

## Association Analysis (Pattern Mining)

2, 16, 3, 5, 20, 4, 3, 12, 3, 5, 20, 21, 19, 18, 9, 21, 19, 18, 23, 36, 21, 19, 18

Another stream of data mining. If we have the above data. We can find some patterns in data (i.e. values that are repeated).

- 3, 5 are repeated

- 21, 19, 18

Finding pattern is imp. for diff. app. If, we want to design a mathematical func. to better represent the data like:

$$f(x) = \begin{cases} x_1 + 2 \\ \text{or} \\ x_1 - 2 \end{cases}$$

These can be used to write algorithms better

So for our T/P we can rep. the patterns.

e.g. In E-commerce website if we buy a shirt then we get recommended the pants by looking at the purchase history pattern of the user. Similarly in scientific analysis (e.g. DNA analysis to prevent diseases), identify similar web doc. for better classification, providing tags to the email for similarity, etc.

If dataset is very large, pattern finding needs some good algo. not manually. Algo. need to efficiently find pattern in data.

### Terminologies

↳

Item: each thing bought by user in a particular transaction

Itemset: set of items

set:

each item is a 1-itemset

→ 1-itemset → {Beer} has only item {nuts}, {milk}

→ 2-itemset → itemset with 2 items e.g. {Beer, nuts}, {Beer, Diaper}

→ k-itemset → itemset with k items

Support Count: sc of an itemset is its freq. of occurrence / appearance in a dataset  
e.g. If we want to find sc of {Beer, Diaper} then it is 3

to judge whether an association rule is good -- we use support and confidence....

freq

This means that itemset occurs more frequently as compared to other items. Customers are more likely to buy Beer and Diaper. So, it is freq pattern. freq patterns are rep. by association rules like Beer  $\rightarrow$  Diaper. So, many pattern and hence association rules exist e.g. Diaper  $\rightarrow$  Beer.

pattern  
finding  
 $\downarrow$   
find useful  
patterns

How to judge which association rule is good? i.e. which patterns are more useful/interesting?

There are 2 measures for it, we compute for the  $X \rightarrow Y$  rule:

support      confidence

If we have an ass. rule  $X \rightarrow Y$

$$\text{support}(X \rightarrow Y) = \frac{\text{#}(X \cup Y)}{N}$$

support count

No. of trans.

out of total transactions how many contains Beer and Diaper

$$= \frac{3}{5} = 60\%$$

$$\text{confidence}(X \rightarrow Y) = \frac{\text{#}(X \cap Y)}{\text{#}(X)}$$

where, X is present out of that Y is also present

$$= \frac{3}{3} = 100\%$$

Clearly

support means e.g. 60% of students are agreed / liking something, tot of entries are liking Beer with Diaper. i.e. the pattern is more likely. & support measures pattern is liked by many customers.

reliability

confidence means e.g. out of 60% are all agreed on same constraint. If 2 vars are there X and Y,  $x = y$ , how much they are associated i.e. has numerical quantification of the (pattern) association.

For any pattern even if both are good then it means it is always purchased & purchased together.

0	0.25	0.5	0.75	1	$\rightarrow$	$0 \leq c \leq 1$
none	low	medium	partial	high		

An itemset is said to be frequent if its support is greater than the min. support threshold...

An itemset is said to be frequent if its support is greater than the min. support threshold of the application.

e.g. if it is 3 and threshold is 2 then  $\{ \text{Beer}, \text{Diaper} \}$  is frequent itemset

If we have an itemset  $\{ a_1, \dots, a_{100} \}$

If in prev. we convert it to standard dataset

	Bear	Nut	Diaper	Latte	Fqgs	Milk
1	1	1	1	0	0	0
2	2	0	1	1	0	0
3	1	0	1	0	1	0
4	0	1	0	0	1	1
5	0	1	2	1	1	1

From trans. dataset the items are represented as feature

In a general store tree standard dataset will have a no. of features. **More length of itemset depends on the no. of items in store.**

In above max length is 6 i.e.  $a_1, \dots, a_6$ .

If itemset is very long  $\{ a_1, \dots, a_{100} \}$  then using it we can form different diff. comb. i.e. sub-itemset e.g.  $\{ a_1 \}$ ,  $\{ a_1, a_2 \}$  and so on. In above,  $2^{100}-1$ .

