

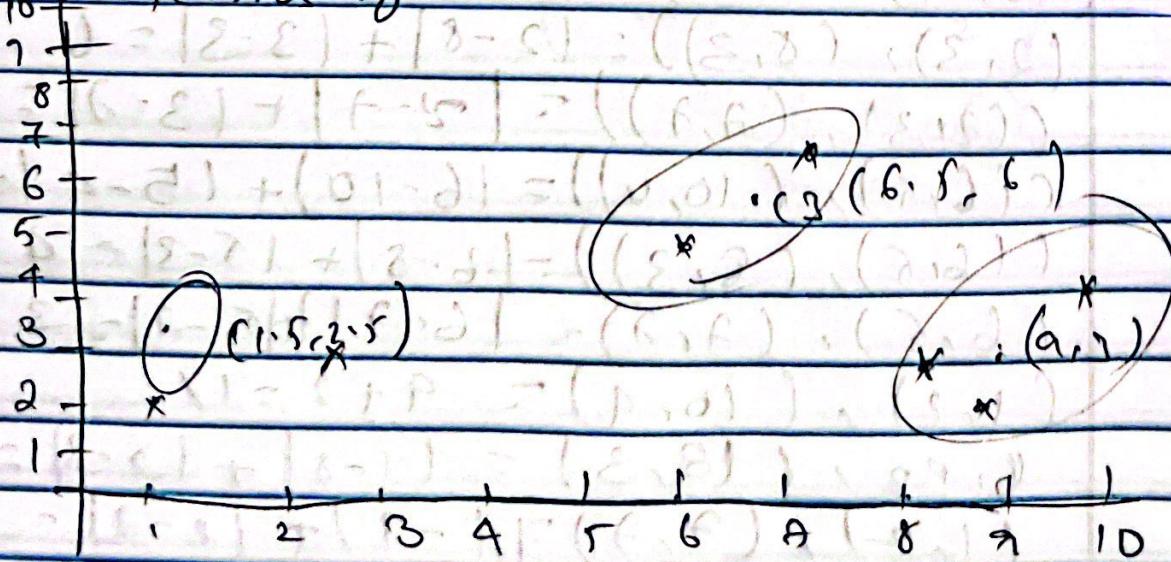
## Data mining: 2

### 2/1 K-means clustering

Given 7 sensor nodes deployed in a network, The location of three sensor are describe in terms of x and y coordinate the given points are

	x	y	
1	10	4	3 (6, 1) 3
2	8	3	4 (5, 0) 5
3	4	1	5 (2, 6) 1
4	2	3	6 (1, 2) 3
5	6	5	7 (0, 1) 2
6	1	2	
7	9	2	

7 sensor nodes deployed in a network locations of these sensors are described in terms of



Distance	$m_1(10, 4)$	$m_2(8, 3)$	$m_3(7, 7)$
1 (10, 4)	0	3	6
2 (8, 3)	3	0	5
3 (7, 7)	6	5	0
4 (2, 3)	9	6	9
5 (6, 5)	5	4	3
6 (1, 2)	11	8	11
7 (9, 2)	3	2	7

