# ASSOCIATION RULES

Market Basket Analyses are a common application of association rules.  One goal of a market basket analysis is to understand the association between items purchases.  The relationship between items purchased at a grocery store will be considered in this handout.



An association rule highlights the fact that some items are more (or less) indicative of the purchase of others.  For example, purchasing cereal increases the likelihood of purchasing milk.  These types of analyses may also reveals that liquor and milk are rarely purchased together.

| Rule | Item | # | # Milk + ___ | Confidence | Lift |
|---|---|---|---|---|---|
| {kitchen utensil} -> {Milk} | kitchen utensil | 4 | 3 | 0.750 | 2.935 |
| {honey} -> {Milk} | honey | 15 | 11 | 0.733 | 2.870 |
| {cereals} -> {Milk} | cereals | 56 | 36 | 0.643 | 2.516 |
| {rice} -> {Milk} | rice | 75 | 46 | 0.613 | 2.400 |
| {rubbing alcohol} -> {Milk} | rubbing alcohol | 10 | 6 | 0.600 | 2.348 |
| {cocoa drinks} -> {Milk} | cocoa drinks | 22 | 13 | 0.591 | 2.313 |
| {pudding powder} -> {Milk} | pudding powder | 23 | 13 | 0.565 | 2.212 |
| {jam} -> {Milk} | jam | 53 | 29 | 0.547 | 2.142 |
| {baking powder} -> {Milk} | baking powder | 174 | 91 | 0.523 | 2.047 |
| {cooking chocolate} -> {Milk} | cooking chocolate | 25 | 13 | 0.520 | 2.035 |
| {preservation products} -> {Milk} | preservation products | 2 | 1 | 0.500 | 1.957 |
| {baby cosmetics} -> {Milk} | baby cosmetics | 6 | 3 | 0.500 | 1.957 |
| {butter} -> {Milk} | butter | 545 | 271 | 0.497 | 1.946 |
| | | | | | |
| {candy} -> {Milk} | candy | 294 | 81 | 0.276 | 1.078 |
| {hair spray} -> {Milk} | hair spray | 11 | 3 | 0.273 | 1.067 |
| {seasonal products} -> {Milk} | seasonal products | 140 | 37 | 0.264 | 1.034 |
| {specialty chocolate} -> {Milk} | specialty chocolate | 299 | 79 | 0.264 | 1.034 |
| {photo/film} -> {Milk} | photo/film | 91 | 23 | 0.253 | 0.989 |
| {frozen fruits} -> {Milk} | frozen fruits | 12 | 3 | 0.250 | 0.978 |
| | | | | | |
| {canned beer} -> {Milk} | canned beer | 764 | 87 | 0.114 | 0.446 |
| {liquor} -> {Milk} | liquor | 109 | 6 | 0.055 | 0.215 |
| {baby food} -> {Milk} | baby food | 1 | 0 | 0.000 | 0.000 |

**Association Rules** are used to uncover associations or relationships that exist between items. Often these rules are constructed to identify relationships between items purchased, i.e. Market Basket Analysis.

Procedural Steps

1.    Determine how often items are purchased
2.    Determine how often items are purchased in conjunction with other items
3.    Identify which purchased items are indicative of others being purchased

Data Technologies

1.    Filtering in Excel
2.    Creating Tables in Excel
3.    Applications of Functions in Excel

Consider the following subset of data from a collection of transactions from a grocery store.

| Transaction ID | Items Purchased |
|---|---|
| 1 | {Bread, Milk} |
| 2 | {Eggs, Ham} |
| 3 | {Bread, Fruit, Milk} |
| 4 | {Beer, Bread, Butter, Fruit, Soda} |
| 5 | {Bread, Fruit, Milk, Soda} |

Association rules are developed under the following guiding principles.

| | | |
|---|---|---|
| 1. | Items should be purchased somewhat often | **Support** |
| 2. | Reliability, i.e. the degree to which one set of items predicts the purchase of another set of items | **Confidence** |

Consider the following association rule – the purchase of Butter indicates the purchase of Milk.

