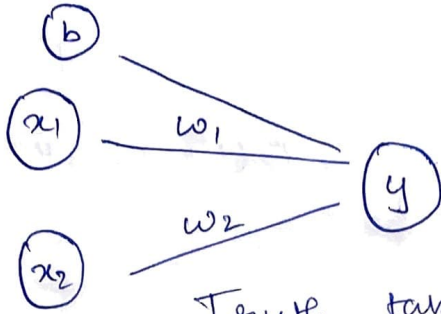


HOMEWORK - 2

CS 541

DEEP LEARNING

1. XOR Problem.



Truth table of XOR

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Input x_1, x_2

Output y

weights w

bias b

for the network

$$\Rightarrow X = (x_1^{(i)}, x_2^{(i)}) = \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$y = [y^{(i)}] = [0 \quad 1 \quad 1 \quad 0]^T$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad \text{or} \quad b = \begin{bmatrix} b \end{bmatrix}$$

Loss function as given in equation 6.1 of Deep learning book.

$$J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(y) - f(x; \theta))^2$$

where θ is ~~the~~ w_1, w_2 & b in our case

& $f^*(y)$ = true value of XOR from truth table.

$$J(w_1, w_2, b) = \frac{1}{4} \sum_{x \in X} [y^{(i)} - (x^{(i)T} w + b)]^2$$

$$= \frac{1}{4} \left[\begin{pmatrix} 0 & 1 & 1 & 0 \end{pmatrix} - \left(\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_{2 \times 1} + b \right) \right]^2$$

$$= \frac{1}{4} \left[\begin{pmatrix} 0 & 1 & 1 & 0 \end{pmatrix} - \left(\begin{bmatrix} 0 \\ w_2 \\ w_1 \\ w_1 + w_2 \end{bmatrix} + b \right) \right]^2$$

$$= \frac{1}{4} \left[\begin{pmatrix} 0 & 1 & 1 & 0 \end{pmatrix} - \begin{bmatrix} b & w_2 + b & w_1 + b \\ w_1 + w_2 + b \end{bmatrix}^T \right]^2$$

$$= \frac{1}{4} \left[\begin{matrix} b & 1 + w_2 + b & 1 + w_1 + b & w_1 + w_2 + b \end{matrix} \right]^2$$

$$= \frac{1}{4} \left[\begin{matrix} -b & 1 - (w_2 + b) & 1 - (w_1 + b) & -(w_1 + w_2 + b) \end{matrix} \right]^2$$

$$\begin{aligned}
 & \text{So } \left(\begin{bmatrix} -b & 1-(\omega_2+b) & 1-(\omega_1+b) & -(\omega_1+\omega_2+b) \end{bmatrix}^T \right)^2 \\
 &= \begin{bmatrix} -b & 1-(\omega_2+b) & 1-(\omega_1+b) & -(\omega_1+\omega_2+b) \end{bmatrix} \begin{bmatrix} -b \\ 1-(\omega_2+b) \\ 1-(\omega_1+b) \\ -(\omega_1+\omega_2+b) \end{bmatrix}
 \end{aligned}$$

$$= \left[b^2 + (1-(\omega_2+b))^2 + (1-(\omega_1+b))^2 + (\omega_1+\omega_2+b)^2 \right]$$

$$\begin{aligned}
 &= \left[b^2 + 1 + (\omega_2+b)^2 - 2(\omega_2+b) + 1 + (\omega_1+b)^2 - 2(\omega_1+b) \right. \\
 &\quad \left. + (\omega_1+\omega_2)^2 + 2(\omega_1+\omega_2)b + b^2 \right]
 \end{aligned}$$

$$\begin{aligned}
 &= \left[b^2 + 2 + \omega_2^2 + b^2 + 2\omega_2 b - 2(\omega_1 + \omega_2 + 2b) + \omega_1^2 + b^2 + 2\omega_1 b \right. \\
 &\quad \left. + \omega_1^2 + \omega_2^2 + 2\omega_1 \omega_2 + 2\omega_1 b + 2\omega_2 b + b^2 \right]
 \end{aligned}$$

$$= \left[4b^2 + 2 + 2(\omega_1^2 + \omega_2^2) + 4b(\omega_1 + \omega_2) - 2(\omega_1 + \omega_2 + 2b) + 2\omega_1 \omega_2 \right]$$

$$\begin{aligned}
 J(\omega, \theta) &= \\
 &\left[4b^2 + 2(\omega_1^2 + \omega_2^2) + 4b(\omega_1 + \omega_2) + 2\omega_1 \omega_2 + 2 - 2(\omega_1 + \omega_2 + 2b) \right]
 \end{aligned}$$

