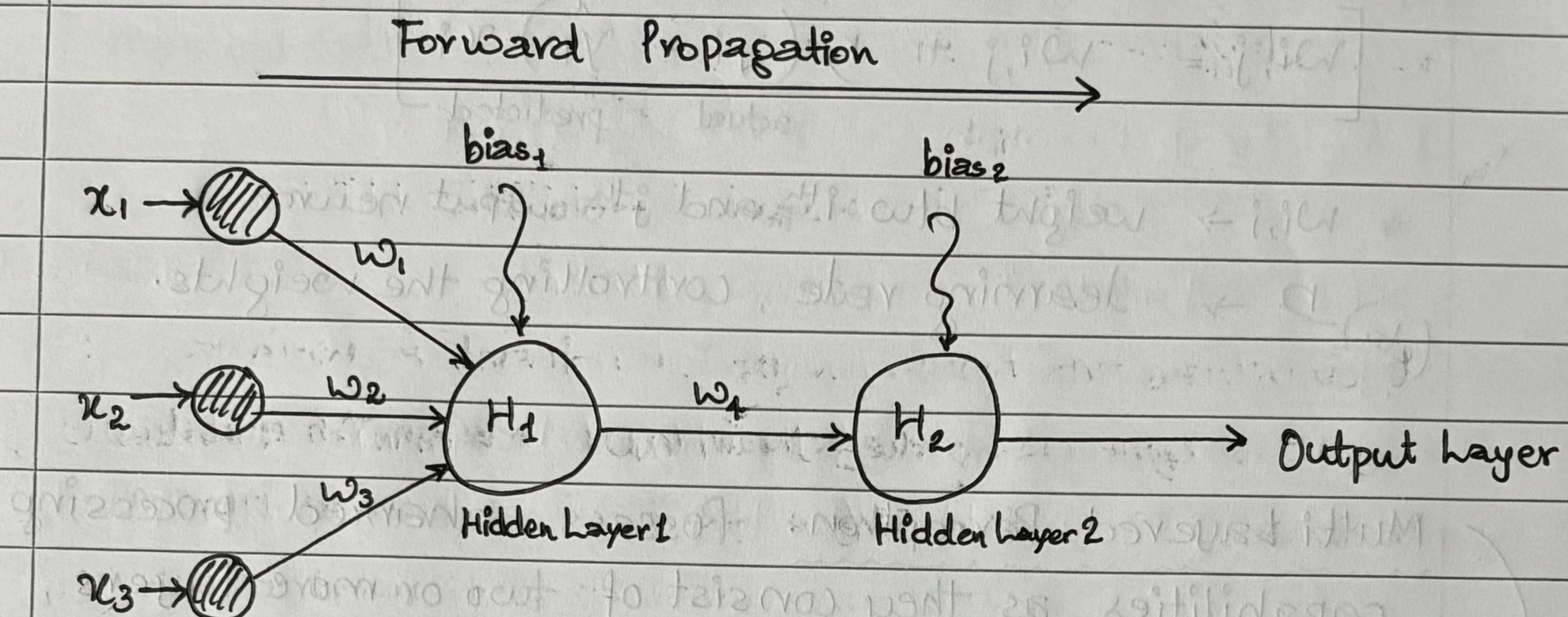


* ① Forward Propagation*: The input data moves through each layer of neural network where each neuron applies weighted sum, add bias, passes the result through an activation function and making predictions. This process is crucial before backpropagation updates weights.



→ Hidden Layer 1: Step 1: $Z = \sum_{i=1}^n x_i w_i + \text{bias}_1$

Step 2: Activation (Z)

