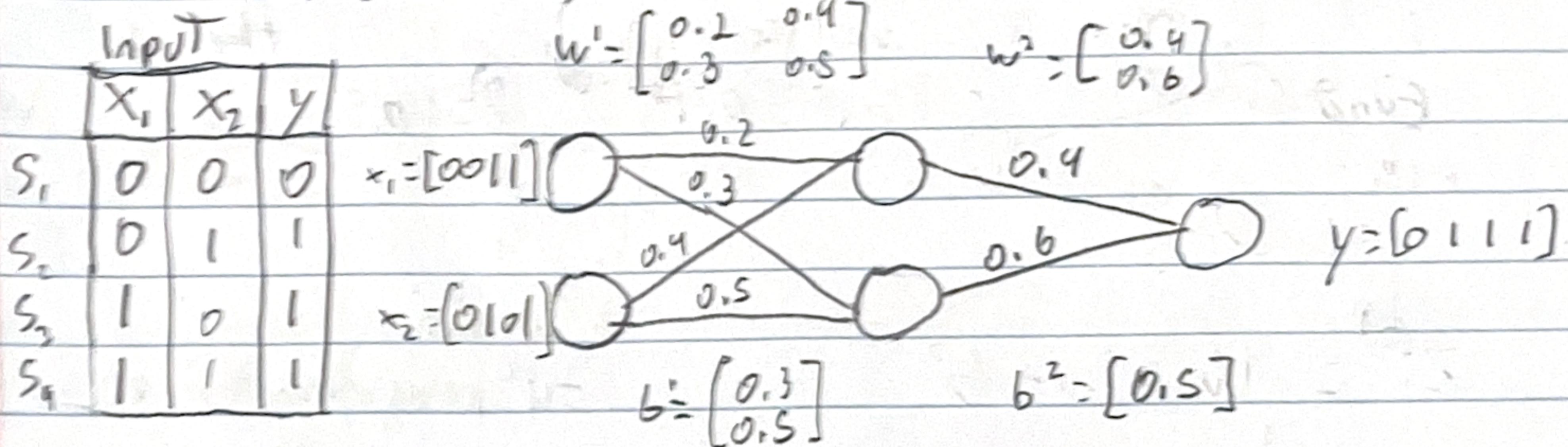


## 2. Math problem.

Compute the forward passes and gradient manually for a neural network with 1 hidden layer (w/ 2 neurons) and 1 output neuron.



\* input layer size is determined by number of features

$$Z^{(0)} = Wx + b$$

$$Z^{(0)} = \begin{bmatrix} 0.2 & 0.4 \\ 0.3 & 0.5 \end{bmatrix} \cdot \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix} + \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}$$

$$= (0.2 \cdot 0 + 0.4 \cdot 0) \quad (0.2 \cdot 0 + 0.4 \cdot 1) \quad (0.2 \cdot 1 + 0.4 \cdot 0) \quad (0.2 \cdot 1 + 0.4 \cdot 1) \\ (0.3 \cdot 0 + 0.5 \cdot 0) \quad (0.3 \cdot 0 + 0.5 \cdot 1) \quad (0.3 \cdot 1 + 0.5 \cdot 0) \quad (0.3 \cdot 1 + 0.5 \cdot 1)$$

$$Z^{(1)} = \begin{bmatrix} 0 & 0.4 & 0.2 & 0.6 \\ 0 & 0.5 & 0.3 & 0.8 \end{bmatrix} + \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.3 & 0.7 & 0.5 & 0.9 \\ 0.5 & 1.0 & 0.8 & 1.3 \end{bmatrix}$$

$$a^{(0)} = \sigma(Z^{(0)}) \text{ using sigmoid} = \frac{1}{1 + e^{-Z}} = \begin{bmatrix} \sigma(0.3) & \sigma(0.7) & \sigma(0.5) & \sigma(0.9) \\ \sigma(0.5) & \sigma(1.0) & \sigma(0.8) & \sigma(1.3) \end{bmatrix}$$

