

21/5/19

PAGE NO.	
DATE	/ /

Association Rule - Correlation between two or more items.

Define relationship between various items

Market basket analysis - Evaluating market strategies based on the sell and applying association rule mining on the data set.

Support denotes the probability that a transaction contains 2 item set while, confidence denotes the probability that a transaction containing item set 1 also contains item set 2.

Basic Concept of Association Rule

Let  $T = \{i_1, i_2, i_3, \dots, i_m\}$

$D$  = Task Relevant Data = Transaction having TID

$A$  = Itemset } one or more items.  
 $B$  = Itemset }

$A \Rightarrow B$  where  $A \subset T$  &  $A \cap B = \emptyset$   
 $B \subset T$

Support =  $P(A \cup B)$

Confidence =  $P\left(\frac{B}{A}\right)$

Ex:  $T_1 : A, B, C$   
 $T_2 : A, C$   
 $T_3 : A, B, D$

→ both A & B present

For  $A \Rightarrow B$  : Support = 2 (66%)

Confidence = 66% ( $\frac{2}{3} \times 100$ )

if A is present then B is present

But, when  $B \rightarrow A$

Support = 66%

Confidence = 100%

Classification of Association Rule (AR)  
i) Based on the type handled by AR

a) Boolean

0 or 1 value

Ex: age(30)  $\rightarrow$  buys (laptop)

b) Quantitative

Ex: age(X, '25-40')  $\wedge$  income(X, '40k-50k')  
 $\rightarrow$  buys (X, computer)

ii) Dimension of data involved in rule

a) Single dimension

Ex: age(40)  $\rightarrow$  buys(car)

b) Multi-dimensional

Ex: age(40)  $\wedge$  gender(male)  $\rightarrow$  buys(car)

iii) Level of abstraction

a) Higher level

Ex: occupation (students)  $\rightarrow$  buys(computer)

b) Lower level

Ex: occupation (students)  $\rightarrow$  buys(laptop)

iv) Based on various expansion

$T_1 : \{a_1, a_2, \dots, a_{50}\}$

$T_2 : \{a_1, a_2, \dots, a_{100}\}$

1-itemset

$\{a_1\}$  support 2

$\{a_2\} : 2$

$\vdots$

$\{a_{99}\} : 1$

$\vdots$

$\{a_{100}\} : 1$

2-itemset

$\{a_1, a_{99}\} : 2$

$\{a_1, a_{98}\} : 1$

$\vdots$

$\{a_{99}, a_{100}\} : 1$

$\vdots$

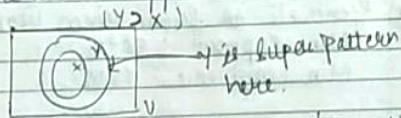
$\{a_{99}, a_{100}\} : 1$

Confidence  
not possible to  
find)

$$\text{Total itemsets (subsets)} = 2^n - 1 = 2^{100} - 1$$

a) Closed pattern (Loss-less)

A pattern (itemset)  $X$  is a closed pattern if  $X$  is frequent & there exists no super pattern  $Y$  with the same support as  $X$ .



Here, (support)

$P_1 = \{a_1, a_2 \dots a_{99}\} : 2$

$P_2 = \{a_1, a_2 \dots a_{100}\} : 1$

$\{a_1, a_{99}\} : 1 \quad \{a_1, a_{98}\} : 2$

b) Max pattern

A pattern  $X$  is a max pattern if  $X$  is frequent and there exists no frequent super pattern  $Y$ .

Here,  $P_1 = \{a_1, a_2 \dots a_{100}\} : 1$

$\{a_1, a_{98}\} : 1$

Thus, lossy representation.

Find AR from large dataset :

Step 1 : Find frequent dataset itemset.

Step 2 : Find a strong association rule from the frequent itemset

→ If the support of a itemset satisfies the minimum <sup>min</sup> threshold, then it is called a frequent item set.

⇒ If the confidence of a item set satisfies the minimum confidence threshold then it is called a strong AR.

Apriori algorithm

Any subset of a frequent itemset should be frequent. This is the principle of apriori

Ex : Generate AR from given Transaction  
 Min. support : 50%  
 Min. confidence : 75%.