In above,  $2^{100}-1$

Out of all, we calc. support & confidence for all and see which is good. This is not practical. We need an algo. to filter out the itemsets.

**Closed itemset:** itemset  $X$  is closed if  $X$  is frequent and there exists no (subset) superset  $Y$  of itemset with the same support as  $X$ .

e.g. if  $\{ \text{Beer}, \text{Diaper} \}$  is closed or not?

→ Frequent? If it is higher than min. support threshold  
 $3 > 2$ . So, it is frequent

Now, a lot of supersets can be formed but none has the same support value as  $\{ \text{Beer}, \text{Diaper} \}$

Itemset  $X$  is max-pattern if  $X$  is freq. and there exists no frequent super-pattern.

Y of X. i.e. X is the only frequent itemset so max is possible.

Support count is maximum for X. e.g. if we find a superset,  $\{ \text{Beer}, \text{Diaper}, \text{Nuts} \}$  then

**Closed itemset -- is frequent and no superset has same support count...**  
**MAX PATTERN ....**

its support count is only 1 which is not  $\geq 2$ . Similarly, for othersupersets we consider  $\{B, D, E, G\}$ ,  $\{B, D, C, F, G\}$ . So,  $\{B, D, G\}$  is both closed and max.

if we have 2 itemsets  $\rightarrow \langle a_1, \dots, a_{100} \rangle$   
 $\rightarrow \langle a_1, \dots, a_{50} \rangle$

Threshold = 1

$$SC(\{a_1, \dots, a_{100}\}) = 1$$

$$SC(\{a_1, \dots, a_{50}\}) = 2$$

Both are frequent. No superset of  $\langle a_1, \dots, a_{100} \rangle$  exists. So, it is closed.

$\langle a_1, \dots, a_{50} \rangle$  has a superset but no superset's support count isn't 2. So, it is not closed.

Only,  $\langle a_1, \dots, a_{100} \rangle$  is max-pattern.

For  $\langle a_1, \dots, a_{50} \rangle$  its superset is also frequent. So, it isn't max-pattern.

no need  
to scan  
DB again  
↓  
complexity  
reduced

Example

set of all patterns  $\rightarrow$  all possible subset

Problem  $\rightarrow$  No of patterns

We need good algo. to id. the pattern & check if more good or not

We have  
→ Apriori  
→ FP-growth  
→ ECLAT

Downward  
closure

"Any subset of a freq. itemset must be frequent". So, if we find a freq. itemset then no need to check its subset.

Freq.  $\rightarrow$  appears a lot  $\rightarrow$  so obviously

e.g., Beer or Diaper will be present in atleast 3 trans. So we don't need to calc the supp & confidence of those

Apriori Algorithm

frequent in pattern mining

If there is any itemset that is infrequent then its superset shouldn't be taken as frequent because it will also be infrequent. So, no need to check

Upward  
property  
(Apriori  
principle)

on check  
steps

- Given the whole DB to get frequent 1-itemset  
Using (Generate) those, generate list of 2-itemset, using 2 generate 3  
candidate itemset  
↓  
complexity: From candidate eliminate the infrequent  
reduced: Repeat until ...



Example  $C_1 \rightarrow$  candidates for 1-itemset given, minsup = 2

$\{A\}$	2	→	$L_1$	$\{A\}$	frequent 1-itemset
$\{B\}$	3		$\{B\}$		
$\{C\}$	3		$\{C\}$		
$\{D\}$	1		$\{D\}$		
$\{E\}$	3		$\{E\}$		

1st scan

From  $L_1$ , prepare (2 i.e.  $(n+1)$ ) candidate itemset

$\{A, B\}$	1	X	→	$L_2$	$\{A, C\}$
$\{A, C\}$	2			$\{B, C\}$	
$\{A, E\}$	1	X		$\{B, E\}$	
$\{B, C\}$	2			$\{B, E\}$	
$\{B, E\}$	3			$\{C, E\}$	
$\{C, E\}$	2				

