



ALY6110: BIG DATA AND DATA MANAGEMENT

METHODS USED IN BIG DATA (2)

MARKET BASKET ANALYSIS

- One source of Big data is the huge amount of data produced at retailing environment.
- Given a large collection of transactions in which each transaction consists of one or more items, association rules go through the items being purchased to see what items are frequently bought together and to discover a list of rules that describe the purchasing behavior.
- The goal with association rules is to discover interesting relationships among the items. The relationships that are interesting depend both on the business context and the nature of the algorithm being used for the discovery.

OVERVIEW

- Each of the uncovered rules is in the form $X \rightarrow Y$, meaning that when item X is observed, item Y is also observed. In this case, the left-hand side (LHS) of the rule is X , and the right-hand side (RHS) of the rule is Y .
- Using association rules, patterns can be discovered from the data that allow the association rule algorithms to disclose rules of related product purchases. The uncovered rules are listed on the right side of this figure.



OVERVIEW

- In the example of a retail store, association rules are used over transactions that consist of one or more items. In fact, because of their popularity in mining customer transactions, association rules are sometimes referred to as **market basket analysis**. Each transaction can be viewed as the shopping basket of a customer that contains one or more items. This is also known as an **itemset**.
- The term itemset refers to a collection of items or individual entities that contain some kind of relationship. This could be a set of retail items purchased together in one transaction, a set of hyperlinks clicked on by one user in a single session, or a set of tasks done in one day. An itemset containing k items is called a **k-itemset**.
- This chapter uses curly braces like $\{item1, item2, \dots, itemk\}$ to denote a k-itemset. Computation of the association rules is typically based on itemsets.

APRIORI ALGORITHM

- The *Apriori* algorithm takes a bottom-up iterative approach to uncovering the **frequent itemsets** by first determining all the possible items (or 1-itemsets, for example {bread}, {eggs}, {milk}, ...) and then identifying which among them are frequent.
- Assuming the minimum **support** threshold (or the minimum support criterion) is set at 0.5, the algorithm identifies and retains those itemsets that appear in at least 50% of all transactions and discards (or “**prunes away**”) the itemsets that have a support less than 0.5 or appear in fewer than 50% of the transactions.
- In the next iteration of the *Apriori* algorithm, the identified frequent *1-itemsets* are paired into *2-itemsets* (for example, {bread, eggs}, {bread, milk}, {eggs, milk}, ...) and again evaluated to identify the frequent 2-itemsets among them.

APRIORI ALGORITHM

One major component of Apriori is support. Given an itemset L , the **support** [2] of L is the percentage of transactions that contain L . For example, if 80% of all transactions contain itemset $\{\text{bread}\}$, then the support of $\{\text{bread}\}$ is 0.8. Similarly, if 60% of all transactions contain itemset $\{\text{bread}, \text{butter}\}$, then the support of $\{\text{bread}, \text{butter}\}$ is 0.6.

A **frequent itemset** has items that appear together often enough. The term “often enough” is formally defined with a **minimum support** criterion. If the minimum support is set at 0.5, any itemset can be considered a frequent itemset if at least 50% of the transactions contain this itemset. In other words, the support of a frequent itemset should be greater than or equal to the minimum support. For the previous example, both $\{\text{bread}\}$ and $\{\text{bread}, \text{butter}\}$ are considered frequent itemsets at the minimum support 0.5. If the minimum support is 0.7, only $\{\text{bread}\}$ is considered a frequent itemset.

If an itemset is considered frequent, then any subset of the frequent itemset must also be frequent. This is referred to as the **Apriori property** (or **downward closure property**). For example, if 60% of the transactions contain $\{\text{bread}, \text{jam}\}$, then at least 60% of all the transactions will contain $\{\text{bread}\}$ or $\{\text{jam}\}$. In other words, when the support of $\{\text{bread}, \text{jam}\}$ is 0.6, the support of $\{\text{bread}\}$ or $\{\text{jam}\}$ is at least 0.6. Figure 5-2 illustrates how the Apriori property works. If itemset $\{B, C, D\}$ is frequent, then all the subsets of this itemset, shaded, must also be frequent itemsets. The Apriori property provides the basis for the Apriori algorithm.

APRIORI ALGORITHM

At each iteration, the algorithm checks whether the support criterion can be met; if it can, the algorithm grows the itemset, repeating the process until it runs out of support or until the itemsets reach a predefined length. The Apriori algorithm [5] is given next. Let variable C_k be the set of candidate k -itemsets and variable L_k be the set of k -itemsets that satisfy the minimum support. Given a transaction database D , a minimum support threshold δ , and an optional parameter N indicating the maximum length an itemset could reach, Apriori iteratively computes frequent itemsets L_{k+1} based on L_k .

APRIORI ALGORITHM

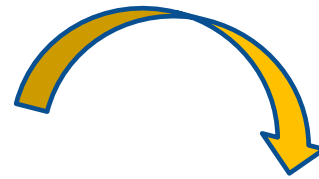
```
1  Apriori ( $D, \delta, N$ )
2     $k \leftarrow 1$ 
3     $L_k \leftarrow \{1\text{-itemsets that satisfy minimum support } \delta\}$ 
4    while  $L_k \neq \emptyset$ 
5      if  $\nexists N \vee (\exists N \wedge k < N)$ 

6       $C_{k+1} \leftarrow$  candidate itemsets generated from  $L_k$ 
7      for each transaction  $t$  in database  $D$  do
8        increment the counts of  $C_{k+1}$  contained in  $t$ 
9       $L_{k+1} \leftarrow$  candidates in  $C_{k+1}$  that satisfy minimum support  $\delta$ 
10      $k \leftarrow k + 1$ 
11  return  $\bigcup_k L_k$ 
```


APRIORI ALGORITHM

| Transaction _ Id | List of Item_Ids |
|------------------|------------------|
| T100 | I2, I1, I5 |
| T200 | I2, I4 |
| T300 | I2, I3 |
| T400 | I2, I1, I4 |
| T500 | I1, I3 |
| T600 | I2, I3 |
| T700 | I1,I3 |
| T800 | I2, I1, I3, I5 |
| T900 | I2,I1, I3 |

Minimum Support count:2



Compare with minimum
support count



| Itemset | Sup.Count |
|---------|-----------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

C1

| Itemset | Sup.Count |
|---------|-----------|
| {I1} | 6 |
| {I2} | 7 |
| {I3} | 6 |
| {I4} | 2 |
| {I5} | 2 |

L1

APRIORI ALGORITHM

| Itemset | Sup.Count |
|----------|-----------|
| {I1,I2} | 4 |
| {I1, I3} | 4 |
| {I1, I4} | 1 |
| {I1, I5} | 2 |
| {I2, I3} | 4 |
| {I2, I4} | 2 |
| {I2, I5} | 2 |
| {I3, I4} | 0 |
| {I3, I5} | 1 |
| {I4, I5} | 0 |

C2

L2

→
Compare with
minimum support
count

| Itemset | Sup.Count |
|----------|-----------|
| {I1,I2} | 4 |
| {I1, I3} | 4 |
| {I1, I5} | 2 |
| {I2, I3} | 4 |
| {I2, I4} | 2 |
| {I2, I5} | 2 |

| Itemset | Sup.Count |
|------------|-----------|
| {I1,I2,I3} | 2 |
| {I1,I2,I5} | 2 |

C3

| Itemset | Sup.Count |
|------------|-----------|
| {I1,I2,I3} | 2 |
| {I1,I2,I5} | 2 |

L3

{I1,I2,I3}: {I1, I2}, {I1, I3}, {I2, I3} \in L2 \rightarrow in C3

{I2,I3,I5}: {I2, I3}, {I2, I5}, {I3, I5} \in L2 \rightarrow not in C3

{I1,I2,I3,I5}: {I1,I2,I3}, {I1,I2, I5}, {I2, I3,I5} \in L3
 \rightarrow not in C4 \rightarrow **Finish**

ASSOCIATION RULES

$$\text{Confidence (A} \rightarrow \text{B)} = P(\text{B}|\text{A}) = \text{Support(A} \wedge \text{B)} / \text{Support(A)}$$