| Rule #1 | $\{Bread\} \rightarrow \{Milk\}$ |
|---------|----------------------------------|

Compute the support and confidence for this rule.

$$Support(Bread\ AND\ Milk) = \frac{\#\ Bread\ AND\ Milk}{\#\ Transactions} =$$

$$Confidence\ of\ Rule\ \#1 = \frac{Support(Bread\ AND\ Milk)}{Support(Bread)} =$$

Questions

1. What is the interpretation of the Support(Bread AND Milk)?

2. What is the interpretation of Confidence of this rule?  Discuss.
   Note:  Confidence is simply a conditional probability, i.e P(Milk | Bread).

Consider a second association rule for the purchase of Milk.

| Rule #2 | $\{Fruit\} \rightarrow \{Milk\}$ |
|---------|----------------------------------|

Compute the support and confidence for this rule.

$$Support(Fruit\ AND\ Milk) =$$

$$Confidence\ of\ Rule =$$

Question

3. Why might Rule #1 be considered "better" than Rule #2 when interest lies in the purchase of Milk?

Consider a third association rule for the purchase of Milk.

| Rule #3 | $\{Bread, Fruit\} \rightarrow \{Milk\}$ |
|---------|------------------------------------------|

Compute the support and confidence for this rule.

$Support(Bread, Fruit, AND\ Milk) =$

$Confidence\ of\ Rule =$

**Lift** is another measure often considered when evaluating rules of association.

$$Lift(\{Bread\} \rightarrow \{Milk\}) = \frac{Confidence(Bread\ AND\ Milk)}{Support(Milk)} = \frac{P(Milk|Bread)}{P(Milk)}$$

For our example, realize that the support for Milk is fairly large. i.e, Milk was purchased in 60% of the transactions. This provides a baseline value for confidence. That is, rules that exceed this value indicate gains when considering the association provided by the rule. When the lift of a rule is near 1, then the rule provides little information to understanding the purchase of the item.

- $Lift > 1$ implies positive association between items
- $Lift \approx 1$ implies no association between items
- $Lift < 1$ implies negative association between items

| Rule | Support | Confidence | Lift |
|------|---------|-----------|------|
| $\{Bread\} \rightarrow \{Milk\}$ | $\dfrac{3}{5}$ | $\dfrac{^3/_5}{^4/_5} = \dfrac{3}{4}$ | $\dfrac{^3/_4}{^3/_5} = 1.25$ |
| $\{Fruit\} \rightarrow \{Milk\}$ | $\dfrac{2}{5}$ | $\dfrac{^2/_5}{^3/_5} = \dfrac{2}{3}$ | $\dfrac{^2/_3}{^3/_5} = 1.11$ |
| $\{Bread, Fruit\} \rightarrow \{Milk\}$ | $\dfrac{2}{5}$ | $\dfrac{^2/_5}{^3/_5} = \dfrac{2}{3}$ | $\dfrac{^2/_3}{^3/_5} = 1.11$ |

Some Comments

- Association rules with no support have zero confidence. E.g. Beer is never purchased with Milk, so the rule $\{Beer\} \rightarrow \{Milk\}$ should not be considered.

- The confidence of a rule should not be considered independent of it's support. For example, the rule $\{Eggs\} \rightarrow \{Ham\}$ has Confidence = 1. That is, 100% of the time eggs were purchased, so was Ham. However, this rule has very low support as Eggs and Ham were only purchased once.

- Association rules are not invariant. For example, the confidence for the rule $\{Bread\} \rightarrow \{Milk\}$ is different than the confidence of the rule $\{Milk\} \rightarrow \{Bread\}$.

Common Data Structure

List

| Transaction ID | Items Purchased |
|---|---|
| 1 | {Bread, Milk} |
| 2 | {Eggs, Ham} |
| 3 | {Bread, Fruit, Milk} |
| 4 | {Beer, Bread, Butter, Fruit, Soda} |
| 5 | {Bread, Fruit, Milk, Soda} |

$\rightarrow$

Binary Representation (Matrix)

| ID | Beer | Bread | Butter | Eggs | Fruit | Ham | Milk | Soda |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 2 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 |
| 4 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 |
| 5 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 1 |

Next, consider the complete grocery dataset. This dataset contains 9835 transactions and 169 unique items. This dataset can be downloaded from the Workshop website.

