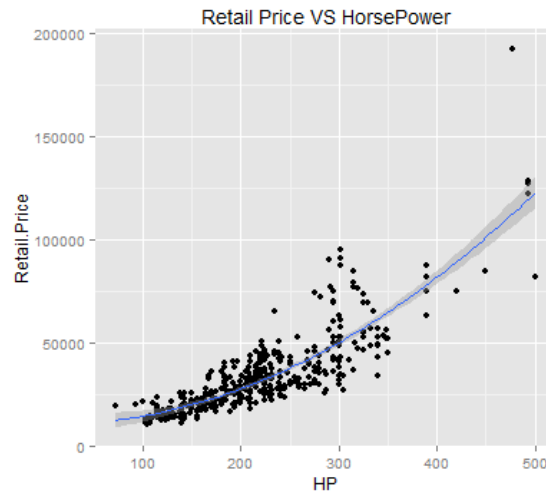


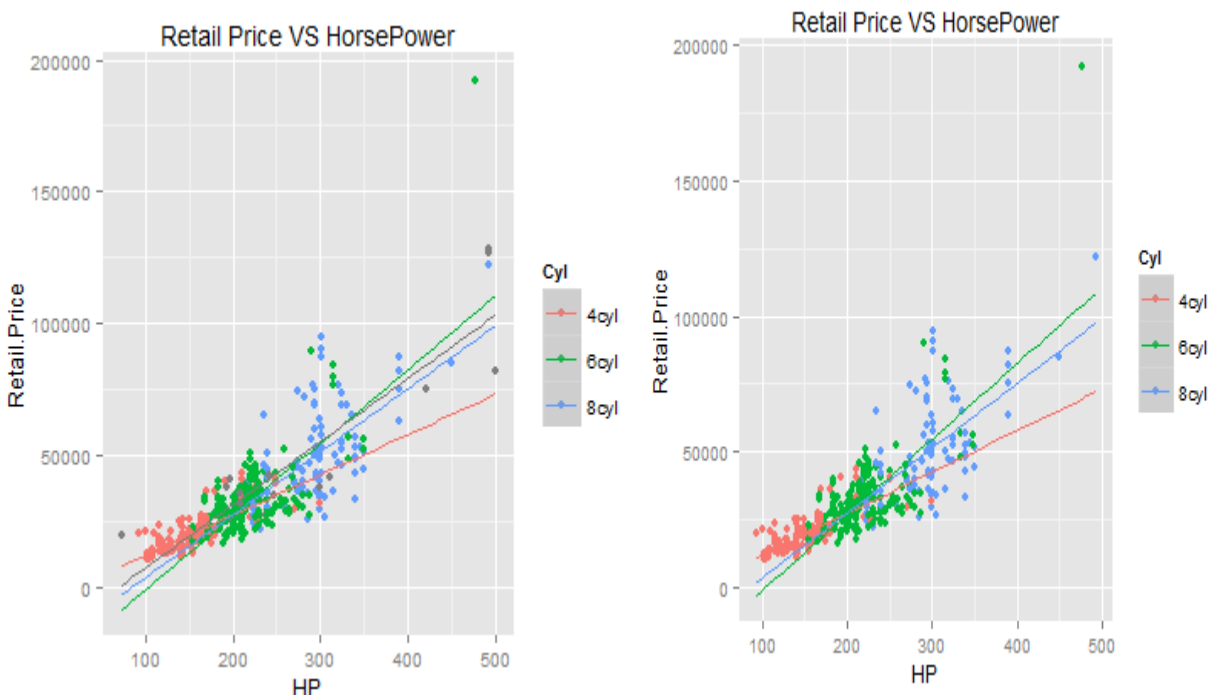
**Lab 4**  
**Steven Lin**

**Exercise 1**

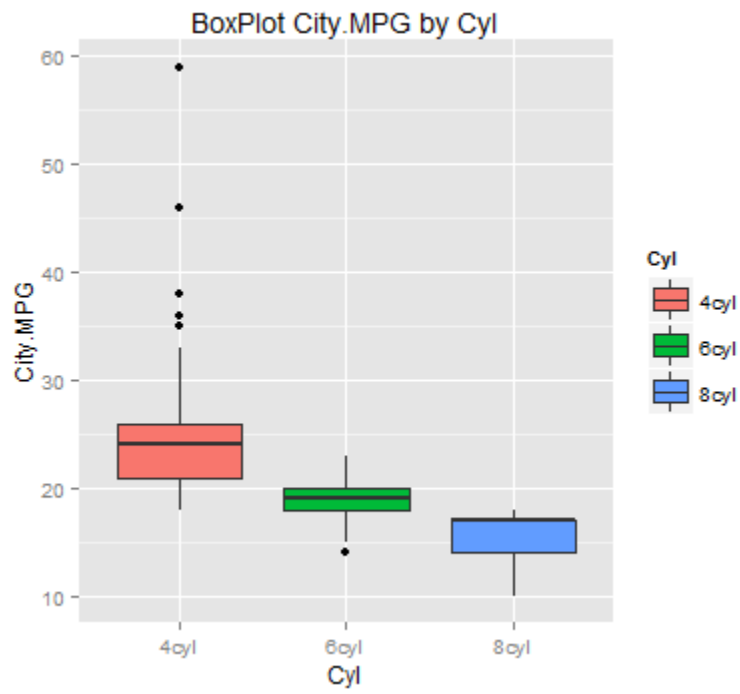
a)



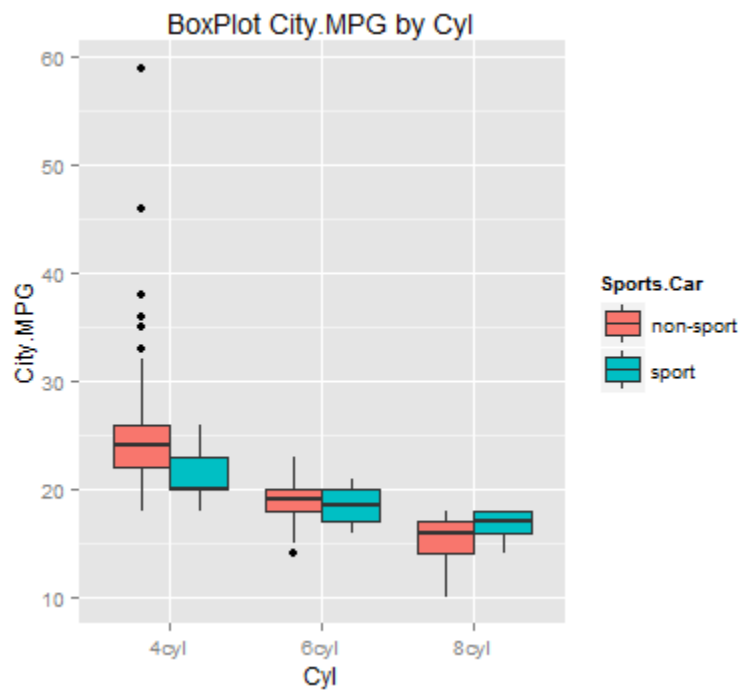
**b)** Note: plot on the left has all the points where point with cylinders other than 4, 6 and 8 are plotted in gray because using factor on cylinders on the entire data converts these points to NA. Plot in the right only plots the subset of data with cylinders 4, 6 and 8. For the next questions I use the subset data.



