

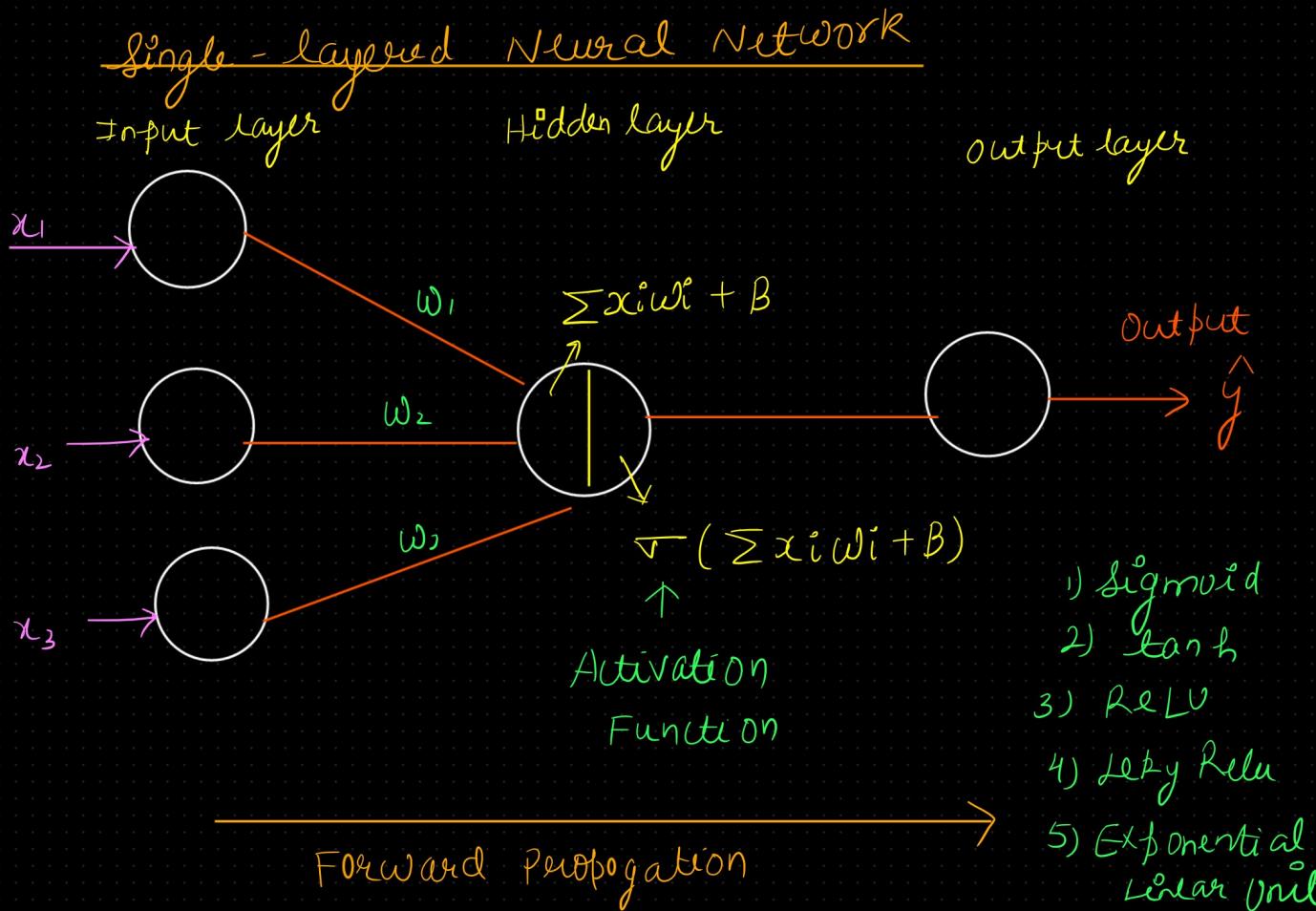
Forward and Backward Propagation

Forward Propagation

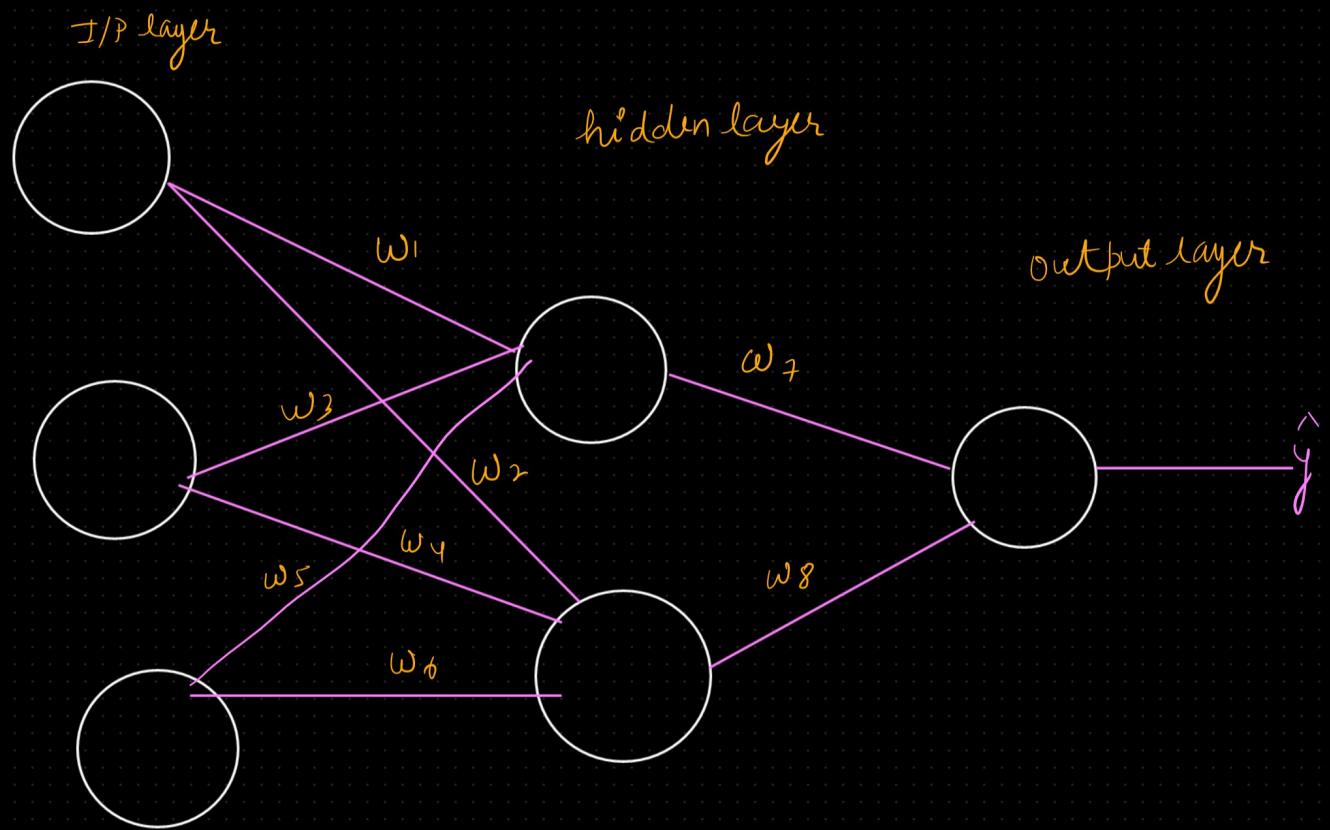
- Forward propagation is the process of moving from the input layer through the hidden layers to the output layer of a neural network to make predictions.
- Each neuron in a neural network layer receives inputs from the previous layer, applies weights and biases, and then passes the result through an activation function to produce an output.

