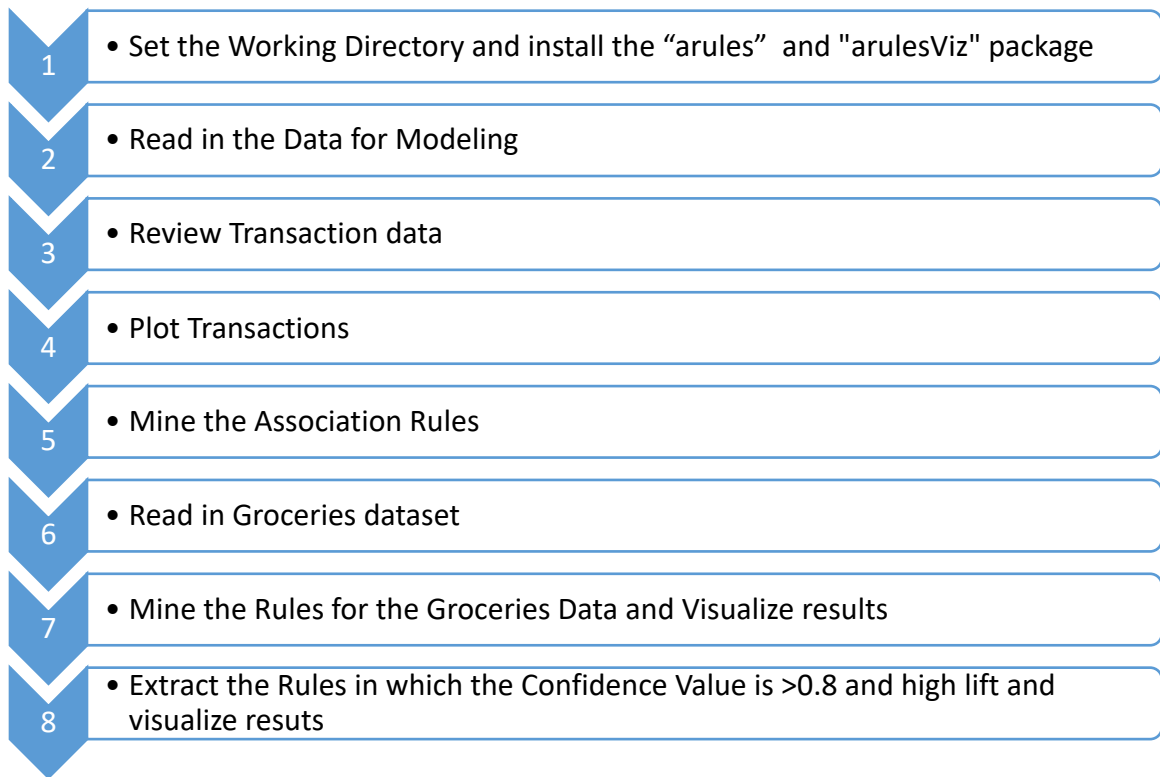


## Lab Exercise 5: Association Rules

<b>Purpose:</b>	<p>This lab is designed to investigate and practice Association Rules. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Association Rule based models</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Use the R –Studio environment to code Association Rule models</li><li>• Apply constraints in the Market Basket Analysis methods such as minimum thresholds on support and confidence measures that can be used to select interesting rules from the set of all possible rules</li><li>• Use R graphics “arules” to execute and inspect the models and the effect of the various thresholds</li></ul>
<b>References:</b>	<ul style="list-style-type: none"><li>• The groceries data set - provided for arules by Michael Hahsler, Kurt Hornik and Thomas Reutterer. <a href="http://rss.acs.unt.edu/Rdoc/library/arules/html/Groceries.html">http://rss.acs.unt.edu/Rdoc/library/arules/html/Groceries.html</a><ul style="list-style-type: none"><li>○ 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, <i>From Data and Information Analysis to Knowledge Engineering, Studies in Classification, Data Analysis, and Knowledge Organization</i>, pages 598–605. Springer-Verlag.</li></ul></li></ul>