- | | |
|--------------------|-------------------------|
| 1: {I1, I2} → {I5} | Confidence: 2/4 = 50% |
| 2: {I1, I5} → {I2} | Confidence : 2/2 = 100% |
| 3: {I2, I5} → {I1} | Confidence : 2/2 = 100% |
| 4: {I1} → {I2, I5} | Confidence : 2/6 = 33% |
| 5: {I2} → {I1, I5} | Confidence : 2/7 = 29% |
| 6: {I5} → {I1, I2} | Confidence : 2/2 = 100% |

**Minimum
Confidence
Threshold: 70%**

Strong Association Rules

- | | |
|--------------------|-------------------------|
| 2: {I1, I5} → {I2} | Confidence : 2/2 = 100% |
| 3: {I2, I5} → {I1} | Confidence : 2/2 = 100% |
| 6: {I5} → {I1, I2} | Confidence : 2/2 = 100% |

$$\text{Confidence}(X \rightarrow Y) = \frac{\text{Support}(X \wedge Y)}{\text{Support}(X)}$$

MEASURES

$$\textit{Confidence}(X \rightarrow Y) = \frac{\textit{Support}(X \wedge Y)}{\textit{Support}(X)}$$

$$\textit{Lift}(X \rightarrow Y) = \frac{\textit{Support}(X \wedge Y)}{\textit{Support}(X) * \textit{Support}(Y)}$$

$$\textit{Leverage}(X \rightarrow Y) = \textit{Support}(X \wedge Y) - \textit{Support}(X) * \textit{Support}(Y)$$

APRIORI ALGORITHM

- The first step of the *Apriori* algorithm is to identify the frequent itemsets by starting with each item in the transactions that meets the predefined minimum support threshold δ . These itemsets are *1-itemsets* denoted as $L1$, as each *1-itemset* contains only one item. Next, the algorithm grows the itemsets by joining $L1$ onto itself to form new, grown *2-itemsets* denoted as $L2$ and determines the support of each *2-itemset* in $L2$.
- Those itemsets that do not meet the minimum support threshold δ are pruned away. The growing and pruning process is repeated until no itemsets meet the minimum support threshold. Optionally, a threshold N can be set up to specify the maximum number of items the itemset can reach or the maximum number of iterations of the algorithm.
- Once completed, output of the *Apriori* algorithm is the collection of all the frequent *k-itemsets*. Next, a collection of candidate rules is formed based on the frequent itemsets uncovered in the iterative process described earlier. For example, a frequent itemset $\{milk, eggs\}$ may suggest candidate rules $\{milk\} \rightarrow \{eggs\}$ and $\{eggs\} \rightarrow \{milk\}$.

EVALUATION OF CANDIDATE RULES

- **Frequent itemsets** from the previous section can form candidate rules such as X implies Y ($X \rightarrow Y$). This section discusses how measures such as **confidence**, **lift**, and **leverage** can help evaluate the appropriateness of these candidate rules. **Confidence** is defined as the measure of certainty or trustworthiness associated with each discovered rule. Mathematically, confidence is the percent of transactions that contain both X and Y out of all the transactions that contain X .
- For example, if $\{bread, eggs, milk\}$ has a support of 0.15 and $\{bread, eggs\}$ also has a support of 0.15, the confidence of rule $\{bread, eggs\} \rightarrow \{milk\}$ is 1, which means 100% of the time a customer buys bread and eggs, milk is bought as well. The rule is therefore correct for 100% of the transactions containing bread and eggs.

EVALUATION OF CANDIDATE RULES

- A relationship may be thought of as interesting when the algorithm identifies the relationship with a measure of confidence greater than or equal to a predefined threshold. This predefined threshold is called the **minimum confidence**.
- A higher confidence indicates that the rule $(X \rightarrow Y)$ is more interesting or more trustworthy, based on the sample dataset. So far, this chapter has talked about two common measures that the *Apriori* algorithm uses: **support** and **confidence**.
- All the rules can be ranked based on these two measures to filter out the uninteresting rules and retain the interesting ones.

EVALUATION OF CANDIDATE RULES

- Even though confidence can identify the interesting rules from all the candidate rules, it comes with a problem. Given rules in the form of $X \rightarrow Y$, confidence considers only the antecedent (X) and the cooccurrence of X and Y ; it does not take the consequent of the rule (Y) into concern.
- Therefore, confidence cannot tell if a rule contains true implication of the relationship or if the rule is purely coincidental. X and Y can be statistically independent yet still receive a high confidence score.
- Other measures such as **lift** and **leverage** are designed to address this issue. **Lift** measures how many times more often X and Y occur together than expected if they are statistically independent of each other. **Lift** is a measure of how X and Y are really related rather than coincidentally happening together.

EVALUATION OF CANDIDATE RULES

- In theory, leverage is 0 when X and Y are statistically independent of each other. If X and Y have some kind of relationship, the leverage would be greater than zero. A larger leverage value indicates a stronger relationship between X and Y . For the previous example, $Leverage(milk \rightarrow eggs) = 0.3 - (0.5 * 0.4) = 0.1$. and $Leverage(milk \rightarrow bread) = 0.4 - (0.5 * 0.4) = 0.2$. It again confirms that milk and bread have a stronger association than milk and eggs. Confidence is able to identify trustworthy rules, but it cannot tell whether a rule is coincidental.
- A high-confidence rule can sometimes be misleading because confidence does not consider support of the itemset in the rule consequent. Measures such as lift and leverage not only ensure interesting rules are identified but also filter out the coincidental rules. This chapter has discussed four measures of significance and interestingness for association rules: support, confidence, lift, and leverage. These measures ensure the discovery of interesting and strong rules from sample datasets. Besides these four rules, there are other alternative measures, such as correlation, collective strength, conviction, and coverage.

APPLICATIONS OF ASSOCIATION RULES

- The term **market basket analysis** refers to a specific implementation of association rules mining that many companies use for a variety of purposes, including these:
 - ✓ **Broad-scale approaches to better merchandising**—what products should be included in or excluded from the inventory each month
 - ✓ **Cross-merchandising** between products and high-margin or high-ticket items
 - ✓ **Physical or logical placement of product within related categories of products**
 - ✓ **Promotional programs**—multiple product purchase incentives managed through a loyalty card program
- *Besides market basket analysis, association rules are commonly used for recommender systems and clickstream analysis.*

APPLICATIONS OF ASSOCIATION RULES

- Many online service providers such as *Amazon* and *Netflix* use recommender systems. **Recommender systems** can use association rules to discover related products or identify customers who have similar interests. For example, association rules may suggest that those customers who have bought product *A* have also bought product *B*, or those customers who have bought products *A*, *B*, and *C* are more similar to this customer. These findings provide opportunities for retailers to cross-sell their products.
- **Clickstream analysis** refers to the analytics on data related to web browsing and user clicks, which is stored on the client or the server side. Web usage log files generated on web servers contain huge amounts of information, and association rules can potentially give useful knowledge to web usage data analysts. For example, association rules may suggest that website visitors who land on page *X* click on links *A*, *B*, and *C* much more often than links *D*, *E*, and *F*. This observation provides valuable insight on how to better personalize and recommend the content to site visitors.

AN EXAMPLE: TRANSACTIONS IN A GROCERY STORE

An example illustrates the application of the Apriori algorithm to a relatively simple case that generalizes to those used in practice. Using R and the `arules` and `arulesViz` packages, this example shows how to use the Apriori algorithm to generate frequent itemsets and rules and to evaluate and visualize the rules.

