

Mean deviation

$$\frac{\sum |x_i - \bar{x}|}{N}$$

Std deviation  $\rightarrow$  RMS

$$\text{var} = \text{MS} = \frac{\sum (x_i - \bar{x})^2}{N}$$

PAGE NO.	/ /
DATE	/ /

MODULE : 3

Q. 45 47 52 52 53 55 56 58 62 80

10 pts total

• Mean = 56

Make 2 sets

[45, 47, 52, 52, 53], [55, 56, 58, 62, 80]

keep separate

• Median:  $(53 + 55)/2 = 54$

• Mode = 52 unimodal

• Max value : 80

• Min value : 45

• Mid range :  $(\text{max} + \text{min})/2 = 62.5$

• Quartile

•  $Q_1$  = median of lower half part of data = 52

•  $Q_3$  = median of upper half part of data = 58

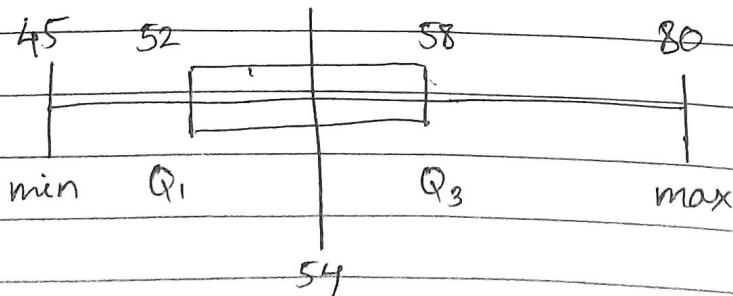
•  $IQR = Q_3 - Q_1 = \text{ht of box} = 58 - 52 = 6$

• Outlier =  $1.5 IQR = 9$  (beyond the quartiles)

• Semi inter quartile range =  $(Q_3 - Q_1)/2 = 6/2 = 3$

• Mid quartile range =  $(Q_3 + Q_1)/2 = 55$

Med Q<sub>2</sub>



if odd no. of nos, make & set exchangeable

PAGE NO.

DATE

Q. 30 36 47 50 52 52 56 60 63 70 70

→ • Mean = 58

• Median = 54

• Mode = 52, 70 bimodal

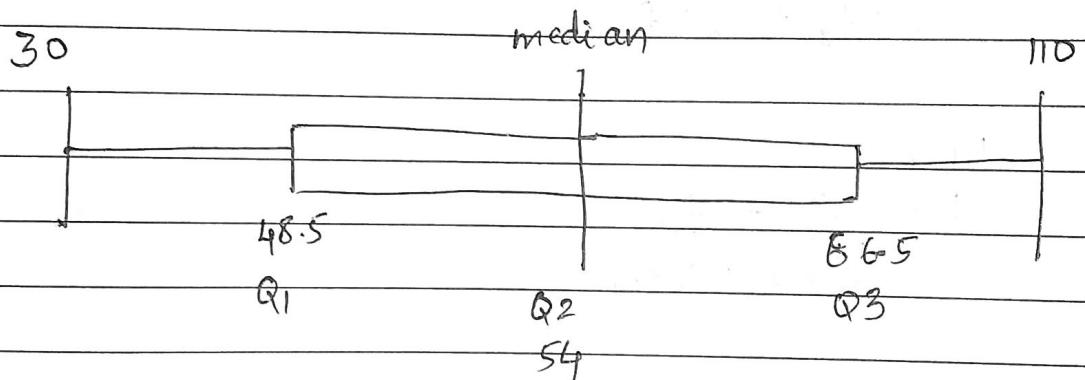
• Max val = 110

• Min val = 30

• Mid range =  $(\text{max} + \text{min}) / 2 = 70$

•  $Q_1 = (47 + 50) / 2 = 48.5$

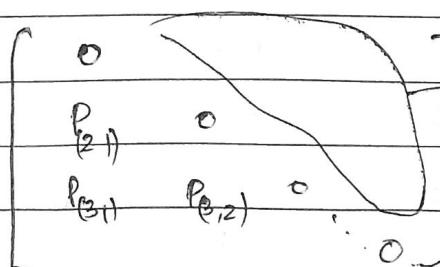
•  $Q_3 = (63 + 70) / 2 = 66.5$



\* (p) Similarity : max [0, 1]