T <sub>1</sub>	Bread, Cheese, Juice, Butter
T <sub>2</sub>	Bread, Cheese, Juice
T <sub>3</sub>	Bread, Milk, Yogurt
T <sub>4</sub>	Bread, Juice, Milk
T <sub>5</sub>	Cheese, Juice, Milk

⇒ Candidate creation:

Item	Support
Bread	4 (80%)
Cheese	3 (60%)



given

Juice	4	(80%)	
Butter	1	(20%)	x (Support < 50%)
Milk	3	(60%)	
Yogurt	1	(20%)	x

L1 Item Support

Bread	4
Cheese	3
Juice	4
Milk	3

L2 Item Set Support

{B, C}	2	(40%)	x
{B, J}	3	(60%)	
{B, M}	2	(40%)	x
{C, J}	3	(60%)	
{C, M}	1	(20%)	x
{J, M}	2	(40%)	x

L2 Item Set Support

{B, J}	3
{C, J}	3

We skip creating C3 from L2 because if have itemsets having 3 items, the support will be only 1

AR possible from L2

Bread $\rightarrow$ Juice	75%
Juice $\rightarrow$ Bread	75%
Cheese $\rightarrow$ Juice	$\frac{3}{3} = 100\%$
Juice $\rightarrow$ Cheese	$\frac{3}{4} = 75\%$

$\therefore$  All are strong AR.

Support  $\approx$  Occurrence of an item.

PAGE NO.	
DATE	/ /

$$\text{confidence } (A \rightarrow B) = \frac{\text{Support } (A \cup B)}{\text{Support } (A)}$$

22/8/19 Confidence of Bread  $\rightarrow$  Juice

Support = 3      Support = 4

$\therefore \text{Confidence} = \frac{3}{4} = 75\%$

Q. Transaction Items

T <sub>1</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>
T <sub>2</sub>	I <sub>2</sub> , I <sub>4</sub>
T <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>4</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>
T <sub>5</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>6</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>7</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>8</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>
T <sub>9</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>

Min. Support = 22%

Min. Confidence = 70%

Min. Support = 22%  $\approx 1.1$   
= 2

$\Rightarrow$  C1.

Item	Support
I <sub>1</sub>	6 (66%)
I <sub>2</sub>	7 (77%)
I <sub>3</sub>	6 (66%)
I <sub>4</sub>	2 (22%)
I <sub>5</sub>	2 (22%)

L1

Item	Support
I <sub>1</sub>	6
I <sub>2</sub>	7
I <sub>3</sub>	6
I <sub>4</sub>	2
I <sub>5</sub>	2

C2.

ItemSupport

$I_1, I_2$	4	
$I_1, I_3$	4	
$I_1, I_4$	1	X
$I_1, I_5$	2	
$I_2, I_3$	4	
$I_2, I_4$	2	
$I_2, I_5$	2	
$I_3, I_4$	0	X
$I_3, I_5$	1	X
$I_4, I_5$	0	X

C2.

Item SetSupport

$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_1, I_5\}$	2
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2
$\{I_2, I_5\}$	2

C3

Item SetSupport

$\{I_1, I_2, I_3\}$	2	
$\{I_1, I_2, I_5\}$	2	
$\{I_1, I_2, I_4\}$	1	X
$\{I_1, I_3, I_5\}$	1	X
$\{I_1, I_2, I_3, I_4\}$		
$\{I_1, I_2, I_3, I_5\}$		
$\{I_1, I_2, I_4, I_5\}$		
$\{I_2, I_3, I_4\}$	0	X
$\{I_2, I_3, I_5\}$	1	X
$\{I_2, I_4, I_5\}$	0	X

don't  
write sets  
with 4 items

L3	Intersect	Support
	$\{I_1, I_2, I_3\}$	2
	$\{I_1, I_2, I_5\}$	2

RULES:

From I2:	$I_1 \Rightarrow I_2$	$4/6 = 66\%$	$I_2 \Rightarrow I_1$	$4/6 = 66\%$
	$I_1 \Rightarrow I_3$	$4/6 = 66\%$	$I_3 \Rightarrow I_1$	$4/6 = 66\%$
	$I_1 \Rightarrow I_5$	$2/6 = 33\%$	$I_5 \Rightarrow I_1$	$2/2 = 100\% \checkmark$
	$I_3, I_2 \Rightarrow I_3$	$2/4 = 50\%$	$I_1, I_2 \Rightarrow I_2$	$2/4 = 50\%$
	$I_3, I_2 \Rightarrow I_1$	$2/4 = 50\%$	$I_1, I_2 \Rightarrow I_5$	$2/4 = 50\%$
	$I_1, I_5 \Rightarrow I_2$	$2/2 = 100\%$		

