

# DATA MINING

## NUMERICALS

### UNIT - 1

- ① Find frequent item sets using Apriori.  
 Min support = 40%. Min. confidence = 60%.

Tid	Items	(2020-SEE)
T <sub>1</sub>	1, 3, 4	
T <sub>2</sub>	2, 3, 5	
T <sub>3</sub>	1, 2, 3, 5	
T <sub>4</sub>	2, 5	
T <sub>5</sub>	1, 3, 5	

$$\text{Support count} = \frac{40}{100} \times \text{no. of transactions}$$

$$= \frac{40}{100} \times 5 = \underline{\underline{2}}$$

C <sub>1</sub> = items	support	L <sub>1</sub> = items	support
{1 3 4}	3	{1 3}	3
{2 3}	3	{2}	3
{3 3}	4	{3}	4
{4 3}	1	{5}	4
{5}	4		
⋮			

$C_2 =$	<u>items</u>	<u>Support</u>	$L_2 =$	<u>items</u>	<u>Support</u>
	$\{1, 2\}$	1		$\{1, 3\}$	3
	$\{1, 3\}$	3		$\{1, 5\}$	2
	$\{1, 5\}$	2		$\{2, 3\}$	2
	$\{2, 3\}$	2		$\{2, 5\}$	3
	$\{2, 5\}$	3		$\{3, 5\}$	3
	$\{3, 5\}$	3			

$C_3 =$	<u>items</u>	<u>Support</u>	$L_3 =$	<u>items</u>	<u>Support</u>
	$\{1, 3, 5\}$	2		$\{1, 3, 5\}$	2
	$\{1, 2, 3\}$	1		$\{2, 3, 5\}$	2
	$\{1, 2, 5\}$	1			
	$\{2, 3, 5\}$	2			

$C_4 =$	<u>items</u>	<u>support</u>	$L_4 = \emptyset$
	$\{1, 2, 3, 5\}$	1	

$\therefore$  Frequent itemsets  $= L_3 = \{1, 3, 5\}, \{2, 3, 5\}$ .

## ② (2019-SEE)

$$L_3 = [\{1, 2, 3\}, \{1, 2, 4\}, \{1, 2, 5\}, \{1, 3, 4\}, \\ \{1, 3, 5\}, \{2, 3, 4\}, \{2, 3, 5\}, \{3, 4, 5\}]$$

- i) Generate  $C_4$  using  $F_{k-1} \times F_k$  strategy.  
ii) Generate  $C_4$  using  $F_{k-1} \times F_{k-1}$  strategy.

i)  $\{1, 2, 3, 4\}, \{1, 2, 3, 5\}$   
 $\{1, 2, 4, 5\}, \{1, 3, 4, 5\}$   
 $\{2, 3, 4, 5\}$

ii)  $\{1, 2, 3, 4\}, \{1, 2, 3, 5\}$   
 $\{1, 2, 4, 5\}, \{1, 3, 4, 5\}$   
 $\{2, 3, 4, 5\}$

Both result the same.

→ Another question had brute force, in that generate all possible combinations.  $F_1 \times F_2 \times F_3$

Unit - 2

① (2020 - SEE)

Trace Apriori &amp; generate rules

Min. Support = 30%.

Min. confidence = 80%.

<u>Sup</u>	<u>Tid</u>	<u>Itemset</u>
	1	1, 2, 5
	2	2, 4
	3	2, 3
	4	1, 2, 4
	5	1, 3
	6	2, 3
	7	1, 3
	8	1, 2, 3, 5
	9	1, 2, 3

$$\text{Min. support count} = \frac{30}{100} \times 9 = 2.7 \Rightarrow 3.$$

## Notes

$C_1 = \text{items}$	Support	$L_1 = \text{items}$	Support
$\{1\}$	6	$\{1\}$	6
$\{2\}$	7	$\{2\}$	7
$\{3\}$	6	$\{3\}$	6
$\{4\}$	2 $\times$		
$\{5\}$	2 $\times$		

$C_2 = \text{items}$	Support	$L_2 = \text{items}$	Support
$\{1, 2\}$	4	$\{1, 2\}$	4
$\{1, 3\}$	4	$\{1, 3\}$	4
$\{2, 3\}$	4	$\{2, 3\}$	4

