



Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

## Association Mining

### Association rule

- Let us assume that the number of items in the shop shop is  $n$ .
- i.e. In the dataset there 6 items in stock, namely Bread, Milk, Diapers, Beer, eggs, Cola.
- The items list is represented  $I$  and its items represented by  $\{i_1, i_2, \dots, i_n\}$ .
- The number of transactions are represented by  $N$  transactions  $N=5$  for the shop data given.

TID	Items
1	{Bread, Milk}
2	{Bread, Diapers, Beer, Eggs}
3	{Milk, Diapers, Beer, Cola}
4	{Bread, Milk, Diapers, Beer}
5	{Bread, Milk, Diapers, Cola}

Each transaction is denoted by  $T(t_1, t_2, \dots, t_n)$  with a unique identifier (TID) and each transaction consists of a subset of items (possibly a small subset) purchased by one customer.

- Let each transaction of  $m$  items  $\{i_1, i_2, \dots, i_m\}$ , where  $m \leq n$  ( $n$  = number of items in a transaction should be less than or equal to total number of items in the shop).
- Typically, transactions differ in the number of items.





Registration No. .... Question .....

Do not write  
in  
this margin

Write on both sides of the paper

From our dataset

$T_1$  = has two items ( $m=2$ ) with  $i_1$  - bread &  $i_2$  - milk  
 $T_4$  =  $m=4$ ,  $i_1$  - Bread,  $i_2$  - Milk,  $i_3$  - Diapers,  $i_4$  - Beer

- Our task will be to find association relationships, given a large number of transactions, such that items that tend to occur together are identified.

Note: Association rules mining does not take into account the quantities of items bought.

- Metrics to Evaluate the Strength of Association Rules

- Support, Confidence, Lift.

- Let  $N$  be the total number of transactions

Support of  $X$  is represented as the number of times  $X$  appears in the database divided by  $N$ .

While support of  $X$  and  $Y$  together, is represented as the number of times they appear together divided by  $N$  as given below.

$$\text{Support}(X) = \frac{\text{Number of times } X \text{ appears}}{N} = P(X)$$

$$\text{Support}(XY) = \frac{\text{Number of times } X \text{ and } Y \text{ appear together}}{N}$$

$$= P(X \cap Y)$$

Thus support of  $X$  is the probability of  $X$  while support of  $XY$  is the probability of  $X \cap Y$ .





Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

Support (Bread) =  $\frac{\text{Number of times Bread appears}}{\text{Total number of transactions}}$

$$= \frac{4}{5} = P(\text{Bread})$$

$$\text{Support (milk)} = \frac{4}{5} = P(\text{milk})$$

$$\text{Support (Diapers)} = \frac{4}{5} = P(\text{Diapers})$$

$$\text{Support (Beer)} = \frac{3}{5} = P(\text{Beer})$$

$$\text{Support (Eggs)} = \frac{1}{5} = P(\text{Eggs})$$

$$\text{Support Cola} = \frac{2}{5} = P(\text{Cola})$$

Support (Bread, milk) =  $\frac{\text{Number of times Bread, Milk appear together}}{\text{Total number of transactions}}$

$$= \frac{3}{5} = P(\text{Bread} \cap \text{milk})$$

Support (Diapers, Beer) =  $\frac{\text{Number of times Diapers, Beer appear together}}{\text{No. of transactions}}$

$$= \frac{3}{5} = P(\text{Diapers} \cap \text{Beer})$$

- A high level of support indicates that the rule is frequent enough for the business to take interest in it.





Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

## ② Confidence.

Suppose the support that support for  $X \rightarrow Y$  is 80 % then

It means that  $X \rightarrow Y$  is very frequent and there are 80 % chances that  $X$  and  $Y$  will appear together in a transaction.

- This would of interest to the sales manager.

- Suppose another pair of items ( $A$  and  $B$ ) and support for  $A \rightarrow B$  is 50 %.

- This is not as frequent as  $X \rightarrow Y$  but if this was high such as when  $A$  appears then is 90% chances that  $B$  also appears, then of course it would be of great interest.

- Thus, not only the probability that  $A$  and  $B$  appear together matters, but also the conditional probability of  $B$  when  $A$  has already occurred play a significant role.

- This conditional probability that  $B$  will follow when  $A$  has already occurred is considered during determining the confidence of the rule.

\* The support and confidence are important metrics to judge the quality of association.

- A high level of confidence shows that the rule is true often enough to justify a decision based on it.

- Confidence for  $X \rightarrow Y$  is defined as the ratio of the support of  $X$  and  $Y$  together to the support for  $X$ . It is same as the conditional probability of  $Y$  when  $X$  has already occurred.





Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

occurred.

Reason If X appears much more frequently than X and Y appearing together, then the confidence will be low.

$$\text{Confidence of } (X \rightarrow Y) = \frac{\text{Support}(XY)}{\text{Support}(X)}$$

$$= \frac{P(X \cap Y)}{P(X)} = P(Y|X)$$

Confidence  $P(Y|X)$  is the probability of Y once X has taken place, also called the conditional probability.

Example of the support measure.

