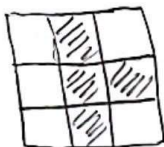


# 1) Convolutional net

Choose the following filter:

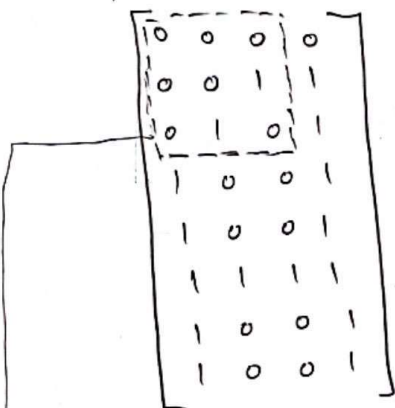


Legend:

□ corresponds to  $w = 0$

▨ corresponds to  $w = 1$

- Pattern  $x^{(1)}$  can be represented as:



The local fields of the feature map of pattern  $x^{(1)}$  can be determined as follows (assume the threshold  $\Theta$  is zero):

Consider the top left  $3 \times 3$  block of the inputs

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow 2 \text{ sum everything}$$

3x3 block of inputs
multiply
kernel

Sweep the kernel over the 2-D input matrix with stride  $[1,1]$ , padding 0.  
 It follows that the local fields of the feature map of pattern  $x^{(1)}$  are:

$$\begin{bmatrix} 2 & 2 \\ 1 & 2 \\ 1 & 1 \\ 1 & 2 \\ 2 & 2 \\ 1 & 2 \end{bmatrix}$$

- Now, pattern  $x^{(2)}$  can be represented as follows:

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

Similarly, the local fields of the feature map of pattern  $x^{(2)}$  can be determined as follows:

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow 2$$

$\begin{matrix} \text{3x3 block} \\ \text{of inputs} \end{matrix}$ 
 $\begin{matrix} \uparrow \\ \text{multiply} \end{matrix}$ 
 $\begin{matrix} \text{kernel} \end{matrix}$ 
 $\begin{matrix} \text{sum everything} \end{matrix}$

Sweep the kernel over the 2-D input matrix with stride  $[1, 1]$ , and padding  $[0, 0, 0, 0]$ . It follows that the local fields of the feature map of pattern  $x^{(2)}$  are:

$$\begin{bmatrix} 2 & 1 \\ 1 & 1 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

The ReLU - activation function does not exert any effect, since all local fields are possible. Thus, the feature maps are therefore equal to the above local fields.

- Apply max-pooling operation on the resulting feature map of pattern  $x^{(1)}$ .

$$\begin{bmatrix} 2 & 2 \\ 1 & 2 \\ 1 & 1 \\ 1 & 2 \\ 2 & 2 \\ 1 & 2 \end{bmatrix}$$

For  $2 \times 2$  max pooling, the maximum element within the top  $2 \times 2$  block is 2

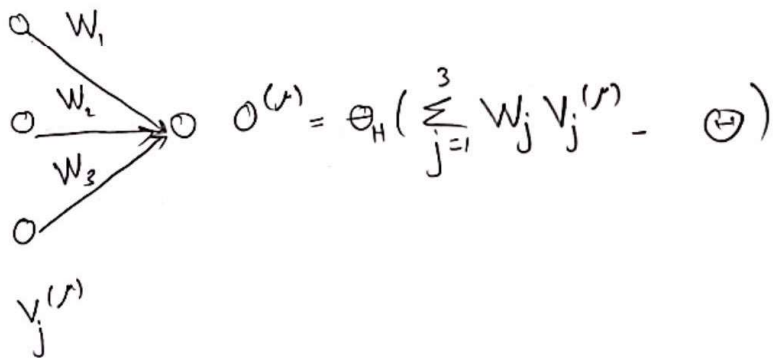
Sweep over the resulting feature map with stride 2, and padding 0  
 The output of the max-pooling layer is:

$$\begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix}$$

- Similarly for pattern  $x^{(2)}$ , the output of the max-pooling layer is

$$\begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$$

- To determine the weights  $W_k$  and threshold  $\Theta$ , the layout of the output of the network can be represented as follows:



where

$$V^{(1)} = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} \quad \text{for pattern } x^{(1)}$$

$$V^{(2)} = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \quad \text{for pattern } x^{(2)}$$

