

## Task 1

- stride:  $s_1 = 1, s_2 = 3$

• Output of convolution layer:  $V (3 \times 2)$

- After max pooling:  $m (2 \times 1)$

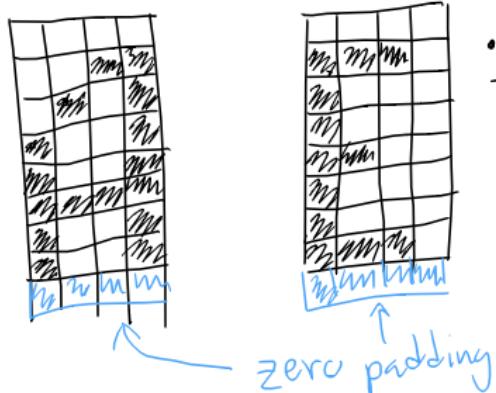
1) Calculate entries of  $V$  matrix.

2) Attempt: set all  $w_{ij} = 1$

3) Get entries for  $m$

4) Determine  $\theta$  and  $\bar{W}$

A



$$V_{11} = w_{23} + w_{32} - \theta$$

$$V_{12} = w_{22} + w_{23} + w_{31} + w_{33} - \theta$$

$$V_{21} = w_{11} + w_{21} + w_{31} + w_{32} + w_{33} - \theta$$

$$V_{22} = w_{13} + w_{23} + w_{31} + w_{32} + w_{33} - \theta$$

$$V_{31} = w_{11} + w_{21} - \theta$$

$$V_{32} = w_{13} + w_{23} - \theta$$

E

$$V_{11} = w_{21} + w_{22} + w_{23} + w_{31} - \theta$$

$$V_{12} = w_{21} + w_{22} - \theta$$

$$V_{21} = w_{11} + w_{21} + w_{22} + w_{31} - \theta$$

$$V_{22} = w_{21} - \theta$$

$$V_{31} = w_{11} + w_{21} + w_{22} + w_{23} - \theta$$

$$V_{32} = w_{21} + w_{22} - \theta$$

$$m(A) = (5 - \theta, 5 - \theta)^T \rightarrow$$

$$\bar{W}_1(5 - \theta) + \bar{W}_2(5 - \theta) < 0$$

$$m(E) = (4 - \theta, 4 - \theta)^T$$

$$\bar{W}_1(4 - \theta) + \bar{W}_2(4 - \theta) > 0$$

ANS with stride  $s_1 = 1, s_2 = 3$ : All  $w_{ij} = 1, \theta = 4, 5,$

note: infinite possible solutions.

$$\bar{W}_1 = -1, \bar{W}_2 = -1$$

5/5

(2)(a) Starting with the KL divergence,

$$D_{KL} = \sum_{n=1}^b P_{\text{data}}(x^n) \log \frac{P_{\text{data}}(x^n)}{P_B(s=x^n)} = - \sum_{n=1}^b P_{\text{data}} \log \frac{P_B}{P_{\text{data}}}$$

use the inequality  $\log(x) \leq x-1$ , where  
 $\log(x) = x-1$   
only if  $x=1$ .  
We find

$$\begin{aligned} D_{KL} &= - \sum_{n=1}^b P_{\text{data}} \log \frac{P_B}{P_{\text{data}}} \geq - \sum_{n=1}^b P_{\text{data}} \left( \frac{P_B}{P_{\text{data}}} - 1 \right) \\ &= - \sum_{n=1}^b \left[ P_B - P_{\text{data}}(x^n) \right] \\ &= - \sum_{n=1}^b [P_B - 1] \end{aligned}$$

$$> 0$$

2)  $D_{KL} > 0$

(b) The equivalence between maximizing the log-likelihood and minimizing the KL divergence can be seen as follows. The log-likelihood reads

$$\log L = \sum_{n=1}^p \log P_B(s=x^{(n)})$$

where, as stated in the lecture notes, the sum runs over a sequence of input patterns  $x^{(n)}$ ,  $n=1, \dots, p$ . Any pattern  $x$  may appear more than once in the sequence, with frequency proportional to  $P_{\text{data}}(x)$ . In the limit of a large number of samples,  $p \gg 1$ ,

