

Assignment No : DMW-3

* Title : Apriori Algorithm

* Problem Statement :

Apply Apriori algorithm to find frequently occurring items from given data and generate strong association rules using support and confident thresholds eg: Market Basket Analysis.

* Objectives : To,

Apply apriori algorithm for market basket analysis and operate strong association rules.

* Outcomes: Students will be able to,

Generate frequent itemset of transactions such that support value is greater than minimum support.

* S/W and H/W requirements :

Linux / windows 10 , intel i5 processor , 64-bit PC , 500 GB HDD.

* Concepts related to theory:

Market Basket Analysis association and correlation between transactions and frequent itemsets.

Association rules -

$A \Rightarrow B$ [support = x% confidence]

A transaction T contains set of item A if $A \subset T$

$A \Rightarrow B$ where $A \subset I$, $B \subset I$ ($I = \text{Itemset}$)

$$A \neq \phi, B = \phi \text{ \& } A \cap B = \phi$$

Support ' s ' = % of transaction in dataset ' D ' containing $(A \cup B)$

Confidence ' c ' = % of transactions in ' D ' that containing $A \& B$.

$$s(A \Rightarrow B) = P(A \cup B)$$

$$c(A \Rightarrow B) = P(B|A) = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

Strong Association - when support and confidence is greater than min-support and min-confidence

A priori property -

All non-empty subsets of a frequent itemset must also be frequent

Generate association rules -

- For each frequent itemset L , generate all non-empty subsets of L .

- For every non-empty subset s of L

output: $s \Rightarrow L-s$ if $\frac{\text{support}(L)}{\text{support}(s)} \geq \text{min-support}$

* Algorithms -

1) calculate the support of items sets of size $k=1$ in the transactional ~~base~~ databases. This is called generating the candidate set.

2) Prune the candidate set by eliminating items with a support less than the given threshold.

3) Join the frequent itemset (of size $k+1$) and repeat the above steps until no more itemsets can be formed. This will happen when the sets formed have support less than given support.

eg: support = 3, confidence = 80%.

Transaction database:

Transaction ID

Items

T₁

I₁, I₂, I₃, I₄

T₂

I₁, I₃

T₃

I₃, I₄

T₄

I₁, I₃, I₄

for rule $x \rightarrow y$

support $(x \& y) / \text{support}(y)$

The following can be obtained from the size of two frequent itemsets.

- 1) $I_2 \rightarrow I_3$ confidence = $3/3 = 100\%$.
- 2) $I_3 \rightarrow I_2$ confidence = $3/4 = 75\%$.
- 3) $I_3 \rightarrow I_4$ confidence = $3/4 = 75\%$.
- 4) $I_1 \rightarrow I_3$ confidence = $3/3 = 100\%$.

Since our required confidence is 80%, only rule 1 & 4 are included in the results.

Therefore, it can be concluded that customers who bought item (I₂) always bought item (I₃) with it. Also the customer who bought (I₁) always bought item (I₃) with it.

Applications -

- 1) Education field
- 2) Medical field
- 3) Business optimization etc.

Advantages:

- 1) It is easy to understand and implement
- 2) Join and prune algorithm are easy to process on large itemsets

Disadvantages:

- 1) Requires high computation, if data of itemsets is very large and minimum support is kept very low.
- 2) The entire database needs to be scanned.

* conclusion:

We have successfully implemented and learned Apriori algorithm.

$$A \cup B = AB = combinations$$

$$A \cap B = AB = combinations$$

$$A \cup B = AB = combinations$$

$$A \cup B = AB = combinations$$

$$A \cup B = AB$$

$$A \cup B = AB$$

$$A \cup B = AB$$

$$A \cup B = AB$$