→ Hidden Layer 2: Step 1: $Z = \text{Output } H_1 * w + \text{bias}_2$

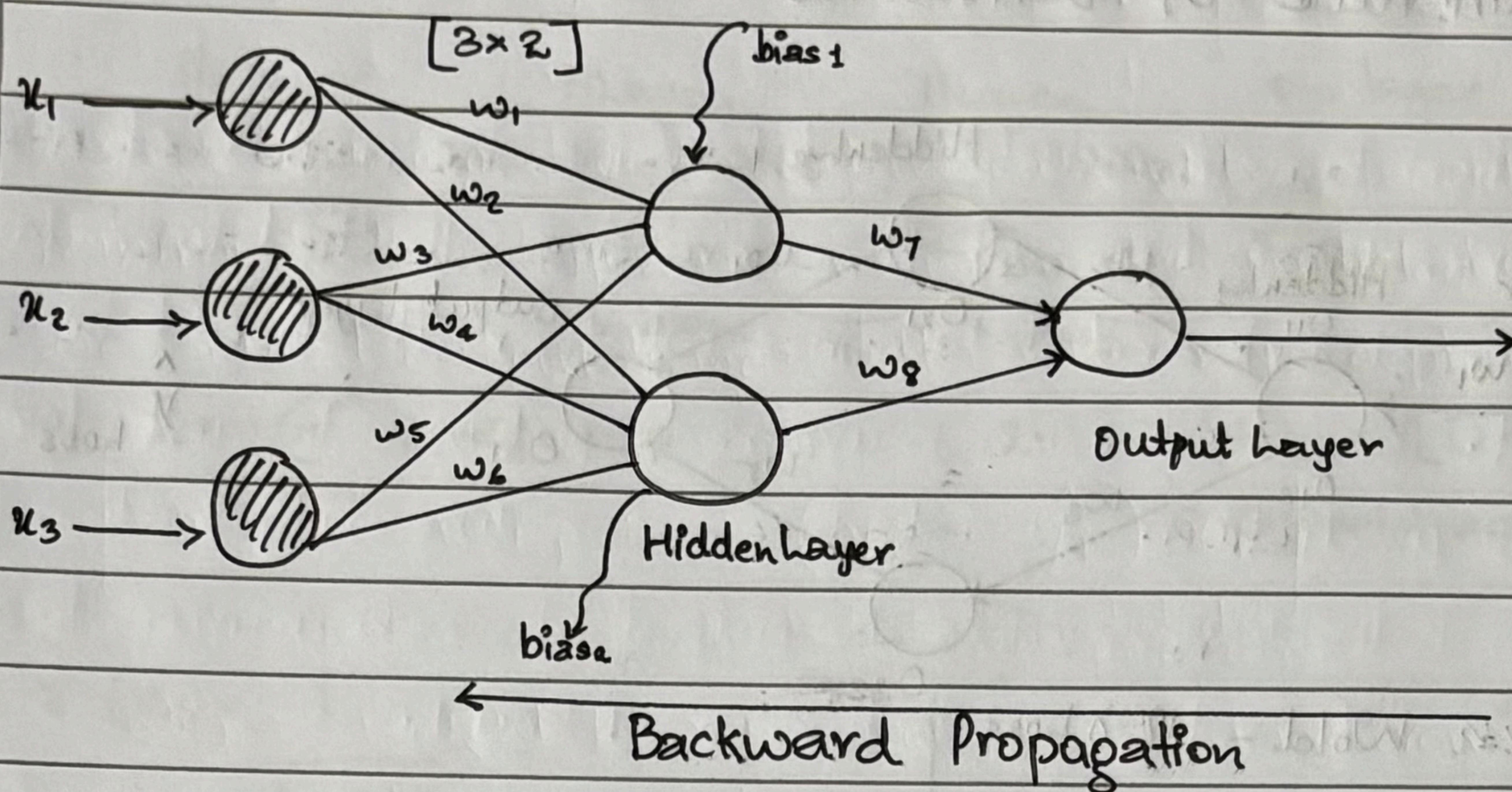
Step 2: Activation (Z) → \hat{y} Output predicted

→ Calculate Loss function → (Actual output - Predicted output)

If loss function is high we will update weights and continue to reduce loss function.

