

# UNIT-1 Association Analysis

10/1/23

\* Association Analysis can be applied to business domain. It can also be used in Healthcare domain, finance.

\* Dataset in Association Analysis can be represented in the form of transactions. Each transaction consists of the items that are bought in that transaction.

Difficulties with Association:-

\* huge data (handling with large data)

\* handling weak transaction (unintentionally, unconditionally)  $\Rightarrow$  milk + pen

\* The entire dataset can be represented in binary format where each row represents a transaction and each column represents an item.

example:-

		milk	sugar	bread	T.p	cp
T <sub>1</sub>	milk + sugar + bread	1	1	1	0	0
T <sub>2</sub>	milk + Teapowder + sugar	1	1	0	1	0
T <sub>3</sub>	milk + sugar + coffee powder	1	1	0	0	1
T <sub>4</sub>	milk + sugar + bread + jam	0	1	1	0	0

\* A collection of zero or more items is known as an item set

\* k-item set, An item set with k items is known as a k-item set

13/1/23

General Rule format will be  $A \rightarrow B$  (A derives B)

where A is known as Antecedent and B is known as consequent

$\Rightarrow$  {milk, coffee}  $\rightarrow$  sugar

for A may contain one or more items.

Support:-

An Association rule can be measured in terms of support:

$\rightarrow$  support determines how often a rule is applicable for a given data set.

Confidence:-

Confidence is a measuring parameter for Association Rule.

$\rightarrow$  confidence determines how frequently the items in B appear in the transactions that contain A.

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

$N \rightarrow$  total no. of transactions in the given dataset.

ex:-

Tid	Items
T1	I <sub>1</sub> , I <sub>2</sub>
T2	I <sub>1</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub>
T3	I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>6</sub>
T4	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub>
T5	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>

calculate the support for the Association  $\{I_2, I_3\} \rightarrow I_4$



$$= \frac{\text{count}(I_2, I_3, I_4)}{N}$$

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

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

$$\Rightarrow \frac{\frac{\text{count}(A \cup B)}{N}}{\frac{\text{count}(A)}{N}} \Rightarrow \frac{\text{count}(A \cup B)}{\text{count}(A)}$$

$$\frac{\text{count}(I_2, I_3, I_4)}{\text{count}(I_2, I_3)} \Rightarrow \frac{2}{3} =$$

$$\text{count}(A) \rightarrow \text{count}\{I_2, I_3\}$$

eg:-  $\{I_1, I_3\} \rightarrow I_5$

$$\text{support}(\{I_1, I_3\} \rightarrow I_5) = \frac{\text{count}(I_1 \cup I_3 \cup I_5)}{N}$$

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

$$\text{confidence}(\{I_1, I_3\} \rightarrow I_5) = \frac{\text{support}(\{I_1, I_3\} \rightarrow I_5)}{\text{support}(\{I_1, I_3\})}$$

$$= \frac{\text{count}(I_1 \cup I_3 \cup I_5)}{\text{count}(I_1, I_3)}$$

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

eg:-  $I_1 \rightarrow \{I_3, I_5\}$

$$\text{support}(\{I_1\} \rightarrow \{I_3, I_5\}) = \frac{\text{count}(I_1 \cup I_3 \cup I_5)}{N}$$

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

$$\text{Confidence}(I_1 \rightarrow \{I_3, I_5\}) = \frac{\text{count}(I_1 \cup I_3 \cup I_5)}{\text{count}(I_1)}$$

$$= \frac{2}{4}$$

Note:-

\* Support is used to eliminate uninteresting rules.

\* Confidence is used to measure the reliability of the Inference made by the rule.

\* **STRONG RULE:-** A rule that satisfy both minimum support threshold and minimum confidence threshold is known as a strong rule.

Apriori Algorithm:-

Step 1:- frequent item set generation

Step 2:- Rule generation

In the step-1 we find all the item sets that satisfy minimum support threshold. These item sets are known as frequent item sets.