TID	Items	Support = Occurrence / Total Support.
1	ABC	Total support = 5
2	ABD	
3	BC	Support{AB} = $2/5 = 40\%$
4	AC	Support{BC} = $3/5 = 60\%$
5	BCD	Support{ABC} = $1/5 = 20\%$

1	Given $X \Rightarrow Y$	
2	Confidence - Occurrence {X, and Y} / occ(X)	
3	$A \Rightarrow B = 2/3$	
4	$B \Rightarrow C = 3/4$	
5	$AB \Rightarrow C = 1/2$	
6		





Registration No. ....

Question .....

Do not write  
in  
this margin

Write on both sides of the paper

Database for identification of association rule.

Antecedent

Consequent

A

0

A

0

A

1

A

0

B

1

B

0

B

1

There are two rules derived from the association of the combined

Rule 1: A implies 0 i.e.  $A \rightarrow 0$ Rule 2: B implies 1 i.e.  $B \rightarrow 1$ The support for rule =  $\frac{\text{Number of times A \& 0 appear}}{\text{Total Number of transactions}}$ 

$$\text{Support of Rule 1} = \frac{3}{7}$$

$$\text{Support for rule 2} = \frac{2}{7}$$

Confidence for Rule 1 is support of (A, 0) / support of A

$$\frac{\text{support of (A, 0)}}{\text{support of A}} = \frac{3/7}{4/7} = \frac{3}{4}$$

$$\text{Support (A)} = \frac{4}{7} = \frac{3}{7} + \frac{1}{7} = \frac{3}{7} \times \frac{1}{4} = \frac{3}{4}$$





Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

Confidence for Rule 2

$$= \text{Support}(B_1) / \text{Support}(B)$$

$$= \text{Support}(B_1) = \frac{2}{7} \quad \text{Support}(B) = \frac{3}{7}$$

$$\frac{2/7}{3/7} = \frac{2}{3}$$

Lift

- It is very important to consider the frequency of  $Y$  or  $P(Y)$  for the effectiveness of the association mining rule  $X \rightarrow Y$

- Confidence is the conditional probability of  $Y$  when  $X$  has already occurred.

- It is very important to consider how frequent  $Y$  is to gauge the effectiveness of the association mining confidence.

Thus

Lift is the ratio of conditional probability of  $Y$  when  $X$  is given to the unconditional probability of  $Y$  in the dataset.

- Simply

It is the confidence of  $X \rightarrow Y$  divided by the probability of  $Y$





Registration No. ....

Question .....

Do not  
in  
this mark

Write on both sides of the paper

$$\text{Lift} = P(Y/X) / P(Y)$$

or

$$\text{Lift} = \text{Confidence of } (X \rightarrow Y) / P(Y)$$

or

$$\text{Lift} = (P(X \cap Y) / P(X)) / P(Y)$$

This lift can be computed by dividing the confidence by the unconditional probability of consequent Y.

- Coke is a very common sales item in a store and appears in most of the transactions.

- Candle  $\rightarrow$  Coke = has support of 20%  
Confidence of 90%.

- logically we think that is very popular and it appears in 95% of transactions.

- then obviously it also appears quite often with the candle as well.

- So the rule of association of candle and coke will not be useful.

• Candle  $\rightarrow$  Matchbox - support 20%  
- Confidence 90%.

Suppose the frequency of Matchbox is very little compared to the sale of Coke.

- Rule suggest that when we make a sale of Candles, 90% chances indicate a matchbox will be sold in the same transaction.





Registration No. .... Question .....

Write on both sides of the paper

Do not write  
in  
this margin

It is more effective and logical to conclude that when we sell a candle then we also sell a Coke (Coke is popular and will appear with every item not just with candle). As support and confidence are unable to handle this case. It is handled by the lift.

In the case probability of Y is very low In the land (Candle  $\rightarrow$  Matchbox)  
(Y Matchbox)

Probability of Y is very high Candle  $\rightarrow$  Coke (Here Y, is high)

$$\text{lift} = P(Y/X) / P(Y)$$

or

$$\text{lift} = \text{Confidence of } (X \rightarrow Y) / P(Y)$$

or

$$\text{lift} = (P(X \cap Y) / P(X)) / P(Y)$$

- High probability of Y i.e.  $P(Y)$  makes lift less effective and low value makes it more effective

lift for Rule 1  $A \rightarrow O$

$$P(A, O) \frac{3}{4}, P(A) \frac{4}{4}, P(O) \frac{4}{4}$$

$$\text{Confidence of } A \rightarrow O = \frac{3/4}{4/4} = \frac{3/4 \times 4/4}{4/4} = \frac{21}{16}$$

$$= 1.3125$$





Registration No. ....

Question .....

Do not write  
in  
this margin

Write on both sides of the paper

$$\text{Lift for Rule 2 } B \rightarrow 1 = \frac{2/3}{3/4} = \frac{2}{3} \times \frac{4}{3} = \frac{8}{9} = 0.88$$
$$\text{Confidence } (2/3) / P_1(3/4)$$

Confidence for Rule 1 is  $3/4 = 0.75$

Rule 2 is  $2/3 = 0.66$

It should be observed that although Rule 1 has a higher confidence as compared to Rule 2, it has a lower lift as compared to Rule 2.

Naturally,

it would appear that Rule 1 is more valuable because of having higher confidence; it appears more accurate. But the accuracy of the rule can be misleading if it is independent of the dataset.

- Lift, as a metric is important because it considers both the confidence of the rule and the overall dataset.