The following commands install these two packages and import them into the current R workspace:

```
install.packages('arules')  
install.packages('arulesViz')
```

```
library('arules')  
library('arulesViz')
```

PREREQUISITE

You need:

- ✓ Rtools R-4.2.0
- ✓ rlang package version 1.1.0 (version 1.0.6 is not compatible with arulesViz)

```
install.packages("pkgbuild") # pkgbuild is not available (for R version 3.5.0)
install.packages("devtools") # make sure you have the latest version from CRAN
library(devtools) # load package
devtools::install_github("r-lib/pkgbuild") # install updated version of pkgbuild from GitHub
library(pkgbuild) # load package
find_rtools()
packageVersion("rlang")
install.packages("rlang")
remove.packages("rlang")
packageVersion("rlang") # '1.0.6'
devtools::install_github("r-lib/rlang")
packageVersion("rlang") # '1.1.0.9000'
install.packages("rlang", dependencies = TRUE)
```

THE GROCERIES DATASET

The example uses the *Groceries* dataset from the R *arules* package. The *Groceries* dataset is collected from 30 days of real-world point-of-sale transactions of a grocery store. The dataset contains 9,835 transactions, and the items are aggregated into 169 categories.

```
data(Groceries)
Groceries
transactions in sparse format with
  9835 transactions (rows) and
  169 items (columns)
```

The summary shows that the most frequent items in the dataset include items such as whole milk, other vegetables, rolls/buns, soda, and yogurt. These items are purchased more often than the others.

```
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 |
| 2513 | 1903 | 1809 | 1715 |
| yogurt | (Other) | | |
| 1372 | 34055 | | |

THE GROCERIES DATASET

```
element (itemset/transaction) length distribution:
```

```
sizes
```

```
   1    2    3    4    5    6    7    8    9   10   11   12   13   14
2159 1643 1299 1005  855  645  545  438  350  246  182  117   78   77
   15   16   17   18   19   20   21   22   23   24   26   27   28   29
   55   46   29   14   14    9   11    4    6    1    1    1    1    3
   32
    1
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
1.000   2.000   3.000   4.409   6.000  32.000
```

```
includes extended item information - examples:
```

```
      labels level2      level1
1 frankfurter sausage meet and sausage
2      sausage sausage meet and sausage
3  liver loaf sausage meet and sausage
```

The class of the dataset is `transactions`, as defined by the `arules` package. The `transactions` class contains three slots:

- **transactionInfo:** A data frame with vectors of the same length as the number of transactions
- **itemInfo:** A data frame to store item labels
- **data:** A binary incidence matrix that indicates which item labels appear in every transaction

THE GROCERIES DATASET

```
class(Groceries)
[1] "transactions"
attr(,"package")
[1] "arules"
```

For the *Groceries* dataset, the `transactionInfo` is not being used. Enter `Groceries@itemInfo` to display all 169 grocery labels as well as their categories. The following command displays only the first 20 grocery labels. Each grocery label is mapped to two levels of categories—`level2` and `level1`—where `level1` is a superset of `level2`. For example, grocery label `sausage` belongs to the `sausage` category in `level2`, and it is part of the `meat` and `sausage` category in `level1`. (Note that “meet” in `level1` is a typo in the dataset.)

```
Groceries@itemInfo[1:20,]
```

| | labels | level2 | level1 |
|----|-------------------|------------|----------------------|
| 1 | frankfurter | sausage | meet and sausage |
| 2 | sausage | sausage | meet and sausage |
| 3 | liver loaf | sausage | meet and sausage |
| 4 | ham | sausage | meet and sausage |
| 5 | meat | sausage | meet and sausage |
| 6 | finished products | sausage | meet and sausage |
| 7 | organic sausage | sausage | meet and sausage |
| 8 | chicken | poultry | meet and sausage |
| 9 | turkey | poultry | meet and sausage |
| 10 | pork | pork | meet and sausage |
| 11 | beef | beef | meet and sausage |
| 12 | hamburger meat | beef | meet and sausage |
| 13 | fish | fish | meet and sausage |
| 14 | citrus fruit | fruit | fruit and vegetables |
| 15 | tropical fruit | fruit | fruit and vegetables |
| 16 | pip fruit | fruit | fruit and vegetables |
| 17 | grapes | fruit | fruit and vegetables |
| 18 | berries | fruit | fruit and vegetables |
| 19 | nuts/prunes | fruit | fruit and vegetables |
| 20 | root vegetables | vegetables | fruit and vegetables |

THE GROCERIES DATASET

The following code displays the 10th to 20th transactions of the `Groceries` dataset. The `[10:20]` can be changed to `[1:9835]` to display all the transactions.

```
apply(Groceries@data[,10:20], 2,  
      function(r) paste(Groceries@itemInfo[r,"labels"], collapse=" ", "  
      )
```

Each row in the output shows a transaction that includes one or more products, and each transaction corresponds to everything in a customer's shopping cart. For example, in the first transaction, a customer has purchased whole milk and cereals.

```
[1] "whole milk, cereals"  
[2] "tropical fruit, other vegetables, white bread, bottled water,  
chocolate"  
[3] "citrus fruit, tropical fruit, whole milk, butter, curd, yogurt,  
flour, bottled water, dishes"  
[4] "beef"  
  
[5] "frankfurter, rolls/buns, soda"  
[6] "chicken, tropical fruit"  
[7] "butter, sugar, fruit/vegetable juice, newspapers"  
[8] "fruit/vegetable juice"  
[9] "packaged fruit/vegetables"  
[10] "chocolate"  
[11] "specialty bar"
```

The next section shows how to generate frequent itemsets from the *Groceries* dataset.

FREQUENT ITEMSET GENERATION

- The *apriori()* function from the *arule* package implements the *Apriori* algorithm to create frequent itemsets. Note that, by default, the *apriori()* function executes all the iterations at once.
- However, to illustrate how the *Apriori* algorithm works, the code examples in this section manually set the parameters of the *apriori()* function to simulate each iteration of the algorithm. Assume that the **minimum support threshold** is set to 0.02 based on management discretion. Because the dataset contains 9,853 transactions, an itemset should appear at least 198 times to be considered a frequent itemset.
- The first iteration of the *Apriori* algorithm computes the support of each product in the dataset and retains those products that satisfy the minimum support. The following code identifies 59 frequent 1-itemsets that satisfy the minimum support.
- The parameters of *apriori()* specify the minimum and maximum lengths of the itemsets, the minimum support threshold, and the target indicating the type of association mined.

FREQUENT ITEMSET GENERATION

```
itemsets <- apriori(Groceries, parameter=list(minlen=1, maxlen=1,  
                                              support=0.02, target="frequent itemsets"))
```

parameter specification:

| confidence | minval | smax | arem | aval | originalSupport | support | minlen |
|------------|--------|------|------|-------|-----------------|---------|--------|
| 0.8 | 0.1 | 1 | none | FALSE | TRUE | 0.02 | 1 |

| maxlen | target | ext |
|--------|-------------------|-------|
| 1 | frequent itemsets | FALSE |

algorithmic control:

| filter | tree | heap | memopt | load | sort | verbose |
|--------|------|------|--------|------|------|---------|
| 0.1 | TRUE | TRUE | FALSE | TRUE | 2 | TRUE |

```
apriori - find association rules with the apriori algorithm  
version 4.21 (2004.05.09)          (c) 1996-2004  Christian Borgelt  
set item appearances ...[0 item(s)] done [0.00s].  
set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].  
sorting and recoding items ... [59 item(s)] done [0.00s].  
creating transaction tree ... done [0.00s].  
checking subsets of size 1 done [0.00s].  
writing ... [59 set(s)] done [0.00s].  
creating S4 object ... done [0.00s].
```

FREQUENT ITEMSET GENERATION

- ✓ The summary of the itemsets shows that the support of *1-itemsets* ranges from 0.02105 to 0.25552.
- ✓ Because the maximum support of the *1-itemsets* in the dataset is only 0.25552, to enable the discovery of interesting rules, the minimum support threshold should not be set too close to that number.

```
summary(itemsets)
set of 59 itemsets

most frequent items:
frankfurter      sausage      ham      meat      chicken
      1          1          1          1          1
(Other)
      54

element (itemset/transaction) length distribution:sizes
1
59

      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
      1      1      1      1      1      1

summary of quality measures:
      support
Min.      :0.02105
1st Qu.:0.03015
Median :0.04809
Mean    :0.06200
3rd Qu.:0.07666
Max.    :0.25552

includes transaction ID lists: FALSE

mining info:
      data ntransactions support confidence
Groceries      9835      0.02      1
```

FREQUENT ITEMSET GENERATION

The following code uses the `inspect()` function to display the top 10 frequent 1-itemsets sorted by their support. Of all the transaction records, the 59 1-itemsets such as {whole milk}, {other vegetables}, {rolls/buns}, {soda}, and {yogurt} all satisfy the minimum support. Therefore, they are called frequent 1-itemsets.

