

Exploratory Data Analysis

1. Principles of Analytic Graphs

Console Terminal Jobs

R 4.1.2 : ~/

```
|-----| How many days of improvement does the median correspond to? | 17%
1: 12
2: 1
3: 4
4: -2
Selection: 2
| Nice work!
|-----| While it's somewhat informative, it's also somewhat cryptic, since the y-axis is | 19%
| claiming to show a change in number of symptom-free days. Wouldn't it be better
| to show a comparison?
...
|-----| Like this? Here's a graphic which shows two boxplots, the one on the left showing | 22%
| the results for a control group that doesn't use an air cleaner alongside the
| previously shown boxplot.
...
|-----| By showing the two boxplots side by side, you can clearly see that using the air | 25%
| cleaner increases the number of symptom-free days for most asthmatic children.
| The plot on the right (using the air cleaner) is generally higher than the one on
| the left (the control group).
...
|-----| what does this graph NOT show you? | 28%
1: Using the air cleaner makes asthmatic children sicker
```

Environment History Connections Tutorial

Files Plots Packages Help Viewer

Figure 2. Percentage Change in Emergency Hospital Admissions Rate for Cardiovascular Diseases per $\mu\text{g}/\text{m}^3$ Increase in Particulate Matter

PM_{2.5}:
1st q: 0.0
Median: 0.0
3rd q: 0.0
IQR: 0.0
PM₁₀:
1st q: 0.0
Median: 0.0
3rd q: 0.0
IQR: 0.0
PM₁₀ - adjusted for PM_{2.5}:
1st q: 0.0
Median: 0.0
3rd q: 0.0
IQR: 0.0
PM_{2.5} - adjusted for PM₁₀:
1st q: 0.0
Median: 0.0
3rd q: 0.0
IQR: 0.0

Percentage Change in Hospital Admissions per
10 $\mu\text{g}/\text{m}^3$ Increase in Particulate Matter

Relative risk estimate 95% posterior intervals

|-----| which of the following is NOT a good principle of graphing? | 86%
1: Having unreadable labels
2: To describe and document evidence
3: To integrate multiple modes of evidence
4: Content is king

Selection: 1

| Your dedication is inspiring!

|-----| which of the following is NOT a good principle of graphing? | 89%
1: Content is king
2: To prove you're always right
3: To demonstrate a causative mechanism underlying a correlation
4: To show two competing hypotheses