To classify the 2 patterns we need the output of the network to be

$$O^{(1)} = 0 \quad \text{for } x^{(1)}$$

$$O^{(2)} = 1 \quad \text{for } x^{(2)}$$

This can be done by setting the output of the network to be

$$O^{(\mu)} = \Theta_H (-V_1^{(\mu)} - V_2^{(\mu)} - V_3^{(\mu)} + 5)$$

$$\text{where } W = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \quad \text{and } \Theta = -5$$

In this case we get

$$\text{For } \mu=1: O^{(1)} = \Theta_H (-2 - 2 - 2 + 5) = \Theta_H (-1) = 0$$

$$\text{For } \mu=2: O^{(2)} = \Theta_H (-2 - 1 - 1 + 5) = \Theta_H (+1) = 1 \quad \text{verified}$$

## 2) Kullback - Leibler divergence

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

Since  $\ln x \leq x-1 \quad \forall x > 0$ , it follows that

$$D_{KL} \geq \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \left( - \left[ \frac{P_B(s=x^{(\mu)})}{P_{\text{data}}(x^{(\mu)})} \right] + 1 \right)$$

$$\geq \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \left[ - \frac{P_B(s=x^{(\mu)}) + P_{\text{data}}(x^{(\mu)})}{P_{\text{data}}(x^{(\mu)})} \right]$$

$$\geq \sum_{\mu=1}^P [P_{\text{data}}(x^{(\mu)}) - P_B(s=x^{(\mu)})]$$

$$\geq \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) - \sum_{\mu=1}^P P_B(s=x^{(\mu)})$$

$$\geq 1 - 1$$

(6)

The total probability is normalized to 1

Thus the Kullback - Leibler divergence is non-negative.

b) For  $P_B(s=x^{(\mu)}) = P_{\text{data}}(x^{(\mu)})$ :

$$D_{KL} = \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log \left[ \frac{P_{\text{data}}(x^{(\mu)})}{P_B(s=x^{(\mu)})} \right]$$

$$= \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log(1)$$

$$= 0 \quad \text{verified}$$

c) Show that minimizing  $D_{KL}$  is equivalent to maximizing log-likelihood function:

$$D_{KL} = \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log \left[ \frac{P_{\text{data}}(x^{(\mu)})}{P_B(s=x^{(\mu)})} \right]$$

$$= \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log P_{\text{data}}(x^{(\mu)}) - \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log P_B(s=x^{(\mu)})$$

$$= \underbrace{\sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log P_{\text{data}}(x^{(\mu)})}_{\text{Fixed by the given data distribution}} - \langle \log P_B(s=x^{(\mu)}) \rangle_{P_{\text{data}}}$$

Fixed by the given  
data distribution (7)

Now,

$$\text{Log } \mathcal{L} = \log \prod_{\mu=1}^P P_B(s=x^{(\mu)}) = \sum_{\mu=1}^P \log P_B(s=x^{(\mu)})$$

using the law of large numbers  $\frac{1}{N} \sum_{i=1}^N f(i) \stackrel{N \gg 1}{\approx} \underbrace{\int dx P(x) f(i)}_{\langle f(i) \rangle_{P(x)}}$

This,  $\text{Log } \mathcal{L}$  can be written as:

$$\begin{aligned} \text{Log } \mathcal{L} &= \sum_{\mu=1}^P \log P_B(s=x^{(\mu)}) \\ &= P \frac{1}{P} \sum_{\mu=1}^P \log P_B(s=x^{(\mu)}) \\ &\stackrel{P \gg 1}{\approx} P \langle \log P_B(s=x^{(\mu)}) \rangle_{P_{\text{data}}} \end{aligned}$$

It follows that:

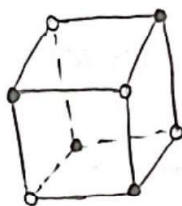
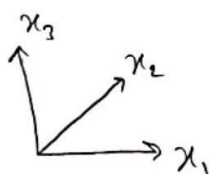
$$\text{Log } \mathcal{L} \stackrel{P \gg 1}{\approx} P \sum_{\mu=1}^P P_{\text{data}}(x^{(\mu)}) \log P_{\text{data}}(x^{(\mu)}) - P D_{\text{KL}}$$

Thus, minimizing  $D_{\text{KL}}$  is equivalent to maximizing  $\text{Log } \mathcal{L}$ .