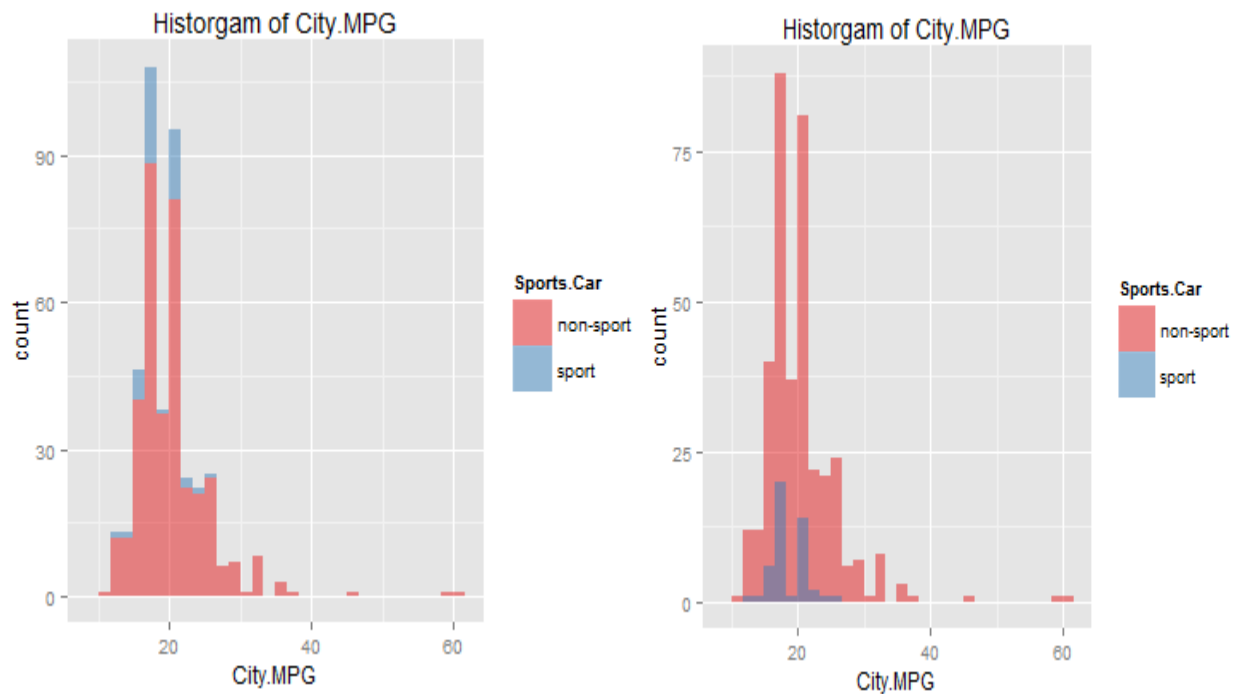
c)



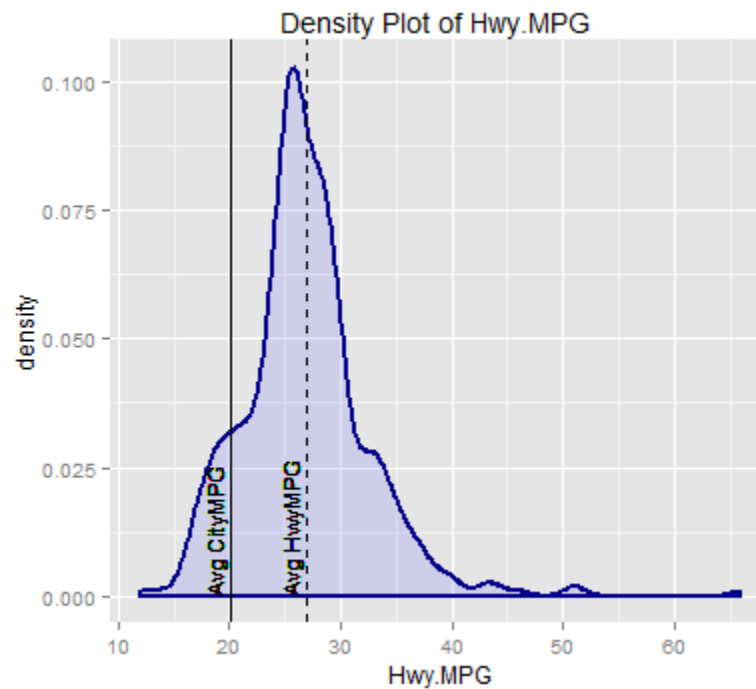
d)



e) Note: not sure what “on top of each other” refers to. So I included “stacked” (left) and “overlaid” (right) histograms. The overlaid histogram makes much more sense to use with transparency since you can clearly compare the distributions of sports vs. non-sports.

