

Association Rules -> Market Basket Analysis

DSL A COURSE

ROHIT PADEBETTU

Market Basket Analysis

Market Basket Example



Image source: deepclimate.org

Applications

Market Basket Analysis

Understand customer shopping habits

Default Risk Analysis

Understand which customers are more likely to default

Customer Churn Analysis

Understand which customers are likely to switch

Medical Diagnosis

Helps doctors diagnose illness or even find a treatment























Crime Investigation

Helps investigators understand patterns and associations in crimes

Hurricane Predictions























Helps forecasters identify and predict the intensity of storms and hurricanes

Transactions -> Associations

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

























Support

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Measure of how popular an item is

$$\text{Support} \{\text{apple}\} = \frac{4}{8}$$

Confidence











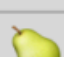











Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Measures how likely the RHS item is purchased given LHS is purchased

$$\text{Confidence} \{ \text{apple} \rightarrow \text{beer} \} = \frac{\text{Support} \{ \text{apple}, \text{beer} \}}{\text{Support} \{ \text{apple} \}}$$

This measure can tend to inflate and show spurious associations when both LHS and RHS are independently popular

Lift

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Measures how likely the RHS item is purchased given LHS is purchased *adjusting for independent popularity of RHS*

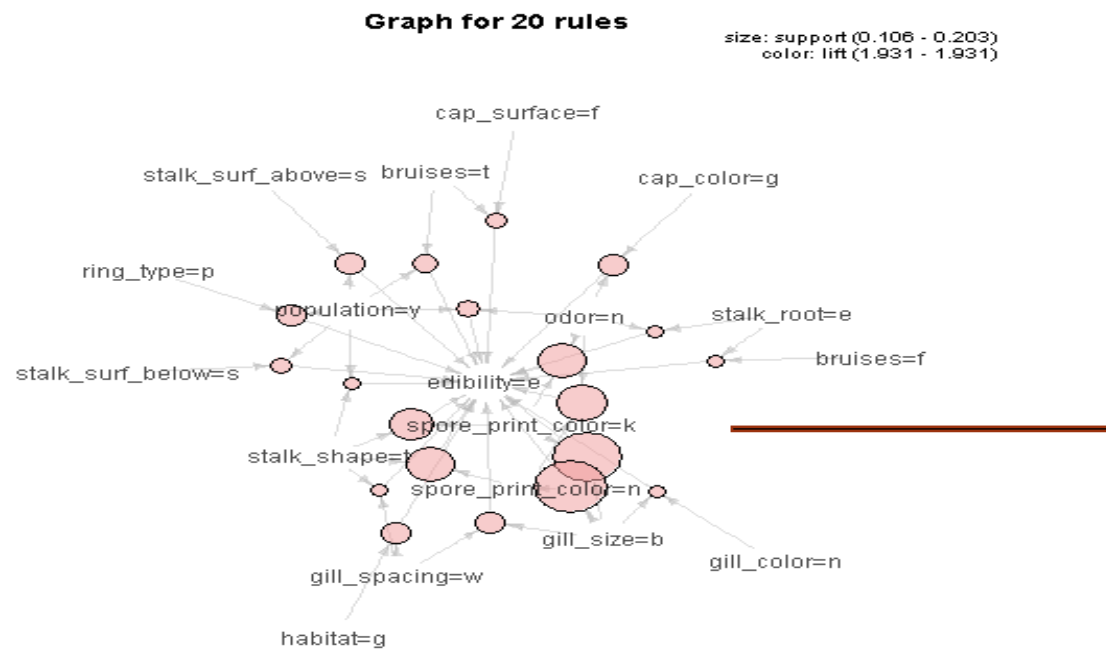
$$\text{Lift} \{ \text{apple} \rightarrow \text{beer} \} = \frac{\text{Support} \{ \text{apple}, \text{beer} \}}{\text{Support} \{ \text{apple} \} \times \text{Support} \{ \text{beer} \}}$$

