Session - 2    (Tashfeen) (Haiba)    Each basket is subset of items

i22-2041                                          ↑ Date: _____

i22-1855    Association rules (market basket analysis)

generates rules from counts

Support (S) => %age of transactions

$$P(A \cap B) \Big\} \text{ measures frequency of association}$$

Confidence (C) →

$$\text{Conf}(I \to j) = \frac{\text{support}(I \cup j)}{\text{support}(I)} \quad \frac{P(A \cap B)}{P(A)} \Big\} \text{ strength of association}$$

Parameters :-

i) Finding all items that appear frequently $\Big\}$ min support count

ii) Find strong associations amoung frequent items $\Big\}$ Confidence

=> Frequent itemsets ( like treshold = 3 )

⇓

Now, we see that in each basket the items which are repeated atleast 3 times.

⇒ Confidence

$$\text{Conf}(I \to s) = \frac{\text{Support}(I \cup j)}{\text{Support}(I)}$$

⇒ Interest $(I \to j) = \text{Conf}(I \to j) - \text{Pr}[j]$

↓

Interesting rules are those with high interest value ( > 0.5 )

## Example:

$B_1 = \{m, c, b\}$

$B_2 = \{m, b\}$

$B_3 = \{m, p, b\}$

$B_4 = \{c, b, j\}$

$B_5 = \{m, p, j\}$

$B_6 = \{c, j\}$

$B_7 = \{m, c, b, j\}$

$B_8 = \{b, c\}$

$\Rightarrow$ Association rule:- $\{m, b\} \to c$

$$\text{Confidence} = \frac{2}{4} = 0.5$$

$$\text{Interest} = \left(\left|0.5 - \frac{5}{8}\right|\right) = \frac{1}{8}$$

$\Rightarrow$ 
- Support threshold $= 3$, Confidence $= 0.75$
- Frequent itemsets :-

$$\{b, m\}, \{b, c\}, \{c, m\}, \{c, j\}, \{m, c, b\}$$

- Generate rules :-

$b \to m : c = 4/6$       $b \to c \, ; \, c = 5/6$       $c \to m$

.....

khilji.

⇒ Compressing the output (Maximal / closed)

i)

| Items | Support | Maximal (S=3) | Closed |
|-------|---------|---------------|--------|
| A | 4 | | No |
| B | 5 | : | Yes -- |
| C | 3 | | No |
| AB | 4 | | Yes |
| AC | 2 | | No -- |
| BC | 3 | | Yes |
| ABC | 2 | | Yes. |

- 'A' ko items wali list ma dakhna hy kis kay set me 'A' hy. (e.g., AB, AC, ABC)
- Ab 'A' ki support 4 hy. Aur Ab A(support) ≠ ≥

AB (support), AC (support), ABC (support)

Closed ⇒⇒ A (support) ≥ Superset (support)

Maximal ⇒
- A (support) ≥ Threshold
- A (supersets → support < Threshold)

Frequent pairs => no. of sets gets slow with size.

Date: _____

=> Naive Algo → fails if $(items)^2$ - exceed main memory

=> Couting pairs in memory

- Approach 1:- Count pairs using matrix (Uses 4 bytes per pair)

- Approach 2:- Table of triplets (12 bytes for pair)

=> A- Priori Algo (monotonicity)

Example:-
- Support = 2

| Tid | Items | $1^{st}$ scan | Itemset | Supp | $L_1$ | Itemset | Sup |
|-----|-------|---------------|---------|------|-------|---------|-----|
| 10 | A,C,D | → | A | 2 | | A | 2 |
| 20 | B,C,E | | B | 3 | | B | 3 |
| 30 | A,B,C,E | | C | 3 | | C | 3 |
| 40 | B,E | | D | 1 | | E | 3 |
| | | | E | 3 | | | |

| Itemset | Supp | | Itemset | Supp | | $C_2$ Itemset |
|---------|------|---|---------|------|---|---------|
| A,C | 2 | | A,B | 1 | | A,B |
| B,C | 2 | | A,C | 2 | $2^{nd}$ scan | A,C |
| B,E | 3 | ← | A,E | 1 | ← | A,E |
| C,E | 2 | | B,C | 2 | | B,C |
| | | | B,E | 3 | | B,E |
| | | | C,E | 2 | | C,E |

OLYMPIC

| Item set | support. |
|----------|----------|
| A, C, B  | 1        |
| A, B, E  | 1        |
| A, C, E  | 1        |
| B, C, E  | 2        |

→

| Itemset | Supp. |
|---------|-------|
| B,C,E   | 2     |

=> **PCY** (Park- Chen- Yu) Algo

S= 3

1. Items = [ milk , coke , pepsi, cookies, juice ]

milk = 1        ,        Cookies =3        ,        juice =5

coke = 2        ,        pepsi =4        )

$B_1 = \{1,2,3\}$        $B_2 = \{1,4,5\}$

$B_3 = \{1,3\}$        $B_4 = \{2,5\}$

$B_5 = \{1,3,4\}$        $B_6 = \{1,2,3,5\}$

$B_7 = \{2,3,5\}$        $B_8 = \{2,3\}$

① 

| Items | frequency |
|-------|-----------|
| 1     | 5         |
| 2     | 5         |
| 3     | 6         |
| 4     | 2         |
| 5     | 4         |

② Remove elements having frequency less or equal than 2

Candidate set = { 1,2,3,4,5 }

③ Map all candidates sets in pairs and calculate frequeny (Sampling)

$b_1 =$ (1,2), (1,3), (2,3)  $=$  (2,4,4)

$b_2 =$ (1,4), (1,5), (4,5)  $=$  (2,2,1)

$b_3 =$ (1,3)  $=$  (4)

$b_4 =$ (2,5)  $=$  (3)

$b_5 =$ (1,3), (1,4), (3,4)  $=$  (4,2,1)

$b_6 =$ (1,2), (1,3), (1,5), (2,3), (2,5), (3,5) $\doteq$ (2,4,2,4,3,2)

$b_7$  (2,3), (2,5), (3,5)  $=$  (4,3,2)

$b_8$  (2,3)  $=$  (4)

Candidate pairs $=$ (1,3) (2,3) (2,5)

④ Apply hash function  $h(i,j) = (i+j) \% 5 = K$

(1,3) $=$ 4%5 $=$ 4

(2,3) $=$ 5%5 $=$ 0

(2,5) $=$ 7%5 $=$ 2

| Bucket no | pair | high freq | Candidate set |
|---|---|---|---|
| 0 | (2,3) | 4 | (2,3) |
| 2 | (2,5) | 3 | (2,5) |
| 4 | (1,3) | 4 | (1,3) |

Candidate pair= 1,2,3,5