Selection: 2

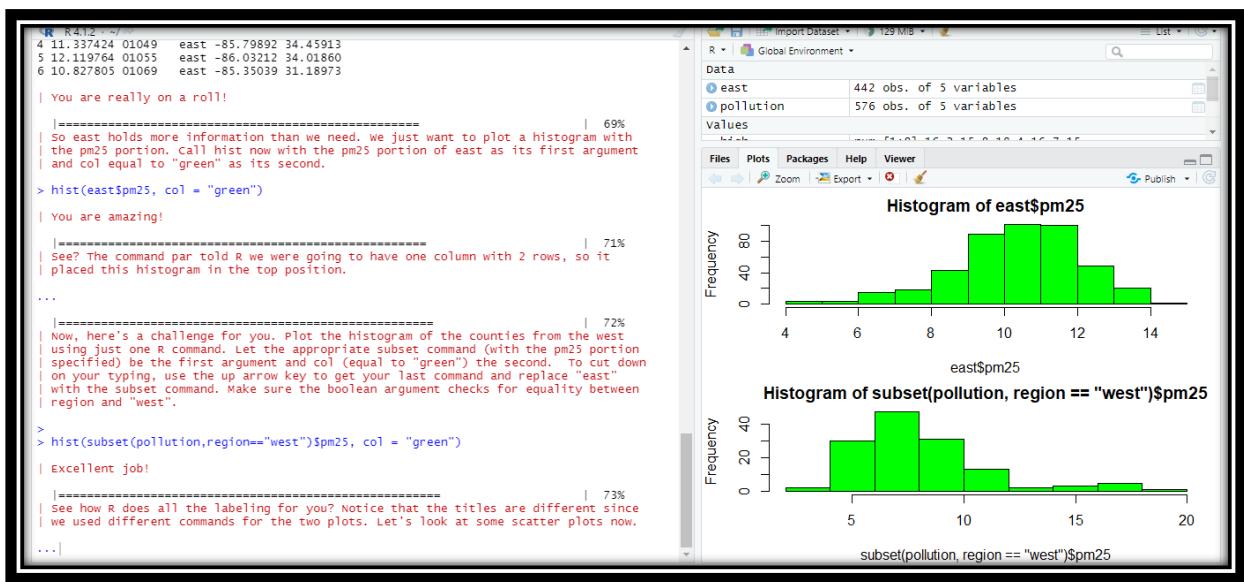
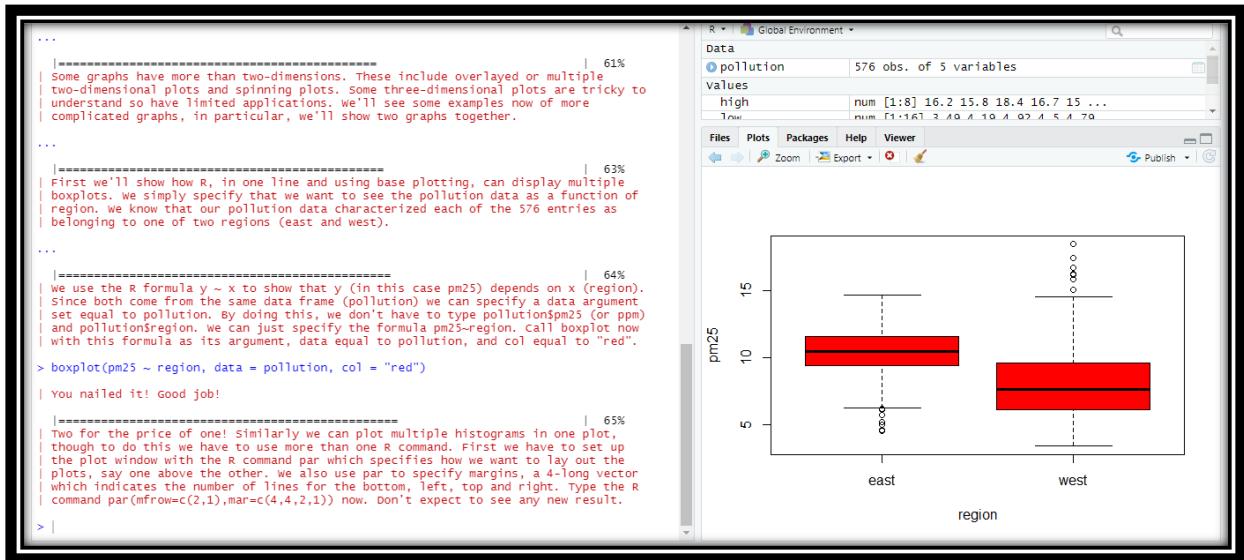
| Perseverance, that's the answer.