### 3) Boolean function I:

#### a) Graphical representation in the input space:



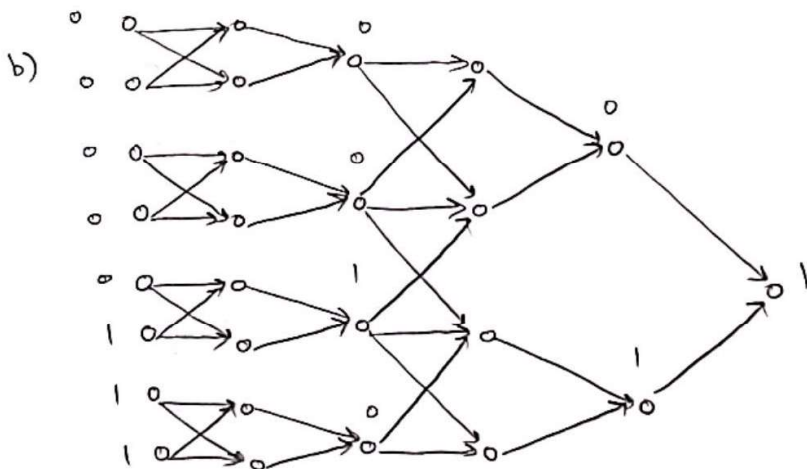
Legend:

•  $t^{(N)} = 1$

○  $t^{(N)} = 0$

As can be seen from the above representation, there is no decision boundary (i.e. plane) that can separate patterns that map to  $t^{(N)} = 1$  from those that map to  $t^{(N)} = 0$ .

Thus, there is no plane that can separate the two classes of outputs. Hence, the problem is not linearly separable.



The network is built from XOR units.

Each XOR unit has a hidden layer with 2 neurons.

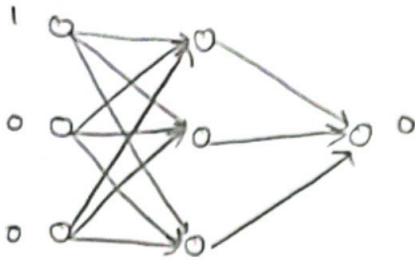
On the figure drawn, only the states of the inputs and outputs of the XOR units are shown, not those of hidden neurons.

In total the network has  $O(N)$  neurons.

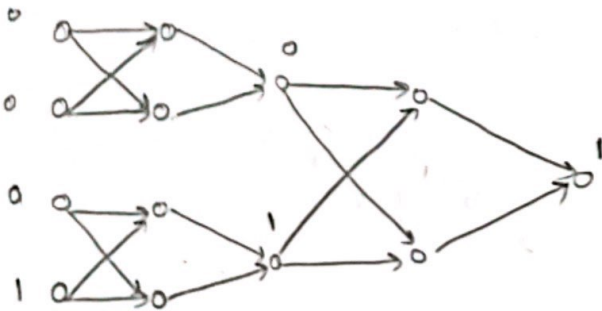
For  $N = 2^k$  input units with  $k = 1, 2, 3, \dots$  the whole network has  $3(N-1)$  neurons.

4) Boolean Function II:

For  $N=3$ :



For  $N=4$ :



The network is built from XOR units.

Each XOR unit has one hidden layer with 2 neurons.

Only the states of the inputs and outputs of XOR units are shown on the above figure.

For  $N=4$ , the network has  $3(N-1) = 3(4-1) = 9$  neurons.