$$\sigma(0.3) = \frac{1}{1 + e^{-0.3}} \approx 0.574$$

$$a^{(1)} = \begin{bmatrix} 0.574 & 0.668 & 0.622 & 0.711 \\ 0.622 & 0.731 & 0.690 & 0.786 \end{bmatrix}$$

$$\sigma(0.7) = \frac{1}{1 + e^{-0.7}} \approx 0.668$$

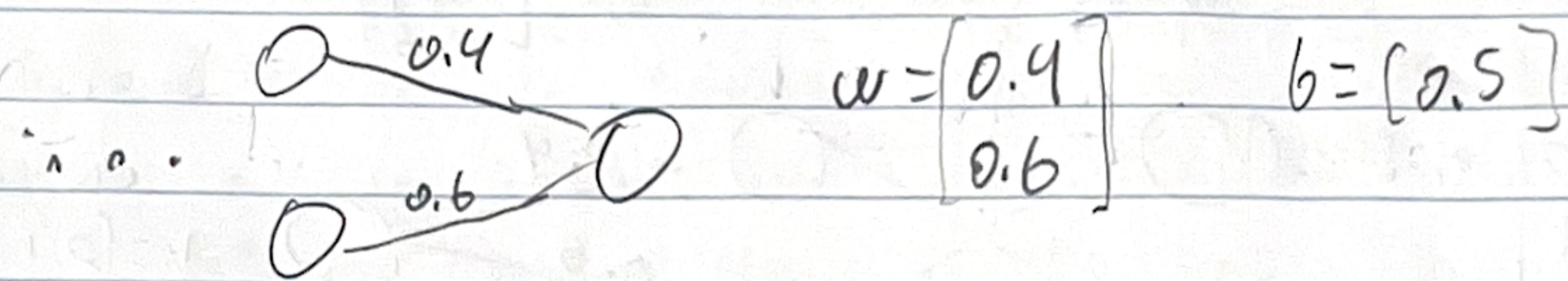
$$x_1 = [0011] \quad x_2 = [0101]$$

$$\sigma(0.5) = \frac{1}{1 + e^{-0.5}} \approx 0.622$$

$$x_1 = [0011] \quad x_2 = [0101]$$

$$\sigma(0.9) = \frac{1}{1 + e^{-0.9}} \approx 0.696$$

$$\alpha^{(1)} = \begin{bmatrix} 0.574 & 0.668 & 0.622 & 0.711 \\ 0.622 & 0.731 & 0.690 & 0.786 \end{bmatrix}$$



$$Z^{(2)} = W^2 X^2 + b^2$$

$$= \begin{bmatrix} 0.4 & 0.6 \end{bmatrix}_{1 \times 2} \cdot \alpha^{(1)}_{2 \times 4} + 0.5$$

$$= (0.4 \cdot 0.574 + 0.6 \cdot 0.622) (0.4 \cdot 0.668 + 0.6 \cdot 0.711) - (0.4 \cdot 0.622 + 0.6 \cdot 0.690) (0.4 \cdot 0.711 + 0.6 \cdot 0.786)$$

$$= [0.603 \quad 0.706 \quad 0.663 \quad 0.756] + [0.5]$$

$$Z^{(2)} = [1.103 \quad 1.206 \quad 1.163 \quad 1.256]$$

• Apply Sigmoid

$$\frac{1}{1+e^{-0.103}} \approx 0.752 \quad \frac{1}{1+e^{-0.206}} \approx 0.770 \quad \frac{1}{1+e^{-0.163}} \approx 0.762 \quad \frac{1}{1+e^{-0.256}} \approx 0.778$$

$$\alpha^{(2)} = [0.752 \quad 0.770 \quad 0.762 \quad 0.778]$$

• Now we can compute the loss using mean square error.

