



Name of the Subject: DATA ANALYSIS & INFO. EXT. Subject Code: IT-704

Seat No: IT076 Student ID: 18ITUBN116 Branch/Sem: IT-VII

[Q2]

[d]

Expected Info

$$I(S_1, S_2) = - \sum_{i=1}^2 \frac{S_i}{S} \log \left(\frac{S_i}{S} \right)$$

$$= - \frac{S_1}{S} \log \left(\frac{S_1}{S} \right) - \frac{S_2}{S} \log \left(\frac{S_2}{S} \right)$$

$$= - \frac{120}{250} \log \left(\frac{120}{250} \right) - \frac{130}{250} \log \left(\frac{130}{250} \right)$$

$$= 0.99088$$

- To check whether "major" is a relevant attribute or not.

Here, V=3 (Science, Engineering, Business)

Expected Info. for Science

$$S_{11} = 84 \text{ (For Grad.)}$$

$$S_{21} = 42 \text{ (For Undergrad.)}$$

$$I(S_{11}, S_{21}) = - \frac{84}{126} \log \left(\frac{84}{126} \right) - \frac{42}{126} \log \left(\frac{42}{126} \right)$$

$$= 0.9183$$



Name of the Subject: DATA ANALYSIS & INFO. EXT. Subject Code: IT-704
Seat No: IT076 Student ID: JSITVBN116 Branch/Sem: IT-VII

For Engg.

$$S_{12} = 36 \text{ (grad.)}$$

$$S_{22} = 46 \text{ (Undergrad.)}$$

$$I(S_{12}, S_{22}) = -\frac{36}{82} \log\left(\frac{36}{82}\right) - \frac{46}{82} \log\left(\frac{46}{82}\right)$$
$$= 0.9892.$$

For Business

$$S_{13} = 0 \text{ (Grad.)}$$

$$S_{23} = 42 \text{ (Undergrad.)}$$

$$I(S_{13}, S_{23}) = -\frac{42}{42} \log(1) = 0$$

$$E(\text{major}) = \frac{126}{250} I(S_{11}, S_{21}) + \frac{82}{250} I(S_{12}, S_{22})$$
$$+ \frac{42}{250} I(S_{13}, S_{23})$$

$$\text{Gain (major)} = 0.9988 - 0.7873$$

$$= \boxed{0.2115}$$



Name of the Subject: DATE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116 Branch/Sem: IT-VII

102

e

Class

 $C_1 = \text{buys_Computer} = \text{'Yes'}$ $C_2 = \text{buys_Computer} = \text{'No'}$

Data to be classified

 $X = (\text{Age} \leq 30, \text{Income} = \text{medium}, \text{Student} = \text{Yes}, \text{Credit_rating} = \text{Fair})$

$$P(C_1) = P(\text{buys_Comp} = \text{'Yes'}) = \frac{9}{14} = 0.643.$$

$$P(\text{buys_Comp} = \text{'No'}) = \frac{5}{14} = 0.357$$

→ Compute $P(X|C_i)$ for each class

$$P(\text{Age} = \text{'}\leq 30\text{'} | \text{buys_Comp} = \text{'Yes'}) = \frac{2}{9} = 0.22.$$

$$P(\text{Age} = \text{'}\leq 30\text{'} | \text{buys_Comp} = \text{'No'}) = \frac{3}{5} = 0.6.$$

$$P(\text{Income} = \text{'medium'} | \text{buys_Comp} = \text{'Yes'}) = \frac{6}{9} = 0.667$$

$$P(\text{Income} = \text{'medium'} | \text{buys_Comp} = \text{'No'}) = \frac{1}{5} = 0.2.$$

$$P(\text{Credit_rating} = \text{'fair'} | \text{buys_Comp} = \text{'Yes'}) = \frac{6}{9} = 0.667$$

$$P(\text{Credit_rating} = \text{'fair'} | \text{buys_Comp} = \text{'No'}) = \frac{2}{5} = 0.4.$$



Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116 Branch/Sem: IT-VII

$X = (\text{age} \leq 30, \text{income} - \text{medium}, \text{Student} = \text{yes}, \text{credit_ratio} = \text{fair})$.

$$P(X|C_i) = P(X \mid \text{buys_Comp} = \text{'Yes'}) = 0.22 \times 0.44 \times 0.667 \\ * 0.667 \\ = 0.044$$

$$P(X|C_i) = P(X \mid \text{buys_Comp} = \text{'No'}) = 0.6 \times 0.4 \times 0.2 \times 0.4 \\ = 0.019$$

$$P(X|C_i) \times P(C_i) = P(X \mid \text{buys_Comp} = \text{'Yes'}) \times P(\text{buys_Comp} = \text{'Yes'}) \\ = 0.028$$

$$P(X \mid \text{buys_Comp} = \text{'No'}) \rightarrow P(\text{buys_Comp} = \text{'No'}) = 0.007.$$

— Therefore, x belongs to class ("Buys.Computer = 'Yes'")

[Q2] [C].

T _{ID}	Item
T ₁₀₀	K, A, D, B.
T ₂₀₀	D, A, C, E, B
T ₃₀₀	C, A, B, E
T ₄₀₀	B, A, D.

↓
Scan D. B

P
↓

PTO.



Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18JTURN116 Branch/Sem: IT-VII

C₁

Item Set	Sup.
{A}	4
{B}	4
{C}	2
{D}	3
{E}	2
{K}	1.

As, Min-Sup = 2



L ₁	Item Set	Sup.
	{A}	4
	{B}	4
	{C}	2
	{D}	3
	{E}	2.



C₂. Item Set

- {A, B}
- {A, C}
- {A, D}
- {A, E}
- {B, C}
- {B, D}
- {B, E}
- {C, D}
- {C, E}
- {D, E}.



→ two over.



Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116 Branch/Sem: IT-VII

Scad.	Itemset	Sup.
C ₂ .	{A, B}	6 0 4
	{A, C}	2
	{A, D}	0 2 3
	{A, E}	2
	{B, C}	2
	{B, D}	3
	{B, E}	2
	{C, D}	1
	{C, E}	2
	{D, E}	1



L ₂	Itemset	Sup.
	{A, B}	6 0 4
	{A, C}	2
	{A, D}	0 2 3
	{A, E}	2
	{B, C}	2
	{B, D}	3
	{B, E}	2
	{C, E}	2



L ₃	Itemset	Sup.
	{A, B, D}	0 2 3
	{A, B, C}	2
	{A, B, E}	2
	{A, C, D}	1
	{A, C, E}	2
	{B, C, D}	1
	{B, C, E}	1
	{B, D, E}	2



Name of the Subject: DAIE

Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116

Branch/Sem: IT-VII

C ₃	Itemset	Sup	%
	{A, B, C}	2	2
	{A, B, D}	3	3
	{A, B, E}	2	2
	{A, C, D}	1	1
	{A C E}	2	2
	{A D E}	1	1
	{B C D}	1	1
	{B C E}	2	2
	{B D E}	1	1
		↓	Scan D

Item Sup L₃.

{A B D}	3
{A C E}	2
{A B E}	2
{B C E}	2.