(q) Dissimilarity : often 0

Proximity: similarity / dissimilarity



These are considered to be same (symmetric) matrix

Sim	Max	Min
Dis	Min	Max

Symmetric binary : if the results  $\Rightarrow$  ~~value~~<sup>attri</sup> can have 2 separate  
 Pass  $\rightarrow$  next class X<sup>sym</sup> binary attri Fail  $\rightarrow$  retain  $\Rightarrow$  PAGE NO. X sym bin  
 DATE

Male Female  $\Rightarrow$  symmetric binary

Obj	Gender	Similarity	Dissimilarity
Ram	M 1	$p(Ram, sita) = 0$	$q = 1$
Sita	F 0	$p(Ram, Laxman) = 1$	$q = 0$
Laxman	M 1		

$$q(i,j) = 1 - p(i,j)$$

\* similarity measure with symmetric binary

$\Rightarrow$  S<sub>b</sub> coeff

$$S = \frac{\text{No. of matching attri values}}{\text{Total nos. of attris}}$$

OR

$$S = \frac{f_{00} + f_{11}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

D<sub>b</sub> c

$$D = \frac{f_{01} + f_{10}}{f_{00} + f_{01} + f_{10} + f_{11}}$$

$$f_{01} + f_{10}$$

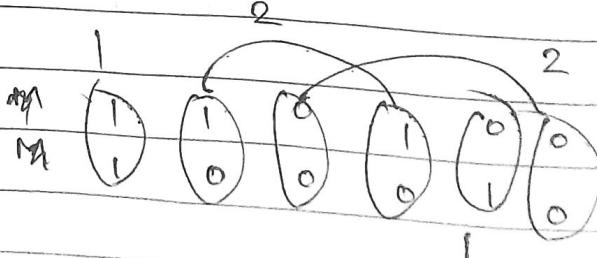
eg

Obj	Gender	Food	Caste	Edu	Hobby	Job
Hari	M	V	M	L	C	N
Ram	M	NV	M	I	T	N
Tomu	F	NV	H	L	C	Y

$$\begin{matrix} M-1 & V-1 & M-0 & L-1 & C-0 & N-0 \\ F-0 & NV-0 & H-1 & I-0 & T-1 & Y-1 \end{matrix}$$

Contingency matrix/table

	R <sub>0</sub>	R <sub>1</sub>
H <sub>0</sub>	2	1
H <sub>1</sub>	2	1



$$S(\text{Hari, Ram}) = \frac{1+2}{1+2+1+2} = 0.5$$

\* Proximity measure with asymmetric coeff  
Jaccard coeff

similarity  $J = \frac{\text{No. of matching presence}}{\text{No. of attri not involved in matching}}$

$$= \frac{f_{11}}{f_{01} + f_{10} + f_{11}}, \quad (\text{dissim}) = 1 - J$$

$$\begin{array}{ccccc} X & R_0 & R_1 & J = & 1 \\ H_0 & 2 & 1 & & 1+2+1 \\ H_1 & 2 & 1 & & \end{array} = 0.25$$