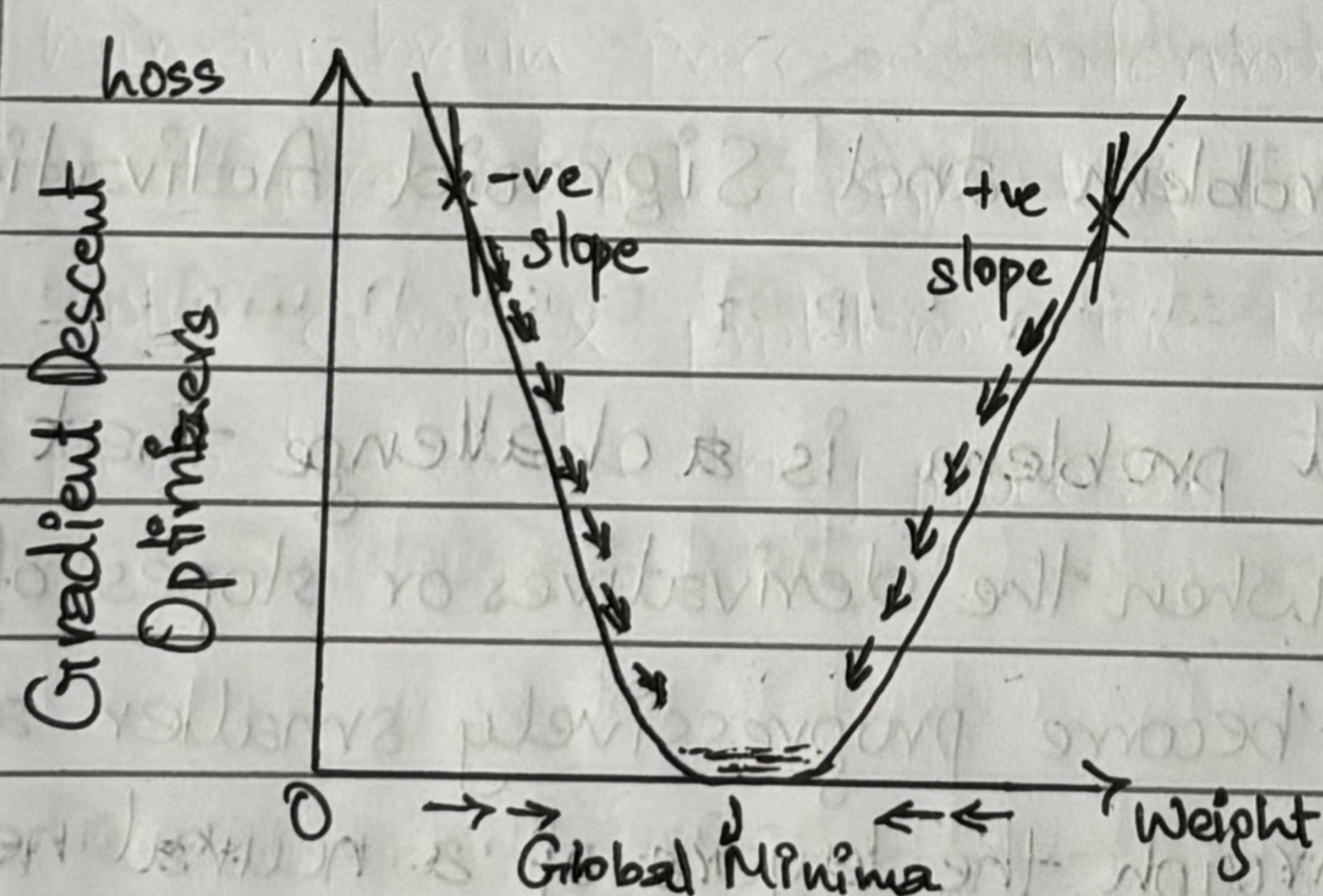
* ② Backpropagation*: Also known as Backward Propagation of Errors is a method used to train neural network.

It's goal is to reduce the difference between the model's predicted output and actual output by adjusting weights and biases.