2nd scan

Now, 3-itemset  $C_3$

A, C, B X

A, B, E X since  $\{A, B\}$  is in freq. 1-itemset is also infrequent

A, C, E X

B, C, E

Then, length 4 no change. So, we stop

3/04/2025

Q. Solve the following.

Tid	Items
T1	I1, I2, I5
T2	I2, I4
T3	I2, I3
T4	I, 2, 4
T5	I, 3
T6	I, 3
T7	I, 3
T8	I, 2, 3, 5
T9	I, 2, 3

### #Growth (Frequent Pattern Growth)

Before this,

in above if we apply Apriori algorithm:

(candidate  
1-itemset)

itemset	support
{I1}	6
{I2}	7
{I3}	5
{I4}	2
{I5}	2

support, min support count = 2

1) O. cuts  
beg.  
for 3侯选  
men stop  
after L3

all the candidate itemset are frequent. L1 is same as above

$L_2 = L_1 \times L_1$	itemset	support	itemset	support	
	{I1, 2}	4		{I4, 5}	0
	{I1, 3}	4			X
	{I1, 4}	1			X
	{I1, 5}	2			
	{I2, 3}	4			
	{I2, 4}	2			
	{I2, 5}	2			
	{I3, 4}	0			X
	{I3, 5}	1			X

if while gen. 3 itemset we get 4. itemset we ignore

else continue

Done

Page

freq itemset  
 $\{1, 2\}$        $\{3, 3\}$   
 $\{1, 3\}$        $\{3, 4\}$   
 $\{1, 5\}$        $\{2, 5\}$

$(l_1 \cup l_2) \cup l_3$        $\{1, 2, 3\}$       2  
 $\{1, 2, 5\}$       2  
 $\{1, 2, 4\}$       need not be tested as  $\{3, 4\}$  is infrequent  
 $\{1, 3, 5\}$        $\{3, 5\}$

$l_3$  is some

$(l_1 \cup l_2) \cup l_3$        $\{1, 2, 3, 5\}$       1  $\times$   
it is infrequent

Now we stop

↳ if  $l_3$  is same to  $l_2$ , we stop  
↳ if we get an empty list, we stop

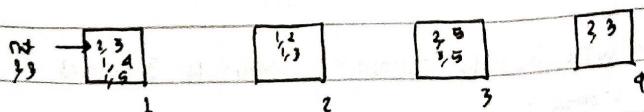
If 0. exists  
freq.  
for 3 itemset  
then stop

Problem:

This algo. is quite complex. e.g. while computing  $l_3$  from  $l_2$  we need to find support values for which we need to check the database i.e. across transactions (length) and within transactions (width). if a trans. has 3 items then  $2^{3-1}$  itemset can be generated. Attempts are made to reduce this complexity.

Hash-based (Support Count) Mechanism

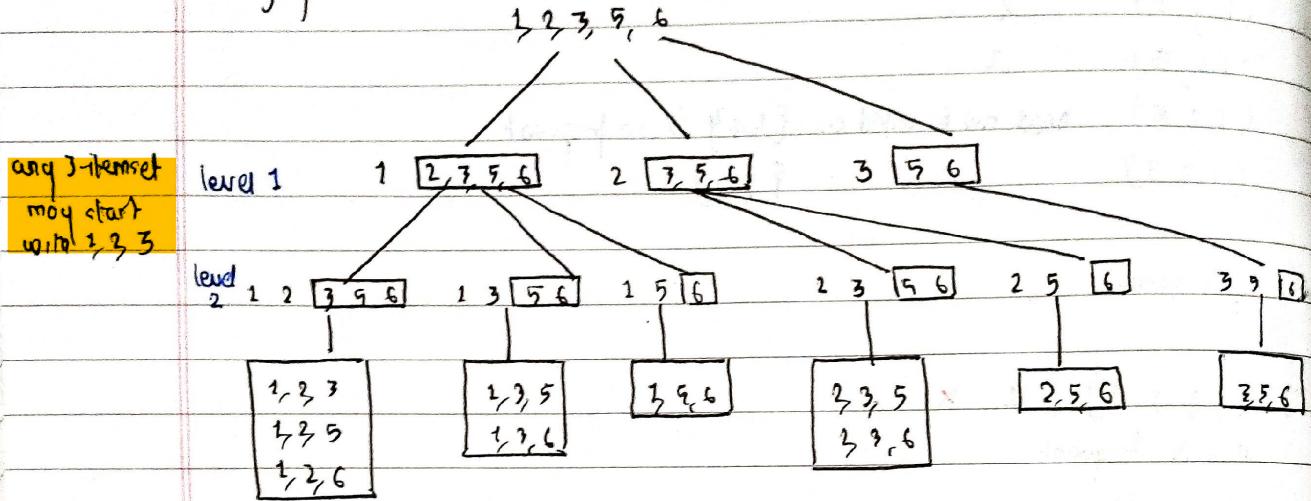
Hash Tree for Support (Count). To reduce complexity, we use a hash func. & diff. hash buckets



