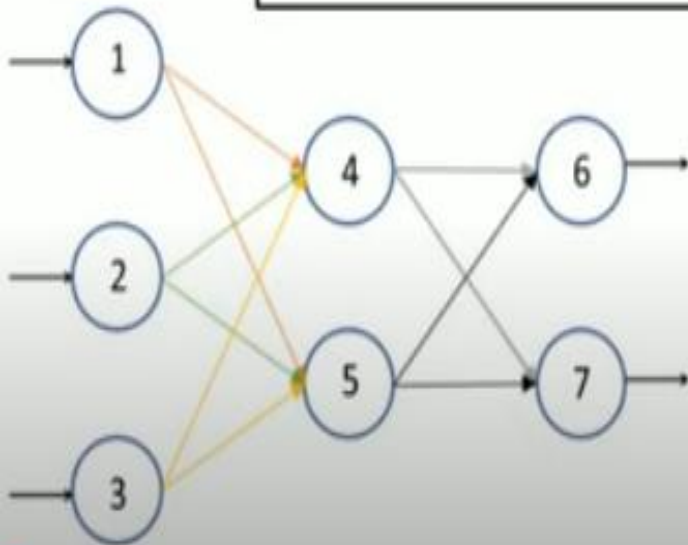


Back Propagation Algorithm

Derivation



**Weight
Update**

**Artificial
Neural
Networks**

**Machine
Learning**

Subscribe

Derivation of Back Propagation Algorithm

- To derive the equation for updating weights in back propagation algorithm, we use Stochastic gradient descent rule.
- Stochastic gradient descent involves iterating through the training examples one at a time, for each training example \mathbf{d} descending the gradient of the error E_d with respect to this single example.
- In other words, for each training example \mathbf{d} every weight w_{ji} is updated by adding to it Δw_{ji} .

- That is,

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

Where

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

Derivation of Back Propagation Algorithm

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

Where

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

- where E_d is the error on training example \mathbf{d} , that is half the squared difference between the target output and the actual output over all output units in the network,

$$E_d(\vec{w}) \equiv \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

- Here outputs is the set of output units in the network, t_k is the target value of unit k for training example \mathbf{d} , and o_k is the output of unit k given training example \mathbf{d} .

Derivation of Back Propagation Algorithm

Notation Used:

x_{ji} = the i^{th} input to unit j

w_{ji} = the weight associated with the i^{th} input to unit j

$net_j = \sum_i w_{ji}x_{ji}$ (the weighted sum of inputs for unit j)

