



# Data Science and Big Data Analytics

## Lab Guide

## Copyright

Copyright © 1996, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012 EMC Corporation. All Rights Reserved. EMC believes the information in this publication is accurate as of its publication date. The information is subject to change without notice.

THE INFORMATION IN THIS PUBLICATION IS PROVIDED “AS IS.” EMC CORPORATION MAKES NO REPRESENTATIONS OR WARRANTIES OF ANY KIND WITH RESPECT TO THE INFORMATION IN THIS PUBLICATION, AND SPECIFICALLY DISCLAIMS IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Use, copying, and distribution of any EMC software described in this publication requires an applicable software license.

EMC2, EMC, Data Domain, RSA, EMC Centera, EMC ControlCenter, EMC LifeLine, EMC OnCourse, EMC Proven, EMC Snap, EMC SourceOne, EMC Storage Administrator, Acartus, Access Logix, AdvantEdge, AlphaStor, ApplicationXtender, ArchiveXtender, Atmos, Authentica, Authentic Problems, Automated Resource Manager, AutoStart, AutoSwap, AVALONidm, Avamar, Captiva, Catalog Solution, C-Clip, Celerra, Celerra Replicator, Centera, CenterStage, CentraStar, ClaimPack, ClaimsEditor, CLARiiON, ClientPak, Codebook Correlation Technology, Common Information Model, Configuration Intelligence, Configuresoft, Connectrix, CopyCross, CopyPoint, Dantz, DatabaseXtender, Direct Matrix Architecture, DiskXtender, DiskXtender 2000, Document Sciences, Documentum, elnput, E-Lab, EmailXaminer, EmailXtender, Enginuity, eRoom, Event Explorer, FarPoint, FirstPass, FLARE, FormWare, Geosynchrony, Global File Virtualization, Graphic Visualization, Greenplum, HighRoad, HomeBase, InfoMover, Infoscapes, Infra, InputAccel, InputAccel Express, Invista, Ionix, ISIS, Max Retriever, MediaStor, MirrorView, Navisphere, NetWorker, nLayers, OnAlert, OpenScale, PixTools, Powerlink, PowerPath, PowerSnap, QuickScan, Rainfinity, RepliCare, RepliStor, ResourcePak, Retrospect, RSA, the RSA logo, SafeLine, SAN Advisor, SAN Copy, SAN Manager, Smarts, SnapImage, SnapSure, SnapView, SRDF, StorageScope, SupportMate, SymmAPI, SymmEnabler, Symmetrix, Symmetrix DMX, Symmetrix VMAX, TimeFinder, UltraFlex, UltraPoint, UltraScale, Unisphere, VMAX, Vblock, Viewlets, Virtual Matrix, Virtual Matrix Architecture, Virtual Provisioning, VisualSAN, VisualSRM, Voyence, VPLEX, VSAM-Assist, WebXtender, xPression, xPresso, YottaYotta, the EMC logo, and where information lives, are registered trademarks or trademarks of EMC Corporation in the United States and other countries.

All other trademarks used herein are the property of their respective owners.

© Copyright 2012 EMC Corporation. All rights reserved. Published in the USA.

Revision Date: February 2012  
Revision Number: MR-1CP-DSBDA .1.2

## Document Revision History

Rev #	File Name	Date
1	DSBDA Integrated Lab Guide 12-01-11 alpha.doc	12-01-11
2	DSBDA Integrated Lab Guide 12-07-11 alpha03.doc  (incorporating the review updates by Noelle and Rashmi)	12-07-11
4	DSBDA Integrated Lab Guide 12-16-11 alpha04.doc  (incorporating the review updates by Noelle)	12-16-11
5	Lab workflow included for labs 3 to 12  Beta version	1-4-12
6	GA Version	1-18-12
7	Made VILT-related edits. Added a disclaimer that applies to the VILT.	4-9-12

## Table of Contents

<b>COPYRIGHT .....</b>	<b>2</b>
<b>DOCUMENT REVISION HISTORY .....</b>	<b>3</b>
<b>THIS PAGE INTENTIONALLY LEFT BLANK. ....</b>	<b>5</b>
<b>LAB EXERCISE 1: INTRODUCTION TO DATA ENVIRONMENT .....</b>	<b>7</b>
<b>1.1 ACCESSING LAB ENVIRONMENT .....</b>	<b>8</b>
<b>1.2 DATABASE ENVIRONMENT – RETAIL DATA .....</b>	<b>9</b>
<b>1.3 DATABASE ENVIRONMENT-CENSUS DATA .....</b>	<b>14</b>
<b>LAB EXERCISE 2: INTRODUCTION TO R .....</b>	<b>21</b>
<b>LAB EXERCISE 3: BASIC STATISTICS, VISUALIZATION, AND HYPOTHESIS TESTS .....</b>	<b>29</b>
<b>PART 1 – BASIC STATISTICS AND VISUALIZATION USING R .....</b>	<b>30</b>
<i>Workflow Overview .....</i>	<i>30</i>
<i>LAB Instructions .....</i>	<i>31</i>
<b>PART 2 – GRAPHICS PACKAGE PLOTS AND HYPOTHESIS TESTS .....</b>	<b>37</b>
<i>Workflow Overview .....</i>	<i>37</i>
<i>Lab Instructions .....</i>	<i>38</i>
<b>LAB EXERCISE 4: K-MEANS CLUSTERING .....</b>	<b>45</b>
<i>Workflow overview .....</i>	<i>46</i>
<i>Lab Instructions .....</i>	<i>47</i>
<b>LAB EXERCISE 5: ASSOCIATION RULES .....</b>	<b>55</b>
<i>Workflow Overview .....</i>	<i>56</i>
<i>LAB Instructions .....</i>	<i>57</i>
<b>LAB EXERCISE 6: LINEAR REGRESSION .....</b>	<b>61</b>
<i>Workflow Overview .....</i>	<i>62</i>
<i>LAB Instructions .....</i>	<i>63</i>
<b>LAB EXERCISE 7: LOGISTIC REGRESSION .....</b>	<b>71</b>
<i>Workflow Overview .....</i>	<i>72</i>
<i>LAB Instructions .....</i>	<i>73</i>
<b>LAB EXERCISE 8: NAÏVE BAYESIAN CLASSIFIER .....</b>	<b>83</b>
<b>PART 1 – BUILDING NAÏVE BAYESIAN CLASSIFIER .....</b>	<b>84</b>
<i>Workflow Overview .....</i>	<i>84</i>
<i>LAB Instructions .....</i>	<i>85</i>
<b>PART 2 – NAÏVE BAYESIAN CLASSIFIER – CENSUS DATA .....</b>	<b>91</b>
<i>Workflow Overview .....</i>	<i>91</i>
<i>LAB Instructions .....</i>	<i>92</i>

<b>LAB EXERCISE 9: DECISION TREES.....</b>	<b>101</b>
<i>Workflow Overview .....</i>	<i>102</i>
<i>LAB Instructions .....</i>	<i>103</i>
<b>LAB EXERCISE 10: TIME SERIES ANALYSIS WITH ARIMA .....</b>	<b>107</b>
<i>Workflow Overview .....</i>	<i>108</i>
<i>LAB Instructions .....</i>	<i>109</i>
<b>LAB EXERCISE 11: HADOOP, HDFS AND MAPREDUCE.....</b>	<b>117</b>
<i>Workflow Overview .....</i>	<i>118</i>
<i>LAB Instructions .....</i>	<i>119</i>
<b>LAB 12: IN-DATABASE ANALYTICS .....</b>	<b>123</b>
<b>PART 1 – IN-DATABASE ANALYSIS OF CLICK-STREAM DATA .....</b>	<b>124</b>
<i>Workflow Overview .....</i>	<i>124</i>
<i>LAB Instructions .....</i>	<i>125</i>
<b>PART 2 – IN-DATABASE COMPUTATION OF MEDIAN WITH ORDERED AGGREGATES .....</b>	<b>132</b>
<i>Workflow Overview .....</i>	<i>132</i>
<i>LAB Instructions .....</i>	<i>133</i>
<b>PART 3: LOGISTIC REGRESSION WITH MADLIB .....</b>	<b>135</b>
<i>Workflow Overview .....</i>	<i>135</i>
<i>LAB Instructions .....</i>	<i>136</i>
<b>FINAL LAB EXERCISE ON BIG DATA ANALYTICS.....</b>	<b>139</b>
<i>Case Study Background and Problem Definition .....</i>	<i>140</i>
<i>Suggested Workflow and Checkpoints for the Lab .....</i>	<i>143</i>

**THIS PAGE INTENTIONALLY LEFT BLANK.**



## Lab Exercise 1: Introduction to Data Environment

Purpose:	<p>The first lab introduces the <i>Analytics Lab Environment</i> you will be working on throughout the course. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Authenticate and access the Virtual Machine (VM) assigned to you for all of your lab exercises</li><li>• Use SQL and Meta commands in PSQL to navigate through the data sets</li><li>• Create subsets of the <i>data</i>, using <i>table joins and filters</i> to analyze subsequent lab exercises</li></ul>
Tasks:	<p>Tasks you will complete in this lab exercise include:</p> <ul style="list-style-type: none"><li>• Exploring databases and datasets</li><li>• Using PSQL statements and Meta commands.</li><li>• Creating subsets of data for use in subsequent lab exercises</li></ul>
References:	<p>References used throughout the labs are located in your <b><i>Student Resource Guide Appendix</i></b>. See the Appendix for:</p> <ul style="list-style-type: none"><li>• PSQL Commands – Quick Reference</li><li>• PSQL Meta Commands – Quick Reference</li><li>• Surviving LINUX – Quick Reference</li><li>• R – Quick Reference</li></ul>

## 1.1 Accessing Lab Environment

**Note:** The lab access is available only in the Instructor Led Training (ILT) program. **For those learning the content via the Video ILT (DVD), the labs are provided only as a reference and are recorded for demonstration purposes only.**

Step	Action
1	<p><b>Accessing from a Public ISP:</b></p> <ol style="list-style-type: none"><li>1. From any Internet connection, open a suitable browser (Internet Explorer strongly recommended) at <a href="https://vdc.emc.com">https://vdc.emc.com</a></li><li>2. Your user name and password details are provided by your instructor.</li></ol> <p><b>Accessing the LAB</b></p> <ol style="list-style-type: none"><li>1. Refer to the instruction sheet for the access details to your assigned windows host - called your “Front-End” (fe) host - and follow the instructions to open a desktop session to your “fe”.</li><li>2. The access to the “Front-End” (fe) of the Lab is now established.</li><li>3. Refer to the instruction sheet for the IP address of the “Back-End” (be) that hosts the databases and the RStudio environment</li><li>4. RStudio may be accessed through the “safari” browser available as a desktop icon on your “fe” host. For RStudio, direct “safari” browser on “fe” to the URL <a href="http://&lt;be host IP Address&gt;:8787/">http://&lt;be host IP Address&gt;:8787/</a></li><li>5. Utilities such as “PuTTY”, WinSCP and PGadmin III are also available on the “fe” to access and update contents in the “be”. Some lab exercises will require that you log in directly on the “be” via a terminal session; you may use “PuTTY” to start a terminal session.</li><li>6. Follow the lab guide for additional instructions that may be associated with individual labs.</li></ol>



## 1.2 Database Environment – Retail Data

Step	Action
1	<p>Log in to the “be” host using “PuTTY” on your “fe” host. Specify “gpadmin” as the username, and the appropriate password (to be supplied by your instructor).</p> <p>Currently you are logged in as <b>GADMIN</b> and you have administrative access to the <i>Greenplum Database Environment</i>, in which you will be working.</p> <p>You must first verify if the database up and running.</p> <ol style="list-style-type: none"> <li>1. Type: <b>gpstate</b></li> <li>2. Review the output; you should be able to see that the database is active with the following output. Please note that because of the large output size we only show selected lines and that your configuration details may slightly differ from the one below.</li> </ol> <pre>[INFO]:-Starting gpstate with args: [INFO]:-local Greenplum Version: 'postgres (Greenplum Database) 4.1.1.1 build 1' [INFO]:-Obtaining Segment details from master... [INFO]:-Gathering data from segments... [INFO]:-Greenplum instance status summary [INFO]:----- [INFO]:- Master instance           = Active [INFO]:- Master standby           = No master standby configured ... [INFO]:- Total primary segments           = 2 [INFO]:- Total primary segment valid (at master) = 2 [INFO]:- Total primary segment failures (at master) = 0 ... [INFO]:- Mirrors not configured on this array [INFO]:-----</pre>

Step	Action
2	<p>Now you're ready to open a PSQL session and check all available databases.</p> <p>Refer to the <i>PSQL Commands – Quick Reference</i>, located in your <b>Student Resource Guide Appendix</b>, for the PSQL meta commands.</p> <p><b>Note:</b> PSQL meta commands start with a backslash (\). To review all available meta commands type backslash and question mark (\?).</p> <p>To review all available databases in your environment:</p> <ol style="list-style-type: none"> <li>1. Type: <code>psql</code> This will open a new PSQL session to the default database.</li> <li>2. Next type: <code>\l</code> Notice a list of databases and record databases named "training*".</li> </ol> <p><b>Note:</b> Another way of listing all available databases (without opening a PSQL session) is to call PSQL executable with parameter (-l): <code>psql -l</code></p>
3	<p><b><u>Connect to the training1 database:</u></b></p> <ol style="list-style-type: none"> <li>1. At the PSQL prompt type : <code>\c training1</code> at the OS level prompt type: <code>psql training1</code></li> </ol> <p>To see the schemas you have in this database:</p> <ol style="list-style-type: none"> <li>2. Type: <code>\dn</code> <ul style="list-style-type: none"> <li>• You should see "ddemo" schema, listed.</li> <li>• You should also ensure that this schema is included in the search path.</li> </ul> </li> <li>3. Execute your first PSQL command, type:           <pre>SET search_path TO ddemo, public;</pre> </li> </ol> <p><b>Note:</b> PSQL commands are terminated with a semi-colon- ";"</p>

Step	Action						
4	<p>You can now view the tables in this database.</p> <ol style="list-style-type: none"><li>1. Type: <code>\dt</code></li><li>2. Record the number of tables in the database: _____</li><li>3. Locate the table, “customers_dim”.</li><li>4. Review the column descriptions for this table:</li><li>5. Type: <code>\d+ customers_dim</code></li><li>6. Record the column descriptions, their types and column name(s) by which the table is distributed (aka: the distribution key):</li></ol> <table><thead><tr><th>Column Descriptions</th><th>Type</th><th>Distribution Key Column(s)</th></tr></thead><tbody><tr><td></td><td></td><td></td></tr></tbody></table>	Column Descriptions	Type	Distribution Key Column(s)			
Column Descriptions	Type	Distribution Key Column(s)					
5	<p><b><u>Analyze the gender distribution of the customer base:</u></b></p> <ol style="list-style-type: none"><li>1. To locate the number of males and females type:  <code>SELECT gender,count(*) FROM customers_dim GROUP BY gender;</code></li><li>2. Record the number of female customers: _____</li><li>3. Record the number of male customers: _____</li><li>4. Record the total number of customers: _____</li></ol>						

Step	Action
6	<p>1. Using PSQL, generate a report on the average spending by gender, Type:</p> <pre> SELECT     c.gender , AVG(o.item_price) AS avg_price FROM     ddemo.order_lineitems AS o JOIN     ddemo.customers_dim AS c     ON o.customer_id = c.customer_id GROUP BY c.gender ; </pre> <p><b>Note:</b> You can find this code in the LAB01 directory. This script can be executed using the following command from the OS prompt:</p> <p>2. To exit the PSQL environment, use the following meta command, type:</p> <pre>\q</pre> <p>You are now at the OS prompt.</p> <p>3. To execute the SQL script type:</p> <pre>psql -d training1 -f lab1p1step6.sql</pre> <p><b>Note 1:</b> In the <i>psql</i> command above option “-d” specifies the database name to connect to (“training1”). This is equivalent to specifying <i>dbname</i> as the first <b>non-option argument</b> on the command line. As a convention we have used the option “-d” throughout this document. However <i>dbname</i> can be specified without option “-d” as long as it is the first argument of the <i>psql</i> command.</p> <p><b>Note 2:</b> This query may take some time to execute as it is processing a million rows of data.</p> <p>4. Record the average expenditures by gender:</p> <p>Male : _____ Female: _____</p>

Step	Action																		
7	<p>Use the script, “lab1p1step7”, with the appropriate modifications to list the top five product categories ordered by men and women.</p> <table><tr><td></td><td><i>Men</i></td><td><i>Women</i></td></tr><tr><td>1</td><td></td><td></td></tr><tr><td>2</td><td></td><td></td></tr><tr><td>3</td><td></td><td></td></tr><tr><td>4</td><td></td><td></td></tr><tr><td>5</td><td></td><td></td></tr></table>		<i>Men</i>	<i>Women</i>	1			2			3			4			5		
	<i>Men</i>	<i>Women</i>																	
1																			
2																			
3																			
4																			
5																			

## 1.3 Database Environment-Census Data

Step	Action
1	Follow the steps detailed in, Lab 1 - Data Set 1, to connect to and inspect another database “training2”.
2	Record the tables in database (Schema – Public)“training2”
3	Describe the type of data in the database.
4	Record the number of rows in each table.
5	<p><b><u>Data Preparation &amp; Cleanup – 1:</u></b></p> <p>(Scenario) You realize that the Intern who loaded the “housing” data has copied records into the table twice. Each different row is represented by a unique combination of “serialno” and “state” columns.</p> <p>1. Execute the following code:</p> <pre> SELECT     SUM(c) AS total_records   , SUM(CASE WHEN c&gt;1 THEN c-1 ELSE 0 END) AS total_dupes   , COUNT(*) AS total_uniques FROM (     SELECT         COUNT(*) AS c     FROM         housing     GROUP BY         serialno         , state     ) AS dupes ; </pre> <p><b>Note:</b> This code is also available at,</p> <p><b><u>/home/gpadmin/LAB01/countdupes.sql,</u></b></p> <p>2. Record the total number of records in the table: _____</p> <p>3. Record the total number of duplicate records: _____</p> <p>4. Record the total number of unique records: _____</p>

Step	Action
6	<p><b><u>Data Preparation &amp; Cleanup – 2:</u></b></p> <p>To prepare and clean the data you need to create a “housing_nodupes” table. Make sure that you are in the PSQL environment if you have previously exited to the OS command line.</p> <ol style="list-style-type: none"> <li>1. Check to see if a table already exists with the name (“housing_nodupes”). Type <code>\dt</code>  Note: the command <code>\dt</code> will list all tables in the database. <code>\dt public.*</code> will list all tables in the public schema.</li> <li>2. If this table already exists execute the following SQL statement:  <pre>DROP TABLE IF EXISTS housing_nodupes;</pre></li> <li>3. Execute the following SQL statement:  <pre>CREATE TABLE housing_nodupes AS SELECT DISTINCT ON     (serialno, state) * FROM     housing DISTRIBUTED BY (serialno, state) ;</pre> <p><b>Note:</b> This code is also available at, <code>/home/gpadmin/LAB01/lab1p2step6.sql</code></p></li> <li>4. Repeat the queries in Step 5 (previous step) to ensure that there are no duplicate records in the housing_nodupes table.</li> </ol>

Step	Action
7	<p data-bbox="337 184 844 220"><b><u>Basic Analytics Using the “Housing” Data:</u></b></p> <ol data-bbox="337 256 1404 1228" style="list-style-type: none"> <li data-bbox="337 256 1404 609">1. Execute the following SQL statement to calculate correlation between household income and number of rooms:   <pre data-bbox="435 367 763 609"> SELECT     corr(hinc, rooms) FROM     housing_nodupes WHERE     state = 25 ; </pre> </li> <li data-bbox="337 651 617 682">2. Record your result:</li> <li data-bbox="337 793 1404 1144">3. Execute the following SQL statement calculate the R-squared of the regression line of household income and number of rooms::   <pre data-bbox="435 903 812 1144"> SELECT     regr_r2(hinc, rooms) FROM     housing_nodupes WHERE     state = 25 ; </pre> </li> <li data-bbox="337 1186 617 1218">4. Record your result:</li> </ol>



Step	Action
8	<p><b><u>Prepare “Housing” Data for Subsequent Analytic Exercises:</u></b></p> <p>You need to prepare data from the, “housing_nodupes” and “persons” tables, for subsequent analysis with “R” in the next module.</p> <p>1. Run the following commands and SQL query to move (pipe) the results into a text file</p> <p><b>Note:</b> Use the meta commands to render your output to a file and remove the white spaces (formatting)</p> <pre> \q \o lab1_01.txt SELECT     serialno , hinc , rooms FROM     housing_nodupes WHERE     hinc &gt; 0     AND state = 25 ; </pre> <p><b>Note:</b> The SQL query is also available at the following location:</p> <p>/home/gpadmin/LAB01/lab1p2step8.sql</p> <p>2. Alternatively you can execute the following command from the OS prompt:</p> <pre>psql -d training2 -f lab1p2step8.sql</pre> <p>Now, your data is ready for the lab exercise in the next module.</p> <p>3. Remove the summary line at the end of the output file lab1_01.txt</p>

Step	Action																		
9	<p><b><u>Prepare “Persons” Data for Subsequent Analytic Exercises:</u></b></p> <p>Prepare a summary table with the number of people by race and by education level.</p> <p><b>Note:</b> Use the following Races: White, Black, American Indian/Alaska Native, Asian, Hawaiian /Pacific Islander, and Others.</p> <pre>(white) White, (black) Black, (aian) American_Indian_Alaska_native, (asian) Asian, (nhpi) Hawaii_pacific_islander, (other) Others</pre> <p><b><u>Use the following Education Levels:</u></b></p> <table><tr><td>01. No schooling completed</td><td>06. 10th grade</td><td>11. One or more years of college, no degree</td></tr><tr><td>02. Nursery school to 4th grade</td><td>07. 11th grade</td><td>12. Associate degree</td></tr><tr><td>03. 5th grade or 6th grade</td><td>08. 12th grade, no diploma</td><td>13. Bachelor’s degree</td></tr><tr><td>04. 7th grade or 8th grade</td><td>09. High school graduate</td><td>14. Master’s degree</td></tr><tr><td>05. 9th grade</td><td>10. Some college, but less than 1 year</td><td>15. Professional degree</td></tr><tr><td></td><td></td><td>16. Doctorate degree</td></tr></table> <p>1. Create a table with columns for Races and rows for Educational Level. (The cells denote the number of “persons” for each category.) Prepare a text file with headers to use in the next module. SQL code necessary for this task is presented below:</p> <pre>\a \o lab1_02.txt SELECT     educ AS Education_Level     , SUM(white) AS White     , SUM(black) AS Black     , SUM(aian) AS American_Indian_Alaska_Native     , SUM(asian) AS Asian     , SUM(nhpi) AS Hawaii_Pacific_Islander     , SUM(other) AS Others FROM     persons WHERE     age &gt; 17     AND educ &gt; 0 GROUP BY educ ORDER BY educ ;</pre>	01. No schooling completed	06. 10th grade	11. One or more years of college, no degree	02. Nursery school to 4th grade	07. 11th grade	12. Associate degree	03. 5th grade or 6th grade	08. 12th grade, no diploma	13. Bachelor’s degree	04. 7th grade or 8th grade	09. High school graduate	14. Master’s degree	05. 9th grade	10. Some college, but less than 1 year	15. Professional degree			16. Doctorate degree
01. No schooling completed	06. 10th grade	11. One or more years of college, no degree																	
02. Nursery school to 4th grade	07. 11th grade	12. Associate degree																	
03. 5th grade or 6th grade	08. 12th grade, no diploma	13. Bachelor’s degree																	
04. 7th grade or 8th grade	09. High school graduate	14. Master’s degree																	
05. 9th grade	10. Some college, but less than 1 year	15. Professional degree																	
		16. Doctorate degree																	

Step	Action
10	<p>The code in step 9 is also available at the following location: /home/gpadmin/LAB01/lab1p2step9.sql</p> <p>Execute the following command from the OS prompt:</p> <pre>psql -d training2 -f lab1p2step9.sql</pre> <p>Remove the last “summary” line as you did in Step 8 and prepare the file “lab1_02.txt” for the lab exercise in the next module.</p>

