

relational terms

30/11/2024

UNIT-3

* Frequent Pattern Analysis:

→ Frequent pattern: a pattern that occurs frequently in a data set.

* → Applications:

↳ Basket Data analysis, cross marketing, catalog design, sale campaign analysis, DNA sequence analysis.

→ Importance of frequent patterns:

- ① Association, correlation & causality analysis.
- ② sequential, structural patterns.
- ③ Pattern analysis in spatiotemporal
- ④ classification
- ⑤ cluster analysis
- ⑥ data warehousing
- ⑦ semantic data compression
- ⑧

→ Basic Concepts:

* Market Basket Analysis:

Process of analysing customer b/w diff items that a customer will place in their baskets.

- Itemsets: a set of 1/more items.
- k-itemset: $X = \{x_1, \dots, x_k\}$
- support / support count of X: Frequency of X
- relative support: fraction of transactions that contains X.
(or) probability that a transaction contains X

Ex:

ID	ITEMS
1	Bread, milk
2	Bread, diaper, beer, eggs
3	milk, diaper, beer, cake
4	Bread, milk, diaper, beer
5	bread, milk, diaper, coke

$$\begin{aligned} \text{Support } (A \rightarrow B) &= SPP(A \rightarrow B) \\ &= \frac{s(A \cup B)}{|T|} \end{aligned}$$

① {milk, diaper}

$= \frac{3}{5} \rightarrow$ no. of times the above combination is seen in the item sets.
 $= 0.6 = 60\%$

$$② \{milk, bread\} = \frac{3}{5} = 0.6 = 60\%.$$

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

① {milk, bread}

$$\text{Support } (A \rightarrow B) = 0.6.$$

$$\begin{aligned} \text{Support } (A) &= \text{Support } (\text{milk}) \\ &= \frac{4}{5} = 0.8. \end{aligned}$$

$$\text{Conf } (A \rightarrow B) = \frac{0.6}{0.8} = 0.75.$$

Ex: {milk, bread, beer}

$$\text{Support} = \frac{1}{5} = 0.2$$

$$\text{Conf} = \frac{0.2}{0.8} = \underline{\underline{0.25}}$$

* Closed and Max Patterns:

→ Closed: An itemset X is closed if X is frequent and there exists no super pattern $Y \supset X$ with the same support as X .

→ Max: An itemset X is a max-pattern if X is frequent and there exists no frequent super pattern $Y \supset X$

→ An itemset is said to be frequent if the corresponding support count is greater than minimum support count.

Ex:

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

$$\text{Support}(E) = \frac{3}{5} = 0.6.$$

$$\text{Support}(CE) = \frac{2}{5} = 0.4$$

$$\text{Supp}(DE) = \frac{2}{5} = 0.4$$

$$\rightarrow \text{min-support} = \underline{\underline{2}} \quad \underline{\underline{2}} \quad \underline{\underline{\text{supp}}}$$

$$\therefore \text{supp}(E) \neq \text{support}(CE)$$

$$\text{and } " \neq \text{supp}(DE)$$

→ CE, DE, E are closed ~~pat~~ itemset

$\because \text{support}(E) \geq \text{min-support} \Rightarrow \text{frequent}$

" < " \Rightarrow not frequent

$\because \text{supp}(E) < \text{min-support} \Rightarrow E \text{ is not maximal.}$

$$\text{Supp}(CDE) = \frac{1}{5} = 0.2$$

*→ Download closure property of frequent patterns:

All subsets of any frequent itemset must also be frequent.

Ex: If $\{milk, bread, butter\}$ is a frequent itemset then $\{milk\}$, $\{Bread\}$, $\{butter\}$, $\{milk, bread\}$, $\{milk, butter\}$, $\{bread, butter\}$ are frequent itemsets.

Ex:

TID	List of items
T1	A, B, C, D
T2	A, B, C, D
T3	A, B, C
T4	B, C, D
T5	C, D

Item	Count
A	3
B	4
C	5
D	4

$$\min_supp = 3$$

∴ all the support_counts of items are $\geq \min_supp$
→ all items are frequent.

Superset

Item	Count
AB	3
AC	3
AD	2
BC	4
BD	3
CD	4

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

→ A(count) is not greater than its immediate superset.

A is not closed.