By calculating manhattan distance  
from points to centroids

$$d_{\text{mny}} = \sum_{i=1}^P |x_i - y_i|^P \quad \text{where } P=1$$

$$(10, 4), (8, 3) : |10-8| + |4-3| = 2+1 = 3$$

$$(10, 4), (7, 7) : |10-7| + |4-7| = 3+3 = 6$$

$$(8, 3), (7, 7) : |8-7| + |3-7| = 1+4 = 5$$

$$(10, 4), (2, 3) : |10-2| + |4-3| = 8+1 = 9$$

$$(2, 3), (5, 3) : |2-5| + |3-3| = 6$$

$$(2, 3), (7, 7) : |2-7| + |3-7| = 9$$

$$(6, 5), (10, 4) : |6-10| + |5-4| = 5$$

$$(6, 5), (8, 3) : |6-8| + |5-3| = 4$$

$$(6, 5), (7, 7) : |6-7| + |5-7| = 3$$

$$(1, 2), (10, 4) : 9+2 = 11$$

$$(1, 2), (8, 3) : |1-8| + |2-3| = 8$$

$$(1, 2), (7, 7) : |1-7| + |2-7| = 11$$

$$(9, 12, 1, 10, 4) = |9-10| + |2-4| = 1+2=3$$

$$(9, 2), (8, 3), |9-8| + |2-3| = 1+4=2$$

$$(9, 2), (7, 7) = |9-7| + |2-7| = 2+5=7$$

$$C_1 = \{10, 4\} \quad y = 214$$

$$C_2 = \{2, 4, 6, 7\}$$

$$C_3 = \{3, 5\}$$

mean of  $C_1$  is same it won't change  
as it only contains only one

$$\text{mean} = \left( \frac{8+2+1+9}{4}, \frac{3+3+2+1}{4} \right)$$

mean  $C_3$  is somewhere

$$\text{mean} = \left( \frac{4+6}{2}, \frac{1+5}{2} \right)$$

$$= (5, 3)$$

$$\left( \frac{4+8+1}{3}, \frac{2+3+1}{3} \right) = (5, 3)$$

$$(5, 3) = 1$$

Distance	$m_1(10, 4)$	$m_2(5, 2, 5)$	$m_3(6.5, 6)$
1(10, 4)	0	6.5	5.5
2(8, 3)	3	3.5	4.5
3(7, 2)	6	6.5	1.5
4(2, 3)	9	3.5	9.5
5(6, 5)	5	3.5	1.5
6(1, 2)	11	4.5	9.5
7(9, 2)	3	4.5	6.5

New centroid

$$C_1 = \{1, 1, 7\}$$

$$C_2 = \{4, 5\}$$

$$C_3 = \{3, 5\}$$

As the Sensors changed to different  
Centroids, the mean of  $C_1$  &  $C_2$   
will change

$C_3$  remains as previous

$$C_3 = \left( \frac{7+6}{2}, \frac{7+5}{2} \right) = (6.5, 6)$$

$$C_1 = \left( \frac{10+8+9}{3}, \frac{4+3+2}{3} \right)$$

$$C_1 = (9, 3)$$

$$C_2 = \left\{ \frac{2+1}{2}, \frac{3+2}{2} \right\}$$

$$C_2 = (1.5, 2.5)$$

Centroids are:  $C_1(9, 3)$

$$C_2 = (1.5, 2.5) \quad C_3 = (6.5, 6)$$

Distance	$m_1(9, 3)$	$m_2(1.5, 2.5)$	$m_3(6.5, 6)$
1 (10, 4)	2	10	5.5
2 (8, 3)	1	7	4.5
3 (7, 3)	6	10	1.5
4 (2, 3)	4	1	7.5
5 (6, 5)	5	7	1.5
6 (1, 2)	9	1	9.5
7 (9, 2)	1	8	6.5

for the new centroid the Sensors are remaining

$$C_1 \text{ has } (1, 2, 3)$$

$$C_2 \text{ has } (4, 6)$$

$$(3 - (3, 5))$$

So, we can lock the new centroid as main centroids

$$C_1(9, 3) \quad C_2(1.5, 2.5) \quad C_3(6.5, 6)$$

(3)

## Hierarchical Clustering

	A	B	C	D	E	F
A	0.00					
B	0.71	0.00				
C	5.66	4.95	0.00			
D	3.61	2.92	2.24	0.00		
E	4.24	3.54	1.41	1.00	0.00	
F	3.20	2.50	2.50	0.50	1.12	0.00

Step 1 (i) Single-link hierarchical clustering

As F-D	F-D	F-D	A	B	C	E
having the	F-D	0.00	3.20	0.71	2.24	1.00
shorter	A		0.00	0.71	5.66	4.24
distance	B		0.00	4.95	3.54	1.41
	C			0.00		0.00
	E					

	F, D	A, B	C	E
F, D	0	2.50	2.24	1.00
A, B		0	4.95	3.54
C			0	1.41
E				0

Hence we are combining A & B because they are having the shortest distance whole distance matrix.

