# UNIT 2 DMA Notes

Data Mining And Analytics (SRM Institute of Science and Technology)

# UNIT-2

**Frequent patterns:-**

The patterns that appear frequently in dataset (include frequent dataitems, → eg: computer buying
Sequences, → computer, mouse, keyboard
substructures) → graph, tree etc

eg: Eg: milk & bread → bought together.

**Market basket analysis :-**

→ Process of analysing customer buying habits by finding the association b/w different items that a customer will place in their baskets.

→ Mainly useful for sellers. (they can understand what type of products customers choosing)

**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 <u>association rules</u>.

(Eg:-) computer & antivirus

Computer ⇒ antivirus software[support – 2%; (comp. + antivirus are purchased together) confidence – 60%

(denotes among the total transactions of store 2% computers are purchased along with the antivirus)

(denotes 60% customers who purchased the comp. purchased anti-virus)

**Support:**
Identifies how frequently a rule is applied to given dataset.

$$S(P \to Q) = \frac{\sigma(P \cup Q)}{N} \qquad N = \text{total transactions}$$
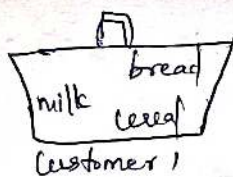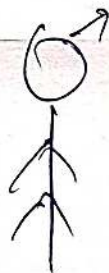
$$\Rightarrow P(A \cup B)$$

association rules {

**Confidence:**
Defines frequent occurrence of items of $Q$ is transactions.

$$\boxed{C(P \to Q) = P(B/A)} \to \text{conditional probability}$$

transactions {



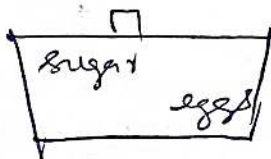which items are frequently purchased

Shopping Baskets

milk  bread
cereal
Customer 1

milk  bread
sugar  eggs.
customer 2

milk  bread
butter
Customer 3

sugar
eggs
Customer n

Market analyst

$*$ Mining frequent itemset using vertical data format.

Given: Horizontal data format.

min. supp = 2.
count

Uses intersection based approach

| TID | list of items ID's |
|---|---|
| $T_{100}$ | $I_1, I_2, I_5$ |
| $T_{200}$ | $I_2, I_4$ |
| $T_{300}$ | $I_2, I_3$ |
| $T_{400}$ | $I_1, I_2, I_4$ |
| $T_{500}$ | $I_1, I_3$ |
| $T_{600}$ | $I_2, I_3$ |
| $T_{700}$ | $I_1, I_3$ |
| $T_{800}$ | $I_1, I_2, I_3, I_5$ |
| $T_{900}$ | $I_1, I_2, I_3$ . |

$\Downarrow$

Vertical data format (1 - Itemset)

| Itemset | TID Set | Sup.count | |
|---|---|---|---|
| $I_1$ | $\{T_{100}, T_{400}, T_{500}, T_{700}, T_{800}, T_{900}\}$ | 6 | ✓ |
| $I_2$ | $\{T_{100}, T_{200}, T_{300}, T_{400}, T_{600}, T_{800}, T_{900}\}$ | 7 | ✓ |
| $I_3$ | $\{T_{300}, T_{500}, T_{600}, T_{700}, T_{800}, T_{900}\}$ | 6 | ✓ |
| $I_4$ | $\{T_{200}, T_{400}\}$ | 2 | ✓ |
| $I_5$ | $\{T_{100}, T_{800}\}$ | 2 | ✓ |