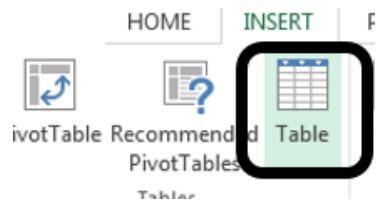| Data Source | |
|---|---|
| Address | http://course1.winona.edu/cmalone/workshops/uscots2015/ |
| Description | Groceries Dataset<br>Michael Hahsler, Kurt Hornik, and Thomas Reutterer (2006) Implications of probabilistic data modeling for mining association rules. In M. Spiliopoulou, R. Kruse, C. Borgelt, A. Nuernberger, and W. Gaul, editors, From Data and Information Analysis to Knowledge Engineering, Studies in Classification, Data Analysis, and Knowledge Organization, pages 598–605. Springer-Verlag. |

Open the Groceries dataset in Excel.  The binary representation of this market basket dataset is provided in this Excel file.  A snippet is shown here.

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | frankfurter | sausage | liver loaf | ham | meat | finished products | organic sausag | chicken | turkey | pork | beef | hamburger meat | fish | citrus fruit | tropical fruit | pip fruit | grapes | berries |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 15 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Spreadsheets consist of rows and columns.  Datasets also consist of rows and columns as well, but also contain information that is not data, e.g. variable names.  Excel does not differentiate the header row from actual data unless you convert the collection of rows and columns into a **Table**.
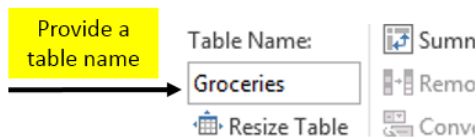
Putting Data into the Table structure in Excel

HOME    INSERT    F

ivotTable Recommended  Table
          PivotTables
          Tables

Give your table a name for easy referencing

Provide a table name

Table Name:     Summ
Groceries        Remo
Resize Table    Conv

Short-cuts for cursor movement in Excel

| Short-cuts in Excel | |
|---|---|
| Ctrl Home | Upper Left Corner |
| Ctrl End | Lower Right Corner |
| Ctrl ← | Move to Left edge |
| Ctrl → | Move to Right edge |
| Ctrl ↑ | Move to top |
| Ctrl ↓ | Move to bottom |

The following snippet shows the Groceries dataset specified as a table.

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | frankfurter | sausage | liver loaf | ham | meat | finished products | organic sausage | chicken | turkey | pork | beef |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 15 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

The drop-down arrows provided for each variable (or field) are called Filters.  Filters in Excel allow you to subset rows.

Filter on Whole Milk               Select Whole Milk = 1 to identify transactions
that purchased whole milk



After a Filter is applied, certain rows are hidden from view.  Excel indicates this fact with changing the row label color to blue.



The status bar in Excel, the bar across the bottom of the Excel file, provides simple summaries for columns of the table.  For example, if the Whole Milk column (column Y) is highlighted, the following summaries are shown.
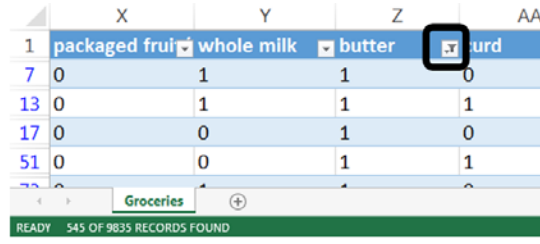
Applying filters to columns Whole Milk and Butter allows one to easy compute the support and confidence for the rule $\{Butter\} \to \{Milk\}$.



| Rule | $\{Butter\} \to \{Milk\}$ |
|------|---------------------------|

Support(Milk AND Butter) = 271/9835

Support(Butter) = 545/9835

- $Support(Butter\ AND\ Milk) = \frac{\#\ Butter\ AND\ Milk}{\#\ Transactions} = \frac{271}{9835} = 0.028$