$C_3 = \text{items}$	Support	$L_{\neq 3} = \emptyset$
$\{1, 2, 3\}$	2	

$\therefore$  Frequent items =  $\{1, 2\}, \{1, 3\}, \{2, 3\}$

~~Confidence ( $A \rightarrow B$ )~~

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

Date:

$$C(1 \rightarrow 2) = \frac{S(1, 2)}{S(1)} = \frac{4}{6} = 66\% < 80\% \times$$

$$C(2 \rightarrow 1) = \frac{4}{7} < 80\% \times$$

$$C(1 \rightarrow 3) = \frac{4}{6} < 80\% \times$$

$$C(3 \rightarrow 1) = \frac{4}{6} < 80\% \times$$

$$C(2 \rightarrow 3) = \frac{4}{7} < 80\% \times$$

$$C(3 \rightarrow 2) = \frac{4}{6} < 80\% \times$$

∴ There are no strong rules.



② (SEE-2020) Draw FP tree & find frequent itemsets.

Min support = 30%.

Tid	Item
1	E, A, D, B
2	D, A, C, E, B
3	C, A, B, E
4	B, A, D
5	D
6	D, B
7	A, D, E
8	B, C

$$\text{Min. support count} = \frac{30}{100} \times 8 = \frac{24}{2.4} = 3.$$

Items'	Support
A	5
B	6
C	3
D	6
E	4

→ Eliminate any that are below min support.