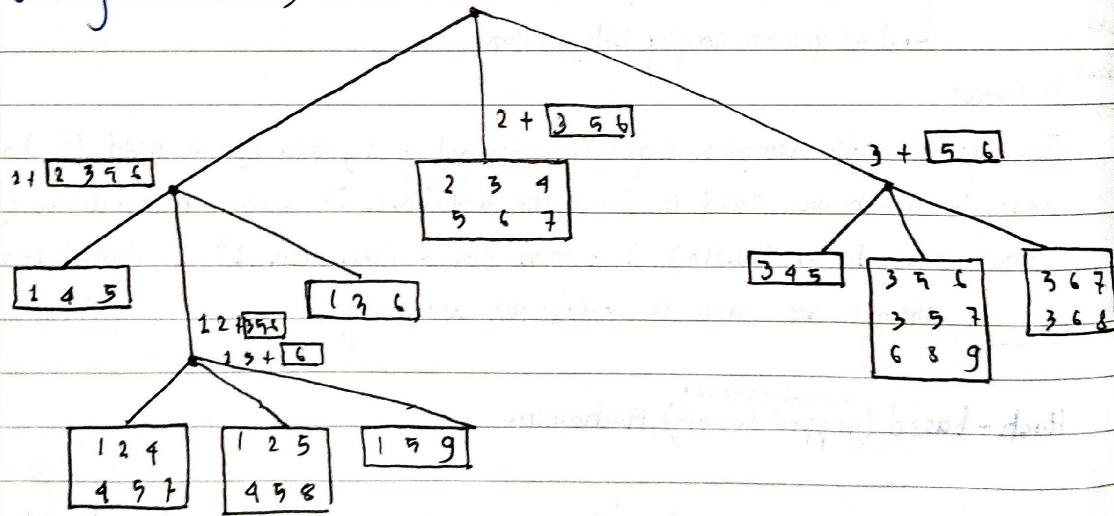
Candidate itemsets generated in each step are placed in buckets. Whatever itemsets belonging to  $T_1$  i.e.  $\{1, 2\}$ ,  $\{1, 5\}$ ,  $\{2, 5\}$ . Previously, a particular transaction was compared against the width of transaction as well.

We hash  $\{1, 2\}$  to find its bucket and we need to only compare with itemsets

of that bucket only to find the support counts.  
 let us say we have in a transaction  $t$  with items  $3, 3, 3, 6$ . We arrange items in lexicographical order. Now we need to check for candidate 3-itemset. we find all 3-itemsets belonging to the transaction.



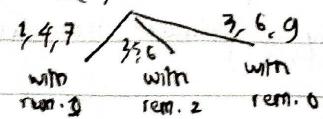
Creating a hash tree,



let us say we have a hash tree with buckets containing candidate 3-itemset

lets say hash func. is  $H(P) = P \bmod 3$

$H(P) \rightarrow P \bmod 3$  where  $P$ , may be 1st item of itemset for 1st level