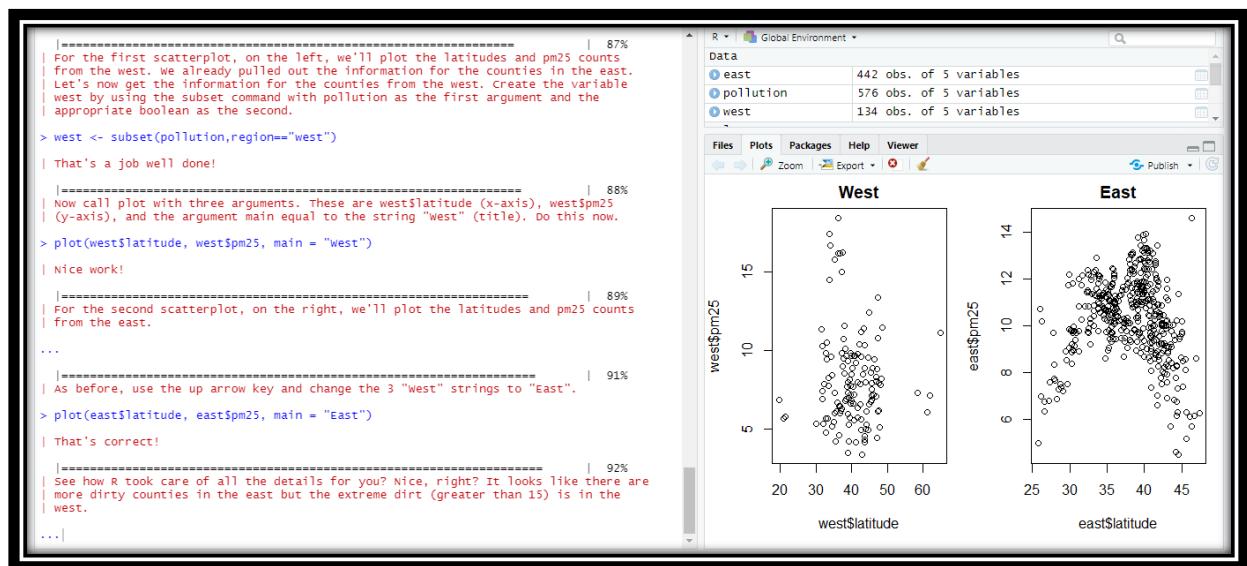
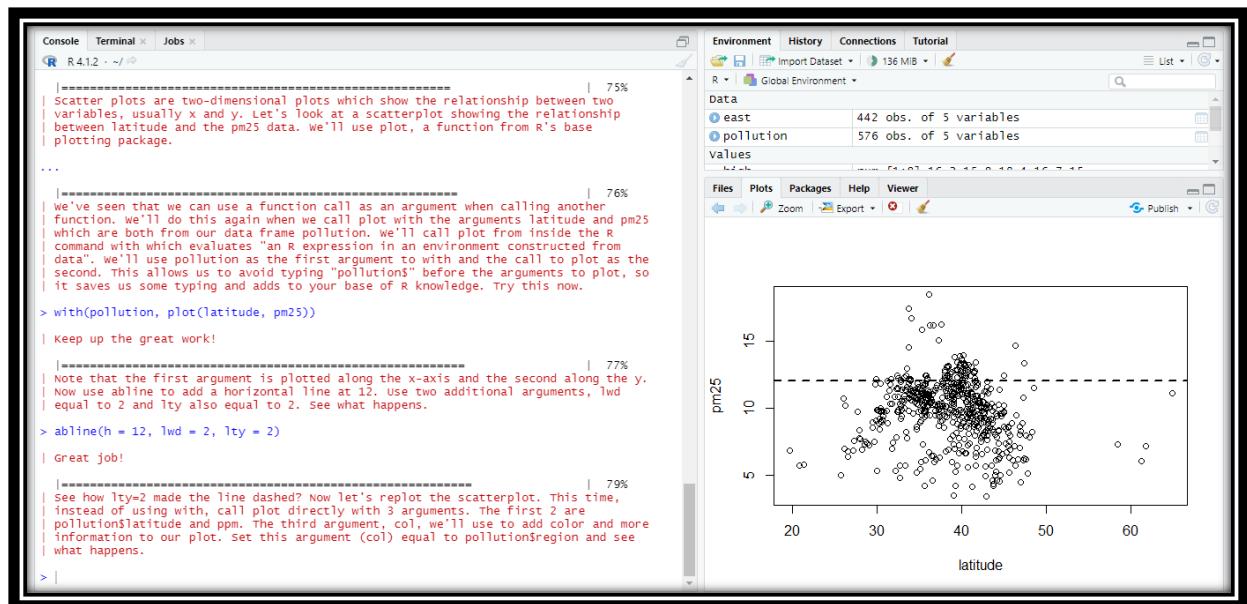
|-----| which of the following is NOT a good principle of graphing? | 92%
1: To integrate different types of evidence
2: To show that some fonts are better than others
3: Content is king
4: To show good labels and scales

