# Association Analysis

## Support:

The support of an association pattern is the percentage of t relevant data transaction for which the pattern is true i.e

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

$$= \frac{\text{No. of tuples containing } A \text{ & } B}{\text{total data set}}$$

## Confidence:-

Confidence is defined as the measure of certainty or trust associated with each discovered pattern i.e.

$$\text{Confidence} (A \Rightarrow B) = P(A \cap B)$$

$$= \frac{\text{No. of tuple containing } A \text{ & } B}{\text{No. of tuple containing } A}$$

## Item set:

A set of item is referred as itemset. An item set containing 'k is called k-itemset.

An itemset satisfies minimum support then it is called fr itemset.

## Association Rule Mining:→

Mining of association rule involves

(i) frequent item generation

(ii) Rules generation.

Given a set of transactions, the goal of association rule mining i find all rules having

(i) support ≥ minimum support threshold.

(ii) confidence ≥ minimum confidence threshold.

## Approaches of Rule Mining:-

### (A) Brute force Approach:

(i) Ast all possible association rules.

(ii) Compute the support and confidence for each rule

(iii) Pure rules that fail the minimum support and minimum confidence threshold

## (5) Apriori Approach:-

→ If an itemset is frequent then all of its subsets must a__
be frequent. or
→ If an itemset is frequent then its subset of non frequent item set is also non-frequ__

Apriori algorithm is an influential algorithm for mining frequ__
itemset.

→ It use a level wise search i.e K item sets are used to explo__
→ K+1 itemset.
→ At first the set is found at level 1, and so on truth no freq__
itemset at level 1 is used to find frequent item set at lev__
2 and so on until no frequent itemset is found.

### Algorithm:-

Step 1 → Read the transaction database and get support from each__
itemset.

Step 2 → compute the support with minimum support to generate at__
of frequent itemset at level 1.

Step 3 → use joint to generate a set of candidate K item set at__
next level.

Step 4 → Generate frequent itemset at next level using minimum__

Step 5 → Repeat. 2 and 3 until no frequent itemset can be generated__

Step 6 → Generate rules from frequent item sets from level 2 onward__
using minimum confidence.

### Example:

| T id | Items |
|---|---|
| 1 | A, E, D |
| 2 | B, C, E |
| 3 | A, B, C, E |
| 4 | B, E, |
| 5 | A, C, E |

Let minimum support $= 38\% = \dfrac{33}{100} \times 5 = 1.65 \approx 2$  (25 to 45% Normally)

minimum confidence $= 70\%$  (65 to 80% normally) → Assume if not given

---

At level 1,
candidate item set

$C_1 =$
| item | count |
|---|---|
| A | 3 |
| B | 3 |
| C | 4 |
| D | 1 |
| E | 4 |

frequent items at level 1

$F_1 =$
| Item | count |
|---|---|
| A | 3 |
| B | 3 |
| C | 4 |
| E | 4 |
| D | 1 |

D is removed as it is less frequent

At level 2,
candidate $C_2 =$
| item | count |
|---|---|
| AB | 1 |
| AC | 3 |
| AE | 2 |
| BC | 3 |
| BE | 3 |
| CE | 3 |

$F_2 =$
| itemset | count |
|---|---|
| AC | 3 |
| AE | 2 |
| BC | 3 |
| BE | 3 |
| CE | 3 |

At level 3,

$C_3 =$
| item | count |
|---|---|
| ACE | 2 |
| ABC | 1 |
| ABE | 1 |
| BCE | 2 |

$F_3 =$
| item | count |
|---|---|
| ACE | 2 |
| BCE | 2 |

At level 4,

$C_4 =$
| item | count |
|---|---|
| ABCE | 1 |

(No frequent item is generated. So, end.)

Now, Generating rules from level 2:

$A \Rightarrow C$, confidence $(A \Rightarrow C) = \dfrac{3}{3} = 100\%$  ← occurrence of A

$C \Rightarrow A$, confidence $(C \Rightarrow A) = \dfrac{3}{4} = 75\%$  → occurrence of A

$A \Rightarrow E$; confidence $(A \Rightarrow E) = \dfrac{2}{3} = 66.67\%$

$E \Rightarrow A$; confidence $(E \Rightarrow A) = \dfrac{2}{4} = 50\%$

Similarly other confidence can be generated

Rules are
$A \Rightarrow C$
$C \Rightarrow A$  } Since minimum confidence is 70%
$A \Rightarrow C$
$C \Rightarrow A$

From Level 3:

$A \Rightarrow CE$, confidence $(A \Rightarrow CE) = \frac{2}{3} = 66.67\%$

$CE \Rightarrow A$, confidence $(CE \Rightarrow A) = \frac{2}{3} = 66.67\%$

$Ac \Rightarrow E$, confidence $(Ac \Rightarrow E) = \frac{2}{3} = 66.67\%$

$E \Rightarrow Ac$, confidence $(E \Rightarrow Ac) = \frac{2}{4} = 50\%$

$C \Rightarrow AE$,

$AE \Rightarrow C$,
.....