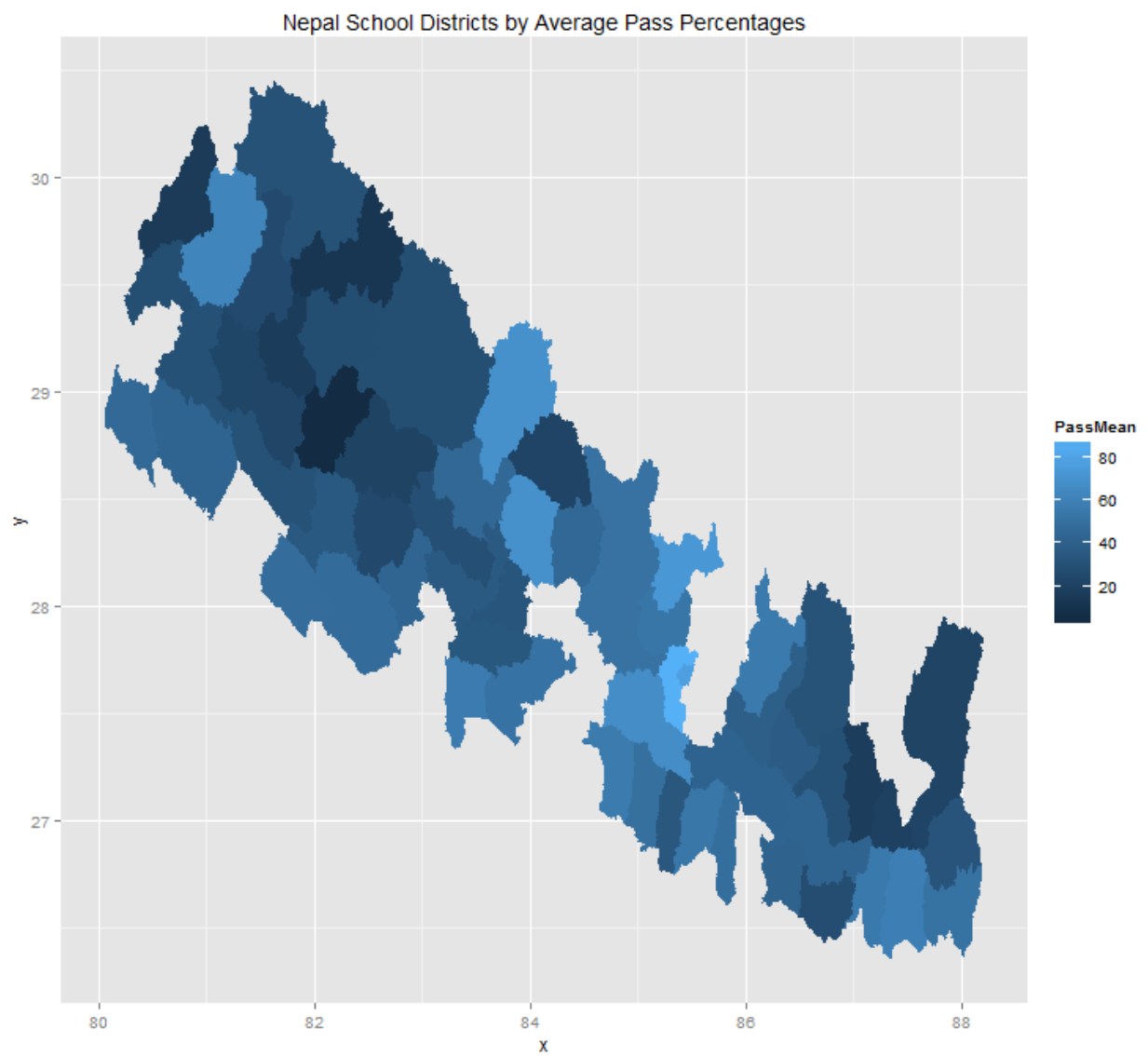


f)