$$C = \frac{1}{2} (\alpha^{(2)} - Y_{true})^2 \quad Y = [0 \ 1 \ 1 \ 1]$$

$$\therefore \begin{array}{l} (0.752 - 0)^2 \\ 0.565 \end{array} + \begin{array}{l} (0.770 - 1)^2 \\ 0.059 \end{array} + \begin{array}{l} (0.762 - 1)^2 \\ 0.056 \end{array} + \begin{array}{l} (0.778 - 1)^2 \\ 0.049 \end{array} = 0.123$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (a_i^{(2)} - y_i)^2$$

$$\text{loss} = \frac{1}{2n} \sum_{i=1}^n (a_i^{(2)} - y_i)^2$$

$$MSE = \frac{1}{4} \times 0.723 \approx [0.181], \text{ loss} = [0.090]$$

- To simplify derivation we modify MSE to make it easier to do back propagation

Backpropagation using Loss Function =  $\frac{1}{2n} \sum (a^{(2)} - y_i)^2$

calculate the output layer error forward by  $\delta^{(2)}$  and  $\delta^{(1)}$

$$\begin{aligned}\delta^{(2)} &= (a^{(2)} - y_{\text{true}}) \cdot \sigma'(z^{(2)}) && \text{derivative of sigmoid} \\ &= [0.752 \ 0.770 \ 0.762 \ 0.778] - \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \\ &= [0.752 \ -.230 \ -.238 \ -.222] \cdot \sigma(z) \cdot (1 - \sigma(z))\end{aligned}$$

$$\begin{aligned}&= [0.752 \ -.230 \ -.238 \ -.222] \cdot \begin{bmatrix} 0.752 \\ 0.770 \\ 0.762 \\ 0.778 \end{bmatrix} \cdot \begin{bmatrix} 1 + .752 = 0.248 \\ 1 - 770 = 0.230 \\ 1 - 762 = 0.238 \\ 1 - 778 = 0.222 \end{bmatrix}\end{aligned}$$

$$= [1.566 \ -.171 \ -.181 \ -.173] \cdot \begin{bmatrix} 0.248 \\ 0.230 \\ 0.238 \\ 0.222 \end{bmatrix}$$

$$\delta^{(2)} = [.140 \ -.091 \ -.043 \ -.039] \cdot \begin{bmatrix} 0.248 \\ 0.230 \\ 0.238 \\ 0.222 \end{bmatrix}$$

- AFTER reviewing, I realized that I should have used the Hadamard Product rather than doing regular matrix multiplication  
- Luckily since matrices here were of a single dimension, the result is the same.

Now The hidden layer error for  $\delta^{(1)}$

$$l = 0.579 \dots$$

$$\begin{aligned}\delta^{(1)} &= (W^{(2)})^\top \cdot \delta^{(2)} \cdot \sigma(z^{(1)}) \cdot (1 - \sigma(z^{(1)})) \\ &= [0.140 \ -0.91 \ -0.43 \ -0.38] \cdot \begin{bmatrix} 0.9 \\ 0.6 \end{bmatrix} \cdot \begin{bmatrix} 0.574 & 0.668 & 0.622 & 0.711 \\ 0.622 & 0.731 & 0.690 & 0.219 \end{bmatrix} \\ &= \begin{bmatrix} 0.9 \\ 0.6 \end{bmatrix} \cdot [0.140 \ -0.91 \ -0.43 \ -0.38]^\top \cdot \sigma(z^{(1)}) \cdot \begin{bmatrix} 0.426 & 0.332 & 0.378 & 0.269 \\ 0.378 & 0.269 & 0.310 & 0.214 \end{bmatrix} \\ &\quad \boxed{\delta^{(1)}} = \boxed{l \times 4} \\ &= \frac{0.9 \times 0.140}{0.6 + 0.140} \ \frac{0.9 \times (-0.91)}{0.6 + (-0.91)} \ \frac{0.9 \times (-0.43)}{0.6 + (-0.43)} \ \frac{0.9 \times (-0.38)}{0.6 + (-0.38)}\end{aligned}$$

$$\textcircled{1} = \begin{bmatrix} -0.056 & -0.016 & -0.017 & -0.015 \\ -0.084 & -0.025 & -0.026 & -0.023 \end{bmatrix}.$$

~~AMM~~

~~MM~~ 9.07

\* We can now multiply these 3 matrices ( $W^{(2)}$ ),  $\sigma(z^{(1)})$ , and  $(1 - \sigma(z^{(1)}))$  using the hadamard product. First I will 2 and 3. Then multiply that by 1.