• Setting priority

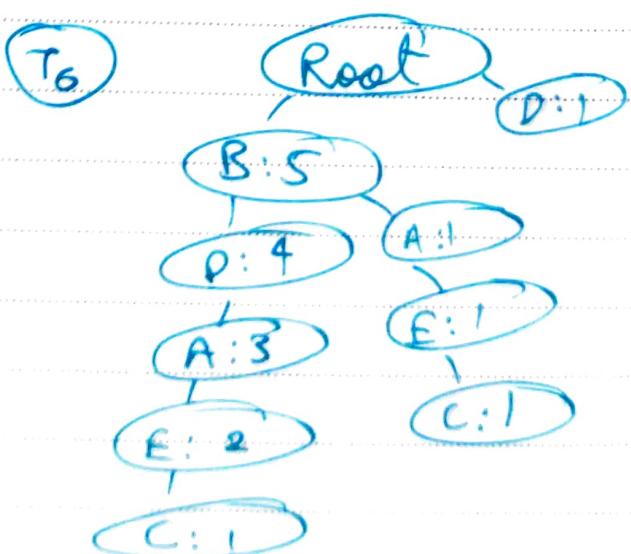
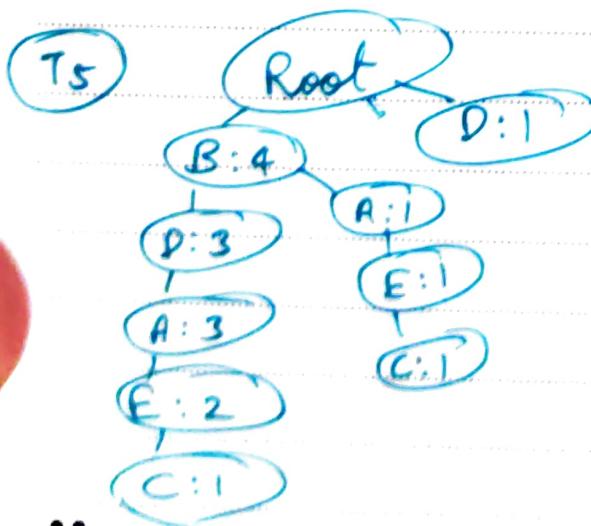
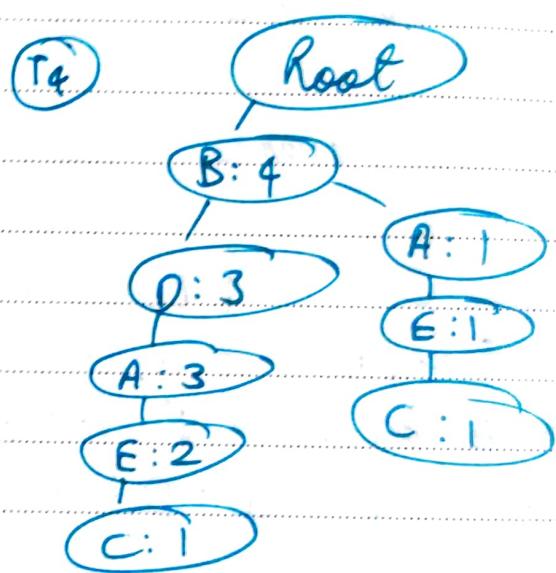
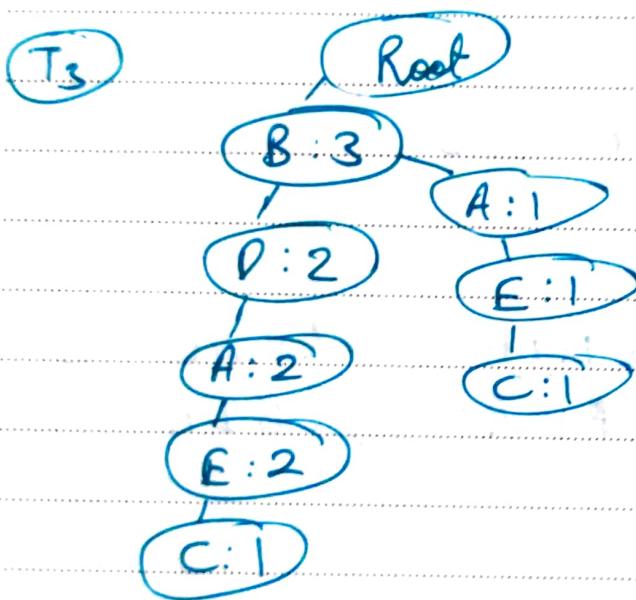
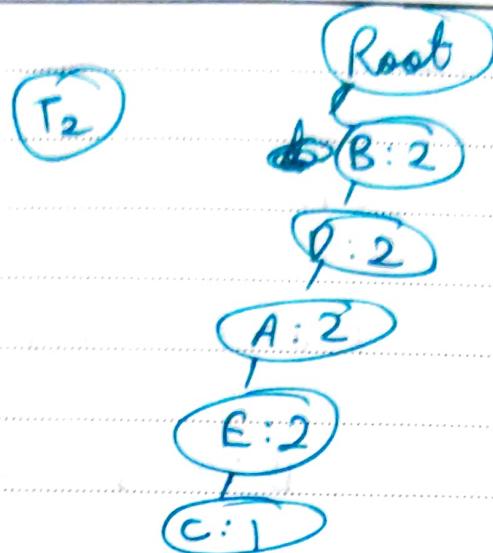
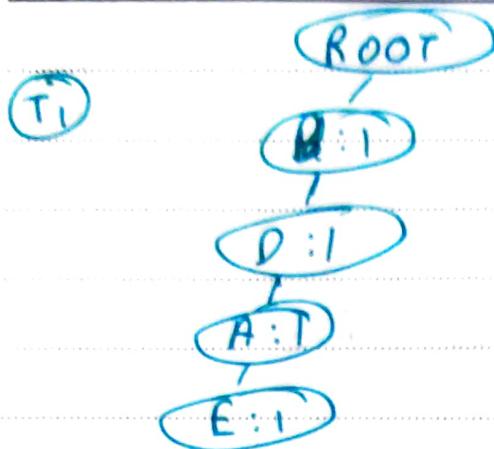
<u>Items</u>	<u>Support</u>
B	6
D	6
A	5
E	4
C	3