$\therefore I_1, I_2, I_3, I_4, I_5 \Rightarrow$ freq. itemsets

$(\because sup. >= 2)$

## 2- Itemsets In Vertical Data format:

| Itemset | TID Set | Sup. count | |
|---|---|---|---|
| $\{I_1, I_2\}$ | $\{T_{100}, T_{400}, T_{800}, T_{900}\}$ | 4 | |
| $\{I_1, I_3\}$ | $\{T_{500}, T_{700}, T_{800}, T_{900}\}$ | 4 | |
| $\{I_1, I_4\}$ | $\{T_{400}\}$ | 1 | ✗ |
| $\{I_1, I_5\}$ | $\{T_{100}, T_{800}\}$ | 2 | |
| $\{I_2, I_3\}$ | $\{T_{300}, T_{600}, T_{800}, T_{900}\}$ | 4 | |
| $\{I_2, I_4\}$ | $\{T_{200}, T_{400}\}$ | 2 | |
| $\{I_2, I_5\}$ | $\{T_{100}, T_{800}\}$ | 2 | |
| $\{I_3, I_4\}$ | $\{-\}$ | 0 | ✗ |
| $\{I_3, I_5\}$ | $\{T_{800}\}$ | 1 | ✗ |
| $\{I_4, I_5\}$ | $\{-\}$ | 0 | ✗ |

## 3- Itemset in Vertical Data format:

| Itemset | TID Set | Sup. count | |
|---|---|---|---|
| $\{I_1, I_2, I_3\}$ | $\{T_{800}, T_{900}\}$ | 2 | |
| $\{I_1, I_2, I_5\}$ | $\{T_{800}, T_{100}\}$ | 2 | |
| $\{I_1, I_3, I_5\}$ | $\{T_{800}\}$ | 1 | ✗ |
| $\{I_2, I_3, I_4\}$ | $\{-\}$ | 0 | ✗ |
| $\{I_2, I_3, I_5\}$ | $\{T_{800}\}$ | 1 | ✗ |
| $\{I_2, I_4, I_5\}$ | $\{-\}$ | 0 | ✗ |

⇒ 3-Itemset format:

| Itemset | ID |
|---|---|
| $\{I_1, I_2, I_3\}$ | $\{T_{800}, T_{900}\}$ |
| $\{I_1, I_2, I_5\}$ | $\{T_{800}, T_{100}\}$ |

4- Itemset .Vertical data format :

| Itemset | ID | Sup.count | |
|---|---|---|---|
| $\{I_1, I_2, I_3, I_5\}$ | $\{T_{800}\}$ | 1 | ✗ |

∴ If 4-itemset is not frequent we can take 3-itemset as output.

⇒ Mined frequent itemsets are,

$\{I_1, I_2, I_3\}$
$\{I_1, I_2, I_5\}$  } output.

# * Improving the efficiency of apriori algorithm :-

Many methods are available for improving the efficiency of apriori algorithm.

## ① Hash-based technique :-

↳ This method uses a hash-based structure called hash table for generating the k-itemsets and their corresponding count.

↳ It uses a hash function for generating a table.

## ② Transaction Reduction :-

↳ This method reduces the no. of transactions scanned in iterations.

↳ The transactions which do not contain frequent items are marked or removed.

## ③ Partitioning :-

↳ This method requires only 2 database scans to mine the frequent itemsets.

↳ It says that for any itemset to be potentially frequent in the database, it should be frequent in at least one of the partitions of database.

## ④ Sampling :-

↳ This method picks the random sample 'S' from Database 'D', and then searches for frequent itemset in 'S'.

↳ It may be possible to lose a global frequent itemset.

↳ This can be reduced by lowering the min_sup.

• Variation of Apriori which tries to reduce the number of passes made over a transactional db while keeping the no of itemsets counted in

⑤ **Dynamic itemset coupling** :— candidate itemsets

↳ This technique can add new a pass relatively low at any marked start point of the db during the scanning of the database.

\* | From association rule to correlation analysis :— |

↳ Association rule algorithms tends to produce too many rules.

Many of them are uninteresting or redundant.

Redundant if,

$\{A, B, C\} \rightarrow \{D\}$ & $\{A, B\} \rightarrow \{D\}$ have same support & confidence.

∴ Asso. rules consist of support & confidence.

But this support & confidence is insufficient at filtering out uninteresting association rules.

To tackle this weakness, a correlation measures can be used.

This leads to the correlation rule of the form;

| $A \Rightarrow B$ [support, confidence, correlation] |

i.e., a correlation rule is measured by not only sup & confidence but also the correlation b/w the itemsets A & B.