It computes the gradient of loss function with respect to each weight using the chain rule making it possible to update weights efficiently. These gradients indicate how much each weight and bias should be adjusted to minimize the error in the next iteration. The activation function through its derivative plays a crucial role in computing these gradients during Back Propagation.

$$\rightarrow \text{Weight Updation formula: } W_{\text{new}} = W_{\text{old}} - \eta \frac{\delta \text{loss}}{\delta W_{\text{old}}} \quad \text{learning rate}$$



* We use optimizers to reduce the loss value

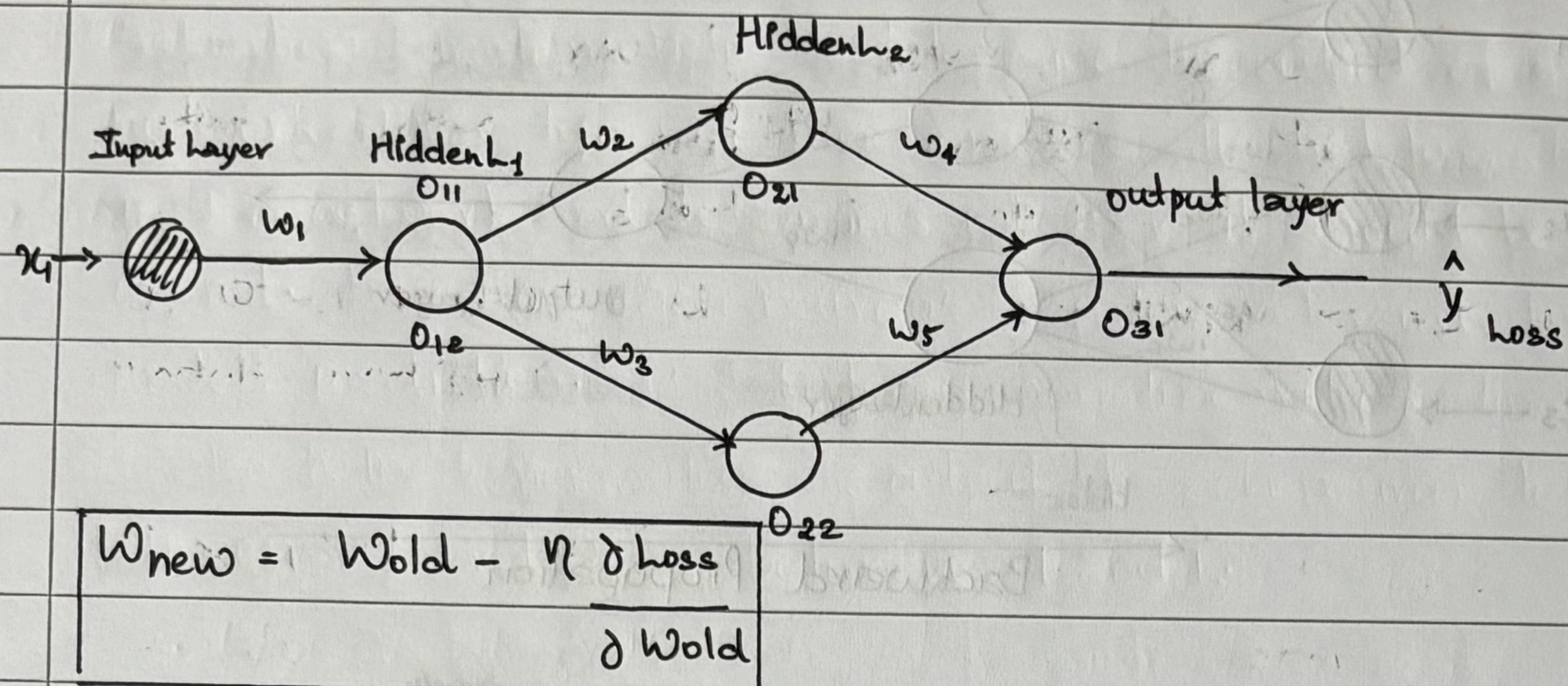
* η → should be always a small value (0.001)

* When $W_{\text{new}} = W_{\text{old}}$, we reach global minimum

$$\begin{aligned} \text{for -ve slope, } W_{\text{new}} &= W_{\text{old}} - \eta (-\text{ve}) \\ &= W_{\text{old}} + \eta (+\text{ve}) \\ W_{\text{new}} &>> W_{\text{old}} \end{aligned}$$

$$\begin{aligned} \text{for +ve slope, } W_{\text{new}} &= W_{\text{old}} - \eta (+\text{ve}) \\ &= W_{\text{old}} - \eta (+\text{ve}) \\ W_{\text{new}} &<< W_{\text{old}} \end{aligned}$$

