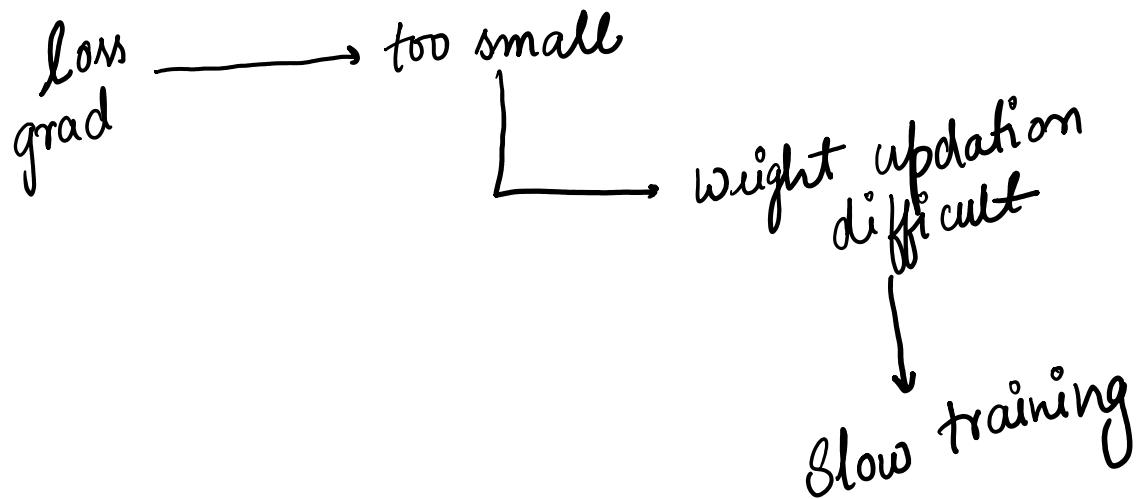


## WHAT YOU WILL STUDY IN TODAY VIDEO ?

-  What is Vanishing Gradient Problem? 
-  What is Exploding Gradient Problem? 
-  Dying ReLu Problem 

# ① Vanishing Gradient

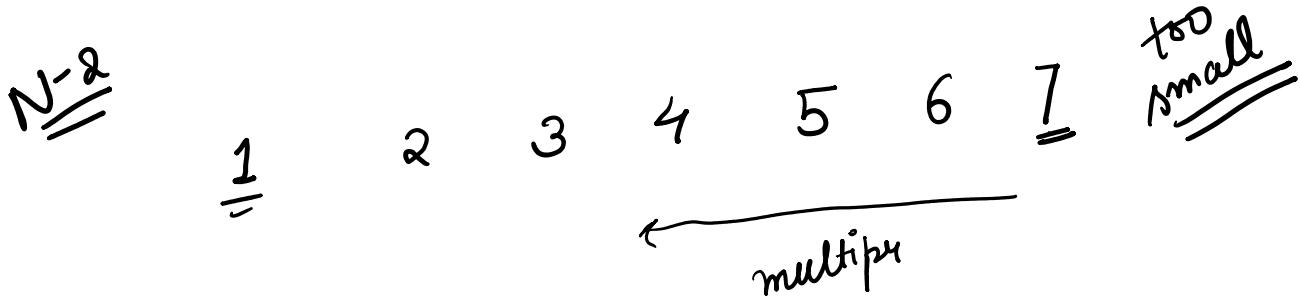
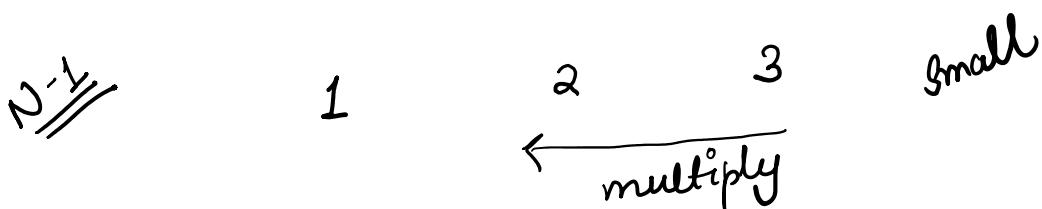
MLP