$I_2 \Rightarrow I_3$	$4/7 = 57\%$	$I_3 \Rightarrow I_2$	$4/6 = 66\%$
$I_2 \Rightarrow I_4$	$2/7 = 28\%$	$I_4 \Rightarrow I_2$	$2/2 = 100\% \checkmark$
$I_2 \Rightarrow I_5$	$2/7 = 28\%$	$I_5 \Rightarrow I_2$	$2/2 = 100\% \checkmark$

From I3:

$\{I_1, I_2\} \Rightarrow I_3$	$2/4 = 50\%$	$I_3 \Rightarrow \{I_1, I_2\}$	$2/6 = 33\%$
$\{I_1, I_3\} \Rightarrow I_2$	$2/4 = 50\%$	$I_2 \Rightarrow \{I_1, I_3\}$	$2/7 = 28\%$
$\{I_2, I_3\} \Rightarrow I_1$	$2/4 = 50\%$	$I_1 \Rightarrow \{I_2, I_3\}$	$2/6 = 33\%$
$\{I_1, I_2\} \Rightarrow I_5$	$2/4 = 50\%$	$I_5 \Rightarrow \{I_1, I_2\}$	$2/2 = 100\%$
$\{I_1, I_5\} \Rightarrow I_2$	$2/2 = 100\% \checkmark$	$I_2 \Rightarrow \{I_1, I_5\}$	$2/7 = 28\%$
$\{I_2, I_5\} \Rightarrow I_1$	$2/2 = 100\% \checkmark$	$I_1 \Rightarrow \{I_2, I_5\}$	$2/6 = 33\%$

Strong AR:

- $I_5 \Rightarrow I_1$
- $I_4 \Rightarrow I_2$
- $I_5 \Rightarrow I_2$
- $\{I_1, I_5\} \Rightarrow I_2$
- $\{I_2, I_5\} \Rightarrow I_1$
- $I_5 \Rightarrow \{I_1, I_2\}$



Improving the efficiency of Apriori's methods.

- i) Shrink the no. of candidates
  - a) Hash-base technique
  - b) Sampling

- ii) Reduce the number of transaction database scan
  - a) Transaction reduction
  - b) Partitioning
  - c) Dynamic item set counting

⇒ Hash-base technique

$$h(x, y) = \text{position of } x \times 10 + (\text{index of } y) \pmod{7}$$

For  $T_1$ , 3 item sets with 2 items are possible

$$h(I_1, I_2) = [(1 \times 10) + 2] \% 7 = 5 \quad \rightarrow \text{place in bucket 5}$$

$$h(I_1, I_6) = [(1 \times 10) + 6] \% 7 = 1$$

$$h(I_2, I_5) = [(2 \times 10) + 5] \% 7 = 4$$

Bucket index	0	1	2	3	4	5	6	$\neq$
Bucket count	2	2	3	2	2	4	4	
Bucket content	$I_1, I_4$	$I_1, I_5$	$I_2, I_3$	$I_2, I_4$	$I_2, I_5$	$I_1, I_2$	$I_1, I_3$	
		$I_3, I_5$	$I_2, I_3$	$I_2, I_4$	$I_1, I_5$	$I_1, I_2$	$I_1, I_3$	
			$I_2, I_5$			$I_1, I_4$	$I_1, I_3$	
						$I_1, I_2$	$I_1, I_3$	

Now, if min. support > bucket count, the candidate can be discarded.

## 26/8/17 Transaction Reduction

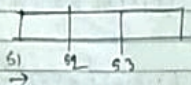
Abscise the trans. which don't have the frequent item sets.

### Partitioning

Divide the entire data set into multiple parts and find freq. item sets from each.

### Dynamic item set counting

Include  $n$ -item sets in the first scan and calculate its support. In subsequent scans, we simply add the pre-calculated support.