*End of Lab Exercise*



## Lab Exercise 2: Introduction to R

<b>Purpose:</b>	<p>This lab introduces you to the use of the R statistical package within the Data Science and Big Data Analytics environment. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Read data sets into R, save them, and examine the contents</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Invoke the R environment and examine the R workspace</li><li>• Read tables created in Lab 1 into the R statistical package</li><li>• Examine, manipulate and save data sets</li><li>• Exit the R environment</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>. See the Appendix for:</p> <ul style="list-style-type: none"><li>• R Commands – Quick Reference</li><li>• Surviving LINUX – Quick Reference</li></ul>

Step	Action
1	<p><b><u>Invoke the R Environment:</u></b></p> <p>Logon to RStudio environment.</p> <ol style="list-style-type: none"> <li>1. The RStudio is accessed through the “safari” browser available as a desktop icon on your “fe” host.</li> <li>2. To start an RStudio session, start “safari” on the “fe” host and direct it to the URL <a href="http://&lt;IP Address of your be host&gt;:8787/">http://&lt;IP Address of your be host&gt;:8787/</a></li> <li>3. RStudio access details are as follows:  User-id : gpadmin  Password : &lt;supplied by your instructor&gt;</li> </ol> <p>With a successful login to the back-end, you should see the standard RStudio four-panel display.</p> <ol style="list-style-type: none"> <li>1. Verify that you see the following text in the lower left-hand pane:</li> </ol> <pre>R version 2.13.1 (2011-04-13) Copyright (C) 2011 The R Foundation for Statistical Computing ISBN 3-900051-07-0 Platform: i386-pc-mingw32/i386 (32-bit)  --- &lt;snip&gt; ---  Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R &gt;</pre>
2	<p><b><u>Examine the Workspace:</u></b></p> <p>Type the following command into the R command panel, and hit [ENTER]</p> <pre>ls()</pre> <p>You should see the following:</p> <pre>character(0)</pre> <p><b>Note:</b> R is telling you that you have nothing in your workspace.</p>

Step	Action
3	<p><b><u>Getting Familiar with R</u></b></p> <ol style="list-style-type: none"> <li>1. Click each tab in each panel. What happens?</li> <li>2. Type the following commands into the R command panel</li> </ol> <pre>help() help.start() demo() demo(graphics)</pre> <p>Hit esc to exit out of the demo</p>
4	<p><b><u>Read-in the Lab1 Script</u></b></p> <ol style="list-style-type: none"> <li>1. Now, in the script window, open the script called “Module3Lab1.R”. (Click on “File”, “Open File” and Navigate to directory LAB02 and click on file “Module3Lab1.R”. )</li> <li>2. All the commands we will be executing in this lab are contained in this script. In order to execute a command, do the following: <ul style="list-style-type: none"> <li>o Position your cursor inside the line that represents the command you wish to execute.</li> <li>o Either click on the “Run” button, or hit [CTRL-Enter]. You can execute many commands at once by selecting a sequence of commands and then issuing the “Run” command.</li> </ul> </li> <li>3. The command will be executed in the command pane. If the command produces graphical output, it will appear in the graphic frame. Note that you can expand this panel by clicking on the “expand window” box. In some instances, this will show more information that has been hidden because of the size of the panel.</li> </ol> <ul style="list-style-type: none"> <li>• The (<i>Module3Lab1.R</i>) file is divided into sections. Each section corresponds to a step in this lab. By selecting an individual line or lines, you can click “Run...” and the command(s) will be executed in the R panel.</li> </ul> <ol style="list-style-type: none"> <li>1. On the 1<sup>st</sup> line in Section 1, put your cursor on the line containing the word ls().</li> <li>2. Click <b>Run</b>. The ls() command will execute in the command window and show you the contents of your workspace.</li> </ol>

Step	Action
5	<p><b><u>Working with R:</u></b></p> <p>Load the .txt files you created in the first lab. Load the first file, <b>lab1_01.txt</b></p> <ol style="list-style-type: none"> <li>1. Set the working directory to LAB01 where we have stored the data. On the console window type:  <pre>setwd("~/LAB01")</pre> </li> <li>2. Select the line and press &lt;ctl&gt;Enter:  <pre>lab1 &lt;- table.read("lab1_01.txt", sep=" ", header=TRUE)</pre> <ul style="list-style-type: none"> <li>• If correct, R will simply return you to the command prompt ("&gt; ").</li> </ul> </li> <li>3. Now load the second .txt file, <b>lab1_02.txt</b>, by modifying the command (using the line of code in the RStudio command panel) you just entered.            (Use the up/down, left/right arrow buttons to move from and within lines; change each occurrence of "lab1" to "lab2".)            The command should read:  <pre>lab2 &lt;- read.table("lab1_02.txt", sep=" ", header=TRUE)</pre> </li> <li>4. When you have completed the edits, make sure that your cursor is within the line, press <b>Enter</b>.</li> </ol> <p><b>Note:</b> R supports copy and paste, as well as up and down arrows for moving to previous commands, left and right arrows to move within/between lines and home/end to move to the beginning or end of a line.</p>
6	<p><b><u>Verify the Contents of the Tables:</u></b></p> <p>It is always a good idea to look at the data to make sure that everything works. You can use the <b>head()</b> command to print out the first 6 lines of a table or the, <b>tail()</b> command to print out the last 6 lines of the table.</p> <ol style="list-style-type: none"> <li>1. Select and run the command:  <pre>head(lab1, n=10)</pre>           Record the value of the 10<sup>th</sup> line here: _____         </li> <li>2. Now do the same for the lab2 table, but use the <b>tail(lab2, n=10)</b> command instead.</li> <li>3. Record the value of the 1<sup>st</sup> line here: _____</li> </ol>



Step	Action
7	<p><b><u>Manipulating Data Tables (data frames) in R:</u></b></p> <p>Examine the contents of the table in more detail.</p> <p>1. Execute the following command:</p> <pre><b>summary(lab1)</b></pre> <p>Ignore the values for the <i>hinc</i> and <i>rooms</i> columns for now. The <i>serialnoid</i> field represents a unique identifier (it's the household identifier) from the Postgres database. You no longer need it and it will interfere with some of the procedures you want to run against this data set, so create a copy of the lab1 table without that column.</p> <p>2. Select and run:</p> <pre><b>nlab1 &lt;- lab1[,2:3]</b></pre> <p>This uses a feature of R that allows us to refer to rows and columns in a dataframe as if they were entries in a matrix. A blank entry in a row or column position means "use all available." This statement says: use all the rows in the table, but only use columns 2 and 3</p> <p>You could have used the following for the same effect (Note that the following code is not part of the script you can see in the source file <i>Module3lab1.R</i>):</p> <pre><b>hinc &lt;- lab1\$hinc</b> <b>rooms &lt;- lab1\$rooms</b> <b>nlab11 &lt;- data.frame(hinc, rooms)</b></pre> <p>You're taking advantage of R behavior that names the columns after the name of the variable. You could have used the following for the same effect:</p> <pre><b>nlab11 &lt;- data.frame(lab1\$hinc, lab1\$rooms)</b> <b>names(nlab11) = c("hinc", "rooms")</b></pre>

Step	Action
7 Cont.	<p>3. The <code>dim(&lt;table&gt;)</code> has the nice property of telling us how many rows exist in the table. Execute the following commands:</p> <pre>dim(nlab1) typeof(nlab1) class(nlab1)</pre> <p>Each of these commands tells us something about this particular object. You may not use these often, but they can be useful when R complains that it doesn't like something about the object that you just used.</p>
8	<p><b><u>Continue to Investigate Your Data:</u></b></p> <p>1. Select and execute the following commands:</p> <pre>summary(nlab1) cor(nlab1)</pre> <p>The summary function for data frames prints out summary statistics.</p> <p>2. Compare the median and the mean. What does it mean if the mean is less than the median? _____</p> <p>3. How about the mean greater than the median? _____</p> <p>4. Does the min and max value for the quartiles make sense to you?</p> <p>Here again you have a chance to do further cleaning of your data sets, but postpone this until you've finished the next few lessons.</p> <p>5. How do the values returned by the <code>cor()</code> function differ from the results obtained in lab 1? _____</p>
9	<p><b><u>Save the Data Sets:</u></b></p> <p>1. Execute the following commands:</p> <pre>rm(lab1) lab1 &lt;- nlab1 save(lab1, lab2, file="Labs.Rdata") rm(lab1, lab2) ls()      # make sure they're not in the workspace</pre>

Step	Action
10	<p><b><u>Examine Your Data:</u></b></p> <ol style="list-style-type: none"> <li>Experiment with some of the examples used in the lecture portion of this lesson. Using the same selection techniques that you used earlier, run each line in the file. <ul style="list-style-type: none"> <li>Some commands don't print their results. If this is the case, type in the value of the variable you created in the command window. If the variable was named "x", you can type "x". You can also type "print(x)" which will do the same thing.</li> </ul> </li> <li>Experiment with R functions that identify the class and data type of a particular variable, type: <p><b><code>typeof(x) , class(x) , attributes(x) , names(x) , dim(x)</code></b></p> </li> <li>Which ones work on which kind of data types? _____</li> <li>Type these values into the RStudio command panel. _____</li> <li>Typing all these commands for each variable is tedious. Alternatively, we will write a function <i>tellme</i> that takes a variable as an argument and performs <code>typeof</code>, <code>class</code>, <code>names</code> and <code>str</code> on that variable. Select and run the lines beginning with "<b><code>tellme &lt;- function(x)</code></b> <b><code>{</code></b> <i>extending</i> through the right curly brace.</li> <li>Now execute the following command <b><code>tellme</code></b> You should see the definition of the function that you just entered! This is because R doesn't interpret a plain <b><code>tellme</code></b> as a function, but rather as an object to be printed out. The default print function for a function is to print its definition. You can try this with any other R function. Type <b><code>mean</code></b> and inspect the results.</li> <li>Try <b><code>tellme ()</code></b> with a series of variables.</li> <li>Which commands actually list something? _____</li> <li>How might you get the other commands to list their return value? [Hint: try <code>print ()</code>]</li> </ol>
11	<p><b><u>Exit R:</u></b></p> <ol style="list-style-type: none"> <li>Execute the following command: <b><code>q()</code></b></li> <li>R will ask you if you want to save your workspace. Answer "<b><code>no</code></b>".</li> </ol>

*End of Lab Exercise*

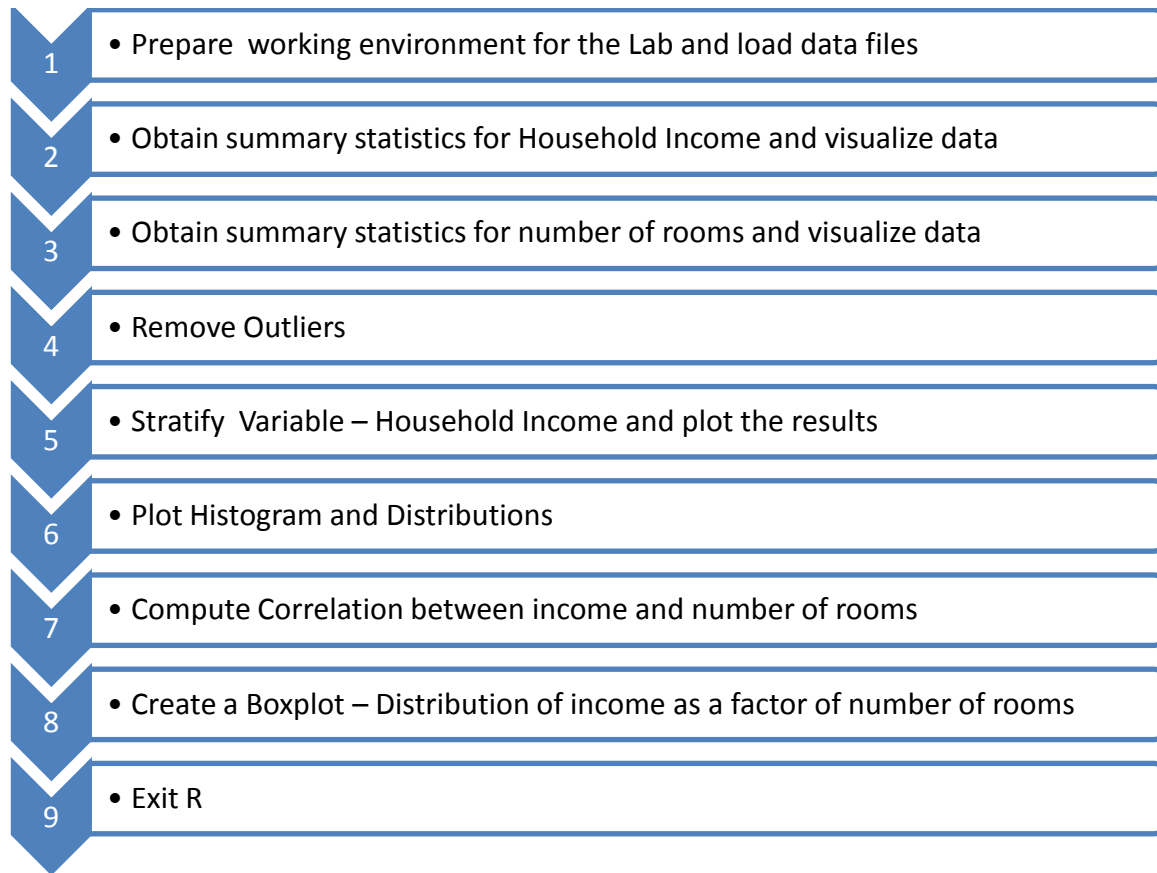


## Lab Exercise 3: Basic Statistics, Visualization, and Hypothesis Tests

<b>Purpose:</b>	<p>The lab introduces you to the analysis of data using the R statistical package within the Data Science and Big Data Analytics environment. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Perform summary (descriptive) statistics on the data sets</li><li>• Create basic visualizations using R both to support investigation of the data as well as exploration of the data</li><li>• Create plot visualizations of the data using a graphics package</li><li>• Test a hypothesis about the data</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Reload data sets into the R statistical package</li><li>• Perform summary statistics on the data</li><li>• Remove outliers from the data</li><li>• Plot the data using R</li><li>• Plot the data using lattice and ggplot</li><li>• Test a hypothesis about the data</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>. See the Appendix for:</p> <ul style="list-style-type: none"><li>• R Commands – Quick Reference</li><li>• Surviving LINUX – Quick Reference</li></ul>

## Part 1 – Basic Statistics and Visualization Using R

### Workflow Overview



## LAB Instructions

Step	Action
1	<p><b><u>Prepare working environment for the Lab and load data files</u></b></p> <ol style="list-style-type: none"> <li>Set the working directory to LAB01 where we have stored the data. On the console window type:  <pre>setwd("~/LAB01")</pre> </li> <li>In the script window, open the script called "Module3Lab2.R". (Click on "File", "Open File" and Navigate to directory LAB03 and click on file "Module3Lab2.R"). Start R and Read the Data Set Back Into Your Workspace:</li> <li>Execute the following commands from the script window:   <pre>options(digits=3) options(width=68)  ls() load(file="Labs.Rdata") ls()  rm(lab2)  ds &lt;- lab1 colnames(ds) &lt;- c("income", "rooms")</pre> </li> </ol>
2	<p><b><u>Obtain summary statistics for Household Income and visualize data:</u></b></p> <ol style="list-style-type: none"> <li>Execute the following commands from the script window:   <pre>summary(ds\$income) range(ds\$income) sd(ds\$income) var(ds\$income)  plot(density(ds\$income)) # left skewed</pre> </li> <li>What is the mean? _____</li> <li>What is the median? _____</li> <li>What is the standard deviation? _____</li> </ol>

Step	Action
3	<p><b><u>Obtain summary statistics for Number of rooms and visualize data:</u></b></p> <p>Execute the following commands from the script window:</p> <pre>summary(ds\$rooms) range(ds\$rooms) sd(ds\$rooms) plot(as.factor(ds\$rooms))</pre> <p>What is the mean?</p> <p>What is the median?</p> <p>What is the standard deviation?</p>
4	<p><b><u>Remove Outliers</u></b></p> <p>In a previous lab, you recorded the range of income. You observed that the minimum household income is 4, and the maximum is 1,620,560.</p> <ol style="list-style-type: none"> <li>Does this make sense to you? Why? *</li> <li>What happens if you throw out the top and bottom 10%? Execute the following line from the script window <pre>(m &lt;- mean(ds\$income, trim=0.10) )</pre> </li> <li>How does this compare to the previous mean of this variable?</li> <li>Execute the following commands from the script window: <pre>ds &lt;- subset(ds, ds\$income &gt;= 10000 &amp; ds\$income &lt; 1000000) summary(ds) quantile(ds\$income, seq(from=0, to=1, length=11))</pre> </li> <li>How do these values vary from the values in the original data set?</li> <li>Do they make more sense?</li> <li>Which data set would you prefer to use?</li> </ol> <hr/> <p>*We might consider the high and low value as outliers, and get rid of them. On the other hand, as we will discover, income is best described via a lognormal distribution, and hence these values are in the extreme ends <math>\pm 3</math> sds from the mean.</p>



Step	Action
5	<p><b><u>Stratify Variable – Household Income and plot the results:</u></b></p> <p>Stratify breaks that occur close to U.S. Guidelines for Poverty, Median Income, Wealth, and Rich (&gt; \$250k @ year)</p> <ol style="list-style-type: none"> <li>Execute the following code (listed under comment heading “step 5” in the script file): <pre>breaks &lt;- c(0, 23000, 52000, 82000, 250000, 999999) labels &lt;- c("Poverty", "LowerMid", "UpperMid", "Wealthy", "Rich") wealth &lt;- cut(ds\$income, breaks, labels) # add wealth as a column to ds ds &lt;- cbind(ds, wealth) # show the 1<sup>st</sup> few lines. head(ds)</pre> </li> <li>Continue to execute the remaining part of the code in Step 5 <pre>wt &lt;- table(wealth) percent &lt;- wt/sum(wt)*100 wt &lt;- rbind(wt, percent) wt plot(wt)</pre> </li> <li>Take another look at the relationship between wealth and income. Execute the following lines: <pre># take another look -- wealth by rooms  nt &lt;- table(wealth, ds\$rooms) print(nt) plot(nt)          # nice mosaic plot</pre> </li> <li>Execute this code from the script file. These lines will remove the variables wealth, breaks and labels, and then save the variables data set and write into a file named “Census.Rdata”. <pre>rm(wealth,breaks,labels) save(ds, wt, nt, file="Census.Rdata")</pre> </li> </ol>

Step	Action
6	<p><b><u>Plot Histogram and Distributions:</u></b></p> <p>Problem: How do you represent income given the range of values?</p> <ol style="list-style-type: none"> <li>1. Select and execute the code under Step 6 Histograms and distributions in the script file.</li> </ol> <pre>library(MASS)  with(ds, {   hist(income, main="Distribution of Household Income",     freq=FALSE)   lines(density(income), lty=2, lwd=2) # line type (lty) 2 is dashed   xvals = seq(from=min(income), to=max(income),     length=100)   param = fitdistr(income, "lognormal")   lines(xvals, dlnorm(xvals, meanlog=param\$estimate[1],     sdlog=param\$estimate[2]), col="blue") })</pre> <ol style="list-style-type: none"> <li>2. Now try the same thing with log10(income)</li> </ol> <pre>logincome = log10(ds\$income) hist(logincome, main="Distribution of Household Income",   freq=FALSE) # line type lty(2) is a dashed line lines(density(logincome), lty=2, lwd=2) xvals = seq(from=min(logincome), to=max(logincome),   length=100) param = fitdistr(logincome, "normal") lines(xvals, dnorm(xvals, param\$estimate[1],   param\$estimate[2]), lwd=2, col="blue")</pre>

Step	Action
7	<p><b><u>Compute Correlation between income and number of rooms:</u></b></p> <ol style="list-style-type: none"> <li>You need to consider your hypothesis. <ul style="list-style-type: none"> <li>Your hypothesis is that the number of rooms in a house is predicted by household income (the rich can buy bigger houses), e.g. <math>lm(\text{rooms} \sim \text{income})</math></li> <li>Therefore, our null hypothesis: no correlation between income and number of rooms.</li> <li>Alternate hypothesis: there is a correlation between income and the number of rooms.</li> </ul> </li> <li>Execute the following code (listed after the comment line "Step7 in the script file"). <pre>with(ds, cor(income, rooms))  with(ds, cor(log(income), rooms)) # this will give a better correlation</pre> </li> <li>For comparison, correlate rooms with a completely unrelated variable. <pre>n = length(ds\$income) with(ds, cor(runif(n), rooms))</pre> </li> </ol>
8	<p><b><u>Create a Boxplot - Distribution of income as a factor of number of rooms:</u></b></p> <ol style="list-style-type: none"> <li>Select and execute the code (Listed after the comment line "Step 8") in the script window.</li> <li>Plot the distribution of income as a factor of # of rooms. 'log="y"' plots income on log scale. We will suppress the outlier points and let the whiskers cover the full range of the data. <pre>boxplot(income ~ as.factor(rooms), data=ds, range=0, outline=F, log="y", xlab="# rooms", ylab="Income")</pre> </li> <li>Plot the # of rooms as a function of wealth level. <pre>boxplot(rooms ~ wealth, data = ds, main="Room by Wealth", Xlab="Category", ylab="# rooms")  # we'll keep the outlier points in this one</pre> </li> </ol>

Step	Action
9	<p><b><u>Exit R:</u></b></p> <ol style="list-style-type: none"> <li>1. Type the following command into the RStudio command window:   <code>q()</code></li> <li>2. R will ask you if you want to save your workspace. Answer “<b>no.</b>”</li> </ol>

*End of Lab Exercise*

## Part 2 – Graphics Package Plots and Hypothesis Tests

### Workflow Overview



## Lab Instructions

Step	Action
1	<p><b><u>Define problem - Analysis of Variance (ANOVA):</u></b></p> <p>Suppose we are evaluating our marketing department's incentive campaign that is trying to increase the amount of money that customers spend when they visit our online site. We ran a short experiment, where visitors to our site randomly received one of two incentive offers or got no offer at all.</p>
2	<p><b><u>Generate the Data:</u></b></p> <pre>offers = sample(c("noffer", "offer1", "offer2"), size=500, replace=T)</pre> <pre>purchasesize = ifelse(offers=="noffer", rlnorm(500, meanlog=log(25)), ifelse(offers=="offer1", rlnorm(500, meanlog=log(50)), rlnorm(500, meanlog=log(55))))</pre> <pre>offertest = data.frame(offer=offers, purchase_amt=purchasesize)</pre>
3	<p><b><u>Examine the Data:</u></b></p> <pre>summary(offertest)</pre> <p>The following command does the equivalent of the SQL command "SELECT avg(purchase_amt) FROM offertest GROUP BY offer",</p> <pre>aggregate(x=offertest\$purchase_amt, by=list(offertest\$offer), FUN="mean")</pre>
4	<p><b><u>Plot and determine how purchase size varies within the three groups:</u></b></p> <p>1. The 'log="y"' argument plots the y axis on the log scale. Does it appear that making offers increases purchase amount?</p> <pre>boxplot(purchase_amt ~ as.factor(offers), data=offertest, log="y")</pre>

Step	Action
5	<p><b><u>Use lm() to do the ANOVA:</u></b></p> <p>1. Execute the following commands:</p> <pre>model = lm(log10(purchase_amt) ~ as.factor(offers), data=offertest)  summary(model)</pre> <p>2. What is the p-value on the F-stat? Can we reject the null hypothesis? _____</p> <p>The intercept of the model is the mean value of <math>\log_{10}(\text{purchase\_amt} \mid \text{no offer})</math>, (call it <math>m_0</math>) and the other coefficients are:</p> <p style="padding-left: 40px;"><math>\text{mean}(\log_{10}(\text{purchase\_amt} \mid \text{offer1})) - m_0</math>, and</p> <p style="padding-left: 40px;"><math>\text{mean}(\log_{10}(\text{purchase\_amt} \mid \text{offer2})) - m_0</math>, respectively.</p> <p>3. What are the p-values on those coefficients? _____</p> <p>4. Can we reject the null hypotheses that the mean purchase amount for offer1 was different from that of no offer, and similarly for offer2 vs. no offer? _____</p>
6	<p><b><u>Use Tukey's test to check all the differences of means:</u></b></p> <p>1. Execute the following command:</p> <pre>TukeyHSD(aov(model))</pre> <p>1. Did offer1 and offer2 increase purchase size to different amounts (to the <math>p &lt; 0.05</math> significance level)? _____</p> <p>2. What would you recommend to the marketing department, based on these results? _____</p>

