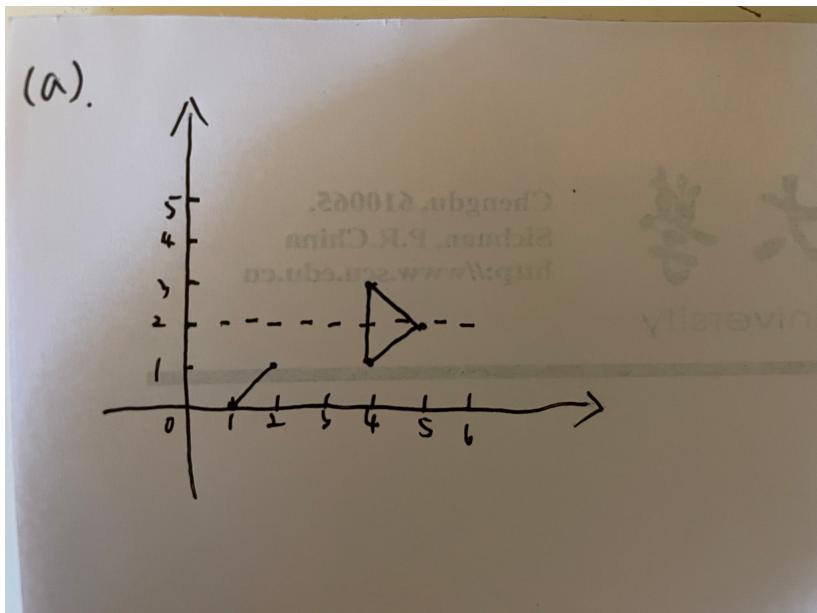


Name: Zhang Tong  
Student ID: z5235242

Question 5

(a)



If the training examples are linearly separable, there would be no intersection of the convex hull of the positive examples and negative examples. By observing the picture above, it is obvious that the range of  $h$  is  $0 \leq h < 4$ .

(b)

If  $h = 1$ , the minimum distance of the positive examples and negative examples would be 2 ( $x_5$  and  $x_1$ ). Thus the functional margin of the training dataset would be  $2/2 = 1$ .

(c)

We just need to ensure that the minimum distance of the positive examples and the negative examples is still 2.

According to the picture in (a), we have  $0 \leq h \leq 2$ .

(d)

$x_5$  and  $x_1$ . If the other points are removed, the boundary would not change since the functional margin would not change. But if  $x_1$  or  $x_5$  is removed, the functional margin would change because the minimum distance of the positive and negative examples would change.

Question 2

(a)

Thus clusterA : {1, 3, 5, 8}

2.

(a)  $d(2,1) = \sqrt{(5.98 - 1.90)^2 + (2.68 - 0.97)^2} = \sqrt{16.6464 + 2.9241} = 4.42$

$d(3,1) = \sqrt{0.6525} = 0.81$        $d(2,10) = 0.74$

$d(4,1) = 3.50$        $d(3,10) = 2.41$

$d(5,1) = 0.90$        $d(4,10) = 2.53$

$d(7,1) = 4.79$        $d(5,10) = 4.14$

$d(8,1) = 1.05$        ~~$d(6,10) = 2.09$~~

$d(9,1) = 4.19$        $d(7,10) = 0.72$

$d(2,6) = 2.82$        $d(8,10) = 3.26$

$d(3,6) = 3.51$        $d(9,10) = 2.88$

$d(4,6) = 6.73$

~~$d(5,6)$~~        $d(5,6) = 3.53$

$d(7,6) = 1.70$

$d(8,6) = 2.97$

$d(9,6) = 0.86$

ClusterB: {4, 6, 9}

ClusterC: {2, 7, 10}

(b)

Center of A:

$$X_1 = (1.9 + 2.68 + 1.54 + 2.46) / 4 = 2.145$$

$$X_2 = (0.97 + 1.18 + 1.80 + 1.86) / 4 = 1.4525$$

$$(2.145, 1.4525)$$

Center of B:

$$X_1 = (3.14 + 3.82 + 3.17) / 3 = 3.3767$$

$$X_2 = (4.24 + 4.50 + 4.96) / 3 = 4.5667$$

$$(3.3767, 4.5667)$$

Center of C:

$$X_1 = (5.98 + 5.74 + 5.44) / 3 = 5.72$$

$$X_2 = (2.68 + 3.84 + 3.18) / 3 = 3.2333$$

(5.72, 3.2333)

(c)

Yes we should. Because that we use Manhattan distance, the lower distance does not always represents that the point is more close to the center.