Step 3:- From the distance matrix as (F, D), E) are having the shortest distance hence we combine (F, D, E)

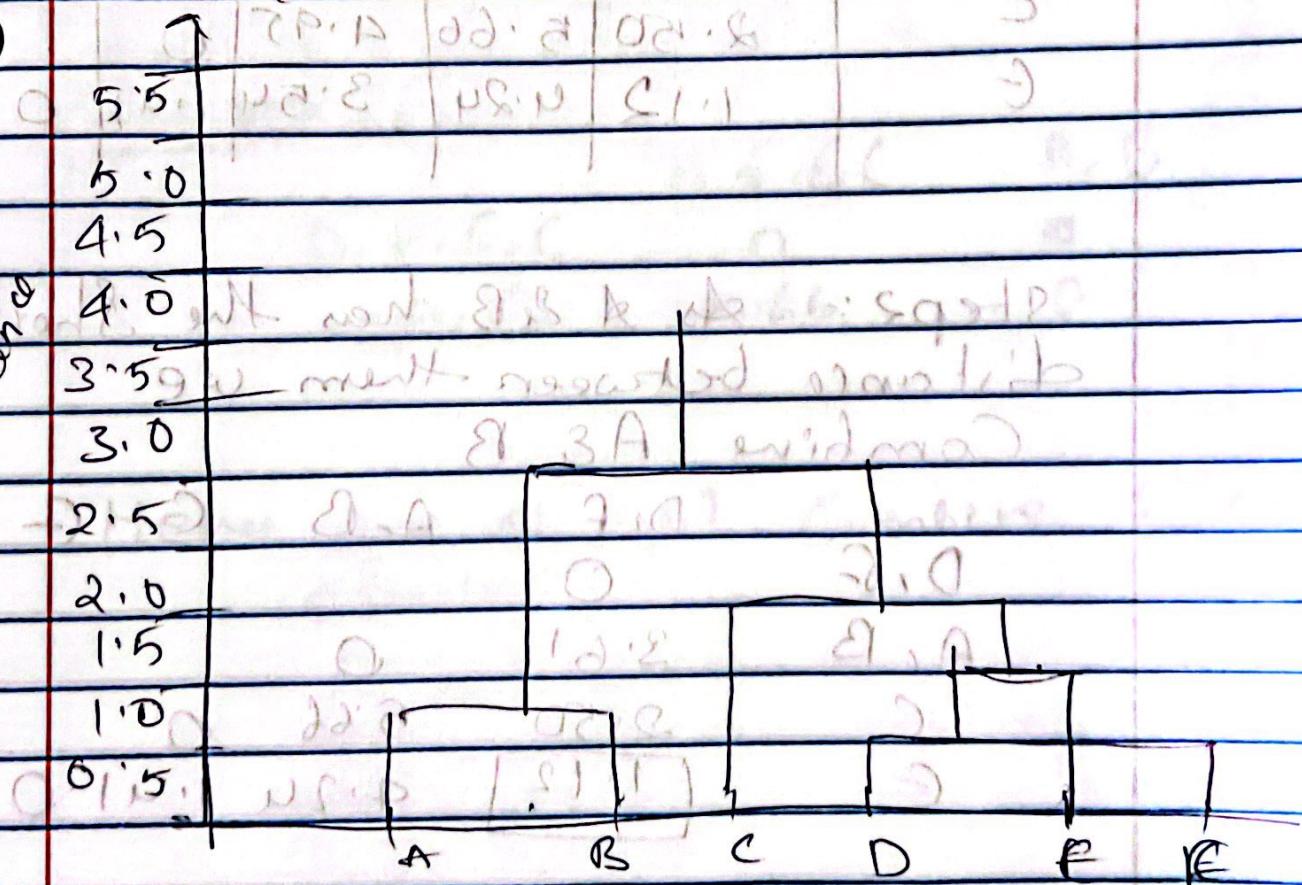
	F, D, E	A, B	C
F, D, E	0	2.50	1.41
(A, B)		0	4.95
C			0

