

1. Forward propagation

In this process we move in forward path from input to the output. In this process we have different hidden layers and one output layer. Each hidden layer consist of many hidden units which are independent to each other and are the outcome of the activation function applied to the some linear function of the previous hidden layer. By calculating these equations we find our predicted output from given input.

Backward propagation

This process is exactly opposite to the above process. Using this we find gradients (or more precisely some refinement term) from the loss function which is the net error in computing predicted output. So, we tune our all parameters to ensure less cost (less error). For this we use different optimizers. Once we tuned our parameters, then finally we use it to predict output for the test input data.

2.

Input $x \{4 \times 1\}$

Output $y \{1 \times 1\}$

$W_1 \{5 \times 4\} \quad b_1 \{5 \times 1\}$] hidden unit parameters
 $W_2 \{1 \times 5\} \quad b_2 \{1 \times 1\}$]

g_1 activation function for layer 1
 g_2 " " layer 2

Forward prop. (Vectorised)

$$z_1 = w_1 \cdot x + b_1$$

$$A_1 = g_1(z_1)$$

$$z_2 = w_2 \cdot A_1 + b_2$$

$$A_2 = g_2(z_2)$$

$$\text{cost} = -y \log A_2 - (1-y) \log(1-A_2)$$

Backward prop. (Vectorised)

$$dz_2 = A_2 - y \quad \left\{ \text{If } g_2 \text{ is sigmoid} \right\}$$

$$dw_2 = dz_2 \cdot A_1$$

$$db_2 = dz_2$$

$$dA_1 = w_2^T \cdot dz_2$$

$$dz_1 = dA_1 * g_1'(z_1) \quad (\text{element-wise product})$$

$$dw_1 = dz_1 \cdot \cancel{A_0} X$$

$$db_1 = dz_1$$

Generalization for MLP: Forward prop.

consider x as $A^{[0]}$.

For every layer l in $[1, L]$

$$\begin{aligned} z^{[l]} &= w^{[l]} \cdot A^{[l-1]} + b^{[l]} && \left\{ \text{store } A^{[l]} \text{ and } z^{[l]} \right\} \\ A^{[l]} &= g^{[l]}(z^{[l]}) && \text{for Back prop.} \end{aligned}$$

Generalization for MLP: Backward prop.

~~$$dz^{[L]} = A^{[L]} - y$$~~

$$dW^{[L]} = dz^{[L]} \cdot A^{[L-1]}$$

$$db^{[L]} = dz^{[L]}$$

For $(L-1)$ to 1 s- ($l \in [1, \dots, L-1]$)

$$dz^{[l]} = dA^{[l]} * g^{[l-1]}(z^{[l]})$$

$$dW^{[l]} = dz^{[l]} \cdot A^{[l-1]}$$

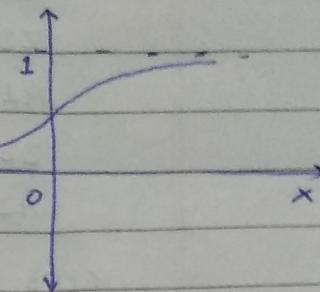
$$db^{[l]} = dz^{[l]}, \quad dA^{[L-1]} = w^{[L]} T \cdot dz^{[L]}$$

3. a) Sigmoid

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

$$\frac{dy}{dx} = y(x)(1-y(x))$$

$$y(x)$$

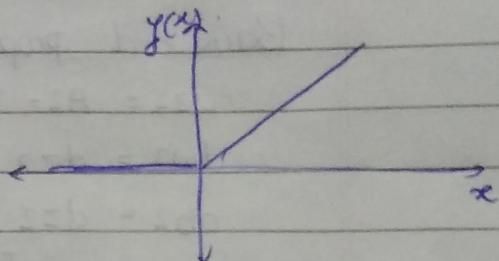


b) Relu

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

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

$$y(x)$$

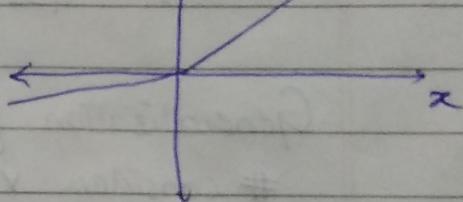


c) Leaky Relu

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

$$\frac{dy}{dx} = \begin{cases} 0.01 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

$$y(x)$$

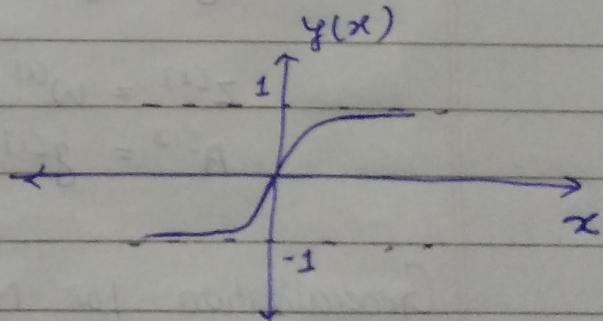


d) Tanh

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

$$\frac{dy}{dx} = 1 - y^2(x)$$

$$y(x)$$



e) Softmax

$Z \sim \text{column vector / Vector}$

$$y(z[i]) = \frac{e^{z[i]}}{\sum_{i=1}^K e^{z[i]}}$$

$K \sim \text{size of } Z \text{ vector}$

$$\frac{dy}{dz[i]} = y(z[i])(1 - y(z[i]))$$