Scan part (S1-S2) and list all  $n$  itemsets with their support.

On next scan (S2-S3), add the support for corresponding item set.

⇒ when we have  $10^4$  1-itemset, we get  $10^7$  2-itemset.

finding freq. itemsets without candidate generation.

Real one: i) It may need to generate huge no. of candidates.

ii) It may repeatedly require to scan the

entire dataset.

Frequent Pattern Growth (FP Tree)  
we can generate 2-itemset using this  
method useful when  $k > 1$

$T_1$	$A_1, A_2, A_5$	$A_1$	6
$T_2$	$A_2, A_4$	$A_2$	7
$T_3$	$A_2, A_3$	$A_3$	6
$T_4$	$A_1, A_2, A_4$	$A_4$	2
$T_5$	$A_1, A_3$	$A_5$	2
$T_6$	$A_2, A_3$		
$T_7$	$A_1, A_3$		
$T_8$	$A_1, A_2, A_3, A_5$		
$T_9$	$A_1, A_2, A_3$		

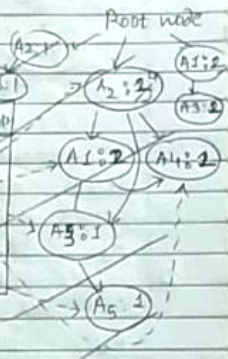
Order

$A_2, A_1, A_3, A_4, A_5$

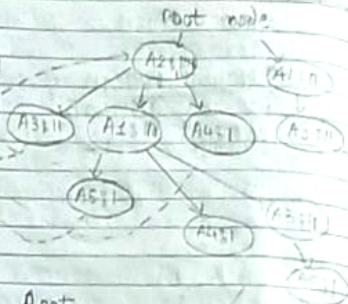
Reorder all trans.

$T_1$	$A_2, A_1, A_5$
$T_2$	$A_2, A_4$
$T_3$	$A_2, A_3$
$T_4$	$A_2, A_1, A_4$
$T_5$	$A_1, A_3$
$T_6$	$A_2, A_3$
$T_7$	$A_1, A_3$
$T_8$	$A_2, A_1, A_3, A_5$
$T_9$	$A_2, A_1, A_3$

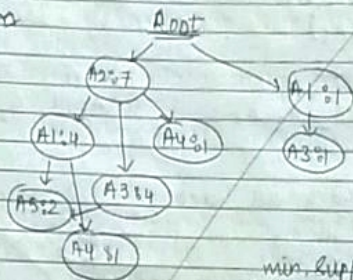
Item	Sup	Header
$A_2$	7	
$A_1$	6	
$A_3$	6	
$A_4$	2	
$A_5$	2	



Item	Supplier	Printer
A <sub>1</sub>	7	-
A <sub>2</sub>	6	-
A <sub>3</sub>	6	-
A <sub>4</sub>	2	-
A <sub>5</sub>	2	-



Ex: 950m



min. supply (given) = 2  
8 9 12 3 has supply  
3

Item	Conditional Pattern Base	Conditional Tree	Frequent Tree
$A_3$	$\{A_2, A_1\}$ $\{A_2, A_3\}$ $\{A_2, A_1, A_3\}$	$\{A_2\}$ $\{A_1\}$ $\{A_1, A_3\}$ $\{A_3\}$	
$A_4$	$\{A_2, A_1\}$ $\{A_2\}$		
$A_3$	$\{A_2, A_1\}$ $\{A_1\}$ $\{A_2\}$		
$A_1$	$\{A_2\}$		



# From Conditional Pattern Base

$A_5$   
 $\{A_2, A_1\} : 1$   
 $\{A_2, A_1, A_3\} : 2$   
 Pattern leading upto  $A_5$  & the count of  $A_5$  node

## Conditional FP Tree

$\{A_1\} : 2$   
 $\{A_1, A_2\} : 2$   
 Freq. item  
 Set (min support = 2)

## Frequent

$\{A_1, A_5\} : 2$   
 $\{A_2, A_5\} : 2$   
 $\{A_1, A_2, A_5\} : 2$   
 Add  $A_5$  to set obtained in CFP