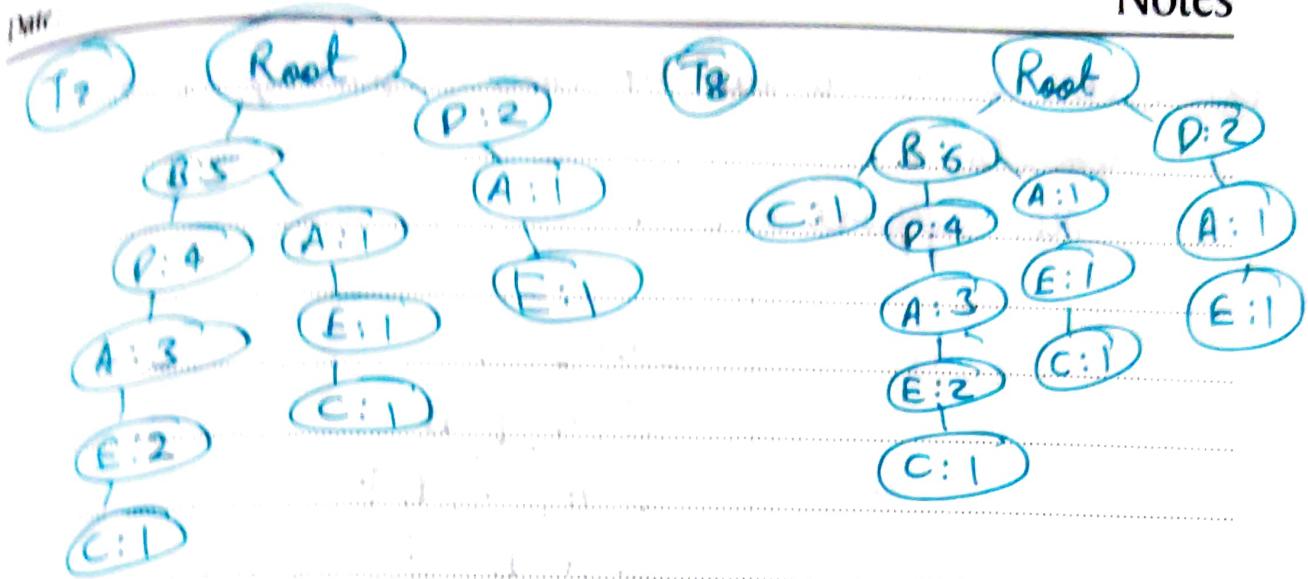
Redefine transactions, do not show items less than min support.

<u>Tid</u>	<u>Items</u>
1	B, D, A, E
2	B, D, A, E, C
3	B, A, E, C
4	B, D, A
5	D
6	B, D
7	D, A, E
8	B, C

# Notes

Date: \_\_\_\_\_





Item	Conditional Pattern base	Conditional FP tree	FP generated
I: C	$\{\{B, D, A, E\} : 1\}$ , $\{B : 1\}$ , $\{BAE : 1\}$	$\{B : 3\}$ , ignored others because min. support.	$\{B, C\}$
E	$\{BDA : 2\}$ , $\{BA : 1\}$ , $\{D, A : 1\}$	<del><math>\{A : 3\}</math>, <math>\{B : 3\}</math></del>	$\{A, E\}$
A	$\{BD : 3\}$ , $\{$		

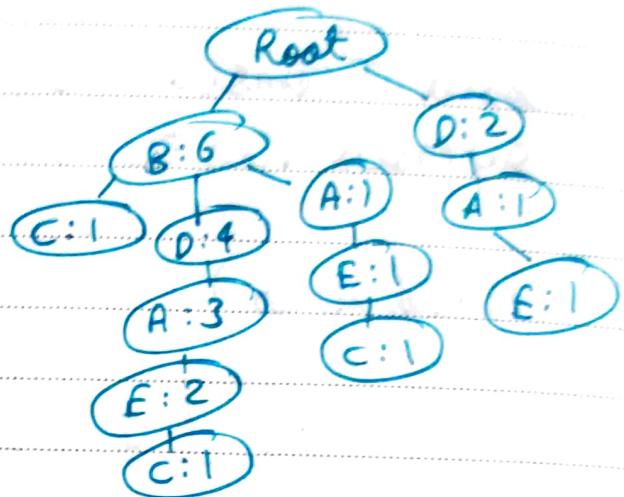
## @ Conditional Pattern Base:

C : 3

E : 4

A : 5

D : 6



→ For C

{ B, D, A, E : 1 }

{ B : 1 }

{ B, A, E : 2 }

→ For E : { B, D, A : 2 }, { B, A : 1 }, { D, A : 1 }

→ For A : { B, D : 3 }, { B : 1 }, { D : 1 }

→ For D : { B : 4 }

① Conditional FP tree

MS = 3

→ For C : {B: 3}

→ For E : {B: 3, A: 3} ~~DX23~~ ~~E: 3~~

→ For A : {B: 4, D: 3}

→ For D : {B: 4}

② Frequent Itemsets Patterns

→ For C : {B, C: 3}

→ For E : {B, E: 3}, {A, E: 3}, {BAE: 3}

→ For A : {B, A: 4}, {D, A: 3}, {BDA: 3}

→ For D : {B, D: 4}

## ③ (SEE 2019)

No. of Transactions =  $N = 100$

" " items = 20

$$S(a) = 25\%. \quad S(b) = 90\%.$$

$$S(a, b) = 20\%. \quad MS = 10\%. \quad MC = 60\%.$$

a)  $C(a \rightarrow b) = ?$

b)  $I(a, b) = ?$

$$\text{a)} \quad C(a \rightarrow b) = \frac{S(a, b)}{S(a)} = \frac{20}{25} = 80\% > 60\%$$

$\therefore$  It is interesting .

$$\text{b)} \quad I(a, b) = \frac{C(A \rightarrow B)}{S(B)} = \frac{80\%}{90\%} = 88\%.$$

→ Interest factor compares the frequency of a pattern to the baseline frequency .

- \* **Maximal Itemset** : If none of its immediate supersets is frequent.
- \* **Closed Itemset** : If ~~not~~ none of its immediate supersets has the same support count.

④ (Makeup 2020)

Tid	Itemset
1	A B D E
2	B C E
3	A B D E
4	A B C E
5	A B C D E
6	B C D

$$MS = 3$$

compute Closed itemset & maximal itemset

item	support
A	4
B	6
C	4
D	4
E	5

## Notes

Date:

$C_2 = \text{item}$	<u>Support</u>	$C_3 = \text{item}$	<u>Support</u>
AB	4	ABD	3
<del>AC</del>	2	<del>ABE</del>	4
AD	3	<del>ABC</del>	2
AE	4	ADE	3
BC	4	<del>ACE</del>	2
BD	4	<del>BCD</del>	2
BE	5	BCE	3
<del>CD</del>	2	<del>BDE</del>	3
CF	3	<del>CDE</del>	1
DE	3		

$C_4 = \text{item}$	<u>Support</u>
ABDE	3 ✓
<del>ABCDE</del>	2
<del>BODE</del>	1

Maximal: ABDE, BCE,

Closed: ABDE, BCE, ABE, BE, BD, BC, B

	Basket Ball	Not BB	$\Sigma$
Cereal	2000	1750	3750
Not cereal	1000	250	1250
$\Sigma$	3000	2000	5000

i) Lift =  $\frac{S(A, B)}{S(A) \times S(B)} = I(A, B) = \frac{S(A, B)}{S(A) S(B)} = \frac{N f_{11}}{f_{1+} f_{+1}}$

$$\frac{f_{11}}{N} = \frac{f_{1+}}{N} \times \frac{f_{+1}}{N} = \frac{f_{1+} f_{+1}}{N}$$

$$\Rightarrow \text{Lift} = \frac{5000 \times 2000}{3750 \times 3000} = 0.88 < 1.$$

A & B are positively correlated.



Unit 3Decision Tree

Age	Competition	Type	Profit
Old	Yes	Software	Down
Old	No	Software	Down
Old	No	Hardware	Down
Mid	Yes	SW	Down
Mid	Yes	HW	Down
Mid	No	HW	Up
Mid	No	SW	Up
New	Yes	SW	Up
New	No	HW	Up
New	No	SW	Up

$$P = 5 \quad N = 5$$

① CALCULATE ENTROPY OF CLASS ATTRIBUTE

$$\text{Entropy (class)} = -\frac{P}{P+N} \log_2 \left( \frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left( \frac{N}{P+N} \right)$$

$$= -\frac{5}{10} \log_2 \left( \frac{5}{10} \right) - \frac{5}{10} \log_2 \left( \frac{5}{10} \right)$$

$$-0.5 \times (-1) - 0.5 (-1)$$

$$\therefore = 0.5 + 0.5 = 1$$

② FOR EACH ATTRIBUTE, calculate

- ① information gain
- ② Entropy (attribute)
- ③ Gain

① Information Gain

$$I(P_i, N_i) = -\frac{P}{P+N} \log_2 \left( \frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left( \frac{N}{P+N} \right)$$

② Entropy of attribute

$$= \sum \frac{P_i + N_i}{P+N} \times (I(P_i, N_i))$$

③ Gain = class entropy - Entropy of attribute.