③ *Chain Rule of Derivatives*



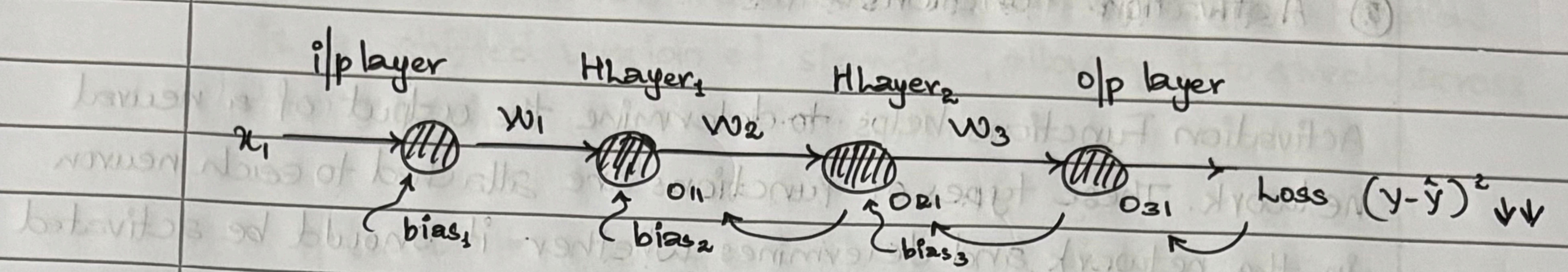
$$\text{Calculating } w_{\text{new}} = w_{\text{old}} - \eta \frac{\delta \text{loss}}{\delta w_{\text{old}}}$$

$$\Rightarrow \frac{\delta \text{loss}}{\delta w_{\text{old}}} = \left[\frac{\delta \text{loss} \times \delta O_{31} \times \delta O_{21} \times \delta O_{11}}{\delta O_{31} \delta O_{21} \delta O_{11} \delta w_{\text{old}}} \right]$$

$$+ \left[\frac{\delta \text{loss} \times \delta O_{31} \times \delta O_{22} \times \delta O_{11}}{\delta O_{31} \delta O_{22} \delta O_{11} \delta w_{\text{old}}} \right]$$

④ *Vanishing Gradient Problem and Sigmoid Activation*

The vanishing gradient problem is a challenge that emerges during backpropagation when the derivatives or slopes of the activation functions become progressively smaller as we move backward through the layers of a neural network.

