

## Forward Propagation

As the name suggests, the input data is fed in the forward direction through the network. Each hidden layer accepts the input data, processes it as per the activation function and passes to the successive layer.

In order to generate some output, the input data should be fed in the forward direction only. The data should not flow in reverse direction during output generation otherwise it would form a cycle and the output could never be generated.

## Backward Propagation

Back Propagation is the essence of neural net training. It is the practice of fine-tuning the weights of a neural net based on the error rate (i.e. loss) obtained in the previous epoch (i.e. iteration). Proper tuning of the weights ensures lower error rates, making the model reliable by increasing its generalisation. Backpropagation is a short form for "backward propagation of errors". It is a standard method of training artificial neural networks. This method helps to calculate the gradient of a loss function with respect to all the weights in the network.

## Question 2: (Vectorised Implementation)

$$\begin{array}{lll} X \in \mathbb{R}^{1 \times 4} & & \} \text{Input} \\ W_1 \in \mathbb{R}^{5 \times 4} & b_1 \in \mathbb{R}^{5 \times 1} & \} \text{Parameters} \\ W_2 \in \mathbb{R}^{1 \times 5} & b_2 \in \mathbb{R}^{1 \times 1} & \\ Y \in \mathbb{R}^{1 \times 1} & & \} \text{Output} \end{array}$$

### Forward Propagation:

# Assume weights and bias are initialized

$$Z_1 = (W_1) \cdot (X^T) + b_1$$

$$A_1 = \text{sigmoid}(Z_1)$$

$$Z_2 = (W_2) \cdot (A_1^T) + b_2$$

$$A_2 = \text{sigmoid}(Z_2)$$

$$\text{Cost} = -Y \log A_2 - (1-Y) \log (1-A_2)$$

### Backward Propagation:

# Say  $dZ$  denotes  $d(\text{cost})/d(Z)$

$$dZ_2 = A_2 - Y$$

$$dW_2 = (dZ_2) \cdot (A_1^T)$$

$$dB_2 = dZ_2$$

$$dZ_1 = [(W_2^T) \cdot (dZ_2)] * (A_1 * (1-A_1))$$

$\text{sigmoid}'(Z_1)$

$$dW_1 = (dZ_1) \cdot (X^T)$$

$$dB_1 = dZ_1$$



## Question 2: (For general MLP)

$$X \in \mathbb{R}^{n[0] \times m}$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$W[l] \in \mathbb{R}^{n[l] \times n[l-1]}$$

$$b[l] \in \mathbb{R}^{n[l] \times m}$$

$n[l]$ : number of features in layer  $l$

$m$ : number of examples

$L$ : Number of layers

### Forward Propagation:

# Assume  $X = A[0]$

for  $l$  in range( $L$ ):

~~$Z[l+1] = W[l+1] \cdot A[l] + b[l+1]$~~

$$Z[l+1] = W[l+1] \cdot A[l] + b[l+1]$$

$$A[l+1] = \text{Sigmoid}(Z[l+1])$$

{ store  $Z[l+1]$  and  $A[l+1]$  }

$$\text{Cost} = \frac{1}{m} \sum_{i=1}^m \left\{ -Y^{(i)} \log(A[L]^{(i)}) - (1 - Y^{(i)}) \log(1 - A[L]^{(i)}) \right\}$$

### Backward Propagation:

$$dZ[l+1] = A[l+1] - Y, \quad dW[l+1] = (dZ[l+1]) \cdot (A[l]^\top), \quad db[l+1] = dZ[l+1]$$

for  $l$  from  $L-1$  to  $0$ :

$$dZ[l] = (W[l+1]^\top \cdot dZ[l+1]) * (A[l] * (1 - A[l]))$$

$$dW[l] = (dZ[l]) \cdot (A[l-1]^\top)$$

$$db[l] = dZ[l]$$

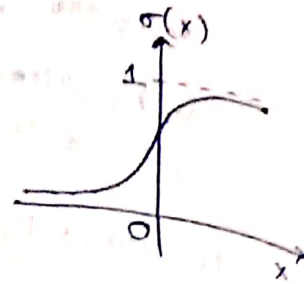
Store all  $dW[l]$  and  $db[l]$

### Question 3

(a) Sigmoid:

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

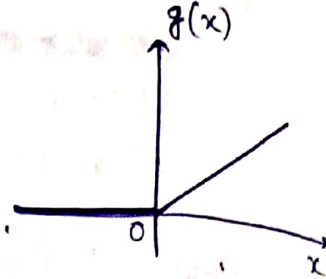
$$\sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \sigma(x) \cdot (1-\sigma(x))$$



(b) Relu:

$$g(x) = \max(0, x)$$

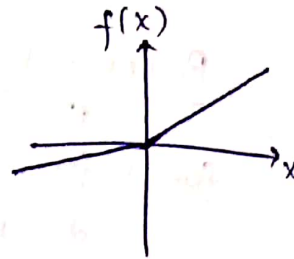
$$g'(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$



(c) Leaky Relu:

$$f(x) = \max(0.01x, x)$$

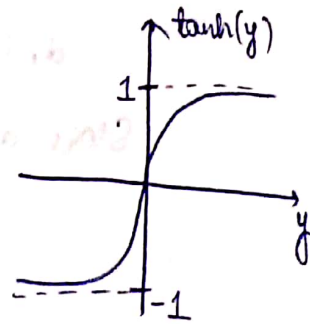
$$f'(x) = \begin{cases} 0.01 & x < 0 \\ 1 & x \geq 0 \end{cases}$$



(d) Tanh:

$$\tanh(y) = \frac{e^y - e^{-y}}{e^y + e^{-y}}$$

$$\tanh'(y) = \frac{4}{(e^y + e^{-y})^2} = 1 - (\tanh(y))^2$$



(e) Softmax:

$$\sigma(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \quad \forall i=1, 2, \dots, k \text{ and } \mathbf{z} = (z_1, z_2, \dots, z_k) \in \mathbb{R}^k$$

$$\frac{\partial \sigma(z_i)}{\partial z_i} = \sigma(z_i) \cdot (1 - \sigma(z_i))$$

$$\frac{\partial \sigma(z_i)}{\partial z_j} = -\sigma(z_i) \cdot \sigma(z_j) \quad \forall i \neq j$$