→ A's supersets are having count of 3 which is equal/greater than min-supp; hence A is not maximal.

By B & D are also not closed and not maximal.

for C:

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

∴ C(count) is greater than and not equal to its supersets; hence C is closed; whereas AC, BC, CD being frequent; makes C non-maximal.

further supersets:

$$\text{Count}(AB) \rightarrow$$

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

$$AB(3) \rightarrow ABC(3), ABD(2)$$

AB is not closed.

∴ ABC & ABD are frequent
→ AB is not maximal.

AC is not closed & maximal.

BC is closed & not maximal.

CD "

BD is not closed & maximal.

ABCD \rightarrow not frequent

ABC is closed & maximal.

BCD is closed & maximal

* Computational Complexity of frequent itemset mining:

→ The no. of frequent itemsets to be generated is sensitive to minsup threshold.

→ worst case: M^N

$M \rightarrow$ distinct items

$N \rightarrow$ max length of transactions

* Frequent Item Set Mining Methods:

① Apriori: A candidate generation & test approach

Improve efficiency of Apriori

② FP Growth: A frequent pattern growth approach.

③ ECLAT: Frequent Pattern Mining with vertical data format.

* Apriori:

→ Principle: If there is any itemset which is infrequent, its superset should not be generated / tested.

→ Steps:

→ Initially; scan database once to get frequent 1-itemset.

→ Generate length $(k+1)$ candidate itemsets from length k frequent itemsets.

→ Test the candidates against DB.

→ Terminate when no frequent / candidate set can be generated.

→ Algorithm: $C_k =$ candidate itemset of size k
 $L_k =$ frequent itemset of "

$L_1 = \{$ frequent items $\}$;

for ($k=1$; $L_k \neq \emptyset$; $k++$) do begin

C_{k+1} = candidates generated from L_k .

for each transaction t in database do

↑ment the count of all candidates in C_{k+1} that are contained in t .

$L_{k+1} =$ candidates in C_{k+1} with min. supp
end

→ Objective is to generate an association.

Ex: min supp = 50%.

Threshold confidence = 70%.

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5.

ItemSet	Support	Min Support
1	2	$2/4 = 50\%$.
2	3	75%.
3	3	75%.
4	1	25%.
5	3	75%.

Item Set = (1, 2, 3, 5)

Pairs: (1, 2); (1, 3), (1, 5), (2, 3), (2, 5)
(3, 5)

② ItemSet	Support	Min-Support
(1, 2)	1	$\frac{1}{4} = 25\%$.
(1, 3)	2	50%.
(1, 5)	1	25%.
(2, 3)	2	50%.
(2, 5)	3	75%.
(3, 5)	2	50%.

Item Sets: (1, 3), (2, 3), (2, 5), (3, 5)

③ (1, 2, 3), (1, 3, 5), (2, 3, 5), (1, 2, 5)

ItemSet	Support	Min-Support
(1, 2, 3)	1	25%.
(1, 3, 5)	1	25%.
(2, 3, 5)	2	50%.
(1, 2, 5)	1	25%.

Item Set = (2, 3, 5)

→ Generate Association Rules for (2,3,5)

<u>A → B</u>	<u>Support</u>	<u>Confidence</u>
$(2^3) \rightarrow 5$	2	$\frac{2}{2} = 100\%$
$(3^5) \rightarrow 2$	2	$\frac{2}{2} = 100\%$
$(2^5) \rightarrow 3$	2	$\frac{2}{3} = 66.6\%$
$2 \rightarrow (3^5)$	2	$\frac{2}{3} = 66.6\%$
$3 \rightarrow (2^5)$	2	$\frac{2}{3} = 66.6\%$
$5 \rightarrow (2^5)$	2	$\frac{2}{3} = 66.6\%$

$$\text{Conf} = \frac{\text{Supp}(A \cup B)}{\text{Supp}(A)}$$

$$\text{Conf} = 70\% \text{ (threshold)}$$

⇒ Item Sets = (2,3,5) and (3,5,2).

* Computational Challenges:

- Multiple scans of transaction database.
- Huge number of candidates.
- Tedious workload of support counting for candidates.

* Ideas to improve:

- Reduce passes of transaction database scans.
- Shrink no. of candidates.
- Facilitate support counting.

