

Association Rules & Item Sets

Frequent pattern :-

Eg : (itemsets, subsequences or substructures) patterns appear frequently in a dataset

Ex

Set of Items such as milk and bread

appear frequently in a data set. (together)

Freq. Itemset = milk & bread

Subsequence :-

Buying first pc, then .

digital camera then

memory card

} occurs frequently
in shopping history
DB.

↓
called (frequent)
sequential pattern.

Substructure :

different structural forms such as subgraphs - subtrees or sublattices which may be combined with itemsets or subsequences.

→ if substructure occurred frequently called (frequent) structured pattern.

MBA (Market Basket Analysis algm)

- Analyze customer buying habit
- discover correlation relationship between business txn. records
 - Helps for Business decision making.
- Rep. frequently associated or purchased together
 - It's rep'td in the form of Association Rule.

Ex Customer

↳ buys Computer

↳ also tend to buy Antivirus at the same time

$\boxed{\text{Computer} \rightarrow \text{Antivirus - sw}}$

Support = 2% (Computer & Antivirus sw purchased together out of all txns of the entire shop)
Confidence = 60% (60% of customer purchased computer also bought the sw).

→ It needs to satisfy minimum support threshold & minimum confidence threshold

→ set by user / domain experts

Frequent Itemsets, Closed Itemsets & Association Rules

Freq. Itemset

Itemset whose support is greater than user specified minimum support.

Closed freq. Itemset

If none of its immediate supersets has the same support as that of the itemset.

Maximal freq. Itemset

If none of its immediate supersets is frequent.

Downward closure property of frequent patterns

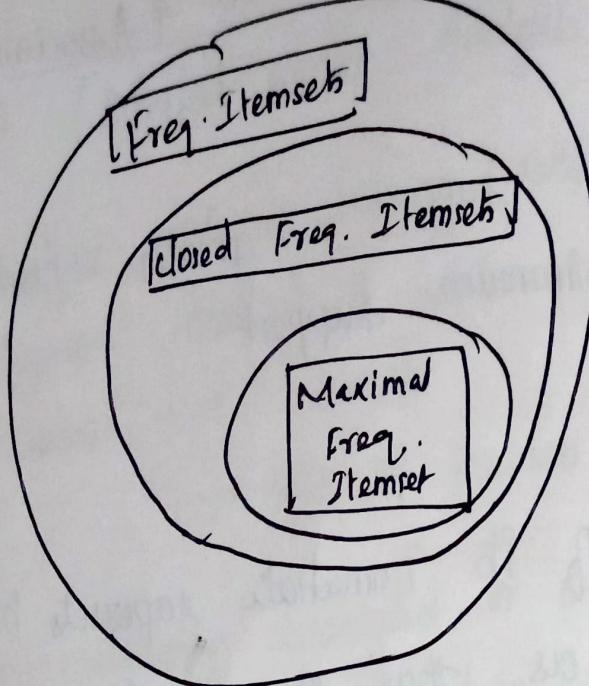
All subset of any frequent itemset must also be freq.

Ex

If Milk, Bread, Butter is a freq. Itemset then the following items are freq.

Milk
Bread
Butter
Milk, Bread
Milk, Butter
Bread, Butter

} smaller subset of freq. pattern Itemset



Need of closed & Maximal Itemset

It's useful when huge amount of data and in association rule mining

If length of freq. itemset is k then.

by downward closure property

→ All of its 2^k subsets are also freq.

Disadv of Maximal freq Itemsets

Even all its subsets are frequent,
we don't know their support

For mining rules, support info is very less

→ Therefore closed freq itemset is preferred

TID	List of Items
T ₁	A, B, C, D
T ₂	A, B, C, D
T ₃	A, B, C
T ₄	B, C, D
T ₅	C, D

Min sup-count = 3

Find freq, closed, Maximal itemset

Item	Count
A	3
B	4
C	5
D	4

All items A, B, C & D are frequent because their sup count is greater than or equal to min.sup count

↳ sup count \geq min sup count

2 Itemset as combination

Item	Count
A, B	3
A, C	3
A, D	2
B, C	4
B, D	3
C, D	4

All 2 items are freq.
except AD

\therefore min sup count = 3

$A(3) \rightarrow AB(3), AC(3), AD(2)$

↳ sup. count Immediate superset

A (count) is not $>$ its immediate
condn not satisfied.