**correlation measures** :—

There are many different correlation measures,

## Lift :-

Lift is a simple correlation measure that is given as follows,

The occurence of itemset A is independent of the occurence of itemset B if $P(A \cup B)$

$$\Rightarrow \boxed{P(A \cup B) = P(A) P(B) ;}$$

itemsets A & B are dependant & correlated.

if there are more than 2 itemset,

$$\boxed{lift(A, B) = P(A \cup B) / P(A) P(B)}$$

i)if value of lift is less than 1 $\Rightarrow$ the occurence of A is negatively correlated with the occurence of B. i.e, the occurence of one likely leads to the absence of the other one.

$\hookrightarrow$ lift value greater than 1 $\Rightarrow$ A & B are positively correlated, meaning that the occurence of one implies the occurence of other.

$\hookrightarrow$ value equal to 1 $\Rightarrow$ A & B are independant & there is no correlation b/w them.

## Eg :-

If the probability of purchasing a computer game is $P(\{game\}) = 0.60$,

the prob. of purchasing a video is $P(\{video\}) = 0.75$,

"       "   both is $P(\{game, video\}) = 0.40$.

$\therefore$ lift value is,

$$\boxed{P(\{game, video\}) / \{P\{game\} \times P(\{video\})\}}$$

$$= 0.40 / \{0.60 \times 0.75\} = \boxed{0.89.}$$

∴ the value is less than 1,
∴ there is a negative correlation b/w the occurence of {game} & {video}.

the second correlation measure is,

$$x^2 = \sum \frac{(observed - expected)^2}{expected}$$

eg'. Contingency table,

| | game | game | $\Sigma_{row}$ |
|---|---|---|---|
| video | 4000 (4500) | 3500 (3000) | 7500 |
| video | 2000 (1500) | 500 (1000) | 2500 |
| $\Sigma$ col. | 6000 | 4000 | 10,000 |

$$x^2 = \sum \frac{(ob - exp)^2}{expected}$$

$$= \frac{(4000 - 4500)^2}{4500} + \frac{(3500 - 3000)^2}{3000} + \frac{(2000 - 1500)^2}{1500}$$

$$+ \frac{(500 - 1000)^2}{1000} = 555 \cdot 6.$$

⟹ $x^2$ value is greater than 1, then observed value of {game, vedio} = 4000 which is less than expected value 4500,

∴ buying game & buying video are negatively correlated.

# Apriori Algorithm

$$\frac{50}{100} \times 4 = 2$$

| TID | Items |
|-----|-------|
| 100 | 1 3 4 |
| 200 | 2 3 5 |
| 300 | 1 2 3 5 |
| 400 | 2 5 |

Min support = 50 %

Threshold Confidence = 70 %

## Iteration 1 : (Step 1 =)

| Itemset | Support | min support |
|---------|---------|-------------|
| 1 | 2 | 2/4 = 50 % |
| 2 | 3 | 3/4 = 75 % |
| 3 | 3 | 3/4 = 75 % |
| 4 | 1 | 1/4 = 25 % (X) |
| 5 | 3 | 3/4 = 75 % |

Itemset (1, 2, 3, 5)

## Iteration 2 (Step 2 =)

Form pairs (1,2) (1,3) (1,5) (2,3) (2,5) (3,5)

| Itemset | Support | min supp |
|---------|---------|----------|
| (1,2) | 1 | 1/4 = 25 % (X) |

| | | |
|---|---|---|
| (1,3) | 2 | 2/4 = 50% |
| (1,5) | 1 | 1/4 = 25% ⊗ |
| (2,3) | 2 | 2/4 = 50% |
| (2,5) | 3 | 3/4 = 75% |
| (3,5) | 2 | 2/4 = 50% |

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

Iteration 3 : Form Triplets (steps)

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

| Itemset | Support | min support |
|---|---|---|
| (1,2,3) | 1 | 1/4 = 25% ✗ |
| (1,2,5) | 1 | 1/4 = 25% ✗ |
| (1,3,5) | 1 | 1/4 = 25% ✗ |
| (2,3,5) | 2 | 2/4 = 50% |

Itemset = (2,3,5)

Now Calcl support & Confidence

Confidence = Support (A∪B) / Support of A
   → freq count of an item

generate Associat<sup>n</sup> rules using (2,3,5)

| 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 % ✗ |

$$(2,3) \underset{A}{} \to \underset{B}{5} \quad \to \quad \text{Confidence} = \frac{\text{Support}(A∪B)}{\text{supp}(A)}$$

$$\frac{S((2,3)∪5)}{S(2,3)} = 2/2 = 100\%$$

$$\underset{A}{2} \rightarrow \underset{B}{(3,5)} = \frac{S(2 \cup (3,5))}{S(2)} = 2/3 \, 566\%$$