Available Association Rules are,

$$A \cap B \Rightarrow D = 3/4 = 75\%$$

$$A \cap D \Rightarrow B = 3/3 = 100\%$$

$$B \cap D \Rightarrow A = 3/3 = 100\%$$

$$D \Rightarrow A \cap B = 3/3 = 100\%$$

$$A \cap D \Rightarrow B. = 3/4 = 75\%$$

$$B \cap D \Rightarrow A. = 3/4 = 75\%.$$

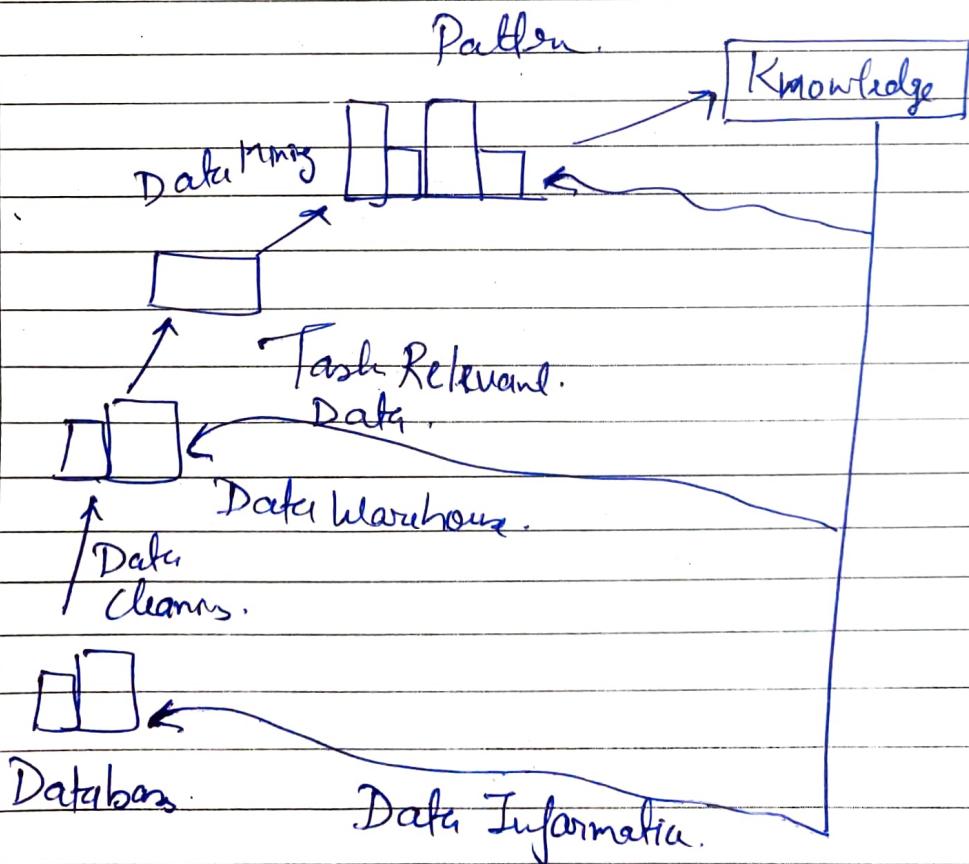
So, $A \cap D \Rightarrow B$
 $B \cap D \Rightarrow A$ } are Strong Association.
 $A \cap B \Rightarrow D$



[Q3] [9].

KDD - Knowledge Discovery in Database.

- Volume of Information is increasing everyday that we can handle business transaction & different data.
- So, we need a System that will be capable of extracting essential of information available & that can generate report & making decision.
- It refers to the nontrivial extraction of implicit provided & potentially useful information from data stored in database.





Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116 Branch/Sem: IT-VII

① Data Cleaning :- Removal of noisy & irrelevant data.

② - Cleaning with Data discrepancy detection & Data Transformation

② Data Integration - Heterogeneous data from multiple source combined in Common Source

- Data Integration using ETL

③ Data Selection - The process where data relevant to analysis selected & retained from data collection.

④ - It uses Clustering.

④ Data Transformation - The process of transforming data into appropriate form required by mining.