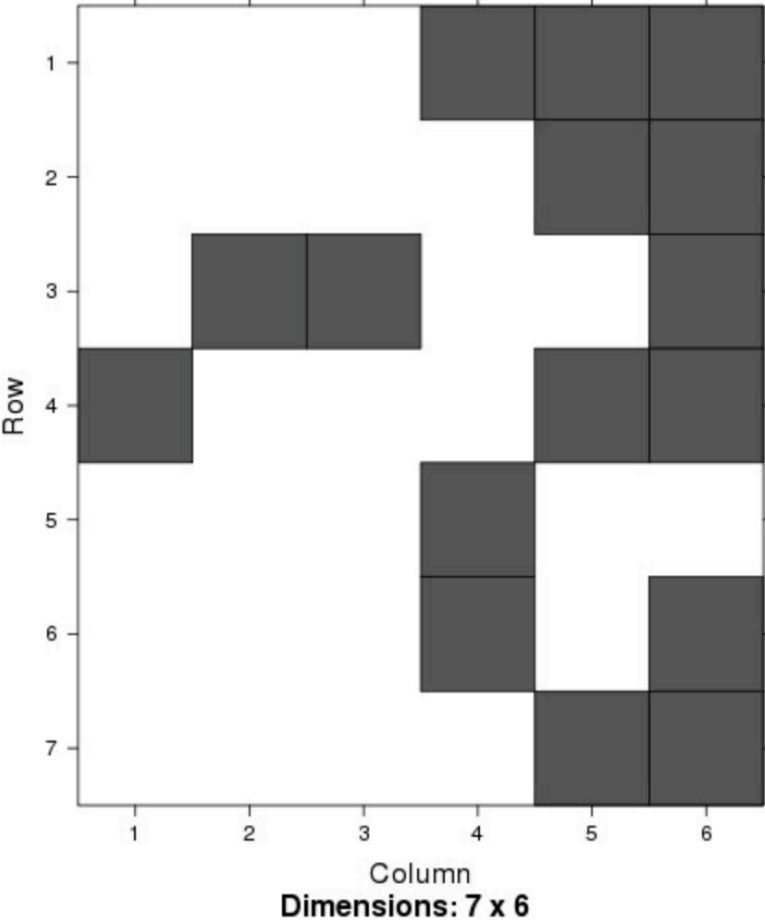
## Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set the Working Directory and install the “arules” package:</u></b></p> <p>To understand Market Basket Analysis and the R package “arules,” use a simple set of transaction lists of “book-purchases”.</p> <ol style="list-style-type: none"><li>Set the working directory to ~/LAB05/ by executing the command: <pre>setwd("~/LAB05")</pre><ul style="list-style-type: none"><li>(Or using the “Tools” option in the tool bar in the RStudio environment.)</li></ul></li><li>Load the package (select the mirror if prompted) and the required libraries: <pre>#Install the packages and load libraries &gt;install.packages('arules') &gt;install.packages('arulesViz') &gt;library('arules') &gt;library ('arulesViz')</pre></li></ol>

Step	Action
3	<p><b><u>Read in the Data for Modeling:</u></b></p> <ul style="list-style-type: none"> <li>• <b>Transaction List</b> is a special data type function in the “arules” package.</li> </ul> <p>1. Read the data in as a Transaction List using the following statement for the states data, “MBAdata.csv”.</p> <pre>&gt; #read in the csv file as a transaction data &gt; txn &lt;- read.transactions ("MBAdata.csv",rm.duplicates = FALSE,format="single",sep=" ",cols=c(1,2))</pre> <p>The arguments for the <b>read.transaction functions</b> are detailed below:</p> <ul style="list-style-type: none"> <li>• <b>file</b> the file name.</li> <li>• <b>format</b> a character string indicating the format of the data set. One of "basket" or "single", can be abbreviated.</li> <li>• <b>Sep</b> a character string specifying how fields are separated in the data file, or NULL (default). For basket format, this can be a regular expression; otherwise, a single character must be given. The default corresponds to white space separators.</li> <li>• <b>Cols</b> For the ‘single’ format, cols is a numeric vector of length two giving the numbers of the columns (fields) with the transaction and item ids, respectively. For the ‘basket’ format, cols can be a numeric scalar giving the number of the column (field) with the transaction ids. If cols = NULL</li> <li>• <b>rm.duplicates</b> a logical value specifying if duplicate items should be removed from the transactions.</li> </ul>
4	<p><b><u>Review Transaction data:</u></b></p> <p>1. First inspect the transaction data</p> <pre>&gt;txn@transactionInfo &gt;txn@itemInfo</pre> <p>2. Review the results on the console</p> <pre>&gt; txn@transactionInfo  &gt; txn@itemInfo transactionID          labels 1           101         1 Harry-Potter-DVD 2           102         2   Jane-Austen 3           103         3   Learn-Spanish 4           104         4   PSQL-basics 5           105         5     R-basics 6           106         6   Stat-intro 7           107         &gt; &gt;</pre>