$$\textcircled{1} \times \textcircled{2} \times \textcircled{3} = \begin{bmatrix} (-0.056 \times 0.574 \times 0.426) & (-0.056 \times 0.668 \times 0.332) \dots \\ (-0.084 \times 0.622 \times 0.378) & (-0.084 \times 0.731 \times 0.269) \dots \end{bmatrix}$$

$$\delta^{(1)} = \begin{bmatrix} -0.0135 & -0.0036 & -0.0040 & -0.0031 \\ -0.0197 & -0.0048 & -0.0055 & -0.0038 \end{bmatrix}$$

• Calculating the gradients for output layer and hidden layer.

Now that we have the loss functions for each layer's weights, calculate the gradient and update the bias.

$$\frac{\partial \text{Cost}}{\partial w^{(2)}} = \frac{1}{n} \delta^{(1)} (\delta^{(2)})^T$$

$$= \begin{bmatrix} 0.47 \\ 0.6 \end{bmatrix} - \frac{1}{4}$$

$$= \frac{1}{4} \begin{bmatrix} -0.00079 \\ -0.00243 \end{bmatrix}$$

$$\text{gradient}_{w^{(2)}} = \begin{bmatrix} -0.0005 \\ -0.0000 \end{bmatrix}$$

Loss =  $\sum_{j=1}^n \delta_j^{(2)} \cdot (y_j - \hat{y}_j)$

$$\frac{\partial \text{Loss}}{\partial w^{(2)}} = \begin{bmatrix} 0.140 \\ -0.041 \\ -0.043 \\ -0.038 \end{bmatrix}$$

$$= \frac{1}{4} \begin{bmatrix} 0.018 \\ 0.0045 \end{bmatrix}$$

\* output max to match slope of weights.

To update the gradient for biases  $b^{(2)}$  we use  $\frac{\partial \text{Cost}}{\partial b^{(2)}} = \frac{1}{n} \sum_{j=1}^n \delta_j^{(2)}$ .

$$\sum [0.140 - 0.041 - 0.043 - 0.038] = \frac{0.018}{4} = \boxed{0.0045 = b^{(2)} \text{ gradient}}$$

$$\frac{\partial \text{Cost}}{\partial w^{(1)}} = \frac{1}{n} \delta^{(1)} \cdot X^T$$

$$= \delta = \begin{bmatrix} 0.135 & -0.0036 & 0.0040 & 0.0031 \\ 0.0197 & 0.0048 & 0.0055 & 0.0038 \end{bmatrix}$$

$$= \frac{1}{4} \begin{bmatrix} -0.0071 & -0.0067 \\ -0.0093 & -0.0066 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$\frac{\partial \text{Cost}}{\partial w^{(1)}} = \begin{bmatrix} -0.0018 & -0.0017 \\ -0.0023 & -0.0022 \end{bmatrix}$$

$$\frac{\partial \text{Cost}}{\partial b^{(1)}} = \frac{1}{n} \sum \delta^{(1)} = \frac{(0.135 + -0.0036 + 0.0040 + 0.0031) + (0.0197 + 0.0048 + 0.0055 + 0.0038)}{4} = \begin{bmatrix} 0.003 \\ 0.0050 \end{bmatrix} \frac{1}{4} = \boxed{\begin{bmatrix} 0.0008 \\ 0.0014 \end{bmatrix}}$$

Update The parameters =  $\theta_{new} = \theta_{old} - lr \times \frac{\partial L(\theta)}{\partial \theta}$

$$w^{(2)} \leftarrow w^{(1)} - lr \frac{\partial L}{\partial w^{(2)}}$$

$$lr = 0.5$$

$$w_{new}^{(2)} = \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix} - 0.5 \times \begin{bmatrix} -0.0002 \\ -0.0006 \end{bmatrix} = \begin{bmatrix} 0.4001 \\ 0.6003 \end{bmatrix}$$