Step 4:- Here we combine (F, D, C) & C, which are having shortest distance from the above distance matrix

(F, D, E, C)  $\rightarrow$  (A, B)

(A, B)  $\rightarrow$  0

Step 3:- At last we combine (C A, D), E  $\rightarrow$  (A, B)



(b) Complex link hierarchical clustering

~~Step 1~~

As from the table

Step 1'. As D & F having shortest distance in whole data matrix

we combine D & F.

	D, F	A	B	C	E
D, F	0				
A	3.61	0			
B	2.92	0.91	0		
C	2.50	5.66	4.95	0	
E	1.12	4.24	3.54	1.41	0

Step 2: As A & B has the shortest distance between them we combine A & B

	D, F	A, B	C, E
D, F	0		
A, B	3.61	0	
C	2.50	5.66	0
E	1.12	4.24	1.41

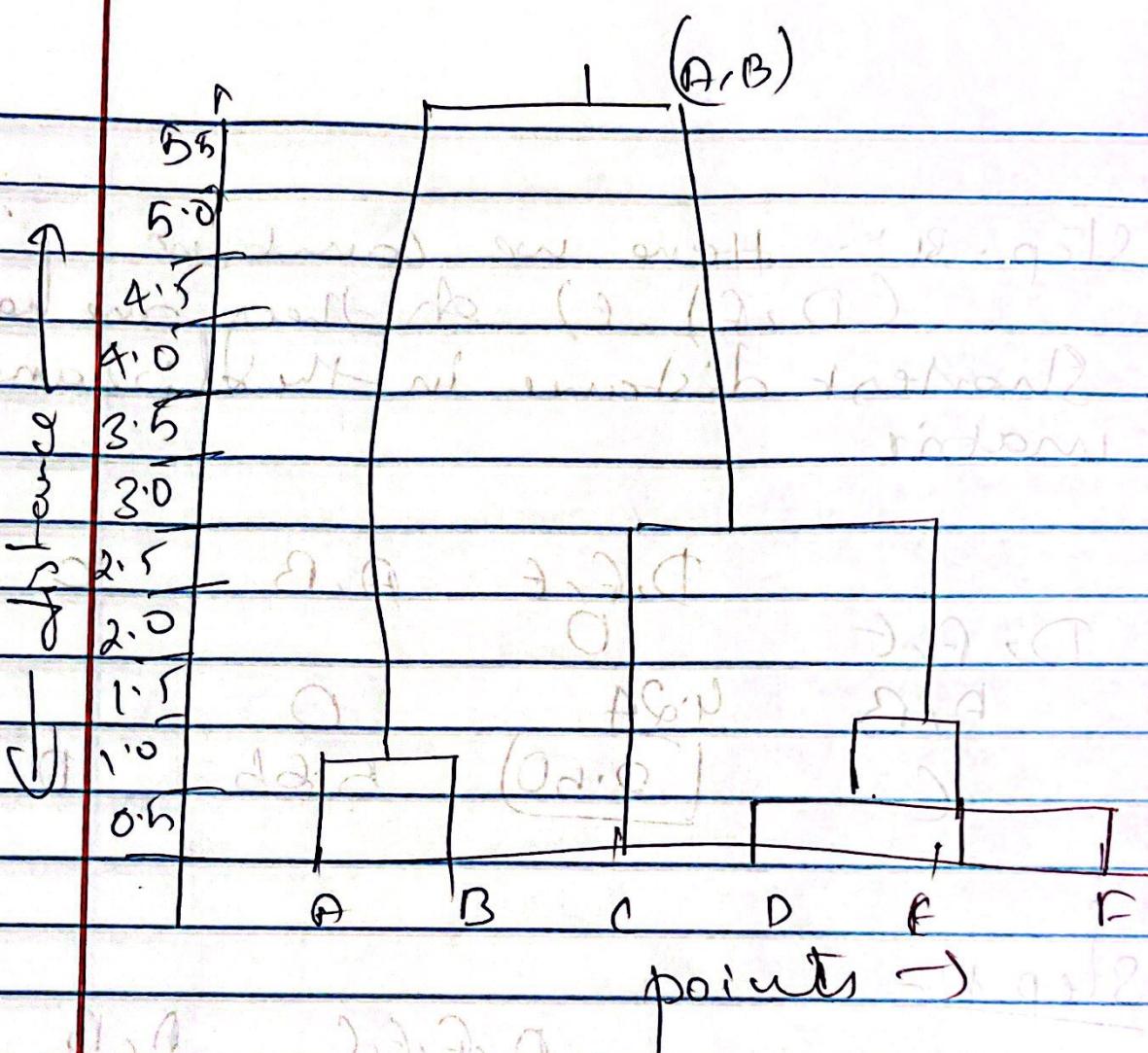
Step - 3 :- Here we combining  
 $(D, F, E)$  as they are having  
 shortest distance in the distance  
 matrix