→ Partition: Scan database only twice:

- * Any itemset that is potentially frequent in DB must be frequent in at least one of the partitions:

- Scan 1: local frequent sets
- Scan 2: global frequent sets.

* Hash Based Techniques:

Ex:	TID	Items
→ Hash	T1	I1, I2, I5
	T2	I2, I4
	T3	I2, I3
	T4	I1, I2, I4
	T5	I1, I3
	T6	I2, I3
	T7	I1, I3
	T8	I1, I2, I3, I5
	T9	I1, I2, I3

$$\text{Order of } \begin{cases} I_1 = 1 \\ I_2 = 2 \end{cases} \left| \begin{array}{l} I_3 = 3 \\ I_4 = 4 \end{array} \right. \left| \begin{array}{l} I_5 = 5 \\ I_6 = 6 \\ I_7 = 7 \\ I_8 = 8 \\ I_9 = 9 \end{array} \right.$$

Itemset	Support count
I1	6
I2	7
I3	6
I4	2
I5	2

$$H(x, y) = (\text{or}(1) * 10 + \text{or}(2)) \bmod 7$$

Item Set	Count	Hash Function
I ₁ , I ₂	4	5
I ₁ , I ₃	4	6
I ₁ , I ₄	1	0
I ₁ , I ₅	2	1
I ₂ , I ₃	4	2
I ₂ , I ₄	2	3
I ₂ , I ₅	2	4
I ₃ , I ₄	0	—
I ₃ , I ₅	1	0
I ₄ , I ₅	0	—

min-supp = 3

Bucket Addr	Bucket count	Bucket contents	L ₂
0	2	I ₁ , I ₄ -1 I ₃ , I ₅ -1	No
1	2	I ₁ , I ₅ -2	No
2	4	I ₂ , I ₃ -4	Yes
3	2	I ₂ , I ₄ -2	No
4	2	I ₂ , I ₅ -2	No
5	4	I ₁ , I ₂ -4	Yes
6	4	I ₁ , I ₃ -4	Yes

Transaction Reduction.

Partitioning

→ Dynamic Itemset Counting:

* algorithm which reduces no. of passes made over the data while keeping the no. of itemsets which are counted in any pass relatively low.

Trans	Items
T ₁	A, B
T ₂	A
T ₃	B, C
T ₄	—

	A	B	C
T ₁	1	1	0
T ₂	1	0	0
T ₃	0	1	1
T ₄	0	0	0

$$\text{Min support} = 25\% \quad M=2 \quad (\text{partition})$$

* Solid box → itemset we have finished counting and exceeds threshold min-supp.



* dashed box → itemset we are still counting that exceeds min-supp.



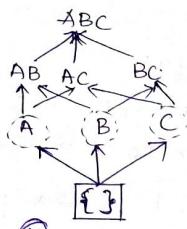
* solid circle → we have finished counting and it is below min-supp.



* dashed circle → we are still counting that is below min-supp.



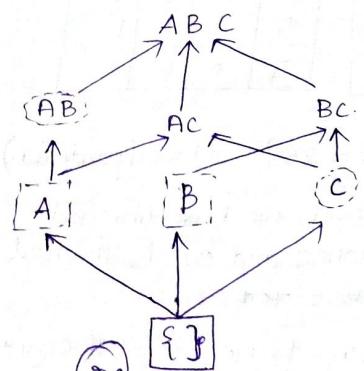
* Before transactions.



$$\begin{array}{l} A = 0 \\ B = 0 \\ C = 0 \end{array}$$

→ After 2 transactions are read: (T_1, T_2)

$$A = 2; B = 1; C = 0; AB = 0$$

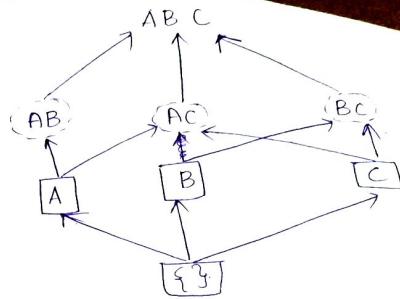


$$\text{min_supp} = 25\% \Rightarrow 1$$

→ After 4 transactions. (T_1, T_2, T_3, T_4)