## Notes

Date:

For Age	P <sub>i</sub>	N <sub>i</sub>	I(P <sub>i</sub> , N <sub>i</sub> )
old	0	3	0
mid	2	2	1
New	3	0	0

$$\text{Entropy of age} = \frac{\sum P_i + N_i (I(P_i, N_i))}{P+N}$$

$$= \frac{0+3}{5+5} \times 0 + \frac{2+2}{5+5} \times 1 + \cancel{3+0} \times 0$$

$$= \frac{2}{5} = \underline{0.4}$$

$$\text{Gain} = 1 - 0.4 = \underline{0.6}$$

For Competition	P <sub>i</sub>	N <sub>i</sub>	I(P <sub>i</sub> , N <sub>i</sub> )
Yes	1	3	0.81
No	4	2	0.92

$$\text{a) Info gain} = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right)$$

$$= 0.5 + 0.31 = 0.81128$$

$$\text{b) } \therefore \frac{-4}{6} \log\left(\frac{4}{6}\right) - \frac{2}{6} \log\left(\frac{2}{6}\right) = \cancel{0.3899} + 0.918$$

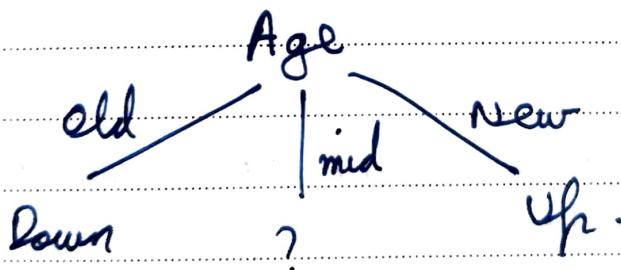
Entropy of competition =  $\frac{(1+3)}{5+5} \times 0.81 + \frac{6}{10} (0.92)$   
 $= 0.8754.$

Grain =  $1 - 0.8754 = 0.124515.$

<u>For Type</u>	$P_i$	$N_i$	$I(P_i, N_i)$
H.W	2	2	1
SW	3	3	1

Entropy of SW type =  $\frac{2+2}{10} \times 1 + \frac{6}{10} \times 1$   
 $= \frac{4}{10} + \frac{6}{10} = 1.$

Grain = 0



Age	competition	Type	Profit
mid	Y	SW	D
mid	X	HW	D
mid	N	SW	U
mid	N	HW	U

For competition:		P <sub>i</sub>	N <sub>i</sub>	I(P <sub>i</sub> , N <sub>i</sub> )
Y		2/4	2	0
N		2	0	0

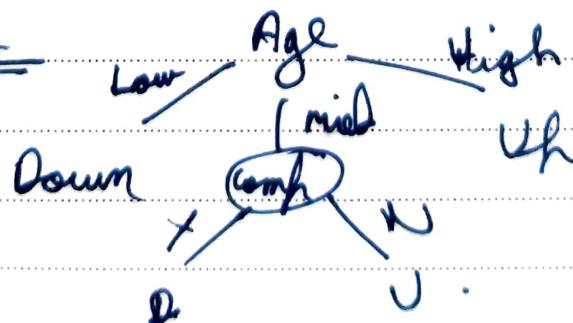
$$\text{Entropy} = \frac{2}{4} \times 0 + \frac{2}{4} \times 0 = 0 = 0.$$

Grain = 1

For type		P <sub>i</sub>	N <sub>i</sub>	I(P <sub>i</sub> , N <sub>i</sub> )
SW		1	1	1
HW		1	1	1

$$\frac{2}{4} \times 1 + \frac{2}{4} \times 1 = 1$$

Grain = 0



∴

## Decision Tree using CART

$$\text{Gini index (attribute = value)} = 1 - \sum_{i=1}^N (p_i)^2$$

$$\text{Gini index (attribute)} = \sum_{v=\text{values}} p_v \times GJ(v)$$