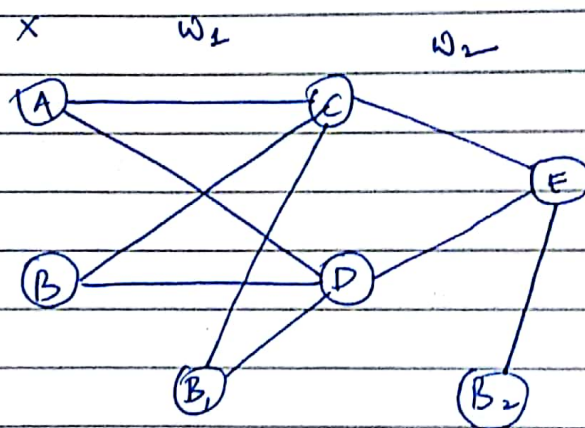


THEORY ASSIGNMENT

QUESTION 1

XOR f^{\wedge} cannot be modelled using a neural network with linear activation function. As a neural network with linear activation f^{\wedge} is equal to a linear model.



Input at hidden layer 1.

$$= w_1 x$$

Output at hidden layer 1

$$= w_1 x + B_1$$

Input at output layer

$$= w_2 (w_1 x + B_1)$$

Output at output layer

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

$$= w_1 w_2 x + w_2 b_1 + b_2$$

$$= \underbrace{w_1 w_2}_{w^*} x + \underbrace{w_2 b_1 + b_2}_{b^*}$$

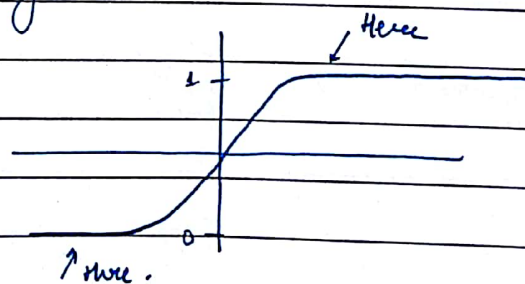
$$= w^* x + b^*$$

which is equivalent to linear models.

QUESTION 2

One of the possible reason that the model didn't get trained is because of 'vanishing gradient problem' of sigmoid function.

Vanishing Gradient Problem ~~is~~ occurs if slope when for sigmoid fn enters a region where slope difference is negligible. i.e.



If ReLU fn is used instead of sigmoid, vanishing gradient problem won't occur but data explosion can occur and also dead unit problem can occur.

Solution

Possible solutions:

- Scale down input by a factor for Vanishing Gradient Problem.
- using leaky Relu or normalising input data.

QUESTION 3

Let

$$y = \text{target}$$

$$\hat{y} = \text{predicted}$$

$$\alpha = \sigma(z)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid})$$

$$\hat{y} = e(z) = \sigma(wx + b)$$

Now

$$MSE = \frac{1}{2} (y - \hat{y})^2$$

By diff.

$$\frac{d(MSE)}{dw} = -(y - \hat{y}) \frac{\partial(\hat{y})}{\partial w} \quad \text{--- (1)}$$

$$\frac{\partial(\hat{y})}{\partial w} = \frac{\partial(\sigma(z))}{\partial w}$$

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

$$= \sigma(z)(1-\sigma(z)) \cdot x \quad \text{--- (2)}$$

Using (1) and (2)

$$\frac{\partial \text{MSE}}{\partial w} = (\hat{y} - y) \sigma(z)(1-\sigma(z)) \cdot x \quad \text{--- (3)}$$

For Cross Entropy Cost fn.

$$C = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

$$\frac{\partial C}{\partial w} = - \left[\frac{y}{\hat{y}} \frac{\partial(\hat{y})}{\partial w} + \left(\frac{-(1-y)}{(1-\hat{y})} \frac{\partial(\hat{y})}{\partial w} \right) \right]$$

$$= - \left[\frac{y}{\hat{y}} + \left(\frac{-(1-y)}{(1-\hat{y})} \right) \right] \frac{\partial(\hat{y})}{\partial w}$$

$$= \frac{y}{\hat{y}} + \left(\frac{-(1-y)}{(1-\hat{y})} \right) \sigma(z) \cdot (1-\sigma(z))$$

$$= \frac{y - \hat{y}}{\hat{y}(1-\hat{y})} \sigma(z)(1-\sigma(z)) \cdot x$$

$$= (\hat{y} - y) \cdot x \quad \text{--- (4)}$$

Comparing 3 and 4, we can say $\sigma(z)(1-\sigma(z))$ acts as slow down factor.