla C \uparrow \uparrow \end{array}$$
$$W = W - \eta \nabla C \uparrow \uparrow$$

$$\nabla C \downarrow \downarrow$$

$$\nabla C \uparrow \uparrow$$

unstable Gradients \Rightarrow you will never reach to a good solution

\Downarrow
solution will not converge

ii) NN requires lot of data to train

2M parameters \rightarrow MNIST

266610 -

solⁿ: Transfer learning or data augmentation.

iii) for complex problem statement:

Increase size of NN \Rightarrow increase in no of hidden layers

\Rightarrow slow training

solⁿ: Choose a better optimizer & activation f.

iv) Risk of overfitting.

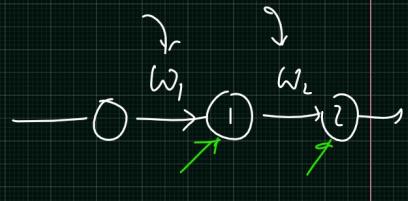


solution : $\{ \text{Regularization} \rightarrow \text{Dropout} \}$

Observation :-

$$\frac{\partial e}{\partial \omega_1} = \frac{\partial e}{\partial a_2} \cdot \underbrace{\left[\frac{\partial a_2}{\partial z_2} \right]}_{\text{gradient}} \cdot \frac{\partial z_2}{\partial \omega_1} - (1)$$

$$\frac{\partial e}{\partial \omega_1} = \frac{\partial e}{\partial a_2} \cdot \underbrace{\left[\frac{\partial a_2}{\partial z_2} \right]}_{\text{gradient}} \cdot \frac{\partial z_2}{\partial a_1} \cdot \underbrace{\left[\frac{\partial a_1}{\partial z_1} \right]}_{\text{gradient}} \cdot \frac{\partial z_1}{\partial \omega_1} - (2)$$



$$z_i = \omega_i a_i + b_i$$

These are product of ratios

1st case :- all the ratio's $\ll 1$

equ. (1) $\frac{\partial e}{\partial \omega_1} = 0.1 \times 0.2 \times 0.3 \approx 0.006$

$$\omega_1 = \omega_1 - \eta \times 0.006$$

$$\omega_1 = \underbrace{\omega_1 - 0.1 \times 0.006}_{\approx 0}$$

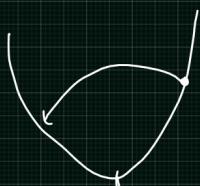
\Rightarrow Vanishing gradient \Rightarrow negligible weight update
 \Rightarrow leads to train lower layer very slow

2nd case

product terms $\gg 1$

$$\frac{\partial e}{\partial \omega_1} = 10 \times 20 \times 30 = 6000$$

$$\omega_1 = \omega_1 - \eta \underbrace{6000}_{\text{gradient}}$$



\Rightarrow Exploding gradient issue
 \Rightarrow Solution will diverge

for Vanishing & Exploding gradient.

{ it depends on choice of activation function.
&
Weight initialization technique

in 2010 by

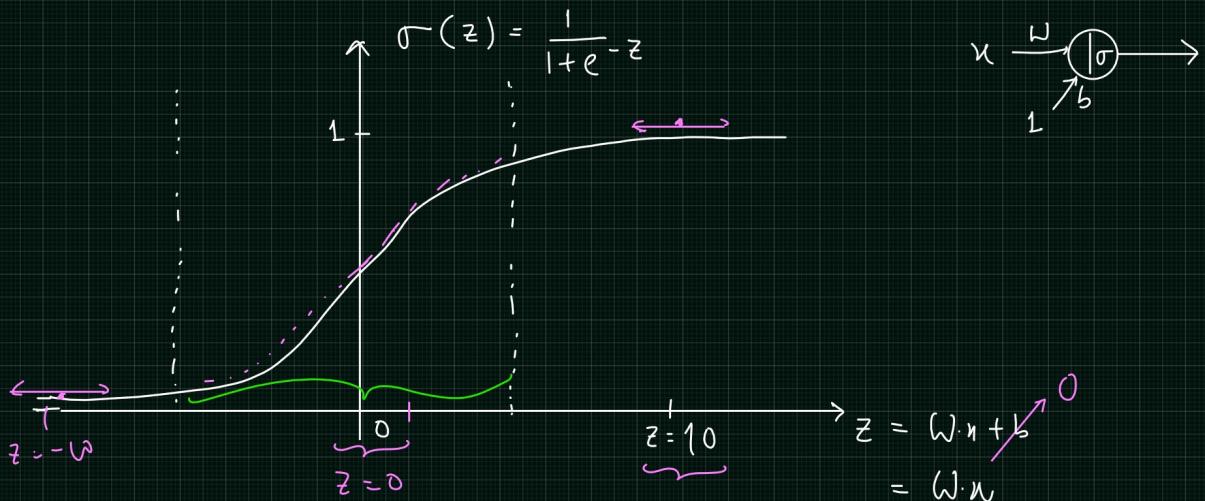
Xavier Glorot & Yoshua Bengio

in their paper

"Understanding the difficulty of training deep feed forward NN"



each layer learns at different speed



Case 1 initialized weight

$$(w = 10 \quad \text{assume } n = 1)$$

$$z = w \cdot n = 10 \times 1 = 10$$

$$\begin{aligned}\sigma(z=10) &= \frac{1}{1+e^{-10}} \\ &= \frac{1}{1+0} = 1\end{aligned}$$

$$e \approx 2.73 \approx 3$$

$$3^{-10} = \frac{1}{3^{10}} \approx 0$$

at $z=10$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z) (1 - \sigma(z))$$

$$\sigma'(z=10) = \sigma(z=10) \left\{ 1 - \sigma(z=10) \right\}$$

$$= 1 \left\{ \underbrace{1 - 1} \right\}$$

$$= 0$$

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial a_2} \cdot \cancel{\frac{\partial a_2}{\partial z_2}}^0 \cdot \frac{\partial z_2}{\partial w_2} = 0$$