Causes

① Activation fn  $\xrightarrow{\text{main}}$  Sigmoid  
Tan h

② Too much hidden layers



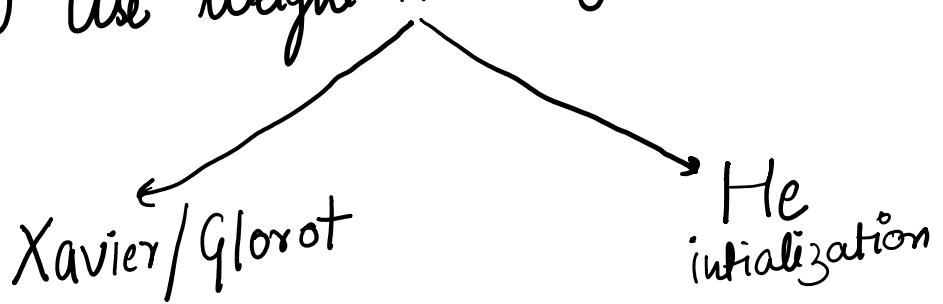
③ Weight initialization, using weights  $\frac{1}{n}$

③ Weight initialization  
Starting weights too small

fix:

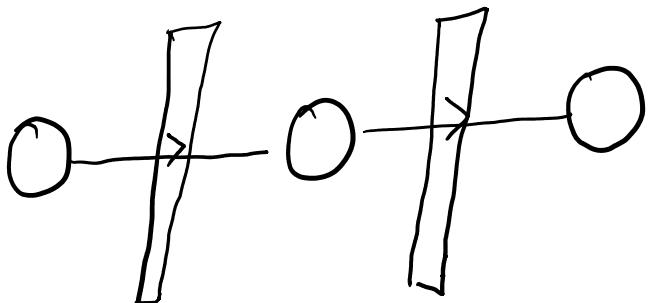
① Better act fm → ReLU  
→ leaky ReLU

② Use weight initialization techniques



③ Batch normalization layer

Gradient  
don't shrink  
much



Mathematical

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

Derivative of Sigmoid

$$\Rightarrow \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

~~of~~

~~3 layers~~ →  $w = 0.5 \equiv$   
 $x = 1$

~~1st layer~~ →  $z_1 = w \times x = 0.5 \times 1 = 0.5$   
 $a_1 = \frac{1}{1+e^{-0.5}} = 0.62 \equiv$

~~2nd layer~~ →  $z_2 = 0.5 \times 0.62 = 0.31 \equiv$   
 $a_2 = \frac{1}{1+e^{-0.31}} = 0.58 \equiv$

~~3rd layer~~ →  $z_3 = 0.5 \times 0.58 = 0.29 \equiv$   
 $a_3 = 0.57 \rightarrow \text{off} \equiv$

~~Backpropagation~~  
~~Compute gradient~~

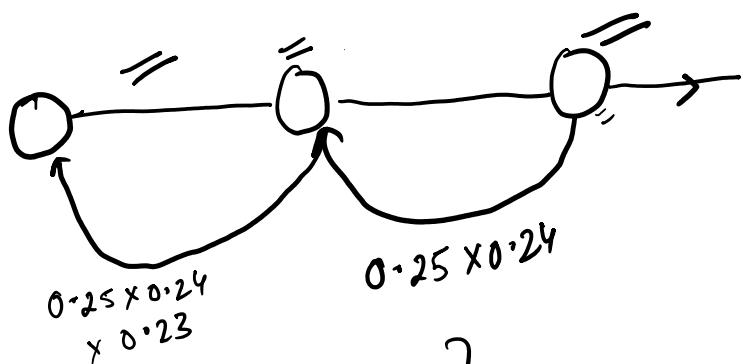
~~3rd layer~~ →  $\sigma'(z_3) = 0.57 \times (1 - 0.57)$   
 $= 0.25 \equiv$