$$\nabla J(\omega_1) = [0 + 4\omega_1 + 4b + 2\omega_2 - 2] = 0 \rightarrow (1)$$

$$\nabla J(\omega_2) = [0 + 4\omega_2 + 4b + 2\omega_1 - 2] = 0 \rightarrow (2)$$

$$\nabla J(b) = [8b + 4(\omega_1 + \omega_2) - 4] = 0 \rightarrow (3)$$

$$4\omega_1 + 2\omega_2 + 4b = 2 \rightarrow (4)$$

$$4\omega_2 + 2\omega_1 + 4b = 2 \rightarrow (5)$$

$$8b + 4\omega_1 + 4\omega_2 = 4 \rightarrow (6)$$

$$(4) - (5)$$

$$2\omega_1 - 2\omega_2 = 0$$

$$\omega_1 = \omega_2 \rightarrow (7)$$

Substituting (7) in (6) & (5)

$$8b + 8\omega_1 = 4 \times 3$$

$$6\omega_1 + 4b = 2 \times 4$$

we get

$$\begin{array}{r} 24b + 24\omega_1 = 12 \\ (-) 24\omega_1 + 16b = 8 \\ \hline 8b = 4 \end{array}$$

$$\boxed{b = 1/2} \rightarrow (8)$$

Substituting $b = 1/2$ in (5)

$$4\omega_1 + 2\omega_1 + 4 \times \frac{1}{2} = 2$$

$$6\omega_1 + 2 = 2$$

$$\boxed{\omega_1 = 0} \rightarrow (9)$$

Since $\omega_1 = \omega_2$
 $\Rightarrow \omega_2 = 0$

∴ The values of w_1, w_2 & b to minimize the function $J(w, b) =$

$$\frac{1}{4} \sum_{x \in X} (f^*(y) - f(x; w, b))^2$$

is

$$\begin{array}{l} w_1 = 0 \\ w_2 = 0 \\ b = 1/2 \end{array}$$

Q3.

L_2 regularization term

$$L_2 = \frac{\alpha}{2n} W^T W.$$

Given input $x = [x_1 \quad x_2]$

Let weights be $w = [w_1 \quad w_2]^T$

Let a matrix S be included in the L_2 regularization term.

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}$$

$$\Rightarrow L_2 \text{ cost} = J = \frac{\alpha}{2n} W^T S W$$

$$= \begin{bmatrix} \omega_1 & \omega_2 \end{bmatrix} \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix}$$

$$= \begin{bmatrix} \omega_1 & \omega_2 \end{bmatrix} \begin{bmatrix} S_{11}\omega_1 + S_{12}\omega_2 \\ S_{21}\omega_1 + S_{22}\omega_2 \end{bmatrix}$$

$$= \omega_1^2 S_{11} + \omega_1 \omega_2 S_{12} + \omega_1 \omega_2 S_{21} + \omega_2^2 S_{22}$$

$$J_{L2} = \omega_1^2 S_{11} + \omega_1 \omega_2 (S_{12} + S_{21}) + \omega_2^2 S_{22}$$

$J_{L2} = 0$ when ω_1 approximately equal to ω_2 , such that we need not penalise the cost function as the weights are reflective over the middle column

$\Rightarrow \omega_1 - \omega_2 \approx 0$

$\Rightarrow (\omega_1 - \omega_2)^2 \approx 0$

$\left[\text{LHS} \mid \text{RHS} \right] \quad \text{LHS} \approx \text{RHS}$

$$\Rightarrow (\omega_1 - \omega_2)^2 = \omega_1^2 S_{11} + \omega_1 \omega_2 (S_{12} + S_{21}) + \omega_2^2 S_{22}$$

$$\omega_1^2 + \omega_2^2 - 2\omega_1 \omega_2 = \omega_1^2 S_{11} + \omega_2^2 S_{22} + \omega_1 \omega_2 (S_{12} + S_{21})$$

comparing coefficients of ω_1^2 , ω_2^2 & $\omega_1 \omega_2$.

$$S_{11} = 1, \quad S_{22} = 1, \quad S_{12} + S_{21} = -2$$

$$S_{12} = S_{21} = -1 \quad (\text{assumed})$$

To prove $S_{12} = S_{21} = -1$.

$$\begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$= \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} w_1 - w_2 \\ w_2 - w_1 \end{bmatrix}$$

$$= w_1(w_1 - w_2) + w_2(w_2 - w_1) \rightarrow \textcircled{1}$$

Equation $\textcircled{1} = 0$.