\Rightarrow zero weight update

Case 2: Weight init-

$$w = -10 \quad \text{assume } n = 1$$

$$z = -10 \times 1 = -10$$

$$\sigma(z=-10) = \frac{1}{1+e^{-(-10)}} = \frac{1}{1+e^{10}} \approx 0$$

$$\sigma'(z=-10) = \underbrace{\sigma(z=-10)}_{=0} \left\{ 1 - \sigma(z=-10) \right\}$$

$$\frac{\partial e}{\partial w_2} = \frac{\partial e}{\partial a_2} \cdot \cancel{\frac{\partial a_2}{\partial z_2}}^0 \cdot \frac{\partial z_2}{\partial w_2} \approx 0$$

\Rightarrow zero weight update

Case 3:

$$w = 0$$

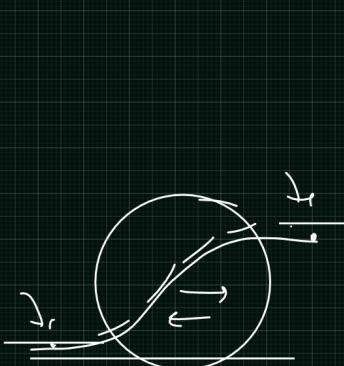
$$z = 0$$

$$\sigma(z=0) = \frac{1}{1+e^0} = \frac{1}{2} = 0.5$$

$$\sigma'(z=0) = \sigma(z=0) \left\{ 1 - \sigma(z=0) \right\}$$

$$= 0.5 \left(1 - \underbrace{0.5}_{0.25} \right)$$

$$= 0.25 \neq 0 \checkmark$$



\xrightarrow{fp}

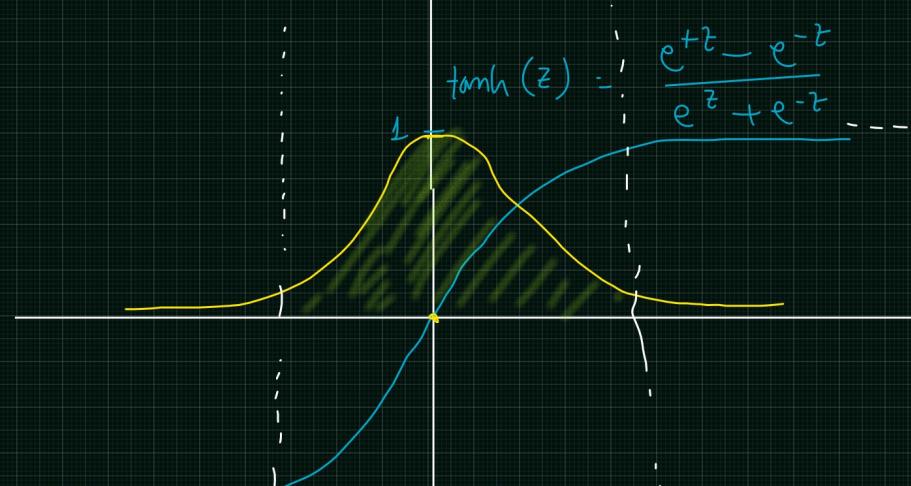
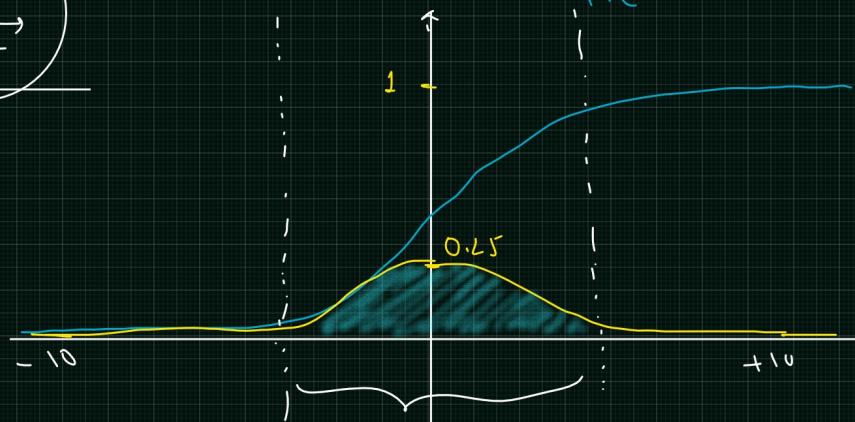
\xleftarrow{bp}

forward

backward prop.

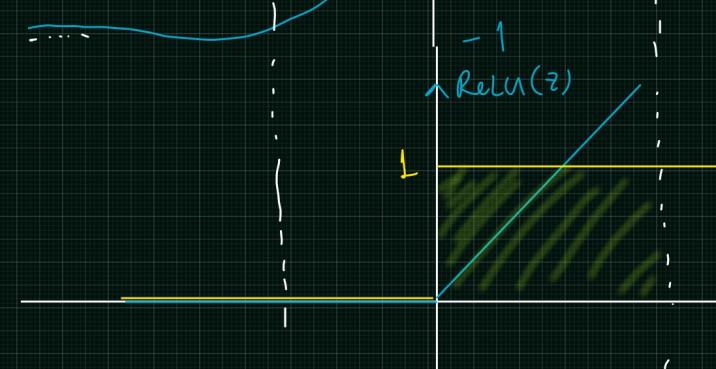
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\sigma'(z) = \underbrace{\sigma(z)}_{0.5} \underbrace{\{1-\sigma(z)\}}_{0.5 \times 0.5} = 0.25$$