Association Rules



$$\begin{array}{l}
 \text{Rule: } X \Rightarrow Y \begin{cases} \nearrow \text{Support} = \frac{\text{freq}(X, Y)}{N} \\ \rightarrow \text{Confidence} = \frac{\text{freq}(X, Y)}{\text{freq}(X)} \\ \searrow \text{Lift} = \frac{\text{Support}}{\text{Supp}(X) \times \text{Supp}(Y)} \end{cases}
 \end{array}$$

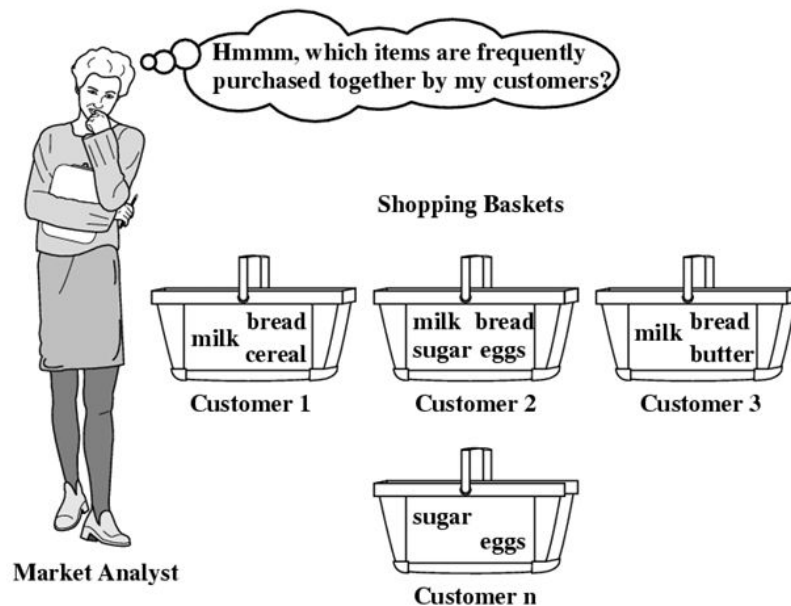
Association Rules - Visualization



Can you guess
the dataset ?


Mushroom Dataset !

Interesting Rules



Association Rules are interesting when they satisfy both a minimum support and a minimum confidence

Useful Rules

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

Transaction	Support
Canned Beer	10%
Soda	20%
Berries	3%
Male Cosmetics	0.5%

Transaction	Support	Confidence	Lift
Canned Beer → Soda	1%	20%	1.0
Canned Beer → Berries	0.1%	1%	0.3
Canned Beer → Male Cosmetics	0.1%	1%	2.6

Association Rules- Caution

Correlation doesn't imply Causation!

The rules below only suggest a strong co-occurrence relationship between items

Causation requires knowledge about Cause and Effect attributes and typically needs information about how relationships evolve over time

Men who purchase diapers also tend to buy beer at the same time !

Transaction	Support	Confidence	Lift
Canned Beer → Soda	1%	20%	1.0
Canned Beer → Berries	0.1%	1%	0.3
Canned Beer → Male Cosmetics	0.1%	1%	2.6

Association Rules - Computation

Brute Force Approach

- List each item in the basket
- List all possible rules from such items
- Count support and confidence of all such rules
- Prune the rules failing minimum thresholds

N item basket $\rightarrow 2^N - 1$ Rules

10 item basket 1023 rules !

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



	Beer	Bread	Milk	Diaper	Eggs	Coke
T_1	0	1	1	0	0	0
T_2	1	1	0	1	1	0
T_3	1	0	1	1	0	1
T_4	1	1	1	1	0	0
T_5	0	1	1	1	0	1

Apriori Algorithm

Mathematical Formulation

```
Apriori( $T, \epsilon$ )  
   $L_1 \leftarrow \{\text{large 1-itemsets}\}$   
   $k \leftarrow 2$   
  while  $L_{k-1} \neq \text{emptyset}$   
     $C_k \leftarrow \{a \cup \{b\} \mid a \in L_{k-1} \wedge b \in \bigcup L_{k-1} \wedge b \notin a\}$   
    for transactions  $t \in T$   
       $C_t \leftarrow \{c \mid c \in C_k \wedge c \subseteq t\}$   
      for candidates  $c \in C_t$   
         $\text{count}[c] \leftarrow \text{count}[c] + 1$   
       $L_k \leftarrow \{c \mid c \in C_k \wedge \text{count}[c] \geq \epsilon\}$   
       $k \leftarrow k + 1$   
  return  $\bigcup_k L_k$ 
```