	$D, F, E$	$A, B$	$C$
$D, F, E$	0		
$A, B$	4.24	0	
$C$	2.50	5.66	0

Step 4 :- Iteration

	$D, F, E, C$	$A, B$
$D, F, E, C$	0	9
$A, B$	15.66	0

Hence we at last combining to one cluster.



4)

## Association Rule.

Transaction id

1  
2  
3  
4  
5  
6

Items Purchased

A, B, D, E

B, C, D

A, D, F

A, B, C, D, F

A, B

C, E; F

Q) the apriori algorithm on the above table with min support = 0.25

Items	Support	min-Support
A	4	0.6
B	4	0.6
C	3	0.5
D	4	0.6
E	2	0.3
F	3	0.5

- i)  $\emptyset$  (All items are generated)
- ii) No 'items' are pruned -  $\emptyset$
- iii)  $\emptyset$ -No items are pruned after scan
- iv) {A, B, C, D, E, G, F}

### 2nd Iteration

Here  $k=2$ , so take all the item set of size 2 from above

Items for support min. Support

{A, B}

3

0.5

{A, C}

0.16 x

{A, D}

0.15

{A, E}

0.16 x

{A, F}

0.3

{B, C}

0.3

{B, D}

0.5

{B, E}

0.16

{B, F}

0.16

{C, D}

0.3

{C, E}

0.16

{C, F}

0.3

{D, E}

0.16

$$\begin{array}{ll} \text{I} D, F \text{Y} & ? \\ \text{II} E, F \text{Y} & 1 \end{array}$$

$\frac{2}{6} = 0.3$   
 $\frac{1}{6} = 0.16 \times$

- i)  $\phi$  - All items are generated
- ii)  $\phi$  - No items are pruned before scan
- iii)  $\{A, C\}, \{A, E\}, \{B, E\}, \{B, F\}, \{C, E\}, \{D, E\}, \{E, F\}$   
 are pruned after scan,  
 as they don't meet the  
 min support of '0.25'.
- iv)  $\{A, B\}, \{A, D\}, \{A, F\}, \{B, C\},$   
 $\{B, D\}, \{C, D\}, \{C, F\}, \{D, F\}$