```
inspect(head(sort(itemsets, by = "support"), 10))
```

| | items | support |
|----|--------------------|------------|
| 1 | {whole milk} | 0.25551601 |
| 2 | {other vegetables} | 0.19349263 |
| 3 | {rolls/buns} | 0.18393493 |
| 4 | {soda} | 0.17437722 |
| 5 | {yogurt} | 0.13950178 |
| 6 | {bottled water} | 0.11052364 |
| 7 | {root vegetables} | 0.10899847 |
| 8 | {tropical fruit} | 0.10493137 |
| 9 | {shopping bags} | 0.09852567 |
| 10 | {sausage} | 0.09395018 |

FREQUENT ITEMSET GENERATION

In the next iteration, the list of frequent 1-itemsets is joined onto itself to form all possible candidate 2-itemsets. For example, 1-itemsets {whole milk} and {soda} would be joined to become a 2-itemset {whole milk, soda}. The algorithm computes the support of each candidate 2-itemset and retains those that satisfy the minimum support. The output that follows shows that 61 frequent 2-itemsets have been identified.

```
itemsets <- apriori(Groceries, parameter=list(minlen=2, maxlen=2,
                                              support=0.02, target="frequent itemsets"))
```

```
parameter specification:
```

```
confidence minval smax arem aval originalSupport support minlen
          0.8    0.1    1 none FALSE              TRUE    0.02    2
maxlen              target  ext
          2 frequent itemsets FALSE
```

```
algorithmic control:
```

```
filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE    2    TRUE
```

```
apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)          (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].
sorting and recoding items ... [59 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 done [0.00s].
writing ... [61 set(s)] done [0.00s].
creating S4 object ... done [0.00s].
```

FREQUENT ITEMSET GENERATION

The summary of the itemsets shows that the support of 2-itemsets ranges from 0.02003 to 0.07483.

```
summary(itemsets)
set of 61 itemsets

most frequent items:
      whole milk other vegetables      yogurt      rolls/buns
           25           17           9           9
           soda           (Other)
           9           53

element (itemset/transaction) length distribution:sizes
  2
61
```

FREQUENT ITEMSET GENERATION

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|---------|--------|------|---------|------|
| 2 | 2 | 2 | 2 | 2 | 2 |

summary of quality measures:

support

| | |
|----------|----------|
| Min. | :0.02003 |
| 1st Qu.: | 0.02227 |
| Median | :0.02613 |
| Mean | :0.02951 |
| 3rd Qu.: | 0.03223 |
| Max. | :0.07483 |

includes transaction ID lists: FALSE

mining info:

| data | ntransactions | support | confidence |
|-----------|---------------|---------|------------|
| Groceries | 9835 | 0.02 | 1 |

FREQUENT ITEMSET GENERATION

- The top 10 most frequent 2-*itemsets* are displayed next, sorted by their *support*. Notice that whole milk appears 6 times in the top 10 2-*itemsets* ranked by *support*. As seen earlier, {*whole milk*} has the highest support among all the 1-*itemsets*.
- These top 10 2-*itemsets* with the highest *support* may not be interesting; this highlights the limitations of using support alone.

```
inspect(head(sort(itemsets, by = "support"), 10))
```

| | items | support |
|----|--|------------|
| 1 | {other vegetables, whole milk} | 0.07483477 |
| 2 | {whole milk, rolls/buns} | 0.05663447 |
| 3 | {whole milk, yogurt} | 0.05602440 |
| 4 | {root vegetables, whole milk} | 0.04890696 |
| 5 | {root vegetables, other vegetables} | 0.04738180 |
| 6 | {other vegetables, yogurt} | 0.04341637 |
| 7 | {other vegetables, rolls/buns} | 0.04260295 |
| 8 | {tropical fruit, whole milk} | 0.04229792 |
| 9 | {whole milk, soda} | 0.04006101 |
| 10 | {rolls/buns, soda} | 0.03833249 |

FREQUENT ITEMSET GENERATION

- Next, the list of frequent 2-*itemsets* is joined onto itself to form candidate 3-*itemsets*.
- For example {*other vegetables*, *whole milk*} and {*whole milk*, *rolls/buns*} would be joined as {*other vegetables*, *whole milk*, *rolls/buns*}.
- The algorithm retains those itemsets that satisfy the *minimum support*.
- The following output shows that only two frequent 3-*itemsets* have been identified.

```
itemsets <- apriori(Groceries, parameter=list(minlen=3, maxlen=3,
                                              support=0.02, target="frequent itemsets"))

parameter specification:
  confidence minval  smax  arem  aval originalSupport  support  minlen
           0.8     0.1    1 none  FALSE              TRUE    0.02      3
  maxlen                target    ext
           3 frequent itemsets FALSE

algorithmic control:
  filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE    2    TRUE

apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)                (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].
sorting and recoding items ... [59 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 3 done [0.00s].
writing ... [2 set(s)] done [0.00s].
creating S4 object ... done [0.00s].
```

FREQUENT ITEMSET GENERATION

The 3-itemsets are displayed next:

```
inspect(sort(itemsets, by = "support"))
  items          support
1 {root vegetables,
  other vegetables,
  whole milk}    0.02318251
2 {other vegetables,
  whole milk,
  yogurt}        0.02226741
```

In the next iteration, there is only one candidate 4-itemset {root vegetables, other vegetables, whole milk, yogurt}, and its support is below 0.02. No frequent 4-itemsets have been found, and the algorithm converges.

```
itemsets <- apriori(Groceries, parameter=list(minlen=4, maxlen=4,
                                              support=0.02, target="frequent itemsets"))
```

```
parameter specification:
confidence minval smax arem  aval originalSupport support minlen
      0.8    0.1    1 none FALSE          TRUE    0.02     4
maxlen          target  ext
      4 frequent itemsets FALSE
```

[illegible]

```
algorithmic control:
  filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE    2    TRUE

apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)          (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].
sorting and recoding items ... [59 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 3 done [0.00s].
writing .. [0 set(s)] done [0.00s].
creating S4 object ... done [0.00s].
```

The previous steps simulate the Apriori algorithm at each iteration. For the *Groceries* dataset, the iterations run out of support when $k = 4$. Therefore, the frequent itemsets contain 59 frequent 1-itemsets, 61 frequent 2-itemsets, and 2 frequent 3-itemsets.

When the `maxlen` parameter is not set, the algorithm continues each iteration until it runs out of support or until `k` reaches the default `maxlen=10`. As shown in the code output that follows, 122 frequent itemsets have been identified. This matches the total number of 59 frequent 1-itemsets, 61 frequent 2-itemsets, and 2 frequent 3-itemsets.

[illegible]

FREQUENT ITEMSET GENERATION

parameter specification:

```
confidence minval smax arem aval originalSupport support minlen
          0.8    0.1    1 none FALSE          TRUE    0.02    1
maxlen          target  ext
          10 frequent itemsets FALSE
```

algorithmic control:

```
filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE    2    TRUE
```

apriori - find association rules with the apriori algorithm

version 4.21 (2004.05.09) (c) 1996-2004 Christian Borgelt

set item appearances ...[0 item(s)] done [0.00s].

set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].

sorting and recoding items ... [59 item(s)] done [0.00s].

creating transaction tree ... done [0.00s]

checking subsets of size 1 2 3 done [0.00s].

writing ... [122 set(s)] done [0.00s].

creating S4 object ... done [0.00s].

Note that the results are assessed based on the specific business context of the exercise using the specific dataset. If the dataset changes or a different minimum support threshold is chosen, the Apriori algorithm must run each iteration again to retrieve the updated frequent itemsets.

RULE GENERATION AND VISUALIZATION

The `apriori()` function can also be used to generate rules. Assume that the minimum support threshold is now set to a lower value 0.001, and the minimum confidence threshold is set to 0.6. A lower minimum support threshold allows more rules to show up. The following code creates 2,918 rules from all the transactions in the *Groceries* dataset that satisfy both the minimum support and the minimum confidence.

```
rules <- apriori(Groceries, parameter=list(support=0.001,
                                           confidence=0.6, target = "rules"))

parameter specification:
  confidence minval  smax  arem  aval originalSupport support  minlen
         0.6      0.1    1 none  FALSE              TRUE   0.001      1
  maxlen target   ext
         10 rules FALSE

algorithmic control:
  filter tree heap memopt load sort verbose
    0.1 TRUE TRUE  FALSE TRUE    2    TRUE

apriori - find association rules with the apriori algorithm
version 4.21 (2004.05.09)          (c) 1996-2004  Christian Borgelt
set item appearances ...[0 item(s)] done [0.00s].
set transactions ...[169 item(s), 9835 transaction(s)] done [0.00s].
sorting and recoding items ... [157 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 3 4 5 6 done [0.01s].
writing ... [2918 rule(s)] done [0.00s].
creating S4 object ... done [0.01s].
```

RULE GENERATION AND VISUALIZATION

The summary of the rules shows the number of rules and ranges of the support, confidence, and lift.

