## * FREQUENT PATTERNS

the patterns that appear frequently in dataset

(include frequent data items, sequences, Substructures)

Example: MILK and bread.

__market Basket Analysis:__
process of Analysing customer buying habits by finding the associations b/w the dif. items that a customer.
    will place in their baskets.
— mainly useful for sellers.

__Strategies Used:__

1. placing them together.

2. placing them at ② different ends.

— This Analysis will help sellers to plan their shelf space for increased Sales.

— frequent patterns are represented by association Rule

Ex: Computer and anti-virus.

__Support:__
    identifies how frequently a rule is applied to given dataset.

$$S(P \to Q) = \frac{\sigma(P \cup Q)}{N} \quad (\because N = Total\ Transactions)$$

$$P(A \cup B)$$

---

Right margin (partial, cut off):

Confidence
de-fir
tro

* MINING

Apriori Alg
   - by R.
   - Shows

objective

Example

If

## Confidence :

defines frequent occurence of items of a In transactions of

$$C(p \rightarrow Q) - p(B/A)$$

## * MINING METHODS

- Apriori Algorithm
- Fp Growth Algorithm

## Apriori Algorithm

- by R. Agarwal and R. srikant.

- Shows how objects are associated with each other

objective : To generate on association.

Example :

minimum Support = 50%.
Threshold confidence = 70%.

| T/D | items |
|-----|-------|
| 100 | ① ③ ④ |
| 200 | ② 3 ⑤ |
| 300 | 1 2 3 |
| 400 | 2 5 |

| Itemset | Support | miniSupport |
|---------|---------|-------------|
| 1 | 2 | 2/4 = 50%. |
| 2 | 3 | 3/4 = 75%. |
| 3 | 3 | 3/4 = 75%.   (x) |
| 4 | 1 | 4/4 = 25%. |
| 5 | 3 | 3/4 = 75%. |

Itemset : (1, 2, 3, 5)

— form pairs

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

| itemset | Support | minimum Support |
|---------|---------|-----------------|
| (1,2)   | 1       | 1/4 = 25%  (×) |
| (1,3)   | 2       | 2/4 = 50%.     |
| (1,5)   | 1       | 1/4 = 25%  (×) |
| (2,3)   | 2       | 2/4 = 50%.     |

itemset = (1,3) (2,3) (2,5) (3,5)

— form triplets

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

| itemset   | Support | minimum Support. |
|-----------|---------|------------------|
| (1,2,3)   | 1       | 1/4 = 25%        |
| (1,2,5)   | 1       | 1/4 = 20%.       |
| (1,3,5)   | 1       | 1/4 = 25%.       |
| (2,3,5)   | 2       | 2/4 = 50%.       |

Itemset : (2,3,5)

— now that lets calculate Support and confidence

Confidence: Support (A∪B) / Support of A

using (2,3,5) we can generate association Rules

| Rules        | Support | confidence. |
|--------------|---------|-------------|
| (2^3) → 5    | 2       | 2/2 = 100%. |
| (3^5) → 2    | 2       | 2/2 = 100%. |
| (2^5) → 3    | 2       | 2/3 = 66%. (×) |
| 2 → (3^5)    | 2       | 2/3 = 66%. (×) |
| 5 → (2^3)    | 2       | 2/3 = 66%. (×) |
| 3 → (2^5)    | 2       | 2/3 = 66% (×) |

$(2^{\wedge}3) \to 5$ - confidence $= \dfrac{\text{Support}(A \cup B)}{\text{Support}(A)}$  2,3,5

$$\dfrac{S \cdot ((2^{\wedge}3) \cup 5)}{S(2^{\wedge}3)} = \dfrac{2}{2} = 100\%.$$

$$2 \to (3^{\wedge}5) = \dfrac{S(2 \cup (3^{\wedge}5))}{S(2)} = 2/3 = 66\%.$$

$\therefore$ $(2^{\wedge}3) \to 5$, $(3^{\wedge}5) \to 2$ are association rules

## * FP GROWTH ALGORITHM

$FP \to$ frequent pattern
- is an efficient and scalable method for mining the complete set of FP using a free Structure for storing Information about FP called FP tree.

Example:
minimum Support = 30%.

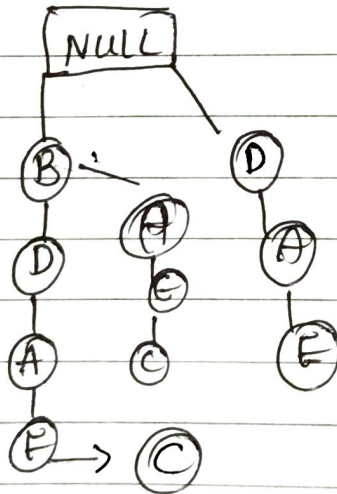| Trans id | items |
|----------|-------|
| 1 | E, A; D B |
| 2 | D, A, E, C, B |
| 3 | C, A, B, E |
| 4 | B, A, D |
| 5 | D |
| 6 | D, B |
| 7 | A, D, E |
| 8 | B, C |

write priorities
More frequency $\to$ More priority
Same frequency $\to$ FCFS

- list out the priorties

| itemset | Frequency | priority |
|---------|-----------|----------|
| A | 5 $\to$ | 3 |
| B | ⑥ $\to$ | 1 |
| C | 3 $\to$ | 5 |
| D | ⑥ $\to$ | 2 |
| E | 4 $\to$ | 4 |