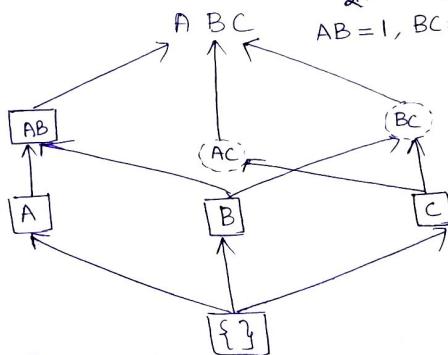
$$A = 2; B = 2; C = 1$$

$$AB = 0; AC = 0; BC = 0$$



→ After 6 transactions: $(T_1, T_2, T_3, T_4, T_5, T_6)$

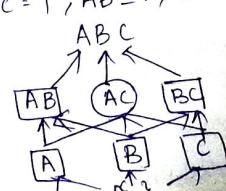
$$A = 2; B = 2; C = 1$$



→ After 8 transactions: $(T_1, T_2, T_3, T_4, T_5, T_6, T_7, T_8)$

$$A = 2; B = 2; C = 1; AB = 1; BC = \emptyset; AC = 0$$

∴ There
are no dashes
left;
we are
done.



Sampling:

Pick up a random sample S of the given data D , then search for frequent item sets in S instead of D .

* Apriori Approach:

→ BFS

→ generates a huge no. of candidates.

* FP growth approach:

→ DFS

→ Avoid explicit candidate generation.

* Benefits of FP tree:

→ Completeness

→ Compactness

* FP tree construction:

① Scan DB once, find frequent 1-itemset

② Sort frequent items in frequency descending Order f-list.

③ Scan DB again; construct FP tree.

* Idea of FP growth method:

→ Recursively grow frequent patterns & database partition.

* Method:

→ Construct conditional pattern-base for each item set, and then its conditional FP-tree.

→ Repeat the process on each newly created conditional FP tree

→ Until the resulting FP tree is empty, or it contains only one path - single path will generate all the combinations of its sub-paths, each of which is a frequent pattern.

* Parallel projection:

→ Project DB in parallel for each frequent item

→ space costly

→ partitions can be processed in parallel.

* Partition projection:

→ partition the DB based on ordered frequent items

→ passing the unprocessed parts to the subsequent partitions

* Advantages of FP growth:

- Decompose both mining task & DB according to the frequent patterns obtained so far.
- focused search of smaller databases.
- compressed database
- No repeated scan

Ex: min-supp = 3

Transaction ID	Items
T1	{E, K, M, N, O, Y}
T2	{D, E, K, N, O, Y}
T3	{A, E, K, M}
T4	{C, K, M, O, Y}
T5	{C, E, I, K, O, O}

Frequency

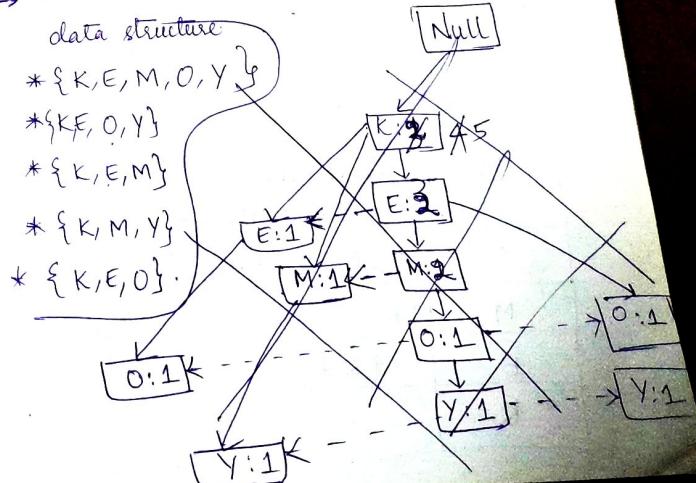
Item	Freq		
A	1	N	2
C	2	O	3
D	1	V	1
E	4	Y	3
I	1		
K	5		
M	3		