Step	Action
7	<p><b><u>Use the lattice package for density plot:</u></b></p> <p>For this course, you are only expected to become familiar with the base graphics capabilities of R; however, there are other graphics packages available for R that makes certain kinds of visualizations easier to produce. If you continue to use R in the future, it will be helpful to be aware of these alternatives to base graphics.</p> <p>The lattice package makes it easy to split data into different groups to highlight the differences between the groups. Here, we split the purchase_amt data by offer, and plot the three offer-specific purchase_amt densityplots on the same graph.</p> <pre>library(lattice)  densityplot(~ purchase_amt, group=offers, data=offertest, auto.key=T)</pre>
8	<p><b><u>Plot the Logarithms of the Data:</u></b></p> <p>1. Because the data is so left-skewed, we may want to plot the logarithms of the data to more clearly see the differences in the distributions, and the different locations of the modes.</p> <pre>densityplot(~ log10(purchase_amt), group=offers, data=offertest, auto.key=T)</pre> <p>2. Also try the plots:</p> <pre>densityplot(~purchase_amt   offers, data=offertest)  densityplot(~log10(purchase_amt)   offers, data=offertest)</pre> <p>3. Which style of graph do you find more helpful?</p>



Step	Action
10	<p><b><u>Use ggplot() package:</u></b></p> <p>The ggplot2 package is based on a theory of the “algebra of graphs”. The syntax is rather complex, but ggplot excels at assembling rich composite graphs that use a variety of different graphic techniques. Here, we show how to produce a variation of the scatterplot + box-and-whisker plot that we saw earlier in the course, to plot the distributions of purchase amounts against offer.</p> <p>1. Execute the following commands:</p> <pre>library(ggplot2)  ggplot(data=offertest, aes(x=as.factor(offers), y=purchase_amt)) + geom_point(position="jitter", alpha=0.2) + geom_boxplot(alpha=0.1, outlier.size=0) + scale_y_log10()</pre> <p>The function geom_point() plots scatterplots. The function geom_boxplot() plots box-and-whisker plots; outlier.size()=0 removes the outlier points beyond the whiskers that normally would be plotted. The function scale_y_log10() plots the y axis on a log10 scale.</p> <p>2. You need to plot at least one geom_xxx to get a graph. Try adding and removing the different terms of the graphing command to create simpler scatterplots or box-and-whisker plots, with and without log scaling.</p> <p>3. Here’s how you would create the densityplots that you created in lattice. Execute the following commands.</p> <pre>ggplot(data=offertest) + geom_density(aes(x=purchase_amt, colour=as.factor(offers)))  ggplot(data=offertest) + geom_density(aes(x=purchase_amt, colour=as.factor(offers))) + scale_x_log10()</pre>

Step	Action
11	<p><b><u>Generate the example data to perform a Hypothesis Test with manual calculations:</u></b></p> <p>Hopefully, you won't have to do this too often. Most statistical packages have functions that calculate a test statistic and evaluate it against the proper distribution, for the most common hypothesis tests. On occasion, you may need to calculate the p-values yourself. For our example, we will calculate the Student's t-test for difference of means (unlike Welch's test, Student's t-test assumes identical variances), under the alternative hypothesis that the means are not equal.</p> <p>1. Select and execute the following commands:</p> <pre>x = rnorm(10) # distribution centered at 0 y = rnorm(10,2) # distribution centered at 2</pre>
12	<p><b><u>Create a function to calculate the pooled variance, which is used in the Student's t statistic:</u></b></p> <p>1. Select and execute the following commands. This will create a function named <i>pooled.var</i>.</p> <pre>pooled.var = function(x, y) {   nx = length(x)   ny = length(y)   stdx = sd(x)   stdy = sd(y)   num = (nx-1)*stdx^2 + (ny-1)*stdy^2   denom = nx+ny-2 # degrees of freedom   (num/denom) * (1/nx + 1/ny) }</pre>
13	<p><b><u>Examine the Data:</u></b></p> <p>Select and execute the following commands:</p> <pre>mx = mean(x) my = mean(y)  mx - my  pooled.var(x,y)</pre>

Step	Action
14	<p><b><u>Calculate the t statistic for Student's t-test:</u></b></p> <p>1. Select and execute the following commands:</p> <pre>tstat = (mx - my) / sqrt(pooled.var(x,y)) tstat</pre>
15	<p><b><u>Calculate the degrees of freedom:</u></b></p> <p>Under the null hypothesis, the t statistic is distributed in a Student's distribution with <math>n_x + n_y - 2</math> degrees of freedom. Calculate the degrees of freedom for our problem.</p> <p>Select and execute the following commands:</p> <pre>dof = length(x) + length(y) - 2 dof</pre>
16	<p><b><u>Compute the area under the curve:</u></b></p> <p>The function <math>pt(x, dof)</math> gives the area under the curve from <math>-\infty</math> to <math>x</math> for the Student's distribution with <math>dof</math> degrees of freedom. Since in this case we have a negative <math>tstat</math>, <math>pt(tstat, dof)</math> will give us the area under the left tail.</p> <p>1. Select and execute the following commands:</p> <pre>tailarea = pt(tstat, dof)</pre> <p>2. Since our null hypothesis is that <math>m_1 \neq m_2</math>, we need the area under both tails.</p> <pre>pvalue = 2*tailarea</pre> <p>3. Are the means different (to the <math>p &lt; 0.05</math> significance level)? _____</p>
17	<p><b><u>Perform Student's t-test directly and compare the results:</u></b></p> <p>1. Execute the following command:</p> <pre>t.test(x, y, var.equal=T)</pre> <p>2. Does <math>t.test()</math> give the same results? _____</p>

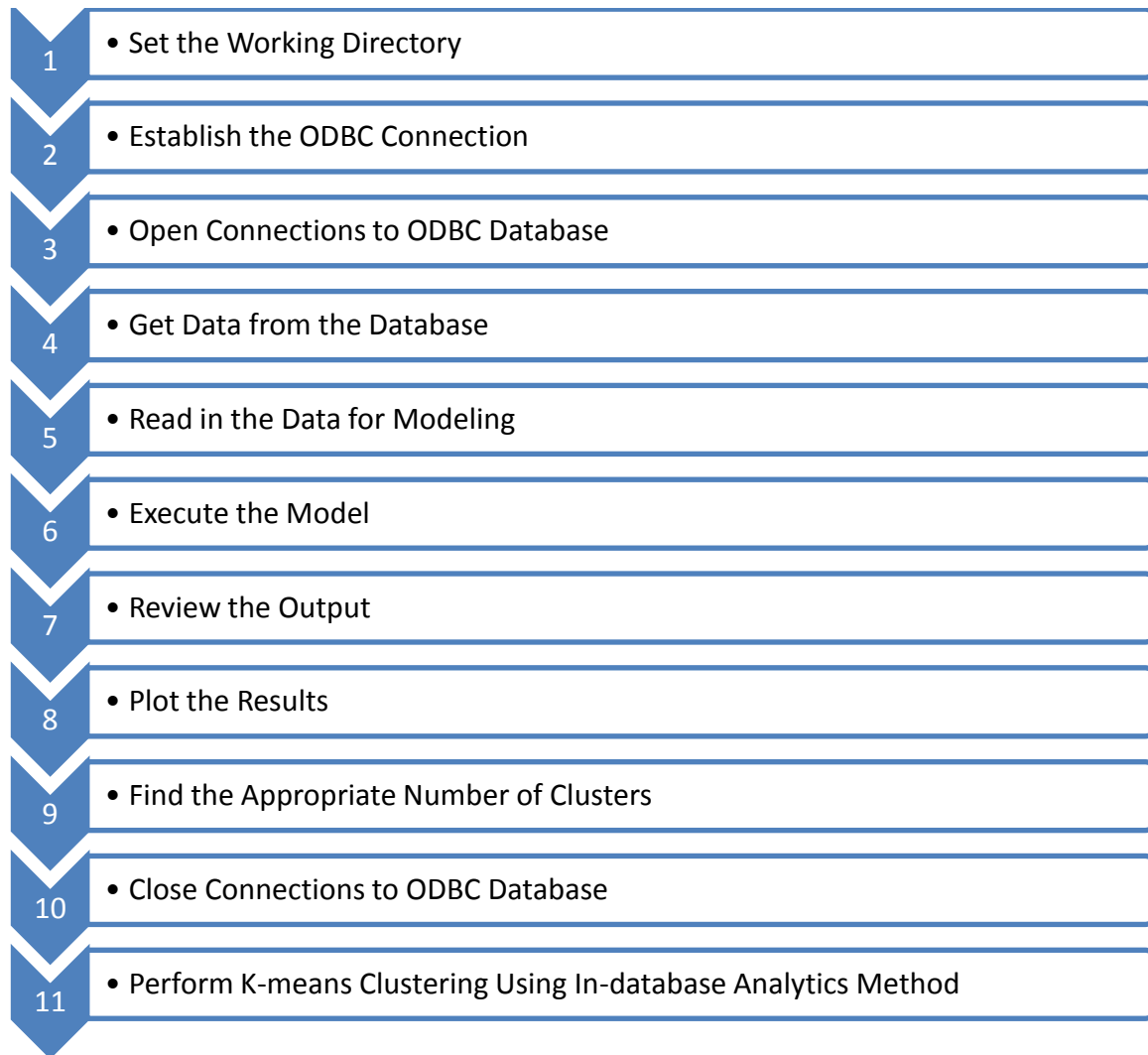
*End of Lab Exercise*



## Lab Exercise 4: K-means Clustering

Purpose:	<p>This lab is designed to investigate and practice K-means Clustering. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions to create K-means Clustering models</li><li>• Use ODBC connection to the database and execute SQL statements and load datasets from the database in an R environment</li><li>• Visualize the effectiveness of the K-means Clustering algorithm using graphic capabilities in R</li><li>• Use MADlib functions for K-means clustering</li></ul>
Tasks:	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Use the R –Studio environment to code K-means Clustering models</li><li>• Use the ODBC connection in the R environment to create the average household income from the census database as test data for K-means Clustering</li><li>• Use R graphics functions to visualize the effectiveness of the K-means Clustering algorithm</li><li>• Use MADlib functions for K-means clustering</li></ul>
References:	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>. <a href="http://www.statmethods.net/advstats/cluster.html">http://www.statmethods.net/advstats/cluster.html</a> (originally from Everitt &amp; Hothorn).</p>

## Workflow overview



## Lab Instructions

Step	Action
1	Log in with GPADMIN credentials onto R-Studio.
2	<p><b><u>Set the Working Directory:</u></b></p> <p>1. Set working directory to ~/LAB04/, execute the command:</p> <pre>setwd("~/LAB04")</pre> <ul style="list-style-type: none"> <li>(Or use the “Tools” option in the tool bar in the RStudio environment.)</li> </ul>
3	<p><b><u>Establish the ODBC Connection:</u></b></p> <p>Load the RODBC package, type in:</p> <pre>library('RODBC')</pre>
4	<p><b><u>Open Connections to ODBC Database:</u></b></p> <p>Before connecting to the ODBC database make sure the file, /etc/odbc.ini, is set to point to database, “training2”.</p> <ol style="list-style-type: none"> <li>If not, edit the line that starts with "Database =" within the file /etc/odbc.ini, to point to “training2”.</li> <li>Ensure the username(uid) and password (pwd) are provided correctly in the following command – note, you need to supply the correct password in the “pwd” field below:</li> </ol> <pre>ch &lt;- odbcConnect("Greenplum",uid="gpadmin",   case="postgresql",pwd="password_of_the_gpadmin_user")</pre>

Step	Action
5	<p><b><u>Get Data from the Database:</u></b></p> <ol style="list-style-type: none"> <li>Before creating the table, “income_state”, you must first delete the table, if it already exists. Type in:   <pre>sqlDrop(ch,"income_state")</pre> </li> <li>Use the sqlQuery command to create the table, “income_state” :   <pre>&gt; sqlQuery(ch, "CREATE TABLE income_state AS SELECT   f.name AS state , round(avg(h.hinc),0) AS income FROM   housing AS h JOIN   fips AS f ON   h.state = f.code WHERE   h.hinc &gt; 0 GROUP BY   f.name DISTRIBUTED BY (income); " )</pre> <p><b>Note:</b> This code creates the table, “income_state”, in database “training2”.</p> </li> <li>Inspect this table using the following command:   <pre>sqlColumns(ch,"income_state")</pre> </li> <li>Review the output on the console.</li> </ol> <p><b>Note:</b> The SQL Query is available for you to copy and paste into the working directory File name: mod4lab4.sql.</p>



Step	Action
6	<p><b><u>Read in the Data for Modeling:</u></b></p> <p>You need the data to be read in as a “matrix”.</p> <ol style="list-style-type: none"> <li>1. Execute the following statement to read in the database table “income_state”. Use the sqlFetch command. The “rownames” attribute ensures the data is rendered as a matrix and the row names are taken from the column “state”.</li> </ol> <pre>income &lt;- as.matrix(sqlFetch(ch,"income_state",                              rownames="state")) &gt; summary(income)</pre> <ol style="list-style-type: none"> <li>2. Review the results of “income” on the console window.</li> <li>3. Ensure that in the “data” window the variable “income” is represented as a 52x1 integer matrix.</li> </ol>
7	<p><b><u>Execute the Model:</u></b></p> <ol style="list-style-type: none"> <li>1. Sort the data “income” before the modeling process. This will make it easier to understand the results and in visualizing.</li> </ol> <pre>income &lt;- sort(income)</pre> <p>The K-Means function, provided by the <i>cluster</i> package, is used as follows:</p> <pre>kmeans(x, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"))</pre> <p>where the arguments are:</p> <ul style="list-style-type: none"> <li>• <b>x:</b> A numeric matrix of data, or an object that can be coerced to such a matrix (such as a numeric vector or a data frame with all numeric columns).</li> <li>• <b>centers:</b> Either the number of clusters or a set of initial (distinct) cluster centers. If a number, a random set of (distinct) rows in x is chosen as the initial centers.</li> <li>• <b>iter.max:</b> The maximum number of iterations allowed.</li> <li>• <b>nstart:</b> If <i>centers</i> is a number, <i>nstart</i> gives the number of random sets that should be chosen.</li> <li>• <b>algorithm:</b> The algorithm to be used. It should be one of the following "Hartigan-Wong", "Lloyd", "Forgy" or "MacQueen". If no algorithm is specified, the algorithm of Hartigan and Wong is used by default.</li> </ul> <ol style="list-style-type: none"> <li>2. Cluster the data into 3 groups (centers = 3) and also specify the number of random sets to start with as, 15.</li> </ol> <pre>&gt; # Fit the k-means cluster with 3 initial cluster centers &gt; km &lt;- kmeans (income,3,15)</pre>

Step	Action
8	<p><b><u>Review the Output:</u></b></p> <ol style="list-style-type: none"> <li>1. Use the following command to display the fitted model on the console:  <code>&gt; km</code></li> <li>2. What are the cluster means?</li> <li>3. What are the available components in the model?</li> <li>4. How many data points cluster into each group?</li> </ol> <p>The output from the model provides the following:</p> <ul style="list-style-type: none"> <li>• <b>cluster</b> A vector of integers (from 1:k) indicating the cluster to which each point is allocated.</li> <li>• <b>centers</b> A matrix of cluster centers.</li> <li>• <b>withinss</b> The within-cluster sum of squares for each cluster.</li> <li>• <b>totss</b> The total within-cluster sum of squares.</li> <li>• <b>tot.withinss</b> Total within-cluster sum of squares, that is, <code>sum(withinss)</code>.</li> <li>• <b>betweenss</b> The between-cluster sum of squares.</li> <li>• <b>size</b> The number of points in each cluster.</li> </ul>
9	<p><b><u>Plot the Results:</u></b></p> <ol style="list-style-type: none"> <li>1. Now plot the results using the following commands:  <pre>&gt; # plot clusters &gt; plot(income, col = km\$cluster) &gt; # plot centers &gt; points(km\$centers, col = 1:3, pch = 8)</pre></li> <li>2. Review the output on the graphic window.</li> </ol>
10	<p><b><u>Find the Appropriate Number of Clusters:</u></b></p> <ol style="list-style-type: none"> <li>1. Plot the within-group-sum of squares and look for an "elbow" of the plot. The elbow (if you can find one) tells you what the appropriate number of clusters probably is. (Adapted from <a href="http://www.statmethods.net/advstats/cluster.html">http://www.statmethods.net/advstats/cluster.html</a> (originally from Everitt &amp; Hothorn)).  <pre>wss &lt;- numeric(15) &gt; for (i in 1:15) wss[i] &lt;- sum(kmeans(income,   centers=i)\$withinss) &gt; plot(1:15, wss, type="b", xlab="Number of Clusters",   ylab="Within groups sum of squares")</pre></li> <li>2. Review the output on the graphic window. Is there an elbow to the plot?</li> <li>3. Repeat the modeling with a few values around the elbow (or 4 and 5 centers if there is no elbow) and review the results.</li> </ol>

Step	Action
11	<p><b><u>Close Connections to ODBC Database:</u></b></p> <p>Use the following command:</p> <pre><b>odbcClose (ch)</b></pre> <p>The R Code for this exercise is available at <a href="/home/gpadmin/LAB04/kmeans1.R">/home/gpadmin/LAB04/kmeans1.R</a></p>

Step	Action
12	<p><b><u>Perform K-means Clustering Using In-database Analytics Method:</u></b></p> <ol style="list-style-type: none"> <li>1. In-database analytics using “MADlib” function calls will be detailed in module 5 later in this course. However this step provides an overview of the flexibility to “cluster” the data points within the database environment. Review the code provided below: <pre> DROP TABLE IF EXISTS myschema.data; CREATE TABLE myschema.data (   pid INT   , position FLOAT8[]) DISTRIBUTED BY (pid);  INSERT INTO myschema.data (pid,position[1]) SELECT   h.state   , round(avg(h.hinc),0) FROM   housing AS h WHERE   h.hinc &gt; 0 GROUP BY   h.state ;  SET SEARCH_PATH to madlib,public,myschema;  SELECT madlib.kmeans('myschema.data', 3, 1, 'mytestrun', 'myschema');  SELECT * FROM myschema.kmeans_out_centroids_mytestrun; SELECT * FROM myschema.kmeans_out_points_mytestrun; </pre> </li> <li>2. The first part of the code is similar to the one created in step 5 of this lab. The K-means function is called by: <pre> SELECT madlib.kmeans( 'input_table', k,                       'goodness', 'run_id', 'output_schema'); </pre> <p>The centroid locations are stored in kmeans_out_centroids_(run_id):  The cluster assignments for each input data point are stored in kmeans_out_points_(run_id):  In the example above  input_table is myshema.data,  number of clusters is 3  goodness is 1  run_id is mytestrun and output schema is myschema.</p> <p>The code is as follows:</p> <pre> SELECT madlib.kmeans('myschema.data', 3, 1, 'mytestrun', 'myschema'); </pre> </li> </ol>

Step	Action
12 Cont.	<p>3. The outputs are available in tables  kmeans_out_centroids_madlib.mytestrun  kmeans_out_points_madlib.mytestrun</p> <p>in schema myschema. The code below will output the results:</p> <pre>SELECT * FROM myschema.kmeans_out_centroids_mytestrun; SELECT * FROM myschema.kmeans_out_points_mytestrun;</pre> <p>4. Review the results and compare them with the results generated with R in step 7.</p>
13	<p>The code is available /home/gpadmin/LAB04/kmeansmadlib.sql</p> <p>You can execute this code with the following command in the appropriate directory:</p> <pre>psql -d training2 -f kmeansmadlib.sql</pre>

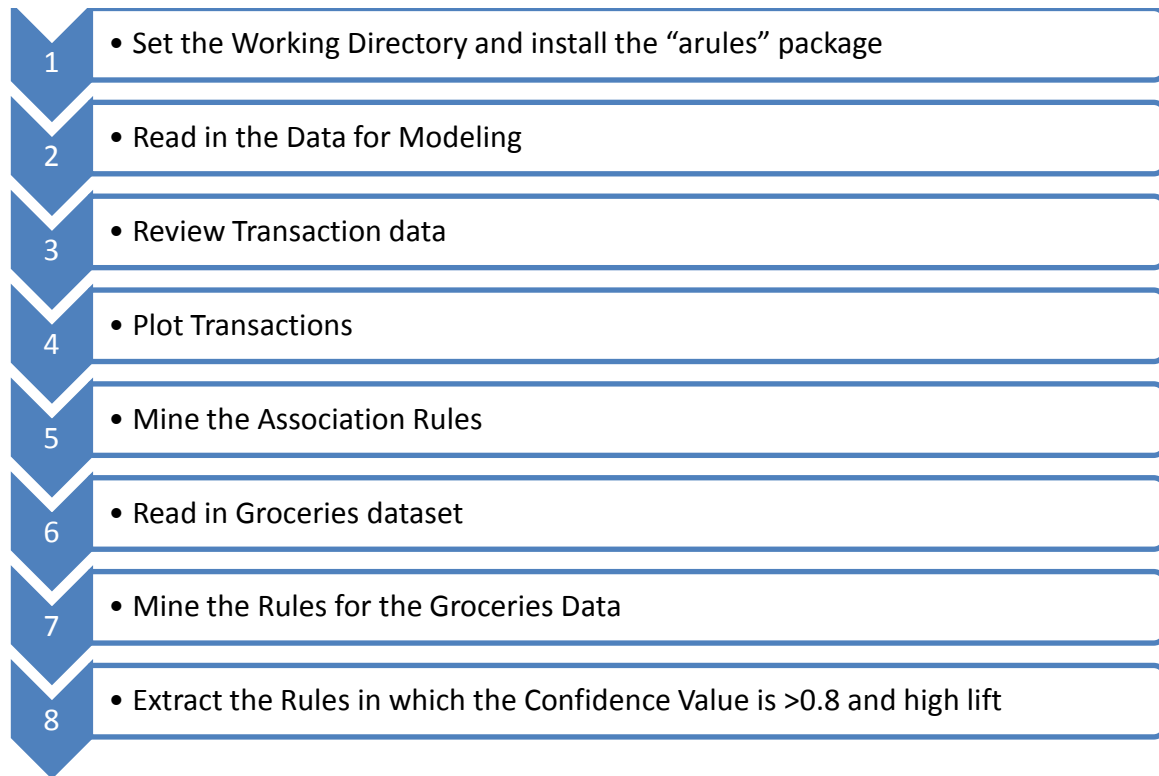
*End of Lab Exercise*



## 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> <p>1. Set the working directory to ~/LAB05/ by executing the command:</p> <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> <p>2. Load the package and the required libraries:</p> <pre>&gt; # Load libraries &gt; library('arules')</pre>

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>
5	<p><b><u>Plot Transactions:</u></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> <pre>&gt; image (txn)</pre> <p>2. Review the output in the graphics window</p>

Step	Action
6	<p><b><u>Mine the Association Rules:</u></b></p> <p>The “apriori” function, provided by the <i>arulesr</i> package, is used as follows:</p> <pre>rules &lt;- 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"> <li>• <b>data</b> object of class transactions or any data structure which can be coerced into transactions (for example, a binary matrix or data.frame).</li> <li>• <b>parameter</b> named list. The default behavior is to mine rules with support 0.1, confidence 0.8, and maxlen 5.</li> </ul> <p>1. Read in the statement for the transaction data:</p> <pre>&gt; #mine association rules &gt; basket_rules &lt;- apriori(txn,parameter=list(sup=0.5,conf=0.9,target="rules"))</pre> <p>2. Review the output on the console. The number of rules generated can be seen in the output and is represented as follows:</p> <pre>writing ... [1 rule(s)] done [0.00s]</pre> <p>3. Inspect the rule using the following statement:</p> <pre>&gt; inspect(basket_rules)</pre> <p>4. Review the output.</p> <p>5. State the generated rule and the support, confidence and the lift thresholds for the rule</p>