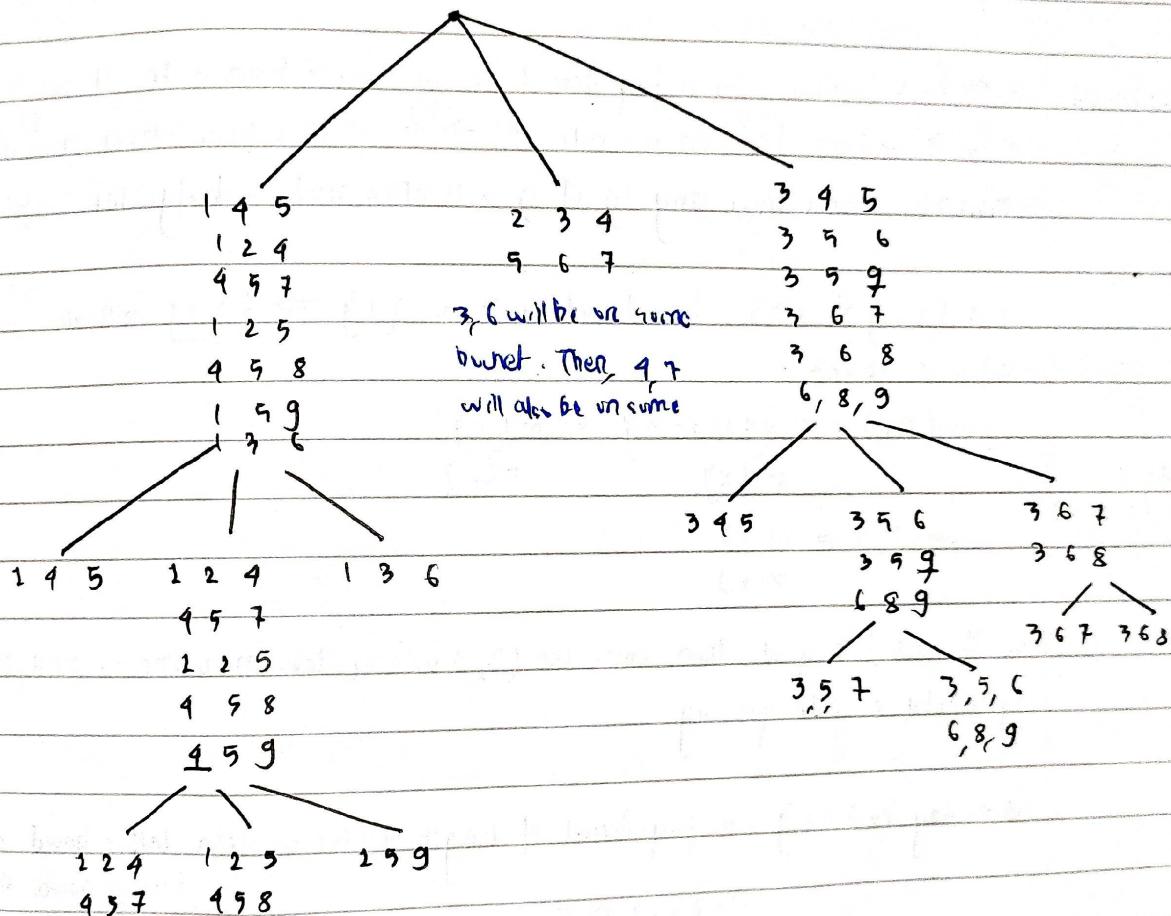


If we hash transactions 1, 2, 3, 4, 5, 6 then more starting with 1 are on left, 2 middle, 3 on (left) right. This is for 1st level.

For 2nd level, 2nd item may be 2, 3, 5. If it is 2 then at middle, 5 men also middle, if 3 then on right. So, if our 2nd element is  $\overset{2 \text{ or } 5}{\text{in}}$  our trans. we need to check the middle. Similarly for others.

Q. To find the support count only few of the itemsets gets compared not all. If  $\{3, 4, 5\}$ ,  $\{3, 2, 4\}$ , ... is our candidate 3-itemset.

- \* 1st hash on 1st item, then 2nd item, then 3rd item



In this way, we form the Hashtree.

### Confidence Based Pruning

If we have an itemset  $y$  of length ' $k$ ' then no. of (associated) association rules (then) is  $2^k - 2$ .

E.g. if our itemset is  $\overset{\text{freq}}{\{1, 2, 3\}}$ . We can write the ass. rules:

$$\{1\} \Rightarrow \{1, 2\} \quad \{1, 2\} \Rightarrow \{2\} \quad \{1, 2\} \Rightarrow \{1, 3\} \quad \{1, 3\} \Rightarrow \{1\}$$

$$\{2\} \Rightarrow \{1, 2\} \quad \{1, 2\} \Rightarrow \{2\} \quad \{2, 3\} \Rightarrow \{2\}$$

not considered

itemset is freq.  
if it satisfies  
support (only  
needs support)

We want to find the useful (itemsets) that support the support & also confidence criteria. If an itemset is freq. then all the ass. rules will satisfy the support rule i.e.  $\sigma(x \cup y)$ . But, we also need it to support confidence. For each rule the

confidence would be different i.e.  $\sigma(x \cup y)$  confidence based ass. because  $\frac{\sigma(y)}{\sigma(x)}$  is useful while doing whole db.