After - Pruning the Item Set  
 Present

items	min-support
$\{A, B\}$	0.5
$\{A, D\}$	0.5
$\{A, F\}$	0.3
$\{B, C\}$	0.3
$\{B, D\}$	0.5
$\{C, D\}$	0.3

$$\frac{f(C, F)}{f(D, F)} = 0.3$$

## II rd iteration

Hence  $K=3$ , so take all the items set of length  $l=3$  from the pruned set

Items	Pruned	min support
-------	--------	-------------

$\{A, B, D\}$	No	0.3
$\{A, B, F\}$	Pruned	
$\{A, D, F\}$	No	0.3
$\{A, B, C\}$	pruned	
$\{B, C, D\}$	No	0.3
$\{C, D, F\}$	No	0.16
$\{A, C, D\}$	pruned	
$\{A, C, F\}$	pruned	
$\{B, C, F\}$	pruned	
$\{B, D, F\}$	pruned	

Items not even generated are

$\{A, B, E\}$ ,  $\{A, C, E\}$ ,  $\{A, D, E\}$ ,  $\{A, E, F\}$   
 $\{B, C, E\}$ ,  $\{D, B, E\}$ ,  $\{B, E, F\}$

$\{C, D, E\}$ ,  $\{C, E, F\}$ ,  $\{D, E, F\}$

ii) Items pruned without Scan are

$\{A, B, F\}$ ,  $\{A, B, C\}$ ,  $\{A, C, D\}$

$\{A, C, F\}$ ,  $\{B, C, F\}$ ,  $\{B, D, F\}$

iii) Items pruned after Scan  
 $\{C, D, F\}$

iv) frequent items sets are

$\{A, B, D\}$ ,  $\{A, D, F\}$ ,  $\{B, C, D\}$

$\{A, B, D\}$ ,  $\{A, D, F\}$ ,  $\{B, C, D\}$

are find frequent sets are  
we can't find frequent  
sets further.

The frequent items are  $\{A, B\}$ ,  
 $\{A, D\}$ ,  $\{A, F\}$ ,  $\{B, C\}$ ,  $\{B, D\}$ ,  
 $\{C, D\}$ ,  $\{C, F\}$ ,  $\{D, F\}$

$$A \rightarrow B = \frac{\text{Confidence}}{\text{Support}(A)} = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

The Association rules are:-

$$i) A \rightarrow B \Rightarrow \text{Confidence} = \frac{\text{Support}(A, B)}{\text{Support}(A)} = 0.75$$

$$ii) B \rightarrow A \Rightarrow \text{Confidence} = \frac{\text{Support}(A, B)}{\text{Support}(B)} = 3/4 = 0.75$$

$$iii) A \rightarrow D \Rightarrow \text{Confidence}$$

$$= \frac{\text{Support}(A, D)}{\text{Support}(A)} = 3/9 = 0.75$$

$$iv) D \rightarrow A \Rightarrow \text{Confidence}$$

$$= \frac{\text{Support}(D, A)}{\text{Support}(D)} = 0.75$$

$$v) F \rightarrow A \Rightarrow \text{Confidence}$$

$$= \frac{\text{Support}(F, A)}{\text{Support}(F)} = 0.66$$

$$vi) A \rightarrow F \Rightarrow \text{Confidence} = \frac{\text{Support}(A, F)}{\text{Support}(A)} = 2/4 = 0.5$$

$$vii) B \rightarrow C \Rightarrow \text{Confidence} = \frac{\text{Support}(B, C)}{\text{Support}(B)} = 0.5$$

viii)  $f(c) \rightarrow q(BY)$

$$\text{Confidence} = \frac{\sigma(c, B)}{\sigma(c)} = 0.6$$

$\neg BY \rightarrow DS$

$$\text{Confidence } \neg(c, D) = \frac{3}{4} = 0.75$$

(x)  $\neg D \rightarrow q(BY)$

$$\text{Confidence } \neg(c, B) = \frac{3}{4} = 0.75$$

(xi)  $f(cY \rightarrow \neg D)$

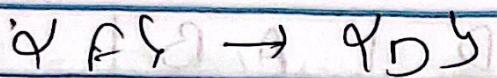
$$\text{Confidence} = \frac{\sigma(c, D)}{\sigma(c)} = 0.66$$