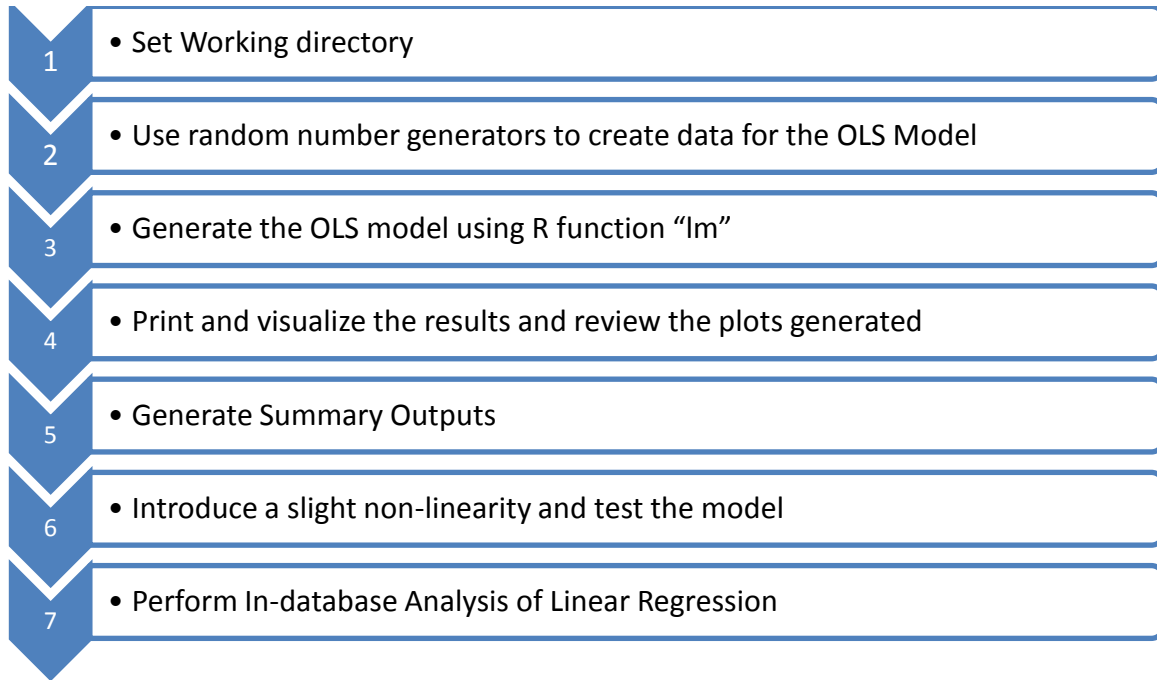
Step	Action
7	<p><b><u>Read in Groceries dataset</u></b></p> <p>Use the standard data set, “Groceries” available with the “arules” package.</p> <ul style="list-style-type: none"> <li>The Groceries data set contains 1 month (30 days) of real-world point-of-sale transaction data from a typical local grocery 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</pre>
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> <p>1. How many rules are generated?</p>
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; inspect(subrules)</pre> <p>2. Review the results.</p> <p>3. How many sub-rules did you extract?</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”.</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> <p>5. List the rules and the value of the parameters associated with these rules:</p> <ul style="list-style-type: none"> <li>The R Code for this exercise is available at <a href="#">home/gpadmin/LAB05/mba.R</a></li> </ul>

*End of Lab Exercise*

## Lab Exercise 6: Linear Regression

<b>Purpose:</b>	<p>This lab is designed to investigate and practice the Linear Regression method. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Linear Regression (Ordinary Least Squares – OLS)</li><li>• Predict the dependent variables based on the model</li><li>• Investigate different statistical parameter tests that measure the effectiveness of the model</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 OLS models</li><li>• Review the methodology to validate the model and predict the dependent variable for a set of given independent variables</li><li>• Use R graphics functions to visualize the results generated with the model</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>.</p>

## Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set Working directory</u></b>  Set the working directory to ~/LAB06/ by executing the command:</p> <pre>setwd("~/LAB06")</pre> <ul style="list-style-type: none"> <li>(Or using the “Tools” option in the tool bar in the RStudio environment).</li> </ul>
3	<p><b><u>Use random number generators to create data for the OLS Model :</u></b></p> <ol style="list-style-type: none"> <li>Run the “runif” function in R which generates random deviates within the specified minimum and maximum range.  <pre>x &lt;- runif(100, 0, 10)</pre> <p>This generates 100 random values for “x” in the range 0 to 10.</p> </li> <li>Create the dependent variable “y” with the “beta” values as 5 and 6 and the “sigma” = 1 (generated with the “rnorm” function, random generation for the normal distribution with mean =0 and SD= 1.)  <pre>y &lt;- 5 + 6*x + rnorm(100)</pre> </li> <li>Plot it  <pre>&gt; plot(x,y)</pre> </li> <li>Review the results in the graphics window</li> </ol>

Step	Action
4	<p><b><u>Generate the OLS model using R function “lm”:</u></b></p> <p>An OLS Model is generated with an R function call “lm”.  You can learn about “lm” with the following command on the console:</p> <pre>?lm</pre> <ol style="list-style-type: none"> <li>1. Generate an OLS Model using the following command:  <pre>d &lt;- lm(y ~ x)</pre></li> <li>2. Use the following command to display the structure of the object “d” created with the function call “lm”  <pre>str(d)</pre></li> <li>3. You can see the details of the fitted model. What are the values of the coefficients (Beta) in the model?</li> </ol>



Step	Action
5	<p><b><u>Print and visualize the results and review the plots generated</u></b></p> <ol style="list-style-type: none"> <li>1. Get the compact results of the model with the following command:  <pre>print(d)</pre> </li> <li>2. Visualize the model with the command:  <pre>par(mfrow=c(2,2)) plot(d)</pre> <p>The explanation of the plots are as follows:</p> <p><b>Residuals vs. Fitted:</b> you want to make sure that the errors are evenly distributed over the entire range of fitted values; if the errors are markedly larger (or biased either positively or negatively) in some range of the data, this is evidence that this model may not be entirely appropriate for the problem.</p> <p><b>Q-Q plot:</b> tests whether or not the errors are in fact distributed approximately normally (as the model formulation assumes). If they are, the Q-Q plot will be along the x=y line. If they aren't, the model may still be adequate, but perhaps a more robust modeling method is suitable. Also, the usual diagnostics (R-squared, t-tests for significance) will not be valid.</p> <p><b>Scale-Location:</b> a similar idea to Residuals v. Fitted; you want to make sure that the variance (or its stand-in, "scale") is approximately constant over the range of fitted values.</p> <p><b>Residuals vs. Leverage:</b> used for identifying potential outliers and "influential" points. Points that are far from the centroid of the scatterplot in the x direction (high leverage) are influential, in the sense that they may have disproportionate influence on the fit (that doesn't mean they are wrong, necessarily). Points that are far from the centroid in the y direction (large residuals) are potential outliers.</p> </li> <li>3. Here are some examples of plots that may be a little more intuitive, Type in the following:  <pre>&gt; ypred &lt;- predict(d) &gt; par(mfrow=c(1,1)) &gt; plot(y,y, type="l", xlab="true y", ylab="predicted y") &gt; points(y, ypred)</pre> <p>Review the results in the graphics window. The plot of predicted vs. true outcome can be seen there. The plot should be near the x=y line. Where it does not run along the x=y line indicates where the model tends to over-predict or under-predict. You can also use this plot to identify ranges where the errors are especially large. This information is similar to the Residuals vs. Fitted plot, but perhaps is more intuitive to the layperson.</p> <p>Note: The “predict” function requires the variables to be named exactly as in the fitted model.</p> </li> </ol>

Step	Action
6	<p><b><u>Generate Summary Outputs:</u></b></p> <ol style="list-style-type: none"> <li>For more detailed results type:  <pre>d1 &lt;- summary(d)</pre> <pre>print(d1)</pre> <p>Read the explanations given below from the summary output and note the values from the output on the console for each statistic detailed:</p> <p><b>coefficients</b> : the estimated value of each coefficient, along with the standard error. coefficient +/- 2*std.error is useful as a quick measure of confidence interval around the estimate.</p> <p><b>t-value</b>: coefficient/std.error, or how tight an estimate this is (compared to 0). If the "true" coefficient is zero (meaning this variable has no effect on the outcome), t-value is "small".</p> <p><b>Pr(&gt; t )</b>: the probability of observing this t-value if the coefficient is actually zero. You want this probability to be small. How small depends on the significance desired. Standard significances are given by the significance codes. So, for example "****" means that the probability that this coefficient is really zero is negligible.</p> <p><b>R-squared</b>: A goodness of fit measure: the square of the correlation between predicted response and the true response. You want it close to 1. Adjusted R-squared compensates for the fact that having more parameters tends to increase R-squared. Since we only have one variable here, the two are the same.</p> <p><b>F-statistic and p-value</b>. Used to determine if this model is actually doing better than just guessing the mean value of y as the prediction (the "null model"). If the linear model is really just estimating the same as the null model, then the F-statistic should be about 1. The p-value is the probability of seeing an F-statistic this large, if the true value is 1. Obviously, you want this value to be very small.</p> </li> <li>Type in the following command:  <pre>&gt; cat("OLS gave slope of ", d1\$coefficients[2,1],</pre> <pre>    "and an R-sqr of ", d1\$r.squared, "\n")</pre> </li> <li>Note the result you see on the console in the space below:</li> </ol>



Step	Action
8	<p><b><u>Perform In-database Analysis of Linear Regression:</u></b></p> <p>To illustrate an in-database execution of linear regression we make use of with approximately 64,000 rows of data aggregated from the census data. The following columns are generated for zip code and sex(male/female) tuples:</p> <p>zip code,</p> <p>mean age in the zip code,</p> <p>mean number of years of education,</p> <p>mean level of employment (Categorical level values were converted to a range from 1 - 3, three being a "high status" job, 1 being "low status"</p> <p>mean household employment in the zip code</p> <p>This data is stored in table "zeta" in the database "training2".</p> <ol style="list-style-type: none"> <li>1. Explore the table with Meta commands or with the pgadmin utility</li> <li>2. Review the MADlib documentation at <a href="http://doc.madlib.net/v0.2beta/group_grp_linreg.html">http://doc.madlib.net/v0.2beta/group_grp_linreg.html</a> and the following code:</li> </ol> <pre> DROP TABLE IF EXISTS zeta1; CREATE TABLE zeta1 (     depvar FLOAT8     , indepvar FLOAT8[]) DISTRIBUTED BY (depvar) ;  INSERT INTO zeta1 (     depvar     , indepvar[1]     , indepvar[2]     , indepvar[3]     , indepvar[4] ) SELECT     ln(meanhouseholdincome + 1)     , 1     , CASE         WHEN sex = 'M' THEN 0         WHEN sex = 'F' THEN 1     END AS sex     , meanage     , meanemployment FROM     zeta ; </pre>

Step	Action
8	<pre>SET SEARCH_PATH to madlib,public,myschema;</pre>
Cont.	<pre>SELECT (linregr(depvar,indepvar)).r2 FROM zetal; SELECT (linregr(depvar,indepvar)).coef FROM zetal; SELECT (linregr(depvar,indepvar)).std_err FROM zetal; SELECT (linregr(depvar,indepvar)).t_stats FROM zetal; SELECT (linregr(depvar,indepvar)).p_values FROM zetal;</pre> <p>We are predicting mean household income with drivers age, sex, and years of employment:</p> <p>Notice that we take the log of income (refer to the discussions in the student resource guide)</p> <p>The code is available at</p> <pre>/home/gpadmin/LAB06/madliblinear.sql</pre> <p>1. You can execute the code with the following command at the command prompt:</p> <pre>psql -d training2 -f madliblinear.sql</pre> <p>2. Review the results and note your observations below:</p>

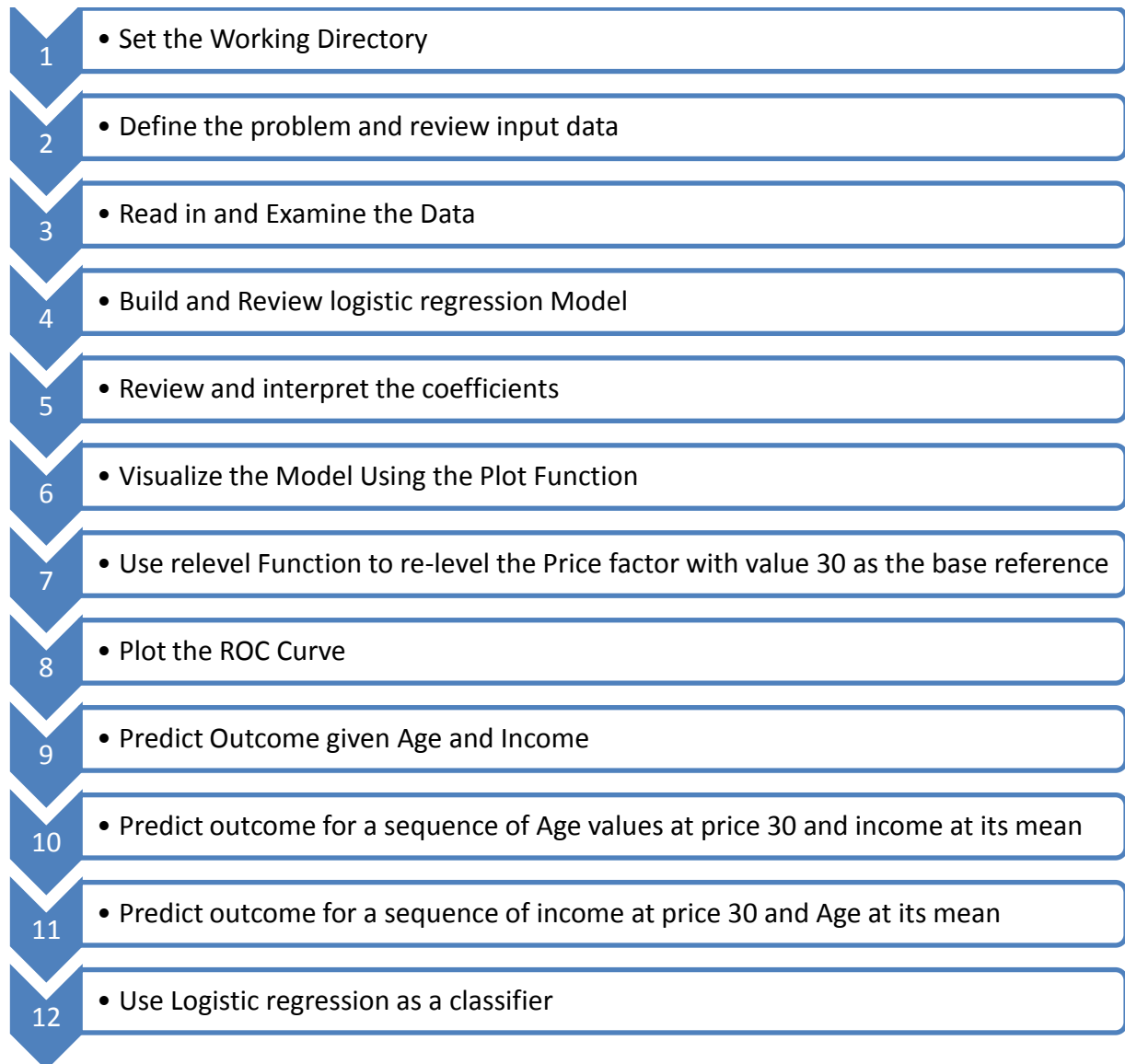
*End of Lab Exercise*



## Lab Exercise 7: Logistic Regression

<b>Purpose:</b>	<p>This lab is designed to investigate and practice the Logistic Regression method. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Logistic Regression – <i>also known as Logit</i>)</li><li>• Predict the dependent variables based on the model</li><li>• Investigate different statistical parameter tests that measure the effectiveness of the model</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Use R –Studio environment to code Logit models</li><li>• Review the methodology to validate the model and predict the dependent variable for a set of given independent variables</li><li>• Use R graphics functions to visualize the results generated with the model</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>.</p>

## Workflow Overview





## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set the Working Directory</u></b> Set the working directory to ~/LAB07/ by executing the command:</p> <pre>setwd("~/LAB07")</pre> <ul style="list-style-type: none"> <li>(Or using the “Tools” option in the tool bar in the RStudio environment).</li> </ul>
3	<p><b><u>Define the problem and review input data</u></b></p> <p><b>Logistic Regression</b>, also known as <b>Logit</b>, is typically used in models where the dependent variables have a binary outcome (True/False, which is coded with 1/0). You model the log odds of the outcome as a linear combination of predictor variables).</p> <p><b><u>Marketing Survey Data</u></b> In this lab you use hypothetical marketing survey data in which customers:</p> <ul style="list-style-type: none"> <li>Responded to the question: <ul style="list-style-type: none"> <li>Would you choose a product based on a “pricing” factor (three “Price” ranges 10, 20 and 30)?</li> </ul> </li> <li>Response options: <ul style="list-style-type: none"> <li>“1” for “yes” and “0” for “no”</li> </ul> </li> <li>The survey also collected information such as “Age” and “Income” of the respondent.</li> </ul> <p><b><u>Business Need</u></b> The marketing campaign team wants to send special offers to those respondents with the highest probability of purchase.</p> <p>This data file “survey.csv” is available in the folder ~/LAB07/survey.csv.</p> <ol style="list-style-type: none"> <li>Review the survey.csv file.</li> <li>How many responses to the survey does the file contain?</li> <li>What is the main purpose of building this model?</li> </ol>

Step	Action
4	<p><b><u>Read in and Examine the Data:</u></b></p> <ol style="list-style-type: none"> <li>1. The first step in the modeling process is to examine the data and determine if there are any outliers. To do this you must read in the survey data, use the following command:   <pre>Mydata &lt;- read.csv("survey.csv",header=TRUE,sep=",")</pre> </li> <li>2. With the following command, explore the data further:   <pre>&gt; table(Mydata\$MYDEPV) &gt; with(Mydata, table(Price,MYDEPV)) &gt; summary(Mydata\$Age) &gt; cor.mat &lt;- cor(Mydata[, -1]) &gt; cor.mat</pre> </li> <li>3. Review the results on the console</li> </ol> <p><b>Note:</b> The general rule <b>is not</b> to include variables in your model that are too highly correlated with other predictors. For example, including two variables that are correlated by 0.85 in your model may prevent the true contribution of each variable from being identified by the statistical algorithm. Confirm that the variables in our survey do not fall in this category.</p>

Step	Action
5	<p><b><u>Build and Review Logistic Regression Model:</u></b></p> <ol style="list-style-type: none"> <li>1. Use the “glm” function for logit modeling. Type in the following command:  <pre>mylogit &lt;- glm(MYDEPV ~ Income + Age + as.factor(Price) ,                data=Mydata, family=binomial(link="logit") ,                na.action=na.pass)</pre> </li> <li>2. Review the model by typing the “summary” and “plot” functions:  <pre>summary(mylogit)</pre> <p>Review the results of the summary command, for the fitted model, on the console. Results you should see:</p> <ul style="list-style-type: none"> <li>• The first line provides the model you specified.</li> <li>• Next, you should see the <b>deviance residuals</b>, which provide the measure of the model fit.</li> <li>• The next part of the output shows the <b>coefficients</b>, their standard errors, the <b>z-statistic</b> (sometimes called a Wald z-statistic), and the associated <b>p-values</b>.</li> <li>• Both <b>Income</b> and <b>Age</b> are statistically significant, as are the two terms for <b>Price</b>.</li> <li>• The <b>logistic regression coefficients</b> show the change in the <b>log odds</b> of the outcome for a one unit increase in the predictor variable.</li> <li>• Residual deviance: analogous to the Residual Sum of Squares of a linear model; that is, it is related to the "total error" of the fit. It is twice the negative log likelihood of the model.</li> <li>• Null deviance: the deviance associated with the "null model" -- that is the model that returns just the global probability of TRUE for every x. The quantity 1 - (Residual deviance/Null deviance) is sometimes called "pseudo-R-squared"; you use it to evaluate goodness of fit in the same way that R-sqr is used for linear models.</li> </ul> <p>The interpretation of the results are as follows:</p> <ol style="list-style-type: none"> <li>1. Review the “Estimate” column. For every one unit change in <b>Income</b>, the log odds of Purchase (versus no-Purchase) increases by 0.12876.</li> <li>2. Record the number that describes how much one unit increase in <b>Age</b> increases the log odds of purchase: The indicator variables for <b>Price</b> are interpreted differently. <b>For</b> example, Purchase decision at a <b>Price</b> of 20, compared with a <b>Price</b> of <b>10</b>, decreases the log odds of admission by 0.74418</li> <li>3. Record the log odds at Price point 30 compared to Price point 10 below:</li> </ol> <p>The summary then shows the table of coefficients that are “fit indices”, including the null and deviance residuals and the AIC.</p> </li> </ol>

Step	Action
6	<p><b><u>Review the results and interpret the coefficients</u></b></p> <ol style="list-style-type: none"> <li>1. Use the “confint” function to obtain the confidence intervals of the coefficient estimates:   <code>confint(mylogit)</code> </li> <li>2. Review the results on the console. <ul style="list-style-type: none"> <li>• You can also exponentiate the coefficients and interpret them as odds-ratios.</li> <li>• To get the exponentiated coefficients, use (<b>exp( )</b>)</li> <li>• The object you want to exponentiate is called coefficients and it is part of mylogit (<b>mylogit\$coefficients</b>).</li> </ul> <code>exp(mylogit\$coefficients)</code> </li> </ol> <p>You can observe that for every unit change in income, the odd-ratio of Purchase increases by a multiplicative factor of 1.137 (and remember a multiplicative factor of 1 corresponds to no change).</p> <p>This is actually a bit more intuitive than the log odds explanation you reviewed in the previous step. Observe that that Age does not appear to be a very strong factor in this model, and the price factor of 30 has a stronger effect than a price factor of 20.</p>
7	<p><b><u>Visualize the Model Using the Plot Function:</u></b></p> <code>plot(mylogit)</code> <p>You should see multiple plots generated on the graphics window.</p>

Step	Action																																										
8	<p><b><u>Use relevel Function to re-level the Price factor with value 30 as the base reference.</u></b></p> <p>In the original model that we fitted with the function call:</p> <pre>mylogit &lt;- glm(MYDEPV ~ Income + Age + as.factor(Price) ,               data= Mydata,family=binomial(link="logit"), na.action=na.pass)</pre> <p>we obtained the results shown below:</p> <p><b>Coefficients:</b></p> <table><thead><tr><th></th><th>Estimate</th><th>Std. Error</th><th>z value</th><th>Pr(&gt; z )</th><th></th></tr></thead><tbody><tr><td>(Intercept)</td><td>-6.02116</td><td>0.53244</td><td>-11.309</td><td>&lt; 2e-16</td><td>***</td></tr><tr><td>Income</td><td>0.12876</td><td>0.00923</td><td>13.950</td><td>&lt; 2e-16</td><td>***</td></tr><tr><td>Age</td><td>0.03506</td><td>0.01179</td><td>2.974</td><td>0.00294</td><td>**</td></tr><tr><td>as.factor(Price) 20</td><td>-0.74418</td><td>0.26439</td><td>-2.815</td><td>0.00488</td><td>**</td></tr><tr><td>as.factor(Price) 30</td><td>-2.21028</td><td>0.31108</td><td>-7.105</td><td>1.2e-12</td><td>***</td></tr><tr><td>---</td><td></td><td></td><td></td><td></td><td></td></tr></tbody></table> <p>What does this tell us?</p> <p>The odds of MYDEPV decreases when price changes from 10 to 20 and decreases even more when we go from 10 to 30.</p> <p>1. Now let's use 30 as the reference price, instead of 10. Type in the following:</p> <pre>Mydata\$pricefactor = relevel(as.factor(Mydata\$Price) , "30")</pre>		Estimate	Std. Error	z value	Pr(> z )		(Intercept)	-6.02116	0.53244	-11.309	< 2e-16	***	Income	0.12876	0.00923	13.950	< 2e-16	***	Age	0.03506	0.01179	2.974	0.00294	**	as.factor(Price) 20	-0.74418	0.26439	-2.815	0.00488	**	as.factor(Price) 30	-2.21028	0.31108	-7.105	1.2e-12	***	---					
	Estimate	Std. Error	z value	Pr(> z )																																							
(Intercept)	-6.02116	0.53244	-11.309	< 2e-16	***																																						
Income	0.12876	0.00923	13.950	< 2e-16	***																																						
Age	0.03506	0.01179	2.974	0.00294	**																																						
as.factor(Price) 20	-0.74418	0.26439	-2.815	0.00488	**																																						
as.factor(Price) 30	-2.21028	0.31108	-7.105	1.2e-12	***																																						
---																																											