```
summary(rules)
```

```
set of 2918 rules
```

```
rule length distribution (lhs + rhs):sizes
```

```
  2    3    4    5    6
  3  490 1765  626   34
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
2.000   4.000   4.000   4.068   4.000   6.000
```

```
summary of quality measures:
```

| support | confidence | lift |
|------------------|----------------|----------------|
| Min. :0.001017 | Min. :0.6000 | Min. : 2.348 |
| 1st Qu.:0.001118 | 1st Qu.:0.6316 | 1st Qu.: 2.668 |
| Median :0.001220 | Median :0.6818 | Median : 3.168 |
| Mean :0.001480 | Mean :0.7028 | Mean : 3.450 |
| 3rd Qu.:0.001525 | 3rd Qu.:0.7500 | 3rd Qu.: 3.692 |
| Max. :0.009354 | Max. :1.0000 | Max. :18.996 |

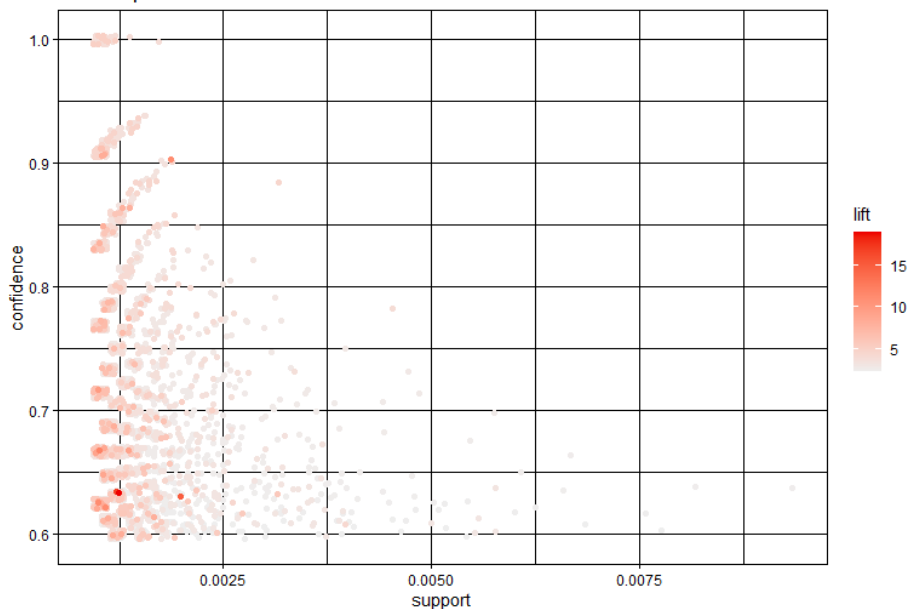
RULE GENERATION AND VISUALIZATION

mining info:

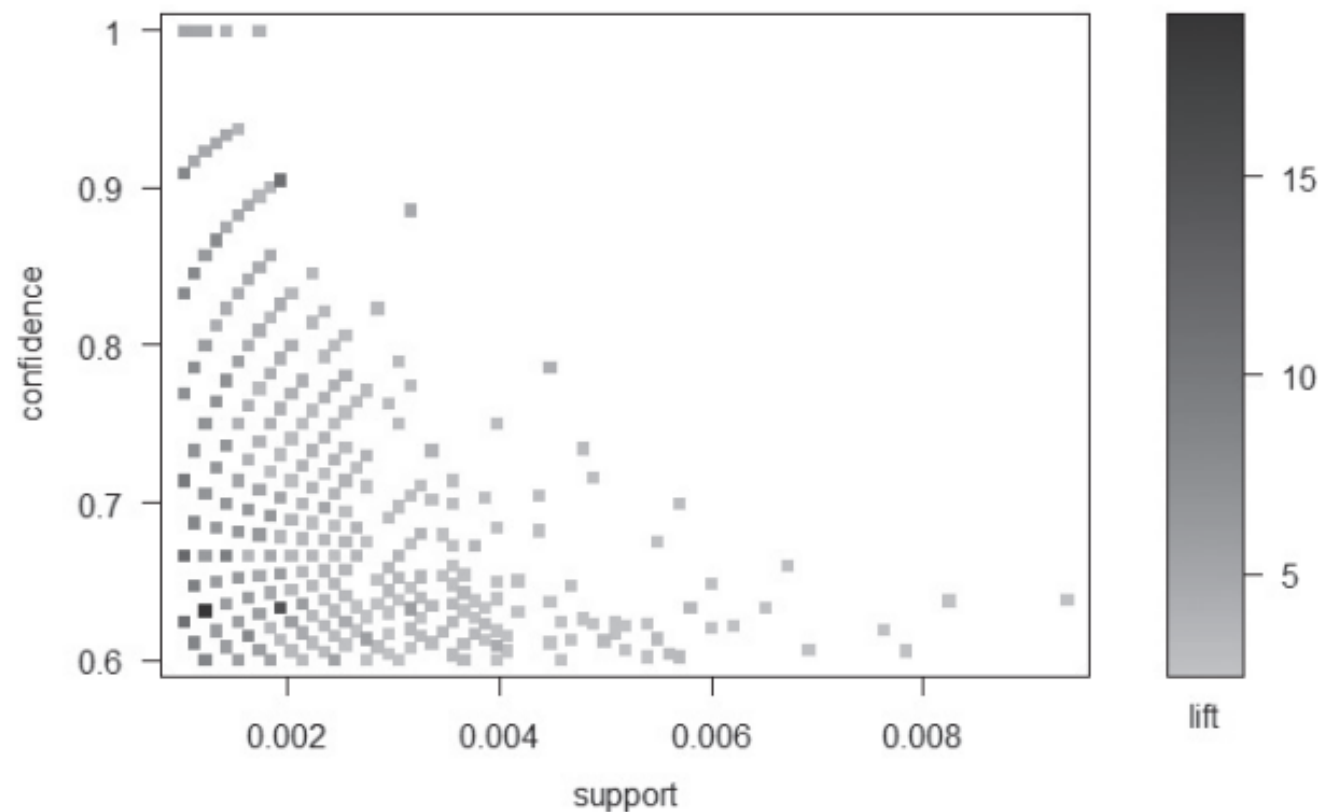
| data | ntransactions | support | confidence |
|-----------|---------------|---------|------------|
| Groceries | 9835 | 0.001 | 0.6 |

Enter `plot(rules)` to display the scatterplot of the 2,918 rules (Figure 5-3), where the horizontal axis is the support, the vertical axis is the confidence, and the shading is the lift. The scatterplot shows that, of the 2,918 rules generated from the *Groceries* dataset, the highest lift occurs at a low support and a low confidence.

Scatter plot for 2918 rules



Scatter plot for 2918 rules



RULE GENERATION AND VISUALIZATION

The `inspect()` function can display content of the rules generated previously. The following code shows the top ten rules sorted by the lift. Rule `{Instant food products, soda} → {hamburger meat}` has the highest lift of 18.995654.