(xii)  $\neg D \rightarrow f(cY)$

$$\text{Confidence} = \frac{\sigma(D, C)}{\sigma(D)} = 0.5$$

(xiii)  $f(cY \rightarrow \neg F)$

$$\text{Confidence} = \frac{\sigma(c, F)}{\sigma(c)} = 0.66$$



Confidence:  $\frac{\sigma(F,D)}{\sigma(F)} = 7^2/3 = 0.66$

# Data mining Assignment

## Classification Evaluation

(1) we have developed two classifiers C1 and C2 which determine the posterior probability of movies nominated for the academy awards. The given we can take nominated as (+) and not nominated as (-) classes

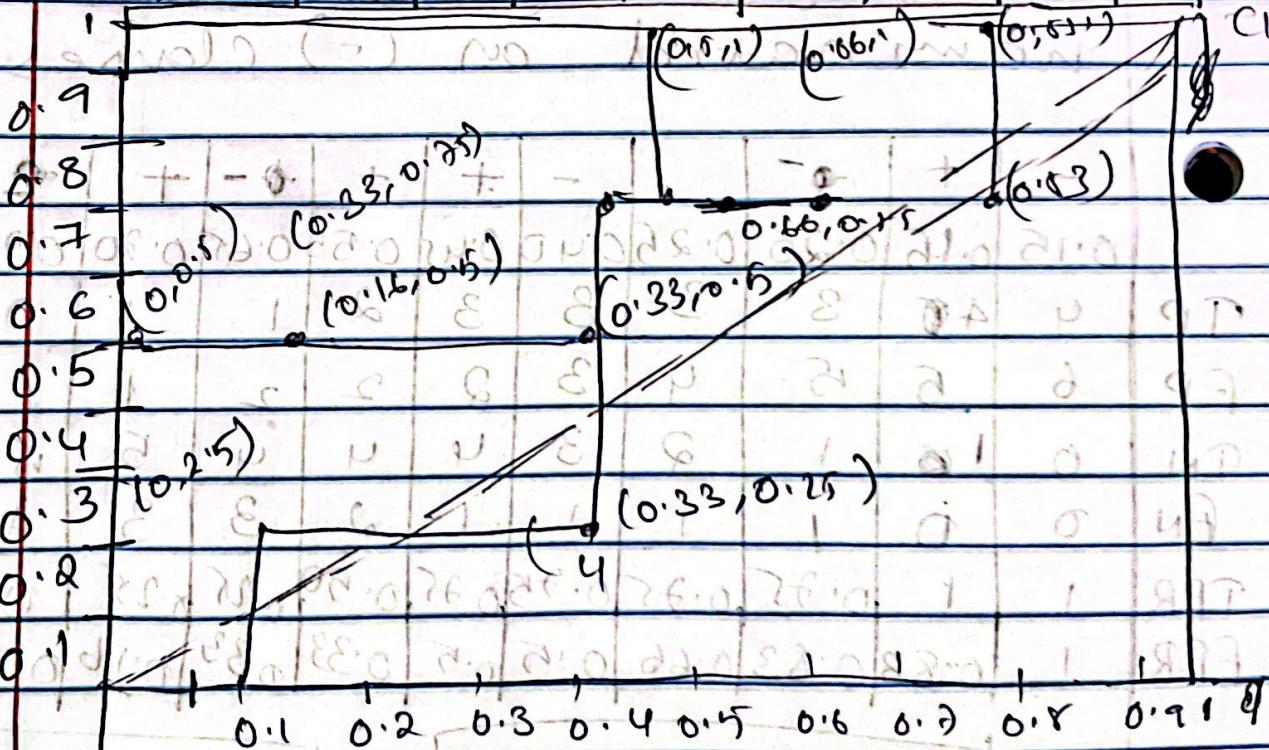
	-	+	-	-	-	+	+	-	-	+	-
	0.15	0.15	0.20	0.25	0.40	0.45	0.55	0.65	0.70	0.75	0.100
TP	4	4	3	3	3	3	2	1	1	0	0
FP	6	5	5	4	3	2	2	2	1	1	0
TN	0	1	1	2	3	4	4	4	5	5	6
FN	0	0	1	1	1	1	2	3	3	4	4
TPR	1	1	0.75	0.75	0.75	0.75	0.50	0.25	0.25	0	0
FPR	1	0.83	0.83	0.66	0.15	0.5	0.33	0.33	0.16	0.16	0