Step	Action
8 Cont.	<p>Fit the Model Again (mylogit2) and Display the Summary:</p> <pre>mylogit2 = glm(MYDEPV ~ Income + Age + pricefactor ,                data= Mydata,family=binomial(link="logit") ,                na.action=na.pass) summary(mylogit2)</pre> <p>You will see the results as follows:</p> <pre>Coefficients:               Estimate Std. Error z value Pr(&gt; z ) (Intercept)  -8.23144    0.66180  -12.438  &lt; 2e-16 *** Income         0.12876    0.00923   13.950  &lt; 2e-16 *** Age           0.03506    0.01179    2.974  0.00294 ** pricefactor10  2.21028    0.31108    7.105  1.20e-12 *** pricefactor20  1.46610    0.29943    4.896  9.76e-07 *** --- </pre> <p>Notice that the intercept has changed (because we changed the reference situation), but the coefficients for Income and Age are the same. The new model tells us that the odds of MYDEPV increase when price decreases from 30 to 10, and less so price decreases from 30 to 20.</p>

Step	Action
9	<p><b><u>Plot the ROC Curve:</u></b></p> <ol style="list-style-type: none"> <li>1. Make sure you have the package ROCR installed and the library included <pre>install.packages("ROCR") library(ROCR)</pre> </li> <li>2. First get all the probability scores on the training data <pre>pred = predict(mylogit, type="response")</pre> </li> <li>3. Every classifier evaluation using ROCR starts with creating a prediction object. This function is used to transform the input data (which can be in vector, matrix, data frame, or list form) into a standardized format. We create the prediction object needed for ROCR as follows: <pre>predObj = prediction(pred, Mydata\$MYDEPV)</pre> </li> <li>4. All kinds of predictor evaluations are performed using the function "performance". Read and understand the parameters of the function with <pre>?performance</pre> </li> <li>5. We now create the ROC curve object and the AUC object with performance function <pre>rocObj = performance(predObj, measure="tpr", x.measure="fpr") # creates ROC curve obj aucObj = performance(predObj, measure="auc") # auc object</pre> </li> <li>6. Extract the value of AUC and display on the console: <pre>auc = aucObj@y.values[[1]] auc</pre> </li> <li>7. What is the value of AUC?</li> <li>8. We will plot the ROC curve now <pre>plot(rocObj, main = paste("Area under the curve:", auc))</pre> </li> <li>9. Review the curve on the plot window. Review the discussions on ROC in the student resources guide. Record your observations below:</li> </ol>

Step	Action
10	<p><b><u>Predict Outcome given Age and Income:</u></b></p> <ol style="list-style-type: none"> <li>1. Use the “predict” function to predict the probability of the purchase outcome given <b>Age</b> and <b>Income</b>. Start with predicting the probability of the purchase decision at different <b>Price</b> points (10, 20, and 30). Create a “data frame” called “newdata1” using the following commands: <pre> Price &lt;- c(10,20,30) Age &lt;- c(mean(Mydata\$Age) ) Income &lt;- c(mean(Mydata\$Income) ) newdata1 &lt;- data.frame(Income, Age, Price) newdata1 </pre> <p>You are predicting with <b>Income</b> and <b>Age</b> both set at their mean value and <b>Price</b> at 10, 20 and 30.</p> <p><b>Note:</b> The values of the data frame “newdata1” displayed on the console. The predict function requires the variables to be named exactly as in the fitted model.</p> </li> <li>2. Create the fourth variable “PurchaseP”. <pre> newdata1\$PurchaseP &lt;- predict (mylogit, newdata=newdata1, type="response") newdata1 </pre> </li> <li>3. What is your observation on the probability of purchase at different <b>Price</b> levels?</li> </ol>



Step	Action
11	<p><b><u>Predict outcome for a sequence of Age values at price 30 and income at its mean:</u></b></p> <ol style="list-style-type: none"> <li>Keep the <b>Price</b> at 30, <b>Income</b> at its mean value and select a sequence of values for <b>Age</b> starting at a minimum age, incrementing by 2 until the maximum age in our dataset: <pre>newdata2 &lt;- data.frame(Age=seq(min(Mydata\$Age),max(Mydata\$Age),2),             Income=mean(Mydata\$Income),Price=30) newdata2\$AgeP&lt;- predict(mylogit,newdata=newdata2,type="response") cbind(newdata2\$Age,newdata2\$AgeP)</pre> <p>Newdata2\$AgeP stores the predicted variables and you just display the sequence for <b>Age</b> you generated and the corresponding probability of the purchase decision using the “cbind” function shown above.</p> </li> <li>Plot and visualize how the “purchase” probability varies with Age: <pre>plot(newdata2\$Age,newdata2\$AgeP)</pre> </li> </ol>
12	<p><b><u>Predict outcome for a sequence of income at price 30 and Age at its mean:</u></b></p> <ol style="list-style-type: none"> <li>Using the same methodology, create a data frame newdata3 with the following characteristics: <ul style="list-style-type: none"> <li><b>Income</b> is a sequence from 20 to 90 in steps of 10</li> <li><b>Age</b> is the mean value for the dataset Mydata</li> <li><b>Price</b> point at 30</li> </ul> </li> <li>Predict <b>newdata3\$IncomeP</b> and display the Income sequence along with the predicted probabilities.</li> <li>Plot the results.</li> </ol>

Step	Action
13	<p><b><u>Use Logistic regression as a classifier:</u></b></p> <p>Recall the problem statement in Step 3, the marketing campaign team wants to send special offers to those respondents with the highest probability of purchase. They have established a threshold of 0.5 and they want to target customers whose probability of purchase are greater than 0.5.</p> <p><b>Note:</b> We are assuming that age and income are uniformly distributed in our customer base, and the price factors of our products are also uniformly distributed. Typically in order to run a scenario like this you should understand the demographic distribution of the customers (and the price distribution of the products).</p> <ol style="list-style-type: none"> <li>1. You want an idea of how many offers will be sent out, using this threshold, so you test it on a 'random' set of data. First, generate this random set using “runif” functions:</li> </ol> <pre>newdata4 &lt;- data.frame (   Age= round(runif(10,min(Mydata\$Age),max(Mydata\$Age))), Income=round(runif(10,min(Mydata\$Income),max(Mydata\$Income))), Price = round((runif(10,10,30)/10))*10) newdata4\$Prob &lt;- predict(mylogit,newdata=newdata4,type="response") newdata4</pre> <ol style="list-style-type: none"> <li>2. How many samples in your random selection qualify for special offers?</li> </ol> <p>The R Code for this exercise is available at <a href="#">home/gpadmin/LAB07/logit.R</a></p>

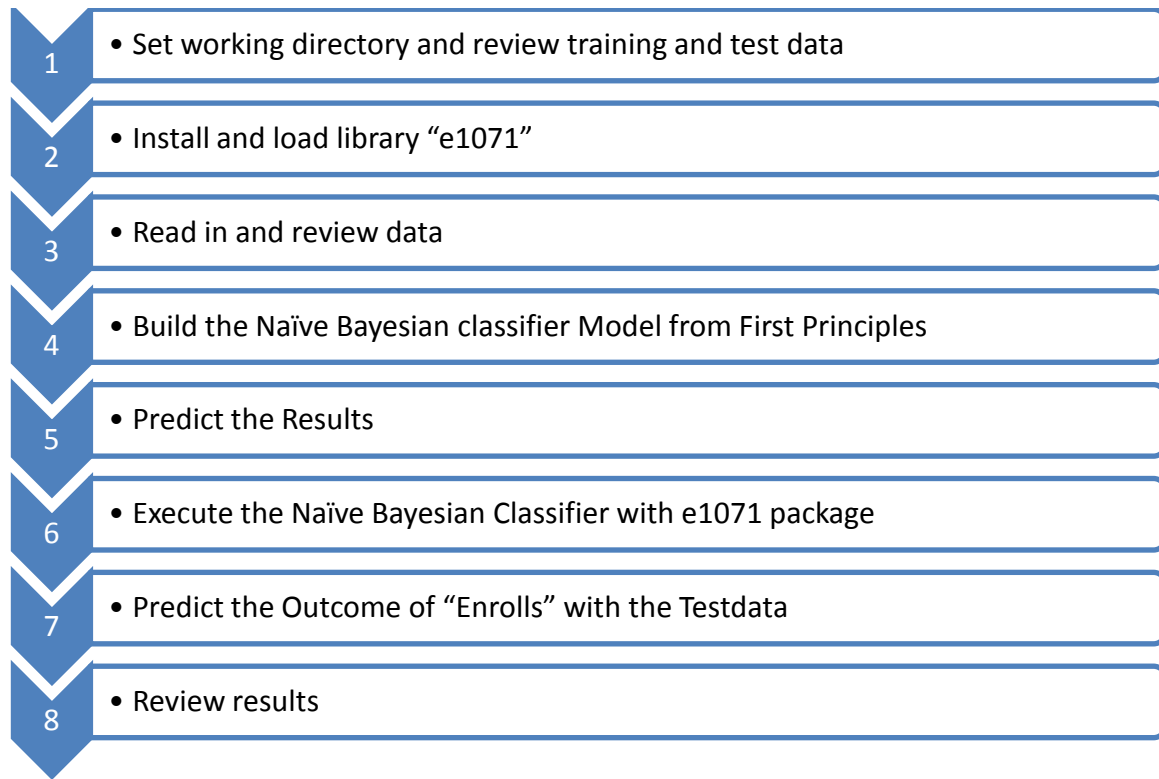
*End of Lab Exercise*

## Lab Exercise 8: Naïve Bayesian Classifier

<b>Purpose:</b>	<p>This lab is designed to investigate and practice the Naïve Bayesian Classifier analytic technique. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Naïve Bayesian Classification</li><li>• Apply the requirements for generating appropriate training data</li><li>• Validate the effectiveness of the Naïve Bayesian Classifier with the big data</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none"><li>• Use R –Studio environment to code the Naïve Bayesian Classifier</li><li>• Use the ODBC connection to the “census” database to create a training data set for Naïve Bayesian Classifier from the big data</li><li>• Use the Naïve Bayesian Classifier program and evaluate how well it predicts the results using the training data and then compare the results with original data</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>.</p>

## Part 1 – Building Naïve Bayesian Classifier

### Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set working directory and review training and test data</u></b></p> <p>1. Set the working directory using the following command:</p> <pre>&gt; setwd("~/LAB08")</pre> <ul style="list-style-type: none"> <li>The “<b>sample1.csv</b>” file in this directory represents the data worked with in the instructor led training session. The file has a header row, followed by 14 rows of training data.</li> <li>The <b>testing data</b> on which you will predict the results should be appended after the <b>training data</b>. The data set should read:</li> </ul> <pre>Age,Income,Jobstaisfaction,Desire,Enrolls      ←-----Header &lt;=30,High,No,Fair,No &lt;=30,High,No,Excellent,No 31 to 40,High,No,Fair,Yes &gt;40,Medium,No,Fair,Yes &gt;40,Low,Yes,Fair,Yes &gt;40,Low,Yes,Excellent,No 31 to 40,Low,Yes,Excellent,Yes &lt;=30,Medium,No,Fair,No &lt;=30,Low,Yes,Fair,Yes &gt;40,Medium,Yes,Fair,Yes &lt;=30,Medium,Yes,Excellent,Yes 31 to 40,Medium,No,Excellent,Yes 31 to 40,High,Yes,Fair,Yes &gt;40,Medium,No,Excellent,No &lt;=30,Medium,Yes,Fair,                          ←-----testing data</pre>
3	<p><b><u>Install and load library “e1071”</u></b></p> <p>Execute the following command to install the required packages and load the libraries:</p> <pre>&gt; install.packages("e1071") &gt; library("e1071")</pre>

Step	Action
4	<p><b><u>Read in and review data</u></b></p> <ol style="list-style-type: none"> <li>1. Execute the following to read in the data. <pre data-bbox="321 289 1425 487">&gt; # read the data into a table from the file &gt; sample &lt;- read.table("sample1.csv",header=TRUE,sep=",") &gt; # we will now define the data frames to use the NB classifier &gt; # we will now define the data frames to use the NB classifier &gt; traindata &lt;- as.data.frame(sample[1:14,]) &gt; testdata &lt;- as.data.frame(sample[15,])</pre> <p data-bbox="321 525 1445 556">You now have two data frame objects “<b>traindata</b>” and “<b>testdata</b>” for running the NB Classifier.</p> </li> <li>2. Execute the following command to display the data frames, to ensure they are loaded properly. <pre data-bbox="321 703 706 798">&gt; #Display data frames &gt; traindata &gt; testdata</pre> </li> <li>3. Review the output on the console window.</li> </ol>

Step	Action
5	<p><b><u>Build the Naïve Bayesian classifier Model from First Principles:</u></b></p> <ol style="list-style-type: none"> <li>1. The first step in building the model is the computation of prior probabilities. The independent variables here are the “Age”, “Income”, “Jobsatisfaction” and “Desire”. The dependent variable is “Enrolls” Compute the prior probabilities of enrollment, P(no), P(yes) first, the counts :   <pre>&gt; tprior &lt;- table(traindata\$Enrolls)</pre> then, normalize over the total number of instances to get the probabilities  <pre>&gt; tprior &lt;- tprior/sum(tprior)</pre> <pre>&gt; tprior</pre>  Review the results of prior probabilities on the console </li> <li>2. Compute the summaries that you need to create a Bayes model: <math>P(A b)</math>, <math>b=\{no, yes\}</math> First, count up "no" and "yes" by Age:  <pre>&gt; ageCounts &lt;- table(traindata[,c("Enrolls", "Age")])</pre> </li> <li>3. Then, normalize by the total number of "no" and "yes" each to get the conditional probabilities  <pre>&gt; ageCounts &lt;- ageCounts/rowSums(ageCounts)</pre>  Display the results on the console and review the conditional probabilities  <pre>&gt; ageCounts</pre> </li> <li>4. Do the same for the other variables.   <pre>&gt; incomeCounts &lt;- table(traindata[,c("Enrolls", "Income")])</pre> <pre>&gt; incomeCounts &lt;- incomeCounts/rowSums(incomeCounts)</pre> <pre>&gt; jsCounts &lt;- table(traindata[,c("Enrolls", "Jobsatisfaction")])</pre> <pre>&gt; jsCounts&lt;-jsCounts/rowSums(jsCounts)</pre> <pre>&gt; desireCounts &lt;- table(traindata[,c("Enrolls", "Desire")])</pre> <pre>&gt; desireCounts &lt;- desireCounts/rowSums(desireCounts)</pre> </li> </ol>

Step	Action
6	<p><b><u>Predict the Results:</u></b></p> <ol style="list-style-type: none"> <li>1. Use the Naïve Bayesian Classifier formula to compute product of <math>P(A b)</math>, for <math>b=\{no, yes\}</math>. The maximum of the two is the “predicted” result of the dependent variable. In the test data we need to predict the “Enrolls” given the for Age<math>\leq</math>30, Income = Medium, Jobsatisfaction = yes and Desire = Fair <pre> &gt; pyes &lt;-   ageCounts["Yes", "&lt;=30"] *   incomeCounts["Yes", "Medium"] *   jsCounts["Yes", "Yes"] *   desireCounts["Yes", "Fair"] *   tprior["Yes"] </pre> <p>followed by</p> <pre> &gt; pno &lt;-   ageCounts["No", "&lt;=30"] *   incomeCounts["No", "Medium"] *   jsCounts["No", "Yes"] *   desireCounts["No", "Fair"] *   tprior["No"] </pre> </li> <li>2. The prediction will be <math>\max(\text{pyes}, \text{pno})</math>. <pre> &gt; pyes &gt; pno &gt; max(pyes, pno) </pre> </li> <li>3. What is the predicted result for “Enrolls”?</li> </ol>



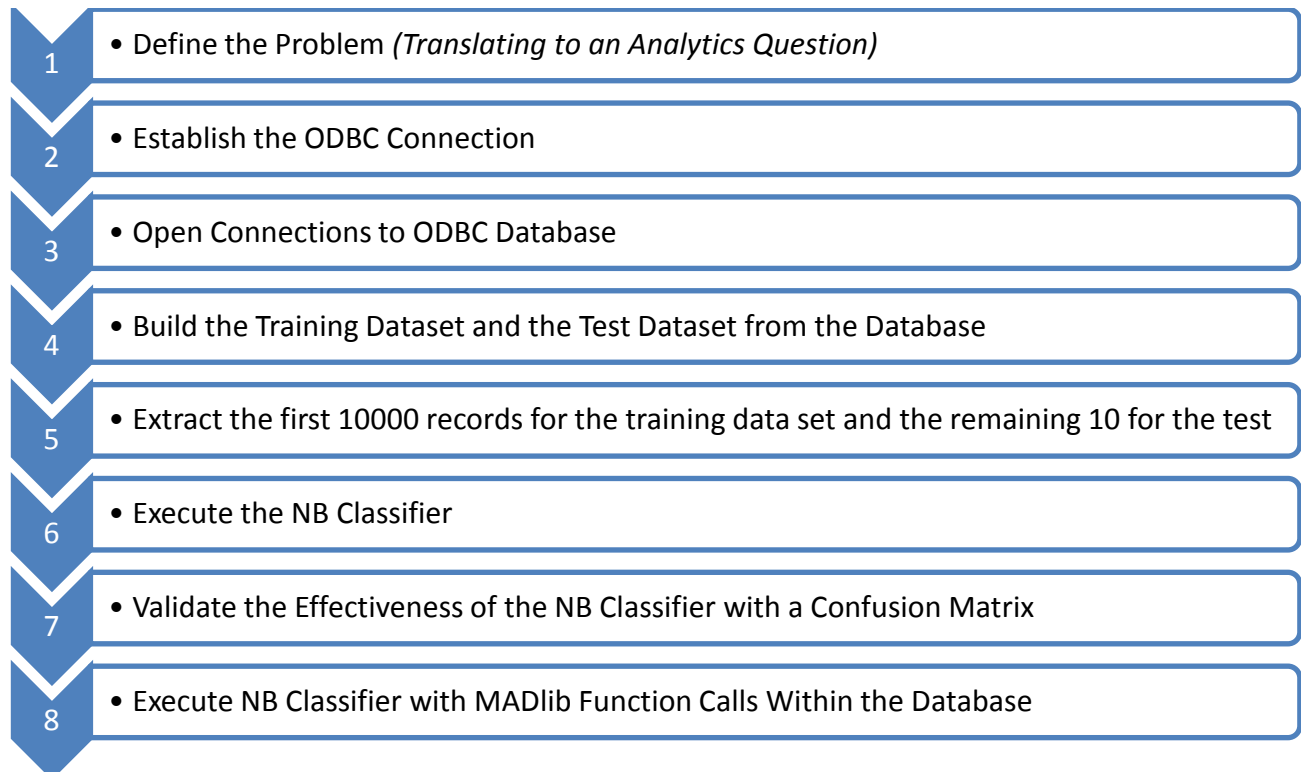
Step	Action
7	<p><b><u>Execute the Naïve Bayesian Classifier with e1071 package:</u></b></p> <p>The Naïve Bayes function computes the conditional a-posterior probabilities of a categorical class variable given independent categorical predictor variables using the Bayes rule. The usage takes the form of <code>naiveBayes(formula, data,...)</code> where the arguments are defined as follows:</p> <ul style="list-style-type: none"> <li>• <b>formula</b> A formula of the form <code>class ~ x1 + x2 + ....</code> Interactions are not allowed.</li> <li>• <b>data</b> Either a data frame of factors or a contingency table.</li> </ul> <ul style="list-style-type: none"> <li>• You are modeling for attribute “Enrolls”.</li> </ul> <ol style="list-style-type: none"> <li>1. Use the following commands to execute the model and display the results.  <pre>&gt; # use the NB classifier &gt; model &lt;- naiveBayes(Enrolls ~., traindata) &gt; # display model &gt; model</pre> </li> <li>2. Review the results on the console and compare these results to the <b>apriori probabilities</b> you manually computed earlier in step 5.</li> </ol>
8	<p><b><u>Predict the Outcome of “Enrolls” with the Testdata:</u></b></p> <ol style="list-style-type: none"> <li>1. To use the predict function, type in the following:  <pre>&gt; # predict with testdata &gt; results &lt;- predict(model, testdata) &gt; # display results &gt; results</pre> </li> <li>2. Review the results (Prediction for “Enrolls”) on the console.</li> </ol>

Step	Action
9	<p><b><u>Review results</u></b></p> <ol style="list-style-type: none"> <li>1. Look at <math>P(\text{age}=31-40 \mid \text{Enrolls} = \text{No})</math>. You will observe a zero probability. Is this a problem?</li> <li>2. Build another NB model, with Laplace smoothing <code>model2 = naiveBayes(Enrolls ~.,traindata, laplace=0.01)</code></li> <li>3. Compare the probabilities here with those of the first model</li> </ol> <p>Note down your observations in the space provided below:</p>

*End of Lab Exercise*

## Part 2 – Naïve Bayesian Classifier – Census data

### Workflow Overview



## LAB Instructions

Step	Action
1	<p><b><u>Define the Problem (Translating to an Analytics Question):</u></b></p> <ol style="list-style-type: none"> <li>Use “Persons” table in the Census dataset to categorize Age, Gender, Educational qualifications and Annual Income from each record as follows: <ul style="list-style-type: none"> <li>3 Age categories: <math>&gt;20</math> and <math>\leq 30</math> , <math>&gt;30</math> and <math>\leq 45</math> and <math>&gt;45</math></li> <li>2 Gender categories: M and F</li> <li>3 Educational Qualifications categories: <math>&gt;14</math> (Professional/Phd), <math>&gt;12</math> and <math>\leq 14</math> (College) and <math>&lt;12</math> (others) –</li> <li>3 Annual Income categories: <math>&gt;10000</math> and <math>\leq 50000</math> , <math>&gt;50000</math> and <math>\leq 80000</math> and <math>&gt; 80000</math></li> </ul> </li> <li>Build an appropriate “training” dataset, which will be a subset of the “categorized” table with four columns age, gender, education and income.</li> <li>Predict the annual income category a person will belong to given the Age, Gender and educational qualifications, using the Naïve Bayesian Classifier.</li> </ol>
2	<p><b><u>Establish the ODBC Connection:</u></b></p> <ol style="list-style-type: none"> <li>Load the RODBC package, using the following command:  <pre>library('RODBC')</pre> </li> </ol>
3	<p><b><u>Open Connections to ODBC Database:</u></b></p> <ol style="list-style-type: none"> <li>Before connecting to the ODBC database make sure the file, /etc/odbc.ini is properly set to point to database “training2”.</li> <li>If not, edit the line that starts with, "Database =" within the file /etc/odbc.ini to point to “training2”</li> <li>Ensure the username(uid) and password (pwd) are provided correctly in the following command :</li> </ol> <pre>ch &lt;- odbcConnect("Greenplum",uid="gpadmin",   case="postgresql",pwd="password_of_the_gpadmin_user")</pre>