A is not closed

superset \rightarrow
closed

In A 's immediate superset,

Itemset are present with min sup. count 3

A is not maximal

If immediate superset is not frequent (not satisfy min support count) \rightarrow called maximal

$B(4) \rightarrow AB(3), BC(4), BD(3)$

B (count) is not greater than its immediate superset

B is not closed

B 's immediate superset, Itemset are present with min sup. count 3

B is not maximal.

C(5) \rightarrow AC(3), BC(4), CD(2)

C(count) > immed. superset(count)

• C is closed

C Not maximal

D(H) \rightarrow AD(2), BD(3), CD(2)

D(count) > im. superset(count)

D is not closed

Im. superset count > min sup count (3)

\rightarrow No

\rightarrow D is not maximal

3 Itemset Generation

Item	count
A B C	3
A B D	2
A C D	2
B C D	3

AB(3) \rightarrow ABC(3), ABD(2), ACD(2), BCD(3)

AB(count) > IS(count)

AB is not closed

IS(count) > min sup. count(3)

No \rightarrow AB not maximal

$AC(3) \rightarrow ABC(3), ACD(2)$

- ↳ Not closed
- ↳ not maximal

$AD(2) \rightarrow ABD(2), ADC(2)$

- ↳ not frequent \therefore if $< \text{sup count}^{\min}$

$BC(4) \rightarrow ABC(3), BCD(3)$

- ↳ closed
- ↳ not maximal

$BD(3) \rightarrow ABD(2), BCD(3)$

- ↳ not closed
- ↳ not maximal

$CD(4) \rightarrow ACD(2), BCD(3)$.

↳ closed

↳ Not maximal

H Itemset Generation	
item	Count
A, B, C, D	2

$ABC(3) \rightarrow ABCD(2)$

↳ closed

↳ itemset not present with min sup count (3)

↳ Maximal

$ABD(2) \rightarrow$ Not freq.

$ACD(2) \rightarrow$ Not freq.

$\overline{BCD}(3) \rightarrow ABCD(2)$

↳ closed

↳ Maximal.

Assignment

TID	List of items
T1	B, J, P
T2	B, P
T3	B, M, P
T4	E, B
T5	E, M

Min sup. count - 40% $\Rightarrow 2/11$

Find frequent, closed, max.
Itemset

All freq except J

Item	count
B	4
E	2
J	1
M	2
P	3

12

$B(4) \rightarrow BE(1), BM(1), BP(3) \rightarrow$ closed, Not max

$E(2) \rightarrow BE(1), EM(1), EP(0) \rightarrow$ closed, max

$M(2) \rightarrow BM(1), EN(1), MP(1) \rightarrow$ closed, max

$P(3) \rightarrow BP(3), EP(0), MP(1) \rightarrow$ not closed, not max.

3 itemset empty so BP is closed & maximal.

sup & conf b/w 0% and 100% rather than 0 to 1.0.

Itemset set of items

$$\text{sup } (A \Rightarrow B) = P(A \cup B)$$

$$\begin{aligned} \text{Confidence } (A \Rightarrow B) &= P(B|A) \\ &= \frac{\text{sup } (A \cup B)}{\text{sup } (A)} \end{aligned}$$

Freq. Itemset

Itemset occurs frequently $\Rightarrow \geq \text{minsup. count}$

Generate strong association rules from freq. itemset

rules must satisfy min support and min. confidence.

Ex

freq itemset of length 100

$\{a_1, a_2, \dots, a_{100}\}$ contain $\binom{100}{1} = 100$

freq. 1 itemsets: $\{a_1\} \{a_2\} \dots \{a_{100}\}$

$\binom{100}{2}$ freq. 2 itemsets: $\{a_1, a_2\} \{a_1, a_3\} \dots \{a_{99}, a_{100}\}$

Tot. no. of freq itemset contains

$$\binom{100}{1} + \binom{100}{2} + \dots + \binom{100}{100} = 2^{100} - 1 \approx 1.2 \times 10^{30}$$

This is too huge number of itemsets for my computer to compute or store.

