



Data Science And Big Data Analytics Course

Lab Exercise 2: Introduction to R

Purpose:	<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">• Read data sets into R, save them, and examine the contents
Tasks:	<p>Tasks you will complete in this lab include:</p> <ul style="list-style-type: none">• Invoke the R environment and examine the R workspace• Read tables created in Lab 1 into the R statistical package• Examine, manipulate and save data sets• Exit the R environment

LAB 2

Introduction to R

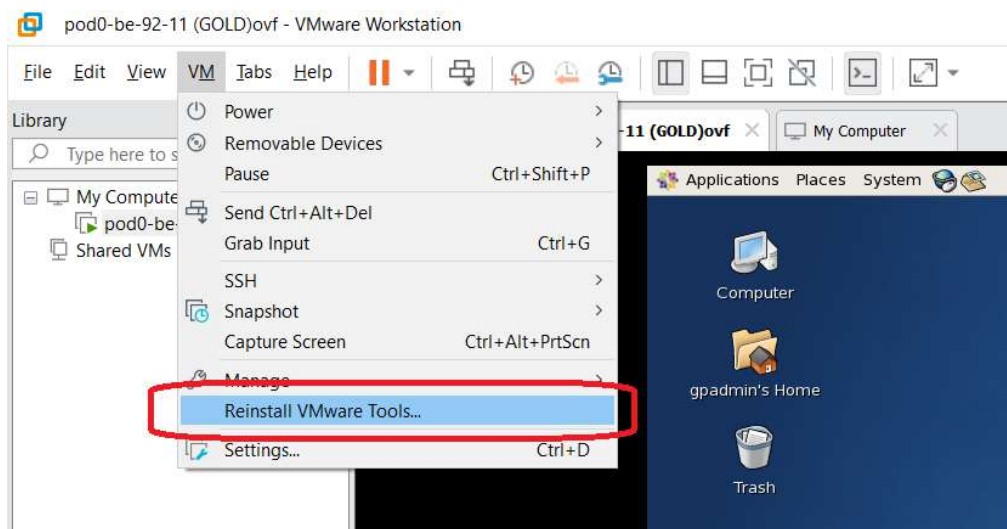
✓ Lab Content:

1. Install/update VMWare Tools
2. Configure the Network Adaptor on the MV
3. Work on RStudio from the host machine

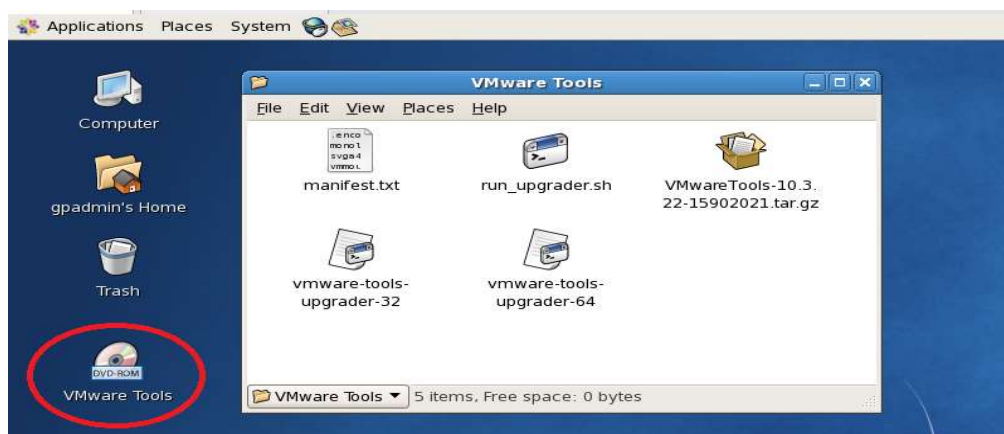
2. 1 Install/update VMWare Tools

We need to install/update VMWare tools to give better usability for the VM. One time job to be done and automated.