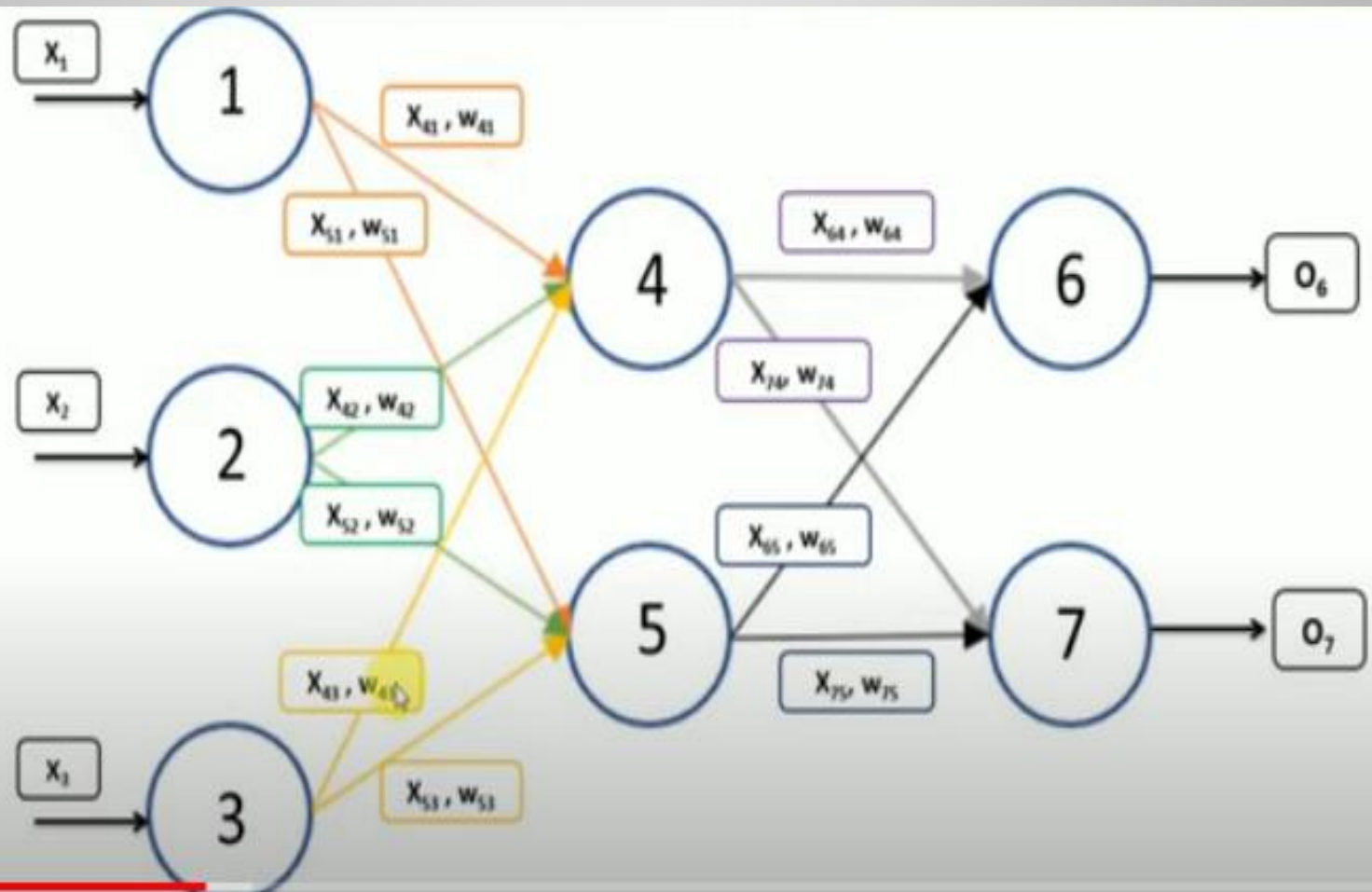
o_j = the output computed by unit j

t_j = the target output for unit j

σ = the sigmoid function

outputs = the set of units in the final layer of the network

Downstream(j) = the set of units whose immediate inputs include the output of unit j



$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

Where

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

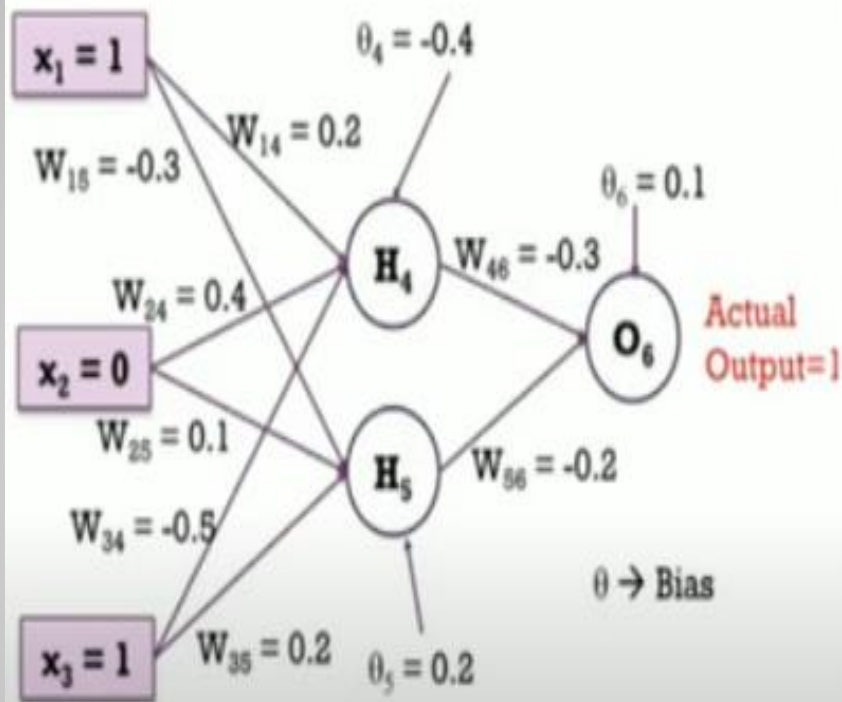
- To begin, notice that weight w_{ji} can influence the rest of the network only through net_j .