$$y = \sigma(\sum x_i w_i + b)$$

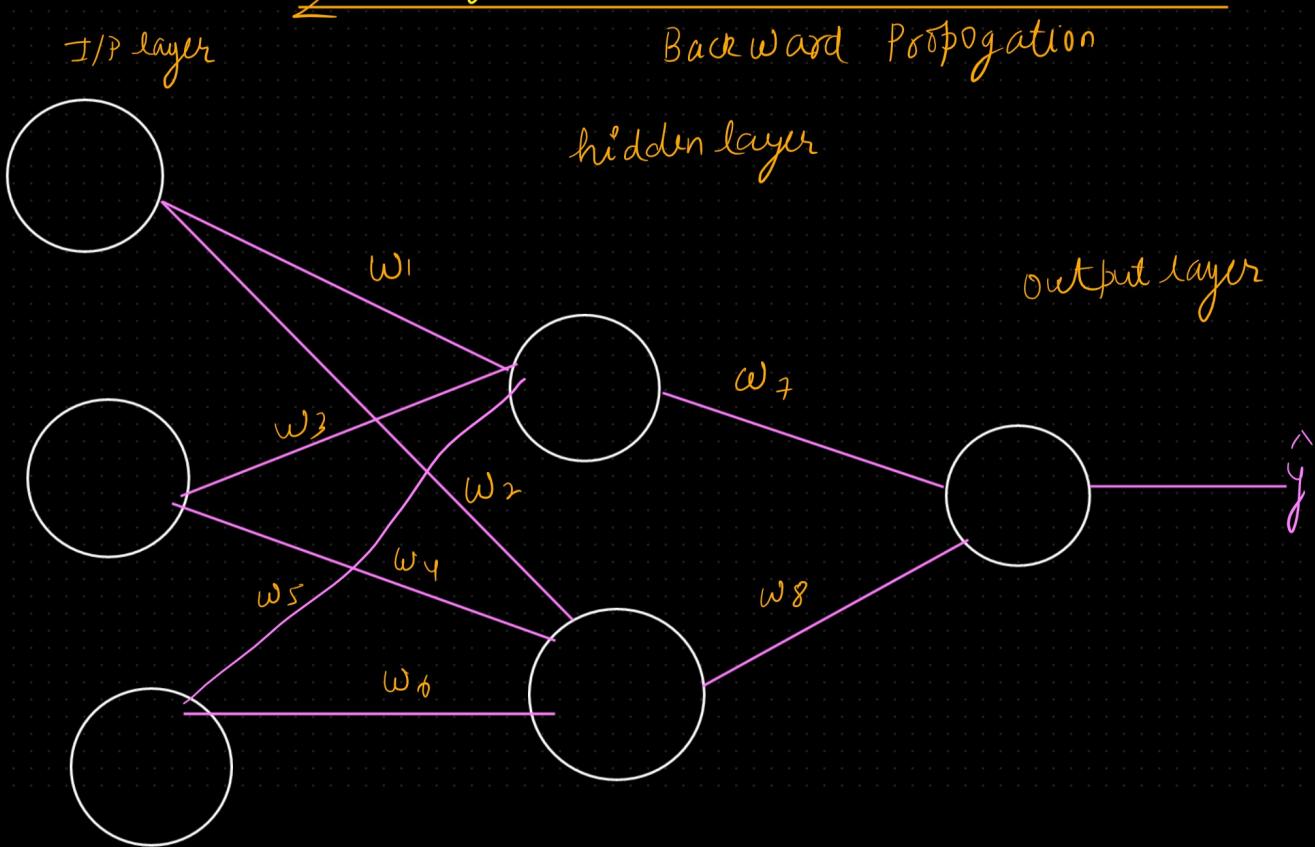
- This process is repeated layer by layer until the output layer is reached, and the final predictions are obtained.



Multi-layered Neural Network



Backward Propagation



① Error → Diff b/w predicted output and actual target output

$$\text{Error} = \hat{y} - y$$

① weights update

② loss funcn ↓↓

→ Backward propagation, also known as backpropagation, is the process of updating the weights of the neural network based on the error between the predicted output and the actual target output

Step 2 : \rightarrow diff. → rate of change

It involves calculating the gradient of the loss function with respect to each weight and bias in the network, and then using this gradient to update the weights and biases in the opposite direction of the gradient to minimize the loss.

single data point loss function $\begin{cases} \text{Regression} \rightarrow \text{MSE} \\ \text{Classification} \rightarrow \text{Log loss} \\ \text{Cross Entropy} \end{cases}$

$$\text{loss function} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 = \text{mse}$$

cost function = $\frac{1}{2n} \sum_{i=1}^n (y - \hat{y})^2$ batches data points
learning rate - step size



Step 3:

The gradient of the loss function with respect to the weights and biases is calculated using the chain rule of calculus, which allows us to propagate the error backward through the network.

Chain Rule

Derivative → Rate of change

$$f(x) = y$$

$$f(x) = x^2$$

- Power Rule

$$f'(x) = nx^{n-1}$$

n = power

x = variable

$$f(x) = x^2$$

$$f'(x) = 2x^{2-1} = 2x$$

$$f(x) = 4x^3$$

$$f'(x) = 12x^2$$

Partial derivative

$$f(x)$$

$$f(x, y)$$

w.r.t x

$$\frac{\partial f}{\partial x} \quad \text{taking } y \text{ as constant}$$

Similarly $\frac{\partial f}{\partial y}$ w.r.t y
taking x as constant

$$\therefore f(x, y) = x^2 + 2xy + y^2$$

$$\frac{\partial f}{\partial x} = \frac{\partial}{\partial x} (x^2 + 2xy + y^2)$$

$$= 2x + 2y$$

$$\frac{\partial f}{\partial y} = \frac{\partial}{\partial y} (x^2 + 2xy + y^2)$$
$$= 2x + 2y$$

Chain Rule

$$f(x) \quad g(x)$$

$$\rightarrow f(g(x))$$

The chain rule says that to find the rate of change of the outer function with respect to the inner function, you multiply the rate of change of the outer function with respect to its input by the rate of change of the inner function with respect to its input.