\* Proximity measure with categorical attri

$$S(x, y) = \frac{\text{No. of matches}}{\text{Total no. of attri}} = \frac{m}{a}$$

$$d(x, y) = \frac{a-m}{a}$$

$m = \text{no. of matches}$

$a =$

w. ordinal attri

$$\hat{a}_i = \frac{i-1}{n-1} \text{ rank}$$

$$S(x, y) = \sqrt{(\hat{a}_i - \hat{a}_j)^2}$$

Obj	Size	Qualitative	size = {S, M, L}
A	S(0.0)	A(0.66)	1 2 3
B	L(1.0)	Ex(1.0)	S < M < L
C	L(1.0)	C(0.0)	Ex > A > B > C
D	M(0.5)	B(0.33)	4 3 2 1
		↓	Ex > A > B > C

$$a_S = \frac{1-1}{3-1} = 0$$

$$a_M = \frac{2-1}{3-1} = \frac{1}{2}$$

$$a_L = \frac{3-1}{3-1} = 1$$

\* with interval scale

i) Manhattan distance ( $L_1$  Norm :  $r=1$ )

$$d = \sum_{i=1}^n |x_i - y_i|$$

$$\text{eg } x = [7, 3, 5], y = [3, 2, 6]$$

$$|7-3| + |3-2| + |5-6| = 6$$

Hamming distance with when attri values  $\in [0, 1]$

2) Euclidean Distance ( $L_2$  Norm :  $r = 2$ )

3) Chebychev Distance ( $L_\infty$  Norm:  $r \in \mathbb{R}$ )

$$d(x, y) = \max_{\forall i} \{ |x_i - y_i| \}$$

$$\max \{ |7-3|, |3-2|, |5-6| \} = 4$$

\*

for ratio-scale

$$X = A e^B$$

Apply the Naive Bayes classifier algorithm for buys computer classification and classify the type

$X = (\text{age} = \text{"young"}, \text{income} = \text{"medium"}, \text{student} = \text{"yes"}, \text{credit-rating} = \text{"fair"})$

Id	Age	Income	Student	Credit rating	Buys Computer
1	young	high	no	fair	no
2	young	high	no	good	no
3	middle	high	no	fair	yes
4	old	medium	no	fair	yes
5	old	low	yes	fair	yes
6	old	low	yes	good	no
7	middle	low	yes	good	yes
8	young	medium	no	fair	no
9	young	low	yes	fair	yes
10	old	medium	yes	fair	yes
11	young	medium	yes	good	yes
12	middle	medium	no	good	yes
13	middle	high	yes	fair	yes
14	old	medium	no	good	no

→ Total samples = 14

$$\text{No. of Yes} = 9$$

$$P(\text{Yes}) = 9/14$$

$$\text{No. of No} = 5$$

$$P(\text{No}) = 5/14$$

### Age

$$P(\text{young} | \text{Yes}) = \frac{2}{9}$$

$$P(\text{young} | \text{No}) = \frac{3}{5}$$

$$P(\text{middle} | \text{Yes}) = \frac{4}{9}$$

$$P(\text{middle} | \text{No}) = \frac{0}{5}$$

$$P(\text{old} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{old} | \text{No}) = \frac{2}{5}$$

### Income

$$P(\text{high} | \text{Yes}) = \frac{2}{9}$$

$$P(\text{high} | \text{No}) = \frac{2}{5}$$

$$P(\text{medium} | \text{Yes}) = \frac{4}{9}$$

$$P(\text{medium} | \text{No}) = \frac{2}{5}$$

$$P(\text{low} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{low} | \text{No}) = \frac{1}{5}$$

### Student

$$P(\text{Yes} | \text{Yes}) = \frac{6}{9}$$

$$P(\text{Yes} | \text{No}) = \frac{1}{5}$$

$$P(\text{No} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{No} | \text{No}) = \frac{4}{5}$$

### Credit rating

$$P(\text{fair} | \text{Yes}) = \frac{6}{9}$$

$$P(\text{fair} | \text{No}) = \frac{2}{5}$$

$$P(\text{good} | \text{Yes}) = \frac{3}{9}$$

$$P(\text{good} | \text{No}) = \frac{3}{5}$$

Data sample :  $X = \text{age} = \text{"young"}$

$\text{income} = \text{"medium"}$

$\text{student} = \text{"Yes"}$

$\text{credit-rating} = \text{"fair"}$

$$P(Y_{\text{Yes}} | X) = P(X|Y_{\text{Yes}}) \cdot P(Y_{\text{Yes}})$$

$$= P(\text{Young} | \text{Yes}) \cdot P(\text{medium} | \text{Yes}) \\ \cdot P(\text{study} | \text{Yes}) \cdot P(\text{Fair} | \text{Yes}) \cdot P(\text{Yes})$$