~~2nd layer~~ →  $\sigma'(z_2) = 0.58 \times (1 - 0.58)$   
 $= 0.24 \equiv$

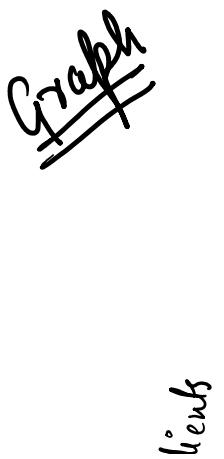
$\rightarrow 1 \rightarrow 1 = 0.62 \times (1 - 0.62)$

1st layer  $\rightarrow \sigma'(z_1) = 0.62 \times (1 - 0.62)$   
 $= \underline{\underline{0.23}}$

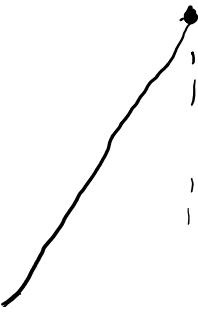
Overall gradient  $\rightarrow 0.25 \times 0.24 \times 0.23$   
 $= \underline{\underline{0.0138}}$   $\rightarrow$  small value

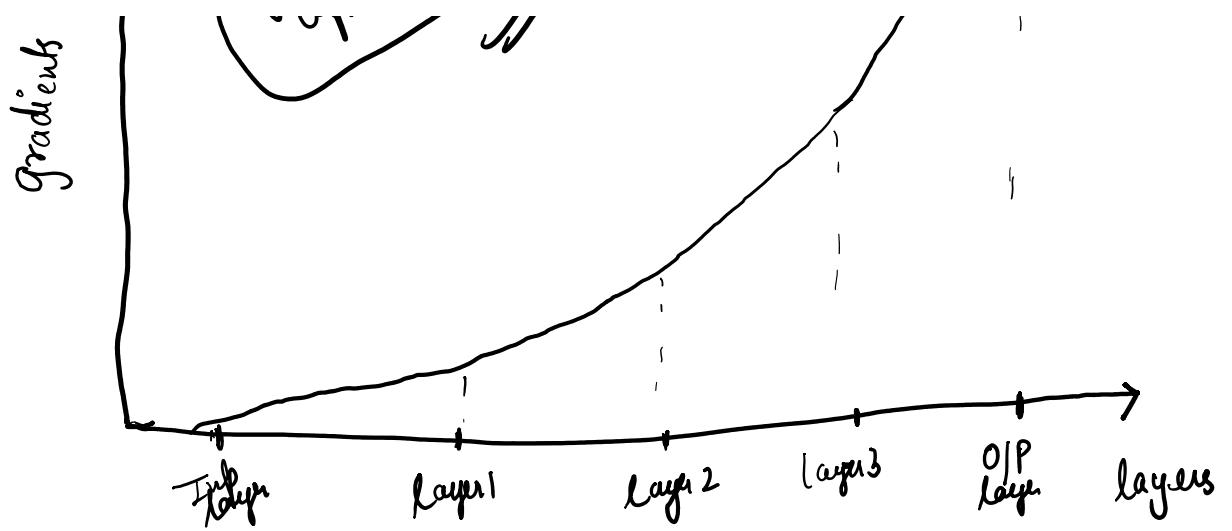


{  $\underline{\underline{\text{No updation}}}$   
 $\underline{\underline{\text{Weight}}}$  }