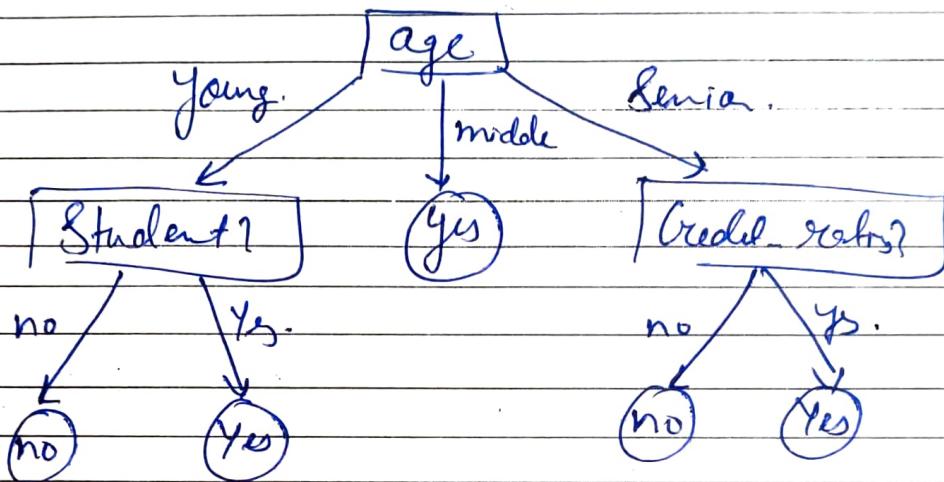
⑤ Data Mining - Techniques that are applied to extract patterns

⑥ Pattern Evaluation - Identifying strictly increasing pattern representing knowledge based

⑦ Knowledge Representation - Technique which utilizes visualization tool to represent data mining result.

Q3 b

A Decision Tree is a structure that include a root node, branches & leaf nodes. Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, & each leaf node holds a class label.



→ Benefits of decision Tree

- It does not require any domain knowledge
- It is easy to comprehend
- The learning & classification steps of a decision tree are simple & fast.

- ~~• Drawbacks~~



Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITURN116 Branch/Sem: IT-VII

— Working.

Create a node N;

if tuples in D are all of the same class, c then
return N as leaf node labeled with c;

If attribute list is empty then

return N as leaf node major class in D;

apply attribute selection method (D, attribute list)

to find the best splitting criteria;
label node N with splitting crit;

If splitting attribute is discrete-valued &
multiway splits allowed

attribute-list = splitting attribute;

for each outcome j of splitting criteria

let D_j be = set of data tuples.

If D_j is empty then
attach a leaf
class in D to node N;
else,

attach the node refined by generate
decision tree

end if.

return N;



Name of the Subject: DAIE Subject Code: IT-704

Seat No: IT076 Student ID: 18ITVBN116 Branch/Sem: IT-VII

103 C

(i) Dimension & Fact Table.

- date_key - calendar_date.
day_no.
day_overall
day_of_week.
month_no_overall
Year
time.

- customer_key - Cust_full_name
Cust_gender.
Cust_title.
Cust_birthdate.
Cust_Shipto_name.
Cust_postalcode
Cust_number
Cust_adol
Cust_payment
Cust_credit_Status.

- product_key - prod_number
prod_size
prod_brand.
prod_descript.
prod_manufact.

- ship_mode - ship_mode_name
Carriar_name.
Ship_Contact_number
Ship_Contact_name



Name of the Subject: DAIT Subject Code: IT-704

Seat No: IT076 Student ID: 18JITURN116 Branch/Sem: IT-VII

date_key INT

Calendar date DATE

day_no INT

day_overall INT

day_week TEXT

month_overall

year TEXT

Product_key INT

prod_no TEXT

product_Sig TEXT

prodlt_brand TEXT

prodle_din TEXT

prod_ma TEXT

Fact Table

date_key

Customer INT.

cust_full_name TEXT

" - gender TEXT

" - title TEXT

" - ShipTo_name TEXT

cust_postalcode TEXT

cust_no INT

cust_addr TEXT

cust_paym TEXT

cust_credit_lim TEXT

Custom_key.

product_key.

Ship_mode.

Ship_mode INT

Total amount

Total quantity

Total Discount.

Ship_mode_no

CARRIER_no

Ship_contact_no

Ship_contact_name

Ship.

Product

date