

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.

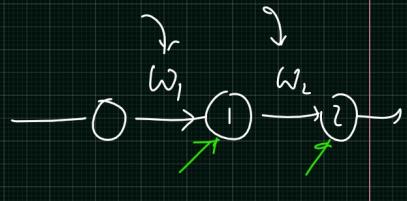


solution : $\{ \text{Regularization} \rightarrow \text{Dropout} \}$

Observation :-

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

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



$$z_i = \omega_L a_i + b$$

These are product of ratios

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

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

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

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

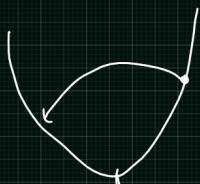
\Rightarrow Vanishing gradient \Rightarrow negligible weight update
 \Rightarrow leads to train lower layer very slow

2nd case

product terms $\gg 1$

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

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



\Rightarrow Exploding gradient is even
 \Rightarrow solution will diverge

for Vanishing & Exploding gradient.

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

in 2010 by

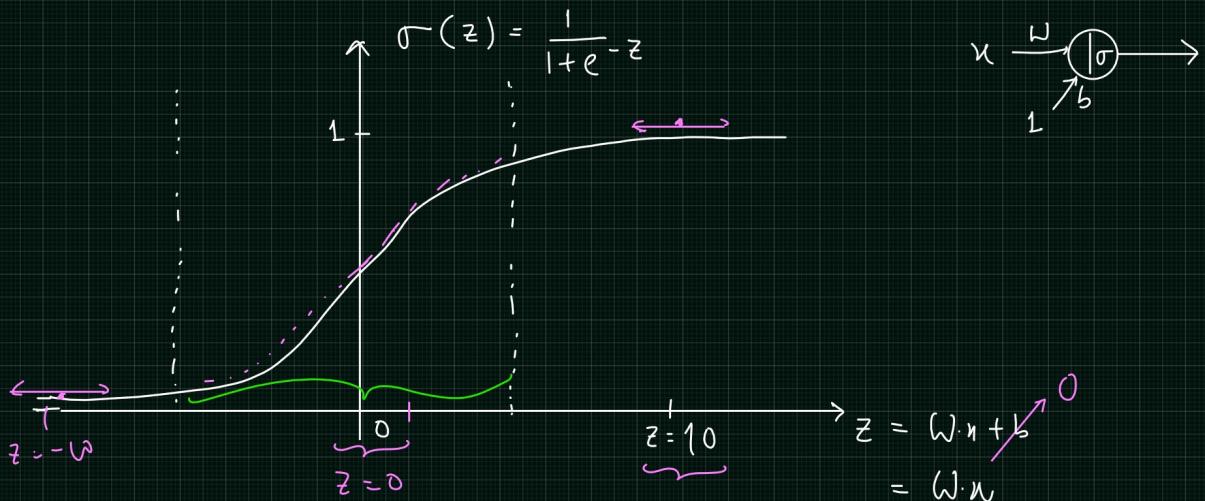
Xavier Glorot & Yoshua Bengio

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-

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

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

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

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

$$= 0$$

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

\Rightarrow zero weight update

Case 3:

$$w = 0$$

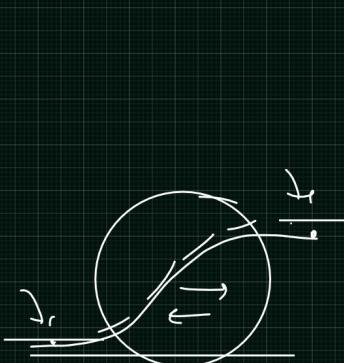
$$z = 0$$

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

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

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

$$= 0.25 \neq 0 \checkmark$$



f_p

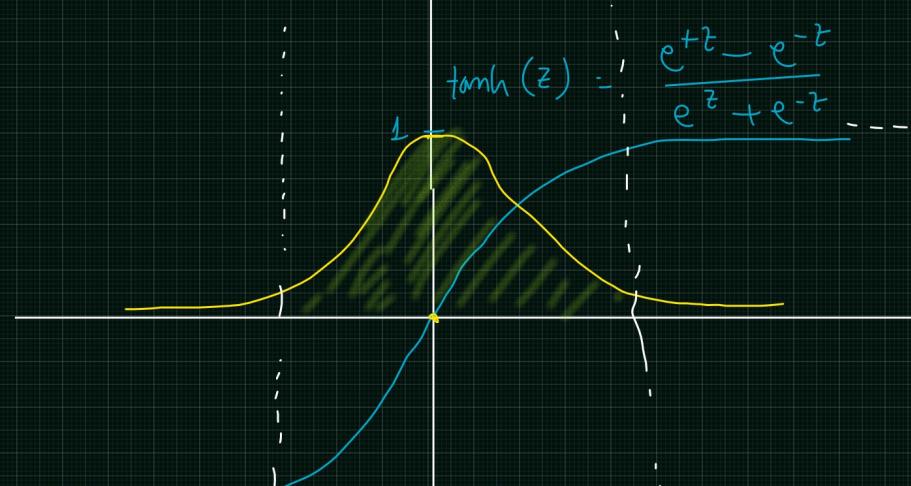
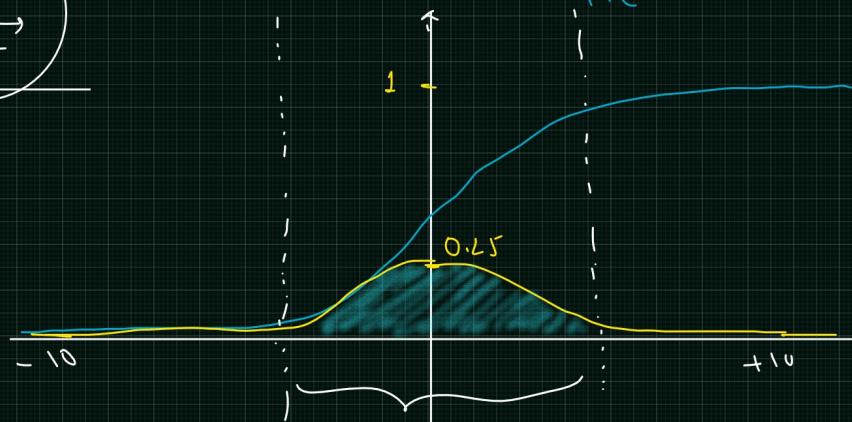
b_p