$$= \frac{2}{9} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9} \times \frac{9}{14} = 0.02821$$

$$P(N_{\text{No}} | X) = P(X|N_{\text{No}}) \cdot P(N_{\text{No}})$$

$$= P(\text{Young} | N_{\text{No}}) \cdot P(\text{medium} | N_{\text{No}}) \cdot P(\text{study} | N_{\text{No}}) \\ \cdot P(\text{Fair} | N_{\text{No}}) \cdot P(N_{\text{No}})$$

$$= \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{5}{14} = 0.006857$$

$$\text{As } P(Y_{\text{Yes}} | X) > P(N_{\text{No}} | X)$$

Thus,  $X$  belongs to the class buys computer  
∴ The student buys a computer

$$E(S) = - \underbrace{p \log_2 p}_{\text{Yes}} - \underbrace{p \log_2 p}_{\text{No}}$$

PAGE NO.	
DATE	/ /

$E(A_i)$  = same (for all possible values)

$$IG(A_i) = E(S) - E(A_i)$$

Q. ID3

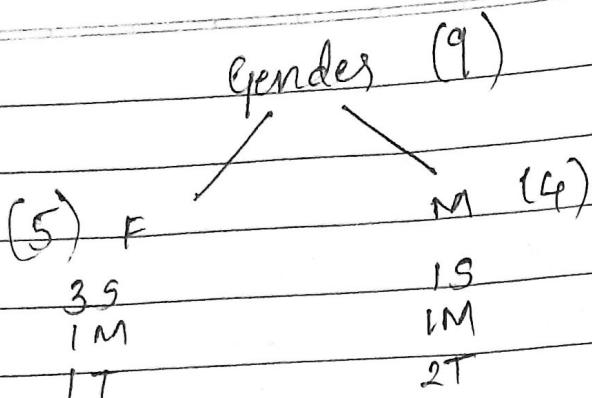
Person	Gender	Height	Class
1	Female	1.6 m	Short
2	Male	2 m	Tall
3	Female	1.9 m	Medium
4	Female	2.1 m	Tall
5	Female	1.7 m	Short
6	Male	1.85 m	Medium
7	Female	1.6 m	Short
8	Male	1.7 m	Short
9	Male	2.2 m	Tall

$$\rightarrow \begin{aligned} \text{Resp } & \leq 1.7 \quad H_1 \\ & > 1.7 = 1.85 \quad H_2 \\ & > 1.85 = 2 \quad H_3 \\ & > 2 \quad H_4 \end{aligned}$$

		Ht	
1	F	$H_1$	S
2	M	$H_3$	T
3	F	$H_3$	M
4	F	$H_4$	T
5	F	$H_{21}$	S
6	M	$H_2$	M
7	F	$H_1$	S
8	M	$H_1$	S
9	M	$H_4$	T

$$E(S) = - \frac{4}{9} \log_2 (4/9) - \frac{2}{9} \log_2 (2/9) - 3/9 \log_2 (3/9)$$

