

Data Science

Association Rule Mining

&

Recommendation Engine







01

Association Rule Mining

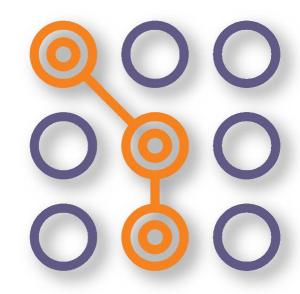
02 Apriori Algorithm

03

Apriori Algorithm in R





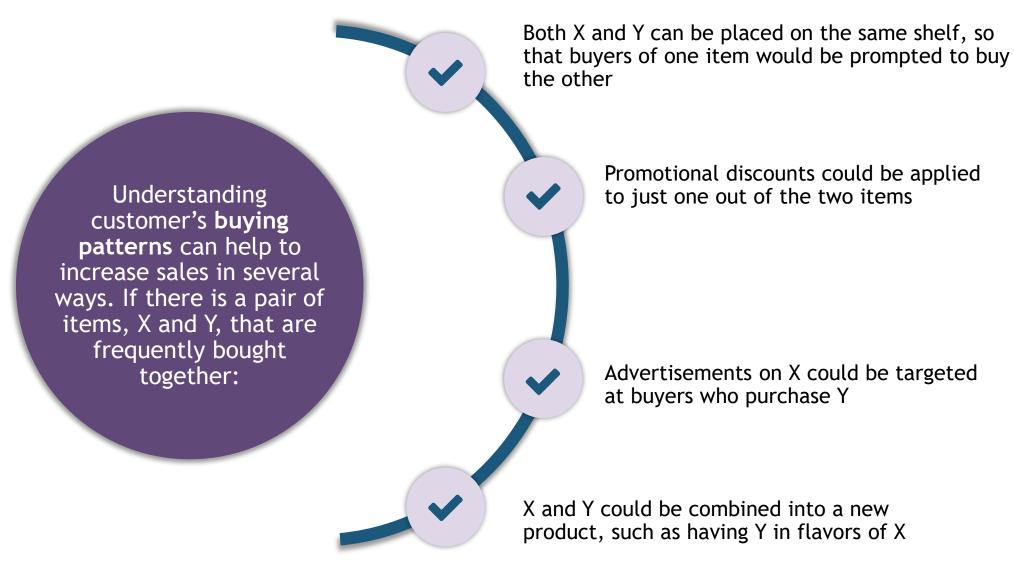


Association rules analysis is a technique to uncover how items are associated to each other



Association mining is usually done on transactions data from a retail market or from an online ecommerce store





3



Since most transactions data is large, the Apriori Algorithm makes it easier to find these patterns or *rules* quickly

So, what is a *rule*?

A rule is a notation that represents which item/s is frequently bought with what item/s

It has an LHS and an RHS part and represented as follows: itemset A => itemset B

This means, the item/s on the right were frequently purchased along with items on the left



Measures of Association Rule Mining

Measures



There are three common ways to measure association

Support

Confidence

Lift







Support



This says how popular an itemset is, as measured by the proportion of transactions in which an itemset appears

$$Support = \frac{Number\ of\ transactions\ with\ both\ A\ and\ B}{Total\ number\ of\ transactions} = P\ (A\cap B)$$

Support



It is an important measure because a rule that has very low support may occur simply by chance

A low support rule is also likely to be uninteresting from a business perspective because it may not be profitable to promote items that customers rarely buy together

For these reasons, Support is often used to eliminate uninteresting rules



Confidence



This says how likely item B is purchased when item A is purchased, expressed as {A -> B}

$$Confidence = \frac{Number\ of\ transactions\ with\ both\ A\ and\ B}{Total\ number\ of\ transactions\ with\ A} = \frac{P\ (A\cap B)}{P\ (A)}$$

Confidence



Confidence measures the reliability of the inference made by a rule

This is measured by the proportion of transactions with item A, in which item B also appears

For a given rule $A \to B$, the higher the confidence, the more likely it is for B to be present in transactions that contain A. Confidence also provides an estimate of the conditional probability of B given A



Lift



Lift says how likely item B is purchased when item A is purchased, while controlling for how popular item B is

$$ExpectedConfidence = \frac{Number\ of\ transactions\ with\ B}{Total\ number\ of\ transactions} = P\left(B\right)$$

$$Lift = \frac{Confidence}{Expected\ Confidence} = \frac{P(A \cap B)}{P(A) \cdot P(B)}$$

Lift



Lift is the factor by which the co-occurrence of A and B exceeds the expected probability of A and B co-occurring, had they been independent

So, higher the lift, higher the chance of A and B occurring together





Apriori Algorithm

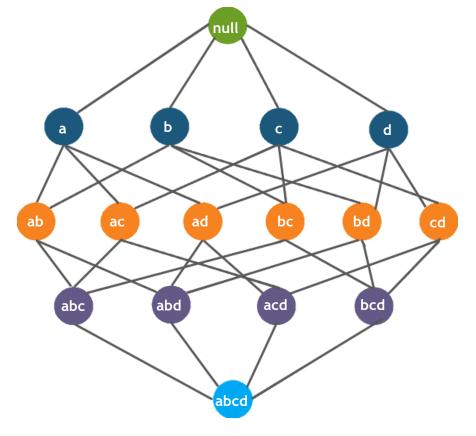
Apriori Algorithm



If an itemset is frequent, then all of its subsets must also be frequent. A key concept in Apriori algorithm is the anti-monotonicity of the support measure