14/7/23

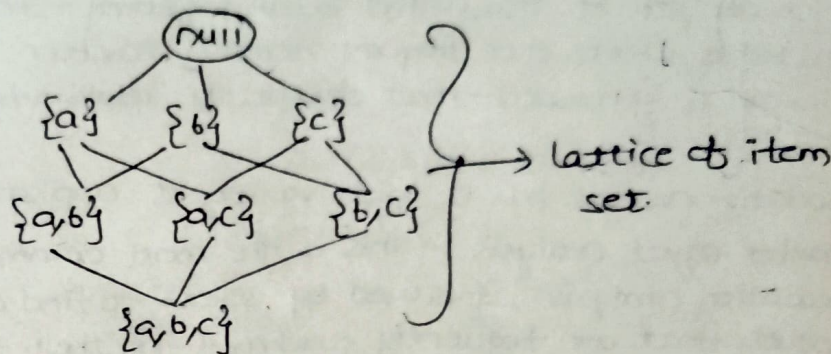


In the step 2 we extract all the  $i$  confidence rules from the frequent item sets. These rules are known as strong rule frequent item set generation:-

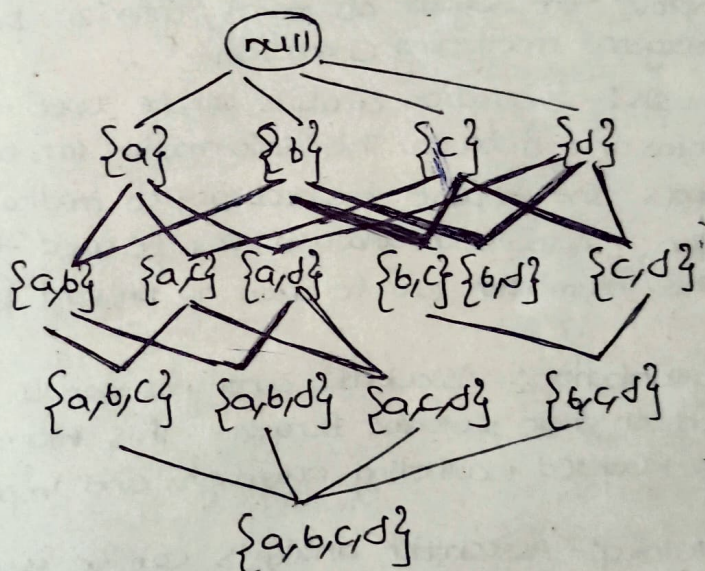
for a dataset with  $k$  items, there will be  $2^k - 1$  frequent item sets.

eg:- a, b, c  $[k=3]$

we get  $\rightarrow \{a\}, \{b\}, \{c\}, \{a,b\}, \{a,c\}, \{b,c\}, \{a,b,c\} \rightarrow$  item sets



eg:-  $k=4$



As  $k$  might be very large the number of frequent item sets that need to be explore may become exponentially large, in order to reduce this complexity we use apriori principle.

**Apriori principle 1:-**

All the non empty subsets of a frequent item set must also be frequent.

for eg:- for a four-item set  $\{b,c,d\}$  is frequent then all its subsets are frequent. ( $\{b\}, \{c\}, \{d\}, \{b,c\}, \{c,d\}, \{b,d\}$ )

**Apriori principle 2:-**

If an item set is not frequent then all its supersets are considered to be non-frequent.

for eg:- if  $\{b\}$  is considered to be non-frequent then all its supersets are non-frequent ( $\{a,b\}, \{b,c\}, \{b,d\}, \{a,b,c\}, \{a,b,d\}, \{b,c,d\}, \{a,b,c,d\}$ ).  
for 4-item data set.



Association Analysis? Explain various applications where association analysis can be applied.

Association analysis is a data mining technique that identifies relationships between variables in a data set. It is often used to find associations between products that are frequently purchased together such as milk and bread.

Association analysis works by identifying frequent itemsets, which are sets of items that occur together more often than would be expected by chance. Once frequent itemsets have been identified, association rules can be generated that describe the relationships between these itemsets.

Association analysis has a wide variety of applications, including:-

- **Market Basket Analysis**:- This is the most common application of association analysis. It is used by sellers to find associations between products that are frequently purchased together. This information can be used to improve the layout of stores, optimize product placement, and develop targeted marketing campaigns.
  - **Medical diagnosis**:- Association analysis can be used to find associations between symptoms and diseases. This information can be used to develop diagnostic tools and improve the accuracy of medical diagnoses.
  - **Fraud detection**:- Association analysis can be used to detect fraudulent transactions. This information can be used to prevent fraud and protect customers.
  - **Customer Segmentation**:- Association analysis can be used to segment customers based on their purchase behaviour. This information can be used to develop targeted marketing campaigns and improve customer service.
  - **Web Usage Mining**:- Association analysis can be used to find associations between web pages that are frequently visited together. This information can be used to improve the design of websites and develop targeted advertising campaigns.
- Some additional examples are product recommendations, social media analysis, logistics etc...