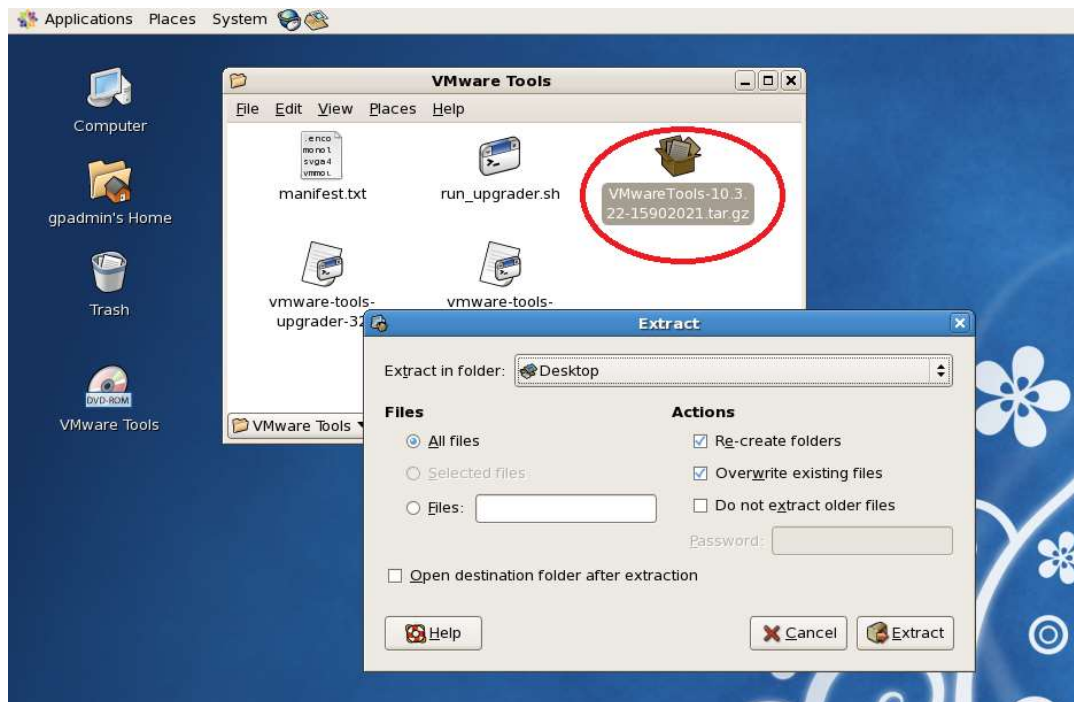
1. Start the VM and open the GUI mode and login with gpadmin.
2. From the VMware station app, select from the menu bar VM → Install VMware Tools...



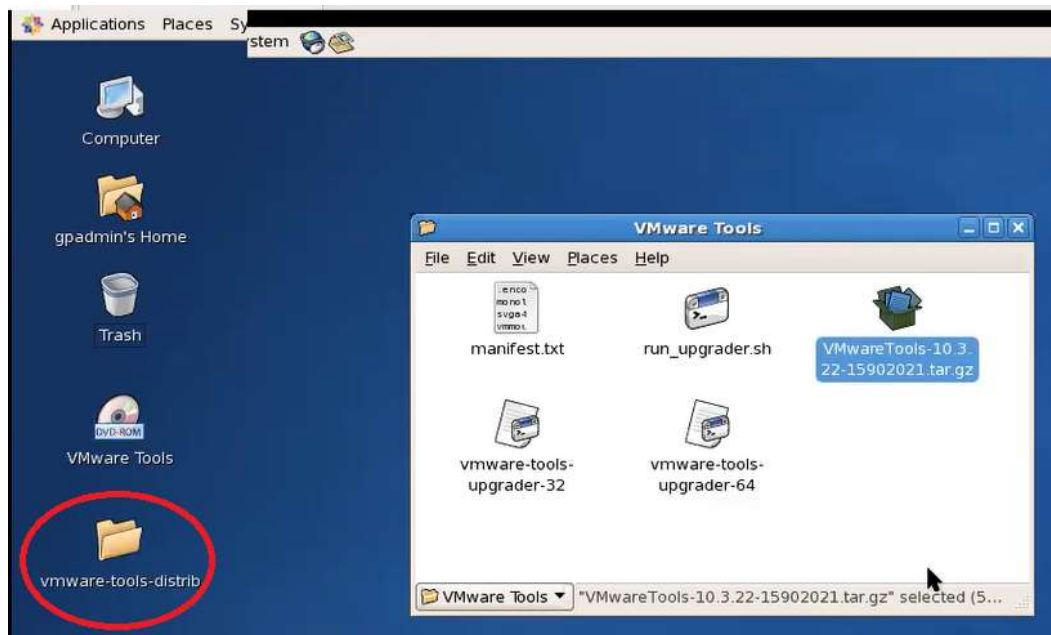
It will mount a virtual DVD the contains the contents of the VMware tools on the guest desktop and will open the drive of the VMware tools folder automatically.