Step	Action																																																								
5	<p><b>Plot Transactions:</b></p> <p>1. Use the “image” function that shows a visual representation of the transaction set in which the rows are individual transactions (identified by transaction ids) and the dark squares are items contained in each transaction.</p> <p><b>&gt; image(txn)</b></p> <div><p>The plot displays a 7x6 matrix where rows represent transactions (1-7) and columns represent items (1-6). Dark squares indicate the presence of an item in a transaction.</p><table border="1"><thead><tr><th>Transaction \ Item</th><th>1</th><th>2</th><th>3</th><th>4</th><th>5</th><th>6</th></tr></thead><tbody><tr><th>1</th><td></td><td></td><td></td><td>1</td><td>1</td><td>1</td></tr><tr><th>2</th><td></td><td></td><td></td><td></td><td>1</td><td>1</td></tr><tr><th>3</th><td></td><td>1</td><td>1</td><td></td><td></td><td>1</td></tr><tr><th>4</th><td>1</td><td></td><td></td><td></td><td>1</td><td>1</td></tr><tr><th>5</th><td></td><td></td><td></td><td>1</td><td></td><td></td></tr><tr><th>6</th><td></td><td></td><td></td><td>1</td><td></td><td>1</td></tr><tr><th>7</th><td></td><td></td><td></td><td></td><td>1</td><td>1</td></tr></tbody></table><p>Dimensions: 7 x 6</p></div>	Transaction \ Item	1	2	3	4	5	6	1				1	1	1	2					1	1	3		1	1			1	4	1				1	1	5				1			6				1		1	7					1	1
Transaction \ Item	1	2	3	4	5	6																																																			
1				1	1	1																																																			
2					1	1																																																			
3		1	1			1																																																			
4	1				1	1																																																			
5				1																																																					
6				1		1																																																			
7					1	1																																																			
	2. Review the output in the graphics window																																																								

**Mine the Association Rules:**

The “**apriori**” function, provided by the *arulesr* package, is used as follows:

```
rules <- apriori(File,
                 parameter = list(supp = 0.5, conf = 0.9,
                                target = "rules"))
```

where the arguments are:

- **data** object of class transactions or any data structure which can be coerced into transactions (for example, a binary matrix or data.frame).
- **parameter** named list. The default behavior is to mine rules with support 0.1, confidence 0.8, and maxlen 5.

1. Read in the statement for the transaction data:

```
> #mine association rules
> basket_rules <-
apriori(txn,parameter=list(sup=0.5,conf=0.9,target="rules"
))

> basket_rules <- apriori(txn,parameter=list(sup=0.5,conf=0.9,target="rules"))

parameter specification:
confidence minval smax arem aval originalSupport support minlen maxlen target ext
0.9 0.1 1 none FALSE TRUE 0.5 1 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 ...[6 item(s), 7 transaction(s)] done [0.00s].
sorting and recoding items ... [2 item(s)] done [0.00s].
creating transaction tree ... done [0.00s].
checking subsets of size 1 2 done [0.00s].
writing ... [1 rule(s)] done [0.00s].
creating S4 object ... done [0.00s].
>
```

2. Review the output on the console. The number of rules generated can be seen in the output and is represented as follows:

```
writing ... [1 rule(s)] done [0.00s]
```

3. Inspect the rule using the following statement:

```
> inspect(basket_rules)
```

Step	Action
	<pre data-bbox="397 317 1023 409">&gt; inspect(basket_rules)   lhs      rhs      support confidence lift 1 {R-basics} =&gt; {Stat-intro} 0.571      1 1.17 &gt;</pre> <p data-bbox="397 443 652 474">4. Review the output.</p> <p data-bbox="397 514 1414 583">5. State the generated rule and the support, confidence and the lift thresholds for the rule</p> <p data-bbox="397 623 558 655">0.571, 1, 1.17</p>



