Assignment - 3 [Theory]

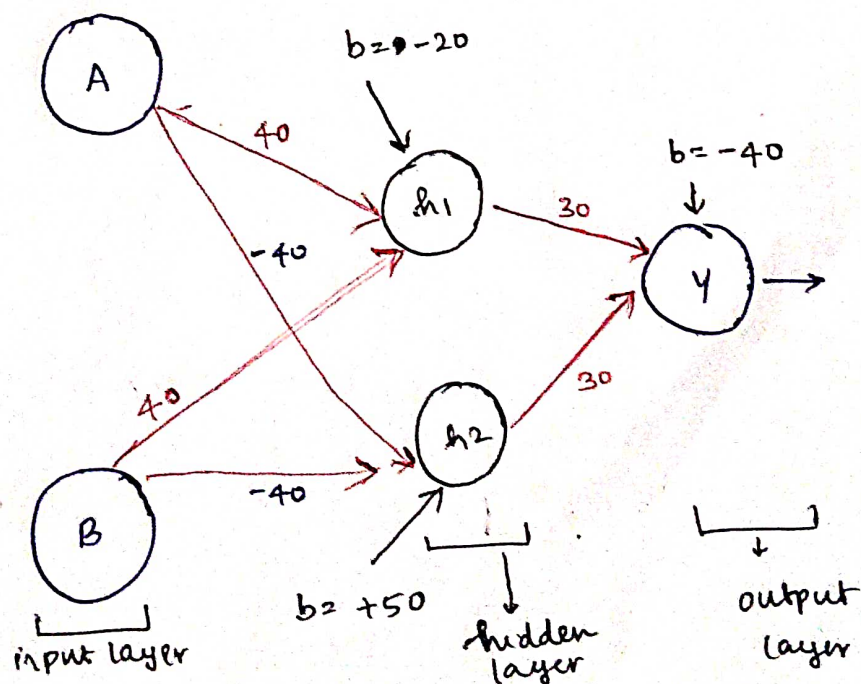1. Neural Networks.

a) The truth table of XOR is

| A | B | Y = A⊕B |
|---|---|---------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

To prove: A two layer perceptron can solve the XOR problem.

solution: Consider the two layer perceptron as



Let

$$h_1 = \sigma(40A + 40B - 20)$$

$$h_2 = \sigma(-40A - 40B + 50)$$

sigmoid function

and $\quad y = \sigma(30h_1 + 30h_2 - 40)$.

Case (i)  when $A = 0$, $B = 0$

$h_1 = \sigma(40(0) + 40(0) - 20) = \sigma(-20) = 0$

$h_2 = \sigma(-40(0) - 40(0) + 50) = \sigma(50) = 1$

$y = \sigma(30(0) + 30(1) - 40) = \sigma(-10) = 0$

$\therefore$ when $A = 0$, $B = 0 \Rightarrow y = 0 \checkmark$.