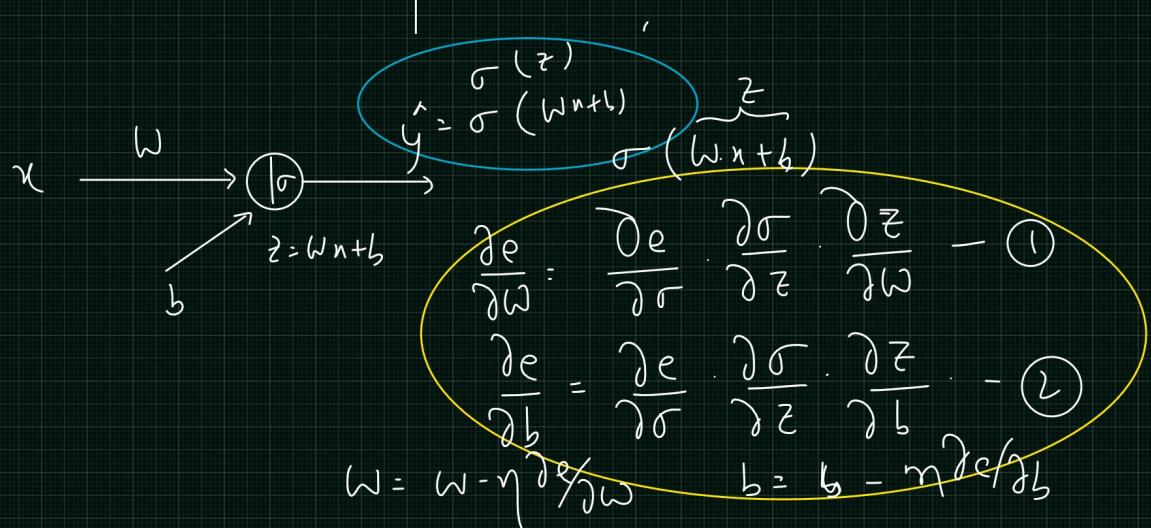
$$\tanh'(z) = 1 - \tanh^2(z)$$

$$\tanh'(0) = 1 - 0 = 1$$



$$ReLU = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$

$$ReLU' = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$$



$$y = \sigma(w \cdot x + b)$$

$$z = w \cdot x + b$$

$$\frac{\partial e}{\partial w} = \frac{\partial e}{\partial r} \cdot \frac{\partial r}{\partial z} \cdot \frac{\partial z}{\partial w} - \textcircled{1}$$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial r} \cdot \frac{\partial r}{\partial z} \cdot \frac{\partial z}{\partial b} - \textcircled{2}$$

$$w = w - \eta \frac{\partial e}{\partial w}$$

$$b = b - \eta \frac{\partial e}{\partial b}$$

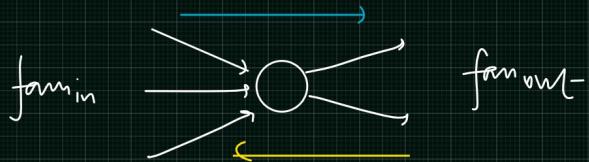
As per the paper,

info ≠ 0

for proper flow of information in forward as well as backward direction :-

$\text{fan}_{\text{in}} \rightarrow$ no. of incoming edges to a layer

$\text{fan}_{\text{out}} \rightarrow$ " " outgoing edges from v → "

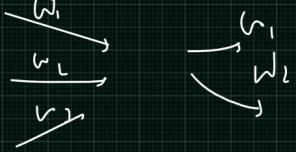


Condition to prevent vanishing & exploding gradient

$$\text{fan}_{\text{in}} = \text{fan}_{\text{out}} \quad \checkmark$$

alternative proposal

$$\text{if } \sigma_{\text{in}}^2 = \sigma_{\text{out}}^2$$



\Rightarrow This can help you to maintain forward & backward info flow
↓
to tackle the vanishing & exploding gradient issues

Calculation :-

$$\text{fan}_{\text{avg}} = \frac{\text{fan}_{\text{in}} + \text{fan}_{\text{out}}}{2}$$

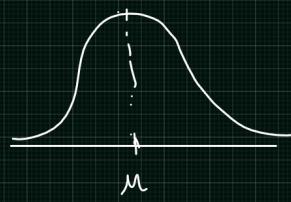
For sigmoid ,

Global initialization → Normal distribution with $\mu=0$ & Variance = $\frac{1}{\text{fan}_{\text{avg}}}$
or
Xavier initialization → Uniform distribution between $-r$ & r where $r = \sqrt{\frac{3}{\text{fan}_{\text{avg}}}}$

Default case of weight initialization in keras

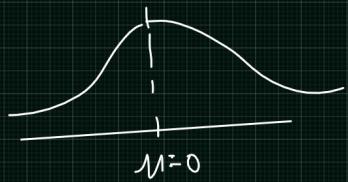
Normal distribution

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



std Normal distribution $\mu=0, \sigma=1$

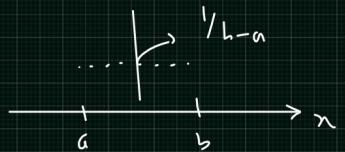
$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



Uniform distribution

a, b

$$f(x) = \frac{1}{b-a}$$



(initialization)

Xavier Glorot
& Yoshua Bengio

Glorot-

None, tanh, sigmoid,
softmax

σ^2 (normal)