$$\frac{dy}{dx} = \text{Derivative of } y \text{ w.r.t } x$$

$$y = f(g(x)) \Rightarrow \frac{dy}{dx} = f'(g(x)) \cdot \frac{dy}{dx}$$

$y = f(u)$ and $u = g(x)$

Chain Rule

$$\boxed{\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}}$$

$$f(u) = u^2 \quad \text{and} \quad g(x) = 3x$$

$$\therefore \frac{dy}{du} = 2u$$

$$\therefore \frac{du}{dx} = 3$$

$$\therefore \frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx} = 2u \cdot 3 = 6u$$

$$= 6(3x) = 18x$$

Step 4: loss function $\rightarrow L$

Weight $\rightarrow w$

$$Z \rightarrow \hat{y} - y$$

By using chain rule

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial Z} \cdot \frac{\partial Z}{\partial w}$$



Gradient of loss function

\rightarrow Step 5

This process is repeated layer by layer in the reverse order until the gradients for all weights and biases are calculated.

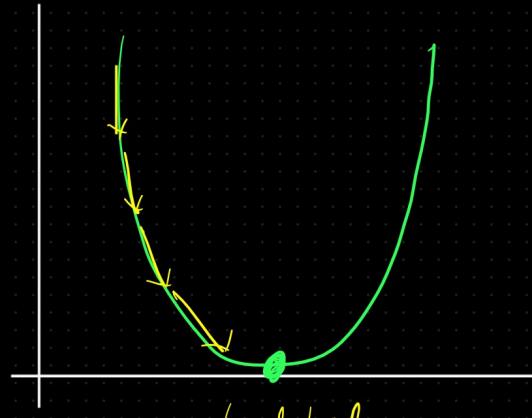
\rightarrow Step 6

Once the gradients are obtained, the weights and biases are updated using an optimization algorithm such as gradient descent:

Weight Updation Formula

$$w_{\text{new}} = w_{\text{old}} - \alpha \cdot \frac{\partial L}{\partial w}$$

α = learning rate \Rightarrow step size



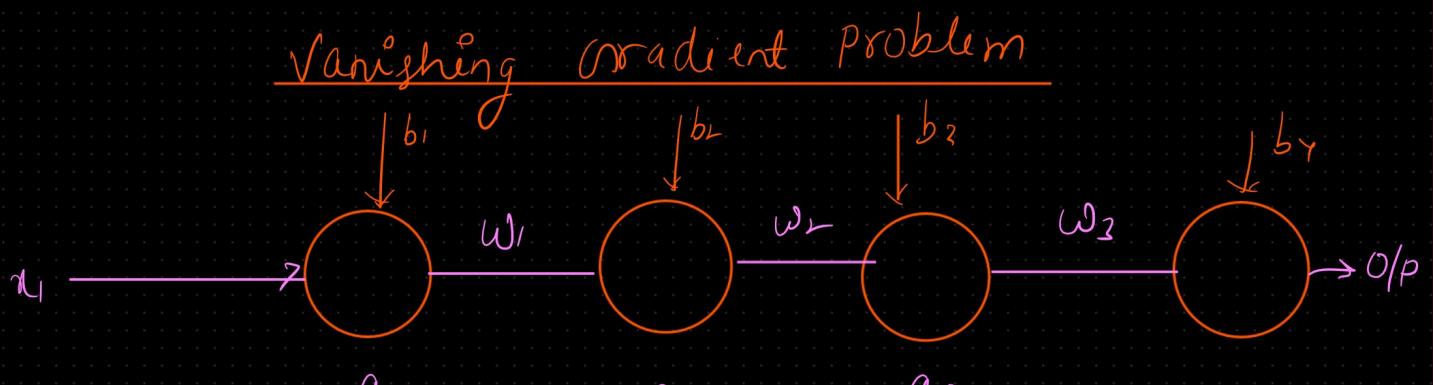
Global
Minima

The intuition behind backward propagation is that it adjusts the weights and biases in such a way that the error between the predicted output and the actual output is minimized. It does this by iteratively updating the parameters in the direction that reduces the error, thus improving the network's ability to make accurate predictions.

λ - learning rate

hyperparameter
↓
step size

The learning rate is a hyperparameter that determines the step size at which the weights of a neural network are updated during training. In the context of training a neural network using optimization algorithms like gradient descent, the learning rate controls how much we are adjusting the weights of our network with respect to the loss gradient.



\rightarrow sigmoid
 $[= 0, 1)$ Funcn'

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a_1} \times \frac{\partial a_1}{\partial a_2} \times \frac{\partial a_2}{\partial a_3} \approx \text{small number}$$

$$\therefore w_{\text{new}} = w_{\text{old}} - \eta (\text{small number})$$

$w_{\text{new}} \approx w_{\text{old}}$ \Downarrow vanishing gradient problem

No change in weight

Activation Func¹
 \downarrow $y = \text{ReLU}(a^4)$
ReLU