Shubham Kumar

2015098

HW-3

A1= No, we cannot build a neural network of arbitary length to learn XOR.

XOR is ~~linear~~ non-linearly seprable. Neural network using just linear activation function is a generalized linear model

∴ XOR can't be classified by it.

$$y = f(w_2 (w_1 x + b_1) + b_2)$$

Its same as SVM with linear function.

$$f(x) = B(Ax + a) + b$$
$$= \underbrace{BA}_{\downarrow} x + \underbrace{Ba + b}_{\downarrow}$$

Constant    Constant.

∴ $f(x) = C_1 x + C_2$

A2 = **Sigmoid**

Reasons :-

<u>Sigmoid Saturation</u> , <u>Killing Gradients</u>

When sigmoid neuron's activation saturates at either 0 or 1, the gradient at these regions is almost 0.

∴ no signal flow through neuron.

<u>Remedy</u> : ① Initialize the weights by extra one.
② Take small value of weights initially.

<u>Relu</u> : It will perform better than sigmoid.

ReLU also have problems like dead neuron. It can be fragile during trit training and die. Large gradients results in update weight that gradient flowing could be zero. This kills neurons.
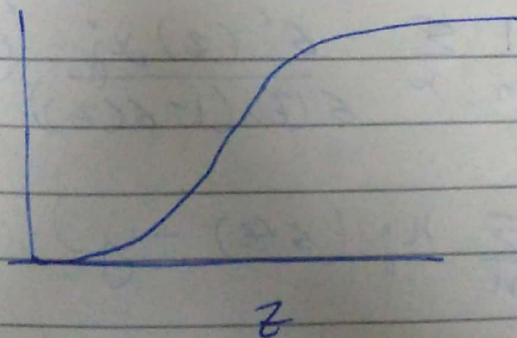
**Remedy** ⇒ Keep the learning rate low.

**B**

A3 = For quadratic, learning is slow when error is large.

$$C = \frac{(y-a)^2}{2}$$

$$a = \sigma(z)$$

$$\frac{\partial C}{\partial w} = (a-y)\sigma'(z) n = a\sigma'(z)$$

illy $\frac{\partial C}{\partial b} = a\sigma'(z)$



z

∴ $\frac{\partial C}{\partial w}$, $\frac{\partial C}{\partial b}$ get very small

∴ slow learning.

Cross entropy cost,

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

I$\frac{}{}$ $\boxed{C > 0}$

II$\frac{}{}$ If neuron's actual output is close to desired o/p , $x$, then cross entropy will be close to zero

III$\frac{}{}$

$$\frac{\partial C}{\partial w_j} = -\frac{1}{n} \sum \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma z} \right) \sigma'(z) x_j$$

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

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

This tells the rate at which weight learns is controlled by
$\sigma(z) - y$

∴ larger error ⟹ more it will learn.