Step	Action
4	<p><b><u>Build the Training Dataset and the Test Dataset from the Database:</u></b></p> <ol style="list-style-type: none"> <li>Drop the table, NBtrain, from the database, use the following command  <pre>sqlDrop(ch, "NBtrain")</pre> </li> <li>Execute a SQL query using the sqlQuery command, creating the table, NBtrain, selecting 10010 random records and categorizing the variables in the categories we defined in Step1 of Part2 of this lab. Use the code below.   <p><b>Note:</b> The code is available in the working directory "NBTrain.sql". To execute, you can copy and paste it into your R script window.</p> <pre>&gt; sqlQuery(ch, " CREATE TABLE NBtrain (   age VARCHAR(8) , sex VARCHAR(8) , educ VARCHAR(8) , income VARCHAR(8) ) DISTRIBUTED BY (age) ; INSERT INTO NBtrain SELECT   t1.age , t1.sex , t1.educ , t1.income FROM (   SELECT     CASE       WHEN age BETWEEN 20 AND 30 THEN '20-30'       WHEN age BETWEEN 31 AND 45 THEN '31-45'       WHEN age &gt; 45 THEN 'GT 45'       ELSE 'unknown age'     END AS age , CASE       WHEN sex = 1 THEN 'M'       WHEN sex = 2 THEN 'F'       ELSE 'unknown sex'     END AS sex , CASE       WHEN educ &gt;14 THEN 'Prof/Phd'       WHEN educ BETWEEN 12 AND 14 THEN 'College'       WHEN educ &lt;12 THEN 'Others'       ELSE 'unknown educ'     END AS educ </pre> </li> </ol>

Step	Action
<p>4 Cont.</p>	<pre> , CASE   WHEN inctot BETWEEN 10000 AND 50000 THEN '10-50K'   WHEN inctot BETWEEN 50000+1 AND 80000 THEN '50-80K'   WHEN inctot &gt; 80000 THEN 'GT 80K'   ELSE 'unknown i' END AS income FROM   persons ) AS t1 WHERE   not (t1.age like 'unk%' or t1.sex like 'unk%' or t1.educ like 'unk%' or t1.income like 'unk%') ORDER BY RANDOM () LIMIT 10010 ; ") </pre>
<p>5</p>	<p><b><u>Extract the first 10000 records for the training data set and the remaining 10 for the test</u></b></p> <ol style="list-style-type: none"> <li>1. Use the sqlFetch command for reading data into an R data frame.  <pre>NBtrain &lt;- (sqlFetch(ch,"NBtrain"))</pre> </li> <li>2. Extract the training dataset  <pre>&gt; NBtrain1 &lt;- NBtrain[1:10000,]</pre> </li> <li>3. Extract the test dataset  <pre>&gt; NBtest &lt;- NBtrain[10001:10010,]</pre> </li> <li>4. Close the ODBC channel  <pre>&gt; odbcClose(ch)</pre> </li> </ol>

Step	Action
6	<p><b><u>Execute the NB Classifier:</u></b></p> <ol style="list-style-type: none"> <li>Run the model as you did in Part 1. Use the following command: <pre># model model &lt;- naiveBayes(income ~.,NBtrain1) model</pre> </li> <li>Review the results of the model on the console.</li> <li>Run the predict function using the following command: <pre># predict with testdata results &lt;- predict (model,NBtest[1:10,-1]) results</pre> </li> <li>Record the results of the predict function: <ul style="list-style-type: none"> <li>1</li> <li>2</li> <li>3</li> <li>4</li> <li>5</li> <li>6</li> <li>7</li> <li>8</li> <li>9</li> <li>10</li> </ul> </li> <li>Use the parameter “type” with a value “raw” and you can see how the scores are very close to 0 or 1 rather than looking like realistic probabilities <pre>&gt; results1 &lt;- predict (model,NBtest[1:10,-1],type="raw") &gt; results1</pre> <p>Note down your observations below:</p> </li> </ol>

Step	Action
7	<p><b><u>Validate the Effectiveness of the NB Classifier with a Confusion Matrix:</u></b></p> <ol style="list-style-type: none"> <li>1. Compare the results of the previous step with the actual data that you have in NBtest. Build a confusion matrix for the predictions vs. actual values:   <pre>&gt; conf &lt;- table(actual=NBtest[1:10,4],predicted=results) &gt; conf</pre> </li> <li>2. The diagonals of “conf” give the count of correctly classified instances by class. The off-diagonals tell how many instances of each class are mis-classified. What class does the model predict best?</li> <li>3. What % of data did the NB Classifier predicted correctly? You can calculate this as the sum of the diagonal elements of the confusion matrix normalized by the total number of test cases   <pre>&gt; accuracy &lt;- sum(diag(conf)) / sum(conf) &gt; accuracy</pre> </li> <li>4. Record any other observations below:</li> </ol> <ul style="list-style-type: none"> <li>• The R code and the data are available at <a href="#">/home/gpadmin/LAB08/NBcoderev.R</a></li> </ul>



Step	Action
8	<p><b><u>Execute NB Classifier with MADlib Function Calls Within the Database:</u></b></p> <p>In this step we will execute a SQL query that uses MADlib function calls for the NB classifier model and assignment of class labels.</p> <ol style="list-style-type: none"> <li>1. Log into the VM assigned to you with “gpadmin” credentials</li> <li>2. Navigate to /home/gpadmin/LAB08/</li> <li>3. Review the file NBmadlib.sql in this directory</li> <li>4. The first part of the code categorizes the data and prepares a table that can be used with MADlib functions for NB classifier:</li> </ol> <pre> DROP TABLE IF EXISTS myschema.NBmdlib; CREATE TABLE myschema.NBmdlib (   attr INTEGER[]   , class INTEGER ) DISTRIBUTED BY (class) ;  INSERT INTO myschema.NBmdlib (attr[1],attr[2],attr[3],class) SELECT   t1.age   , t1.sex   , t1.educ   , t1.income FROM (   SELECT     CASE       WHEN age BETWEEN 20 AND 30 THEN 1       WHEN age BETWEEN 31 AND 45 THEN 2       WHEN age &gt; 45 THEN 3       ELSE 0     END AS age   , CASE       WHEN sex = 1 THEN 1       WHEN sex = 2 THEN 2       ELSE 0     END AS sex   , CASE       WHEN educ &gt; 14 THEN 1       WHEN educ BETWEEN 12 AND 14 THEN 2       WHEN educ &lt; 12 THEN 3       ELSE 0     END AS educ   , CASE       WHEN inctot BETWEEN 10000 AND 50000 THEN 1       WHEN inctot BETWEEN 50000+1 AND 80000 THEN 2       WHEN inctot &gt; 80000 THEN 3       ELSE 0     END AS income </pre>

Step	Action
8 Cont.	<pre> FROM     persons ) AS t1 WHERE     not (income = 0 OR age = 0 OR sex = 0) ; </pre> <p>5. The second part of the code builds the classifier model. Note here we have all the eligible rows in the table taken in for the training data set. Please review NB classifier documentation for the MADlib function at <a href="http://doc.madlib.net/v0.2beta/group_grp_bayes.html">http://doc.madlib.net/v0.2beta/group_grp_bayes.html</a></p> <pre> DROP TABLE IF EXISTS myschema.nb_feature_probs; DROP TABLE IF EXISTS myschema.nb_class_priors;  SELECT madlib.create_nb_prepared_data_tables(     'myschema.NBmdlib', 'class' , 'attr', 3, 'myschema.nb_feature_probs' , 'myschema.nb_class_priors');  SELECT * FROM myschema.nb_feature_probs; SELECT * FROM myschema.nb_class_priors; </pre> <p>6. Step 5 will take some time to execute. We now prepare some data to score the model. We essentially use the same code in step 4 but select 10 random records.</p> <pre> DROP IF EXISTS TABLE myschema.NBmdlib_test; CREATE TABLE myschema.NBmdlib_test (     id SERIAL , attr INTEGER[] , original_data INTEGER ) DISTRIBUTED BY (id) ; INSERT INTO myschema.NBmdlib_test(     attr[1] , attr[2] , attr[3] , original_data ) SELECT     t1.age , t1.sex , t1.educ , t1.income FROM (     SELECT         CASE             WHEN age BETWEEN 20 AND 30 THEN 1             WHEN age BETWEEN 31 AND 45 THEN 2             WHEN age &gt; 45 THEN 3             ELSE 0         END AS age </pre>

Step	Action
<p>8 Cont.</p>	<pre> , CASE     WHEN sex = 1 THEN 1     WHEN sex = 2 THEN 2     ELSE 0 END AS sex , CASE     WHEN educ &gt;14 THEN 1     WHEN educ BETWEEN 12 AND 14 THEN 2     WHEN educ &lt;12 THEN 3     ELSE 0 END AS educ  , CASE     WHEN inctot BETWEEN 10000 AND 50000 THEN 1     WHEN inctot BETWEEN 50000+1 AND 80000 THEN 2     WHEN inctot &gt; 80000 THEN 3     ELSE 0 END AS income FROM     persons ) AS t1 WHERE     NOT (income = 0 OR age = 0 OR sex = 0) ORDER BY RANDOM () LIMIT 10 ; </pre> <p>7. We run the SQL code that predicts the classes for the test data we created in the previous step:</p> <pre> SELECT * from myschema.NBmdlib_test ORDER BY id;  DROP TABLE IF EXISTS myschema.nb_classify_view_fast; DROP TABLE IF EXISTS myschema.nb_probs_view_fast;  SELECT madlib.create_nb_classify_view (     'myschema.nb_feature_probs' , 'myschema.nb_class_priors' , 'myschema.NBmdlib_test' , 'id', 'attr', 3 , 'myschema.nb_classify_view_fast' ); </pre>

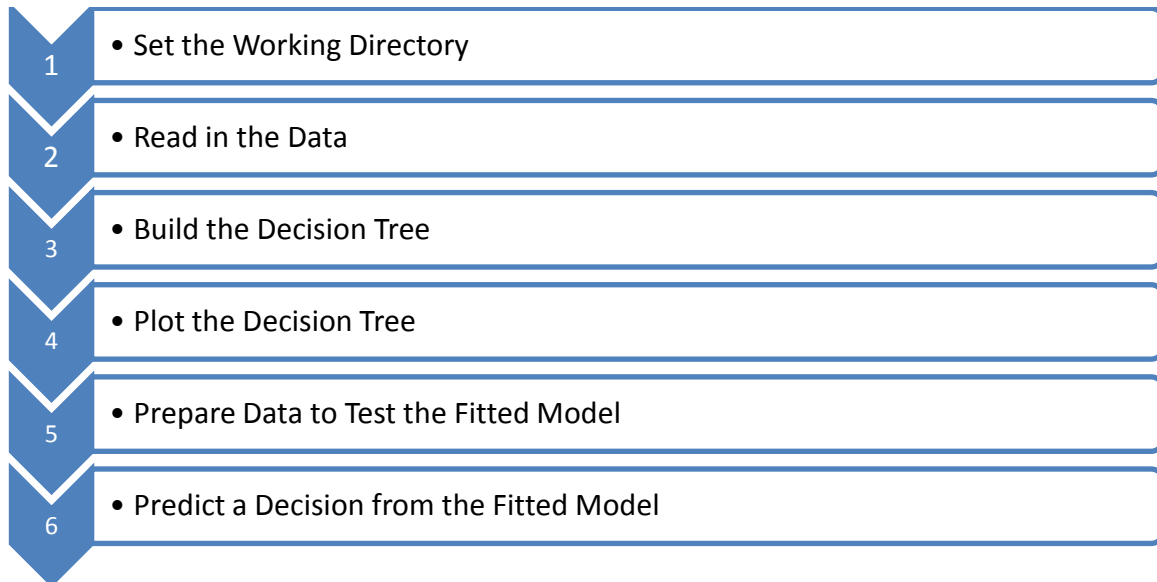
Step	Action
8 Cont.	<pre>SELECT * FROM myschema.nb_classify_view_fast ORDER BY key;</pre> <pre>SELECT madlib.create_nb_probs_view (   'myschema.nb_feature_probs' , 'myschema.nb_class_priors' , 'myschema.NBmdlib_test' , 'id', 'attr', 3 , 'myschema.nb_probs_view_fast');</pre> <pre>SELECT * FROM myschema.nb_probs_view_fast ORDER BY key,class;</pre> <p>All these code are available at /home/gpadmin/LAB08/NBmadlib.sql</p> <p>You can execute them with the following command</p> <pre>psql -d training2 -f NBmadlib.sql</pre> <p>Record your observations below:</p>

*End of Lab Exercise*

## Lab Exercise 9: Decision Trees

<b>Purpose:</b>	<p>This lab is designed to investigate and practice Decision Tree (DT) models covered in the course work. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for Decision Tree models</li><li>• Predict the outcome of an attribute based on the model</li></ul>
<b>Tasks:</b>	<p>Tasks you will complete in this lab exercise include:</p> <ul style="list-style-type: none"><li>• Use the R –Studio environment to code Decision Tree Models</li><li>• Build a Decision Tree Model based on data whose schema is composed of attributes</li><li>• Predict the outcome of one attribute based on the model</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>.</p>

## Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set the Working Directory:</u></b></p> <p>1. Execute the command:</p> <pre>setwd("~/LAB09")</pre> <p>(Or use the “Tools” option in the tool bar in the RStudio environment.)</p> <p>2. Load the package rpart.plot and the associated libraries. If prompted for the location to download select any integer representing a location nearest to you.</p> <pre>&gt; install.packages("rpart.plot") &gt; library("rpart") &gt; library("rpart.plot")</pre>
3	<p><b><u>Read in the Data:</u></b></p> <ul style="list-style-type: none"> <li>• Use a data table with columns for data attributes : Play, Outlook, Temperature, Humidity and Windy</li> <li>• A Decision Tree allows for predicting the values of the attribute Play, given that you know the values for attributes like Outlook, Humidity and Windy.</li> </ul> <p>1. Read in the data from the “Dtdata.csv” file in the working directory and display the contents:</p> <pre>&gt; #Read the data &gt; play_decision &lt;- read.table("DTdata.csv",header=TRUE,sep=",") &gt; play_decision</pre> <p>2. How many observations did you read in?</p> <p>3. How many variables (attributes) did you read in?</p> <p>4. Use the command “summary” for a detailed list of the table object you read in</p> <pre>summary(play_decision)</pre> <p>5. Review the results. (The Summary is located in the console window.)</p>

Step	Action
4	<p><b><u>Build the Decision Tree:</u></b></p> <p>Use the “rpart” package in R for classification by Decision Trees. The RPart Programs build classification or regression models of a very general structure using a two stage procedure; the resulting models can be represented as binary trees.</p> <p>1. Use the following rpart commands to grow a Decision Tree:</p> <pre>rpart (formula, data=, method=, control=)</pre> <div data-bbox="337 537 1492 865"> <ul style="list-style-type: none"> <li>• <b>formula</b> is in the format: outcome ~ predictor1+predictor2+predictor3+ect.</li> <li>• <b>data=</b> specifies the dataframe</li> <li>• <b>method=</b> "class" for a classification tree "anova" for a regression tree</li> <li>• <b>control=</b> optional parameters for controlling tree growth. For example, control=rpart.control(minsplit=30, cp=0.001) requires that the minimum number of observations in a node be 30 before attempting a split and that a split must decrease the overall lack of fit by a factor of 0.001 (cost complexity factor) before being attempted.</li> </ul> </div> <p>The "Play" attribute is the outcome that will be predicted.</p> <p>2. Use the command:</p> <pre>&gt; fit &lt;- rpart(Play ~ Outlook + Temperature + Humidity + Wind, method="class", data=play_decision, + control=rpart.control(minsplit=1))</pre> <p>3. You can now display “fit” and review the results:</p> <pre>&gt; summary(fit)</pre> <p>Note that the leaf nodes information includes both the class label and the class probabilities (P(no), P(yes))</p>
5	<p><b><u>Plot the Decision Tree:</u></b></p> <p>1. Review the arguments for rpart.plot function. Type in:</p> <pre>&gt; ?rpart.plot</pre> <p>We will use the arguments “type” and “extra” in our plot.</p> <p>2. Type in the following :</p> <pre>&gt; rpart.plot(fit, type=4, extra=1)</pre> <p>3. Review the Decision Tree plot on the graphics window.</p>



Step	Action																		
6	<p><b><u>Prepare Data to Test the Fitted Model:</u></b></p> <p>You must use “fit” for a new set of data to create predictions from the DT:</p> <table><tr><td>Play Decision</td><td>Outlook</td><td>Temperature</td><td>Humidity</td><td>Wind</td></tr><tr><td>?</td><td>rainy</td><td>mild</td><td>high</td><td>FALSE</td></tr></table> <p>1. “newdata” is a data frame object and can be built for our test data. Type in the following statement:</p> <pre>newdata &lt;- data.frame(Outlook="rainy",Temperature="mild",Humidity="high",Wind=F ALSE)</pre> <p>2. Review the “newdata” displaying the dataframe</p> <pre>&gt; newdata</pre> <p>3. The data displayed as follows:</p> <table><tr><td>Outlook</td><td>Temperature</td><td>Humidity</td><td>Wind</td></tr><tr><td>1 rainy</td><td>mild</td><td>high</td><td>FALSE</td></tr></table>	Play Decision	Outlook	Temperature	Humidity	Wind	?	rainy	mild	high	FALSE	Outlook	Temperature	Humidity	Wind	1 rainy	mild	high	FALSE
Play Decision	Outlook	Temperature	Humidity	Wind															
?	rainy	mild	high	FALSE															
Outlook	Temperature	Humidity	Wind																
1 rainy	mild	high	FALSE																

7	<p><b><u>Predict a Decision from the Fitted Model:</u></b></p> <p>The “predict” function is used to generate predictions from a fitted rpart object.</p> <ul style="list-style-type: none"> <li>• “type” is a character string denoting the type of the predicted value</li> <li>• Use both “prob” and “class” to predict from a Decision Tree model</li> </ul> <pre>predict(object, newdata = list(),         type = c("vector", "prob", "class", "matrix"))</pre> <p>1. The <b>type=“prob”</b> gives the class probabilities for the prediction for newdata Type in  <b>&gt; predict(fit,newdata=newdata,type="prob")</b></p> <p>2. Repeat the prediction with type=“class”</p> <pre><b>&gt; predict(fit,newdata=newdata,type="class")</b></pre> <p>Review the results.</p> <p>3. What is the prediction for the “newdata”?</p>
8	<p>The code is available at /home/gpadmin/LAB09/DT.R</p>

*End of Lab Exercise*

## Lab Exercise 10: Time Series Analysis with ARIMA

<b>Purpose:</b>	<p>This lab is designed to investigate and practice Time Series Analysis with ARIMA models (Box-Jenkins-methodology). After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use R functions for ARIMA models</li><li>• Apply the requirements for generating appropriate training data</li><li>• Validate the effectiveness of the ARIMA 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 ARIMA models</li><li>• Use the ODBC connection to the database to create the weekly sales data from the retail database</li><li>• Prepare the data (sorting and rendering the data as a Time series)</li><li>• Generate a model and evaluate how well it predicts the results and compare the results with original data</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide Appendix</i></b>.</p>

## Workflow Overview



## LAB Instructions

Step	Action
1	Log in with GPADMIN credentials on to R-Studio.
2	<p><b><u>Set the working directory</u></b>  Set working directory to ~/LAB10/ , execute the command:  <code>setwd("~/LAB10")</code></p> <ul style="list-style-type: none"> <li>(Or using the “Tools” option in the tool bar in the RStudio environment.)</li> </ul>
3	<p><b><u>Establish the ODBC Connection:</u></b></p> <p>Load the RODBC package using the following command:  <code>library('RODBC')</code></p>
4	<p><b><u>Open Connections to ODBC Database:</u></b></p> <ol style="list-style-type: none"> <li>Before connecting to the ODBC database make sure the file, /etc/odbc.ini, is properly set to point to database “training1”.</li> <li>If not, edit the line that starts with, "Database =" within the file /etc/odbc.ini, to point to, “training1”.</li> <li>Ensure the username(uid) and password (pwd) are provided correctly in the following command:</li> </ol> <pre>ch &lt;- odbcConnect("Greenplum",uid="gpadmin",   case="postgresql",pwd="password_of_the_gpadmin_user")</pre>

Step	Action
5	<p><b><u>Get Data from the Database:</u></b></p> <ol style="list-style-type: none"> <li>Drop the table, weekly_sales, from the schema ddemo:</li> </ol> <pre>sqlDrop(ch, "ddemo.weekly_sales")</pre> <ol style="list-style-type: none"> <li>Execute an SQL query using the sqlQuery command, creating a table, weekly_sales, in which the weekly sales are grouped first by year and the week within a year:</li> </ol> <pre>sqlQuery(ch, "CREATE TABLE     ddemo.weekly_sales (         sale INTEGER,         Y1 INTEGER,         W1 INTEGER )     DISTRIBUTED BY (sale); INSERT INTO ddemo.weekly_sales SELECT     SUM((o.item_price*o.item_quantity)) as sale ,     EXTRACT(YEAR FROM o.order_datetime) as y1,     CASE         WHEN (EXTRACT(WEEK FROM o.order_datetime) = 53) THEN 52         ELSE             EXTRACT(WEEK FROM o.order_datetime)         END w1 FROM     ddemo.order_lineitems o GROUP BY     y1,w1 ; "</pre> <ul style="list-style-type: none"> <li><b>Note:</b> Note the use of the <b>EXTRACT</b> function to obtain the year and the week within the year from the order_datetime field</li> <li>The sales number is accumulated for each week</li> <li>The ISO Standard for numbering weeks within a year may lead to a year containing 52 or 53 weeks.</li> <li>In order to work with Time Series data you need a consistent “periodicity” with the data. You must accumulate the vales for week 53 with that of week 52 in the same year. We use the CASE statement to designate the week 53 as 52 and cumulate the sales amount for week 53 into week 52 of the same year.</li> </ul> <ol style="list-style-type: none"> <li>Get the results from the table into data frame msales. Execute the command:</li> </ol> <pre>msales &lt;- (sqlFetch(ch, "ddemo.weekly_sales"))</pre> <ol style="list-style-type: none"> <li>Close the ODBC channel.</li> </ol> <pre>odbcClose(ch)</pre>

Step	Action
6	<p><b><u>Review, Update and Prepare DataFrame "sales" for ARIMA Modeling:</u></b></p> <ol style="list-style-type: none"> <li>Sort the data in the order of Year and Week: Use the R function "order" :</li> </ol> <pre>attach(msales) msales &lt;- msales[order(y1,w1),] detach(msales)</pre> <ol style="list-style-type: none"> <li>Extract 300 values from column 1 of "msales" for modeling and 12 values to compare with the predictions done by the model. Store them in two different vectors: "sales" and "csales". Use the following command:</li> </ol> <pre>&gt;sales &lt;- c(rep(0,300)) &gt;csales &lt;- c(rep(0,12)) &gt;sales[1:300] &lt;- msales[1:300,1] &gt;csales[1:12] &lt;- msales[301:312,1]</pre>
7	<p><b><u>Convert "sales" into Time Series Type Data:</u></b></p> <p>Convert the "sales" into a time series.</p> <ul style="list-style-type: none"> <li>This "transformation" is required for most of the time-series functions, since a time series contains more information than the values themselves, namely information about dates and frequencies at which the time series has been recorded.</li> </ul> <pre>&gt; sales &lt;- ts(sales,start=2005,frequency=52)</pre>
8	<p><b><u>Plot the Time Series:</u></b></p> <ol style="list-style-type: none"> <li>Plot the Time Series using the following command:</li> </ol> <pre>plot(sales,type="l")</pre> <ol style="list-style-type: none"> <li>Review the plot of the Time Series.</li> <li>Identify the seasonality features in the graph.</li> <li>Is the data Seasonal (Do you see patterns that repeat at a particular frequency)?</li> <li>Is the data stationary?</li> <li>Is there a trend to the data?</li> </ol>