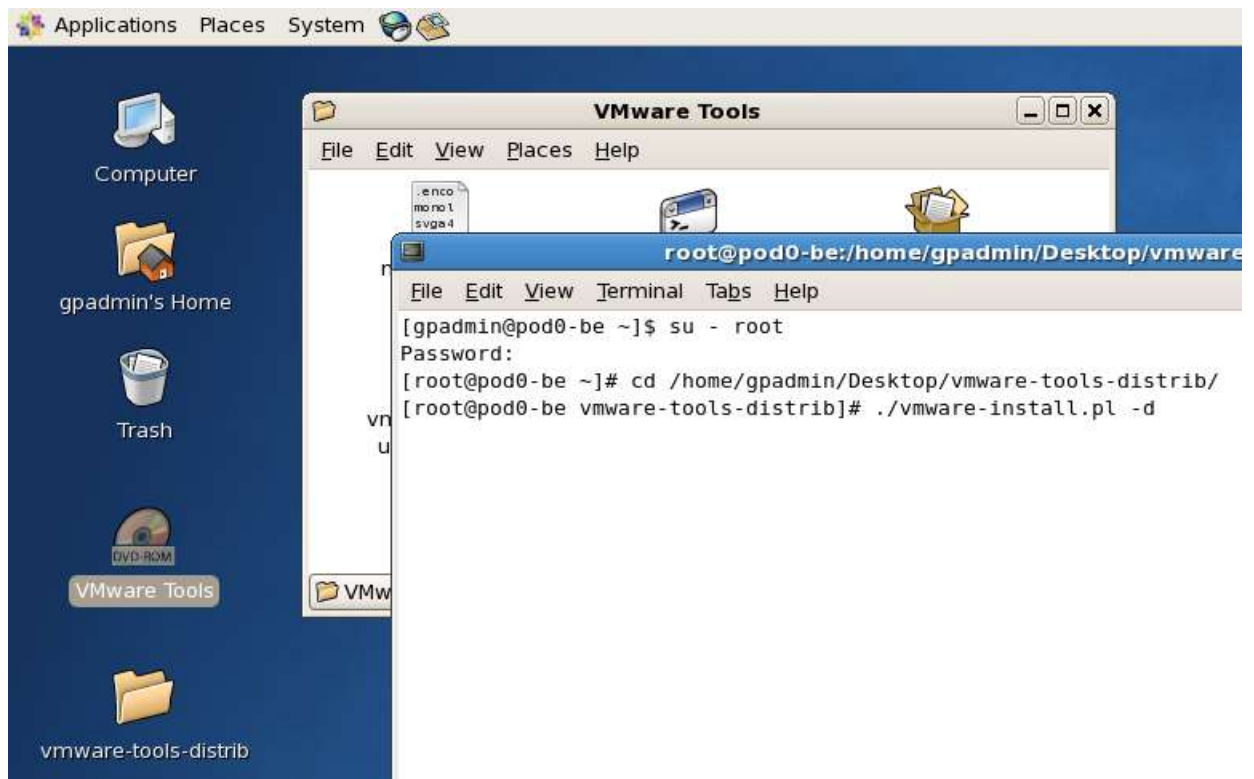
3. Right click on the “VMware Tools-10.3.22.....tar.gz” and select “Extract To...” and click extract.



After finishing the extracted folder will appear on the desktop



4. Open terminal (right click on desktop and choose open terminal) and write the following commands



After pressing enter it will take few seconds, wait till see this message:

```
Stopping VMware Tools services in the virtual machine:
  Guest operating system daemon:          [ OK ]
  VMAuthService:                          [ OK ]
  VMware User Agent (vmware-user):        [ OK ]
  Blocking file system:                   [ FAILED ]
  Unmounting HGFS shares:                 [ OK ]
  Guest filesystem driver:                 [ OK ]
  Guest memory manager:                   [ OK ]
  VM communication interface socket family: [ OK ]
  VM communication interface:             [ OK ]
Unable to stop services for VMware Tools

Execution aborted.

Found VMware Tools CDRom mounted at /media/VMware Tools. Ejecting device
/dev/hdc ...
Enjoy,

--the VMware team

[root@pod0-be vmware-tools-distrib]#
```

5. When reaching this previous message, then you are done and finished installing, you need now to restart your VM. From System → Shut Down, then click restart.