$A_1$   
 $\{A_2\} : 1$   
 $\{A_2, A_1\} : 1$

$\{A_2\} : 2$

$\{A_2, A_5\} : 2$

$A_3$   
 $\{A_2\} : 2$   
 $\{A_2, A_1\} : 2$   
 $\{A_1\} : 2$

$\{A_2\} : 4$   
 $\{A_1\} : 4$   
 $\{A_1, A_2\} : 2$

$\{A_2, A_5\} : 4$   
 $\{A_1, A_5\} : 4$   
 $\{A_1, A_2, A_5\} : 2$

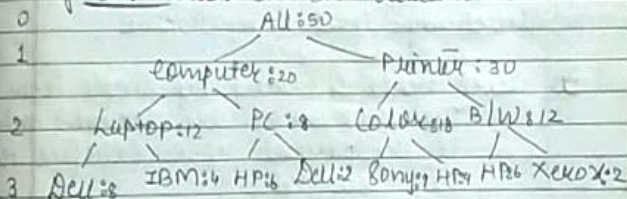
$A_1$   
 $\{A_2\} : 2$

$\{A_2\} : 4$

$\{A_2, A_5\} : 4$

$A_2$

## Mining multi-level association rules from transactional databases



Thus, we have 3 levels of abstraction here.

buy(x, "computer")  $\Rightarrow$  buy(x, "printer")  
support = 50 (High Level AR)

buy(x, "laptop")  $\Rightarrow$  buy(x, "color")  
support = 30 (Lower Level AR)

buy(x, "Dell")  $\Rightarrow$  buy(x, "Sony")  
support = 17

Approaches to mine multi-level AR

i) Using uniform min. support for all the levels

- Need to declare only one global min. support
- The items at lower level can be missed

ii) Using the reduced min support at lower level

- Decreasing the min support at every level

iii) Level by level independence

1) Level wise filtering by single item

- The child nodes are scanned only if the parent node satisfies the min. support
- No. of scans are reduced
- We can miss out on the child nodes which satisfy the min. support but, its parent does not.

Level wise filtering by # item

iv) Controlled level wise filtering  
The concept of Passage min support is used

$$\min \text{ support of Parent} \leq \text{Passage min support} < \min \text{ support of child}$$

Then for every node in the level, we check for the Passage min support

$$\begin{aligned} \Rightarrow \text{buys}(X, \text{computer}) &\Rightarrow \text{buys}(X, \text{printer}) & [s=20, c=70] \\ \text{buys}(X, \text{laptop}) &\Rightarrow \text{buys}(X, \text{printer}) & [s=7, c=42] \end{aligned}$$

are redundant AR

Thus, we can discard the lower level rule if it gives the expected output

Expected values:

From first rule laptop accounts for almost  $\frac{1}{3}$ <sup>rd</sup> of total sale. Therefore an expected support for laptop would be  $\frac{20}{3}$  (6.66) approx

and confidence near to 70, which is satisfied. Thus, second rule can be discarded.

Multidimensional rules

Finding Multidimensional AR from RDB & DW

$$\text{age}(X, 25-30) \wedge \text{buys}(X, \text{computer}) \Rightarrow \text{buys}(X, \text{printer})$$

is a hybrid multidimensional AP as  
 'buys' predicate / dimension is repeated  
 $\text{age}(x, 25-30) \wedge \text{income}(x, 25K-55K) \Rightarrow$   
 $\text{buys}(x, \text{"laptop"})$   
 (No predicate / dimension repeated)

Methods for mining MDAR

- i) Quantitative attributes are discretized using concept hierarchy.
- ii) Quantitative attributes are discretized into fields based on the distribution of data.
- iii) Quantitative attributes are discretized to capture the semantic meaning of such data.

Purchase of IBM PC: 50

min. support = 14

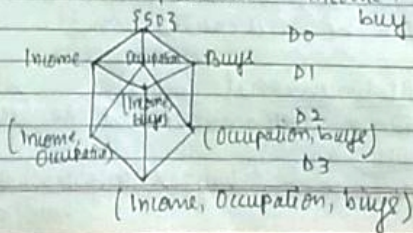
age : 12 x (based on min. support)

Income : 14

Gender : 6 x

Occupation : 18

$\therefore$  Total 3 predicates : Income + Occupation + buys





## Mining Quantitative AR (Binning method)

- 1) Binning
- 2) finding freq predicate set
- 3) clustering the association rule