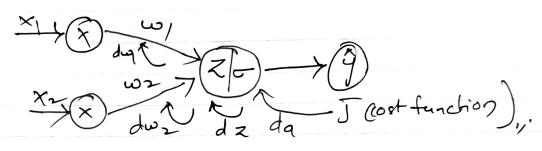# Question-1

For a 2 layer Neural Network with 1 hidden layer and 1 output layer with sigmoid as a activation function but in the hidden layer can be pictured as



Here is the steps for forward propagation

$$Z = W^T X + b.$$

we can derive forward propagation as

$$z^1 = W^{(1)} x + b^1$$
$$A^1 = \sigma(z^1) \qquad \sigma \text{ is sigmoid.}$$
$$z^2 = W^{(2)} A^{(1)} + b^1$$
$$A^2 = \sigma(z^2).$$

Gradient descent to update the weights
$$W = W1 - \alpha \frac{dj}{dw}.$$

Gradient descent can be calculated using backpropagation by applying 2nd order derivative. We need to pass output of the forward propagation to the back propagation and calculate the error which is ~~dy~~

$$z^2 = |w|^2 A^1 + b^2$$

derivative of $z^2$ w.r.t $w^2$

$$\frac{dj}{dw^2} = \left[ -y + y\cancel{A^2} + A^2 - \cancel{yA^2} \right]\left[ A^1 \right]$$

$$= \left[ A^2 - y \right]\left[ A^1 \right]$$

$$= dz^2 A^1$$

$$\frac{dj}{db^2} = \frac{dj}{dA^2} \cdot \frac{dA^2}{dz^2} \cdot \frac{dz^2}{db^2}$$

$$= dz^2$$

$$\frac{dj}{dw^1} = \frac{dj}{dA^2} \cdot \frac{dA^2}{dz^2} \cdot \frac{dz^2}{dA^1} \cdot \frac{dA^1}{dz^1} \cdot \frac{dz^1}{dw^1}$$

$$= \frac{d\hat{s}}{dz^2} \cdot \frac{dz^2}{dA^1} \cdot \frac{dA^1}{dz^1} \cdot \frac{dz^1}{dw^1}$$

$$= \left[ A^2 - y \right] \left[ w^2 \right] \left[ g'(z^1) \right] \left[ A^0 \right]$$

$$= dz^2 w^2 g'(z) A^0$$

$$= dz^1 A^0$$

$$\frac{dz}{dA^1} = \frac{d\hat{s}}{dA^2} \cdot \frac{dA^2}{dz^2} \cdot \frac{dz^2}{dA^1} \Rightarrow \frac{d\hat{s}}{dz^2} w^2 \Rightarrow dz^2 \cdot w^2$$

$$\frac{d\hat{s}}{db^1} = \frac{d\hat{s}}{dA^2} \cdot \frac{dA^2}{dz^2} \cdot \frac{dz^2}{dA^1} \cdot \frac{dA^1}{dz_1} \cdot \frac{dz^1}{db^1}$$

$$= \frac{d\hat{s}}{dz^2} \cdot \frac{dz^2}{dA^1} \cdot \frac{dA^1}{dz_1} \cdot \frac{dz^1}{db^1}$$

$$= \left[ A^2 - y \right] \left[ w^2 \right] \left[ g'(z^1) \right] \left[ 1 \right]$$

$$= dz^2 w^2 g'(z^1) = dz^1.$$

Now, we need to add activation function for every derivative of backpropagation.

$$\frac{\partial G}{\partial x} = \frac{\partial G}{\partial y} f'(\cdot x) \quad \text{where } f'(x) \text{ is derived}$$

from every neuron $\Rightarrow y = [f(x^1) + f(x^2) \cdots f(x^i)]$

$$= f(x)_{\prime\prime}'$$

Sigmoid activation function:-

$$\sigma(x) = 1/1 + e^{-x}$$

after derivating this function

$$\sigma'(x) = (1 + e^{-x})(1 - (1 + e^{-x})).$$

which can be termed as
$$\sigma'(x) = x(1 - x)$$

Now by updating weights for every layer

we need to calculate error using mean squared error.

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

→ If we use sigmoid in the binary classification problem it will be hard to interpret the result as sigmoid will return log values and binary classification will always expect binary value of either '0' or '1.'

Sigmoid will always gives probability ranging from 0 to 1.