After restarting, Now you are done. You can delete these folders:



2. 2 Configure the Network Adaptor on MV

The VM we are working on is taken as a copy from a VM that was deployed on a cloud, so it has the IP addresses and configurations and MAC address of the physical machine that it was deployed on; and so we will not be able to access the internet from this VM.

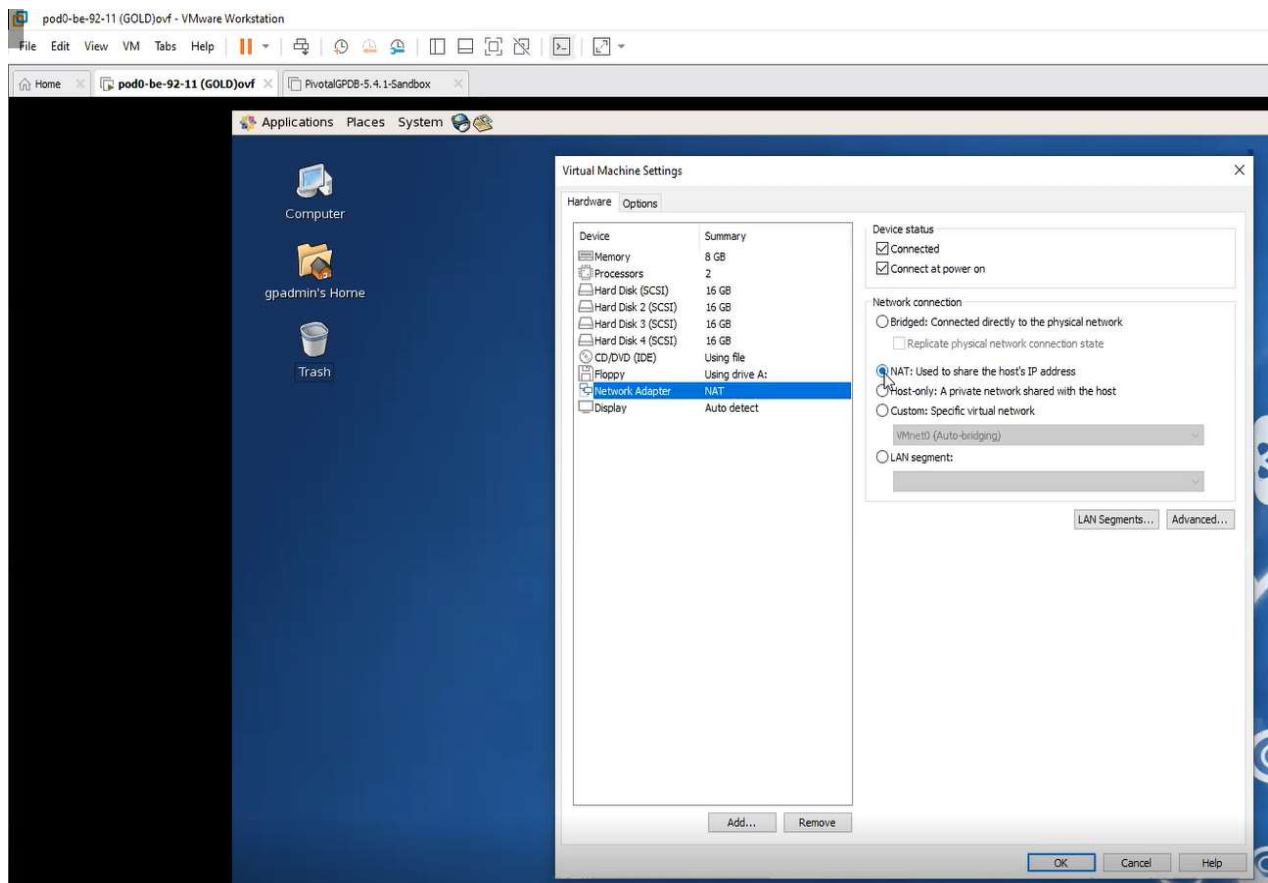
We need to configure the Network Adaptor by these steps:

1- Select System→Administration→Network



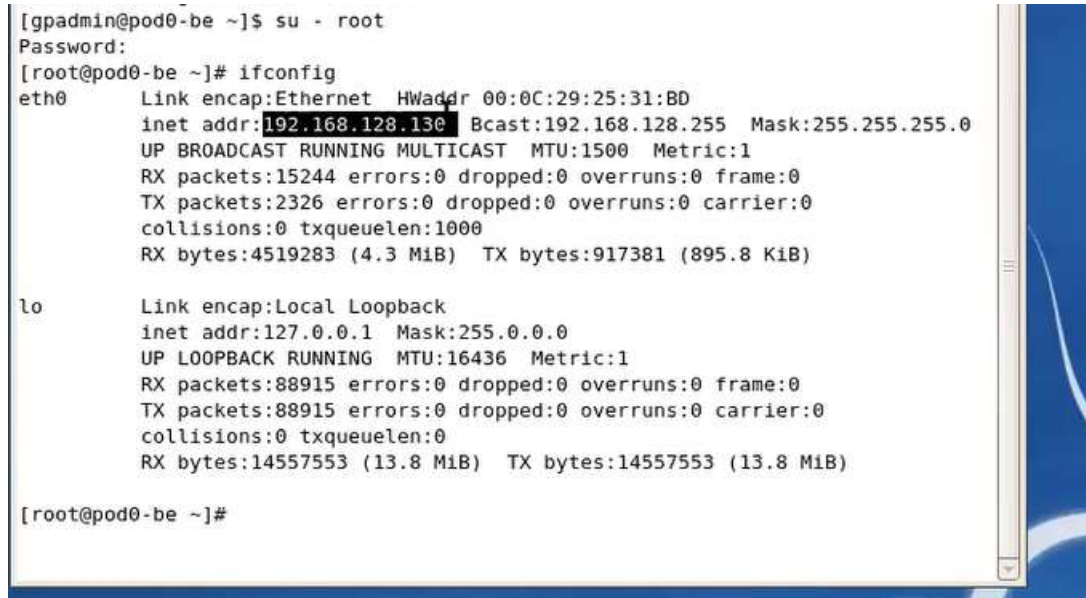
2- “Network configuration” window will open and the Status is “InActive”

- 3- Click on “Edit” .. “Ethernet Device” configuration window will open.
- 4- Select the radio button “Automatically obtain IP address settings with ..” then select the “Hardware Device” tab and click the “Probe” button .. You will notice that it will take new MAC address from your physical machine.. Then click “OK”.
- 5- From the VMware workstation menu bar, click VM → Settings.. From the left Select “Network Adapter” then from right select “NAT” then press “OK”



- 6- Back in your VM, in the Network Configuration window, click on “Activate” .. It will take few time then it should show Active Status.

- 7- Now the VM took an IP address, we need to know it. Open the terminal and switch user to root and type ifconfig



```
[gpadmin@pod0-be ~]$ su - root
Password:
[root@pod0-be ~]# ifconfig
eth0      Link encap:Ethernet  HWaddr 00:0C:29:25:31:BD
          inet addr:192.168.128.130  Bcast:192.168.128.255  Mask:255.255.255.0
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:15244 errors:0 dropped:0 overruns:0 frame:0
          TX packets:2326 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:4519283 (4.3 MiB)  TX bytes:917381 (895.8 KiB)

lo        Link encap:Local Loopback
          inet addr:127.0.0.1  Mask:255.0.0.0
          UP LOOPBACK RUNNING  MTU:16436  Metric:1
          RX packets:88915 errors:0 dropped:0 overruns:0 frame:0
          TX packets:88915 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:0
          RX bytes:14557553 (13.8 MiB)  TX bytes:14557553 (13.8 MiB)

[root@pod0-be ~]#
```

Remember the IP address appeared to you. Now you can test that its working by pinging on it from the cmd from your host machine.

- 8- To open the RStudio from the host machine, open your internet browser and visit RStudio: <IP>:8787

Login by username: gpadmin, password: changeme

- 9- Now you are logged on the RStudio server that contains your datasets. (N.P. The setup of the VM is made as a RStudio server, and we now can access it through our browser using it IP then port 8787)

2. 3 Work on RStudio from host machine

Step	Action
1	<p><u>Invoke the R Environment:</u></p> <p>Logon to RStudio environment.</p> <ol style="list-style-type: none"> 1. The RStudio is accessed through the “safari” browser available as a desktop icon 2. Refer to the access details provided by your instructor. RStudio is accessed with URL <a href="http:<Your assigned IP-address for the backend system>:8787/">http:<Your assigned IP-address for the backend system>:8787/ 3. The RStudio Userid and Password are gpadmin and changeme, respectively. <p>This will start your web browser and connects to the “be” server. You should see the RStudio four- panel display.</p> <p>Verify that you see the following text in the lower left-hand pane:</p> <pre>R version 2.15.0 (2012-03-30) Copyright (C) 2012 The R Foundation for Statistical Computing ISBN 3-900051-07-0 Platform: x86_64-redhat-linux-gnu (64-bit) --- <snip> --- 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 ></pre>
2	<p><u>Examine the Workspace:</u></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>Note: R is telling you that you have nothing in your workspace.</p>