20/1/23

\*We have 2 phases in Apriori:-

- 1) Join
- 2) Prune ~~\*\*\*\*~~ Example-1

Transactional Data for an All Electronics Branch

TID

T100

T200

T300

T400

T500

T600

T700

T800

List of Item IDs

11, 12, 15

12, 14

12, 13

11, 12, 14

11, 13

12, 13

11, 13

11, 12, 13, 15

T900

11, 12, 13



minimum support count is 2.

candidate item set (C)

(selected list)

$$C_1 = L_1 \bowtie L_1$$

$$C_2 = L_1 \bowtie L_1$$

$L_1 \bowtie L_1 = (1-1)$  Items similarity in  $L_1$

$$L_2 = \{(I_1, I_2), (I_1, I_3), (I_2, I_3)\}$$

$$L_2 = \{(I_1, I_2, I_3)\}$$

$$C_{1+2} = L_1 \bowtie L_1$$

Item	support
✓ $I_1$	6
✓ $I_2$	4
✓ $I_3$	6
✓ $I_4$	2
✓ $I_5$	2

Item	support
$I_1$	6
$I_2$	4
$I_3$	6
$I_4$	2
$I_5$	2

let us suppose  $L_2$

$$\{(I_1, I_2), (I_1, I_3), (I_2, I_3), (I_1, I_4)\}$$

$$= \{(I_1, I_2, I_3), (I_1, I_2, I_4)\}$$

min<sup>m</sup> count/support  
\* lexicographic order

$$C_2 = L_1 \bowtie L_1$$

$$L_2$$

$$C_3 = L_2 \bowtie L_2$$

Item	support
$(I_1, I_2)$	4
$(I_1, I_3)$	4
x $(I_1, I_4)$	1
$(I_1, I_5)$	2
$(I_2, I_3)$	4
$(I_2, I_4)$	2
$(I_2, I_5)$	2
x $(I_3, I_4)$	0
x $(I_3, I_5)$	1
x $(I_4, I_5)$	0

Item	support
$(I_1, I_2)$	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

Item	support
$(I_1, I_2, I_3)$	2
<del><math>(I_1, I_2, I_4)</math></del>	
$(I_1, I_2, I_5)$	2
<del><math>(I_1, I_3, I_4)</math></del>	
<del><math>(I_1, I_3, I_5)</math></del>	
<del><math>(I_1, I_4, I_5)</math></del>	
<del><math>(I_2, I_3, I_4)</math></del>	
<del><math>(I_2, I_3, I_5)</math></del>	
<del><math>(I_2, I_4, I_5)</math></del>	

\* We need to eliminate certain item sets based on the apriori principle.

(frequency)

for the item set  $\{I_1, I_3, I_5\}$

all its subsets are not frequent.  $\{(I_1, I_3), (I_1, I_5), (I_3, I_5)\}$  i.e., they are not present in  $L_2$ . so the item  $\{I_1, I_3, I_5\}$  is considered to be nonfrequent

$$C_3 = L_2 \bowtie L_2$$

$$C_3 = \{(I_1, I_2, I_3), (I_1, I_3, I_5), (I_2, I_3, I_4), (I_2, I_3, I_5), (I_2, I_4, I_5)\}$$

$$L_3$$

Item	support
$(I_1, I_2, I_3)$	2
$(I_1, I_2, I_5)$	2

## 11/7/23 JOIN STEP

→ In this step we find  $L_k$  a set of candidate k item sets by joining  $L_{k-1}$  with itself. The resulting set after performing join is denoted as  $C_k$ . The join  $L_{k-1} \bowtie L_{k-1}$  is performed where the members of  $L_{k-1}$  are joinable.

→ The members of  $L_{k-1}$  are said to be joinable if their first  $k-2$  items are common.

## PRUNE STEP

→  $C_k$  is the superset of  $L_k$ , so it may even contain non-frequent item sets. All the candidates having count greater than the minimum support count are included in  $L_k$ .



→ The item sets in  $C_k$  are included if all its subsets are frequent.

$$C_4 = L_3 \bowtie L_3$$

$= \{(I_1, I_2, I_3, I_4, I_5)\} \times \rightarrow$  becoz its subset  $\{I_1, I_3, I_5\}$  is not frequent.

when only the  $C_k = \emptyset$  or {empty} the algorithm terminates.

