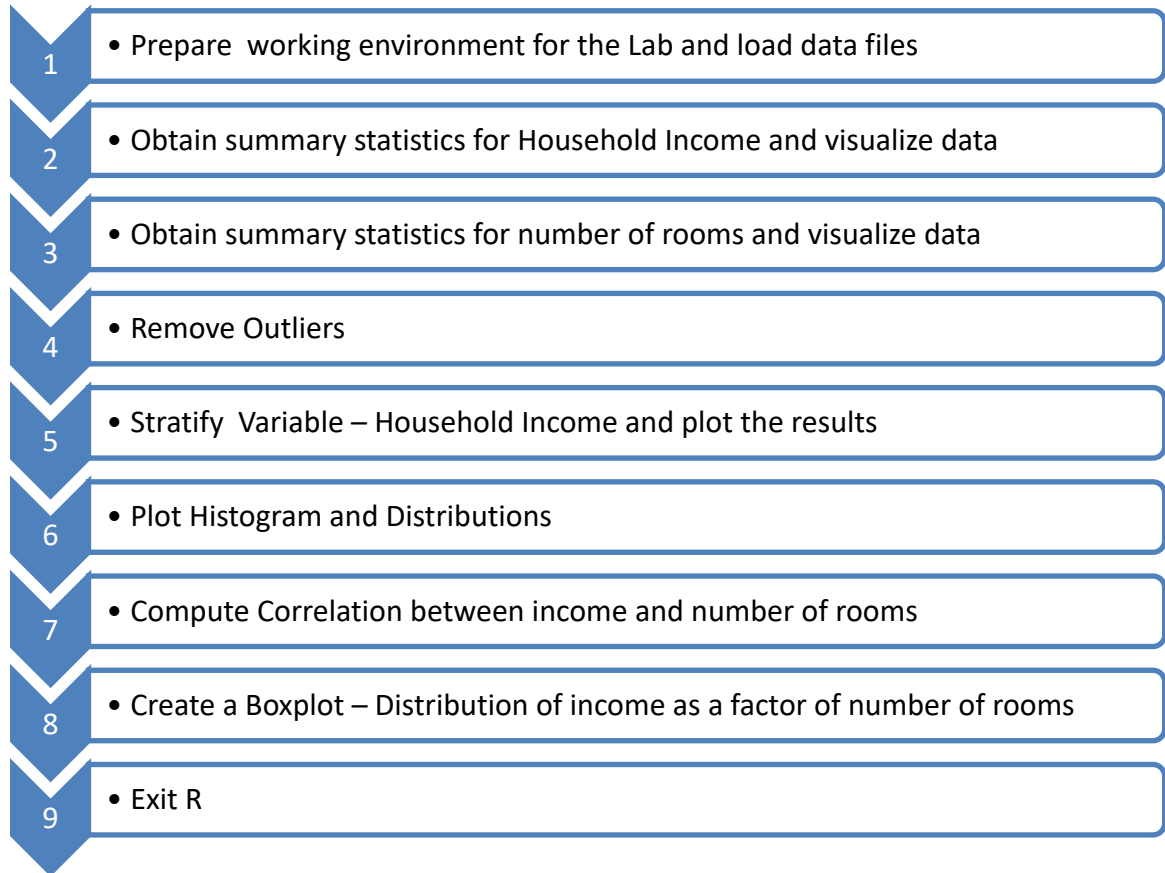Hannah Roach
2/8/2019

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

| Purpose: | 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 able to: |
|---|---|
| | <ul><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> |

| Tasks: | Tasks you will complete in this lab include: |
|---|---|
| | <ul><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> |

| References: | References used in this lab are located in your **Student Resource Guide Appendix**. See the Appendix for: |
|---|---|
| | <ul><li>R Commands – Quick Reference</li><li>Surviving LINUX – Quick Reference</li></ul> |

# Part 1 – Basic Statistics and Visualization Using R

# Workflow Overview

1 • Prepare working environment for the Lab and load data files

2 • Obtain summary statistics for Household Income and visualize data

3 • Obtain summary statistics for number of rooms and visualize data

4 • Remove Outliers

5 • Stratify Variable – Household Income and plot the results

6 • Plot Histogram and Distributions

7 • Compute Correlation between income and number of rooms

8 • Create a Boxplot – Distribution of income as a factor of number of rooms

9 • Exit R

# LAB Instructions

| Step | Action |
|------|--------|
| 1 | **Prepare working environment for the Lab and load data files**<br>1. Set the working directory to LAB01 where we have stored the data. On the console window type:<br>`setwd("~/LAB01")`<br><br>2. 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").<br>Start R and Read the Data Set Back Into Your Workspace:<br><br>3. Execute the following commands from the script window:<br>`options(digits=3)`<br>`options(width=68)`<br><br>`ls()`<br>`load(file="Labs.Rdata")`<br>`ls()`<br><br>`rm(lab2)`<br><br>`ds <- lab1`<br>`colnames(ds) <- c("income", "rooms")` |
| 2 | **Obtain summary statistics for Household Income and visualize data:**<br><br>1. Execute the following commands from the script window:<br><br>`summary(ds$income)`<br>`range(ds$income)`<br>`sd(ds$income)`<br>`var(ds$income)`<br><br>`plot(density(ds$income))   # left skewed`<br><br>2. What is the mean? 67200<br>3. What is the median? 50300<br>4. What is the standard deviation? 68178 |

| Step | Action |
|---|---|
| 3 | **Obtain summary statistics for Number of rooms and visualize data:**<br><br>Execute the following commands from the script window:<br>`summary(ds$rooms)`<br>`range(ds$rooms)`<br>`sd(ds$rooms)`<br>`plot(as.factor(ds$rooms))`<br><br>What is the mean? 5.63<br><br>What is the median? 6.00<br><br>What is the standard deviation? 1.99 |