(11)

5) Backpropagation:

$$\begin{aligned} a) \delta \omega^{(l)} &= -\eta \frac{\partial H}{\partial \omega^{(l)}} \\ &= -\eta 2 \cdot \frac{1}{2} (t - V^{(l)}) \left( -\frac{\partial V^{(l)}}{\partial \omega^{(l)}} \right) \\ &= \eta (t - V^{(l)}) \left( \frac{\partial V^{(l)}}{\partial \omega^{(l)}} \right) \\ &= \eta (t - V^{(l)}) \frac{\partial V^{(l)}}{\partial V^{(l)}} \frac{\partial V^{(l)}}{\partial \omega^{(l)}} \end{aligned}$$

where  $V^{(l)} = g(\omega^{(l)} V^{(l-1)} - \theta^{(l)}) = g(b^{(l)})$

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

It follows that:

$$\begin{aligned} \delta \omega^{(l)} &= \eta (t - V^{(l)}) \frac{\partial V^{(l)}}{\partial V^{(l)}} g'(b^{(l)}) V^{(l-1)} \\ &= \eta \delta^{(l)} V^{(l-1)} \end{aligned}$$

such that  $\delta^{(l)} = (t - V^{(l)}) \frac{\partial V^{(l)}}{\partial V^{(l)}} g'(b^{(l)})$

in other words  $\delta^{(l-1)} = (t - V^{(l)}) \frac{\partial V^{(l)}}{\partial V^{(l-1)}} g'(b^{(l-1)})$  verified

b) Evaluate  $\frac{\partial V^{(L)}}{\partial V^{(L-1)}}$  :

We have:  $V^{(L)} = g(\omega^{(L)} V^{(L-1)} - \Theta^{(L)}) = g(b^{(L)})$

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

- $\frac{\partial V^{(L)}}{\partial V^{(L-2)}} = g'(b^{(L)}) \omega^{(L)} \frac{\partial V^{(L-1)}}{\partial V^{(L-2)}} = g'(b^{(L)}) \omega^{(L)} g'(b^{(L-1)}) \omega^{(L-1)}$

- $\frac{\partial V^{(L)}}{\partial V^{(L-3)}} = g'(b^{(L)}) \omega^{(L)} \frac{\partial V^{(L-1)}}{\partial V^{(L-3)}}$

$$= g'(b^{(L)}) \omega^{(L)} \frac{\partial V^{(L-1)}}{\partial V^{(L-2)}} \frac{\partial V^{(L-2)}}{\partial V^{(L-3)}}$$

$$= g'(b^{(L)}) \omega^{(L)} \cdot g'(b^{(L-1)}) \omega^{(L-1)} g'(b^{(L-2)}) \omega^{(L-2)}$$

It follows that:

$$\frac{\partial V^{(L)}}{\partial V^{(L)}} = \prod_{k=L}^{L+1} [g'(b^{(k)}) \omega^{(k)}]$$

6) Reinforcement learning:

$$a) P(y|x) = \prod_{i=1}^M \begin{cases} p(b_i) & \text{for } y_i = +1 \\ 1 - p(b_i) & \text{for } y_i = -1 \end{cases}$$

$$= (1) p(b_i) + (-1) (1 - p(b_i))$$

$$= \frac{1}{1 + e^{-2\beta b_i}} - 1 + \frac{1}{1 + e^{-2\beta b_i}}$$

$$= \frac{2}{1 + e^{-2\beta b_i}} - 1$$

$$= \frac{2 - 1 - e^{-2\beta b_i}}{1 + e^{-2\beta b_i}}$$

$$= \frac{1 - e^{-2\beta b_i}}{1 + e^{-2\beta b_i}} \times \frac{e^{\beta b_i}}{e^{\beta b_i}}$$

$$= \frac{e^{\beta b_i} - e^{-\beta b_i}}{e^{\beta b_i} + e^{-\beta b_i}}$$

$$= \tanh(\beta b_i)$$

$$b) \langle \delta \omega_{mn} \rangle = \eta \frac{\partial \langle r \rangle}{\partial \omega_n}$$

Compare with  $\frac{\partial \langle r \rangle}{\partial \omega_n} = \langle r(x, y) [y - \tanh(b)] \rangle x_n$  leads to

$$\delta \omega_{mn} = r [y - \tanh(b)] x_n$$

- unbiased estimator of a parameter is the estimator whose expected value is equal to the value of the parameter.