- $Confidence = \frac{Support(Butter\ AND\ Milk)}{Support(Butter)} = \frac{271/9835}{545/9835} = \frac{271}{545} = 0.497$

- $Lift = \frac{Confidence}{Support(Milk)} = \frac{0.497}{2513/9835} = \frac{0.497}{0.256} = 1.946$

The =COUNT() function in Excel can used to count the number of nonblank rows in a column. Excel functions also work with tables and variable names. The following will provide a count of the number of transactions in the Groceries dataset, i.e. 9835. The use of the table and variable names is preferred as this avoids the need to highlight an exact range of cells in Excel.

| Counting the number of rows in table | |
|---|---|
| Using Range | =COUNT(Y:Y) |
| Using Table | =COUNT(Groceries[whole milk]) |

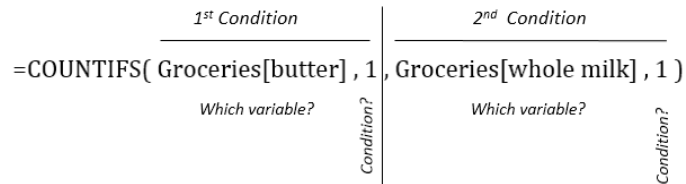The =COUNTIF() function in Excel provides a count of only the cells that satisfy some condition. The following can be used to compute the support for butter.

=COUNTIF( Groceries[butter] , 1 )

If more than one condition is needed, the =COUNTIFS() function can be used. COUNTIFS is necessary to compute Support(Butter AND Whole Milk).

=COUNTIFS( Groceries[butter] , 1 , Groceries[whole milk] , 1 )

8

A brief description of the COUNTIFS function in Excel is provide here.

|  | 1st Condition | | 2nd Condition | |
| --- | --- | --- | --- | --- |
| =COUNTIFS( | Groceries[butter] , 1 | , | Groceries[whole milk] , 1 | ) |
|  | Which variable? | Condition? | Which variable? | Condition? |

Move to far right of the Groceries table in Excel. You can use Ctrl → to move quickly to the far right edge. Enter the following function in Excel to compute the counts necessary for measuring support for the rule $\{Butter\} \rightarrow \{Milk\}$.

|  | FL | FM | FN | FO | FP |
| --- | --- | --- | --- | --- | --- |
| 1 | shop | bags | | | |
| 2 | 0 | 0 | | | Automating Counting in Excel |
| 3 | 0 | 0 | | Number of Transactions | =COUNT(Groceries[whole milk]) |
| 4 | 0 | 0 | | # (Butter) | =COUNTIF(Groceries[butter],1) |
| 5 | 0 | 0 | | | |
| 6 | 0 | 0 | | # (Butter AND Whole Milk) | =COUNTIFS(Groceries[butter],1,Groceries[whole milk],1) |
| 7 | 0 | 0 | | | |
| 8 | 0 | 0 | | | |

Use the value computed above to compute the Confidence and Lift for this rule.

|  | FL | FM | FN | FO | FP |
| --- | --- | --- | --- | --- | --- |
| 1 | shop | bags | | | |
| 2 | 0 | 0 | | | Automating Counting in Excel |
| 3 | 0 | 0 | | Number of Transactions | =COUNT(Groceries[whole milk]) |
| 4 | 0 | 0 | | # (Butter) | =COUNTIF(Groceries[butter],1) |
| 5 | 0 | 0 | | | |
| 6 | 0 | 0 | | # (Butter AND Whole Milk) | =COUNTIFS(Groceries[butter],1,Groceries[whole milk],1) |
| 7 | 0 | 0 | | | |
| 8 | 0 | 0 | | **Confidence** | =  FP6 /  FP4 |
| 9 | 0 | 0 | | | |
| 10 | 0 | 0 | | # (Milk) | =COUNTIF(Groceries[whole milk],1) |
| 11 | 0 | 0 | | **Lift** | = FP8 / ( FP10 / FP3 ) |
| 12 | 0 | 0 | | | |