- Order items according to priority.

| Trans ID | Items | Ordered Items |
|----------|-------|---------------|
| 1 | E A D B | B D A E |
| 2 | D A E C B | B D A E C |
| 3 | C A B E | B A E C |
| 4 | B A D | B D A |
| 5 | D | D |
| 6 | D B | B D |
| 7 | A D E | D A E |
| 8 | B C | Bc |



B — 1, 2, 3
D — 1, 2
A — 1, 2
E — 1, 2
C — 1
A — 1
E — 1
C — 1
D — 1, 2
A — 1

\* MINING VARIOUS KINDS OF Association Rules

- ④ types

1. mining Multilevel association Rule
X. mining Uniform Support for all levels

- Using reduced minimum Support at lower
levels.

- using item or group based minimum Support

2. mining Multidimensional association rules from Relational database or data warehouse.

3. mining Multi dimensional association Rules using Static discretisation of Quantitave attributes

4. minining quantitative association Rules

\* CORRELATION ANALYSIS
Used to measure the relationship b/w ② variables.

$$r_{A/B} = \frac{\sum (A - A^1)(B - B^1)}{(n-1) \sigma_A \sigma_B}$$

$r_A$, $r_{A,B}$ = karle pearson correlation cofficient.

$A^1, B^1$ = mean of A and B.
$\sigma_A, \sigma_B$ = Standard deviation of A and B.
$m$ = no of tuples in db.
$r \to$ ③ values $(0, -1, +1)$
$r \to +1 \Rightarrow$ perfect positive correlation
$r \to 0 \Rightarrow$ NO correlation (no dependence)
$r \to -1 \Rightarrow$ perfect Negative Correlation.

Example:

| A | B |
|---|---|
| 20 | 8 |
| 12 | 34 |
|  | 4 |

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

$$A^1 = \frac{20+12+9}{3} = 13.66 \qquad B^1 = \frac{8+34+4}{3} = 15.33$$

$$\sigma A = \sqrt{\frac{\Sigma(A-A^1)^2}{n-1}}$$

$$= \sqrt{\frac{(20-13.66)^2 + (12-13.66)^2 + (9-13.66)^2}{2}} = 5.68$$

$$\sigma B = \sqrt{\frac{\Sigma(B-B^1)^2}{n-1}}$$

$$= \sqrt{\frac{(8-15.33)^2 + (34-16.33)^2 + (4-15.33)^2}{2}} = 16.28$$

$$rA_1B = \frac{(20-13.66)(8-15.33) + (12-13.66)(34-15.33) + (9-13.66)(4-15.33)}{2 \times 5.68 \times 16.28}$$

$$= -1 \cdots$$
$$\approx -1$$

i.e negative correlation

## ＊ CONSTRAINT BASED ASSOCIATION MINING!

Constraint – Condition

– association Rules are generated based on conditions

### ＊ Types of constraints.

#### 1. knowledge Type;

– Specifies the types of knowledge you want to do mining – association, Correlation, Regression etc,

#### 2. Data Constraints

– Specifies the type of data on which you want to generate the Rules.
– only task relevant data

#### 3. dimension level Constraints.
Specifies the dimension or level concept hierarchy.

#### 4. Interestingness Constraints
Support, confidence are used to Identify.

#### 5. Rule Constraints.
Specifies the form of rules to be mined.
② ways

1. metarules guided mining
2. constraint pushing.

## * GRAPH PATTERN MINING

set of tools techniques used to mine frequent Subgroups Subgraphs.

- Used to Analyse the properties of real world graphs

- used to Analyse how structure of graph will effect the rules

### 2 ways

1. Apriori based approaches
2. Pattern growth approaches

### Algorithms used:

1. Gspan → all types
2. closed Graph → closed Subgraphs.

## Applications
1. in xmL Structures
2. anomaly detection
3. Network Analysis
4. control flow Analysis
5. Biological Structures etc,

\* SEQUENTIAL PATTERN MINING: (Spm)

sequence = set of ordered events

Ex: $S = \{ e_1, e_2, e_3, e_4, e_5 \}$

Spm → process of finding frequent subsequences from a set of sequences.

Sequences are represented by "< >"

| Normal transaction data | | | Sequential data. |
|---|---|---|---|
| CID | TID | Transactions | CID Sequences |
| ① | 100 | a,b,c,d | |
| ⑤ 3 | 111 | a,f,d,e | 1. < (abcd) (dep), |
| ① | 122 | d,e,f | (bcde) (aep)> |
| 3 | 133 | b,f,s,a | |
| ① | 144 | b,c,d,e | 3. < (afde) (bfsa), |
| 3 | 155 | a,f,d,c | (afdc). |
| ① | 166 | a,e,p | |

min-Sup =2

Challenges in Spm:
- finding all subsequences

| Sid | Sequence |
|---|---|
| 10 | < a (abc) (ac) d (cf) > |
| 20 | < (ad) c (bc) (ae) > |
| 30 | < (ef) (ab) (df c b)> |
| 40 | < eg (af) (bc > |

min-Sup =2

<(ab)c> (✓)

<eg> (✗)

Algorithms used:
1. GSp (Generalised Sequential patterns)
2. SPADE (vertical format based mining)
3. prefixspam
4. clospam - for closed patterns.