$$= 1.5305$$

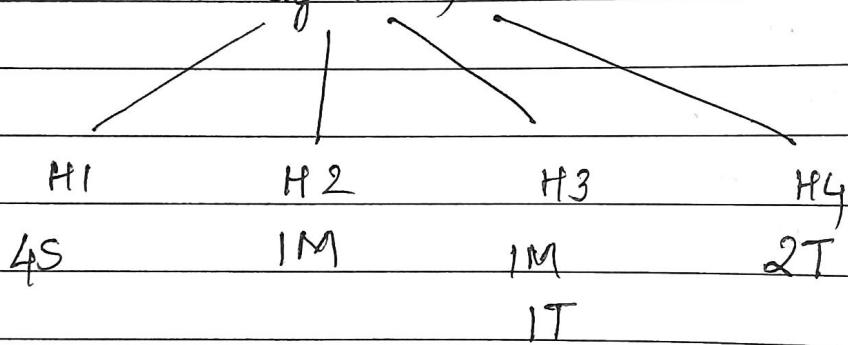


$$E(g) = \frac{5}{9} \left[ -\frac{3}{5} \log_2 \frac{3}{5} - \frac{1}{5} \log_2 \frac{1}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right] + \frac{4}{9} \left[ -\frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right]$$

$$E(g) = 1.4283$$

$$Ig(g) = 0.1022$$

Height (9)



$$E(H) = \frac{4}{9} \left[ -\frac{4}{4} \log_2 \frac{4}{4} \right] + \frac{1}{9} \left[ -\frac{1}{1} \log_2 \frac{1}{1} \right] + \frac{2}{9} \left[ -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] + \frac{2}{9} \left[ -\frac{1}{1} \log_2 \frac{1}{1} \right]$$

$$E(H) = 0.2222$$

$$Ig(H) = 1.3083$$

$$Ig(H) > Ig(g)$$

$\therefore H$  is the root element

\* Kmeans

$$D = \{1, 2, 6, 7, 8, 10, 15, 17, 20\}$$

$$K_1 = \{1, 6, 8, 15, 20\}$$

$$K_2 = \{2, 7, 10, 17\}$$



$$K_1 \text{ mean} = 10$$

$$K_2 \text{ mean} = 9$$

$$K_1(10) = \{10, 15, 17, 20\}$$

$$\text{mean } 15.5$$

closer to 10

$$K_2(9) = \{1, 2, 6, 7, 8\}$$

$$\text{mean } 4.8$$

$$K_1(15.5) = \{15, 17, 20\}$$

$$17.333$$

$$K_2(4.8) = \{1, 2, 6, 7, 8, 10\}$$

$$5.667$$

$$K_1(17.3) = \{15, 17, 20\}$$

$$17.333$$

$$K_2(5.6) = \{1, 2, 6, 7, 8, 10\}$$

$$5.667$$

same means  $\Rightarrow$  stop

\* Agglomerative  
Q. Single link

G/P dist matrix

	BA	FI	MI	NA	RM	TO
BA	0	662	887	225	412	996
FI	662	0	295	468	268	400
MI	887	295	0	754	564	138
NA	225	468	754	0	219	869
RM	412	268	564	219	0	669
TO	996	400	138	869	669	0

Single Link

Min dist

MI  $\xrightarrow{\text{138}}$  TO

	BA	FI	MI/TO	NA	RM
BA	0	662	877	255	412
FI	662	0	295	468	268
MI/TO	887	295	0	754	564
NA/RM	255	468	754	0	219
RM	412	268	564	219	0

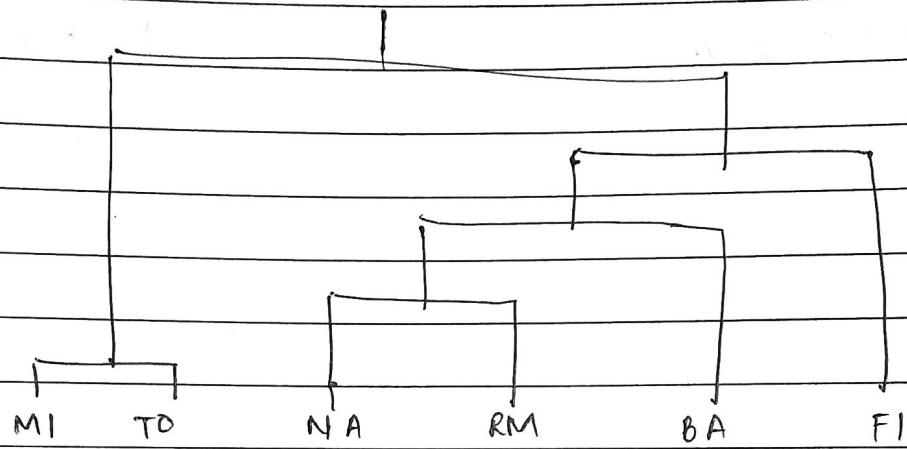
	BA	FI	MI/TO	NA/RM
BA	0	662	877	255
FI	662	0	295	268
MI/TO	877	295	0	564
NA/RM	255	268	564	0