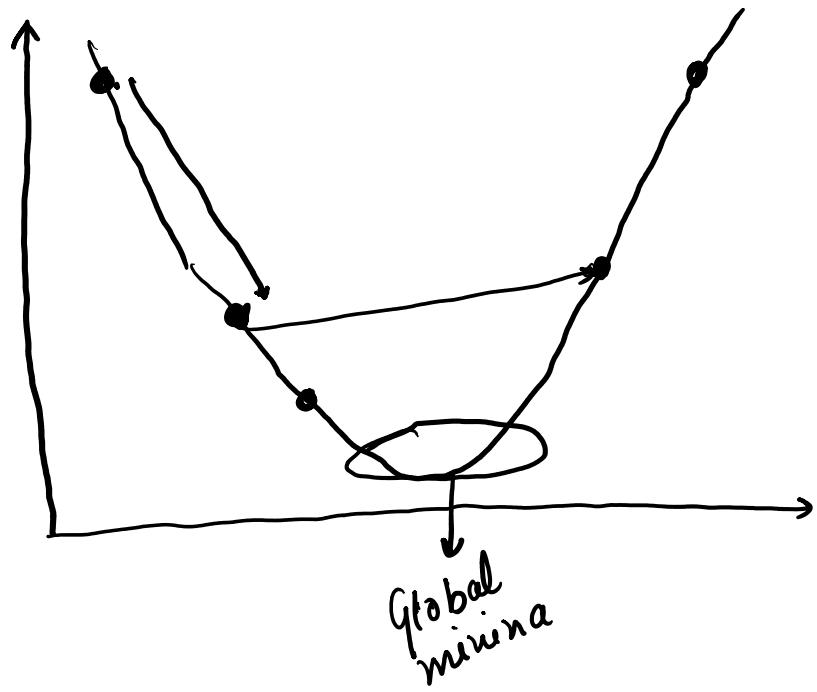
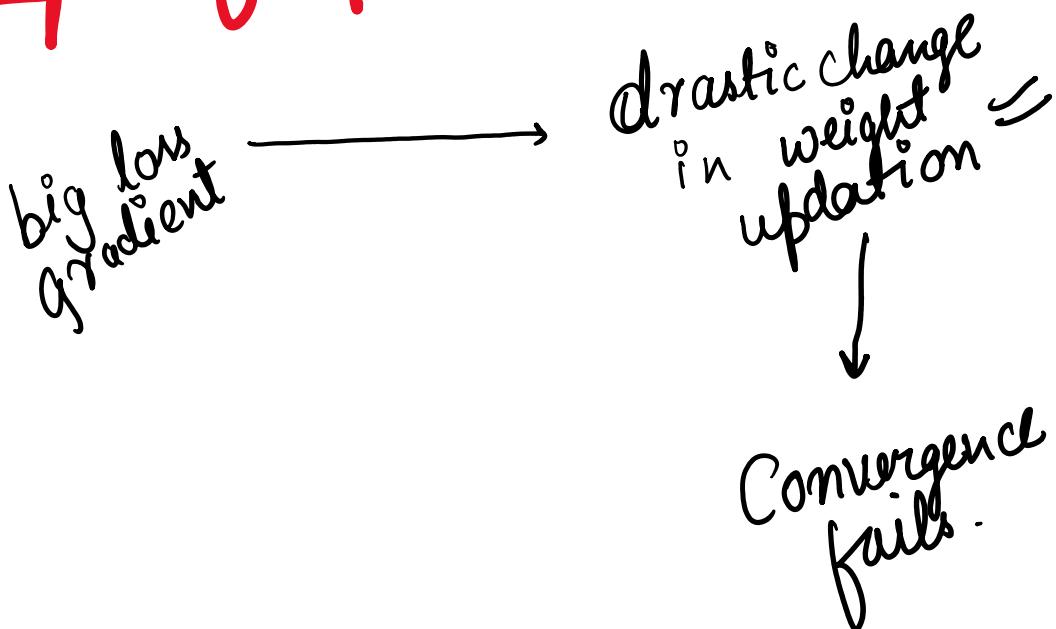


gradient approach zero





## ② Exploding Gradient



comes ↴

① Weight initialization ↴  
Starting weights are big ↴

RNN

② Too many layers =

③ High learning rate

$$\omega_{\text{new}} = \omega_{\text{old}} - \frac{\partial L}{\partial \omega} =$$

fix

① Weight initialization



② Batch normalization layer

Mathematical

Assume:

$$w = 3$$

$$x = 1$$

3 layered network

Act  $\rightarrow$  ReLU  $\rightarrow f(x) = \max(0, x)$

forward  
prop.