The apriori algo. we must have already found the necessary support values. So, for confidence we don't need to rescan the DB.

Statement:  $x \Rightarrow^i y - x$ , where  $y$  is a freq. itemset, is an association rule. If we know an  $x' \subset x$  and we find an ass. rule  $x' \Rightarrow^{ii} y - x'$ ; Then if we know i) satisfies the confidence criteria then any ii) of i) will also not satisfy the confidence criteria.

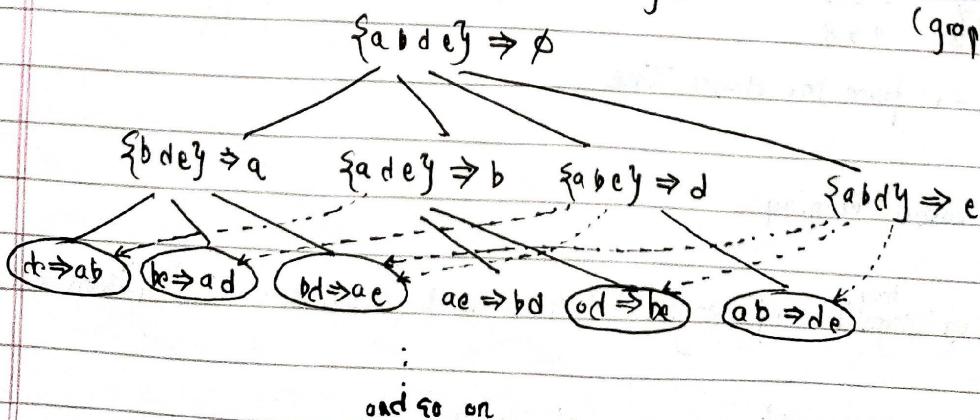
e.g. if  $\{1, 2, y\} \Rightarrow \{3, y\}$  doesn't satisfy then  $\{1, y\} \Rightarrow \{2, 3, y\}$  will also not satisfy the confidence criteria.

$$\text{con}(R_2) = \frac{\sigma(x \cup y - x)}{\sigma(xy)} = \frac{\sigma(y)}{\sigma(x)}$$

$$\text{con}(R_1) = \frac{\sigma(y)}{\sigma(x')}$$

so, if  $\sigma(x')$  is greater then sub.  $\text{con}(R_2)$  will be less than the  $\text{con}(R_1)$ . This is very helpful in rule proving.

Let us say  $\{a, b, c, d, e\}$  is a freq. itemset of length 5. We create a lattice based structure (graph based structure).



1st step:

2nd step:

Once we get a frequent itemset, then we can use the lattice method to find out all the association rules... if the node does not satisfy the confidence then all the rules rooted on it. (in its subtree) will not satisfy so we eliminate them..

This is called lattice based structure for freq. itemset. If we find freq. 4-itemset using apriori, then our application may need the rules to support support, confidence, correlation acc. to its need.

This shows the possible ass. rles corresponding to the freq. itemset. With this, all operations of graph are applicable. Suppose we write any rule  $\{b\} \rightarrow \{a\}$ . If this doesn't satisfy the confidence then all the rules rooted on it don't satisfy the confidence so we eliminate.

We check 1st level, if all rules satisfy confidence then we go to next level. Then if any rule doesn't satisfy (them) then the sub-tree can be removed.

If we know support values of each, we can also find closed & max itemset. If we say  $e \Rightarrow abd$  then we need to check the immediate supersets (parent). If they don't have same support value then  $e \Rightarrow abd$  isn't a closed. Similarly, for more pattern too.

4/4/2025

## FP Growth (Frequent Pattern)

Same Q as before

another method to find freq. itemset. Trying to improve the apriori algorithm.

1st step: 1st we scan our dataset, <sup>find 1-itemset & their SC</sup> same as apriori algorithm.