```
inspect(head(sort(rules, by="lift"), 10))
```

| | lhs | | rhs |
|---|----------------------------------|------------|---------------------|
| | support | confidence | lift |
| 1 | {Instant food products, soda} | | => {hamburger meat} |
| | 0.001220132 | 0.6315789 | 18.995654 |
| 2 | {soda, popcorn} | | => {salty snack} |
| | 0.001220132 | 0.6315789 | 16.697793 |
| 3 | {ham, processed cheese} | | => {white bread} |
| | 0.001931876 | 0.6333333 | 15.045491 |

RULE GENERATION AND VISUALIZATION

```
4 {tropical fruit,  
  other vegetables,  
  yogurt,  
  white bread}      => {butter}  
0.001016777  0.6666667 12.030581  
5 {hamburger meat,  
  yogurt,  
  whipped/sour cream}  => {butter}  
0.001016777  0.6250000 11.278670  
6 {tropical fruit,  
  other vegetables,  
  whole milk,  
  yogurt,  
  domestic eggs}      => {butter}  
0.001016777  0.6250000 11.278670  
7 {liquor,  
  red/blush wine}      => {bottled beer}  
0.001931876  0.9047619 11.235269  
8 {other vegetables,  
  butter,  
  sugar}              => {whipped/sour cream}  
0.001016777  0.7142857  9.964539  
9 {whole milk,  
  butter,  
  hard cheese}         => {whipped/sour cream}  
0.001423488  0.6666667  9.300236  
10 {tropical fruit,  
    other vegetables,  
    butter,  
    fruit/vegetable juice} => {whipped/sour cream}  
0.001016777  0.6666667  9.300236
```

RULE GENERATION AND VISUALIZATION

The following code fetches a total of 127 rules whose confidence is above 0.9:

```
confidentRules <- rules[quality(rules)$confidence > 0.9]
confidentRules
set of 127 rules
```

The next command produces a matrix-based visualization (Figure 5-5) of the LHS versus the RHS of the rules. The legend on the right is a color matrix indicating the lift and the confidence to which each square in the main matrix corresponds.

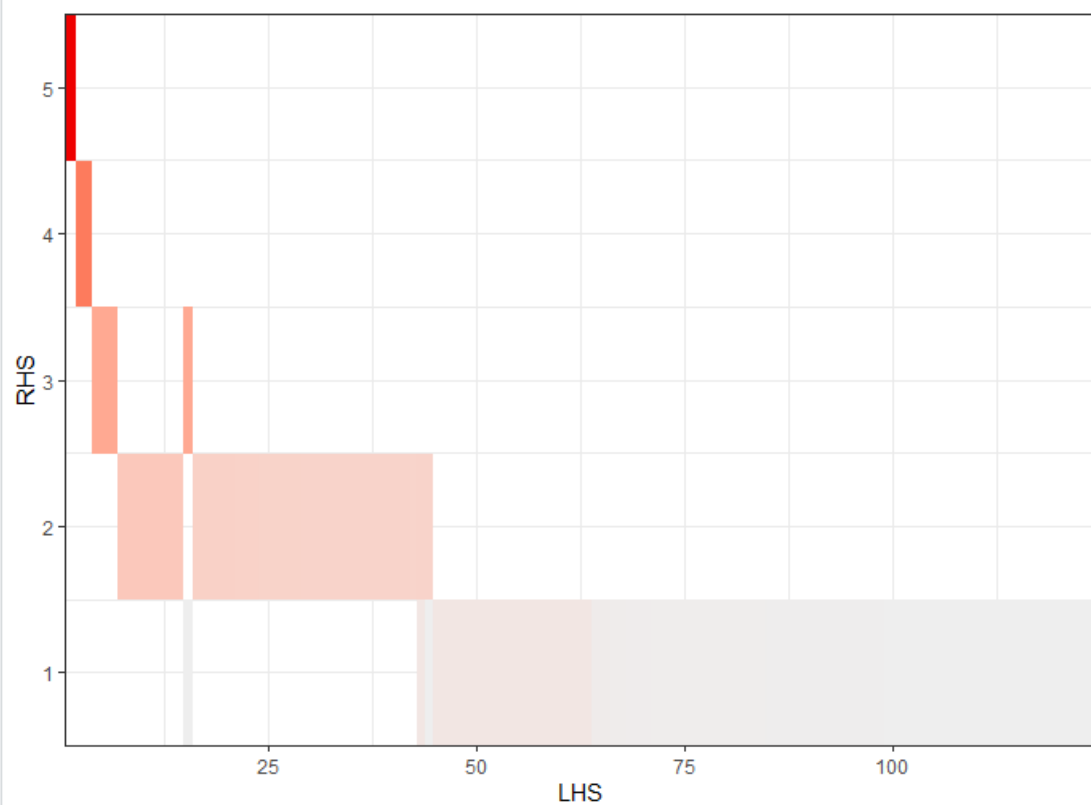
```
plot(confidentRules, method="matrix", measure=c("lift", "confidence"),
control=list(recorder=TRUE))
```

As the previous `plot()` command runs, the R console would simultaneously display a distinct list of the LHS and RHS from the 127 rules. A segment of the output is shown here:

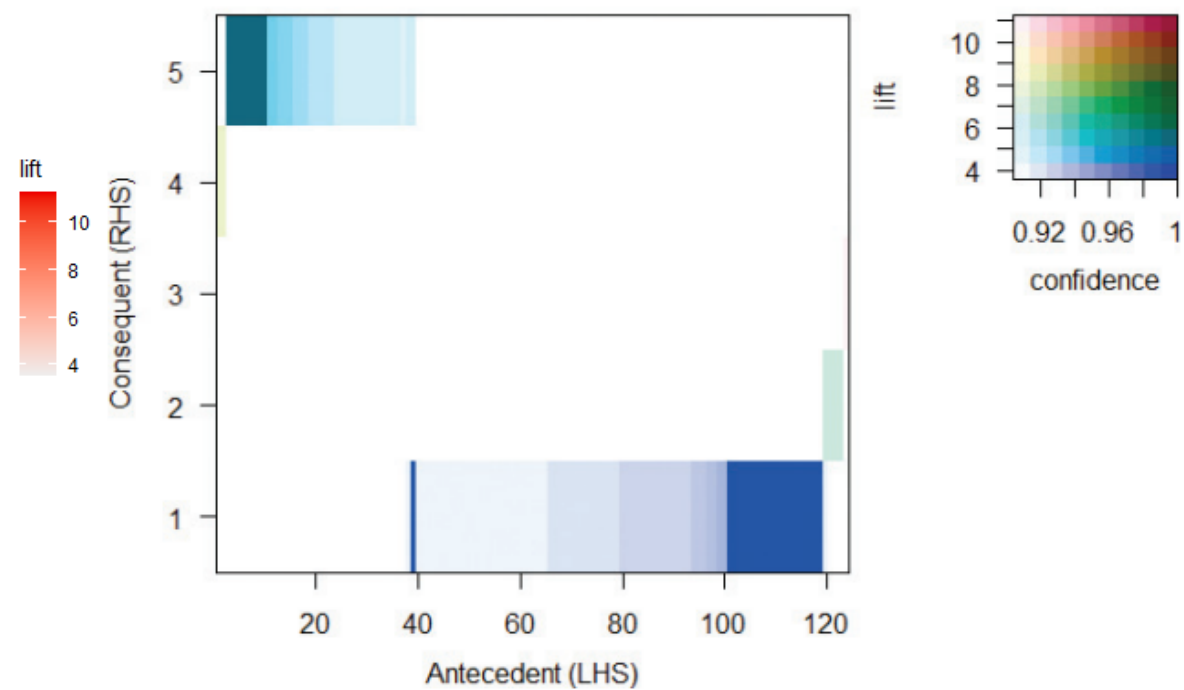
```
Itemsets in Antecedent (LHS)
[1] "{citrus fruit,other vegetables,soda,fruit/vegetable juice}"
[2] "{tropical fruit,other vegetables,whole milk,yogurt,oil}"
```

RULE GENERATION AND VISUALIZATION

Matrix for 127 rules



Matrix with 127 rules



RULE GENERATION AND VISUALIZATION

```
> plot(confidentRules, method="matrix", measure=c("lift","confidence"))
```

Itemsets in Antecedent (LHS)

```
[1] "{liquor,red/blush wine}"
[2] "{citrus fruit,other vegetables,soda,fruit/vegetable juice}"
[3] "{tropical fruit,other vegetables,whole milk,yogurt,oil}"
[4] "{tropical fruit,whole milk,butter,sliced cheese}"
[5] "{other vegetables,curd,whipped/sour cream,cream cheese }"
[6] "{tropical fruit,other vegetables,butter,white bread}"
[7] "{citrus fruit,root vegetables,soft cheese}"
[8] "{pip fruit,whipped/sour cream,brown bread}"
[9] "{tropical fruit,grapes,whole milk,yogurt}"
[10] "{ham,tropical fruit,pip fruit,yogurt}"
[11] "{ham,tropical fruit,pip fruit,whole milk}"
[12] "{tropical fruit,butter,whipped/sour cream,fruit/vegetable juice}"
[13] "{whole milk,rolls/buns,soda,newspapers}"
[14] "{citrus fruit,tropical fruit,root vegetables,whipped/sour cream}"
```