Question 3

(a)

Q3.

(a)  $(P_1, P_2)$  has the lowest value (0.1)

So we have

$\underline{P_{12} \quad P_3 \quad P_4 \quad P_5}$

$P_{12} \quad 0 \quad 0.39 \quad 0.57 \quad 0.48$

$P_3 \quad 0 \quad 0.44 \quad 0.85$

$P_4 \quad 0 \quad 0.76$

$P_5 \quad 0$

$$\min = (P_{12}, P_4).$$

$\underline{P_{124} \quad P_3 \quad P_5}$

$P_{124} \quad 0 \quad 0.44 \quad \cancel{0.54}$

$P_3 \quad 0 \quad 0.85$

$P_5 \quad 0$

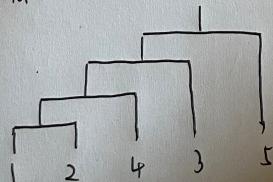
$$\min = (P_{124}, P_3)$$

$\underline{P_{1243} \quad P_5}$

$P_{1243} \quad 0 \quad 0.56$

$P_5 \quad 0.56 \quad 0$

Then we have



Question 4

Q4. (a) id	$a_1$	$a_2$	$a_3$	class
1	T	T	L	Y
2	F	T	H	Y
3	T	F	H	N
4	F	F	M	Y
5	F	T	H	N
6	F	T	M	N
7	F	F	H	N
8	T	F	H	Y
9	F	T	M	N

$$\begin{aligned}
 \text{(b)} \quad \text{gini split}(a_1) &= \frac{4}{9} \times \left[ 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \right] + \frac{5}{9} \times \left[ 1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2 \right] = 0.34 \\
 \text{gini split}(a_2) &= \frac{5}{9} \times \left[ 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 \right] + \frac{4}{9} \times \left[ 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right] = 0.49 \\
 \text{gini split}(a_3) &= \frac{1}{9} \times (1 - 1^2 - 0^2) + \frac{4}{9} \times \left[ 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 \right] + \frac{4}{9} \times \left[ 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 \right] \\
 &= 0.39
 \end{aligned}$$

So  $a_1$  is the first splitting node (Min Gini).

$$\text{(c)} \quad P(Y) = \frac{4}{9}, \quad P(N) = \frac{5}{9}$$

Add 0.5 smoothing,

$$P(a_1=T | \text{Class}=Y) = (3+0.5)/(4+0.5 \times 2) = 0.7$$

$$P(a_1=T | \text{Class}=N) = (4+0.5)/(5+0.5 \times 2) = 0.25$$

$$P(a_2=F | \text{Class}=Y) = (2+0.5)/(4+0.5 \times 2) = 0.5$$

$$P(a_2=F | \text{Class}=N) = (2+0.5)/(5+0.5 \times 2) = 0.42$$

$$\text{So } P(a_1=T, a_2=F, a_3=x | Y) = 0.7 \times 0.5 = 0.35$$

$$P(a_1=T, a_2=F, a_3=x | N) = 0.25 \times 0.42 = 0.105$$

$$P(a_1=T, a_2=F, a_3=x | Y) P(Y) = 0.35 \times \frac{4}{9} = 0.16,$$

$$P(a_1=T, a_2=F, a_3=x | N) P(N) = 0.105 \times \frac{5}{9} = 0.06, \quad 0.16 > 0.06$$

Thus the prediction should be Y.

## Question 1

Q1.

(a) 8 cuboids.

None, A, B, C, AB, AC, BC, ABC.

(b) The table contains 3 tuples, if ~~the~~ the 3 tuples are same,  
they could be merged as one. Thus the minimum ~~possible~~ number of  
tuples in the complete data cube should be  $2 \times 2 \times 2 = 8$ .  
Considering ~~that~~ if A, B, C have different values in the ~~the~~ 3 tuples,  
the maximum number of tuples would be  $4^3 = 64$ .

(c) A=1, B=1  $\Rightarrow$  one tuple

A=1, B=2  $\Rightarrow$  one tuple

A=2  $\Rightarrow$  one tuple

So 3 times.