↓ overcome this difficulty

closed freq. Itemset & maximal freq itemset introduced.

↓

Maximal closed freq Itemset

if none of its immediate supersets is frequent.

Closed Maximal freq Itemset

If none of its immediate supersets has the same support as the itemset.

TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE

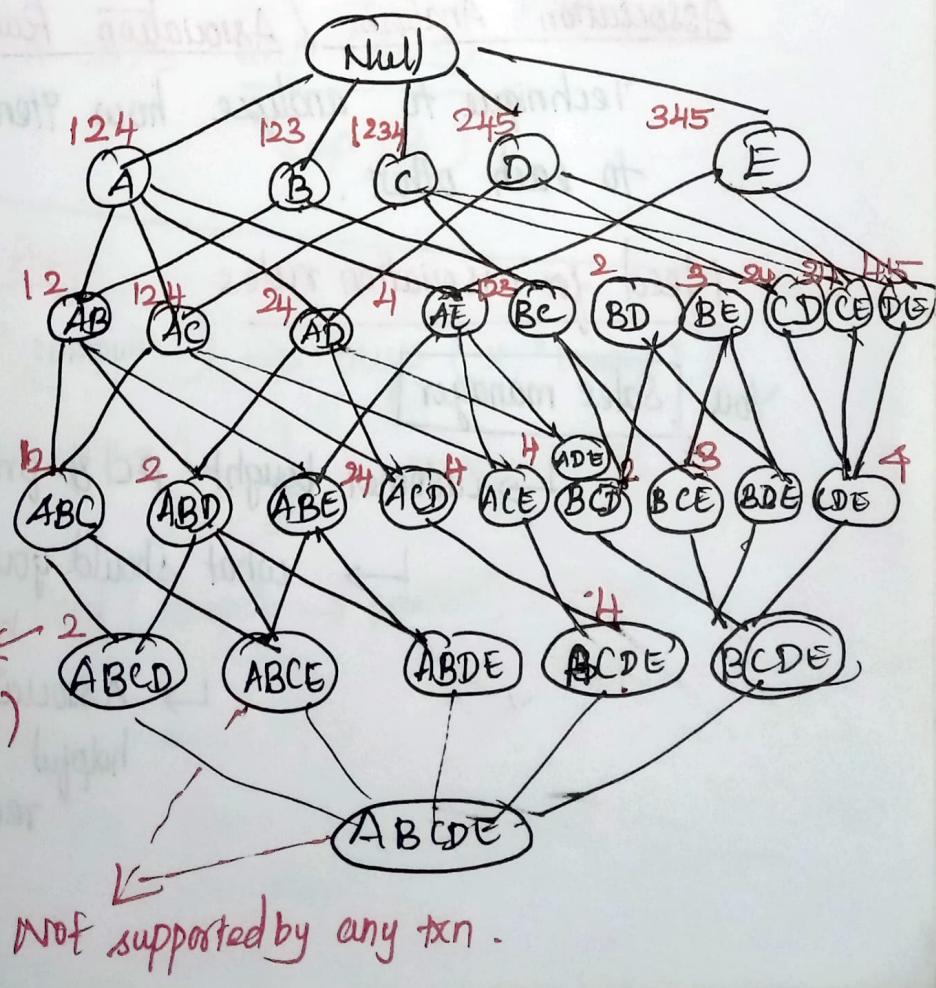
min sup = 2

Closed, not max

~~AC, BC, CD, DE, AE, BE, CE, DE, AC, BC, ACD~~

CE, DE, ABC, ACD

closed = 9
max = 4



Apriori Algorithm : Finding frequent Itemsets by confined Candidate Generation

→ It uses frequent Itemsets to generate association rules.

Itemset :- set of items together

Ex {Milk, Bread, Egg}

→ k Itemset

→ Itemset contains k items

freq:

↳ satisfy min TH value for sup & conf.

Association Analysis / Association Rules

Technique to analyze how items are associated to each other.

Need for association rules

You [Sales manager]

↳ customer bought PC & printer recently

↳ what should you recommend to her next?