MI/T  $\circ$  NA/RM

BA/NA/RM

FI

MI/



Q. Complete link

	1	2	3	4	5
1	1	0.90	0.10	0.65	0.20
2	0.90	1.00	0.70	0.60	0.50
3	0.10	0.40	1.00	0.40	0.30
4	0.65	0.60	0.40	1.00	0.80
5	0.20	0.90	0.30	0.80	1.00

1,3

	1,3	2	4	5
1,3	1.0	0.9	0.65	0.83
2	0.9	1.0	0.46	0.5
4	0.65	0.6	0.4	0.8
5	0.3	0.5	0.8	1.00

note

1, 3, 5

1, 3, 5	2	4	
1, 3, 5	1.0	0.9	0.8
2	0.9	1.0	0.6
4	0.8	0.6	1.0

2, 4

1, 3, 5	2, 4	
1, 3, 5	1.0	0.9
2, 4	0.9	1.0

1, 3, 5, 2, 4

1, 3, 5, 2, 4

## \* Market Basket Analysis

if you buy  $x$ , you are more likely to buy  $y$ .

eg Bread  $\rightarrow$  Milk

if you buy bread, you are more/less likely  
to buy milk

$$x \rightarrow y$$

$$x \wedge y \neq \emptyset \text{ ie } x \neq y$$

for single item

### \* Support

min support threshold = 50%.

$$1: 1, 3, 5$$

$$\text{support of } \{8, 12\} = 2 \text{ (or } 50\%)$$

$$2: 1, 8, 14, 17, 12$$

$$\{1, 5\} = 1 \text{ ( } 25\%)$$

$$3: 4, 8, 6, 12, 9, 10, 4$$

$$\{1\} = 3 \text{ ( } 75\%)$$

$$4: 2, 1, 8$$

$$50 = 50$$

$75 > 50 \Rightarrow \{8, 12\}, \{1\}$  are frequent.  
(high support)

### \* Confidence

min confi thi = 50%.

Dr.

$$\text{confi of } \{8 \rightarrow 12\} = \frac{2}{3} \quad 66\%$$

assoc rule

$$\text{confi of } \{1 \rightarrow 5\} = \frac{1}{3} \quad 33\%$$

$\{8, 12\}$  is strong (high confidence)

Sup → how many times bought

PAGE NO.	/ /
DATE	/ /

confi → associativity

Q. Min Confid there = 50%.

$$\begin{array}{l} \text{sup} = 30\% \\ 1: 3, 5, 8 \\ 2: 268 \\ 3: \\ 4: \end{array}$$

support of  $\{5, 8\} = \frac{4}{10} (40\%)$

confi of  $\{5, 8\} = \frac{4}{5} (80\%)$

(meaningful) strong and frequent

Q. min = 50, 50%.

$$\text{sup of } \{9, 3\} = \frac{1}{10} (10\%)$$

$$\text{confi of } \{9 \Rightarrow 3\} = \frac{1}{1} (100\%)$$

↑ confi ↓ sup → not meaningful

Q \* Apriori Algorithm

Q. Min support = 30%, Min confidence = 75%.

Database D	TID	Items	Itemset	Sup
→	1	A, B, D	C <sub>1</sub>	A 3
	2	B, C, D		B 5
	3	A, B	Scan D	C 2
	4	B, D		D 3
	5	A, B, C		