Verify that these formulas are correct by comparing them to the output provided below.

|  | FL | FM | FN | FO | FP |
| --- | --- | --- | --- | --- | --- |
| 1 | s | ba | | | |
| 2 | 0 | 0 | | Automating Counting in Excel | |
| 3 | 0 | 0 | | Number of Transactions | 9835 |
| 4 | 0 | 0 | | # (Butter) | 545 |
| 5 | 0 | 0 | | | |
| 6 | 0 | 0 | | # (Butter AND Whole Milk) | 271 |
| 7 | 0 | 0 | | | |
| 8 | 0 | 0 | | **Confidence** | 0.497 |
| 9 | 0 | 0 | | | |
| 10 | 0 | 0 | | # (Milk) | 2513 |
| 11 | 0 | 0 | | **Lift** | 1.946 |

Evaluating Several Rules

The procedure provided above lack efficiencies and does not scale well when several rules need to be evaluated.  For example, to evaluate the rule $\{Yogurt\} \rightarrow \{Milk\}$, the formulas for support will need to be changed.   The =INDIRCT() function in Excel will help increase the efficiency in computing the support, confidence, and lift for several rules.

---

**INDIRECT() Function**

Consider the following data in Excel.

| | A | B |
|---|---|---|
| 1 | Cell ID | Value |
| 2 | B3 | Cell B2 |
| 3 | B4 | 2 |
| 4 | | Text in B4 |

Step 1: Obtain value from another cell

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Cell ID | Value | | |
| 2 | B3 | Cell B2 | | |
| 3 | B4 | 2 | | = INDIRECT(A3) |
| 4 | | Text in B4 | | |

Step 2: Use value in specified cell in formula

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Cell ID | Value | | |
| 2 | B3 | Cell B2 | | |
| 3 | B4 | 2 | | = ~~INDIRECT(A3)~~ B4 |
| 4 | | Text in B4 | | |

---

The =INDIRECT() function can be used in the following manner to automatically update the variable names when computing the support for several rules.

$$=\text{COUNTIF( INDIRECT ( " Groceries[ " \& B2 \& " ] " ), 1)}$$

The following setup is used to evaluate six different association rules for Milk.

| | A | B | C |
|---|---|---|---|
| 1 | Rule | Item | # |
| 2 | {butter} -> {Milk} | butter | =COUNTIF(INDIRECT("Groceries["&B2&"]"),1) |
| 3 | {yogurt} -> {Milk} | yogurt | |
| 4 | {whipped/sour cream} -> {Milk} | whipped/sour cream | |
| 5 | {cereals} -> {Milk} | cereals | |
| 6 | {canned beer} -> {Milk} | canned beer | |
| 7 | {make up remover} -> {Milk} | make up remover | |

This formula can be copied down in Excel to evaluate the support for the remaining rules.   The confidence and lift are computed for these rules as well.

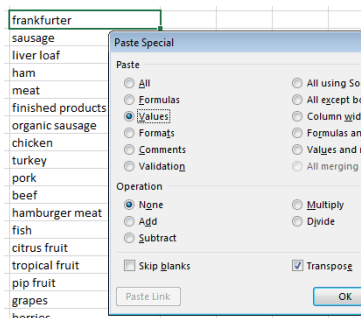| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| 1 | Rule | Item | # | # Milk + ___ | Confidence | | Lift |
| 2 | {butter} -> {Milk} | butter | 545 | 271 | 0.497 | | 1.946 |
| 3 | {yogurt} -> {Milk} | yogurt | 1372 | 551 | 0.402 | | 1.572 |
| 4 | {whipped/sour cream} -> {Milk} | whipped/sour cream | 705 | 317 | 0.450 | | 1.760 |
| 5 | {cereals} -> {Milk} | cereals | 56 | 36 | 0.643 | | 2.516 |
| 6 | {canned beer} -> {Milk} | canned beer | 764 | 87 | 0.114 | | 0.446 |
| 7 | {make up remover} -> {Milk} | make up remover | 8 | 2 | 0.250 | | 0.978 |

Questions

4. The Lift for $\{Cereal\} \rightarrow \{Milk\}$ is about 2.5 which is fairly high. Thus, given that the transaction includes cereal, there is 2.5 fold increase in the likelihood of milk being purchased.

   a. Compute Support(Cereal AND Milk).

   b. This value is fairly low. Why does a low support value negate the usefulness of a rule?

