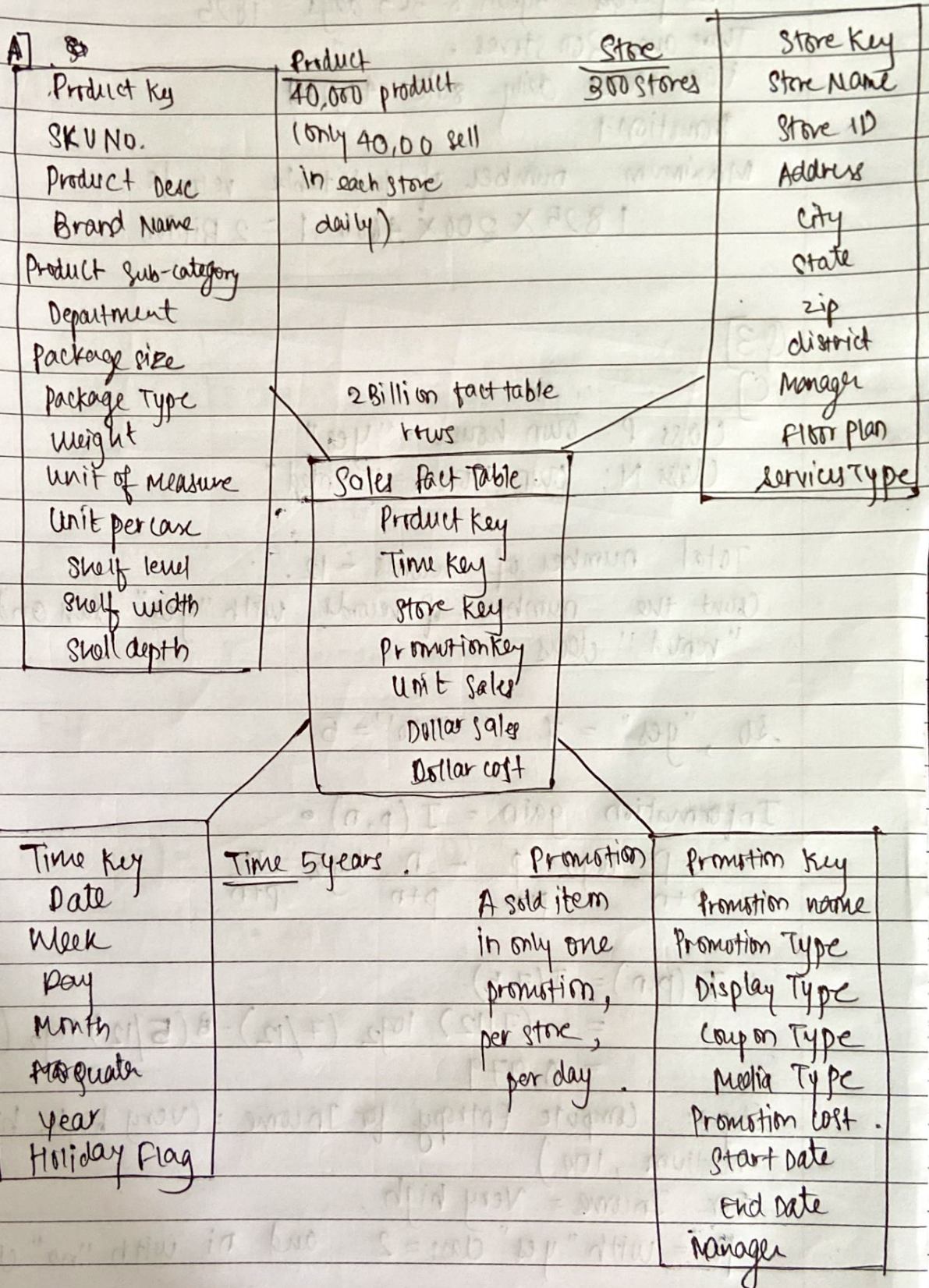


Q.3] Star schema



Time period = 5 years \times 365 days = 1825

There are 300 stores,

Each stores daily sale = 40,000

Promotion = 1

Maximum number of fact table records:

$$1825 \times 300 \times 4000 \times 1 = 2 \text{ Billion}$$

Q3]

c)

Class P: own house = "yes"

Class N: own House = "rented"

Total number of records = 12

Count the number of records with "yes" class and "rented" class

So, "yes" = 7 "no" = 5

Information gain = $I(p, n)$

$$= -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$\therefore I(p, n) = I(7, 5)$$

$$= - (7/12) \log_2 (7/12) - (5/12) \log_2 (5/12)$$

$$= 0.979$$

Step 1: Compute Entropy for Income: (Very high, high, medium, low)

For Income = Very high,

p_i = with "yes" class = 2 and n_i with "no" class = 0

Refers

7/9

$$\text{Therefore, } I(p_i, n_i) = I(2, 0) = 0$$

For Income = High,

$$p_i = \text{yes} = 4, n_i = \text{no} = 0$$

$$\therefore I(p_i, n_i) = I(4, 0) = 0$$

For Income = Medium

$$p_i = \text{yes} = 1, n_i = \text{no} = 2$$

$$\therefore I(p_i, n_i) = I(1, 2) = - \left(\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \left(\frac{2}{3} \right) \log_2 \left(\frac{2}{3} \right) \right)$$

$$= 0.918$$

For Income = Low

$$p_i = \text{yes} = 0, n_i = \text{no} = 3$$

$$\therefore I(p_i, n_i) = I(0, 3) = 0$$

$$E(A) = \sum_{i=1}^n \frac{p_i + n_i}{p + n} I(p_i, n_i)$$

$$E(\text{Income}) = \frac{2}{12} \times I(2, 0) + \frac{4}{12} \times I(4, 0) + \frac{3}{12} \times I(0, 3)$$

$$= 0.229$$

Hence

$$\text{Gain}(S, \text{Income}) = I(p, n) - E(\text{Income})$$

$$= 0.979 - 0.229$$

$$= 0.75$$

8/9

Step 2: Compute entropy for age:
(Young, medium, old)

Similarly for different age ranges $I(p_i, n_i)$
is calculated as

Age	p_i	n_i	$I(p_i, n_i)$
Young	3	1	0.811
Medium	3	2	0.971
Old	1	2	0.918

Cal. Entropy

$$E(\text{Age}) = \frac{4}{12} \times I(3,1) + \frac{5}{12} \times I(3,2) + \frac{3}{12} \times I(1,2)$$

$$= \frac{4}{12} \times 0.811 + \frac{5}{12} \times 0.971 + \frac{3}{12} \times 0.918$$

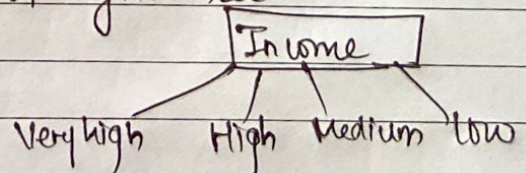
$$= 0.904$$

$$\begin{aligned} \text{Hence, Gain}(S, \text{age}) &= I(p, n) - E(\text{age}) \\ &= 0.979 - 0.904 \\ &= 0.075 \end{aligned}$$

Income attribute has the highest gain, therefore it is used as the decision attribute in the root node

Step 3: Consider income = "very high" and count the number of tuples from the original given set

$$S_{\text{very high}} = 2$$



9/9

Since both the tuples have class label = "yes"
so directly give "yes" as a class label below "very high"

Check for income = "high" and "low" having "yes"
and "rented"

for income = "medium" we have class 'yes' is 1
and rented is 2

so put age label below income = "medium"

Final decision tree is :