Step	Action
9	<p><b><u>Analyze the ACF and PACF :</u></b></p> <p>The next step in analyzing time series is to examine the <b>autocorrelations (ACF)</b> and <b>partial autocorrelations (PACF)</b>. R provides the functions <code>acf()</code> and <code>pacf()</code> for computing and plotting of ACF and PACF.</p> <ol style="list-style-type: none"> <li>1. Use the <b>parameter function “par”</b> to set the plot window to display both the ACF and PACF plots. Plot ACF and PACF on the same graph using the command: <pre>&gt; par(mfrow=c(2,1)) &gt; acf(sales) &gt; pacf(sales)</pre> </li> <li>2. Using the plot generated in the plot window: <ul style="list-style-type: none"> <li>• Does the ACF tail off quickly?</li> <li>• What does the ACF indicate with respect to stationarity of the data?</li> <li>• What will you do to make the data stationary?</li> </ul> </li> </ol>
10	<p><b><u>Difference the Data to Make it Stationary:</u></b></p> <p>To difference the data and make it stationary you should use the “diff” function in R. The “diff” function takes the pair of each observation and differences it from the one previous to it.</p> <ol style="list-style-type: none"> <li>1. Run the following code: <pre>&gt; #Difference the series and plot it &gt; sales1 &lt;- diff(sales) &gt; m &lt;- length(sales1) &gt; par(mfrow=c(1,1)) &gt; plot(1:m,sales1,type="l")</pre> </li> <li>1. Do you see any trend to the data now?</li> <li>2. Is seasonality still shown in the data?</li> <li>3. How does the trend of oscillating spikes change as you move from left to right on this plot?</li> </ol>



Step	Action
11	<p><b><u>Plot ACF and PACF for the Differenced Data:</u></b></p> <ol style="list-style-type: none"> <li>1. Run the following code: <pre data-bbox="321 325 987 457">&gt; #Plot ACF and PACF on the same graph &gt; par(mfrow=c(2,1)) &gt; acf(sales1) &gt; pacf(sales1)</pre> </li> <li>2. Do you see the ACF tailing off quickly?</li> </ol>

Step	Action
12	<p><b><u>Fit the ARIMA Model:</u></b></p> <p>Once you configure the data and review the seasonality elements, you are ready to fit an ARIMA model. Selecting the correct model to fit an ARIMA model is a bit of an art. Algorithms are used to select the correct parameters.</p> <ul style="list-style-type: none"> <li>• Use the data configuration and the principle of “parsimony” to fit a basic model</li> <li>• Use the statement ARIMA</li> <li>• Use the time series</li> <li>• Use “sales” and not the “diff(sales)” that we computed in step 10. ARIMA will automatically invoke the “diff” function based on the parameters specified.</li> <li>• You need to specify the order (p,d,q) where “p” is the order of AR , “q” order of MA and “d” is the number of differences.</li> <li>• Fit (1,1,0) so the model will difference once, use AR and MA parameters as 1 and 0.</li> <li>• Use a seasonal statement since the data seems to have a seasonal component to it as you see spikes every 52 weeks.</li> <li>• In the seasonal statement, you have a “list” where you put in the “order” and a (time) “period” which will be “52”.</li> <li>• Use “include.mean = false” - R will force a mean and continue any trend seen in the past for the model. This automatically defaults to “True”. You turn the automatic default off with this statement.</li> <li>• <b>Note:</b> You can also experiment with not using this statement during this lab.</li> </ul> <p>1. Type in the following code:</p> <pre>sales.fit &lt;- arima (sales,                     order=c(1,1,0) ,                     seasonal = list(order=c(1,1,0) ,period=52) ,                     include.mean=FALSE) &gt; sales.fit</pre> <p>2. Review the output displayed.</p> <p>3. Record the coefficients for the AR term and the seasonal AR terms and the standard errors.</p> <p>4. What is your observation on the standard errors compared to the coefficients?</p> <ul style="list-style-type: none"> <li>• <b>Note:</b> The model gives you the “log likelihoods” that provide important input on model selection.</li> </ul>

Step	Action
13	<p><b><u>Generate Predictions:</u></b></p> <ol style="list-style-type: none"> <li>1. Use the fitted model “sales.fit” and the “predict” statement for 12 periods ahead:  <pre>sales.predict &lt;- predict (sales.fit, n.ahead=12)</pre> <ul style="list-style-type: none"> <li>• <b>Note:</b> You can see the predictions by typing “sales.fit”.</li> </ul> </li> <li>2. Plot the predictions using plot statements:  <pre>&gt; par(mfrow=c(1,1)) &gt; plot (sales,xlim=c(2009.50,2011)) &gt; lines (sales.predict\$pred,col="blue") &gt; lines (sales.predict\$pred+2*sales.predict\$se,col="red") &gt; lines (sales.predict\$pred-2*sales.predict\$se,col="red")</pre> <ul style="list-style-type: none"> <li>• <b>Note:</b> The first plot statement plots the original “sales” data</li> <li>• EMC has restricted the “xlim” value to zoom in on the values just prior to the predicted values</li> <li>• The “blue” line indicates the mean of the predictions</li> <li>• The “red” lines denote the upper and lower bounds of the prediction</li> </ul> </li> </ol>
14	<p><b><u>Compare predicted values with actual values</u></b></p> <ol style="list-style-type: none"> <li>1. Compare the predicted values with the actual values and plot the actual Values Predicted as a barplot. Use the following code:  <pre>&gt; #comparing the predictions with the actual values &gt; par(mfrow=c(1,1))  &gt; forbar &lt;- matrix(x,ncol=12,byrow=TRUE) &gt; colnames(forbar)&lt;- asales[1:12,3] &gt; rownames(forbar)&lt;- c("Actual", "Predicted") &gt; barplot(forbar,beside=TRUE, +         main="Actual Vs Predicted", +         ylab="Weekly Sales", +         col=rainbow(2), +         xlab="Weeks 2010")</pre> </li> <li>2. Review the plot.</li> </ol> <ul style="list-style-type: none"> <li>• <b>Note:</b> The complete code is available in the folder /home/gpadmin/LAB10/arma.R</li> </ul>

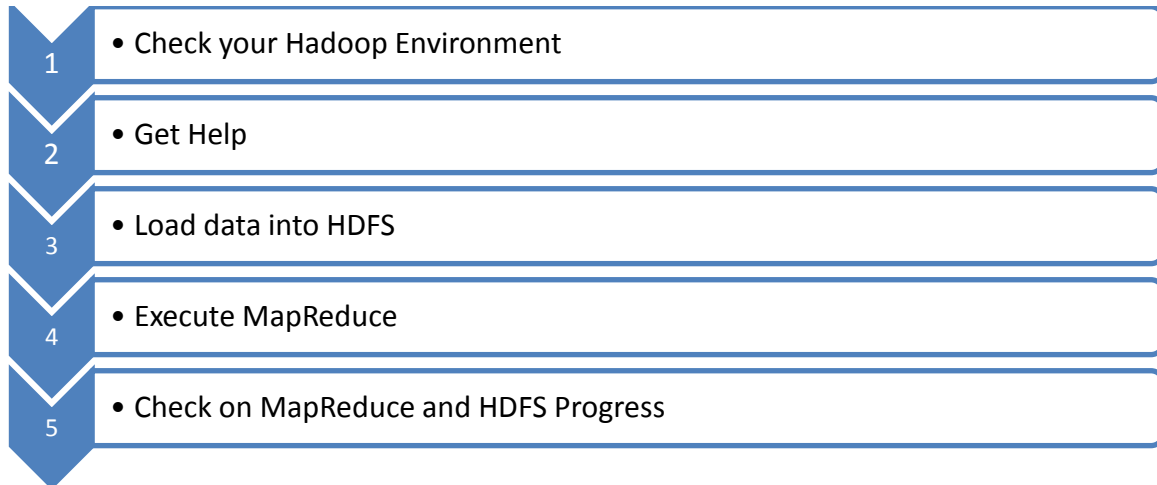
*End of Lab Exercise*



## Lab Exercise 11: Hadoop, HDFS and MapReduce

<b>Purpose:</b>	<p>This lab introduces the <i>Hadoop and MapReduce environment</i> that you will be working on for the next lab. After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Get help on the various Hadoop commands</li><li>• Observe a MapReduce job in action</li><li>• Query various Hadoop servers regarding status</li></ul>
<b>Tasks:</b>	<p>Tasks you'll be completing in this lab include:</p> <ul style="list-style-type: none"><li>• Run Hadoop and Hadoop fs and collect help information</li><li>• Run a shell script to perform a word count activity</li><li>• Run a MapReduce job to produce similar output</li><li>• Investigate the UI for MapReduce/HDFS components to track system behavior</li></ul>
<b>References:</b>	<p>References used in this lab are located in your <b><i>Student Resource Guide</i></b>. See the Guide for:</p> <ul style="list-style-type: none"><li>• Hadoop Commands</li><li>• HDFS Commands</li></ul>

## Workflow Overview



## LAB Instructions

Step	Action
1	<p><b><u>Check your Hadoop Environment:</u></b></p> <ol style="list-style-type: none"> <li>1. Start a terminal session to your “be” host – you may use PuTTY on your “fe” host for this. Log in with username “gpadmin” and the gpadmin password supplied by your instructor. Make sure you are connected to the right directory. Execute the command:   <code>cd ~/LAB11</code></li> <li>2. Execute the command:   <code>printenv   grep '^H'</code></li> <li>3. Are there any Hadoop variables defined? Which ones?</li> <li>4. Execute the command:   <code>hadoop fs -ls</code></li> <li>5. What do you see?</li> <li>6. Execute the command:   <code>ls</code></li> </ol> <p>This should be different</p>
3	<p><b><u>Get Help:</u></b></p> <ol style="list-style-type: none"> <li>1. Execute the command:   <code>hadoop -help 2&gt;&amp;1   tee Hadoop.hlp</code></li> <li>2. Execute the command:   <code>hadoop fs -help   tee HDFS.hlp</code></li> <li>3. Now list the files by executing the commands:   <code>clear &amp;&amp; more Hadoop.hlp</code>  <code>clear &amp;&amp; more HDFS.hlp</code></li> </ol> <p>These files are also contained in your Student Resource Guide. These commands can be run from your command shell whenever you need them.</p> <ol style="list-style-type: none"> <li>4. Now we create an alias so we don’t have to type “hadoop fs” every time   <code>alias hdfs="hadoop fs"</code></li> </ol>

Step	Action
4	<p><b><u>Load data into HDFS:</u></b></p> <p>In this step, we will be loading an input data file into HDFS that we will be using in later activities.</p> <ol style="list-style-type: none"> <li>1. Execute the following commands.</li> </ol> <pre>hdfs -copyFromLocal speech.txt input/speech.txt</pre> <pre>hdfs -ls input</pre> <p>The file <i>speech.txt</i> is now in HDFS.</p>
5	<p><b><u>Execute MapReduce:</u></b></p> <p>In this step, you will run a MapReduce job and observe its output. This job is identical to the example job discussed as part of the lecture.</p> <ol style="list-style-type: none"> <li>1. First, execute the command</li> </ol> <pre>time ./wf.sh speech.txt</pre> <ol style="list-style-type: none"> <li>2. How long did it take to produce its output? _____</li> <li>3. List the file “MRwordcount.sh” by executing</li> </ol> <pre>more MRwordcount.sh</pre> <p>This should be identical to the file used as an example in the lecture.</p> <ol style="list-style-type: none"> <li>4. Execute the following command:</li> </ol> <pre>time ./MRWordCount.sh</pre> <ol style="list-style-type: none"> <li>5. What do you see? _____</li> <li>6. How long did it take to execute this script? _____</li> <li>7. Is this command slower than the Unix script file? Why do you think that is? _____</li> <li>8. Looking at the output</li> </ol> <p>The output of the MapReduce job is stored in a subdirectory of the output directory in HDFS. This directory is named “d” followed by a string of numbers (and is listed in the output of the MapReduce command). You can see the content of the directory by executing the following command:</p> <pre>hdfs -cat output/d*/part-*</pre>



Step	Action
6	<p><b><u>Check on MapReduce and HDFS Progress:</u></b></p> <p>In this step, we will look at some of the status information about MapReduce and Hadoop. First we look at the administrative interfaces for the JobTracker and the TaskTracker components of the MapReduce framework, and then NameNode User Interface (UI) for HDFS.</p> <p>0.1. Create a shortcut on your desktop for each UI that you will be investigating. For each of the following steps, right-click on your Windows desktop and select New&gt;&gt;Shortcut.</p> <p>Enter the following URL: <a href="http://&lt;IP-ADDRESS-OF-SERVER&gt;:&lt;PORT_NUMBER&gt;/">http://&lt;IP-ADDRESS-OF-SERVER&gt;:&lt;PORT_NUMBER&gt;/</a> For the JobTracker, PORT_NUMBER will be 50030. IP-ADDRESS-OF-SERVER is the IP address for your server that you received at the beginning of the course. Click “Next.”</p> <p>Name the shortcut “JobTracker” and click “Finish”</p> <p>0.2. Do the same for the NameNode (port number 50070) and for the TaskTracker (port # 50060).</p> <ol style="list-style-type: none"> <li>1. Click on the desktop icon labeled “Jobtracker”. This will bring up the UI for the JobTracker MapReduce node.</li> <li>2. Examine the output. Do you see anything similar to what you saw in the output of the script you just ran? Anything different?</li> <li>3. Click on the desktop icon labeled “TaskTracker” This will bring up the UI for the MapReduce TaskTracker node.</li> <li>4. Examine the output. Do you see anything similar to what you saw in the output of the script you just ran? Anything different?</li> </ol> <p>Now we look at the UI for the NameNode node in HDFS.</p> <p>5. Click on the desktop icon labeled “NameNode” This invokes the UI for the NameNode node in an HDFS implementation. Examine the output.</p> <p>If you are strictly working as an analyst, you may never need to look at the administrative interface to these components of the Hadoop framework. On the other hand, if things aren’t working out as you might have expected, you can use these interfaces to take a deeper look “under the hood” at the mechanics of Hadoop.</p>

*End of Lab Exercise*

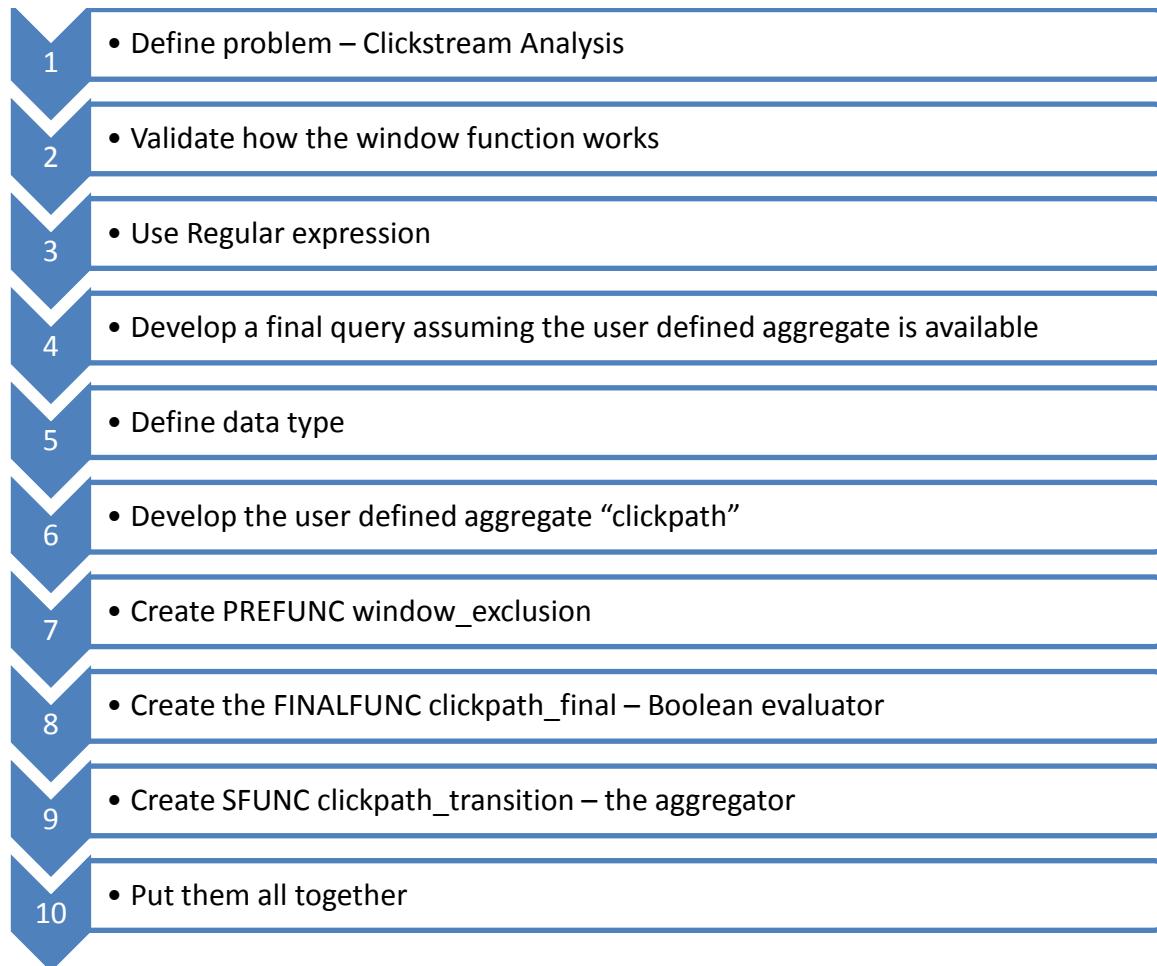


## Lab 12: In-database Analytics

<b>Purpose:</b>	<p>This lab is designed to familiarize you and give you practice with the in-database analytics methods covered in lessons three and four of Module 5</p> <p>After completing the tasks in this lab you should be able to:</p> <ul style="list-style-type: none"><li>• Use window functions</li><li>• Implement user defined aggregates and user defined functions</li><li>• Use ordered aggregates</li><li>• Use Regular Expressions (Regex) in SQL for text filtering</li><li>• Use MADlib functions and plot results from MADlib function outputs</li></ul>
<b>Tasks:</b>	<p>Tasks you'll be completing in this lab include:</p> <ul style="list-style-type: none"><li>• Process Clickstream analysis data using window functions, User defined functions, User defined aggregates and regular expressions</li><li>• Compute median household income using ordered aggregates</li><li>• Use MADlib functions for logistic regression and direct output to plot the results</li></ul>
<b>References:</b>	<p>Student resource guide</p> <p><a href="http://doc.madlib.net/v0.2beta/group_grp_logreg.html">http://doc.madlib.net/v0.2beta/group_grp_logreg.html</a></p>

## Part 1 – In-database analysis of Click-Stream data

### Workflow Overview



## LAB Instructions

Step	Action
1	<p><b><u>Define problem - Clickstream Analysis</u></b></p> <p>Problem Definition: A users' click-stream is defined as the aggregate of all activity a user has through a website via their clicks, derived through the web logs. This has become an important view of the data, as it enables insights into typical paths a user takes to navigate a website of interest. Analysis of click-streams can help to improve the usability of the websites, identify hacking attempts on the website, etc. In this lab, we will be constructing and analyzing click-streams from pre-processed weblog data. You are provided with data in a database table called "clicks". The table is defined as follows:</p> <p>TABLE clicks(user_id BIGINT, timestamp BIGINT, page_type VARCHAR)</p> <p>Where</p> <p>userid : user session number,</p> <p>page_type : identification of the page visited</p> <p>timestamp : the time of the visit.</p> <p>We want to determine which users:</p> <ul style="list-style-type: none"><li>– Start at the home page,</li><li>– <b>then</b> Click on an auction,</li><li>– <b>then</b> View <b>at least one</b> help page</li><li>– <b>then</b> Place a bid</li></ul> <p><b>In this lab you will connect to "module5indb" database and all the data tables required for this lab are available in this database.</b></p>

Step	Action
2	<p><b><u>Validate how the window function works</u></b></p> <p>Key in the following code and test how the windows function works:</p> <pre> SELECT     sid   , page_type   , time   , count(*) OVER (prefix) AS seq_length   , count(*) OVER (PARTITION BY sid) AS max_seq_length FROM     clicks WINDOW prefix AS (PARTITION BY sid ORDER BY time ASC) LIMIT 50 ; </pre> <p>The SELECT statement selects from table “clicks”, Session_id, Page_type, and the time stamp.</p> <p>Two standard “count” aggregate functions (which return the count of all records), are also included in the SELECT statement. The first one is defined as “sequence length” and the second one is defined as maximum sequence length.</p> <p>The first “count” aggregate is cumulated over a window defined as “prefix”; “prefix” is partitioned by variable “session id” and ordered by “time” (in a ascending order).</p> <p>For example if “session_id” = “1” had 10 different clicks at different times, your output for seq_length will be the sequence number of the clicks in the ascending order of time in session_id = “1”.</p> <p>If there are 10 clicks that session you will have 10 rows in the output.</p> <p>The second aggregate is cumulated over the partition defined by session id, the second “count” aggregate in our example will be 10 as there are 10 clicks for “session_id” =1.</p> <p>Execute the code and observe the results. We have limited the output to 50 rows.</p>

Step	Action
3	<p><b><u>Use Regular expression</u></b></p> <p>Check through the window defined as “prefix” and determine if the user went through a particular sequence of “page_types”. We want to know if the user (defined by the session_id):</p> <ul style="list-style-type: none"> <li>a) Starts at the home page</li> <li>b) Then clicks on an auction</li> <li>c) Then views at least one help page</li> <li>d) Then places a bid</li> </ul> <p>Define an aggregate that will step through the window “prefix” and look at the page types at every record in the window.</p> <p>If we call our pages with notation S,A,H,B we are looking for a sequence in regular expression terms “^SAH+B”. (defined with a variable “pattern”)</p> <p>Extract the first character of page_type (use upper case) and build a sequence of the page_type characters and compare this with our regular expression string “^SAB+H”.</p> <p>The code to perform the above mentioned tasks:</p> <pre> SELECT   sid , page_type , time , count(*) OVER (prefix) AS seq_length , count(*) OVER (PARTITION BY sid) AS max_seq_length , upper(substring(page_type for 1)) AS mystring , '^SAH+B' AS pattern FROM   clicks WINDOW prefix AS (PARTITION BY sid ORDER BY time ASC) LIMIT 50 ; </pre> <p>Review the output.</p>

Step	Action
4	<p><b><u>Develop a final query assuming the user defined aggregate is available</u></b></p> <p>The output of the column is the first character of “page id”. As you step through each time stamp of the “preview” window, aggregate the first characters at each pass. This aggregated character set is compared with the “pattern” “^SAH+B”.</p> <p><b>Write a user defined aggregate</b> that will accumulate the text string on each step it traverses in the window and return a Boolean value “true” or “false” based on the match with the pattern.</p> <p>Call this function “clickpath” and the arguments for this function are</p> <ul style="list-style-type: none"> <li>• the upper cased first character of the page_type and</li> <li>• the regular expression “pattern”</li> </ul> <pre>clickpath(upper(substring(page_type for 1)), '^SAH+B' )</pre> <p>This function should work as an aggregate over the window “prefix” accumulating the first character and determining the Boolean value of match.</p> <p>Our final query code (assuming clickpath works the way it is intended) will be:</p> <pre> SELECT   sid FROM (   SELECT     sid     , page_type     , time     , clickpath(upper(substring(page_type for 1)), '^SAH+B'   ) OVER (prefix) AS match     , count(*) OVER (prefix) AS seq_length     , count(*) OVER (PARTITION BY sid) AS max_seq_length   FROM     clicks   WINDOW prefix AS (PARTITION BY sid ORDER BY time ASC)   ) AS subq WHERE   seq_length = max_seq_length   AND match = true ; </pre>