$$w_1^2 - w_1 w_2 + w_2^2 - w_1 w_2 = 0.$$

$$w_1^2 - 2w_1 w_2 + w_2^2 = 0$$

$$(w_1 - w_2)^2 = 0$$

is the condition when the weights are symmetric & this can be derived only when $S_{12} = S_{21} = -1$

Explanation for matrix S .

when weights w_1 & w_2 are not symmetric

$$\begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$= w_1(w_1 - w_2) + w_2(w_2 - w_1)$$

$$= \frac{\alpha}{2n} [(\omega_1 - \omega_2)(\omega_1 - \omega_2)]$$

$$= \frac{\alpha}{2n} [(\omega_1 - \omega_2)^2]$$

which penalises the cost function

when $\omega_1 = \omega_2$, i.e. the weights are symmetric (in the sense reflective along the middle column)
 (not $x^T = x$)

$$J_{L2} = \frac{\alpha}{2n} [(\omega_1 - \omega_2)]^2$$

$$= \frac{\alpha}{2n} \times 0$$

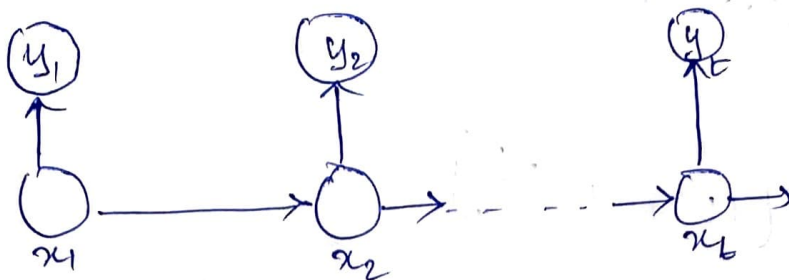
$$= 0$$

$$\left[\begin{array}{c|c} \text{LHS} & \text{RHS} \end{array} \right]$$

Since the weights are already symmetric it does not penalise the cost function
 ~~LHS = 255 - R1~~

Hence $L2$ can be regularization can be used to discourage asymmetric weights ~~or~~ if the data are symmetric in nature.

$$S = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$



$y_1, y_2, y_3, \dots, y_t$ are the observable random variables, ~~here~~ ^{here} the student's behaviour which is observable ~~which is~~ determined by that student's current state which is thought process inside his brain which is not observable given by x_1, x_2, \dots, x_t (corresponding y_i by x_i, x_{i-1}, \dots, x_1 $i \rightarrow 1 \text{ to } t$)

Hence x_1, \dots, x_t are known as hidden units & y_1, \dots, y_t are known as observable units.

Given. $\left\{ \begin{array}{l} P(x_t | x_1, \dots, x_{t-1}) = P(x_t | x_{t-1}) \\ \quad \quad \quad \downarrow \\ \quad \quad \quad \text{markov property.} \\ P(y_t | x_t, y_1, \dots, y_{t-1}) = P(y_t | x_t) \end{array} \right.$

The goal of the teacher is to estimate the current state x_t given the observations y_1, \dots, y_t and update her belief about the student.

ie; to ~~find~~ ^{prove}

$$\begin{cases} P(x_t | y_1, \dots, y_{t-1}, y_t) \propto \\ P(x_t | x_{t-1}) P(x_{t-1} | y_1, \dots, y_{t-1}) \end{cases}$$

Applying Bayes' Theorem.

$$P(x_t | y_t, y_{t-1}, \dots, y_1) = \frac{P(y_t | x_t, y_1, \dots, y_{t-1}) P(x_t | y_1, \dots, y_{t-1})}{P(y_t | y_1, \dots, y_{t-1})} \quad (1)$$

$$P(y_t | y_1, \dots, y_{t-1}) \rightarrow (1)$$

Since the denominator of (1) does not involve x_t , equation (1) can be re-written as.

$$P(x_t | y_1, \dots, y_t) \propto P(y_t | x_t, y_1, \dots, y_{t-1}) P(x_t | y_1, \dots, y_{t-1})$$

$$P(y_t | x_t, y_1, \dots, y_{t-1}) = P(y_t | x_t) \rightarrow \text{Given } x_t \text{ can be inferred from the chain.}$$

$$\Rightarrow P(x_t | y_1, \dots, y_t) \propto P(y_t | x_t) P(x_t | y_1, \dots, y_{t-1}) \rightarrow (2)$$

x_t depends on x_{t-1} as given in the chain and x_{t-1} can be any state of mind of the student, hence according to law of total probability, $P(x_t) = \sum_{x_{t-1}} P(x_t, x_{t-1})$