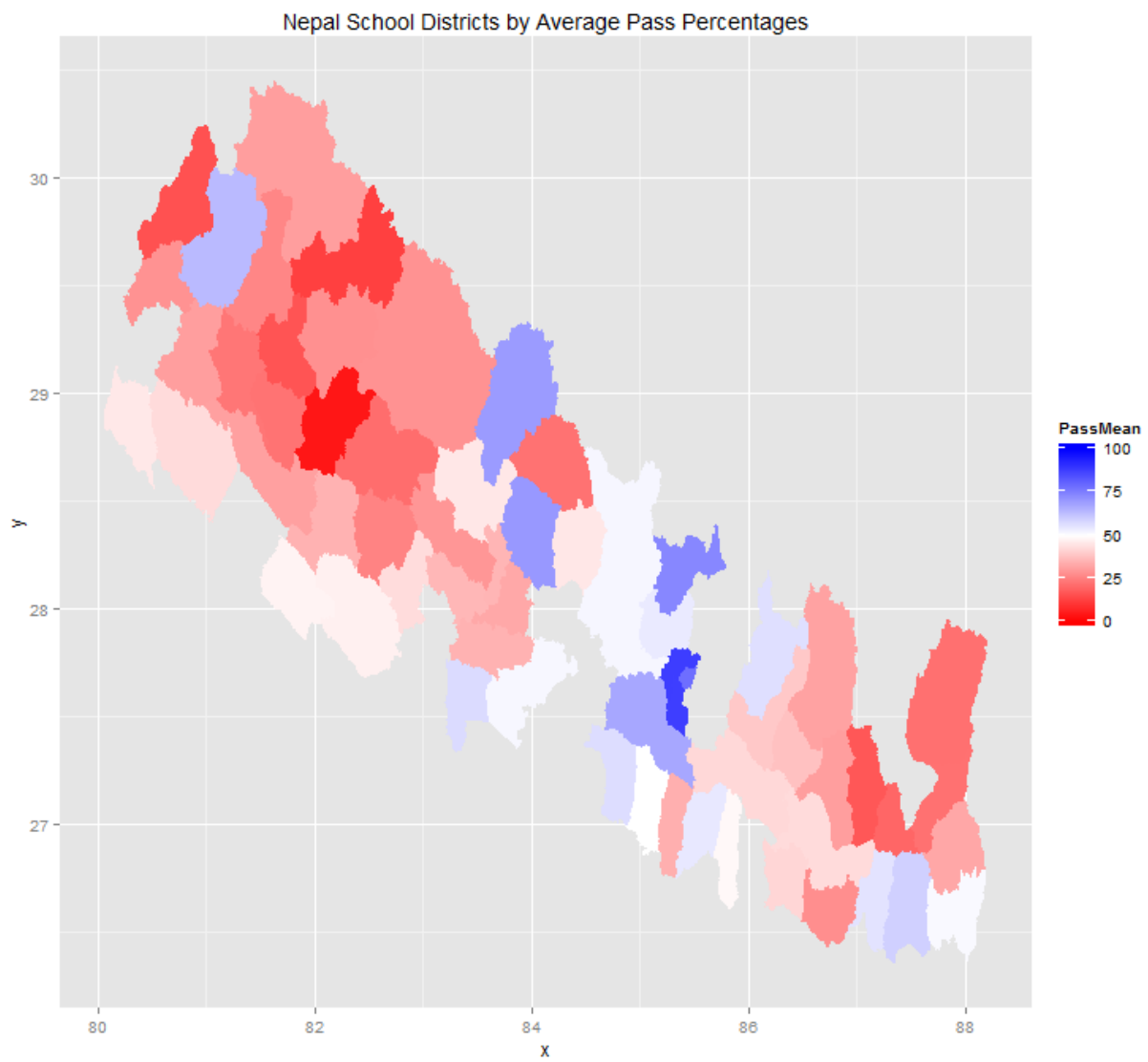


## Exercise 