Step	Action
3	<p><u>Getting Familiar with R:</u></p> <ol style="list-style-type: none"> 1. Click each tab in each panel. What happens? 2. Type the following commands into the R command panel <pre>help() help.start() demo() demo(graphics)</pre> <p>Hit esc to exit out of the demo</p>
4	<p><u>Read-in the Lab Script:</u></p> <ol style="list-style-type: none"> 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 the file “Module3Lab1.R”.) 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"> ○ Position your cursor inside the line that represents the command you wish to execute. ○ 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. 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. <p>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.</p> <ol style="list-style-type: none"> 1. On the 1st line in Section 1, put your cursor on the line containing the word <code>ls()</code>. 2. Click Run. The <code>ls()</code> command will execute in the command window and show you the contents of your workspace.

Step	Action
5	<p><u>Reading external data:</u></p> <p>Load the .txt files you created in the first lab. Load the first file, lab1_01.txt</p> <ol style="list-style-type: none"> 1. Set the working directory to LAB01 where we have stored the data. On the console window type: <pre>setwd("~/LAB01")</pre> 2. Select the line and press <ctl>Enter: <pre>lab1 <- read.table("lab1_01.txt", sep=" ", header=TRUE)</pre> <ul style="list-style-type: none"> • If correct, R will simply return you to the command prompt ("> "). 3. Now load the second .txt file, lab1_02.txt, 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 <- read.table("lab1_02.txt", sep=" ", header=TRUE)</pre> 4. When you have completed the edits, make sure that your cursor is within the line, press Enter. <p>Note: 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><u>Verify the Contents of the Tables:</u></p> <p>It is always a good idea to look at the data to make sure that everything works. You can use the head() command to print out the first 10 lines of a table or the, tail() command to print out the last 10 lines of the table.</p> <ol style="list-style-type: none"> 1. Select and run the command: <pre>head(lab1, n=10)</pre> Record the value of the 10th line here: _____ 2. Now do the same for the lab2 table, but use the tail(lab2, n=10) command instead. 3. Record the value of the 1st line here: _____

Step	Action
7	<p><u>Manipulating data frames in R:</u></p> <p>Examine the contents of the table in more detail.</p> <ol style="list-style-type: none"> 1. Execute the following command: <pre>summary(lab1)</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> 2. Select and run: <pre>nlab1 <- lab1[,2:3]</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>hinc <- lab1\$hinc rooms <- lab1\$rooms nlab1 <- data.frame(hinc, rooms)</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>nlab1 <- data.frame(lab1\$hinc, lab1\$rooms) names(nlab1) = c("hinc", "rooms")</pre>

Step	Action
7 Cont.	<p>3. The <code>dim(<table>)</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><u>Investigate Your Data:</u></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><u>Save the Data Sets:</u></p> <p>Execute the following commands:</p> <pre>rm(lab1) lab1 <- nlab1 save(lab1, lab2, file="Labs.Rdata") rm(lab1, lab2) ls() # make sure they're not in the workspace</pre>

Step	Action
10	<p><u>Continue Investigating the Data:</u></p> <ol style="list-style-type: none"> 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"> 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. Experiment with R functions that identify the class and data type of a particular variable, type: <pre>typeof(x), class(x), attributes(x), names(x), dim(x)</pre> Which ones work on which kind of data types? _____ Type these values into the RStudio command panel. _____ 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 <code>"tellme <- function(x)"</code> {extending through the right curly brace. Now execute the following command tellme You should see the definition of the function that you just entered! This is because R doesn't interpret a plain tellme 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 mean and inspect the results. Try <code>tellme ()</code> with a series of variables. Which commands actually list something? _____ How might you get the other commands to list their return value? [Hint: try <code>print()</code>]
11	<p><u>Exit R:</u></p> <ol style="list-style-type: none"> Execute the following command: q() R will ask you if you want to save your workspace. Answer "no".

End of Lab Exercise