# VANISHING GRADIENT PROBLEM :-
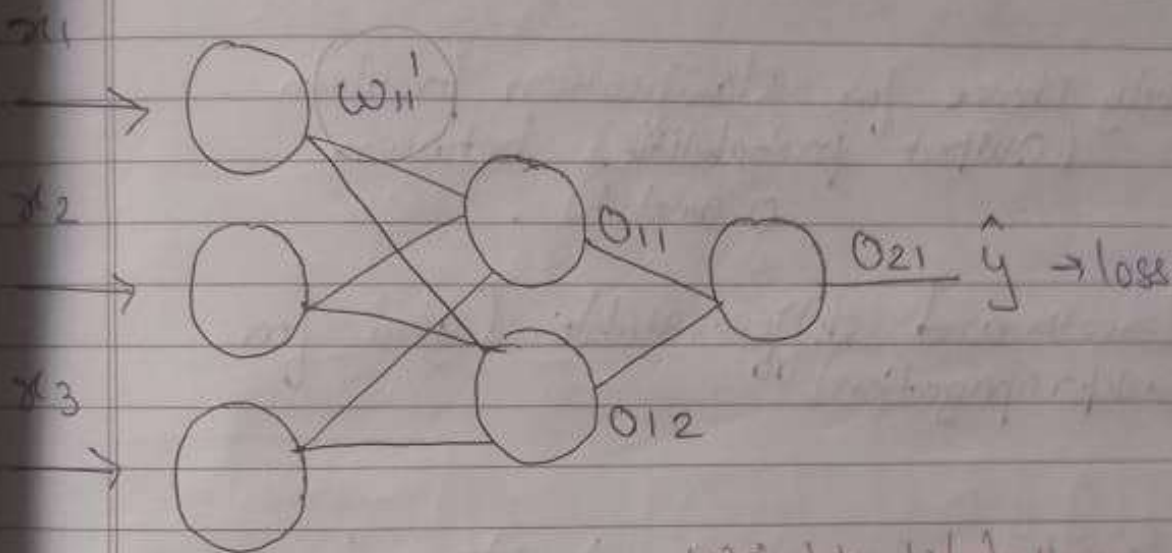
Previously we were using Sigmoid at that time ReLU was not invented. the problem faced was Vanishing Gradient problem.
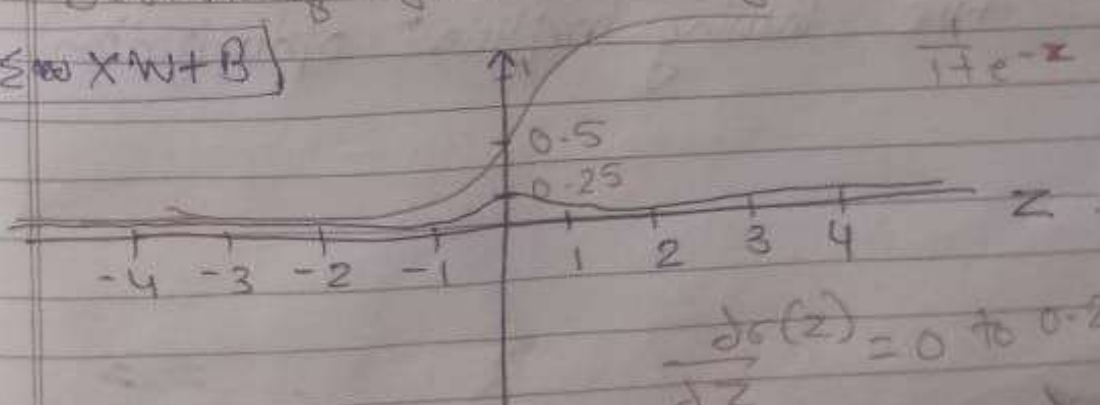
$x_1$ →
$x_2$ →
$x_3$ →

$w_{11}'$
$O_{11}$
$O_{12}$
$O_{21}$ $\hat{y}$ → loss

### Weight Updation:-

$$W_{11} \, new = W_{11} \, old - \eta \frac{\partial L}{\partial W_{11}' \, old}$$

$$\frac{\partial L}{\partial W_{11}} = \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial W_{11}'} \quad \rightarrow \text{chain Rule}$$

Note :- Derivative of Sigmoid will Range between 0 to 0.25

$$Z = \sum X \cdot W + B$$

$$\frac{1}{1 + e^{-x}}$$

0.5
0.25

-4  -3  -2  -1    1   2   3   4    Z

$$\frac{d\sigma(z)}{dz} = 0 \text{ to } 0.25$$

$$0 \leq \sigma(z) \leq 0.25$$

Sigmoid Activation →

$$f(x) = \frac{1}{1 + e^{-x}}$$

- Early choice for classification problem.
  (Output probabilities between
  0 and 1).

- Smooth and differentiable (good for
  Backpropagation)

- Vanishing gradient problem :
  Gradients become tiny, slowing
  learning in deep networks.

  for large/small inputs gradients are
  almost zero.

  Not zero-Centred (all outputs are
  +ve, causing inefficient optimization)

ReLU Activation →

$$f(x) = max(0, x)$$

- Common in deep networks to avoid vanishing gradient.

- No Vanishing gradient for +ve Inputs (faster learning)

- Simple & Efficient to Compute

- Sparse activation :- Most neurons are inactive (output 0)

- Dying ReLU problem → Neurons can get stuck outputting 0 and stop learning.