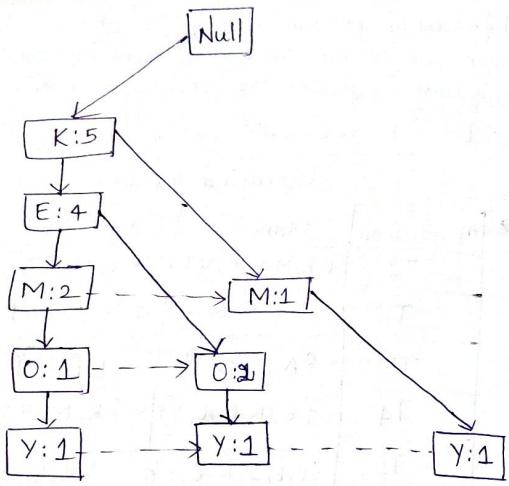
→ A frequency pattern set (L) is built which will contain all the elements whose frequency is greater than or equal to min-supp
 $\therefore L = \{K:5, E:4, M:3, O:3, Y:3\}$

decreasing order

TransactionID	Items	Ordered-Item Set
T1	{E, K, M, N, O, Y}	{K, E, M, O, Y}
T2	{D, E, K, N, O, Y}	{K, E, O, Y}
T3	{A, E, K, M}	{K, E, M}
T4	{C, K, M, O, Y}	{K, M, Y}
T5	{C, E, I, K, O, O}	{K, E, O}

→ ordered item sets are inserted into a tree data structure





* Conditional Pattern Base is computed which is path labels of all the paths which lead to any node of the given item in the frequent pattern tree.

Items	Conditional Pattern Bag	Cond. Frequent Pattern Tree
y	{K, E, M, O: 1}, {K, E, O: 1}, {K, M: 1}, {K: 3}	
O	{K, E, M: 1}, {K, E: 2}	{K, E: 3}
M	{K, E: 2}, {K: 1}	{K: 3}
E	{K: 4}	{K: 4}
K		

Frequent Pattern Tree

Items	Frequent Pattern Generated
y	<K, Y: 3>
O	<K, O: 3>, <E, O: 3>, <E, K, O: 3>
M	<K, M: 3>
E	<E, K: 4>
K	

* ECLAT: Equivalence Class Clustering and Bottom-Up Lattice Traversal

Ex:

TID	Bread	Butter	Milk	Coke	Jam
T1	1	1	0	0	1
T2	0	1	0	1	0
T3	0	1	1	0	0
T4	1	1	0	1	0
T5	1	0	1	0	0
T6	0	1	1	0	0
T7	1	0	1	0	0
T8	1	1	1	0	1
T9	1	1	1	0	0

1 → true 0 → false

min_supp = 2.

① $k=1$ (1 item at a time.)

Item	Tidset
Bread	$T_1, T_4, T_5, T_7, T_8, T_9$
Butter	$T_1, T_2, T_3, T_4, T_6, T_8, T_9$
Milk	$T_3, T_5, T_6, T_7, T_8, T_9$
Coke	T_2, T_4
Jam	T_1, T_8

② $k=2$

Item	Tidset
{Bread, Butter}	T_1, T_4, T_8, T_9
{Bread, Milk}	T_5, T_7, T_8, T_9
{Bread, Coke}	T_4
{Bread, Jam}	T_1, T_8
{Butter, Milk}	T_3, T_6, T_8, T_9
{Butter, Coke}	T_2, T_4
{Butter, Jam}	T_1, T_8
{Milk, Coke}	T_8
{Milk, Jam}	T_8
{Coke, Jam}	

③ $k=3$

Item	Tidset
{Bread, Butter, Milk}	$\{T_8, T_9\}$
{Bread, Butter, Jam}	$\{T_1, T_8\}$
{Bread, Butter, Coke}	$\{T_4\}$
{Bread, Milk, Jam}	$\{T_8\}$

④ $k=4$

Bread, Butter, Milk, Jam $\Rightarrow \underline{\underline{T_8}} < \underline{\underline{1}}$

→ Rules

Items Bought	Products
Bread	butter
bread	MILK
Bread	Jam
Butter	Milk
Butter	Coke
Butter	Jam
Bread & Butter	Milk
Bread & Butter	Jam

Advantages:

- ① Less memory
- ② faster speed
- ③ less no. of computations.

* Pattern Evaluation Methods:

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

* Imbalance Ratio: (IR)

$$IR(A, B) = \frac{|sup(A) - sup(B)|}{sup(A) + sup(B) - sup(A \cup B)}$$

(*)

* Mining multiple-level association rules:

- Hierarchy
- Flexible support settings — lower level has lower support.