```
[118] "{other vegetables,butter,whipped/sour cream,napkins}"
[119] "{pork,root vegetables,other vegetables,butter}"
[120] "{root vegetables,other vegetables,rolls/buns,brown bread}"
[121] "{citrus fruit,other vegetables,butter,bottled water}"
[122] "{citrus fruit,tropical fruit,other vegetables,domestic eggs}"
[123] "{tropical fruit,root vegetables,yogurt,pastry}"
[124] "{tropical fruit,other vegetables,butter,yogurt,domestic eggs}"
```

Itemsets in Consequent (RHS)

```
[1] "{whole milk}"           "{other vegetables}" "{yogurt}"           "{root vegetables}"
[5] "{bottled beer}"
```

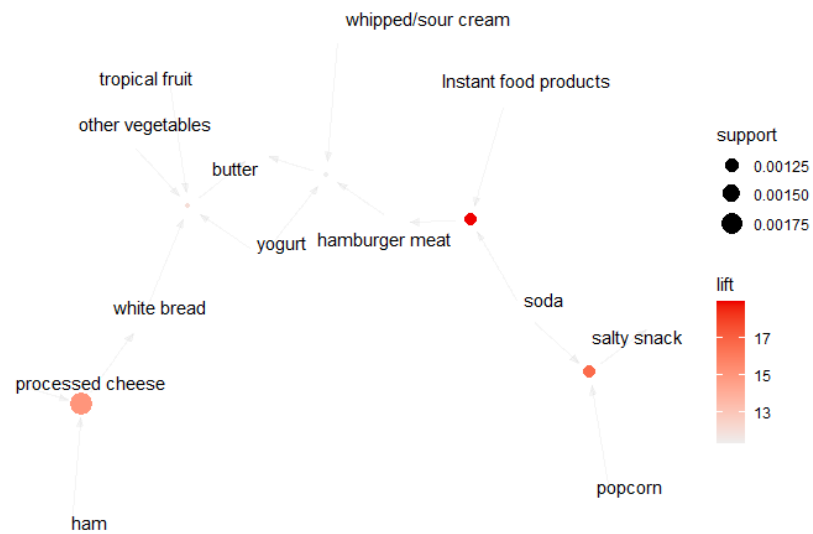
```
> |
```

RULE GENERATION AND VISUALIZATION

The following code provides a visualization of the top five rules with the highest lift. The plot is shown in Figure 5-6. In the graph, the arrow always points from an item on the LHS to an item on the RHS. For example, the arrows that connect ham, processed cheese, and white bread suggest rule $\{\text{ham, processed cheese}\} \rightarrow \{\text{white bread}\}$. The legend on the top right of the graph shows that the size of a circle indicates the support of the rules ranging from 0.001 to 0.002. The color (or shade) represents the lift, which ranges from 11.279 to 18.996. The rule with the highest lift is $\{\text{Instant food products, soda}\} \rightarrow \{\text{hamburger meat}\}$.

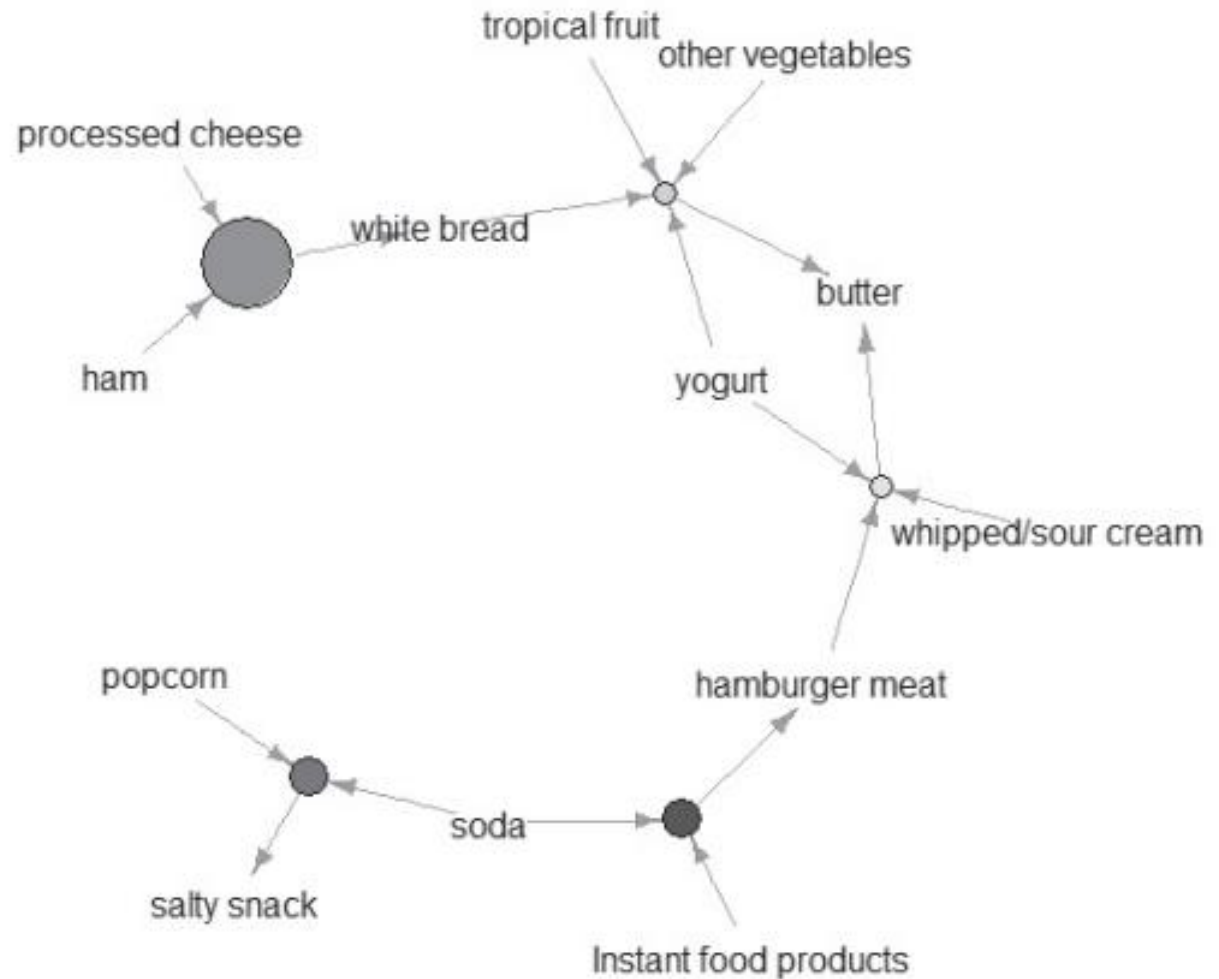
```
highLiftRules <- head(sort(rules, by="lift"), 5)
plot(highLiftRules, method="graph", control=list(type="items"))
```

RULE GENERATION AND VISUALIZATION



Graph for 5 rules

size: support (0.001 - 0.002)
color: lift (11.279 - 18.996)



VALIDATION AND TESTING

- After gathering the output rules, it may become necessary to use one or more methods to validate the results in the business context for the sample dataset. The first approach can be established through statistical measures such as *confidence*, *lift*, and *leverage*.
- Rules that involve mutually independent items or cover few transactions are considered uninteresting because they may capture spurious relationships. As mentioned before, *confidence* measures the chance that X and Y appear together in relation to the chance X appears.
- *Confidence* can be used to identify the interestingness of the rules. *Lift* and *leverage* both compare the support of X and Y against their individual support. While mining data with association rules, some rules generated could be purely coincidental. For example, if 95% of customers buy X and 90% of customers buy Y , then X and Y would occur together at least 85% of the time, even if there is no relationship between the two.
- Measures like *lift* and *leverage* ensure that interesting rules are identified rather than coincidental ones.

VALIDATION AND TESTING

- Another set of criteria can be established through subjective arguments. Even with a high confidence, a rule may be considered subjectively uninteresting unless it reveals any unexpected profitable actions. For example, rules like $\{paper\} \rightarrow \{pencil\}$ may not be subjectively interesting or meaningful despite high support and confidence values.
- In contrast, a rule like $\{diaper\} \rightarrow \{beer\}$ that satisfies both minimum support and minimum confidence can be considered subjectively interesting because this rule unexpected and may suggest a **cross-sell opportunity** for the retailer.
- This incorporation of subjective knowledge into the evaluation of rules can be a difficult task, and it requires collaboration with domain experts. The domain experts may serve as the business users or the business intelligence analysts as part of the Data Science team. Then, the team can communicate the results and decide if it is appropriate to operationalize them.

DIAGNOSTICS

- Although the *Apriori* algorithm is easy to understand and implement, some of the rules generated are uninteresting or practically useless. Additionally, some of the rules may be generated due to coincidental relationships between the variables.
- Measures like *confidence*, *lift*, and *leverage* should be used along with human insights to address this problem. Another problem with association rules is that, in phases of the Data Analytics Lifecycle, the team must specify the minimum support prior to the model execution, which may lead to too many or too few rules. In related research, a variant of the algorithm can use a predefined target range for the number of rules so that the algorithm can adjust the minimum support accordingly.
- The *Apriori* algorithm reduces the computational workload by only examining itemsets that meet the specified minimum threshold. However, depending on the size of the dataset, the *Apriori* algorithm can be computationally expensive. For each level of support, the algorithm requires a scan of the entire database to obtain the result.

DIAGNOSTICS

Accordingly, as the database grows, it takes more time to compute in each run. Here are some approaches to improve *Apriori*'s efficiency:

- **Partitioning:** Any itemset that is potentially frequent in a transaction database must be frequent in at least one of the partitions of the transaction database.
- **Sampling:** This extracts a subset of the data with a lower support threshold and uses the subset to perform association rule mining.
- **Transaction reduction:** A transaction that does not contain frequent k -itemsets is useless in subsequent scans and therefore can be ignored.
- **Hash-based itemset counting:** If the corresponding hashing bucket count of a k -itemset is below a certain threshold, the k -itemset cannot be frequent.
- **Dynamic itemset counting:** Only add new candidate itemsets when all of their subsets are estimated to be frequent.