Step	Action
7	<p><b><u>Read in Groceries dataset</u></b></p> <p>Use the standard data set, “Groceries” available with the “arules” package.</p> <ul style="list-style-type: none"> <li>The Groceries data set contains 1 month (30 days) of real-world point-of-sale transaction data from a typical local grocery outlet. The data set contains 9835 transactions and the items are aggregated to 169 categories.</li> </ul> <p>1. Read in the data set and inspect the item information</p> <pre> &gt; #Read in Groceries data &gt; data(Groceries) &gt; Groceries@itemInfo  &gt; Groceries@itemInfo       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 </pre>

Step	Action
8	<p><b><u>Mine the Rules for the Groceries Data:</u></b></p> <pre> &gt; #mine rules &gt; rules &lt;- apriori(Groceries, parameter=list(support=0.001, confidence=0.5)) </pre> <ul style="list-style-type: none"> <li>Note the values used for the parameter list.</li> </ul> <pre> &gt; rules &lt;- apriori(Groceries, parameter=list(support=0.001, confidence=0.5)) </pre> <pre> parameter specification:  confidence minval  smax  arem  aval originalSupport support  minlen maxlen target  ext       0.5      0.1    1 none  FALSE              TRUE   0.001     1     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 [0.02s]. writing ... [5668 rule(s)] done [0.00s]. creating S4 object ... done [0.00s]. &gt; </pre> <p>1. How many rules are generated? 5668 rules</p>

Step	Action																				
9	<p><b><u>Extract the Rules in which the Confidence Value is &gt;0.8 and high lift:</u></b></p> <p>1. Execute the following commands:</p> <pre>&gt; subrules &lt;- rules[quality(rules)\$confidence &gt; 0.8] &gt; plot(subrules, control = list(jitter=2)) &gt; inspect(subrules)</pre> <p><b>Screenshot of scatterplot</b></p> <pre>&gt; plot(subrules, control=list(jitter=2)) Error in as.double(y) :   cannot coerce type 'S4' to vector of type 'double' &gt;</pre> <p><b>Note: This is the datatype error that we discussed.</b></p> <p>2. Review the results.</p> <p>3. <b>How many sub-rules did you extract? 371</b></p> <ul style="list-style-type: none"><li>• These rules are more valuable for the business.</li></ul> <p>4. Extract the top three rules with high threshold for the parameter “lift” and plot.</p> <pre>&gt; #Extract the top three rules with high lift &gt; rules_high_lift &lt;- head(sort(rules, by="lift"), 3) &gt; inspect(rules_high_lift)</pre> <table><thead><tr><th>lhs</th><th>rhs</th><th>support</th><th>confidence</th><th>lift</th></tr></thead><tbody><tr><td>1 {Instant food products, soda}</td><td>=&gt; {hamburger meat}</td><td>0.00122</td><td>0.632</td><td>19.0</td></tr><tr><td>2 {soda, popcorn}</td><td>=&gt; {salty snack}</td><td>0.00122</td><td>0.632</td><td>16.7</td></tr><tr><td>3 {flour, baking powder}</td><td>=&gt; {sugar}</td><td>0.00102</td><td>0.556</td><td>16.4</td></tr></tbody></table> <pre>&gt; plot(rules_high_lift, method="graph", control=list(type="items")) Error in as.double(y) :   cannot coerce type 'S4' to vector of type 'double'</pre> <p><b>Note: This is the datatype error that we discussed.</b></p>	lhs	rhs	support	confidence	lift	1 {Instant food products, soda}	=> {hamburger meat}	0.00122	0.632	19.0	2 {soda, popcorn}	=> {salty snack}	0.00122	0.632	16.7	3 {flour, baking powder}	=> {sugar}	0.00102	0.556	16.4
lhs	rhs	support	confidence	lift																	
1 {Instant food products, soda}	=> {hamburger meat}	0.00122	0.632	19.0																	
2 {soda, popcorn}	=> {salty snack}	0.00122	0.632	16.7																	
3 {flour, baking powder}	=> {sugar}	0.00102	0.556	16.4																	

*End of Lab Exercise*