

# Vanishing Gradient Problem

- Vanishing gradient problem is a common problem that we face while training deep neural networks. Gradients of neural networks are found during back propagation.
- Generally, adding more hidden layers will make the network able to learn more complex arbitrary functions, and thus do a better job in predicting future outcomes. This is where Deep Learning is making a big difference.

Now during back-propagation i.e moving backward in the Network and calculating gradients, it tends to get smaller and smaller as we keep on moving backward in the Network. Below is Just a simple demonstration of Vanishing Gradient Problem in single layer.

weight updation formula

$$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \frac{\partial L}{\partial (w_{ij}^k)_{\text{old}}}$$

let give an example in single layer.

let we update  $w'_{11}$

let  $(w'_{11})_{\text{old}} = 2.5$   
 $\eta = 1$

Diagram of a single-layer neural network with 3 input nodes ( $x_1, x_2, x_3$ ), 2 hidden nodes ( $f_{11}, f_{12}$ ), and 1 output node ( $o_1$ ). Weights are  $w'_{11}, w'_{12}, w'_{21}, w'_{22}, w'_{31}, w'_{32}$ . The output is  $o_1/p$ . The loss is  $L$ .

weight updation formula

$$(w'_{11})_{\text{new}} = (w'_{11})_{\text{old}} - \eta \frac{\partial L}{\partial (w'_{11})_{\text{old}}}$$

by chain rule

$$\frac{\partial L}{\partial (w'_{11})_{\text{old}}} = \frac{\partial L}{\partial o_1} \times \frac{\partial o_1}{\partial f_{11}} \times \frac{\partial f_{11}}{\partial (w'_{11})_{\text{old}}}$$

here (only take one route for better underst.)

o/p layer

$$= [0.20 \times 0.5 \times 0.02]$$

as we backpropagate the derivative value decrease

$$= 0.002$$

[we know derivative of sigmoid is lies 0 to 0.25 so here all three derivative lies 0 to 0.25]

$(w'_{11})_{\text{new}} = 2.5 - 1 \times 0.002$   
 $\approx 2.49$

As we add more hidden layer the value become reduce & going to zero at a point.

by this eqn become

$$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}}$$

this is called vanishing gradient problem. as here gradient vanish.

- This Happen because of we use sigmoid and tanh activation function in hidden layer. As sigmoid and tanh deriative 0.25,1 respectively. so by calculating number of hidden layer the derivative becomes 0 so avoid it we use RELU activation function in hidden layer.

## Exploding gradient Problem

- We have discussed about vanishing gradient problem. Now we will get in to exploding gradient problem. Earlier we discussed what happens when our gradient becomes very small. Now we will discuss what will happen if it gets large.
- In deep networks or recurrent neural networks, error gradients can accumulate during an update and result in very large gradients.
- These in turn result in large updates to the network weights, and in turn, an unstable network. The explosion occurs through exponential growth by repeatedly multiplying gradients through the network layers that have values larger than 1.0. This will ultimately lead to a total unstable network.