γ_{fanin}

Kaiming He

He

ReLU & its variants

$\gamma/fanin$

Yann LeCun

LeCun

SELU
activation fn

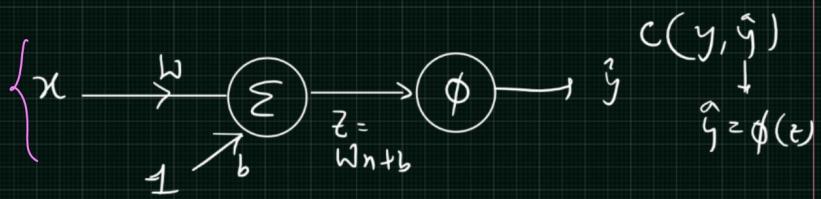
$\gamma/fanin$

tf.kern. Layers.Dense (unit, activation = "relu",
kernel_initializer = "he_normal")

$$\omega_{100} \rightarrow$$

$$\sigma^2 = \gamma/fanin$$

Activation Functions



forward pass ✓

$$\boxed{\hat{y} = \underline{\phi}(z)}$$

where $z = \underbrace{w_n}_\uparrow x + \underbrace{b}_\uparrow$

domain $\rightarrow (-\infty, \infty)$

range $\rightarrow (0, 1)$ sigmoid

Backward pass

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial a} \cdot \frac{\partial a}{\partial z}, \frac{\partial z}{\partial w}$$

$$a' = \underline{\phi'(z)} = \frac{\partial \phi}{\partial z}$$

domain $\rightarrow (-\infty, \infty) \leftarrow$

range $\rightarrow (0, 0.5) \leftarrow$

sigmoid

$$\sigma'(z) = \sigma(z) \{1 - \sigma(z)\}$$

$$\text{at } z=0 \quad \sigma'(z=0) = 0.5$$

$$z = \underset{-\infty}{+\infty} \quad \sigma'(z) = 0$$

for forward pass :-

use $\phi(z)$ $\phi \rightarrow \text{antif}$

consider its
domain & range

for backward pass

use $\underline{\phi'(z)}$

Consider its domain &
range.

1) Sigmoid activation f^n

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

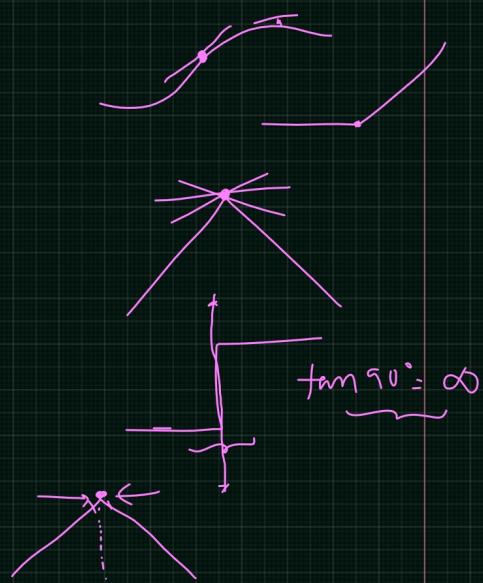
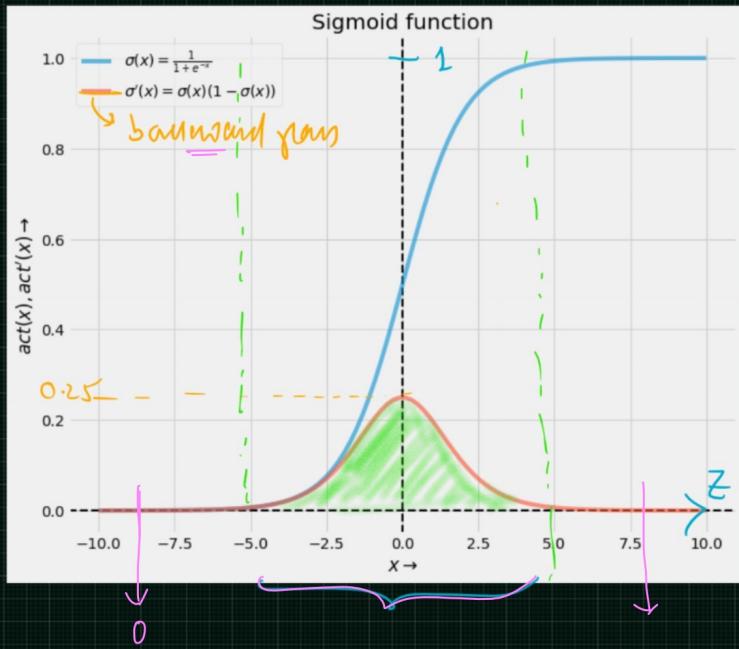
range $\sigma(z) \in (0, 1)$

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

$$\sigma'(z) = \underbrace{\sigma(z)}_{\text{range}} \left\{ 1 - \underbrace{\sigma(z)}_{\text{range}} \right\}$$

$\sigma'(z) \in (0, 0.25)$

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



Advantages:-

- 1) It has smooth gradient \Rightarrow prevents "jumps" in output.
continuous $f^n \Rightarrow$ derivable everywhere.
- 2) Output value or Range $\sigma(z)$ is between 0 & 1.
 \Rightarrow normalizing the o/p of each neuron
 \Rightarrow helps us to reduce the solution space.

Disadvantages:-

- 1) Prone to gradient vanishing

$$\sigma'(z) \in (0, 0.25) \quad 0.25 < 1$$

- 2) Power operation \Rightarrow exponential terms
makes calculation time consuming

latency ↑

$$\left\{ \frac{1}{1+e^{-z}} \right\}$$

2. Hyperbolic Tangent activation function

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

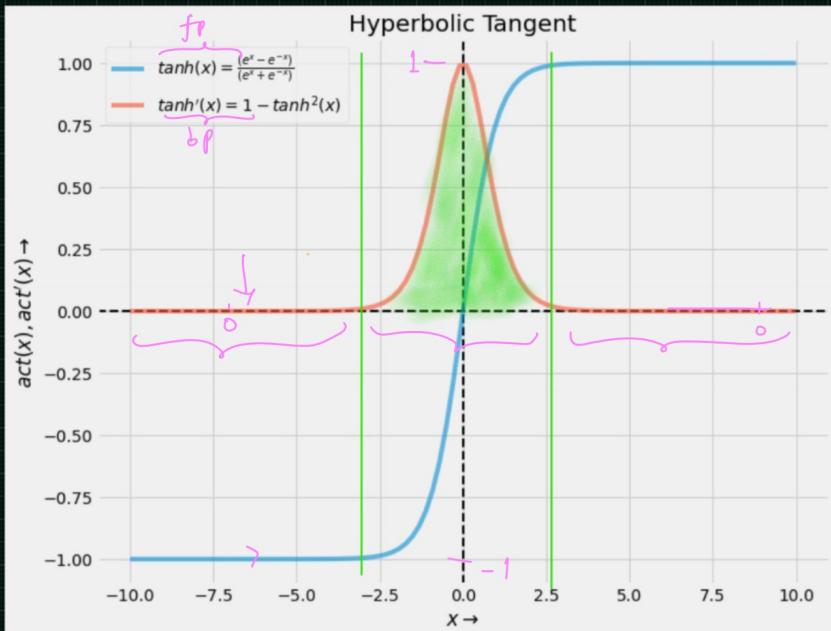
range $\rightarrow \tanh(z) \in (-1, 1)$

