

Q1.

Attribute 1: $H(y) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$

According to $H(y|A) = \sum p(x) H(y|A=x)$, we get

$$H(y|A_1) = \frac{4}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{1}{5} (-1 \log_2 1) = 0.8$$

$$H(y|A_2) = \frac{2}{5} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{2}{5} (-1 \log_2 1) = 0.551$$

$$H(y|A_3) = \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{3}{5} \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) = 0.951$$

So $IG(y|A_1) = H(y) - H(y|A_1) = 0.171$

$IG(y|A_2) = H(y) - H(y|A_2) = 0.42$

$IG(y|A_3) = H(y) - H(y|A_3) = 0.02$

So we should choose A_2 as the first split attribute

Attribute 2:

1) $A_2 = 0$

Example	A_1	A_2	Output y
x_1	1	0	0
x_2	1	1	0

Since x_1 and x_2 have the same output values 0, then we return 0.

2) $A_2 = 1$

Example	A_1	A_2	Output y
x_3	0	0	0
x_4	1	1	1
x_5	1	0	1

As shown above, x_3, x_4, x_5 do not have the same output values.

$$H(y) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$$

$$H(y|A_1) = \frac{1}{3} (-1 \log_2 1) + \frac{2}{3} (-1 \log_2 1) = 0$$

$$H(y|A_2) = \frac{1}{3} (-1 \log_2 1) + \frac{2}{3} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) = 0.667$$

So $IG(y|A_1) = H(y) - H(y|A_1) = 0.918$. $IG(y|A_2) = H(y) - H(y|A_2) = 0.251$

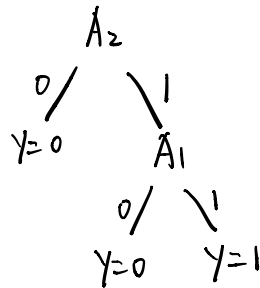
So we choose A_1 as the second split attribute for this subtree.

Attribute 3 =

1) $A_2 = 1$, $A_1 = 0$ Since x_2 is the only input and has output 0, then we return 0.

2) $A_2 = 1$, $A_1 = 1$ Since x_4 and x_5 have the same output values 1, then we return 1.

So the decision tree is:

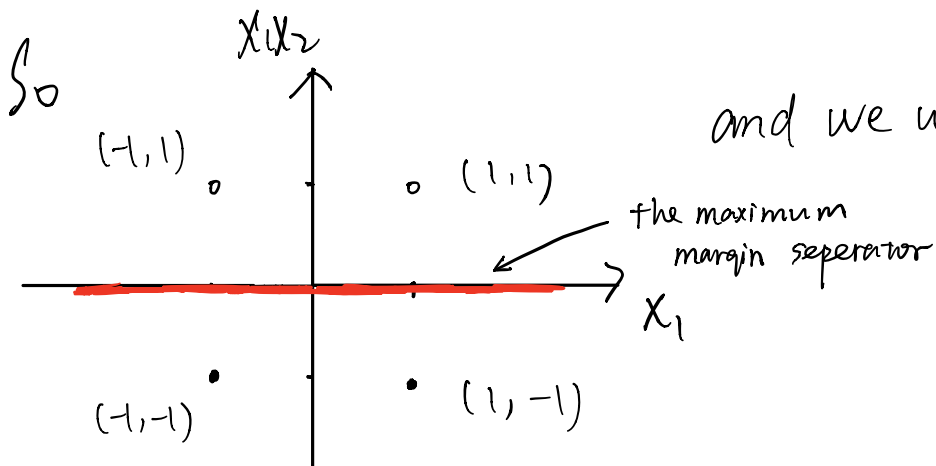


Q2.

(a)

It's easy to get, $[x_1, x_2]$

	x_1	$x_1 x_2$	Xor output
$[-1, -1]$	-1	1	0
$[-1, 1]$	-1	-1	1
$[1, -1]$	1	-1	1
$[1, 1]$	1	1	0

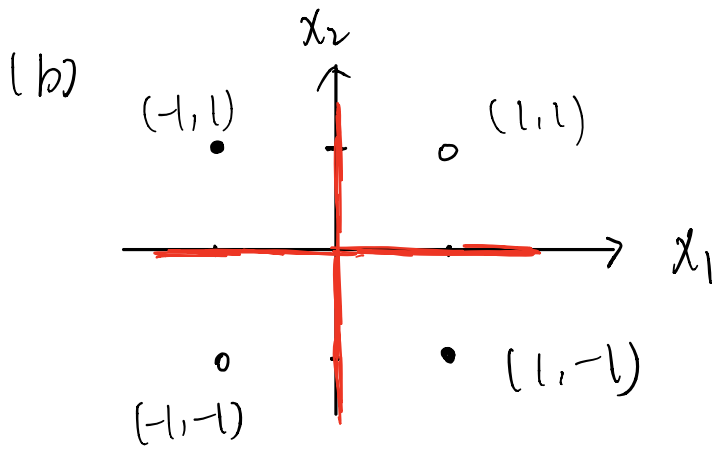


and we use "•" for output 1

"o" for output -1

the maximal margin separator is: $x_1 x_2 = 0$ which is emphasized by red.

The margin is 2 or 1 for only half of the distance.



From $x_1 x_2 = 0$, we get $x_1 = 0$ or $x_2 = 0$

The separating line is drawn by the red line.

Q3.

a. We use w^1, w^2 to represent the weights for the hidden layer and Output layer respectively.

For the one hidden layer:

the i -th input: $\sum_k w_{ik}^1 x_k$

the i -th output $h_i = g(\sum_k w_{ik}^1 x_k) = c \sum_k w_{ik}^1 x_k + b$.

For the output layer:

the j -th input: $\sum_n w_{jn}^2 h_n$

the j -th output $Q_j = g(\sum_n w_{jn}^2 h_n) = c \sum_n w_{jn}^2 (c \sum_k w_{nk}^1 x_k + b) + b$

$$= c^2 \sum_n \sum_k w_{jn}^2 w_{nk}^1 x_k + b (c \sum_n w_{jn}^2 + 1)$$

So there is a network with no hidden units that computes the same function. The parameters are shown as below.

the j -th input: $\sum_n \sum_k w_{jn}^2 w_{nk}^1 x_k$

$c_{\text{new}} = c^2$. $b_{\text{new}} = b (c \sum_n w_{jn}^2 + 1)$.

$$g(x)_{\text{new}} = c^2 x + b \left(c \sum_{\bar{n}} W_{\bar{j}\bar{n}}^2 + 1 \right)$$

b. For a network with n hidden layers, we can use conclusion from part a. to turn it into a network with $(n-1)$ hidden layers. In this way, by induction, we can turn this n -hidden-layer network into a network with no hidden layers.

And the new parameters will be:

$$g_{\text{new}}(x) = c_{\text{new}} x + b_{\text{new}} \quad \text{where}$$

$$c_{\text{new}} = c^n. \quad b_{\text{new}} = b \left(\sum_{\bar{z}=1}^{n-1} c^{\bar{z}} \cdot \sum_{\bar{j}=1}^{\bar{z}} W^{n+1-\bar{j}} + 1 \right).$$

c.

$$\text{For one hidden layer} = h \times n + n \times h = \geq 2hn$$

$$\text{For no hidden layer} = n \times n = n^2$$

Since $h \ll n$, the total number of weights becomes much larger, thus after transformation, the performance of the network becomes greatly worse.