$$\log L = p \left[ \frac{1}{p} \sum_{n=1}^p \log P_B(s=x^{(n)}) \right]$$

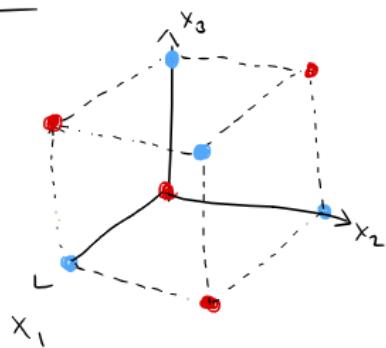
interpreting as a Using the law of large numbers, the quantity in the square brackets reads

$$\log L = \frac{1}{p} \sum_{n=1}^p \log P_B(s=x^{(n)}) \xrightarrow{p \gg 1} \sum_{n=1}^p P_{\text{data}}(x^{(n)}) \log P_B(x)$$

thus,  $\log L = p \left[ \underbrace{(\sum P_{\text{data}} \log P_{\text{data}})}_{\text{constant}} - D_{KL} \right]$

In conclusion, maximising  $\log L$  is equivalent to minimising  $D_{KL}$ .

### Task 3



- 1
- 0

No way to construct a plane that separates the points.

A network which can separate the blue from the red can be found by first constructing an XOR network for the first two variables and then using the output of that network and the third variable as input to a second XOR network.

$$w^{(1)} = \begin{pmatrix} 1 & 1 & 0 \\ -1 & -1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$w^{(2)} = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$w^{(3)} = \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix}$$

$$w^{(4)} = \begin{pmatrix} 1 & 1 \end{pmatrix}$$

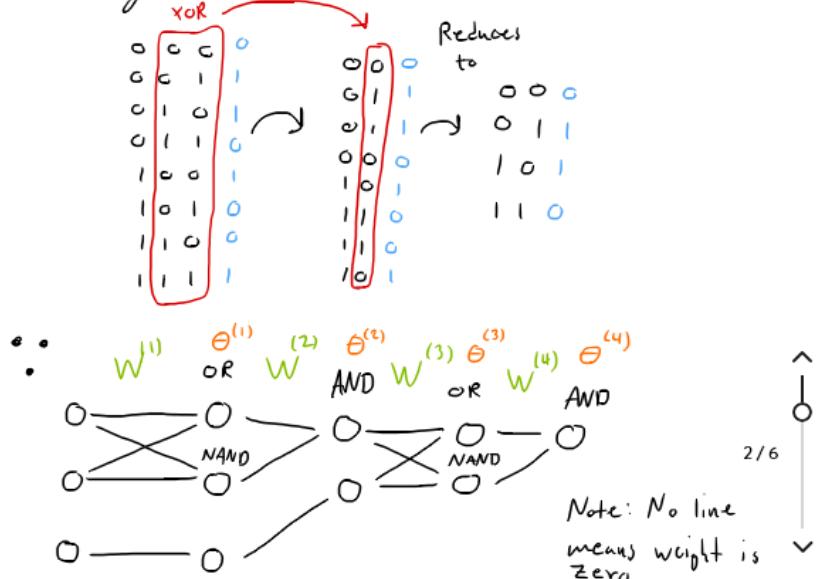
$$\theta^{(1)} = (0.5, -1.5, 0)^T$$

$$\theta^{(2)} = (1.5, 0)^T$$

$$\theta^{(3)} = (0.5, -1.5)^T$$

$$\theta^{(4)} = 1.5$$

This can be seen by looking at the logic table:



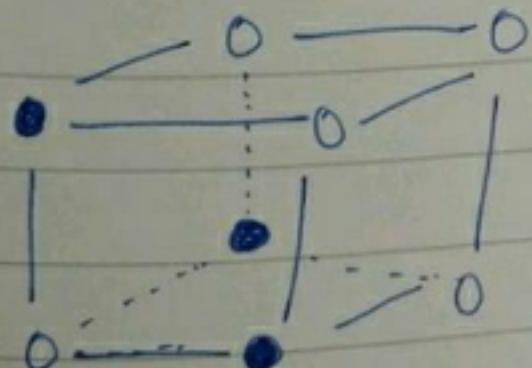
#### ④ Boolean function II

Truth table

$x_1$	$x_2$	$x_3$	$t$
0	0	0	0
0	0	1	1
0	1	0	1
1	0	0	1
0	1	1	0
1	0	1	0
1	1	0	0
1	1	1	0

Empty circle = 0

Filled circle = 1



Not linearly separable = requires hidden layer.

The A neural network representing these functions in  $N (= 3 \text{ or } 4)$  dimensions can be constructed using the winning neuron construction, discussed in section 7.1 of the lecture notes.