Support Count

$\{1\}$	6
$\{2\}$	7
$\{3\}$	6
$\{4\}$	2
$\{5\}$	2

2nd step: Sort items acc. to SC in decreasing order

$\{2\}$	7
$\{1\}$	6

same from any order

$$1, 2, 3, 5 \rightarrow \underbrace{2, 1, 3, 5}_{\text{order}}$$

$\{3\}$  6  
 $\{2\}$  2  
 $\{5\}$  2

called L-list

3rd step: Construct an FP tree

Root is considered as NULL from the dataset & add branches to the tree as per the above order. 1st trans. is  $\{1, 2, 5\}$  so branch is  $T_2, T_1, T_5$  acc. to the L-list. Separate tree, then 2nd and so on acc. to above order.

72 is 1st  
on it is  
the bottom  
of L-list

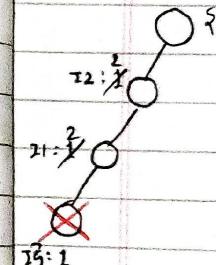
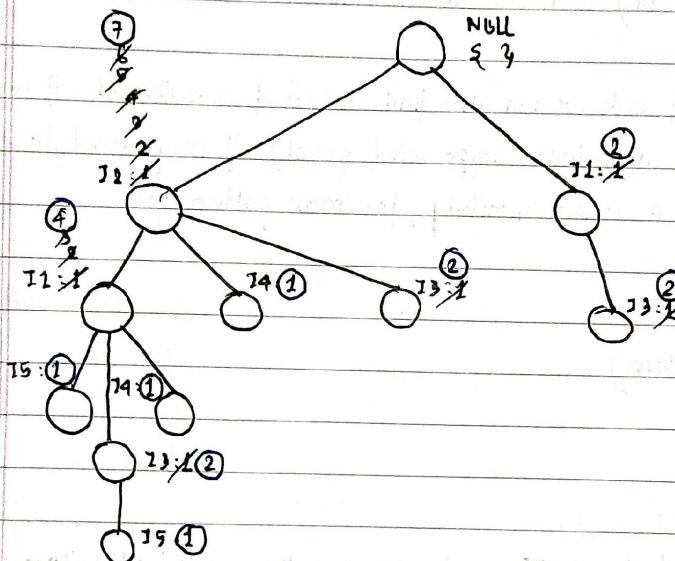
with head.  
pattern tree  
we construct  
a word.  
FP-tree

11  
IS

74

73

For 1st item  
support is 1



it doesn't  
support the  
min support  
count of  
2

71

4th step: Create a table as follows:

Item Id	Support Count	Node with least support
T_2	7	
T_1	6	
T_3	6	
T_4	2	
T_5	2	

72: 4

i) and ii)  
This de

Strong

5th step: Create the following:

Approach is to partition the dataset and find out corresponding to each partition. (Computing acc. to L-list order)  
 To pattern we want the cond. FP-tree. Reduction of complexity

N = 10,  
Computer  
video

ways to reach the item, a kind of partition of the dataset

Item  
75 is 1st  
on list  
use bottom  
of L-list

Per.  
with 1nd.  
pattern tree  
we construct  
a word.

FP-tree

Conditional Pattern Tree  
 $\{72, 71; 1\}$  if 25 un  
 $\{73, 71, 73; 1\}$

conditional FP-tree  
 $\{72: 2, 71: 2\}$

Using mine we gen.  
freq. pattern by simply  
insert 75 with two  
items

Freq. Pattern Generator

$\{72, 75: 2\}$

$\{71, 75: 2\}$

i)  $\{73, 71, 75: 2\} \rightarrow$  os in  
some branch

74

$\{72, 71, (74): 1\}$   
 $\{72: 1\}$

$\{72: 2\}$

$\{73, 74: 2\}$

73

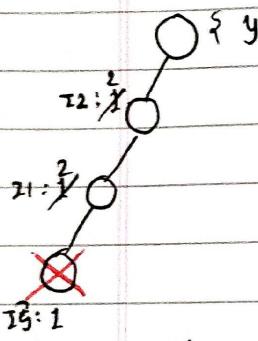
$\{72, 71: 2\}$   
 $\{72: 2\}$   
 $\{71: 2\}$