forward

backward prop.

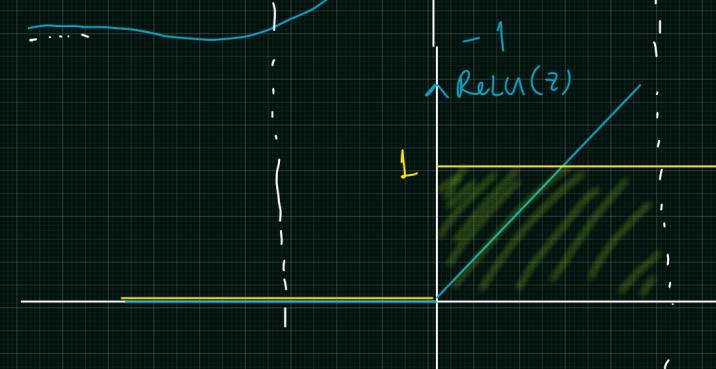
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

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



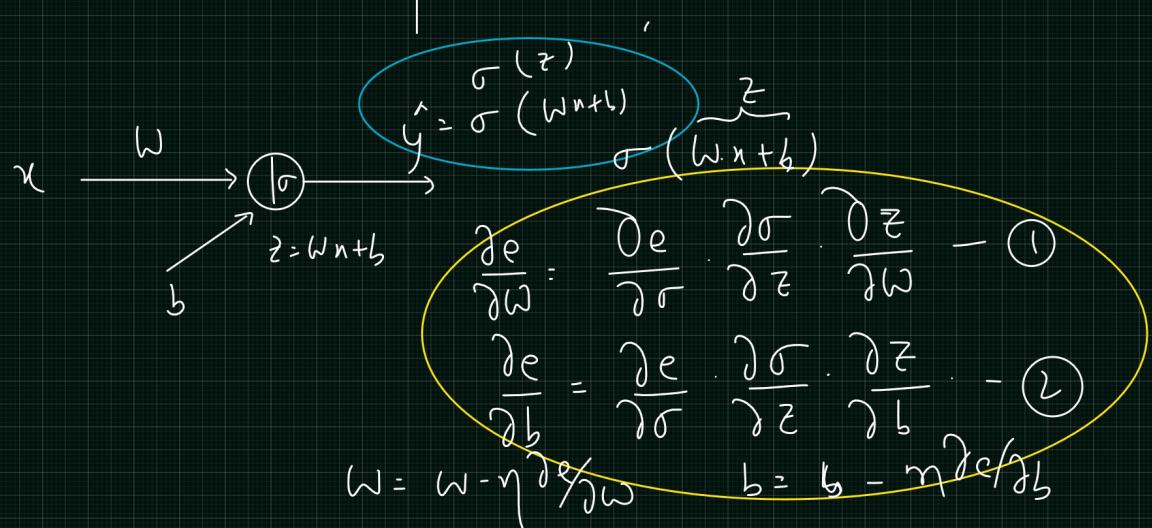
$$\tanh'(z) = 1 - \tanh^2(z)$$

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



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

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



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

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

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

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

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

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

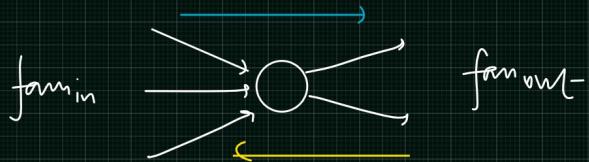
As per the paper,

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 → "

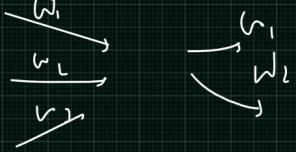


Condition to prevent vanishing & exploding gradient

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

alternative proposal

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



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

Calculation :-

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

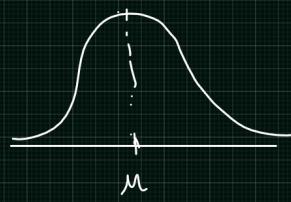
For sigmoid,

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

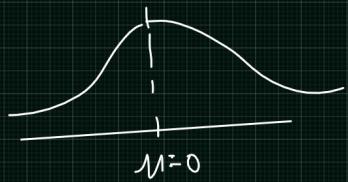
Normal distribution

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



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

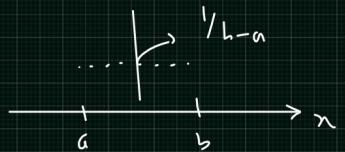
$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



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")

$$\omega_{100} \rightarrow$$

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