$$b^{(2)} \leftarrow b^{(1)} - lr \frac{\partial L}{\partial b^{(2)}}$$

$$= 0.5 - 0.5 \times 0.0005 = 0.4978$$

$$w^{(1)} \leftarrow w^{(0)} - lr \frac{\partial L}{\partial w^{(1)}}$$

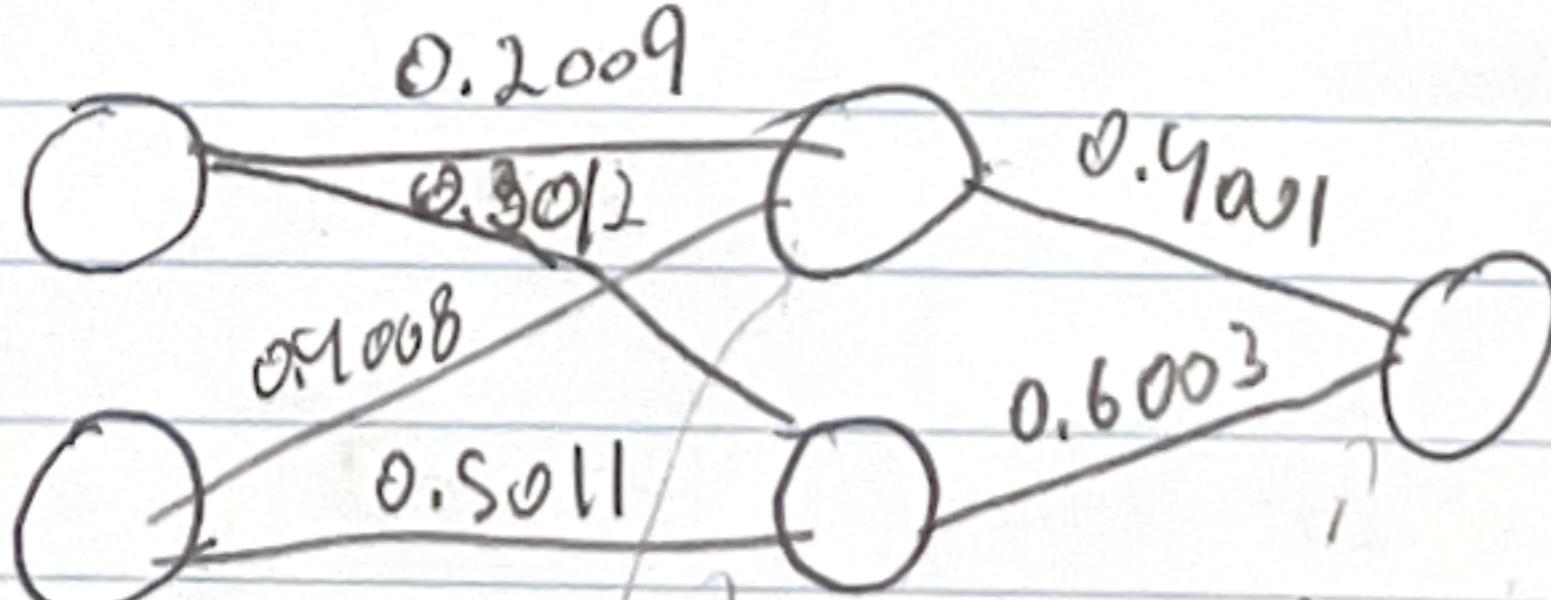
$$w^{(0)} = \begin{bmatrix} 0.2 & 0.4 \\ 0.3 & 0.5 \end{bmatrix} - 0.5 \times \begin{bmatrix} -0.00178 & -0.00017 \\ -0.0023 & -0.0022 \end{bmatrix} = \begin{bmatrix} 0.2009 & 0.4008 \\ 0.3012 & 0.5011 \end{bmatrix}$$

$$b^{(1)} \leftarrow b^{(0)} - lr \frac{\partial L}{\partial b^{(1)}}$$

$$= \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix} - 0.5 \times \begin{bmatrix} 0.0008 \\ 0.0014 \end{bmatrix} = \begin{bmatrix} 0.2996 \\ 0.4993 \end{bmatrix}$$

Read for another forward pass.

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1



$$bias = [0.4978]$$

$$bias = \begin{bmatrix} 0.3 \\ 0.5 \end{bmatrix}$$