Therefore, we can use the chain rule to write,

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}}$$

$$net_j = \sum_i w_{ji} x_{ji}$$

$$\frac{\partial net_j}{\partial w_{ji}} = x_{ji}$$



Assume that the neurons have a sigmoid activation function, perform a forward pass and a backward pass on the network. Assume that the actual output of y is 1 and learning rate is 0.9. Perform another forward pass.

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

Where

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

- To begin, notice that weight w_{ji} can influence the rest of the network only through net_j .

Therefore, we can use the chain rule to write,

$$\begin{aligned} \frac{\partial E_d}{\partial w_{ji}} &= \frac{\partial E_d}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} \\ &= \frac{\partial E_d}{\partial net_j} x_{ji} \end{aligned}$$

$$net_j = \sum_i w_{ji} x_{ji}$$

$$\frac{\partial net_j}{\partial w_{ji}} = x_{ji}$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} x_{ji}$$

- Our remaining task is to derive a convenient expression for $\frac{\partial E_d}{\partial net_j}$

To derive a convenient expression for $\frac{\partial E_d}{\partial net_j}$

We consider two cases in turn:

- Case 1, where unit j is an output unit for the network, and
- Case 2, where unit j is an internal unit of the network.

Case 1: Training Rule for Output Unit Weights

- Just as wji can influence the rest of the network only through net_j , net_j can influence the network only through o_j . Therefore, we can invoke the chain rule again to write,

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial net_j} \quad \frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2 \quad \frac{\partial o_j}{\partial (net_j)} = \frac{\partial \sigma(net_j)}{\partial (net_j)} \quad \frac{\partial \sigma(x)}{\partial (x)} = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} (t_j - o_j)^2 \quad = \sigma(net_j) (1 - \sigma(net_j))$$

$$= \frac{1}{2} 2(t_j - o_j) \frac{\partial (t_j - o_j)}{\partial o_j}$$

$$= -(t_j - o_j)$$

$$= o_j (1 - o_j)$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j (1 - o_j)$$

Case 1: Training Rule for Output Unit Weights

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j(1 - o_j)$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} x_{ji}$$

$$\Delta w_{ji} = \eta \underline{(t_j - o_j) o_j(1 - o_j)} x_{ji}$$

$$\delta_j = (t_j - o_j) o_j(1 - o_j)$$

$$\Delta w_{ji} = \eta \delta_j x_{ji}$$

Case 2: Training Rule for Hidden Unit Weights

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}(j)} \frac{\partial \check{E}_d}{\partial net_k} \frac{\partial \check{net}_k}{\partial net_j}$$

$$= \sum_{k \in \text{Downstream}(j)} \boxed{-\delta_k} \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in \text{Downstream}(j)} -\delta_k \frac{\partial \check{net}_k}{\partial o_j} \frac{\partial o_j}{\partial net_j}$$

$$= \sum_{k \in \text{Downstream}(j)} -\delta_k \underline{w_{kj}} \frac{\partial o_j}{\partial net_j}$$

$$= \sum_{k \in \text{Downstream}(j)} -\delta_k w_{kj} o_j (1 - o_j)$$

$$\frac{\partial E_d}{\partial net_j} = -\underline{(t_j - o_j) o_j (1 - o_j)}$$

$$\delta_j = (t_j - o_j) o_j (1 - o_j)$$

$$\frac{\partial net_k}{\partial o_j} = \frac{\partial x_{kj} w_{kj}}{\partial o_j} = \frac{\partial o_j w_{kj}}{\partial o_j}$$

$$\frac{\partial o_j}{\partial (net_j)} = \frac{\partial \sigma(net_j)}{\partial (net_j)}$$