$\{72: 4, 71: 3, 71: 2\}$

$\{72, 75: 4\}$

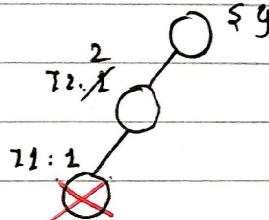
$\{71, 73: 1\}$

ii)  $\{72, 71, 73: 2\}$



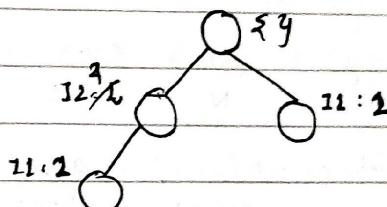
73: 1

exit doesn't  
support the  
min support.  
Count of  
2



71

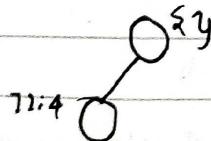
$\{72: 4\}$   
 $\{73: 1\}$



$\{72: 4\}$

$\{72, 71: 4\}$

Here 1st take  
2 as we have  
 $\{72, 71: 2\}$



i) and ii) are freq. 3-itemset. We also have  $\frac{72}{2}$  items

This doesn't generate freq. test. Patterns are freq. grown based on the position of data.

### Strong Association Rule

N: 10,000

Computer Games = 6000

Video = 7500

once we identify the frequent itemset, we can write all possible association rules ... by using lattice based structure and prune to identify the rules satisfying the both support and confidence criteria...

it might be the case that rule is satisfying both the criterias but it is still not strong association.....

Computer Games & Video = 4000

Once we id-the freq-itemset, we can write all possible ass.rules wh by using the lattice based structure & prune to id-the rules satisfying the both support and confidence criteria.

If we have a dataset with 10,000 trans.. 6000 customers buy Game, 4000 customers purchase from computer games and video.

If we write an ass.rule  $X \rightarrow Y$  like,

Computer Games  $\rightarrow$  Video

i.e. a customer buy (g) also buys Video,  
Suppose, min-support = 40% } for weak rules  
min-confidence = 66%.

For above of case find,

$$\text{Support} = \frac{\sigma(X \cup Y)}{N} = \frac{4000}{10000} = 40\%$$

$$\text{confidence} = \frac{\sigma(X \cup Y)}{\sigma(X)} = \frac{4000}{6000} = 66\%$$

✓  $\sigma(\text{Computer Games})$

So, the above is a useful rule or it satisfies both criteria. But it isn't a strong association.  
As, we find probability of purchase of video is  $\frac{7500}{10000} = 75\%$ . But if video is bought

along with (Video) games then the confidence is 66%. So, the prob. decreases. This is known as -ve correlation.

↳ 2 items are dependent or not of how much

Using this we can find better rules. Using the correlation we can ↑ sales. -ve correlation then the sales decreases.

So, to judge the usefulness of rule we have correlation as a criteria.

Correlation

Methods to compute:

i) Lift

$$\text{lift}(A, B) = \frac{P(A \cup B)}{P(A)P(B)}$$

intersection....

i) If 2 events are independent then,  $P(A \cup B) = P(A) \cdot P(B)$   
 $A, B$  are two itemsets.

- lift may be ( $= 0$ ) ( $\geq 1$ ) ( $> 1$ )
- if they are independent ( $\approx 1$ )
- $= 1$  -ve correlation
- $(\geq 1)$  true correlation

$$\text{In prev., lift}(G, V) = \frac{0.4}{0.6 \times 0.75} = 0.88$$

i.e. it is  $< 1$ , so we have a -ve correlation. This much we need to lift to get the good RSS.

### ii) Using Chi-Squared Test

$\chi^2$ : correlation

We prepare a contingency table

		$\rightarrow$ Games / Games		
		4000	(7500) 3500	7500
Video	Video	4000	(7500)	7500
	$\rightarrow$ Video	(6000)	900	2900
		6000	4000	10000

These are observed values we can compute free expected value

$$\chi^2 = \sum_{\text{Expected}} (\text{Observed} - \text{Expected})^2$$

$$\text{Games} \times \text{Video} = \frac{6000 \times 7500}{10000} = 4500 \text{ Similarly for all}$$

Consider to be independent, how if hypothesis is rejected they are dependent. In above, the hypo. is rejected and they are corr. The value we get from above is the amount of correlation.