$(2,3) \rightarrow 5, \ (3,5) \rightarrow 2$ are

association rules

# Frequent Pattern growth (FP growth)

MS = 2

| TID | List of items ID's |
|---|---|
| $T_{100}$ | $I_1, I_2, I_5$ |
| $T_{200}$ | $I_2, I_4$ |
| $T_{300}$ | $I_2, I_3$ |
| $T_{400}$ | $I_1, I_2, I_4$ |
| $T_{500}$ | $I_1, I_3$ |
| $T_{600}$ | $I_2, I_3$ |
| $T_{700}$ | $I_1, I_3$ |
| $T_{800}$ | $I_1, I_2, I_3, I_5$ |
| $T_{900}$ | $I_1, I_2, I_3$ |

50)

| 1 temset | Sup Cant |
|---|---|
| $\{I_1\}$ | 6 |
| $\{I_2\}$ | 7 |
| $\{I_3\}$ | 6 |
| $\{I_4\}$ | 2 |
| $\{I_5\}$ | 2 |

→

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

us called FP tree



| Itemid | SC | node link |
|--------|-----|-----------|
| $I_2$ | 7 | — |
| $I_1$ | 6 | — |
| $I_3$ | 6 | — |
| $I_4$ | 2 | |
| $I_5$ | 2 | |

| item | Conditional pattern Base | Conditional FP tree | FP's generated |
|------|--------------------------|---------------------|----------------|
| $I_5$ | $\{\{I_2, I_1 : 1\},$ $\{I_2, I_1, I_3 : 1\}\}$ | $\{I_2 : 2, I_1 : 2\}$ | $\{I_2, I_5 : 2\},$ $\{I_1, I_5 : 2\},$ $\{I_2, I_1, I_5 : 2\}$ |
| $I_4$ | $\{\{I_2, I_1 : 1\}, \{I_2 : 1\}\}$ | $\{I_2 : 2\}$ | $\{I_2, I_4 : 2\}$ |
| $I_3$ | $\{\{I_2, I_1 : 2\} \{I_2 : 2\}$ $\{I_1 : 2\}\}$ | $\{I_2 : 4, I_1 : 2\}$ $\{I_1 : 2\}$ | $\{I_2, I_3 : 4\},$ $\{I_1, I_3 : 4\},$ $\{I_2, I_1, I_3 : 2\}$ |
| $I_1$ | $\{\{I_2 : 4\}\}$ | $\{I_2 : 4\}$ | $\{I_2, I_1 : 4\}$ |

1) Mining Frequent patterns/itemsets/Market Basket Analysis

Synopsis

* What is itemset?

* What is Freql- itemset/Pattern? (FP) with example

* Application of FP (Market Basket Analysis)(MBA)
  → What is MBA? with diagram

* Purpose of MBA

* How MBA works?
  → (works on Association rule mining)

→ What is Association rule mining with example

* Strategies used on MBA with example (Computer & antivirus)

* Support

* Confidence

* Mining methods
  a) Apriori Algorithm (defn)
  b) FP growth Algorithm (def)

Apriori Algorithm
## Synopsis

* Definition of Apriori Algorithm
* What is Support? with formula
* What is Confidence? with formula
* what is association rule? with
* Any 2 techniques of improving efficiency of Apriori Algorithm.

Frequent itemset & Closed itemset

### Synopsis

* What is itemset?
* " " Freql itemset?
* Support with formula
* Confidence. " "
* What is closed itemset?
* Example of Closed itemset?

## Pattern growth approach | Frequent pattern (FP) growth Algorithm.

### Synopsis

* Definition of Frequent pattern growth.
* Own example with FP tree.

5) Mining Frequent itemsets Using Vertical data format

Synopsis

* What is itemset?

* What is Freql-itemset?

* What is mining freql-itemset?

* own Example.

6. Association analysis to Correlation Analysis

Refer the pdf notes. for answer.