It assumes that:

- All subsets of a frequent itemset must be frequent
- Similarly, for any infrequent itemset, all its supersets must be infrequent too



Working of Apriori Algorithm



The entire algorithm can be divided into two steps:

Step 1 Frequent itemset generation

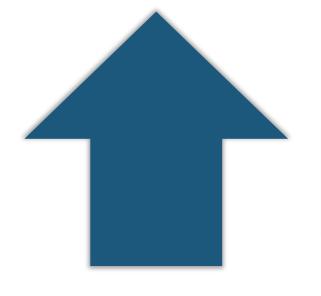
 Find all itemsets for which the support is greater than the threshold support following the process we have already seen earlier in this article

Step 2 Rule generation

- Create rules for each frequent itemset using the binary partition of frequent item-sets and look for the ones with high confidence. These rules are called candidate rules
- For finding association rules, we need to find all rules having support greater than the threshold support and confidence greater than the threshold confidence

Pros & Cons of Apriori Algorithm





Pros

- It is an easy-to-implement and easy-tounderstand algorithm
- It can be used on large itemsets



Cons

- Sometimes, we may need to find a large number of candidate rules which can be computationally expensive
- Calculating support is also expensive because it has to go through the entire database



Association Rule Mining Applications

Association Rule Mining Applications





Identifying value-added services on mobile subscription and recommending what products can be offered to a customer (parallelism)



Understanding which symptoms tend to co-morbid to help improve patient care and medicine prescription



Understanding customer transition so that they can tie up with other vendors for additional discounts



Planning efficient public services and helping public businesses for setting up new factories, shopping malls and so on, as well as marketing particular products





Building the Apriori algorithm on top of the "Groceries" dataset

```
summary(Groceries)
transactions as itemMatrix in sparse format with
 9835 rows (elements/itemsets/transactions) and
 169 columns (items) and a density of 0.02609146
most frequent items:
      whole milk other vegetables
                                       rolls/buns
                                                              soda
                                                                            yogurt
                                                                                            (Other)
                                                                              1372
            2513
                            1903
                                             1809
                                                              1715
                                                                                              34055
element (itemset/transaction) length distribution:
sizes
                                                  11
                         645 545 438 350 246 182 117
                                                                                29
                                                            78
                     26
  Min. 1st Qu. Median
                          Mean 3rd Qu.
  1.000 2.000 3.000 4.409 6.000 32.000
includes extended item information - examples:
       labels level2
1 frankfurter sausage meat and sausage
      sausage sausage meat and sausage
  liver loaf sausage meat and sausage
```

Tasks to be performed



1

Build an Apriori rule on the "Groceries" dataset where the support value is "0.002" & the confidence value is "0.5"

2

Build an apriori rule on the "Groceries" dataset where the support value is "0.007" & the confidence value is "0.6"



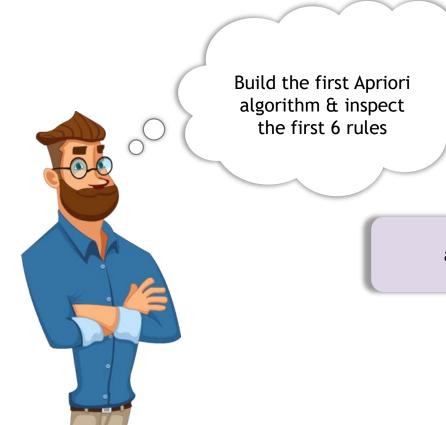
Load the required packages and have a glance at the summary of "Groceries" dataset



library(arules)
library(arulesviz)
data("Groceries")

summary(Groceries)



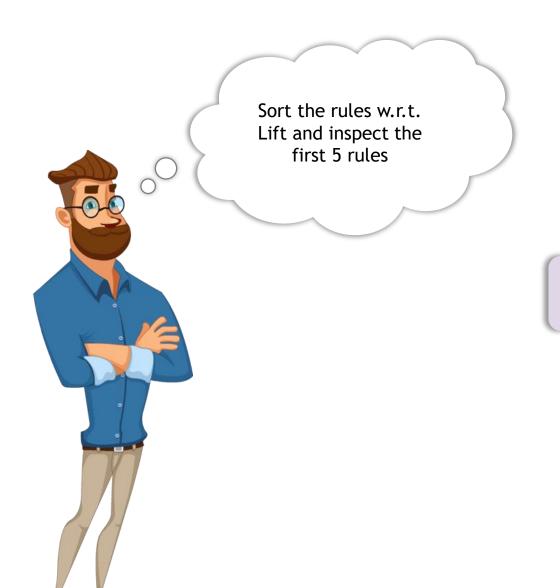


apriori(Groceries,parameter = list(support=0.002, confidence=0.5)) -> rule1



inspect(head(rule1))





inspect(head(sort(rule1,by="lift"),5))





plot(rule1)

plot(rule1,method="grouped")





Build the second apriori algorithm & inspect the first 4 rules

apriori(Groceries,parameter = list(support=0.007, confidence=0.6)) -> rule2



inspect(head(rule2),4)





plot(rule2)

plot(rule2,method="grouped")



Quiz

Quiz



Which package would you have to load to work with the 'apriori()' function?

- 1. Association
- 2. Apriori
- 3. Arules
- 4. Support



Thank You









US: 1-800-216-8930 (TOLL FREE)



sales@intellipaat.com



24X7 Chat with our Course Advisor