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

$$\tanh'(z) = 1 - \tanh^2(z)$$

range $\rightarrow \tanh'(z) \in (0, 1)$

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



Advantages:-

1) Smooth gradients.

2) Symmetric to origin \Rightarrow zero mean

3) Solves vanishing gradient issue to some extent provided there is a proper wt. initialization and also because its range for $\tanh'(x)$ is between $(0, 1)$

4) Suitable function for hidden layers

Disadvantage.

1) It also has saturation region

2) It introduces latency because of exponential terms.

3) Rectified Linear unit (ReLU) activation fn.

$$\text{ReLU}(z) = \max(z, 0)$$

— or —

$$\text{ReLU}(z) = \begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$$

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

range $\text{ReLU}(z) \in [0, \infty)$

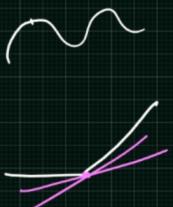
$$\text{ReLU}'(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

{ derivative is not defined
at $z = 0$

domain $z \in (-\infty, \infty)$
except $z = 0$

range

$\text{ReLU}'(z) \in 0 \text{ and } 1$.



Advantages:-

- 1) If input is +ve \Rightarrow There is no gradient division problem
- 2) Calculation is much faster {forward or backward} as compared to tanh & Sigmoid.

Disadvantage:-

- 1) For -ve inputs ReLU is completely inactive

↓
Dying ReLU problem.

- 2) Not a zero centric function.

- 3) At zero its derivative is not defined \Rightarrow jump

4) Leaky ReLU function

$$\frac{d\alpha z}{dz} = \alpha \underset{-\alpha}{\cancel{1}}$$

$$\text{Leaky ReLU}(z) = \max(\tilde{z}, \alpha z)$$

— or —

$$\text{ReLU}(z) = \begin{cases} z & z \geq 0 \\ \underline{\alpha z} & z < 0 \end{cases}$$

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

Range of LeakyReLU(z) $\in \underline{(-\infty, \infty)}$

$$\alpha = \underline{0.01} \quad 0.01z$$

\downarrow

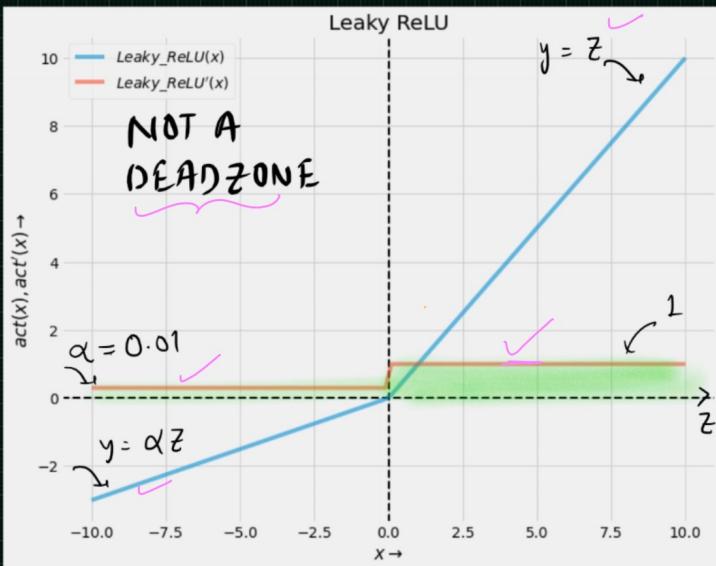
$$\text{LeakyReLU}'(z) = \begin{cases} 1 & z > 0 \\ \underline{\alpha} & z < 0 \end{cases}$$

derivative is not defined
at $z = 0$

domain $z \in (-\infty, \infty)$
except $z = 0$

range $0.01 \text{ or } 1$

$$\text{LeakyReLU}'(z) \in \underline{\alpha} \text{ or } \underline{1}.$$



Advantages:-

for all input

- 1) If input is +ve \Rightarrow There is no gradient admission problem
- 2) Calculation is much faster {forward or backward} as compared to tanh & sigmoid.

Disadvantage:-

1) for -ve inputs ReLU is completely inactive

Dying ReLU problem.

- 2) Not a zero centric function.
- 3) At zero its derivative is not defined \Rightarrow jump

5) PReLU (Parametric ReLU)

$$\text{PReLU}(z) = \begin{cases} z & z \geq 0 \\ \alpha z & z < 0 \end{cases}$$

$\rightarrow \alpha \rightarrow \text{learnable parameter}$

or
learnable

$$\omega = \omega - \eta \frac{\partial e}{\partial \omega}$$

$$\rightarrow \alpha = \alpha - \eta \frac{\partial e}{\partial \alpha}$$

if $\alpha = 0 \Rightarrow \text{ReLU}$

if $\alpha > 0 \Rightarrow \text{Leaky ReLU}$

if α is learnable parameter $\Rightarrow \text{PReLU}$

Advantage:

1) α is now intelligently updated } \checkmark
learnt during the training

↓
No manual fine-tuning required } \checkmark
to set α

↓
 α will be adaptable as per } \checkmark
the data

Disadvantage:

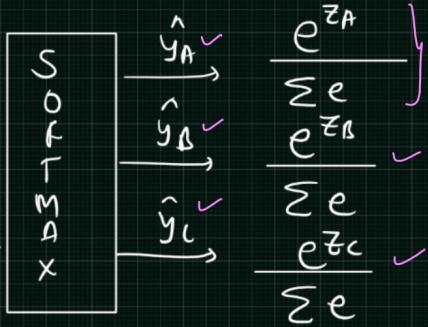
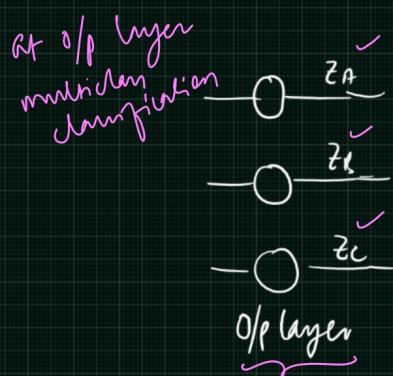
1) You have entire parameter $\underline{\alpha}$

↓
to train

↓
learning time \uparrow

SOFTMAX ACTIVATION FUNCTION :-

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



$$\hat{y}_A = \frac{e^{z_A}}{\sum e}$$

$$\hat{y}_B = \frac{e^{z_B}}{\sum e}$$

$$\hat{y}_C = \frac{e^{z_C}}{\sum e}$$

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

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

Outcome : Probability distribution

$$\sum y = 1$$

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

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

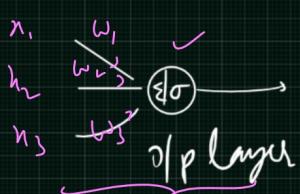
$$\hat{y}_C = 0.0$$

Case Study :-

Case 1: Binary classification { sigmoid }

60%

0.6 →
0, 1
0 → 0.5
↓
A
B



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

$$\begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix}^T \begin{pmatrix} n_1 \\ n_2 \\ n_3 \end{pmatrix}$$

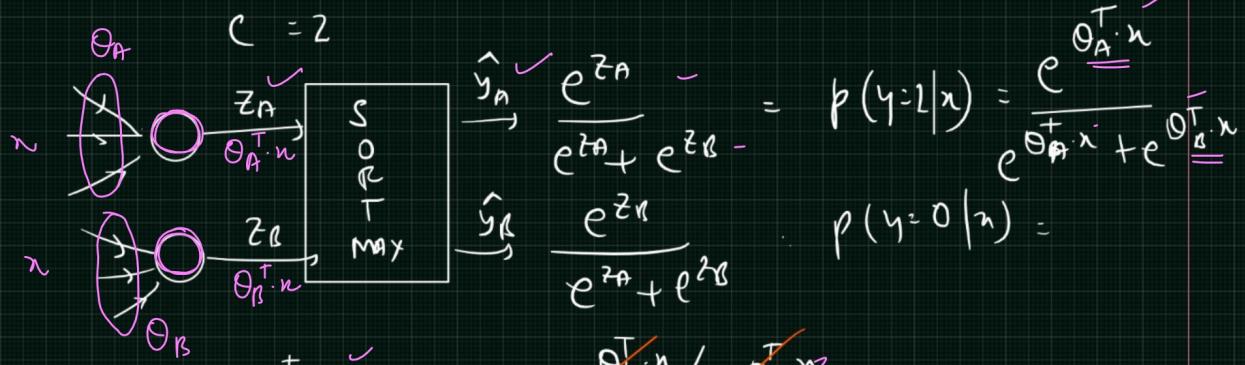
$$z = w_1 n_1 + w_2 n_2 + w_3 n_3 = \underbrace{w^T x}_{(\omega_1, \omega_2, \omega_3)} = \underbrace{\Theta^T n}_{\left(\begin{matrix} n_1 \\ n_2 \\ n_3 \end{matrix} \right)}$$

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

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

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

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 $\theta_A^T - \theta_B^T = \underline{\theta^T}$