$$\begin{aligned} \text{Let } z_1 &= 3 \times 1 = 3 \\ \text{layer } 1 &a_1 = \max(0, 3) = 3 \end{aligned}$$

2nd layer  $\longrightarrow z_2 = 3 \times 3 = 9$   
 $a_2 = \underline{\underline{9}} =$

3rd layer  $\longrightarrow z_3 = 3 \times 9 = 27$   
 $a_3 = \underline{\underline{27}} \longrightarrow \underline{\underline{O/P}}$

Back Prop.

loss fn = MSE  
 $L = (y - \hat{y})^2$

3rd layer  $\frac{dL}{dz_3} = 1 =$   $\longrightarrow$  loss grad at layer 3.

loss grad  
 $w \times \frac{dL}{dz_3}$   
 weighted sum

2nd layer

Apply chain rule

$$\begin{aligned}\frac{dL}{dz_2} &= w \times \frac{dL}{dz_3} \\ &= 3 \times 1 \\ &= \underline{\underline{3}}.\end{aligned}$$

loss grad  
 for layer 2

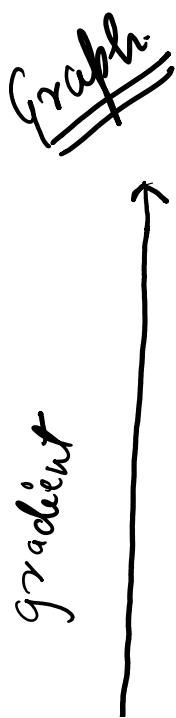
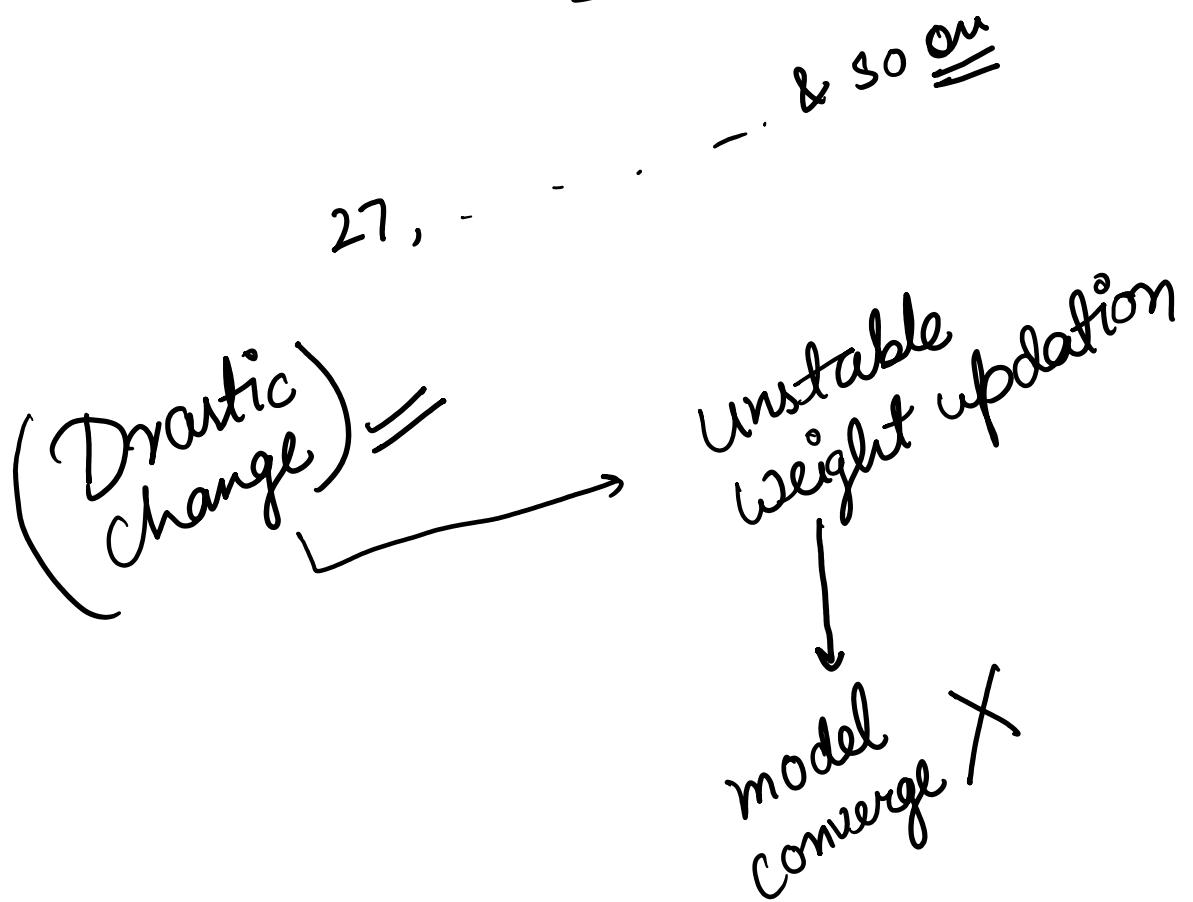
1st