```
1 library(arules)
2 library(arulesViz)
3 library(recommenderlab)
4 library(Matrix)
5
6 m <- matrix(sample(c(0,1),10000,replace=TRUE,prob=c(0.4,0.6)),nrow=2500, ncol=4,dimnames=list(users=paste("u",1:2500,sep=' '),items=paste("i",1:4,sep=' ')))
7 head(m)
8
9 b <- as(m,"binaryRatingMatrix")
10 b
11 as(b,"matrix")
12
13 count <- as.data.frame(rowCounts(b))
14 head(count)
15 sum(count)
16 colCounts(b)
17
18 as(b,"data.frame")
19 rep<- getData.frame(b,ratings=FALSE)
20 head(rep)
21 aa <- as(b,"list")
22 head(aa)
23
24 transaction1 <- as(aa,"transactions")
25 inspect(transaction1)
26
27 AllRules <- apriori(transaction1,parameter=list(supp=0.5,conf=0.5,target="rules"))
28
29 q(AllRules) <- round(quality(AllRules),digits=3)
30 SortedRules <- sort(AllRules, by="lift")
31 inspect(SortedRules)
```

MORE PRACTICE

```
> head(m)
```

| | items | | | |
|-------|-------|----|----|----|
| users | i1 | i2 | i3 | i4 |
| u1 | 1 | 1 | 0 | 0 |
| u2 | 0 | 0 | 1 | 0 |
| u3 | 1 | 0 | 1 | 1 |
| u4 | 0 | 1 | 0 | 1 |
| u5 | 1 | 1 | 1 | 0 |
| u6 | 1 | 0 | 1 | 1 |

```
> as(b,"matrix")
```

| | i1 | i2 | i3 | i4 |
|-----|-------|-------|-------|-------|
| u1 | TRUE | TRUE | FALSE | FALSE |
| u2 | FALSE | FALSE | TRUE | FALSE |
| u3 | TRUE | FALSE | TRUE | TRUE |
| u4 | FALSE | TRUE | FALSE | TRUE |
| u5 | TRUE | TRUE | TRUE | FALSE |
| u6 | TRUE | FALSE | TRUE | TRUE |
| u7 | FALSE | FALSE | FALSE | FALSE |
| u8 | TRUE | FALSE | TRUE | TRUE |
| u9 | TRUE | FALSE | FALSE | FALSE |
| u10 | FALSE | TRUE | TRUE | TRUE |
| u11 | FALSE | TRUE | FALSE | FALSE |
| u12 | TRUE | FALSE | FALSE | TRUE |
| u13 | FALSE | FALSE | TRUE | FALSE |
| u14 | TRUE | FALSE | TRUE | TRUE |
| u15 | TRUE | TRUE | TRUE | FALSE |

```
> head(count)
```

| | rowCounts(b) |
|----|--------------|
| u1 | 2 |
| u2 | 1 |
| u3 | 3 |
| u4 | 2 |
| u5 | 3 |
| u6 | 3 |

```
> sum(count)
```

```
[1] 5981
```

```
> colCounts(b)
```


| | i1 | i2 | i3 | i4 |
|--|------|------|------|------|
| | 1481 | 1473 | 1495 | 1532 |

```
> |
```

MORE PRACTICE

```
> as(b,"data.frame")
```

| | user | item | rating |
|------|-------|------|--------|
| 1 | u1 | i1 | 1 |
| 1482 | u1 | i2 | 1 |
| 1485 | u10 | i2 | 1 |
| 2960 | u10 | i3 | 1 |
| 4454 | u10 | i4 | 1 |
| 60 | u100 | i1 | 1 |
| 3013 | u100 | i3 | 1 |
| 4510 | u100 | i4 | 1 |
| 2068 | u1000 | i2 | 1 |
| 5063 | u1000 | i4 | 1 |
| 603 | u1001 | i1 | 1 |
| 2069 | u1001 | i2 | 1 |
| 3552 | u1001 | i3 | 1 |



```
> aa <- as(b,"list")
> head(aa)
```


```
$u1
[1] "i1" "i2"
```

```
$u2
[1] "i3"
```

```
$u3
[1] "i1" "i3" "i4"
```

```
$u4
[1] "i2" "i4"
```

```
$u5
[1] "i1" "i2" "i3"
```



```
> transaction1 <- as(aa,"transactions")
> inspect(transaction1)
```

| | items | transactionID |
|------|--------------|---------------|
| [1] | {i1, i2} | u1 |
| [2] | {i3} | u2 |
| [3] | {i1, i3, i4} | u3 |
| [4] | {i2, i4} | u4 |
| [5] | {i1, i2, i3} | u5 |
| [6] | {i1, i3, i4} | u6 |
| [7] | {} | u7 |
| [8] | {i1, i3, i4} | u8 |
| [9] | {i1} | u9 |
| [10] | {i2, i3, i4} | u10 |
| [11] | {i2} | u11 |
| [12] | {i1, i4} | u12 |
| [13] | {i3} | u13 |
| [14] | {i1, i3, i4} | u14 |

MORE PRACTICE

```
> AllRules <- apriori(transaction1,parameter=list(supp=0.5,conf=0.5,target="rules"))
```

```
Apriori
```

```
Parameter specification:
```

| confidence | minval | smax | arem | aval | originalSupport | maxtime | support | minlen | maxlen | target | ext |
|------------|--------|------|------|-------|-----------------|---------|---------|--------|--------|--------|------|
| 0.5 | 0.1 | 1 | none | FALSE | TRUE | 5 | 0.5 | 1 | 10 | rules | TRUE |

```
Algorithmic control:
```

| filter | tree | heap | memopt | load | sort | verbose |
|--------|------|------|--------|------|------|---------|
| 0.1 | TRUE | TRUE | FALSE | TRUE | 2 | TRUE |

```
Absolute minimum support count: 1250
```

```
set item appearances ...[0 item(s)] done [0.00s].  
set transactions ...[4 item(s), 2500 transaction(s)] done [0.00s].  
sorting and recoding items ... [4 item(s)] done [0.00s].  
creating transaction tree ... done [0.00s].  
checking subsets of size 1 2 done [0.00s].  
writing ... [4 rule(s)] done [0.00s].  
creating S4 object ... done [0.00s].
```

```
> |
```

REFERENCE

- Data Science and Big Data analytics, Wiley, 2015. Chapter 5: Advanced analytical theory and methods: Association Rules.