The hidden layer weights and thresholds need to be modified from those discussed in the lecture notes. The lecture notes define  $\pm$  hidden layer weights,

$$w_{jk} = \begin{cases} \delta & \text{if the } k^{\text{th}} \text{ digit of binary representation of } j \text{ is 1} \\ -\delta & \text{otherwise.} \end{cases}$$

and thresholds  $\theta_j = N(\delta - 1)$ .

for the above weights and thresholds are for  $\pm 1$  neurons. Let the weights and thresholds for  $0,1$  neurons be  $\tilde{w}_{ij}$  and  $\tilde{\theta}_{ij}$ . Then

$$\tilde{w}_{jk} = 2 w_{jk}$$

$$\tilde{\theta}_j = \sum_k w_{jk} + \theta_j.$$

## Task 5

Back-propagation:

$$\omega^{(l)} \leftarrow \omega^{(l)} - \eta \frac{\partial H}{\partial \omega^{(l)}}$$

$$\frac{\partial H}{\partial \omega^{(l)}} = \frac{\partial}{\partial \omega^{(l)}} \frac{1}{2} (t - V^{(L)})^2 = -(t - V^{(L)}) \frac{\partial V^{(L)}}{\partial \omega^{(l)}}$$

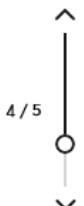
$$= -(t - V^{(L)}) g'(b^{(L)}) \omega^{(L)} \frac{\partial}{\partial \omega^{(L)}} V^{(L-1)}$$

$$= \boxed{-(t - V^{(L)}) \left[ \prod_{i=0}^{L-1} g'(b^{(L-i)}) \omega^{(L-i)} \right] g'(b^{(L)}) V^{(L-1)}}$$

$$\frac{\partial V^{(L)}}{\partial V^{(l)}} = \frac{\partial}{\partial V^{(l)}} g(\omega^{(L)} V^{(L-1)} - \theta^{(L)}) = g'(b^{(L)}) \omega^{(L)} \frac{\partial V^{(L-1)}}{\partial V^{(l)}}$$

$$= \boxed{\prod_{i=0}^{L-1} g'(b^{(L-i)}) \omega^{(L-i)}}$$

$$\therefore \boxed{\delta^{(L-1)} = (t - V^{(L)}) \frac{\partial V^{(L)}}{\partial V^{(L-1)}} g'(b^{(L-1)})}$$



(6.)

(a) Start with by differentiating the average immediate reward w.r.t  $w_n$ ,

$$\frac{\partial \langle r \rangle}{\partial w_n} = \sum_{y=\pm 1} \left\langle r(\bar{x}, y) \frac{\partial P(y, \bar{x})}{\partial w_n} \right\rangle$$

use,

$$P(y, \bar{x}) = \begin{cases} p(b) & \text{if } y = +1 \\ 1 - p(b) & \text{if } y = -1 \end{cases}$$

$$\Rightarrow \frac{\partial \langle r \rangle}{\partial w_n} = \left\langle r(\bar{x}, 1) \frac{\partial p(b)}{\partial w_n} \right\rangle + \cancel{\left\langle r(\bar{x}, -1) \frac{\partial [1 - p(b)]}{\partial w_n} \right\rangle}$$

$$+ \left\langle r(\bar{x}, -1) \frac{\partial [1 - p(b)]}{\partial w_n} \right\rangle$$

$$= \left\langle r(\bar{x}, 1) p(b) [1 - \tanh(b)] \right\rangle \cancel{x_n} \frac{\partial b}{\partial w_n}$$

$$+ \left\langle r(\bar{x}, -1) [1 - p(b)] [-1 - \tanh(b)] \right\rangle \frac{\partial b}{\partial w_n}$$

$$= \sum_{y=\pm 1} \left\langle r(\bar{x}, y) P(y, \bar{x}) [y - \tanh(b)] \right\rangle x_n$$

(b) choosing  $\delta w_n = -\eta r(\bar{x}, y) [y - \tanh(b)]_x$  gives.

$$\begin{aligned}\langle \delta w_n \rangle &= -\eta \langle r(\bar{x}, y) [y - \tanh(b)] \rangle_{x_n} \\ &= -\eta \frac{\partial \langle r \rangle}{\partial w_n}.\end{aligned}$$

(c) An unbiased estimator is an estimator whose expected value equals the parameter being estimated. Suppose the  $X$  is an estimator for  $\theta$ . Then  $X$  is an unbiased estimator if

$$E[X] = \theta.$$