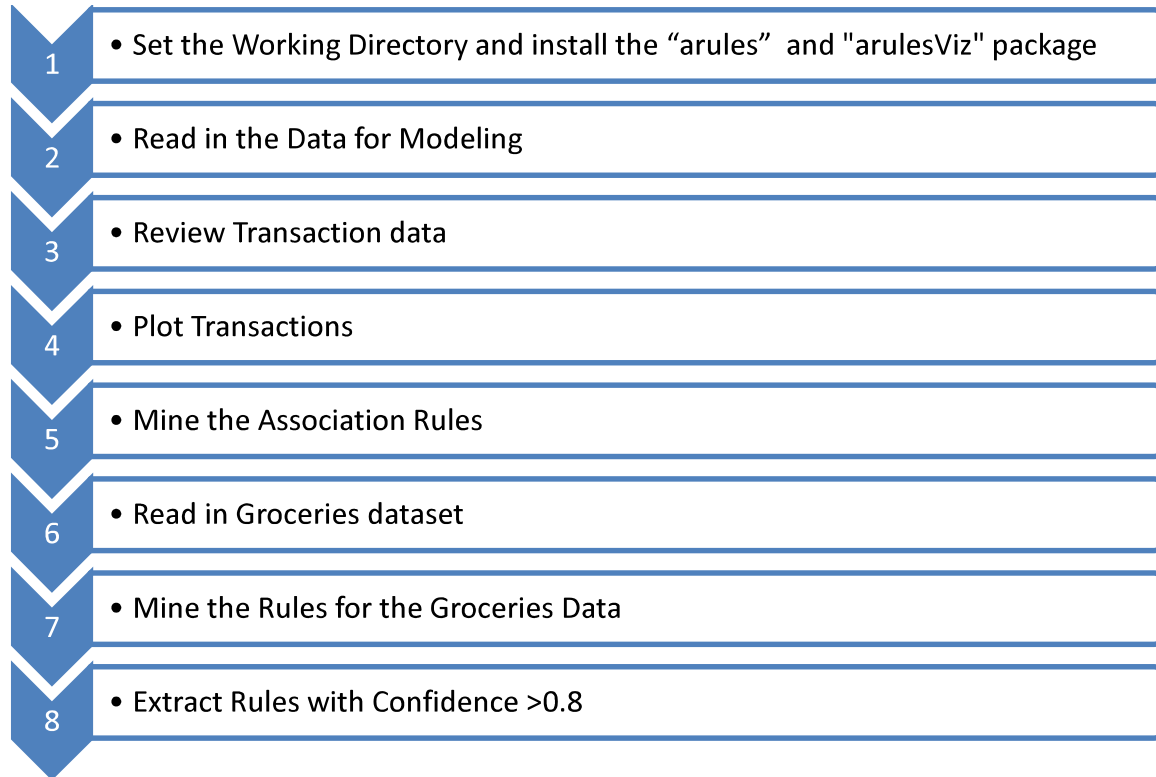


Lab Exercise 5: Association Rules

Purpose:	<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">• Use R functions for Association Rule based models
Tasks:	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none">• Use the R –Studio environment to code Association Rule models• 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• Use R graphics “arules” to execute and inspect the models and the effect of the various thresholds
References:	<ul style="list-style-type: none">• The groceries data set - provided for arules by Michael Hahsler, Kurt Hornik and Thomas Reutterer. http://rss.acs.unt.edu/Rdoc/library/arules/html/Groceries.html<ul style="list-style-type: none">○ 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.

Workflow Overview



LAB Instructions

Step	Action
1	<p>Log in with gpadmin credentials on to R-Studio.</p> <p>The R Code for this exercise is available at <code>home/gpadmin/LAB05/mba.R</code></p>
2	<p><u>Set the Working Directory and install the packages:</u></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">1. Set the working directory to <code>~/LAB05/</code> by executing the command: <pre>setwd("~/LAB05")</pre>2. Load the package (select the mirror if prompted) and the required libraries: <pre>#Install the packages and load libraries >install.packages('arules') >install.packages('arulesViz') >library('arules') >library ('arulesViz')</pre>

Step	Action
3	<p><u>Read in the Data for Modeling:</u></p> <ul style="list-style-type: none"> • Transaction List is a special data type function in the “arules” package. <p>Read the data in as a Transaction List using the following statement for the states data, “MBAdata.csv”.</p> <pre>> txn <- read.transactions ("MBAdata.csv",rm.duplicates = FALSE,format="single",sep=" ",cols=c(1,2))</pre> <p>The arguments for the read.transaction functions are detailed below:</p> <ul style="list-style-type: none"> • file the file name. • format a character string indicating the format of the data set. One of "basket" or "single", can be abbreviated. • Sep 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. • Cols 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 • rm.duplicates a logical value specifying if duplicate items should be removed from the transactions.
4	<p><u>Review Transaction Data:</u></p> <ol style="list-style-type: none"> 1. First inspect the transaction data <pre>>txn@transactionInfo</pre> <pre>>txn@itemInfo</pre> 2. Review the results on the console
5	<p><u>Plot Transactions:</u></p> <ol style="list-style-type: none"> 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. <pre>> image (txn)</pre> 2. Review the output in the graphics window

Step	Action
6	<p><u>Mine the Association Rules:</u></p> <p>The <code>help(apriori)</code> command provides documentation of the <i>apriori</i> function from the <i>arulesr</i> package. The syntax of the <i>apriori</i> function is as follows:</p> <pre>rules <- apriori(File, parameter = list(supp = 0.5, conf = 0.9, target = "rules"))</pre> <p>where the arguments are:</p> <ul style="list-style-type: none"> • 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. <ol style="list-style-type: none"> 1. Build the rules from the transaction data: <pre>> basket_rules <- apriori(txn,parameter=list(sup=0.5,conf=0.9,target="rules"))</pre> 2. Review the output on the console. The number of rules generated can be seen in the output. 3. Inspect the rule using the following statement: <pre>> inspect(basket_rules)</pre> 4. Review the output. Observe the generated rule and the support, confidence and the lift thresholds for the rule.

Step	Action
7	<p><u>Read in Groceries dataset:</u></p> <p>Use the standard data set, “Groceries” available with the “arules” package.</p> <p>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.</p> <p>Read in the data set and inspect the item information</p> <pre>> data(Groceries) > Groceries > Groceries@itemInfo</pre>
8	<p><u>Mine the Rules for the Groceries Dataset:</u></p> <pre>> rules <- apriori(Groceries, parameter=list(support=0.001, confidence=0.5))</pre> <p>Note the values used for the parameter list.</p> <p>How many rules are generated?</p>
9	<p><u>Extract Rules with Confidence >0.8:</u></p> <ol style="list-style-type: none"> Execute the following commands: <pre>> subrules <- rules[quality(rules)\$confidence > 0.8] > plot(subrules, control = list(jitter=2)) > inspect(subrules)</pre> Review the results. How many sub-rules did you extract? Extract the top three rules with high threshold for the parameter “lift” and plot. <pre>> #Extract the top three rules with high lift > rules_high_lift <- head(sort(rules, by="lift"), 3) > inspect(rules_high_lift) > plot(rules_high_lift, method="graph", control=list(type="items"))</pre> List the rules and the value of the parameters associated with these rules:

End of Lab Exercise