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

Table for C2

	-	-	+	-	+	-	-	+	+	+	
θ	0.25	0.30	0.35	0.55	0.60	0.65	0.75	0.80	0.85	0.90	1.00
TP	4	4	4	4	3	3	2	2	2	1	0
FP	6	5	4	3	3	2	2	1	0	0	0
TN	0	1	2	3	3	4	4	5	6	6	.6
FN	0	0	0	0	1	1	2	2	2	3	4
TPR	1	1	1	1	0.75	0.75	0.5	0.5	0.5	0.25	0
FPR	0.83	0.66	0.5	0.5	0.33	0.33	0.16	0	0	0	0



b) Calculate the F-measure for all the threshold values in U

As we know that

$$F\text{-measure} = \frac{2(TP)}{\frac{2(TP)}{TP+FN+FP}}$$

At  $\theta$  is 0.15

$$F\text{-measure} = \frac{2(4)}{1+2+4} = \frac{8}{7} = 0.57$$

At  $\theta$  is 0.45

$$F\text{-measure} = \frac{2(3)}{1+2+3+0} = \frac{6}{9} = 0.66$$

Similarly calculate for other values  
 ' $\theta$ ' value at  $0.45$  given as  
 the highest F-measure

$$\theta = 0.97, 8 - 0.247 = 0.75$$

Confusion matrix at  $\theta = 0.45$

		Predicted class	
		Nominated	not-nominated
Actual class	Nominated (+)	TP(a) 3	FN(b) 1
	class=NOT-nominated (-)	FP(c) 2	TCN(d) 4

Precision:

$$P = \frac{a}{a+c} = \frac{3}{3+2} = \frac{3}{5}$$

recall

$$R = \frac{a}{a+b} = \frac{3}{3+1} = \frac{3}{4}$$

F-measure ( $r$ ):

$$F = \frac{2a}{2a+b+c} = \frac{2(3)}{2(3)+1+2} = \frac{6}{9}$$

① At  $\theta = 0.65$  in C1

The matrix values are

$$FP(a) = FN(b) = 3, FP(c) = 2$$

Given the weight matrix  
with

$$w_1=1, w_2=10, w_3=4, w_4=1$$

We know that

$$\text{Weighted accuracy} = \frac{w_1a + w_4d}{w_1a + w_2b + w_3c + w_4d}$$

$$\frac{1(1) + 4(1)}{1(1) + 3(10) + 2(4) + 4(1)}$$

5  
43  
 $\underline{= 0.116}$

At  $\theta = 0.65$  in C2

The confusion matrix values are  
 $TP(a) = 3$ ,  $FN(b) = 1$      $FP(c) = 1$

$TN(d) = 4$

River

$$w_1 = 1, w_2 = 10, w_3 = 4, w_4 = 1$$

Now,

$$w_A = \frac{1(3) + 4(1)}{(3) + 10(1) + 4(1) + 4(1)}$$

$$\Rightarrow \frac{7}{25} \Rightarrow 0.28$$

C2 has provided the better accuracy at  $\theta = 0.65$

We can conclude 'C2' classifier is better.