Applying law of probability to equation (2) we get.

$$P(x_t | y_1, \dots, y_t) \propto P(y_t | x_t) \sum_{x_{t-1}} P(x_t, x_{t-1} | y_1, \dots, y_{t-1}) \quad \rightarrow (3)$$

Applying conditional probability distribution concept to equation (3) we get.

$$P(x_t | y_1, \dots, y_t) \propto P(y_t | x_t) \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1} | y_1, \dots, y_{t-1})$$

where $P(x_{t-1} | y_1, \dots, y_{t-1})$ is the teacher's belief from the time stamp of $t-1$ or the summation is over the all possible values of x_{t-1} i.e., all the thought process in the student's mind.

$$\therefore \boxed{P(x_t | y_1, \dots, y_t) \propto P(y_t | x_t) \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1} | y_1, \dots, y_{t-1})}$$

Hence Proved

Cost fn.

5.

$$P(y | x, \omega, \sigma^2) = \mathcal{N}(y; x^T \omega, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - x^T \omega)^2}{2\sigma^2}\right).$$

$$P(\mathcal{D} | \omega, \sigma^2) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \omega, \sigma^2) \quad \text{--- (1)}$$

where $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$

Taking log on both sides of eqn. (1).

$$\log P(\mathcal{D} | \omega, \sigma^2) = \log \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \omega, \sigma^2)$$

$$= \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}, \omega, \sigma^2)$$

$$= \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2} \right) \right)$$

$$= \sum_{i=1}^n \left[\log \frac{1}{\sqrt{2\pi\sigma^2}} + \log \left(\exp \left(-\frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2} \right) \right) \right]$$

$$= \sum_{i=1}^n \left[\log \frac{1}{\sqrt{2\pi}} + \log \frac{1}{\sqrt{\sigma^2}} + \left(-\frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2} \right) \right]$$

$$= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}} + \sum_{i=1}^n \log \frac{1}{\sqrt{\sigma^2}} - \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2}$$

$$= -\frac{1}{2} \sum_{i=1}^n \log 2\pi - \frac{1}{2} \sum_{i=1}^n \log \sigma^2 - \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2}$$

$$= -\frac{n(\log 2\pi)}{2} - \frac{n}{2} \log(\sigma^2) - \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^2}$$

Differentiating eqn. (2) w.r.t σ^2 , equating it to zero

$$0 = 0 - \frac{n}{2} \times \frac{1}{\sigma^2} - \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^4}$$

$$\frac{n}{2\sigma^2} = \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T} \omega)^2}{2\sigma^4}$$

$$\therefore \sigma^2 = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - x^{(i)T} \omega)^2$$

Taking (-1) common from

$$\Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)T} \omega - y^{(i)})^2$$

Differentiating equation (2) w.r.t ω and equating it to zero.

$$0 = 0 + 0 - 2 \sum_{i=1}^n (y^{(i)} - x^{(i)T} \omega) \times (-x^{(i)})$$

since $\boxed{\nabla (x^{(i)T} \omega) = x^{(i)}}$

$$0 = - \sum_{i=1}^n y^{(i)} x^{(i)} + \sum_{i=1}^n x^{(i)} x^{(i)T} \omega :$$

$$\sum_{i=1}^n x^{(i)} x^{(i)T} \omega = \sum_{i=1}^n x^{(i)} y^{(i)}$$

$$\omega \sum_{i=1}^n x^{(i)} x^{(i)T} = \sum_{i=1}^n x^{(i)} y^{(i)}$$

$$\omega = \left(\sum_{i=1}^n (x^{(i)} x^{(i)T}) \right)^{-1} \sum_{i=1}^n x^{(i)} y^{(i)}$$

$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)T} \omega - y^{(i)})^2$

Q2.

Set of Learning Rates = $[0.001, 0.005, 0.01, 0.0005]$

Set of no. of epochs = $[50, 100, 200, 400]$

Set of batch sizes = $[128, 256, 512, 1024]$

Set of regularization constant = $[0.01, 0.2, 0.5, 0.005]$

Result

Best Epochs = 400

Best Learning Rate = 0.001

Best batch size = 512

Best regularization constant = 0.5

Best Validation Loss = 118.10845242

~~Best~~
TEST MSE LOSS = 120.87438984

~~(Report Attached)~~
~~Summary~~