Case (ii)  when $A = 0$, $B = 1$

$h_1 = \sigma(40(0) + 40(1) - 20) = \sigma(20) = 1$

$h_2 = \sigma(-40(0) - 40(1) + 50) = \sigma(10) \approx 1$

$y = \sigma(30(1) + 30(1) - 40) = \sigma(20) = 1$

$\therefore$ when $A = 0$, $B = 1 \Rightarrow y = 1 \checkmark$

Case (iii)  when $A = 1$, $B = 0$

$h_1 = \sigma(40(1) + 40(0) - 20) = \sigma(20) \approx 1$

$h_2 = \sigma(-40(1) - 40(0) + 50) = \sigma(10) \approx 1$

$y = \sigma(30(1) + 30(1) - 40) = \sigma(20) \approx 1$

$\therefore$ when $A = 1$, $B = 0 \Rightarrow y = 1 \checkmark$

Case (iv)  when $A = 1$, $B = 1$
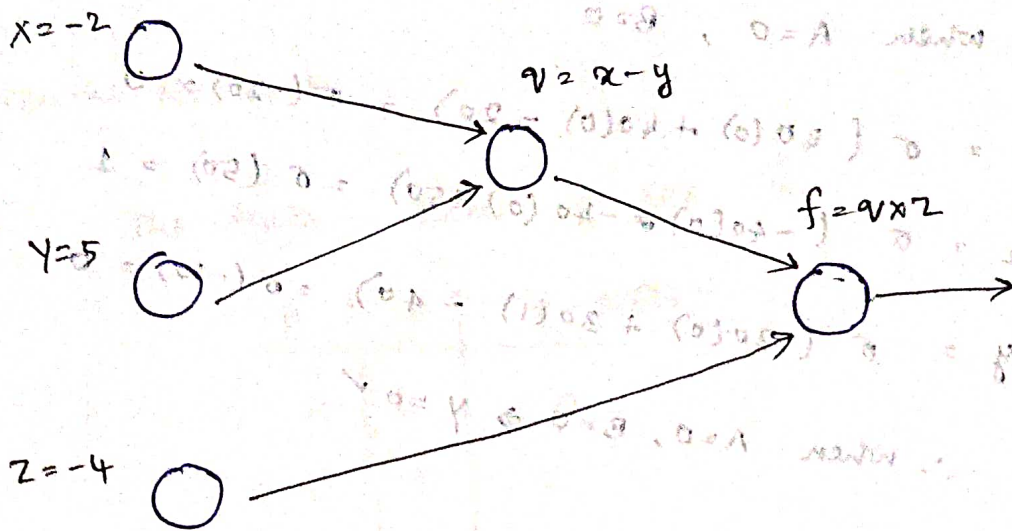
$h_1 = \sigma(40(1) + 40(1) - 20) = \sigma(60) = 1$

$h_2 = \sigma(-40(1) - 40(1) + 50) = \sigma(-30) \approx 0$

$y = \sigma(30(1) + 30(0) - 40) = \sigma(-10) \approx 0$

$\therefore$ when $A = 1$, $B = 1$, $y = 0 \checkmark$   Thus XOR truth table is verified.

(b)

$X = -2$

$Y = 5$

$Z = -4$

$V = x - y$

$f = v \times z$

Gradient f with respect to x $= \dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial v} \cdot \dfrac{\partial v}{\partial x}$

$= z (1)$

$= -4$

Gradient f with respect to y $= \dfrac{\partial f}{\partial y} = \dfrac{\partial f}{\partial v} \cdot \dfrac{\partial v}{\partial y}$

$= z(-1) = -z$

$= -(-4) = 4$

Gradient f with respect to z $= \dfrac{\partial f}{\partial z} = v = x - y = -7$

$$\therefore \dfrac{\partial f}{\partial x} = -4 \ , \qquad \dfrac{\partial f}{\partial y} = 4 \ , \qquad \dfrac{\partial f}{\partial z} = -7$$

2)

Given: cross entropy error function is,

$$E(W) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{kn} \ln y_k(x_n, W) \rightarrow \text{①}$$

where $N$ = no. g data samples, $K$ = number g classes

$$y(x_n, w) = p(t_k = 1/x) = \frac{\exp(a_k(x,w))}{\sum_{j} \exp(a_k(x,w))}$$

$0 \le y \le 1$, and $\sum_k y_k = 1$ ; $a_k \rightarrow$ pre softmax activation of output layer neurons

We know that by chain rule

$$\frac{\partial E}{\partial W_{ki}} = \sum_{j} \frac{\partial E}{\partial y_j} \cdot \frac{\partial y_j}{\partial a_k} \cdot \frac{\partial a_k}{\partial W_{ki}} \rightarrow \text{②}$$

$E(W) = +$ ½ Partial derivative of ①

with respect to $y_j$ for a particular value of $j$

is

$$\frac{\partial E}{\partial y_j} = -\frac{t_j}{y_j} \rightarrow \text{③}$$

Two cases:

$$\text{Let us find } \frac{\partial y_j}{\partial a_k} \begin{cases} j = k \\ j \ne k. \end{cases}$$

when $j = k$

$$\frac{\partial y_j}{\partial a_k} = \frac{\sum \exp(a_i) \exp(a_j) - \exp(a_j) \exp(a_j)}{(\sum \exp(a_i))^2}$$

$$\frac{\partial y_j}{\partial a_k} = y_i - y_i^2 = y_j(1-y_j) \quad \hookrightarrow ④$$

when $j \neq k$,

$$\frac{\partial y_j}{\partial a_k} = - \frac{\exp(a_j)\exp(a_k)}{(\sum \exp(a_i))^2} = -y_j y_k. \quad \hookrightarrow ⑤$$