↳ Association rules are helpful in making your recommendation

Association rules

determine which items are frequently purchased together within the same txn.

$$A \Rightarrow B$$

↳ implies / determines

some value of itemset A

↳ determines the value of itemset B

under condition in which min sup & conf are met

$$\text{sup}(A) = \frac{\text{No. of txn in which } A \text{ appears}}{\text{Total no. of Txs}}$$

$$\text{conf}(A \rightarrow B) = \frac{\text{sup}(A \cup B)}{\text{sup}(A)}$$

Customer Buys

(X "computer") \Rightarrow buys (X, "printer")

$$\text{sup}(x) = \frac{1}{9} = \cancel{50\%} \quad 1\%$$

$$\text{conf}(\text{computer} \rightarrow \text{printer}) = 50\%$$

Support count (σ):

Freq of occurrence of Itemset

Ex $\sigma(\{\text{Bread, Milk, Sugar}\}) = 2.$

T. ID	Items bought
I ₁	Bread, Milk.
I ₂	Bread, Egg, sugar
I ₃	Milk, sugar, Jam.
I ₄	Bread, Milk, sugar, Jam.
I ₅	Bread, Milk, sugar, Eggs

Apriori algm
finding freq. Itemsets by confined Candidate Generation

2 step process

1. Join step

generate $(k+1)$ itemset from k itemsets
by joining each item with itself.

2. prune step

- scans the count of each item in the DB
- If the candidate item doesn't meet min sup.
→ infrequent → remove it
- step to reduce the size of the candidate itemsets

Ex

TID	Item set
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , 4
T300	2, 3
T400	1, 2, 4
T500	1, 3
T600	2, 3
T700	1, 3
T800	1, 2, 3, 5
T900	1, 2, 3

Step 1

Generate 1^o itemset [C₁]

algm count occurrence of each item.

Itemset	sup count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Step 2

$$\text{min sup count} = 2 //$$

Pruning step: No need of pruning step (No need to del. items having min sup count $\neq 2$) \because all items had atleast min sup count

After pruning also same supcount maintained along with items

I _S	Sup. Comt
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Step ③ Joining - To discover the set of freq 2-itemsets, L₂
 the algm uses the join L₁ \bowtie L₁ to generate a
 candidate set of 2 itemsets C₂

Step ④

L ₁	\bowtie	L ₁	C ₂	Itemset	sup. Comt
I ₁		I ₁		{I ₁ , I ₂ }	4
I ₂		I ₂		{I ₁ , I ₃ }	4
I ₃		I ₃		{I ₁ , I ₄ }	1
I ₄		I ₄		{I ₁ , I ₅ }	2
I ₅		I ₅		{I ₂ , I ₃ }	4
				{I ₂ , I ₄ }	2
				{I ₂ , I ₅ }	2
				{I ₃ , I ₄ }	0
				{I ₃ , I ₅ }	1
				{I ₄ , I ₅ }	0

→ Single set converted
 into a set using join
 operation

Step ⑤ From 2 itemsets L_2 determined

Pruning \rightarrow Remove itemset which has $< \text{minsupcount}$

$\boxed{\text{min sup count} = 2}$

L_2

Itemset	Sup. count
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

Step ⑥ Generate 3 itemsets candidate C_3 using join step

$\boxed{C_3 = L_2 \bowtie L_2}$

L_2	\bowtie	L_2
I_1, I_2		I_1, I_2
I_1, I_3		I_1, I_3
I_1, I_5		I_1, I_5
I_2, I_3		I_2, I_3
I_2, I_4		I_2, I_4
I_2, I_5		I_2, I_5



Itemset	Supcount
I_1, I_2, I_3	2
I_1, I_2, I_5	2
I_1, I_2, I_4	1
I_1, I_3, I_5	1
I_2, I_3, I_4	0
I_2, I_3, I_5	1
I_2, I_4, I_5	0
I_1, I_2, I_4, I_5	1

Condition Analysis classification

- Study of degree of linear relationship between 2 or more variables.

⇒ Relation b/w variables

Ex A, B
x, y, z.