Ex:- T100 {M, O, N, K, E, Y} minimum support is 3.

T200 {O, O, N, K, E, Y}

T300 {M, A, K, E}

T400 {M, U, C, K, Y}

T500 {C, O, O, K, I, E}

$$\frac{1+1+1+1+1+1}{+1+1+1+1}$$

candidate Itemset

<u>L<sub>1</sub></u>		<u>L<sub>2</sub></u>		<u>C<sub>2</sub> = L<sub>1</sub> ⋈ L<sub>1</sub></u>		<u>L<sub>2</sub></u>	
Item	support	Item	support	Item	support	Item	support
M	3	M	3	X(M, O)	1	(M, K)	3
O	4	O	4	(M, K)	3	(O, K)	4
X N	2	K	5	X(M, E)	2	(O, E)	3
K	5	E	4	X(M, Y)	2	(K, E)	4
E	4	Y	3	(O, K)	4	(K, Y)	3
Y	3			(O, E)	3		
X D	1			X(O, Y)	2		
X A	1			(K, E)	4		
X U	1			✓(K, Y)	3		
X C	2			X(E, Y)	2		
X I	1						

$$C_3 = L_2 \bowtie L_2$$

$\{(O, K, E), (K, E, Y)\}$  subset is not frequent  $(K, E, Y) \rightarrow (E, Y)$

<u>L<sub>2</sub></u>		<u>L<sub>3</sub></u>	
Item	support	Item	support
(O, K, E)	3	(O, K, E)	3

24/7/23

\*Consider the minimum confidence as 50% and compute the association rules. (Example-1)

let us suppose the frequent item set as  $L_1$  where  $L_1$  can be written as  $(S; L)$  i.e.

$S$  is the starting element of  $L_1$  and  $L$  is list of items such that  $L = L_1 - S$ . we can frame Association rules for  $L_1$  such that  $S \rightarrow (L - S)$  or any other combination which satisfies the minimum confidence.

$$L_1 (I_1, I_2, I_5)$$

$$\begin{aligned} I_1 &\rightarrow (I_2, I_5) & (I_1, I_2) &\rightarrow I_5 \\ I_2 &\rightarrow (I_1, I_5) & (I_1, I_5) &\rightarrow I_2 \\ I_5 &\rightarrow (I_1, I_2) & (I_2, I_5) &\rightarrow I_1 \end{aligned}$$



$$\text{confidence}(I_1 \rightarrow \{I_2, I_3\}) = \frac{\text{count}(I_1 \cup I_2 \cup I_3)}{\text{count}(I_1)}$$

$$\Rightarrow \frac{2}{6} = \frac{1}{3} \times 100 = 33.66\% \quad \times$$

$$\text{confidence}(I_2 \rightarrow \{I_1, I_3\}) = \frac{\text{count}(I_2, I_1, I_3)}{\text{count}(I_2)}$$

$$= \frac{2}{7} \times 100 = 28\% \quad \times$$

$$\text{confidence}(I_3 \rightarrow \{I_1, I_2\}) = \frac{\text{count}(I_3, I_1, I_2)}{\text{count}(I_3)}$$

$$\Rightarrow \frac{2}{2} = 1 = 100\% \quad \checkmark$$

$$\text{confidence}(\{I_1, I_2\} \rightarrow I_3) = \frac{\text{count}(I_1, I_2, I_3)}{\text{count}(I_1, I_2)} = \frac{2}{4} = \frac{1}{2} = 50\%$$

$$\text{confidence}(\{I_1, I_3\} \rightarrow I_2) = \frac{\text{count}(I_1, I_3, I_2)}{\text{count}(I_1, I_3)} = \frac{2}{2} = 1 = 100\% \quad \checkmark$$

$$\text{confidence}(\{I_2, I_3\} \rightarrow I_1) = \frac{\text{count}(I_2, I_3, I_1)}{\text{count}(I_2, I_3)} = \frac{2}{2} = 1 = 100\% \quad \checkmark$$