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

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

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

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

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

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

α = scale value | its a slope of -ve section

at $z=0$, let see continuity of elu ✓

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

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

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

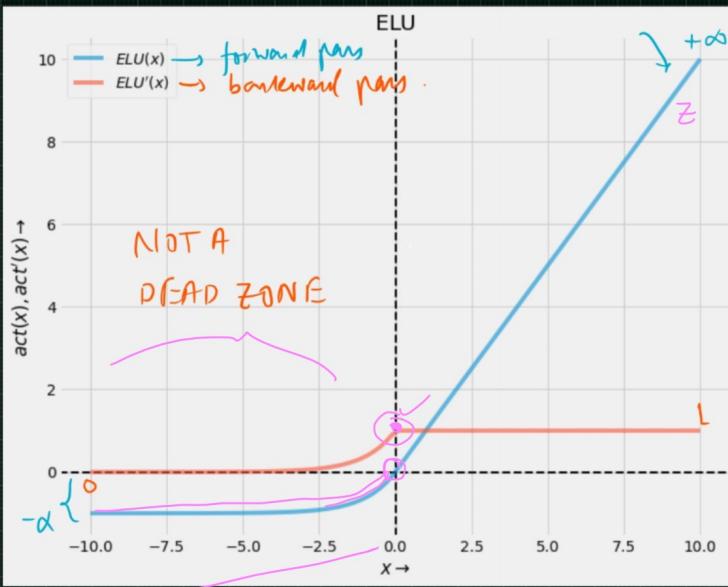
\Rightarrow elu is continuous. ✓

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

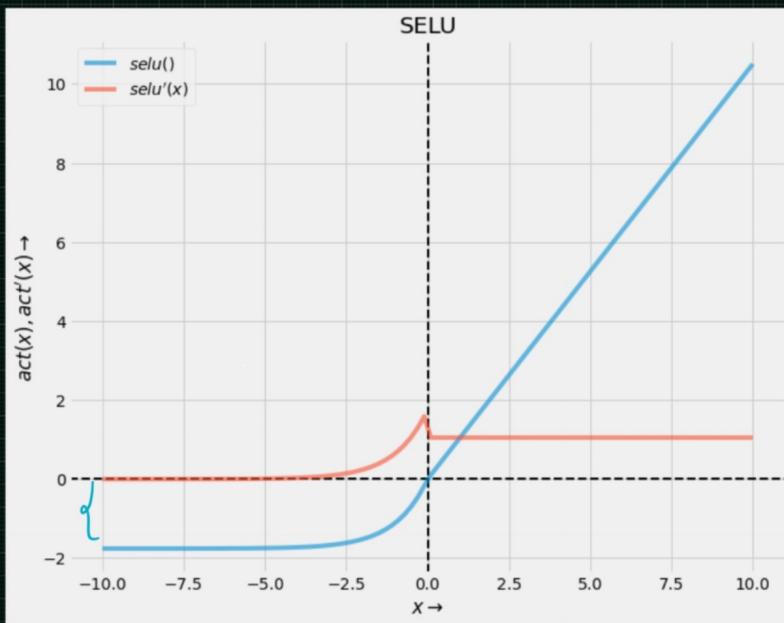
Disadvantage:

- 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

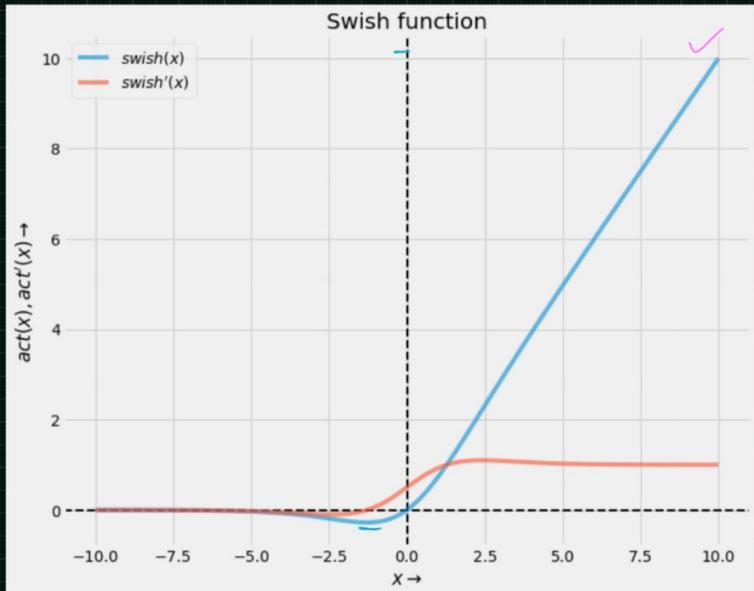
$$\text{swish}(z) = z \cdot \sigma(\beta \cdot z) = \frac{z}{1 + e^{-\beta z}}$$

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

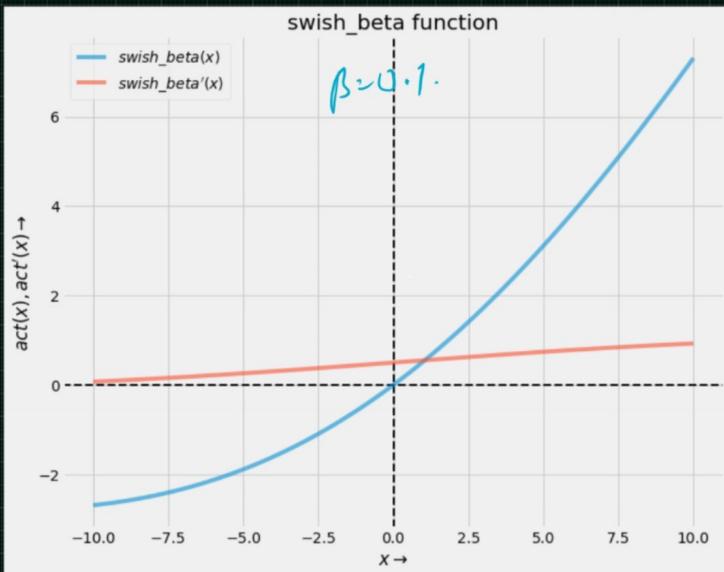
β can be a learnable parameter

for $\beta = 1$.

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



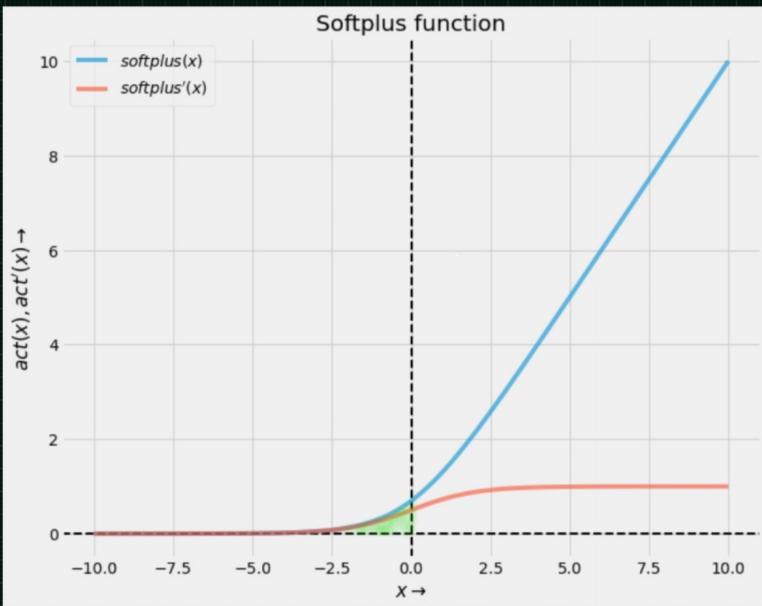
derivative of swish
domain $z \in (-\infty, \infty)$
range $\in (0, 1)$