2

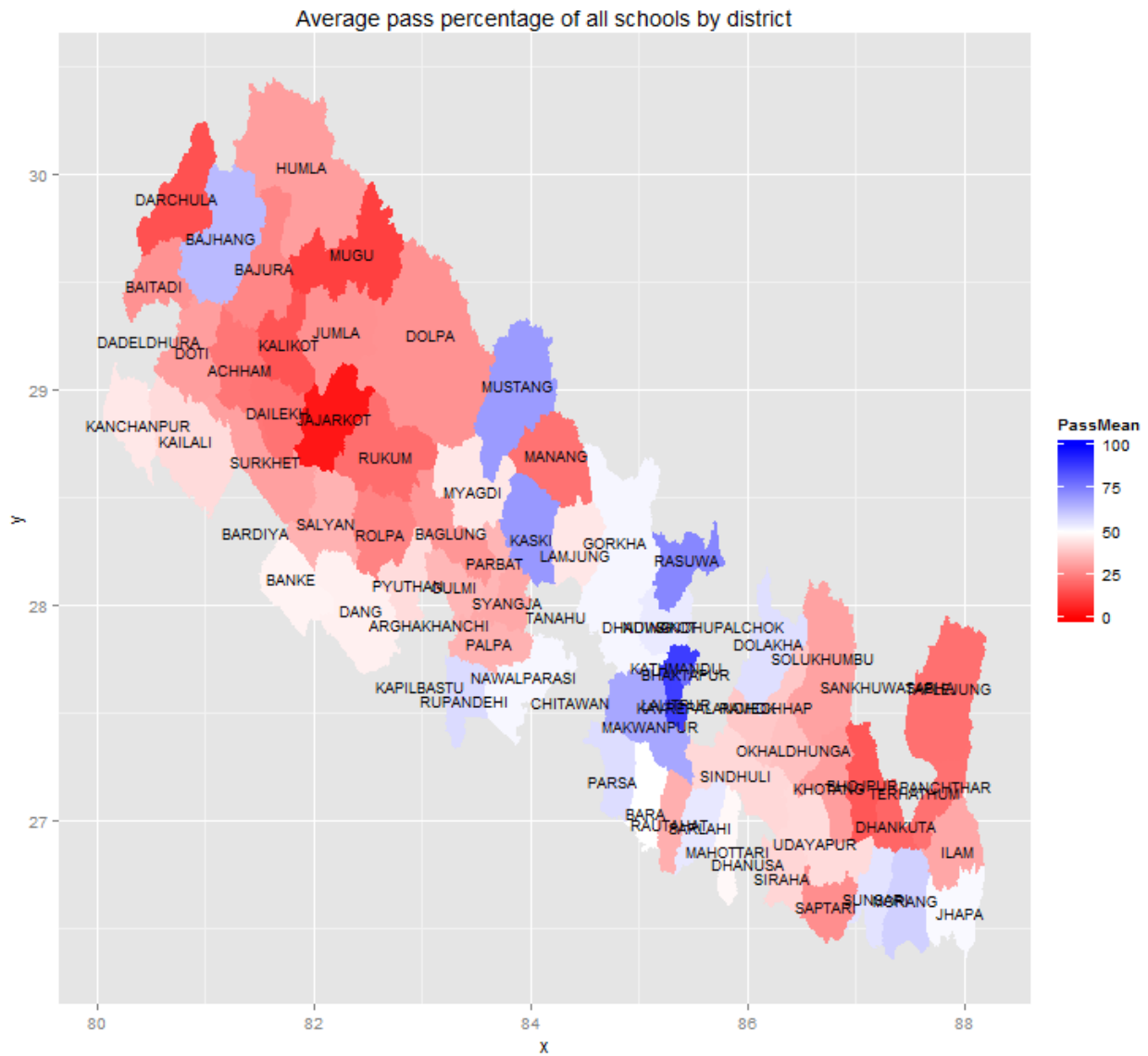
a)



b)