$$\frac{dL}{dz_1} = w \times \frac{dL}{dz_2} = 3 \times 3$$

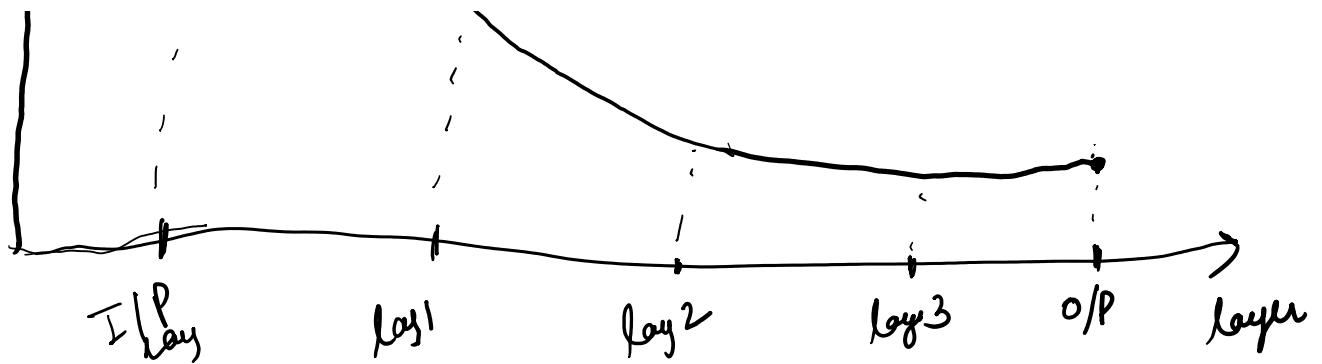
$$\frac{\partial L}{\partial z_1} = \underline{q.}$$

1st layer

27,



Loss gradient  
drastic increase



## ③ Dying ReLU

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

Cause → ReLU act fn

$$\begin{array}{rcl} 4 & \rightarrow & 4 \\ 5 & \rightarrow & 5 \\ -1 & \rightarrow & 0 \\ -2 & \rightarrow & 0 \\ -4 & \rightarrow & 0 \end{array} \quad =$$

Fix → Leaky ReLU    Parametric ReLU    ELU    =

