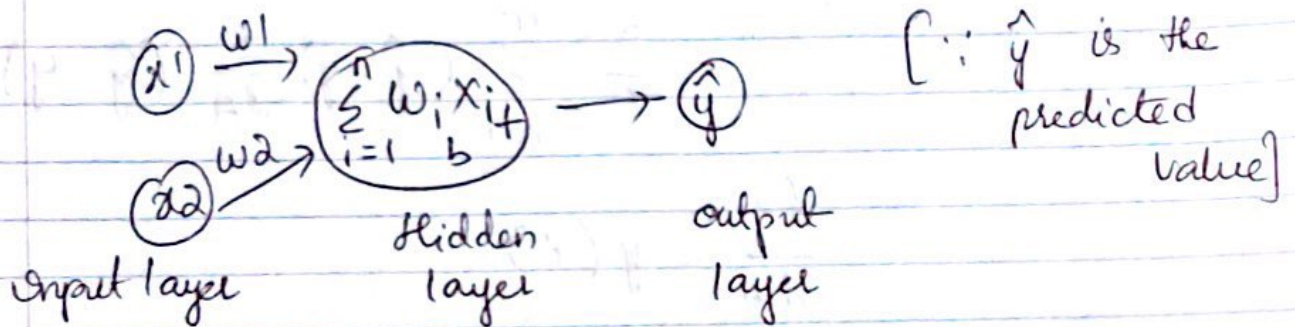


Deepthi Mikkilineni
CUID: A20349205.

Feed forward:

⇒



So as of the question and concepts

$$Z(i) = w(i)x + b \quad \left[\because \text{Here } w = \text{weight} \text{ and } b = \text{bias} \right]$$

$$A(i) = g(Z(i)) \quad \left[\because A \text{ is the activation} \right]$$

By saying that mean square error = $\frac{1}{n} \sum_{i=1}^n (\hat{y} - y)^2$
So the sigmoid activation function being used for the hidden layers then as of the below function
→ The MSE is the Average of (predicted value - actual value)²

Now the sigmoid activation will be

$$\text{Loss } L = \frac{1}{m} \sum_{i=1}^m (\hat{y} - y)^2$$

Chain rule in order to get the output is

$$\rightarrow \frac{dL}{dW} = \frac{dL}{dA} \times \frac{dA}{dz}$$

$$\begin{aligned} \frac{dL}{dA} &= \frac{1}{n} \frac{d}{dA} \sum_{i=1}^m (\hat{y} - y)^2 \\ &= \frac{2}{n} \sum_{i=1}^m (\hat{y} - y) \frac{d}{dA} (\hat{y} - y) \end{aligned}$$

$$\frac{dz}{dW} = g'(z)$$

We know the sigmoid function

$$g'(z) = g(z) (1 - g(z))$$

$$\frac{dL}{dW} = \frac{2}{n} (\hat{y} - y) \sum_{i=1}^m \frac{d}{dA} (\hat{y} - y) \left(\frac{dz}{dW} \right)$$

Assume $\hat{y} = a$

$$\frac{dL}{dW} = \frac{2}{n} (0_i - y) \frac{d0_i}{dA} \frac{dz}{dW}$$

⇒ While coming to the concept of regression, the output layer uses the linear activation and the difference between both using log loss and HSE as the loss function, is that HSE cannot be used to measure the probability to find an outcome to the defined problem. It's evident from the above that updates for regression is twice as much as the classification.