$$= \sigma(net_j) (1 - \sigma(net_j))$$

$$= o_j (1 - o_j)$$

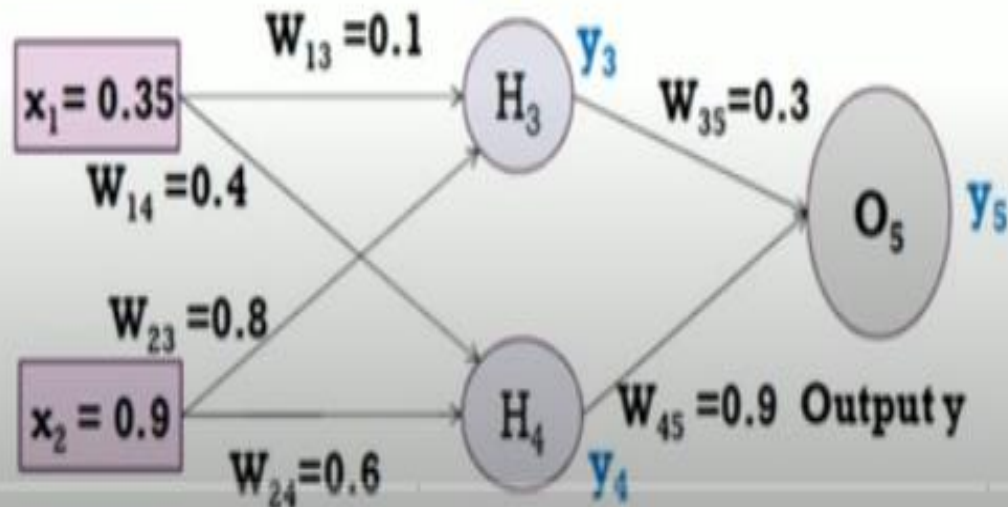
Case 2: Training Rule for Hidden Unit Weights

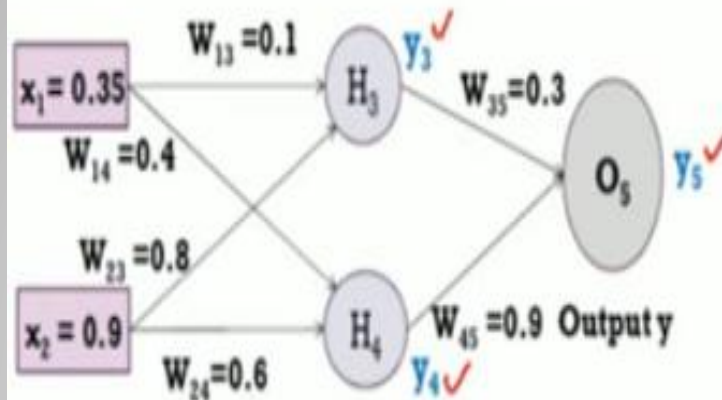
$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} x_{ji} \quad \frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} o_j(1 - o_j)$$

$$\Delta w_{ji} = \eta o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} x_{ji}$$

$$\Delta w_{ji} = \eta \delta_j x_{ji} \quad \delta_j = o_j(1 - o_j) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj}$$

- Assume that the neurons have a sigmoid activation function, perform a forward pass and a backward pass on the network. Assume that the actual output of y is 0.5 and learning rate is 1. Perform another forward pass.





- Forward Pass: Compute output for y_3 , y_4 and y_5 .