Types of Correlation

- positive → x, y ⇒ Both \uparrow se (Corr) & \downarrow se.
↳ changes equal. Ex Yield, fertilizer
- Negative → x, y ⇒ opposite changes
 $x \uparrow$ se $\Rightarrow y \downarrow$ se (or) $x \downarrow$ se $\Rightarrow y \uparrow$ se
Ex: item purchase based on prize (less)
- Simple → x, y → Using only 2 variable, do correlation simply
- partial → x, y, z, a, b → more variables gn but we calculate only 2 var.
- Multiple → x, y, z, a, b → — do —
Finding correlation b/w 2 but Pts based on all other variables

4 methods under simple linear correlation

- 01 - Scatter diagram.
- 02 - Karl Pearson correlation coefficient Best
- 03 - Spearman's Rank correlation coeff.
- 04 - Correlation coeff. by concurrent method.

Ex Karl Pearson correlation coefficient (r)

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

X :	1	3	5	7	9	1st order	2
Y :	2	4	6	8	10	1st order	positive correlation

calculate coefficient of correlation b/w x & y

x	y	xy	x^2	y^2
1	2	2	1	4
3	4	12	9	
5	6	30		16
7	8	56	25	
9	10	90	49	36
				64
				100
$\sum =$	25	30	190	165
				220

$N = 5$
Total no. of values

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} - \sqrt{N \sum y^2 - (\sum y)^2}}$$

$$= \frac{5 \times 190 - 25 \times 80}{\sqrt{5(625) - 625} - \sqrt{5(900) - (900)}}$$

$(\sum x)^2 = 625$
 $(\sum y)^2 = 900$

$$= \frac{950 - 750}{\sqrt{825 - 625} - \sqrt{1100 - 900}}$$

$$= \frac{200}{\sqrt{200} - \sqrt{200}}$$

$2 \Rightarrow \sqrt{2} \times \sqrt{2}$

$$= \frac{\cancel{200}}{\cancel{200}}$$

$$\boxed{r = 1}$$

Correlation Analysis classification :

Basic Learning / mining tasks

Used to measure the relationship b/w 2 variables

var $x \uparrow \rightarrow y \uparrow \Rightarrow$ +ve correlation

var $x \uparrow \rightarrow y \downarrow \Rightarrow$ -ve correlation

$$\gamma_{A,B} = \frac{\sum (A - A') (B - B')}{(n-1) \sigma_A \sigma_B}$$

γ Karle pearson correlation coefficient

$A' B'$ \Rightarrow Mean of A & B

σ_A, σ_B \Rightarrow std. deviation of A & B

n = no. of tuples in DB.

$r \rightarrow$ 3 values (0, -1, +1)

$r \rightarrow +1 \rightarrow$ perfect positive correlation

$r \rightarrow 0 \rightarrow$ No correlation (no dependence)

$r \rightarrow -1 \rightarrow$ perfect -ve correlation

b/w 2 var.

\hookrightarrow A doesn't affect B
B " A

Ex

A	B
20	8
12	34
9	4
ΣA	ΣB

Σ

$$A' = \frac{20 + 12 + 9}{3} = 13.66$$

(Mean of A)

$$B' = \frac{8 + 34 + 4}{3} = 15.33$$

(Mean of B)

$$\sigma_A = \sqrt{\frac{\sum (A - A')^2}{n-1}}$$

Sigma

$$= \sqrt{\frac{(20 - 13.66)^2 + (12 - 13.66)^2 + (9 - 13.66)^2}{2}}$$

$$= \sqrt{\frac{40.20 + 2 \cdot 7.6 + 21.72}{2}} = \sqrt{\frac{44.68}{2}}$$

$$= \sqrt{32.34}$$

$$= 5.68 \text{ // Ans}$$

$$\sigma_B = \sqrt{\frac{\sum (B - B')^2}{n-1}}$$

$$= \sqrt{\frac{(8 - 15.33)^2 + (34 - 15.33)^2 + (4 - 15.33)^2}{2}}$$

$$= 16.28 \text{ // Ans}$$

$$r_{A,B} = \frac{\sum [(A - A') (B - B')]}{(n-1) \sigma_A \sigma_B}$$

$$= \frac{(20-13.66)(8-15.33) + (12-13.66)(34-15.33) + (9-13.66)(4-15.33)}{2 \times 5.68 \times 16.28}$$

$$= -1.***$$

≈ -1 \rightarrow approximately equal to -1

i.e. Negative Correlation

\hookrightarrow one value \uparrow see, other value \downarrow see

\hookrightarrow Inverse proportion