| Step | Action |
|------|--------|
| 4 | **Remove Outliers**<br><br>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.<br><br>1. <mark>Does this make sense to you? Why? Yes, this</mark> make sense because the median is 50,300 and the 3$^{rd}$ quarter is 84200 and the plot appears to be left skewed.<br>2. <mark>What happens if you throw out the top and bottom 10%? Execute the following line from the script window</mark><br><br>`(m <- mean(ds$income, trim=0.10) )`<br><br><mark>The new mean becomes 55,347</mark><br><br>3. <mark>How does this compare to the previous mean of this variable? The new mean with the outliers removed is lower than the original mean.</mark><br>4. Execute the following commands from the script window:<br>`ds <- subset(ds, ds$income >= 10000 & ds$income < 1000000)`<br>`summary(ds)`<br>`quantile(ds$income, seq(from=0, to=1, length=11))`<br><br>5. <mark>How do these values vary from the values in the original data set? These values only include incomes between $10,000 and $1,000,000.</mark><br>6. <mark>Do they make more sense? Yes</mark><br>7. <mark>Which data set would you prefer to use? The second dataset because it excludes the outliers.</mark><br><br>_____<br><br>*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. |

| Step | Action |
|------|--------|
| 5 | **Stratify  Variable – Household Income and plot the results:**<br><br>Stratify breaks that occur close to U.S. Guidelines for Poverty, Median Income, Wealth, and Rich (> $250k @ year)<br><br>1.  Execute the following code (listed under comment heading "step 5" in the script file):<br><br>```r<br>breaks <- c(0, 23000, 52000, 82000, 250000, 999999)<br>labels <- c("Poverty", "LowerMid", "UpperMid",<br>"Wealthy", "Rich")<br>wealth <- cut(ds$income, breaks, labels)<br># add wealth as a column to ds<br>ds <- cbind(ds, wealth)<br># show the 1st few lines.<br>head(ds)<br>```<br><br>2.  Continue to execute the remaining part of the code in Step 5<br><br>```r<br>wt <- table(wealth)<br>percent <- wt/sum(wt)*100<br>wt <- rbind(wt, percent)<br>wt<br>plot(wt)<br>```<br><br>3.  Take another look at the relationship between wealth and income. Execute the following lines:<br><br>```r<br># take another look -- wealth by rooms<br><br>nt <- table(wealth, ds$rooms)<br>print(nt)<br>plot(nt)          # nice mosaic plot<br>```<br><br>4.  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".<br><br>```r<br>rm(wealth,breaks,labels)<br>save(ds, wt, nt, file="Census.Rdata")<br>``` |

| Step | Action |
|------|--------|
| 6 | **Plot Histogram and Distributions:**<br><br><mark>Problem: How do you represent income given the range of values? Given this rage of values, you could illustrate the data as a histogram.</mark><br><br>1. Select and execute the code under Step 6 Histograms and distributions in the script file.<br><br>```r<br>library(MASS)<br><br>with(ds, {<br>  hist(income, main="Distribution of Household Income",<br>freq=FALSE)<br>  lines(density(income), lty=2, lwd=2)<br># line type (lty) 2 is dashed<br>  xvals = seq(from=min(income), to=max(income),<br>length=100)<br>  param = fitdistr(income, "lognormal")<br>  lines(xvals, dlnorm(xvals, meanlog=param$estimate[1],<br>sdlog=param$estimate[2]), col="blue")<br>})<br>```<br><br>2. Now try the same thing with log10(income)<br><br>```r<br>logincome = log10(ds$income)<br>hist(logincome, main="Distribution of Household Income",<br>freq=FALSE)<br># line type lty(2) is a dashed line<br>lines(density(logincome), lty=2, lwd=2)<br>xvals = seq(from=min(logincome), to=max(logincome),<br>length=100)<br>param = fitdistr(logincome, "normal")<br>lines(xvals, dnorm(xvals, param$estimate[1],<br>param$estimate[2]), lwd=2, col="blue")<br>``` |

| Step | Action |
|---|---|
| 7 | **Compute Correlation between income and number of rooms:**<br><br>1. You need to consider your hypothesis.<br><br>    • 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(rooms ~ income)*<br>    • Therefore, our null hypothesis: no correlation between income and number of rooms.<br>    • Alternate hypothesis: there is a correlation between income and the number of rooms.<br><br>4.   Execute the following code (listed after the comment line "Step7 in the script file).<br><br>`with(ds, cor(income, rooms))`<br><br>`with(ds, cor(log(income), rooms))) # this will give a better correlation`<br><br>5. For comparison, correlate rooms with a completely unrelated variable.<br>`n = length(ds$income)`<br>`with(ds, cor(runif(n), rooms))` |
| 8 | **Create a Boxplot - Distribution of income as a factor of number of rooms:**<br><br>1.   Select and execute the code (Listed after the comment line "Step 8") in the script window.<br>2.   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.<br>`boxplot(income ~ as.factor(rooms), data=ds, range=0,`<br>`outline=F, log="y",xlab="# rooms", ylab="Income")`<br><br>3.   Plot the # of rooms as a function of wealth level.<br>`boxplot(rooms ~ wealth, data = ds,main="Room by`<br>`Wealth", Xlab="Category", ylab="# rooms")`<br><br>`# we'll keep the outlier points in this one` |

| Step | Action |
|------|--------|
| 9 | **Exit R:**<br><br>1. Type the following command into the RStudio command window:<br>`q()`<br><br>2. R will ask you if you want to save your workspace. Answer "**no**." |

*End of Lab Exercise*