**Mining frequent item set without candidate generation.**

**Frequent pattern Growth ( FP-Growth )**

→ FP is divide and conquer strategy

→ It compresses the database representing frequent pattern tree (fp-tree), which retain Itemset association information.

→ Divides the compressed database into a set of conditional database each associated with one frequent item or pattern fragment and mines each such database separately.

**FP-tree Algorithm**

step 1: Creat root node of tree labeled with null

step 2: Scan the complete dataset

step 3: The items in each transactions are processed in sorted order (descending to descending) and branch is created for each transaction.

**FP-tree Mining:**

→ start from each frequent length pattern as an initial suffix pattern.

→ Construct it conditional pattern base (A conditional pattern base is a sub database which consis of the set of prefix paths in the fp-tree co-occurring with suffix pattern)

→ construct each FP-tree and perform mining recursively on such tree.

Example!

| T id | List of Items |
|------|---------------|
| 1. | $I_1, I_2, I_5$ |
| 2. | $I_2, I_4$ |
| 3. | $I_2, I_3$ |
| 4. | $I_1, I_2, I_4$ |
| 5. | $I_1, I_3$ |
| 6. | $I_2, I_3$ |
| 7. | $I_1, I_3$ |
| 8. | $I_1, I_2, I_3$ |
| 9. | $I_1, I_2, I_3, I_5$ |

Let, minimum support $= 2$

| Item | Count |
|------|-------|
| $I_1$ | 6 |
| $I_2$ | 7 |
| $I_3$ | 6 |
| $I_4$ | 2 |
| $I_5$ | 2 |

Sorted order : descending →

| Item | count |
|------|-------|
| $I_2$ | 7 |
| $I_1$ | 6 |
| $I_3$ | 6 |
| $I_4$ | 2 |
| $I_5$ | 2 |



fig. FP-tree.

| item | conditional pattern | conditional FP tree | frequent pattern |
|------|--------------------|--------------------|------------------|
| $I_5$ | $\{I_2,I_1:1\}\{I_2,I_1,I_3:1\}$ | $(I_2:2)(I_1:2)$ | $(I_2,I_5:2)(I_1,I_5:2)(I_2,I_1,I_5:2)$ |
| $I_4$ | $\{I_2,I_1:1\}\{I_2:1\}$ | $(I_2:2)$ | $(I_2,I_4:2)$ |
| $I_3$ | $\{I_2,I_1:2\}\{I_2:2\}\{I_1:2\}$ | $(I_2:4)(I_1:2)(I_1:2)$ | $(I_2,I_3:4)(I_1,I_3:4)(I_2,I_1,I_3:4)$ |
| $I_1$ | $\{I_2:4\}$ | $(I_2:4)$ | $(I_2,I_1:4)$ |

(path to an item $I_5$) (i.e same branch so add)

($I_5$ short first write seperately)

(different branch write seperately)

(make combination with item )

$$F_2 =$$

| item | count |
|------|-------|
| $I_1,I_5$ | 2 |
| $I_1,I_5$ | 2 |
| $I_2,I_4$ | 4 |
| $I_1,I_3$ | 4 |
| $I_2,I_4$ | 4 |

$$F_3 =$$

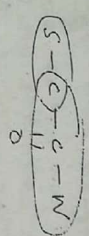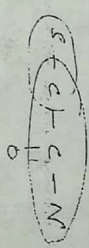| item | count |
|------|-------|
| $I_1,I_2,I_5$ | 2 |
| $I_1\,I_2\,I_3$ | 4 |

**Advantages of fp-growth:-**

(1) It transform the problem of finding long frequent pattern to searchin for shorter ones recursively then concatenating the suffix.

(2) It uses the least frequent item as a suffix.

(3) Has good sensitivity

(4) Reduce the search cost

(5) faster than Apriori

**Disadvantages:-**

(3) when the database is large, It is sometime unrealistic to construct a main memory-based fp-tree.

---

**Sequence pattern:-**

Eg:- $a$ $d$ $c$ $a$ $b$ $c$ $a$ $c$ $d$ $c$ $a$ $c$ $a$

**Subgraph pattern:-**

$$S - C - C - N$$
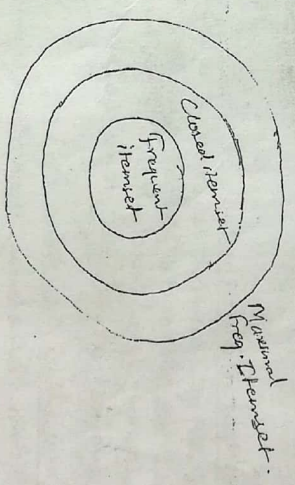$$\underset{O}{\| }$$

$$S - C - C - N$$
$$\underset{O}{\| }$$

**Infrequent pattern:-**

→ Not common items
→ Can hold significant information
→ usefull in scientific research prediction
→ useful in rare item identification.

**Maximal frequent Itemset:-**

An itemset is maximal if non of its immediate superset is freque
(superset होगा)

**Closed item set:**

An itemset is closed if non of its immediate superset has the sam as of the item set.



Closed itemset
Frequent itemset
Maximal freq. Itemset.