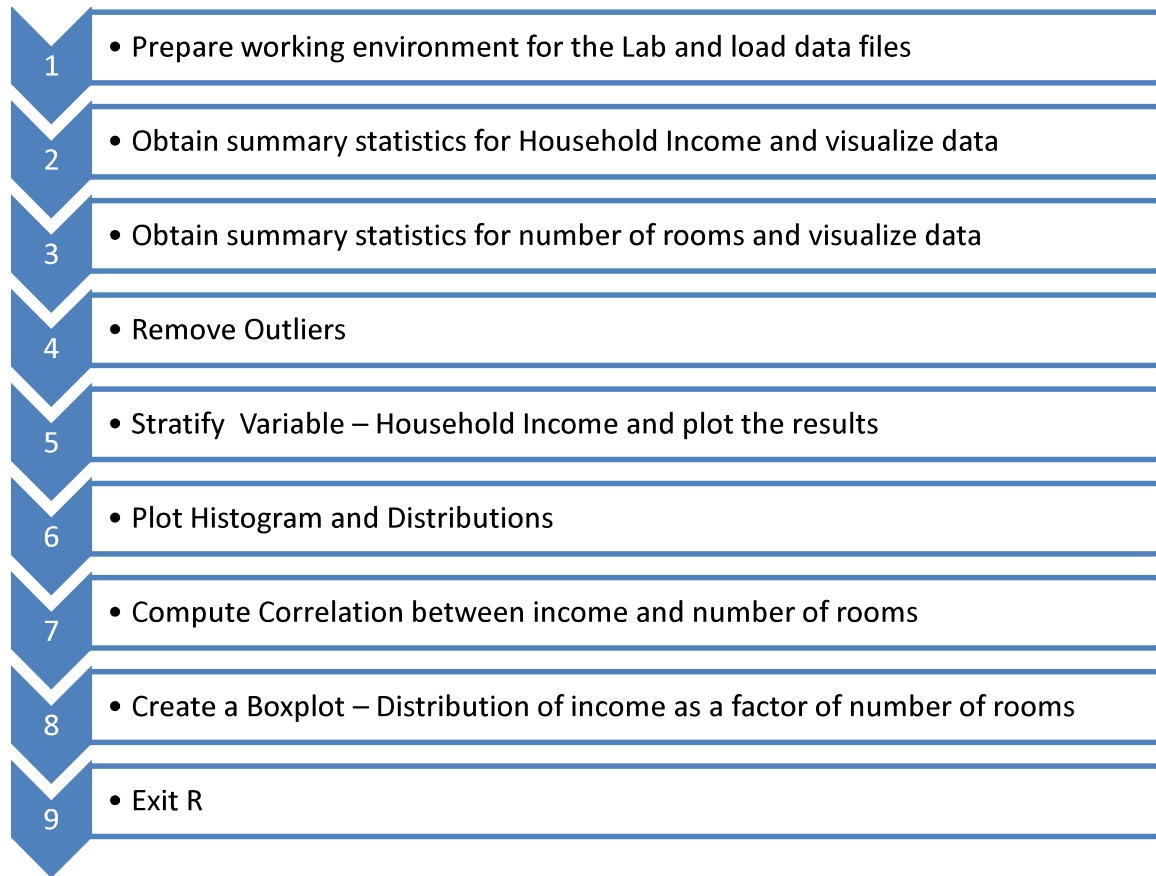


Lab Exercise 3: Basic Statistics, Visualization, and Hypothesis Tests

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

Part 1 – Basic Statistics and Visualization Using R

Workflow Overview



LAB Instructions

Step	Action
1	<p><u>Prepare working environment for the Lab and load data files:</u></p> <ol style="list-style-type: none"> Set the working directory to LAB01 where we have stored the data. On the console window type: <pre>setwd("~/LAB01")</pre> 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: Execute the following commands from the script window: <pre>options(digits=3) ls() load(file="Labs.Rdata") ls() rm(lab2) ds <- lab1 colnames(ds) <- c("income", "rooms")</pre>
2	<p><u>Examine Household Income:</u></p> <ol style="list-style-type: none"> 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)) # right skewed</pre> What is the mean? _____ What is the median? _____ What is the standard deviation? _____

Step	Action
3	<p><u>Examine the Number of Rooms:</u></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><u>Remove Outliers:</u></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"> Does this make sense to you? Why? * What happens if you throw out the top and bottom 10%? Execute the following line from the script window <pre>(m <- mean(ds\$income, trim=0.10))</pre> How does this compare to the previous mean of this variable? Execute the following commands from the script window: <pre>ds <- subset(ds, ds\$income >= 10000 & ds\$income < 1000000) summary(ds) quantile(ds\$income, seq(from=0, to=1, length=11))</pre> How do these values vary from the values in the original data set? Do they make more sense? Which data set would you prefer to use? <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 ± 3 sds from the mean.</p>

Step	Action
5	<p><u>Stratify a Variable – Household Income:</u></p> <p>Stratify breaks that occur close to U.S. Guidelines for Poverty, LowerMid, UpperMid, Wealthy, and Rich (> \$250k @ year)</p> <ol style="list-style-type: none"> Execute the following code (listed under comment heading “step 5” in the script file): <pre>breaks <- c(0, 23000, 52000, 82000, 250000, 999999) labels <- c("Poverty", "LowerMid", "UpperMid", "Wealthy", "Rich") wealth <- cut(ds\$income, breaks, labels) # add wealth as a column to ds ds <- cbind(ds, wealth) # show the 1st few lines. head(ds)</pre> Continue to execute the remaining part of the code in Step 5 <pre>wt <- table(wealth) percent <- wt/sum(wt)*100 wt <- rbind(wt, percent) wt plot(wt) #This does not seem to give good results, why?</pre> Take another look at the relationship between wealth and income. Execute the following lines: <pre># take another look -- wealth by rooms nt <- table(wealth, ds\$rooms) print(nt) plot(nt) #Nice mosaic plot</pre> 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>

Step	Action
6	<p><u>Plot Histogram and Distributions:</u></p> <p>Problem: How do you represent income given the range of values?</p> <ol style="list-style-type: none"> 1. Select and execute the code under Step 6 Histograms and distributions in the script file. <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"> 2. Now try the same thing with log10(income) <pre>logincome = log10(ds\$income) hist(logincome, main="Distribution of Household Income", freq=FALSE) lines(density(logincome), lty=2, lwd=2) #Line type lty(2) is dashed 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><u>Compute Correlation between income and number of rooms:</u></p> <ol style="list-style-type: none"> You need to consider your hypothesis. <ul style="list-style-type: none"> Your hypothesis is that the number of rooms in a house is predicted by household income (the rich can buy bigger houses), e.g. $lm(\text{rooms} \sim \text{income})$ Therefore, our null hypothesis: no correlation between income and number of rooms. Alternate hypothesis: there is a correlation between income and the number of rooms. 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> For comparison, correlate rooms with a completely unrelated variable. <pre>n = length(ds\$income) with(ds, cor(runif(n), rooms))</pre>
8	<p><u>Create a Boxplot - Distribution of income as a factor of number of rooms:</u></p> <ol style="list-style-type: none"> Select and execute the code (Listed after the comment line "Step 8") in the script window. 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> Plot the # of rooms as a function of wealth level. <pre>boxplot(rooms ~ wealth, data = ds, main="Room by Wealth", xlab="Category", ylab="# rooms")</pre>

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

End of Lab Exercise