- * Some items are more valuable but less frequent.

↳ use non-uniform; group based support.

Eg: {diamond, watch, camera}: 0.05%/
 {bread, milk}: 5%.

* Redundancy filtering: Since mining may find redundant rules due to ancestor relationships.

$$\begin{aligned} &\rightarrow \text{milk} \Rightarrow \text{wheat bread} [\text{support} = 9\%, \\ &\quad \text{cnf} = 72\%] \\ &\rightarrow 2\% \text{ milk} \Rightarrow \text{wheat bread} [\text{support} = 2\%, \\ &\quad \text{cnf} = 72\%] \end{aligned}$$

* A rule is said to be redundant if support is close to expected value.

→ using uniform support for all levels.
 → giving a min-supp. value which is same for all levels.

→ using reduced min-support at lower levels.
 ↳ min supp is given.

→ using item (or) group based min. support.

* Mining multi-dimensional association rules from relational database (or) data warehouse:

↳ Single dimensional rules:

Ex: buys(x, "milk") \Rightarrow buys(x, "bread")

↳ Multi-dimensional rules:

① Integer dimension:

age(x, "19-25") \wedge occup(x, "student") \Rightarrow buys(x, "coca")

② hybrid dimension: (repeated predicates)

$$\text{age}(x, "19-25") \wedge \text{buys}(x, \text{"popcorn"}) \\ \Rightarrow \text{buys}(x, \text{"coke"})$$

↓
predicate (= task)

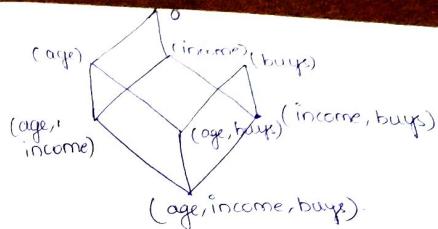
* Mining Quantitative Associations:

→ Based on numerical attributes; methods are divided:

- * static discretization (data cube methods)
- * dynamic discretization (quantitative rules)
- * Clustering
- * Deviation

→ numerical values are replaced by ranges.

- The cells of an n -dimensional cuboid correspond to the predicate sets.
- Mining from data cubes is much faster.



* Negative and Rare patterns:

→ Rare patterns: Very low support but interesting.

Mining: Setting individual-based/group-based support threshold for items.

→ Negative patterns:

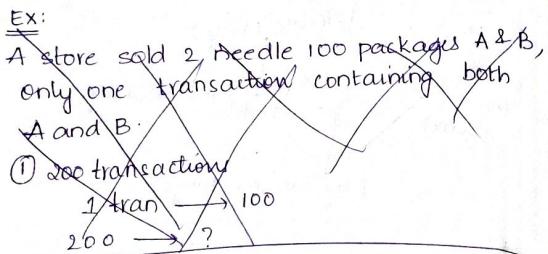
Ex. :: It is unlikely that one buys Ford EXP (SUV) & Toyota Prius (hybrid) together.

→ Ford EXP & Toyota Prius are -vely correlated patterns.

① If itemsets X and Y are both frequent but rarely occur together.

∴ $\text{SUP}(X \cup Y) < \text{SUP}(X) * \text{SUP}(Y)$

→ then X & Y are negatively correlated.



* Constraint Based Frequent Pattern Mining:

→ Mining is done using data mining query language.

→ Advantages of Constraint-based mining :

- * User flexibility
- * Optimization

→ Constraints

① Knowledge type constraint

↳ association, classification, regression

② Data constraint : using SQL like queries (task)

③ dimension/level constraint : hierarchy

④ Rule/pattern constraint ↳ Metarules

⑤ Interestingness constraint : Supp/Confidence ↳ Constraint pushing

decides the type of variables to be used together.

→ Pattern space pruning constraints:

① Anti-monotonic: if constraint c is violated, further mining is stopped

② monotonic: if c is satisfied; no recheck

③ succinct: c must be satisfied; so one can start with datasets satisfying c

④ convertible: c is not monotonic/anti-monotonic; it can be converted by proper ordering.

→ Data space pruning constraint:

① data succinct: Space is pruned with initial pattern

② data anti-monotonic: if transaction t doesn't satisfy c; t can be pruned.

→ Meta-rules guided mining allows the user to give syntactical form of any rule in the form of constraint.