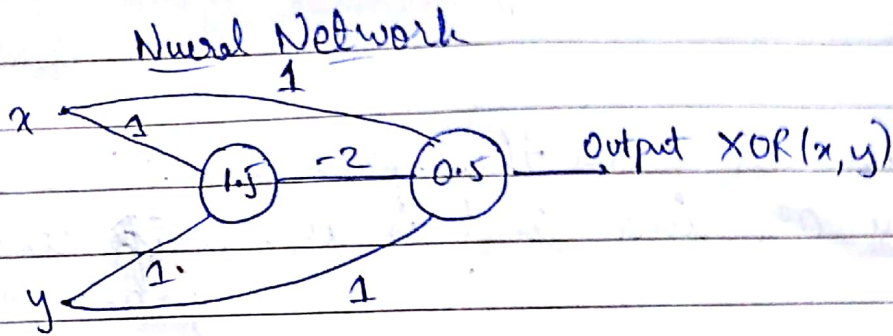


Theory - Ass - 3

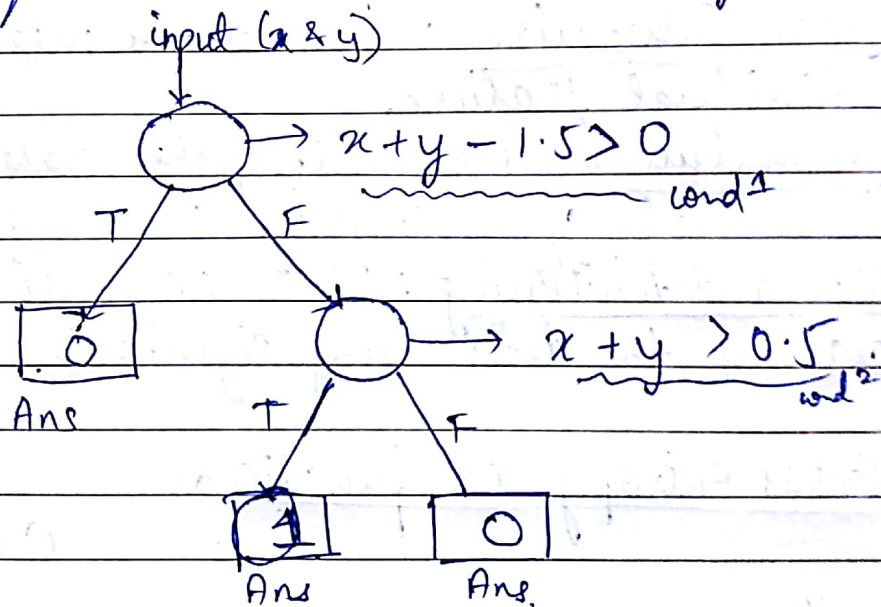
Shabham Khanna

1 YES



Eq of network $x + y - 2(x + y - 1.5) - 0.5 > 0$

The above equation can be also be made using a decision tree.



Both the models represent the same ~~neural~~ mathematical equation.

Sigmoid saturate and kill gradients. The problem is if the sigmoid neuron's activation saturates at either tail of 0 or 1, the gradient at the regions is almost zero, thus during backpropagation, the small gradient will kill itself.

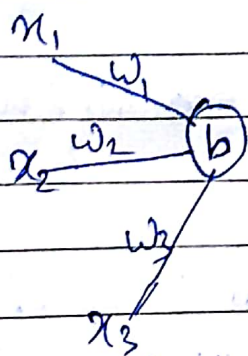
Sigmoid outputs are not zero-centered. This creates problems on the dynamics during gradient decent.

ReLU does not have these problems since it doesn't kill the gradient and is also ~~not~~ zero-centred.

We can use some Pre-processing techniques to counter the problem.

- Mean Subtraction: Subtract the mean across every individual feature.
- Normalization: Normalising the data dimensions so they scale approximately the same way.
- PCA and Whitening: This normalises the data, after centering the data using eigenbasis.

3. Cross entropy loss function.



$$C = -\frac{1}{n} \sum_x y \ln a + (1-y) \ln (1-a)$$

- $n \rightarrow$ no. of items of training data
- $x \rightarrow$ input
- $y \rightarrow$ o/p, $\sigma \rightarrow$ o/p form of Neuron

$$\frac{\partial C}{\partial w_i} = -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \frac{\partial \sigma(z)}{\partial w_i}$$

$$= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \sigma'(x) x_j$$

$$\frac{dC}{dw_j} = \frac{1}{n} \sum_x \frac{\sigma'(z) x_j}{\sigma(z)(1-\sigma(z))} (\sigma(z) - y) = \frac{1}{n} \sum_x (\sigma(z) - y) x_j$$

Thus we observe that the rate of learning is dependent on $(\sigma(z) - y)$ i.e. the error of the output. It avoids the learning slowdown caused by $\sigma'(z)$ in the quadratic cost function.