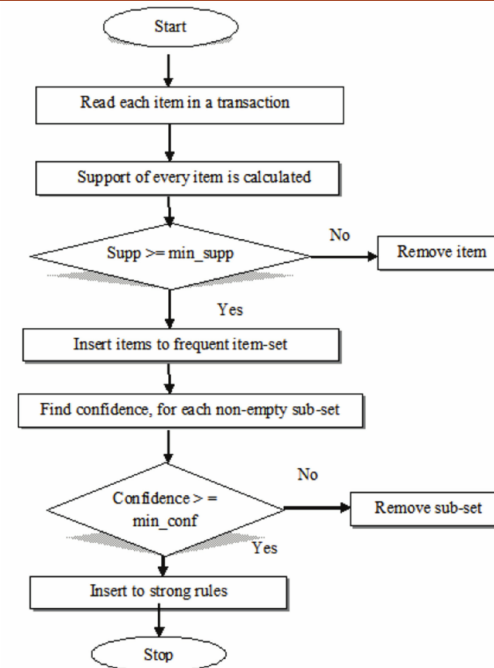
Apriori Algorithm

Pseudo Code & Flow Chart

```

1: procedure APRIORI_FREQUENTITEMSETS( $min\_sup, S$ )
2:    $L_1 \leftarrow itemsets$ 
3:   for  $k = 2; L_{k-1} \neq \emptyset; k++$  do
4:      $C_k = aprioriGen(L_{k-1})$   $\triangleright$  Create the candidates
5:     for each  $c \in C_k$  do
6:        $c.count \leftarrow 0$ 
7:     end for
8:     for each  $I \in S$  do
9:        $C_r \leftarrow subset(C_k, I)$   $\triangleright$  Identify candidates that
       belong to I
10:      for each  $c \in C_r$  do
11:         $c.count++$   $\triangleright$  Counting the support values
12:      end for
13:    end for
14:    if  $c.count \geq min\_sup$  then
15:       $L_k = L_k \cup c$ 
16:    end if
17:  end for
18:  return  $L_k$ 
19: end procedure

```



Apriori Algorithm - Principle

RULE1 : If an “Itemset” is frequent, then all of its subsets must also be frequent

If $\{A,B\}$ is frequent, then both $\{A\}$ & $\{B\}$ are frequent

RULE2 : If an “Itemset” is infrequent, then all of its supersets must also be infrequent

If $\{A\}$ is infrequent, then $\{A,B\}$, $\{A,C\}$ & $\{A,B,C\}$ are infrequent

Anti-Monotonicity

These principles are useful to prune candidates.

Apriori Algorithm - Principle



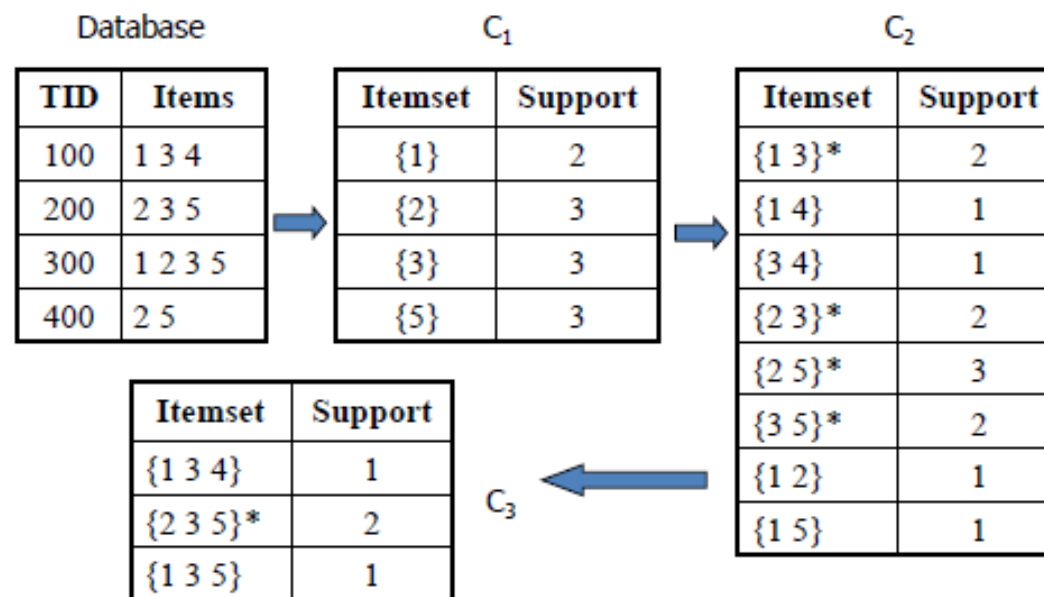
Step 0. Start with itemsets containing just a single item, such as {apple} and {pear}