④ & ⑤ ⟹ Generalised as

$$\boxed{\frac{\partial y_j}{\partial a_k} = y_j(\delta_{jk} - y_k)} \quad \hookrightarrow ⑥$$

$\delta_{jk}$ ⟹ Kronecker delta defined as

$\delta_{jk} = 1$ when $j = k$

$\delta_{jk} = 0$ when $j \neq k$.

Substituting ⑥, ③ in ②

$$\frac{\partial E}{\partial w_{ki}} = \sum_j - \frac{t_j}{y_j} \cdot y_j(\delta_{jk} - y_k) \cdot \frac{\partial a_k}{\partial w_{ki}}$$

$$= \sum_j t_j (y_k - \delta_{jk}) \cdot \frac{\partial a_k}{\partial w_{ki}}$$

Given in question $t_n = [0, 0, 1, 0 \ldots, 0]$

$$\therefore \sum_j t_j = 1$$

for an w one given input,

$$\frac{\partial E}{\partial W_{Ki}} = \left[ \left( \sum_{m=1}^{M} (y_k - \delta_{jk}) \cdot \frac{\partial a_k}{\partial W_{Ki}} \right) \right]$$

$$\Rightarrow \frac{\partial E}{\partial W_{Ki}} = \frac{\partial E}{\partial a_k} \cdot \frac{\partial a_k}{\partial W_{Ki}} = (y_k - \delta_{jk}) \cdot \frac{\partial a_k}{\partial W_{Ki}}$$

$$\delta_{jk} \to t_k \ (target)$$

$$\Rightarrow \left[ \boxed{\frac{\partial E}{\partial a_k} = (y_k - t_k)} \right]$$

**3)**

$$f(x) = x^2$$

$$E_{AV} = \frac{1}{M} \sum_{m=1}^{M} E_x \left[ (y_m(x) - f(x))^2 \right]$$

$$E_{ENS} = E_x \left[ \left[ \frac{1}{M} \sum_{m=1}^{M} y_m(x) - f(x) \right)^2 \right]$$

To prove: $E_{ENS} \leq E_{AV}$

Given

$$E_{ENS} = E_x \left[ \left( \frac{1}{M} \sum_{m=1}^{M} y_m(x) - f(x) \right)^2 \right]$$

$$= E_x \left[ \frac{1}{M^2} \left( \sum_{m=1}^{M} y_m(x) - f(x) \right)^2 \right]$$

$$= E_x \left[ \frac{1}{M^2} \left\{ (y_1(x) - f(x)) + (y_2(x) - f(x) + \ldots \right. \right.$$
$$\left. \left. \ldots + (y_m(x) - f(x)) \right]^2 \right]$$

$$= E_x \left[ \frac{1}{M^2} \left( \sum_{m=1}^{M} \varepsilon_m \right)^2 \right]$$

$$= E_x \left[ \left( \frac{\sum\limits_{m=1}^{m} \varepsilon_m}{M} \right)^2 \right]$$

By Jensen Inequality,

$$f\left( E(x) \right) \leq E\left( f(x) \right)$$

$$E_{ENS} = E_x \left[ \left( \frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots \varepsilon_m}{M} \right)^2 \right]$$

$$\leq \frac{1}{M} E_x(\varepsilon_1^2) + \frac{E_x}{M}(\varepsilon_2^2) + \dots + \frac{E_x(\varepsilon_m^2)}{M}$$

$$\leq \frac{1}{M} \sum_{m=1}^{M} E_x(\varepsilon_m^2)$$

$$\rightarrow \boxed{E_{ENS} \leq E_{AV}}$$

The Jensen's inequality mainly cares about the convexity of the function rather than error function $E(y)$. Hence the result stands true for any error function