Step	Action
5	<p><b><u>Define data type</u></b></p> <p>Define a composite data type that you will use with the aggregation function. Our composite data type will consists of</p> <ul style="list-style-type: none"> <li>• the sequence we are aggregating and</li> <li>• a regular expression “pattern” (which does not change) that we will use for comparison.</li> </ul> <p>Create data type with the following code:</p> <pre> DROP TYPE IF EXISTS clickstream_state CASCADE; CREATE TYPE clickstream_state AS (     sequence VARCHAR     , pattern VARCHAR ); </pre>
6	<p><b><u>Develop the user defined aggregate “clickpath”</u></b></p> <p>There are two major functions of “clickpath”</p> <ul style="list-style-type: none"> <li>• It should aggregate the characters (transition function that aggregates)</li> <li>• It should compare and return a Boolean function (the final function that returns the Boolean value)</li> </ul> <p>Key in the following code:</p> <pre> DROP AGGREGATE IF EXISTS clickpath(     /* Symbol */ CHAR     , /* regex */ TEXT ); CREATE AGGREGATE clickpath(     /* Symbol */ CHAR     , /* regex */ TEXT) (     STYPE = clickstream_state,     SFUNC = clickpath_transition,     FINALFUNC = clickpath_final,     PREFUNC = window_exclusion ); </pre> <p><b>Note:</b> The STYPE is the data type we defined in step 5.</p> <p>We have to now create two functions:</p> <ul style="list-style-type: none"> <li>• clickpath_transition (the aggregator) and</li> <li>• clickpath_final (the Boolean evaluator)</li> </ul> <p>Notice that we also defined a PREFUNC, a function required to enable the function clickpath to be called as a window function.</p>

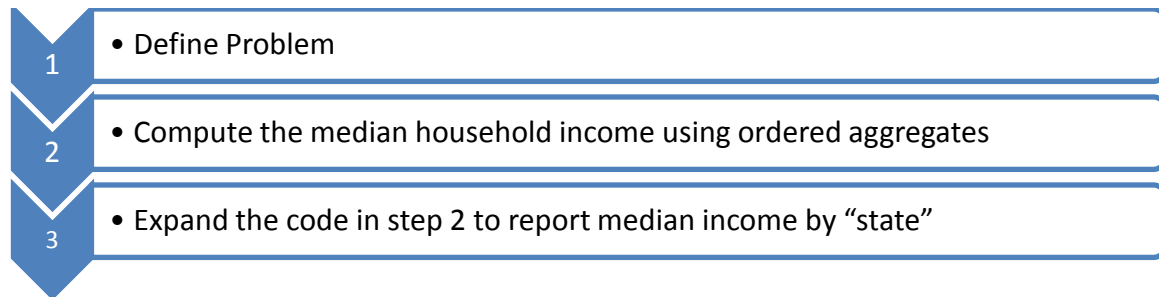
Step	Action
7	<p><b><u>Create PREFUNC window exclusion:</u></b></p> <pre> CREATE OR REPLACE FUNCTION window_exclusion(clickstream_state, clickstream_state) RETURNS clickstream_state AS \$\$ BEGIN     RAISE EXCEPTION 'aggregate may only be called from a window function'; END; \$\$ LANGUAGE PLPGSQL STRICT; </pre>
8	<p><b><u>Create the FINALFUNC clickpath final - Boolean evaluator:</u></b></p> <p>The Boolean evaluator is the simpler of the two remaining functions.</p> <pre> CREATE OR REPLACE FUNCTION clickpath_final(state clickstream_state) RETURNS BOOLEAN AS \$\$     SELECT \$1.sequence ~ \$1.pattern; \$\$ LANGUAGE SQL STRICT; </pre> <p>The sequence and the pattern are matched and the Boolean value is returned. \$1 refers to the first and the only argument in the function call. Recall the composite data type we created has both sequence and pattern.</p>

Step	Action
9	<p><b><u>Create SFUNC clickpath transition – the aggregator</u></b></p> <p>The next and the last function to define is the aggregator. This function has three arguments.</p> <ul style="list-style-type: none"> <li>• The “state” which aggregates with every step,</li> <li>• The “symbol”, the character we read in from the current row</li> <li>• The pattern to match</li> </ul> <p>When you step into a new window, the “state” will be NULL and it will take in the first character. As we step through each row within the window the aggregation will be carried out.</p> <p>Code the function as follows:</p> <pre> CREATE OR REPLACE FUNCTION clickpath_transition(     state clickstream_state, symbol CHAR(1), pattern VARCHAR) RETURNS clickstream_state AS \$\$     SELECT CASE         WHEN \$1 IS NULL THEN (\$2, \$3)::clickstream_state         ELSE (\$1.sequence    \$2, \$3)::clickstream_state     END; \$\$ LANGUAGE SQL CALLED ON NULL INPUT; </pre>
10	<p><b><u>Put them all together:</u></b></p> <ol style="list-style-type: none"> <li>1. Start with the definition of data type (step 5)</li> <li>2. Code the three functions SFUNC, FINALFUNC and PREFUNC (Steps 8,9,7)</li> <li>3. Complete the user defined aggregate (step6)</li> <li>4. Run the query (step4)</li> </ol> <p>The segments of this code are available in <b>/home/gpadmin/LAB12/clickstream_step*.sql</b> (* represents the steps in the document).</p> <p>If you have not completed the code as you reviewed the Lab, compile them in to one file and execute the query with the following command:</p> <pre> psql -d module5indb -f your_code.sql </pre>

*End of Lab Exercise*

## Part 2 – In-database computation of Median with Ordered Aggregates

### Workflow Overview



## LAB Instructions

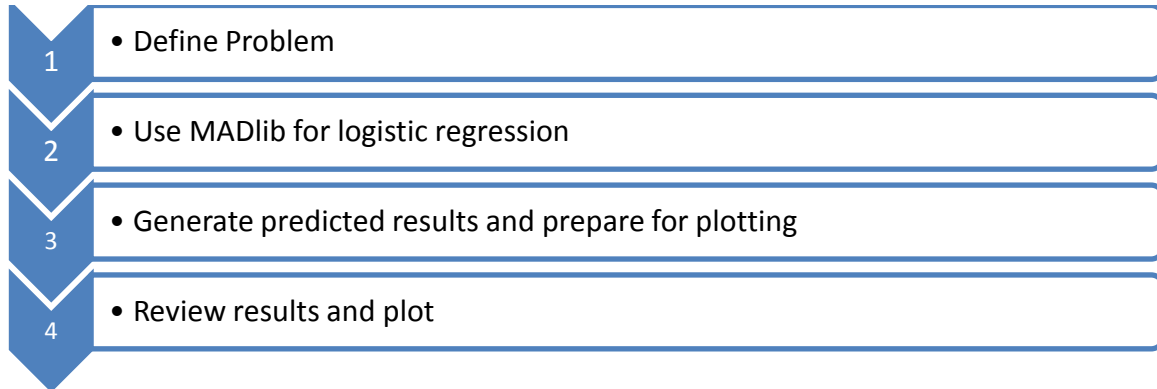
Step	Action
1	<p><b><u>Define Problem:</u></b></p> <p>Use the housing table in training2 database (census) to compute the median household income for each state.</p>
2	<p><b><u>Compute the median household income using ordered aggregates</u></b></p> <p>Use ordered aggregates for the computation of median household income. Code suggestion:</p> <pre> SELECT   ( arr[ length/2 + 1 ] + arr[ (length + 1)/2 ] ) / 2.0 AS median_income FROM(   SELECT     array_agg(hinc ORDER BY hinc) AS arr   , count(*) AS length   FROM     housing   ) AS q </pre> <p>What is the overall median household income in the US?</p>

Step	Action
3	<p><b><u>Expand the code in step 2 to report median income by “state”:</u></b></p> <p>Execute the following code:</p> <pre> SELECT     f.name   , ( arr[ length/2 + 1 ] + arr[ (length + 1)/2 ] ) / 2.0 AS     median_income FROM(     SELECT         state AS s       , array_agg(hinc ORDER BY hinc) AS  arr       , count(*) AS length     FROM         housing     GROUP BY         state     ) AS q JOIN     fips f ON     s = f.code ORDER BY     f.name ; </pre> <p>What is the median income of Massachusetts and Alaska?</p> <p>The code in step 4 is available at <a href="/home/gpadmin/LAB12/median.sql">/home/gpadmin/LAB12/median.sql</a></p>

*End of Lab Exercise*

## Part 3: Logistic Regression with MADlib

### Workflow Overview



## LAB Instructions

Step	Action
1	<p><b><u>Define Problem:</u></b></p> <p>In this exercise you will use the MADlib function for logistic regression and generate the model and plot the predicted results. Synthetic data is available in the table “artificiallogreg” in database “module5indb”</p>
2	<p><b><u>Use MADlib for logistic regression</u></b></p> <p>Execute the following code to generate the model and store the results in a table “logr_coef”</p> <pre> DROP TABLE IF EXISTS logr_coef; CREATE TABLE logr_coef AS     SELECT 0::INT AS bla     , NULL::FLOAT8[] AS coef DISTRIBUTED BY (bla) ;  UPDATE logr_coef SET coef = (SELECT coef FROM madlib.logregr('artificiallogreg', 'y', 'x', 20, 'irls', 0.001) AS coef) ; </pre>
3	<p><b><u>Generate predicted results and prepare for plotting</u></b></p> <p>Generate the predicted results; organize them in ascending order of value of x. Pipe the results using meta commands “\o” to a file called “graphics.txt” that we can use to plot in the next step:</p> <pre> \a \o graphics.txt SELECT     DISTINCT rank::FLOAT8/total_count AS x     , count::FLOAT8/total_true AS y FROM (     SELECT         y         , rank() OVER (ORDER BY prediction DESC)         , count(*) OVER () total_count         , count(*) FILTER (WHERE y = TRUE) OVER (ORDER BY prediction DESC)         , count(*) FILTER (WHERE y = TRUE) OVER () AS total_true </pre>



Step	Action
3 cont'd	<pre> FROM (   SELECT     r.*     , 1. / (1. + exp(-dotProduct(r.x, c.coef))) AS prediction   FROM     artificiallogreg AS r   CROSS JOIN     logr_coef as c ) q ) p ; \o </pre>
4	<p><b><u>Review results and plot</u></b></p> <p>Review the results in the file “graphics.txt”. You can use “GNUplot” or “R” or even EXCEL to plot the results. This task is left un-scripted as a student exercise.</p>
5	<p>The code in step 2 and 3 above is available in file /home/gpadmin/LAB12/logistic.sql</p>

*End of Lab Exercise*



## Final Lab Exercise on Big Data Analytics

<b>Purpose:</b>	This lab allows students to apply what they have learned from the analytical methods and tools to a big data problem using the Analytics Lab Environment.
<b>Tasks:</b>	<p>Tasks you will complete in this lab exercise include:</p> <ul style="list-style-type: none"> <li>• Explore the big data set provided and prepare the data for analysis</li> <li>• Assess data quality, outliers and training sets</li> <li>• Conduct model selection, code, execute and score the model</li> <li>• Use R and PSQL statements during your analysis of big data</li> <li>• Create a narrative summary of your findings, using the methods covered earlier in this module</li> </ul>
<b>References:</b>	<p>References used throughout the labs are located in your <b><i>Student Resource Guide Appendix</i></b>. See the Appendix for:</p> <ul style="list-style-type: none"> <li>• <a href="http://www.ffiec.gov/hmda/">http://www.ffiec.gov/hmda/</a></li> </ul>
<b>Working directory:</b>	<p>The directory /home/gpadmin/FINAL_LAB in your lab environment will be your working directory for the final lab exercise. Following files are pre-loaded in this lab:</p> <p>Analyst.ppt – Analyst presentation template</p> <p>Sponsor.ppt – Sponsor presentation template</p> <p>*.asc – encrypted files with suggested code for the solution. The decrypting of these files are performed with the following command at the \$ prompt in the FINAL_LAB directory:</p> <pre>gpg -o *.* -d *.*.asc</pre> <p>(*.* represents the filename with extension name)</p> <p>You will be prompted for a passphrase. Your instructor will provide the pass phrase</p>

## Case Study Background and Problem Definition

<b>Scenario</b>	<p>A financial planning company, FPC would like to expand the set of services they offer by creating an online site for loan advice. Potential home loan borrowers can enter information about their personal finances and the kind of home loan they want, and the site will return the probability of getting such a loan, along with some general advice about how to increase their likelihood of success. For example, the advice could be: "Increase the down payment so as to decrease the loan amount by X dollars"; or "Consider a home in the price range Y"; "Are you eligible for a particular type of loan?", or "Can you add a co-signer to the loan?".</p> <p>The company hopes this online service will be a lead-in for customers to come to FPC for more focused, personal financial planning to achieve their life goals. FPC would also consider partnering with a real estate broker to showcase houses to potential homebuyers. FPC realizes that the customers are looking for fast responses and the online-service must provide an answer within 45 seconds. FPC plans to enter into a service level agreement with the partner websites such as those managed by the real estate brokers.</p> <p>Ideally, the model behind this advice site can give reasonable, grounded predictions. Of course, the site cannot ask applicants to fill out an entire loan application and the sensitive data it contains, such as credit scores, employment history, or existing debt. The FPC project stakeholders want to stick to basic, easy to enter information such as applicant income, loan type, loan size, and the location of the property (ZIP code).</p> <p>They recognize that a model with only that information can only give general advice, rather than truly precise predictions. We have a set of data that can support this approach, and allow making predictions based on the information above.</p>
-----------------	--

<b>Issues to Address</b>	<p>A number of issues came up during the kick-off meeting for the project:</p> <ol style="list-style-type: none"> <li>1. Should there be one big model, or separate models for different types of loans?</li> <li>2. Someone in the group wondered if personal demographic information (sex, gender, and ethnicity) would improve the prediction. The others are hesitant about the idea of asking such questions on the site, but agreed to explore whether knowing that information would improve the model.</li> <li>3. Someone else offered the opinion that giving the users raw probabilities would not be meaningful to them. She suggested that the model should set thresholds, and deliver qualitative messages instead, such as the following: <ul style="list-style-type: none"> <li>• If the model reports that the probability of getting a loan were greater than 75%, then the system would send the user a message such as: "Congratulations! You have a very good chance of getting your loan!"</li> <li>• For probability less than 50%: "Sorry. Looks like the chances aren't so good," with a link to FPC's advice page.</li> <li>• For probability between 50-75%: "Your chances aren't the strongest. Come talk to us about developing a plan to improve your chances of getting financing."</li> </ul> </li> </ol> <p>This work led the group to a metric for measuring model performance. Of the people who score &gt; 75%, do more than 75% of them actually get a loan? Likewise for people who score less than 50%, how many of them are actually get loan? Also, how many people in the general population get each message? For instance, does the entire population score more than 75%?</p>
<b>Data Scientist Goals</b>	<p>Your goal, as the data scientist on this project, is to answer the following questions:</p> <ol style="list-style-type: none"> <li>1. Would it be more effective to develop different models or one model? Why? If different models, focus on a single one for the initial study.</li> <li>2. Should we ask for personal demographic information, or can we build a good enough model without it?</li> <li>3. How accurate is the model, in terms of the thresholds that the stakeholders set in their discussion (75% and 50%)? What is the coverage of the threshold regimes?</li> <li>4. Provide suggestions for the kind of general advice FPC can put on their advice page.</li> </ol>

<p><b>Considerations for developing an Analytic plan</b></p>	<ul style="list-style-type: none"> <li>• Consider the scope of the data you will need to include in the analysis, and filters you may need to set to construct the data set for your analysis.</li> <li>• Consider the types of models best suited to perform the analysis needed for the new website engine. Does this scenario represent a classification, clustering, or prediction problem?</li> <li>• Examine the distribution of data, such as loan data for home improvement, home purchase, and refinancing loans, to identify the influences on how you will select and create the model</li> <li>• Look at creating several models and compare them in terms of ROC/AUC, or other performance metrics.</li> <li>• Find ways to examine how robust the model is with the help of a confusion matrix or similar diagnostic technique.</li> <li>• Give thought to how you would portray this information to business stakeholders as well as an analytical audience.</li> <li>• Consider the Service Level Targets that FPC can offer to their end users when they score the model with their inputs</li> <li>• Consider Service Level Targets that you can provide to FPC in terms of computational resources required for model generation and validation.</li> <li>• Provide some suggestions for the kind of general advice FPC can put on their advice page, based on the results from your modeling exercise. For instance, mention the types of things an applicant can do to increase their likelihood of success when applying for a loan on the website.</li> </ul>
--	--

## Suggested Workflow and Checkpoints for the Lab

### Suggested Workflow

#### Checkpoint 1

1. The data for this lab is the housing loan database assembled by federal agencies pursuant to the Home Mortgage Disclosure Act (HMDA). This database identifies the census tract location of almost every housing loan and housing loan application made in the United States each year. The data provided for analysis in this lab is an extract for the year 2010. The data is organized in three database tables larDB1, larDB2, larDB3 (in database “hmdalab”) for different states as follows:

larDB1	larDB2	larDB3
AK	AL	CT
AZ	AR	DC
CA	CO	DE
HI	GA	FL
ID	IA	MA
MN	IL	MD
MT	IN	ME
ND	KS	Na
NM	KY	NC
NV	LA	NH
OR	MI	NJ
SD	MO	NY
UT	MS	PA
WA	NE	RI
WI	OH	SC
WY	OK	VA
	PR	VT
	TN	
	TX	
	WV	

2. The tables provide the HMDA Loan Application Registration (lar) details and they have the following structure:

```
As_of_Year INTEGER,  
Respondent_Id VARCHAR(10),  
Agency_Code VARCHAR(1),  
Loan_Type INTEGER,  
Property_Type VARCHAR(1),  
Loan_Purpose INTEGER,  
Occupancy INTEGER,  
Loan_Amount_inK INTEGER,  
Preapproval VARCHAR(1),  
Action_Type INTEGER,  
MSAMD VARCHAR(5),  
State_Code VARCHAR(2),  
County_Code VARCHAR(3),  
Census_Tract_Number VARCHAR(7),  
Applicant_Ethnicity VARCHAR(1),  
Co_Applicant_Ethnicity VARCHAR(1),  
Applicant_Race_1 VARCHAR(1),  
Applicant_Race_2 VARCHAR(1),  
Applicant_Race_3 VARCHAR(1),  
Applicant_Race_4 VARCHAR(1),  
Applicant_Race_5 VARCHAR(1),  
Co_Applicant_Race_1 VARCHAR(1),  
Co_Applicant_Race_2 VARCHAR(1),  
Co_Applicant_Race_3 VARCHAR(1),  
Co_Applicant_Race_4 VARCHAR(1),  
Co_Applicant_Race_5 VARCHAR(1),  
Applicant_Sex INTEGER,  
Co_Applicant_Sex INTEGER,  
Applicant_Income_inK VARCHAR(4),  
Purchase_Type VARCHAR(1),  
Denial_Reason_1 VARCHAR(1),  
Denial_Reason_2 VARCHAR(1),  
Denial_Reason_3 VARCHAR(1),  
Rate_Spread VARCHAR(5),  
HOEPA_Status VARCHAR(1),  
Lien_Status VARCHAR(1),  
Edit_Status VARCHAR(1),  
Sequence_Number VARCHAR(7),  
Population VARCHAR(8),  
Minority_Population_pct VARCHAR(6),  
HUD_Median_Family_Income VARCHAR(8),  
Tract_To_MSAMD_Income_pct VARCHAR(6),  
Number_of_Owner_occupied_units VARCHAR(8),  
Number_of_1_to_4_Family_units VARCHAR(8),  
Application_Date_Indicator INTEGER);
```



3. All the required codes for the modeling exercise are made available in different tables as detailed below:

Table name	variable defined
action	<b>Action_Type</b>
counties	<b>County_Code</b>
ethnicity	<b>Applicant_Ethnicity</b>
fips	<b>State_Code</b>
inst	Institution Record format
lienstatus	<b>Lien_Status</b>
loanpurpose	<b>Loan_Purpose</b>
loantype	<b>Loan_Type</b>
msamd	MSAMD office format
preapproval	<b>Preapproval</b>
race	<b>Applicant_Race_1</b>
sex	<b>Applicant_Sex</b>

Property type is not coded in a table, but has code definitions as follows:

- 1: 1 to 4 family
- 2: Manufactured housing
- 3: Multi-family

Occupancy = 1 indicates owner occupied housing (our focus of analysis)

4. For your analysis you are required to select
- a. A single state
  - b. Occupancy = 1
  - c. Property\_Type = 1
  - d. Action\_Type <= 4
5. Extract data from the "lar" table (with the conditions in step 4) and create a table with the following variables:

```
Loan_Type VARCHAR(20) ,  
Loan_Purpose VARCHAR(25) ,  
Loan_Amount_inK INTEGER,  
Preapproval VARCHAR(25) ,  
Action_Type VARCHAR(25) ,  
County_Name VARCHAR(50) ,
```

	<pre> <b>Applicant_Ethnicity</b> VARCHAR(25) , Co_Applicant_Ethnicity VARCHAR(1) , <b>Applicant_Race_1</b> VARCHAR(25) , <b>Applicant_Sex</b> VARCHAR(25) , Applicant_Income_inK VARCHAR(4) , Rate_Spread VARCHAR(5) , HOEPA_Status VARCHAR(1) , <b>Lien_Status</b> VARCHAR(25) , Minority_Population_pct VARCHAR(6) , HUD_Median_Family_Income VARCHAR(8) , Tract_To_MSAMD_Income_pct VARCHAR(6) , Number_of_Owner_occupied_units VARCHAR(8) </pre> <p>The highlighted variables must be expanded to the values corresponding to the codes in the “lar” table.</p>
<p><b>Suggested Workflow</b></p> <p><b>Checkpoint 2</b></p>	<ol style="list-style-type: none"> <li>1. You can read the table created in Checkpoint 1 in your RStudio environment. Ensure the database “hmdalab” is specified in the “odbc.ini” file. When you read in the table with sqlFetch make sure that as.is=T is specified so that you have control in specifying the variables you want to treat as factors.</li> <li>2. In R, use the factor function to specify those variables highlighted in step 5.</li> <li>3. Convert the income, rate_spread and other variables to Numeric. Check for “NA” in the data</li> <li>4. Use Table Function to tabulate values (number of records) with specific value combinations. You can also plot these data to begin exploring the variables. For instance, compare action_type vs. Rate_spread.</li> <li>5. Remove records without income info (income=“NA”)</li> <li>6. Code appropriate “releveling” for variables to explain the models in subsequent analysis</li> <li>7. Define HUD_Median_Family_Income as numeric. Check for NAs</li> <li>8. Remove rows with nulls (NA) in Tract_To_MSAMD_Income_pct, Minority_Population_pct, Tract_To_MSAMD_Income_pct</li> <li>9. Visualize the variables. You should visualize all (or many, at least) of them. Check whether the loan amount has any odd, multi-modal distribution. This may suggest to us that we might want to build separate models for the different loan purposes.</li> <li>10. Look for spikes in the outputs. You may need to develop one model with all the data below that spike and develop a separate model with the data beyond that spike.</li> <li>11. Optionally eliminate some of the very small loans that trail out on the right.</li> </ol> <p>Subset data based on only those “action_types” you need for your model</p>

<b>Suggested Workflow</b>  <b>Checkpoint 3</b>	<ol style="list-style-type: none"> <li>1. Finalize your dependent variable (Hint: approved implied by action_type==originated) and the drivers. First model with all drivers</li> <li>2. Use the Log of monetary variables</li> <li>3. Model with 10% of data (small set) Hint use:  <pre> <b>probldata\$gp = runif(dim(probldata)[1])</b> <b>smallset = subset(probldata, probldata\$gp &lt;</b> <b>0.1) # 10% of data</b> </pre> </li> <li>4. Execute the chosen model with the smallset</li> <li>5. Analyze the results</li> <li>6. Drop/change variables (such as including or excluding personal demographic information) and repeat steps 3,4 and 5</li> <li>7. Predict using different models (with and without personal demographic info)</li> <li>8. Plot the ROC curves for all the models predicted</li> <li>9. Compare the AUC for the different models</li> </ol>
<b>Suggested Workflow</b>  <b>Final Presentation</b>	<ol style="list-style-type: none"> <li>1. Look at the original questions you had when you started. Can you answer the questions we started with?</li> <li>2. Examine the probability thresholds suggested by marketing. Do the evaluations on the “hold out” set.</li> <li>3. Use a confusion matrix to determine the probability of approval in each bin (high, medium and low probability)</li> <li>4. Use the presentation template to compile your results and develop your narrative summary</li> <li>5. Present your findings to the class.</li> </ol>

*End of Lab Exercise*