**c)**



## Code

```
1. # Lab session 4 exercise
2.
3. ##### Load data #####
4.
5. # My PC
6. main = "C:/Users/Steven/Documents/Academics/3_Graduate School/2014-2015 ~ NU/"
7.
8. # Aginity
9. #main = "\\nas1/labuser169"
10.
11. course = "MSIA_411_Data_Visualization"
12. datafolder = "/Lab/Data"
13. setwd(file.path(main,course, datafolder))
14.
15. carsdata =read.csv("04cars data.csv",header=TRUE,
16.                   na.strings=c("","*", "NA"))
17.
18.
19. ##### Ex1 #####
20.
21. library(ggplot2)
22.
23. ## Part a ###
24.
25. # Using qplot() create a scatter plot of "HP" against "Retail.Price."
26. # Fit a second degree polynomial through the data, i.e. Retail.Price
27. # = HP2. Make sure that the datapoints and the fitted line are
28. # shown on the same figure.
29.
30. qplot(HP,Retail.Price, data = carsdata, geom =c("point","smooth"),
31.       formula=y~poly(x,2), method="lm",main="Retail Price VS HorsePower")
32.
33. ## Part b ###
34.
35. # Select all the cars with 4, 6 and 8 cylinders. For each cylinder
36. # category regress (using "lm") "HP" on Retail.Price." In total you
37. # need to have three regression lines. Remove the 95% confidence
38. # interval around the lines.
39.
40. table(carsdata$Cyl)
41.
42. dim(carsdata) # 428 rows
43. mydata = carsdata
44.
45. mydata$Cyl = factor(mydata$Cyl,levels=c(4,6,8),labels=c("4cyl","6cyl","8cyl"))
46.
47. # data points with no cyl = 4,6,8
48. index = which(!(carsdata$Cyl %in% c(4,6,8)))
49.
50. carsdata[index,]
51. mydata[index,]
52.
53. dim(mydata) # 428 rows
54. table(mydata$Cyl)
55.
56. # when using factors, cyl other than 4,6 and 8 get NA, so when plotting they
57. # are ignored
58.
```

```

59. sub = carsdata[(carsdata$Cyl %in% c(4,6,8)),]
60. sub$Cyl = factor(sub$Cyl,levels=c(4,6,8),labels=c("4cyl","6cyl","8cyl"))
61.
62. # points that are not 4,6 and 8 are plotted in gray color
63. qplot(HP,Retail.Price, data = mydata,
64.       geom = c("point","smooth"),formula=y~x, level =0, fullrange=T,
65.       color=Cyl,method="lm",main="Retail Price VS HorsePower")
66.
67. # subset data 4,6 and 8
68. qplot(HP,Retail.Price, data = sub,
69.       geom = c("point","smooth"),formula=y~x, level =0, fullrange=T,
70.       color=Cyl,method="lm",main="Retail Price VS HorsePower")
71.
72. ## Part c ####
73.
74. # For the same subset of the data as in part (b) (cars with 4,
75. # 6 and 8 cylinders), use ggplot() to a create boxplot for
76. # "City.MPG" for all three categories.
77. # (The final figure should have 3 boxplots and "City.MPG" on
78. # the y-axis.) Give a different fill color to each boxplot.
79.
80. ggplot(sub, aes(x=Cyl, y=City.MPG)) +
81.   geom_boxplot(aes(fill = Cyl))+
82.   ggtitle("BoxPlot City.MPG by Cyl")
83.
84. ## Part d ####
85.
86. # On the final figure from part (c) add another layer of boxplots
87. # in the following way: for every category of cylinders create
88. # two boxplots, one corresponding to sports cars ("Sports.Car" = 1),
89. # and the other corresponding to non-sports cars ("Sports.Car" = 0).
90. # Your resulting figure should have 6 boxplots.
91. # Use two fill colors; one color for sports cars boxplots
92. # and another for non-sports cars boxplots. You must use ggplot().
93.
94. table(sub$Sports.Car)
95. sub$Sports.Car = factor(sub$Sports.Car,levels=c(0,1), labels=c("non-sport","sport"))
96.
97. ggplot(sub, aes(x=Cyl, y=City.MPG)) +
98.   geom_boxplot(aes(fill = Sports.Car))+
99.   ggtitle("BoxPlot City.MPG by Cyl")
100.
101.
102.   ## Part e ####
103.
104.   # Using ggplot() create a histogram for "City.MPG" conditional
105.   # on whether a car is a sports car or not (so you have to end
106.   # up with 2 histograms). Both histograms must be shown on the same
107.   # figure on top of each other. Use two different colors and
108.   # set the transparency control alpha = 0.5
109.
110.   mydata = carsdata
111.   mydata$Sports.Car = factor(mydata$Sports.Car,levels=c(0,1),
112.                             labels=c("non-sport","sport"))
113.
114.   # default = stacked
115.   ggplot(mydata, aes(x = City.MPG)) +
116.     geom_histogram(aes(fill=Sports.Car ),alpha = 0.5)+
117.     ggtitle("Historgam of City.MPG")+
118.     scale_fill_brewer(palette="Set1")
119.