5. The Lift value for the rule $\{Canned\ Beer\} \rightarrow \{Milk\}$ is lowest on this list. What can be said about the purchase of Canned Beer AND Milk?

6. Which of these rules is least useful in the prediction of Milk? Explain how you made this determination.

Task

Use Excel to obtain the Confidence and Lift for all association rules for Whole Milk where only single items are considered on the left.

- Copy all variable names and paste them into a single column. This can be done using Paste Special – specify Values and Transpose when pasting.
- The =CONCATENATE() function can be used to create the Rule column, i.e. =CONCATENATE("{",B2,"} -> {Milk}").

Specify Values and Transpose under Paste Special

Output for Rules

| Rule | Item | # | # Milk + ___ | Confidence | Lift |
|------|------|---|------|------------|------|
| {kitchen utensil} -> {Milk} | kitchen utensil | 4 | 3 | 0.750 | 2.935 |
| {honey} -> {Milk} | honey | 15 | 11 | 0.733 | 2.870 |
| {cereals} -> {Milk} | cereals | 56 | 36 | 0.643 | 2.516 |
| {rice} -> {Milk} | rice | 75 | 46 | 0.613 | 2.400 |
| {rubbing alcohol} -> {Milk} | rubbing alcohol | 10 | 6 | 0.600 | 2.348 |
| {cocoa drinks} -> {Milk} | cocoa drinks | 22 | 13 | 0.591 | 2.313 |
| {pudding powder} -> {Milk} | pudding powder | 23 | 13 | 0.565 | 2.212 |
| {jam} -> {Milk} | jam | 53 | 29 | 0.547 | 2.142 |
| {baking powder} -> {Milk} | baking powder | 174 | 91 | 0.523 | 2.047 |
| {cooking chocolate} -> {Milk} | cooking chocolate | 25 | 13 | 0.520 | 2.035 |
| {preservation products} -> {Milk} | preservation products | 2 | 1 | 0.500 | 1.957 |
| {baby cosmetics} -> {Milk} | baby cosmetics | 6 | 3 | 0.500 | 1.957 |
| {butter} -> {Milk} | butter | 545 | 271 | 0.497 | 1.946 |
| | | | | | |
| {candy} -> {Milk} | candy | 294 | 81 | 0.276 | 1.078 |
| {hair spray} -> {Milk} | hair spray | 11 | 3 | 0.273 | 1.067 |
| {seasonal products} -> {Milk} | seasonal products | 140 | 37 | 0.264 | 1.034 |
| {specialty chocolate} -> {Milk} | specialty chocolate | 299 | 79 | 0.264 | 1.034 |
| {photo/film} -> {Milk} | photo/film | 91 | 23 | 0.253 | 0.989 |
| {frozen fruits} -> {Milk} | frozen fruits | 12 | 3 | 0.250 | 0.978 |
| | | | | | |
| {canned beer} -> {Milk} | canned beer | 764 | 87 | 0.114 | 0.446 |
| {liquor} -> {Milk} | liquor | 109 | 6 | 0.055 | 0.215 |
| {baby food} -> {Milk} | baby food | 1 | 0 | 0.000 | 0.000 |