$$a_j = \sum_i (w_{ij} * x_i) \quad y_j = F(a_j) = \frac{1}{1 + e^{-a_j}}$$

$$a_1 = (w_{13} * x_1) + (w_{23} * x_2) \checkmark$$

$$= (0.1 * 0.35) + (0.8 * 0.9) = 0.755$$

$$y_3 = f(a_1) = 1 / (1 + e^{-0.755}) = 0.68$$

$$a_2 = (w_{14} * x_1) + (w_{24} * x_2) \checkmark$$

$$= (0.4 * 0.35) + (0.6 * 0.9) = 0.68$$

$$y_4 = f(a_2) = 1 / (1 + e^{-0.68}) = 0.6637$$

$$a_3 = (w_{35} * y_3) + (w_{45} * y_4) \checkmark$$

$$= (0.3 * 0.68) + (0.9 * 0.6637) = 0.801 \checkmark$$

$$y_5 = f(a_3) = 1 / (1 + e^{-0.801}) = 0.69 \text{ (Network Output)}$$

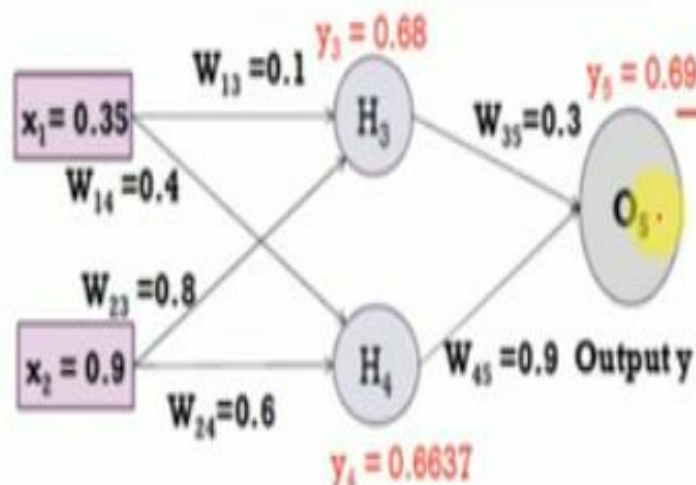
- Each weight changed by:

$$\Delta w_{ji} = \eta \delta_j o_i$$

$$\delta_j = o_j(1 - o_j)(t_j - o_j) \quad \text{if } j \text{ is an output unit}$$

$$\delta_j = o_j(1 - o_j) \sum_k \delta_k w_{kj} \quad \text{if } j \text{ is a hidden unit}$$

- where η is a constant called the learning rate
- t_j is the correct teacher output for unit j
- δ_j is the error measure for unit j



Backward Pass: Compute δ_3 , δ_4 and δ_5 .

For output unit:

$$\delta_5 = y(1-y) (y_{\text{target}} - y) \\ = 0.69 * (1 - 0.69) * (0.5 - 0.69) = -0.0406$$

For hidden unit:

$$\delta_3 = y_3(1-y_3) w_{35} * \delta_5 \\ = 0.68 * (1 - 0.68) * (0.3 * -0.0406) = -0.00265$$

$$\Delta w_{ji} = \eta \delta_j o_i$$

$$\delta_j = o_j(1-o_j)(t_j - o_j)$$

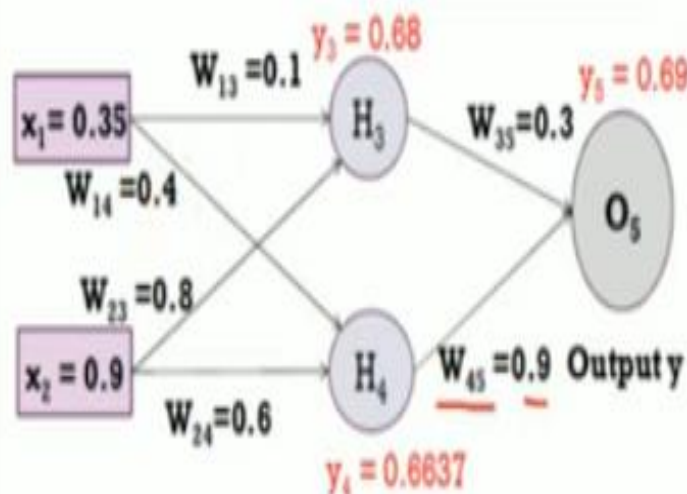
if j is an output unit

$$\delta_4 = y_4(1-y_4) w_{45} * \delta_5$$

$$= 0.6637 * (1 - 0.6637) * (0.9 * -0.0406) = -0.0082$$

$$\delta_j = o_j(1-o_j) \sum_k \delta_k w_{kj}$$

if j is a hidden unit



• Backward Pass: Compute δ_3 , δ_4 and δ_5 .