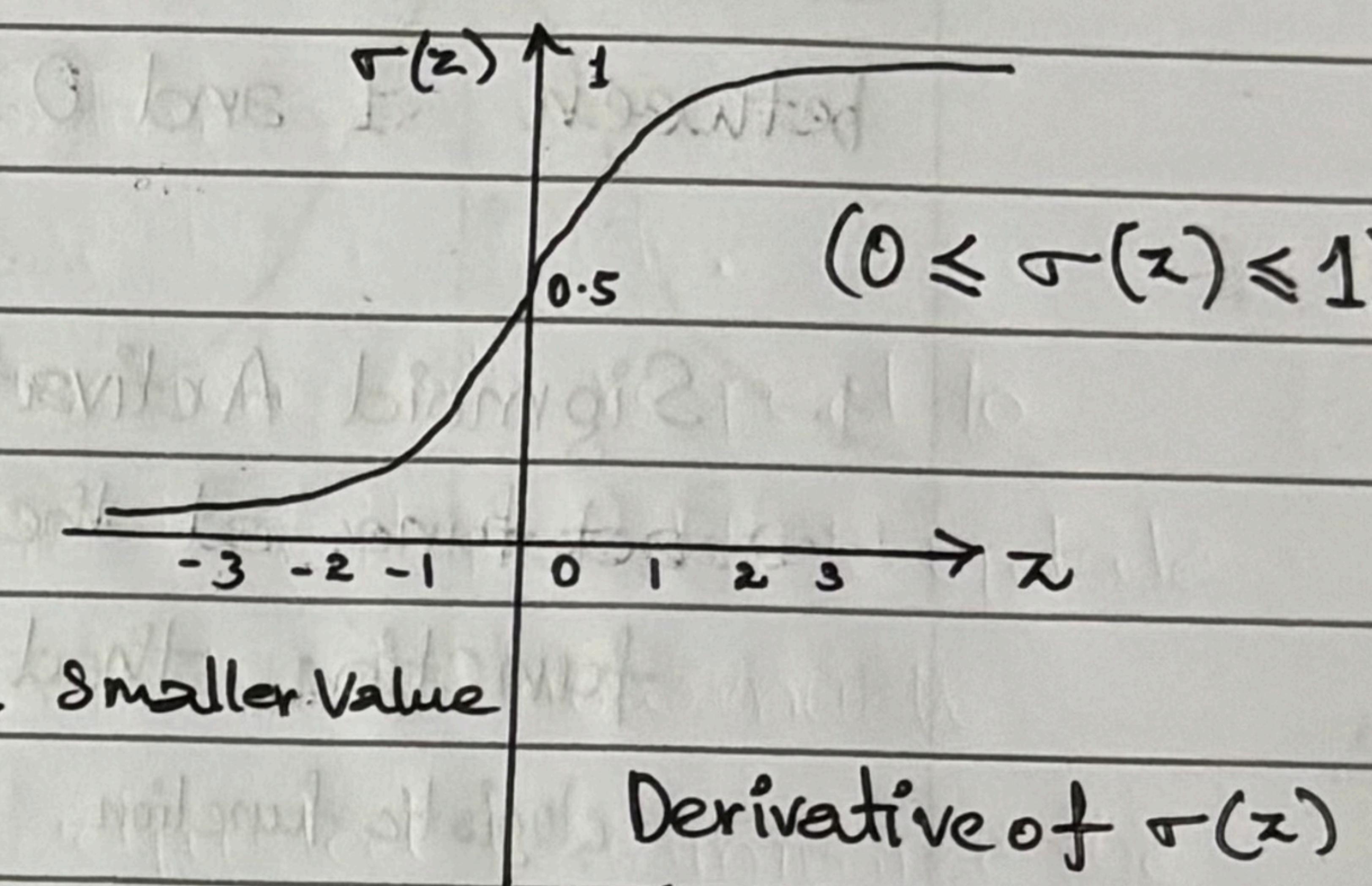


$$z = \sum_{i=1}^n x_i w_i + b$$

$$\sigma(z) \rightarrow [0 \text{ to } 1]$$

↳ Sigmoid activation function

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



$$w.r.t w_{1new} = w_{1old} - \eta \frac{\delta \text{loss}}{\delta w_{1old}}$$

$$\frac{\delta \text{loss}}{\delta w_{1old}}$$

Small Value

$$\text{Derivative of } \sigma(z) \quad (0 \leq \delta(\sigma(z)) \leq 0.25)$$

$$\frac{\delta \text{loss}}{\delta w_{1old}} = \frac{\delta \text{loss}}{\delta O_{31}} * \frac{\delta O_{31}}{\delta O_{21}} + \frac{\delta \text{loss}}{\delta O_{21}} * \frac{\delta O_{21}}{\delta O_{11}}$$

$$O_{31} = \sigma(w_3 * O_{21} + b_3)$$

$$O_{31} = \sigma(z) \quad z \text{ (input going before the activation)}$$

$$\rightarrow (\text{Derivative of sigmoid}) \Rightarrow 0 \leq \sigma(z) \leq 0.25$$

$$\frac{\delta O_{31}}{\delta O_{21}} = \frac{\delta(\sigma(z))}{\delta(z)} * \frac{\delta z}{\delta O_{21}}$$

$$\Rightarrow 0 \leq \sigma(z) \leq 0.25 * \frac{\delta((w_3 * O_{21}) + b_3)}{\delta(O_{21})}$$

So, as δ is around 0.01 and $\delta \text{Loss} / \delta w_{1old} \approx 0.25$,

$w_{1new} \approx w_{1old}$, the weights are never getting updated and we are not moving to the convergence zone. It is because of sigmoid.