## Association Rules in R

The **arules** package in R can be used to expand upon what we've done in Excel. The following code can be used to recreate the table above. The maxlen=2 specification in the apriori() function restricts rules to single items. The subset() function reduces the rules – here to include only rules for which whole milk is on the right-hand side.

```
#Load arules package
library(arules)

#Forcing Groceries to be transactions object
gr.trans = as(Groceries,"transactions")

#Rule Development via apriori function
gr.rules = apriori(gr.trans,parameter = list(supp = 0.0,conf = 0.0, maxlen=2))

#Rules that have whole milk on right-hand side
gr.subset = subset(gr.rules, subset=rhs %in% "whole milk")

#Print rules to screen - sorted by lift
inspect(sort(gr.subset,by="lift"))

     lhs                        rhs                    support confidence      lift
1    {kitchen utensil}       => {whole milk} 0.0003050330 0.75000000 2.9352368
2    {honey}                 => {whole milk} 0.0011184545 0.73333333 2.8700093
3    {cereals}               => {whole milk} 0.0036603965 0.64285714 2.5159172
4    {rice}                  => {whole milk} 0.0046771734 0.61333333 2.4003714
5    {rubbing alcohol}       => {whole milk} 0.0006100661 0.60000000 2.3481894
6    {cocoa drinks}          => {whole milk} 0.0013218099 0.59090909 2.3126108
7    {pudding powder}        => {whole milk} 0.0013218099 0.56521739 2.2120625
8    {jam}                   => {whole milk} 0.0029486528 0.54716981 2.1414306
9    {baking powder}         => {whole milk} 0.0092526690 0.52298851 2.0467935
10   {cooking chocolate}     => {whole milk} 0.0013218099 0.52000000 2.0350975
11   {preservation products} => {whole milk} 0.0001016777 0.50000000 1.9568245
12   {baby cosmetics}        => {whole milk} 0.0003050330 0.50000000 1.9568245
13   {butter}                => {whole milk} 0.0275546518 0.49724771 1.9460530
```

The following parameter specification in the apiori() function will limit rules with support larger than 0.01 and confidence larger than 0.25.
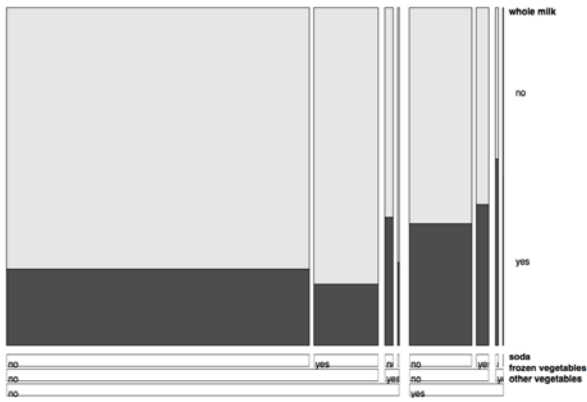
```
#Rule Development via apriori function
gr.rules = apriori(gr.trans,parameter = list(supp = 0.01,conf = 0.25))

#Print rules to screen - sorted by lift
inspect(sort(gr.rules,by="lift"))

     lhs                         rhs                     support confidence      lift
1    {citrus fruit,
      other vegetables}      => {root vegetables}  0.01037112   0.3591549 3.2950455
2    {tropical fruit,
      other vegetables}      => {root vegetables}  0.01230300   0.3427762 3.1447798
3    {beef}                   => {root vegetables}  0.01738688   0.3313953 3.0403668
4    {citrus fruit,
      root vegetables}       => {other vegetables} 0.01037112   0.5862069 3.0296084
5    {tropical fruit,
      root vegetables}       => {other vegetables} 0.01230300   0.5845411 3.0209991
6    {other vegetables,
      whole milk}            => {root vegetables}  0.02318251   0.3097826 2.8420820
7    {whole milk,
      curd}                  => {yogurt}           0.01006609   0.3852140 2.7613555
8    {other vegetables,
      yogurt}                => {root vegetables}  0.01291307   0.2974239 2.7286977
9    {other vegetables,
      yogurt}                => {tropical fruit}   0.01230300   0.2833724 2.7005496
10   {other vegetables,
      rolls/buns}            => {root vegetables}  0.01220132   0.2863962 2.6275247
```

Plotting limited subsets of the association rules can be done in R.  A couple of examples are shown here.

Mosaic type plot of the association rules        Network type plot of the association rules