- 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

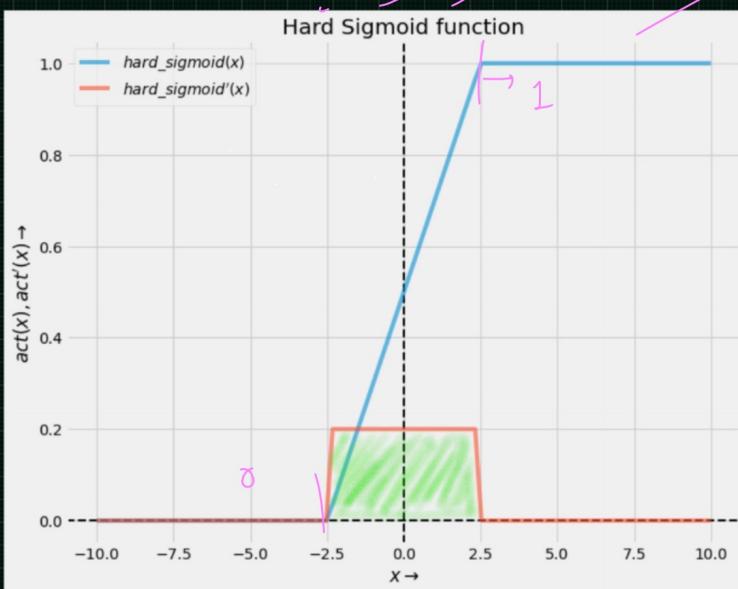
Advantages :-

- i) it's differentiable
everywhere

Hard Sigmoid

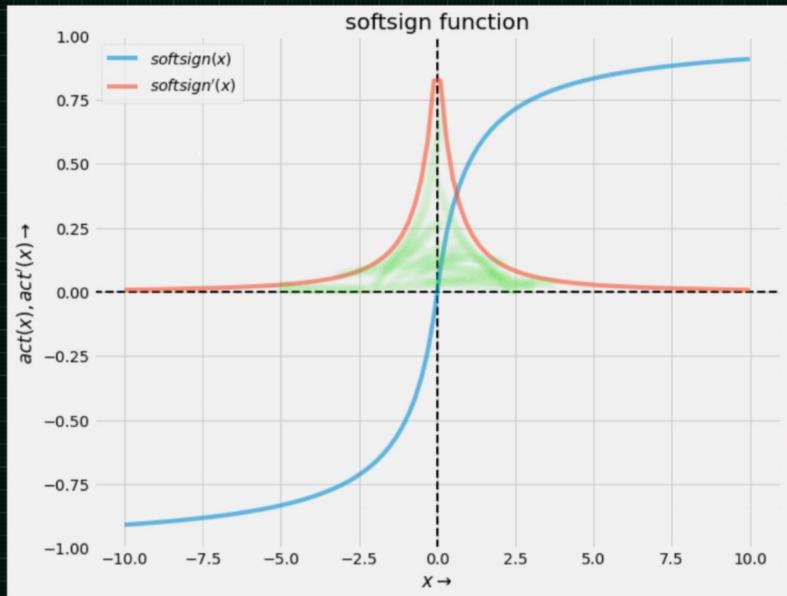
$$\phi(z) = \begin{cases} 0 & z < -2.5 \\ 1 & z > 2.5 \\ 0.2z + 0.5 & -2.5 \leq z \leq 2.5 \end{cases}$$

} if else

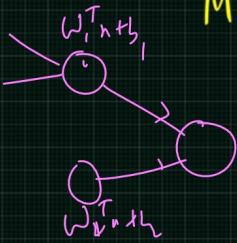


SOFTSIGN

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

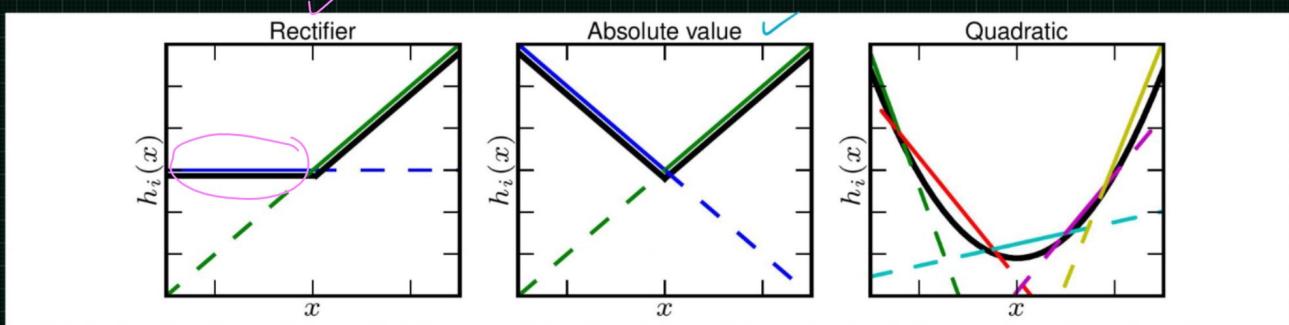


MAXOUT



Generalization of ReLU & Leaky ReLU

$$f(x) = \max(\underbrace{\underbrace{\omega_1^T x + b_1}_{(0)}, \underbrace{\underbrace{\omega_2^T x + b_2}_{(+)}})$$



Drawbacks :-

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

→ How to start / begin selection of activation function?

for any problem,

Go with ReLU for hidden layers

version of relu {as a future break}

for classification :-

binary :-

Sigmoid / softmax at output layer
o/p unit = 1 o/p unit = 2

Multi classification

softmax only at output layer

for Regression:-

No activation function

Result

DL → Architecture

CNN

std. data

Image net
CIFAR } previous