Step 1. Determine the support for itemsets. Keep the itemsets that meet your minimum support threshold, and remove itemsets that do not

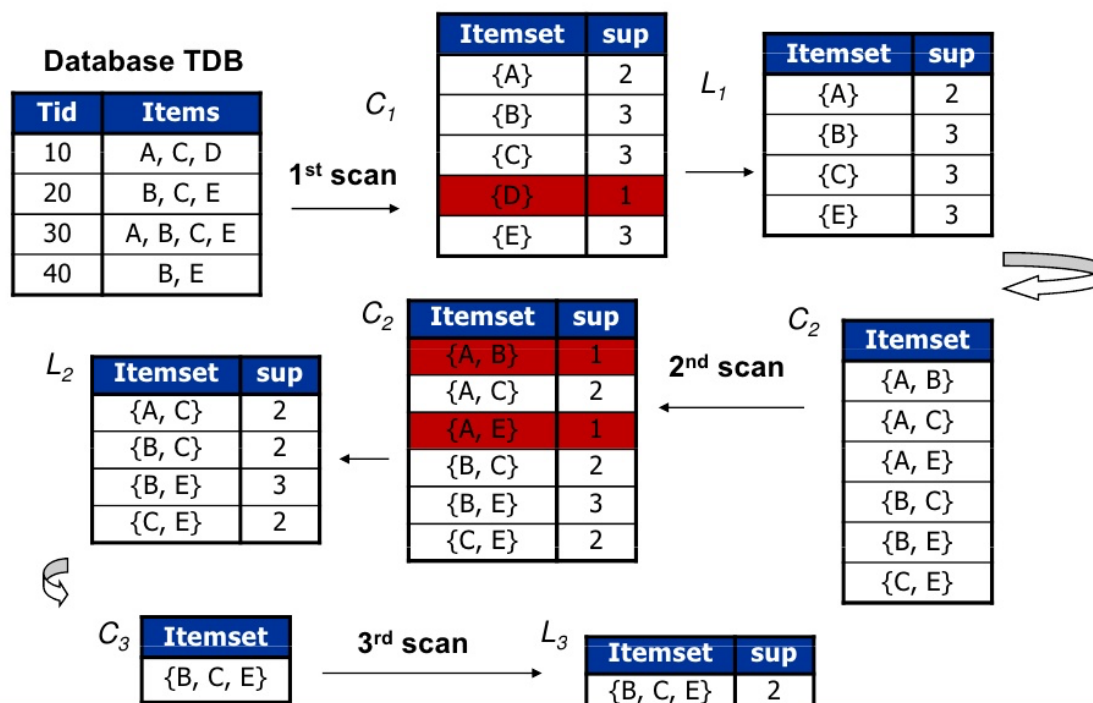
Step 2. Using the itemsets you have kept from Step 1, generate all the possible itemset configurations.

Step 3. Repeat Steps 1 & 2 until there are no more new itemsets

Apriori Algorithm - Example



Apriori Algorithm - Example



Association Rules

Demo

Association Rules

Have good rest of weekend!