L <sub>1</sub>	Itemset	Sup	(K+1)	L <sub>2</sub>	Itemset	Sup	Min sup 30%
{A, B}	{A, B}	3			{A}		
{B, C}	{A, C}	0	X		{B}		
{B, D}	{A, D}	1	X		{C}		
	{B, C}	2			{D}		
	{B, D}	3					
	{C, D}	1	X				

		Sup
Scan D	{A, B, C}	1 x
→	{A, B, D}	1 x
C <sub>3</sub>	{B, C, D}	1 x

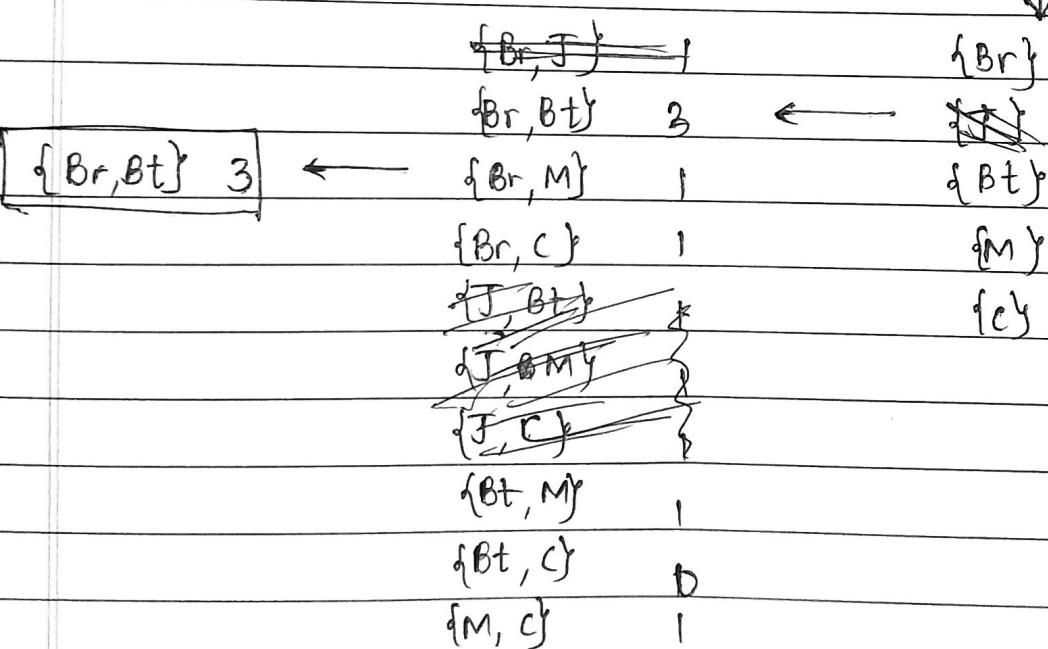
∴ All combinations are < 30% support, so no rules for association can be generated

$$Q. \text{ min support} = 30\%$$

$$\text{min confi} = 80\%$$

Transaction	Items			Sup
T <sub>1</sub>	Bread, Jelly, Butter	C <sub>1</sub>	Bread	1
T <sub>2</sub>	Bread, Butter	→	Jelly	1
T <sub>3</sub>	Bread, Milk, Butter	Scan D	Butter	3
T <sub>4</sub>	Coke, Bread		Bread	2
T <sub>5</sub>	Coke, Milk		Milk	2
			Coke	2

↓ Min sup 30%



$$\text{Bread} \rightarrow \text{Butter} = 3/4 = 75\% \times 280$$

$$\text{Butter} \rightarrow \text{Bread} = 3/3 = 100\%.$$

If item comes twice, count only 1

PAGE No.	
DATE	/ /

Q. Slide 29 Chpt 5

O, K, E - 3

Associa:

$O \rightarrow K, E$       3/3      ✓

$K \rightarrow O, E$       3/5

$E \rightarrow K, O$       3/4

$O, K \rightarrow E$       3/4      ✓

$O, E \rightarrow K$       3/3      ✓

$K, E \rightarrow O$       3/4