```



```

120.     # overlaid
121.     ggplot(mydata, aes(x=City.MPG, fill=Sports.Car)) +
122.       geom_histogram(position="identity", alpha=0.5)+
123.       ggtitle("Histogram of City.MPG") +
124.       scale_fill_brewer(palette="Set1")
125.
126.     ## Part f ####
127.
128.     # Using ggplot() create a density plot for "Hwy.MPG" and draw
129.     # a dashed line showing the mean of "Hwy.MPG" and a straight
130.     # line showing the mean of "City.MPG".
131.
132.     meanHwyMPG = mean(mydata$Hwy.MPG,na.rm=T)
133.     meanCityMPG = mean(mydata$City.MPG,na.rm=T)
134.
135.
136.     ggplot(mydata, aes(x=Hwy.MPG))+
137.       geom_density(colour="darkblue", size=1, fill="blue", alpha=0.1) +
138.       geom_vline(aes(xintercept=meanHwyMPG ), linetype="dashed")+
139.       geom_vline(aes(xintercept=meanCityMPG))+
140.       ggtitle("Density Plot of Hwy.MPG ") +
141.       geom_text(data=mydata, mapping=aes(x=c(meanHwyMPG,meanCityMPG),
142.                                             y=0, label=c("Avg HwyMPG", "Avg CityMPG")),
143.
144.               size=4, angle=90, vjust=-0.4, hjust=0)
145.
146.     ##### Ex1 #####
147.
148.     # Consider the "nepal.csv" data set.
149.
150.     # Part a ####
151.     # Go to GADM.org and download the shape file for Nepal's map (
152.     # you will download a zip file that contains "NPL_adm3.shp").
153.     # Re-run the code from lab 4 to set up Nepal's map such that each
154.     # region's color corresponds to the average pass percentage of all
155.     # schools in that district.
156.
157.     library(plyr)
158.     library(rgeos)
159.     library(maptools)
160.     library(sp)
161.     library(rgdal)
162.     library(ggplot2)
163.     library(Rcpp)
164.     #library(gpclib)
165.
166.     # http://gadm.org/download
167.     nepal = read.csv("nepal.csv",header=TRUE,na.strings=c("", "*", "NA"))
168.
169.     str(nepal)
170.
171.     #download shapefile
172.     #go to GADM.org, select the country you are interested in, and download
173.     #the shp file (you will download a zip file
174.     #that contains the shp file you're interested in)
175.
176.     np_dist = readShapeSpatial("NPL_adm/NPL_adm3.shp")
177.     plot(np_dist)
178.

```

```

179.     #This function turns a map into a data frame than can more easily be plotted with
180.     h ggplot2.
181.     np_dist = fortify(np_dist,region="NAME_3")
182.     #When I was using fortify function, I didn't specify the region and later I set
183.     it as "NAME_3". Here is the reason why I can only use "NAME_3".
184.     #In the unzipped folder, you can see such files: NPL_adm1, NPL_adm2, NPL_adm3, NPL_adm4. #"adm1" means administrative level 1 such as "West", "East", "Central". "adm2"
185.     can refer to #"State". As the data I read into R is NPL_adm3, I should specify the region to the column #which represents this level. So, "NAME_3" is for "NPL_adm3".
186.     #In other words, if you load dataset "NPL_adm2", you have to specify the region="NAME_2".
187.
188.     # since each row contains data about 1 school, we want to take the average
189.     # of schools in the same district.
190.     # use ddply to do averaging
191.
192.     # Take the mean of PASS.PERCENT by District
193.     # Note the use of the '.' function to allow District to be used without quoting
194.     distrpassave = ddply(nepal, .(District), summarize, PassMean = mean(PASS.PERCENT
195.     ))
196.     distrpassave = ddply(nepal, ~District, summarize, PassMean = mean(PASS.PERCENT))
197.
198.     np_dist$id = toupper(np_dist$id) #change ids to uppercase
199.
200.     ggplot() +
201.     geom_map(data = distrpassave, aes(map_id = District, fill = PassMean), map = n
202.     p_dist)+
203.     expand_limits(x = np_dist$long, y = np_dist$lat) +
204.     ggtitle("Nepal School Districts by Average Pass Percentages")
205.
206.     # Part b ####
207.     # Change the colors on the map. Make the low values red, the mid
208.     # values white, and the high values blue. Set the midpoint to 50.
209.
210.     ggplot() +
211.     geom_map(data = distrpassave, aes(map_id = District, fill = PassMean), map = n
212.     p_dist) +
213.     expand_limits(x = np_dist$long, y = np_dist$lat) +
214.     scale_fill_gradient2(low="red", mid = "white", midpoint = 50, high = "blue",li
215.     mits=c(0,100)) +
216.     ggtitle("Nepal School Districts by Average Pass Percentages")
217.
218.     # Part c ####
219.     # Add the district names to the map.
220.
221.     distlabels = ddply(np_dist, .(id), summarize, long = mean(long), lat = mean(lat)
222.     )
223.
224.     ggplot() +
225.     geom_map(data = distrpassave, aes(map_id = District, fill = PassMean), map = n
226.     p_dist) +
227.     expand_limits(x = np_dist$long, y = np_dist$lat) +
228.     scale_fill_gradient2(low="red", mid = "white", midpoint = 50, high = "blue",li
229.     mits=c(0,100)) +
230.     ggtitle("Average pass percentage of all schools by district")+
231.     geom_text(data = distlabels,aes(long,lat,label=id),size=3)

```