$L_2(I_1, I_2, I_3)$

$$I_1 \rightarrow (I_2, I_3) \quad \text{confidence} = \frac{2}{6} \times 100 = 33.66\% \quad \times$$

$$I_2 \rightarrow (I_1, I_3) \quad \text{confidence} = \frac{2}{7} \times 100 = 28\% \quad \times$$

$$I_3 \rightarrow (I_1, I_2) \quad \text{confidence} = \frac{2}{6} \times 100 = 33.66\% \quad \times$$

$$(I_2, I_3) \rightarrow I_1 \quad \text{confidence} = \frac{2}{4} \times 100 = 50\% \quad \times$$

$$(I_1, I_3) \rightarrow I_2 \quad \text{confidence} = \frac{2}{4} \times 100 = 50\% \quad \times$$

$$(I_1, I_2) \rightarrow I_3 \quad \text{confidence} = \frac{2}{4} \times 100 = 50\% \quad \times$$

$$L_1 \rightarrow (I_3 \rightarrow \{I_1, I_2\}), (\{I_1, I_2\} \rightarrow I_3), (\{I_2, I_3\} \rightarrow I_1)$$

IID	<u><math>I_1</math></u> Rice	<u><math>I_2</math></u> pulses	<u><math>I_3</math></u> Milk	<u><math>I_4</math></u> Apple	<u><math>I_5</math></u> oil
t1	1	1	1	0	0
t2	0	1	1	1	0
t3	0	0	0	1	1
t4	1	1	0	1	0
t5	1	1	1	0	1
t6	1	1	1	1	1

Support = 50%

Confidence = 60%

$$\frac{50}{26} = \frac{100}{52}$$

$C_2 = L_1 \times L_1$

Item	support
$I_1$	4 = 67%
$I_2$	5 = 77%
$I_3$	4 = 67%
$I_4$	4 = 67%
$I_5$	3 = 50%

Item	support
$I_1$	4
$I_2$	5
$I_3$	4
$I_4$	4
$I_5$	3

Item	support	
$(I_1, I_2)$	4	67% $(I_3, I_4)$ 2 33%
$(I_1, I_3)$	3	50% $(I_3, I_5)$ 2 33%
$(I_1, I_4)$	2	33% $(I_4, I_5)$ 2 33%
$(I_1, I_5)$	2	33%
$(I_2, I_3)$	4	67%
$(I_2, I_4)$	3	50%
$(I_2, I_5)$	2	33%



L<sub>2</sub>

Item support	
(I <sub>1</sub> , I <sub>2</sub> )	4
(I <sub>1</sub> , I <sub>3</sub> )	3
(I <sub>2</sub> , I <sub>3</sub> )	4
(I <sub>2</sub> , I <sub>4</sub> )	3

C<sub>3</sub> = L<sub>2</sub> ∪ L<sub>2</sub>

Item support	
(I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> )	3 = 50%
X (I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> )	

↓  
Subsets are not frequent

L<sub>3</sub>

Item support	
(I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> )	3

(I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>)

$$I_1 \rightarrow \{I_2, I_3\} = \frac{\text{Count}(I_1 \cup I_2 \cup I_3)}{\text{Count}(I_1)} = \frac{3}{4} \times 100 = 75\% \checkmark$$

$$I_2 \rightarrow \{I_1, I_3\} = \frac{3}{5} \times 100 = 60\% \checkmark$$

$$I_3 \rightarrow \{I_1, I_2\} = \frac{3}{4} \times 100 = 75\% \checkmark$$

$$\{I_2, I_3\} \rightarrow I_1 = \frac{3}{4} \times 100 = 75\% \checkmark$$

$$\{I_1, I_3\} \rightarrow I_2 = \frac{3}{3} \times 100 = 100\% \checkmark$$

$$\{I_1, I_2\} \rightarrow I_3 = \frac{3}{4} \times 100 = 75\% \checkmark$$

### 27/7/23 FP-Growth Algorithm:- (frequency-pattern)

Step-1:- Condense the dataset by constructing a FP-tree

Step-2:- To extract association rules from the FP-Tree

from example 1; Support = 2

descending order of support count

L<sub>1</sub> = L

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

⇒

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

⇒

T <sub>100</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>5</sub> ✓
T <sub>200</sub>	I <sub>2</sub> , I <sub>4</sub> ✓
T <sub>300</sub>	I <sub>2</sub> , I <sub>3</sub> ✓
T <sub>400</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>4</sub> ✓
T <sub>500</sub>	I <sub>1</sub> , I <sub>3</sub> ✓
T <sub>600</sub>	I <sub>2</sub> , I <sub>3</sub> ✓
T <sub>700</sub>	I <sub>1</sub> , I <sub>3</sub> ✓
T <sub>800</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub> ✓
T <sub>900</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>4</sub>

FP Tree