For output unit:

$$\delta_5 = y(1-y) (y_{\text{target}} - y) \\ = 0.69 * (1 - 0.69) * (0.5 - 0.69) = -0.0406$$

For hidden unit:

$$\delta_3 = y_3(1-y_3) w_{35} * \delta_5 \\ = 0.68 * (1 - 0.68) * (0.3 * -0.0406) = -0.00265$$

Compute new weights

$$\Delta w_{ji} = \eta \delta_j o_i$$

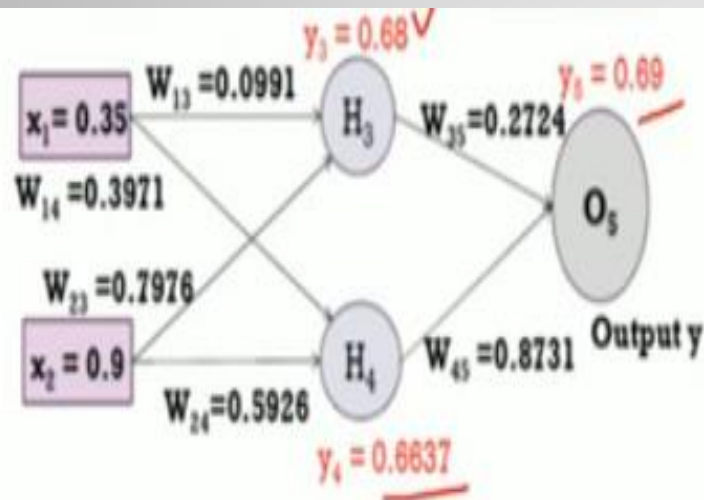
$$\Delta w_{45} = \eta \delta_5 y_4 = 1 * -0.0406 * 0.6637 = -0.0269 \\ w_{45}(\text{new}) = \Delta w_{45} + w_{45}(\text{old}) = -0.0269 + (0.9) = 0.8731$$

$$\delta_4 = y_4(1-y_4) w_{45} * \delta_5 \\ = 0.6637 * (1 - 0.6637) * (0.9 * -0.0406) = -0.0082$$

$$\Delta w_{14} = \eta \delta_4 x_1 = 1 * -0.0082 * 0.35 = -0.00287 \\ w_{14}(\text{new}) = \Delta w_{14} + w_{14}(\text{old}) = -0.00287 + 0.4 = 0.3971$$

- Similarly, update all other weights

i	j	w_{ij}	δ_i	x_i	η	Updated w_{ij}
1	3	0.1	-0.00265	0.35	1	0.0991
2	3	0.8	-0.00265	0.9	1	0.7976
1	4	0.4	-0.0082	0.35	1	0.3971
2	4	0.6	-0.0082	0.9	1	0.5926
3	5	0.3	-0.0406	0.68	1	0.2724
4	5	0.9	-0.0406	0.6637	1	0.8731



$$\text{Error} = y_{\text{target}} - y_5 = -0.182$$

- Forward Pass: Compute output for y_3 , y_4 and y_5 .

$$a_j = \sum_i (w_{ij} * x_i) \quad y_j = F(a_j) = \frac{1}{1 + e^{-a_j}}$$

$$\begin{aligned} a_1 &= (w_{13} * x_1) + (w_{23} * x_2) \\ &= (0.0991 * 0.35) + (0.7976 * 0.9) = 0.7525 \\ y_3 &= f(a_1) = 1 / (1 + e^{-0.7525}) = \underline{0.6797} \end{aligned}$$

$$\begin{aligned} a_2 &= (w_{14} * x_1) + (w_{24} * x_2) \\ &= (0.3971 * 0.35) + (0.5926 * 0.9) = 0.6723 \\ y_4 &= f(a_2) = 1 / (1 + e^{-0.6723}) = \underline{0.6620} \end{aligned}$$

$$\begin{aligned} a_3 &= (w_{35} * y_3) + (w_{45} * y_4) \\ &= (0.2724 * 0.6797) + (0.8731 * 0.6620) = 0.7631 \\ y_5 &= f(a_3) = 1 / (1 + e^{-0.7631}) = \underline{0.6820} \text{ (Network Output)} \end{aligned}$$