Selection: 2

| You are quite good my friend!

2. Exploratory Graphs





3. Graphics Devices in R

Console Terminal Jobs

R 4.1.2 - ~/

```

|-----| 26%
| As an example, run the R command with with 2 arguments. The first is a dataset,
| faithful which comes with R, and the second is a call to the base plotting function
| plot. Your call to plot should have two arguments, eruptions and waiting. Try this
| now to see what happens.

> with(faithful, plot(eruptions, waiting))
| You are doing so well!
|-----| 29%
| See how R created a scatterplot on the screen for you? This shows that relationship
| between eruptions of the geyser old Faithful and waiting time. Now use the R
| function title with the argument main set equal to the string "Old Faithful Geyser
| data". This is an annotation to the plot.

> title(main = "old Faithful Geyser data")
| You are amazing!
|-----| 32%
| Simple, right? Now run the command dev.cur(). This will show you the current
| plotting device, the screen.

> dev.cur()
RStudioGD
2

| Your dedication is inspiring!
|-----| 35%
| The second way to create a plot is to send it to a file device. Depending on the
| type of plot you're making, you explicitly launch a graphics device, e.g., a pdf
| file. Type the command pdf(file="myplot.pdf") to launch the file device. This will
| create the pdf file myplot.pdf in your working directory.

>

```

Environment History Connections Tutorial

R west 134 obs. or 5 variables

Values

high	num [1:8]	16.2 15.8 18.4 16.7 15 ...
low	num [1:16]	3.49 4.19 4.92 4.5 4.79 ...
path_to_course	'C:/Users/bania/Documents/R/win-library/4.1/swirl...'	

Files Plots Packages Help Viewer

Old Faithful Geyser data

The scatterplot shows the relationship between 'eruptions' (x-axis, ranging from 1.5 to 5.0) and 'waiting' (y-axis, ranging from 50 to 90). The data points form a clear positive correlation, with most points clustered between 1.5 and 4.5 on the x-axis and 50 and 80 on the y-axis.

Console Terminal Jobs

R 4.1.2 - ~/

```

|-----| 100%
| would you like to receive credit for completing this course on Coursera.org?
1: No
2: Yes

selection: 1
| You got it right!
| You've reached the end of this lesson! Returning to the main menu...
| Please choose a course, or type 0 to exit swirl.

1: Exploratory Data Analysis
2: Take me to the swirl course repository!

selection: 1
| Please choose a lesson, or type 0 to return to course menu.

1: Principles of Analytic Graphs 2: Exploratory Graphs
3: Graphics Devices in R 4: Plotting Systems
5: Base Plotting System 6: Lattice Plotting System
7: Working with Colors 8: GGPlot2 Part1
9: GGPlot2 Part2 10: GGPlot2 Extras
11: Hierarchical Clustering 12: K Means Clustering
13: Dimension Reduction 14: Clustering Example
15: CaseStudy

selection: 3
|-----| 0%
| Graphics_Devices_in_R. (Slides for this and other Data Science courses may be found
| at https://github.com/datasciencespecialization/courses/. If you care to use
| them, they must be downloaded as a zip file and viewed locally. This lesson
| corresponds to 04_ExploratoryAnalysis/Graphics_Devices_in_R.)
```

Environment History Connections Tutorial

R west 134 obs. or 5 variables

Values

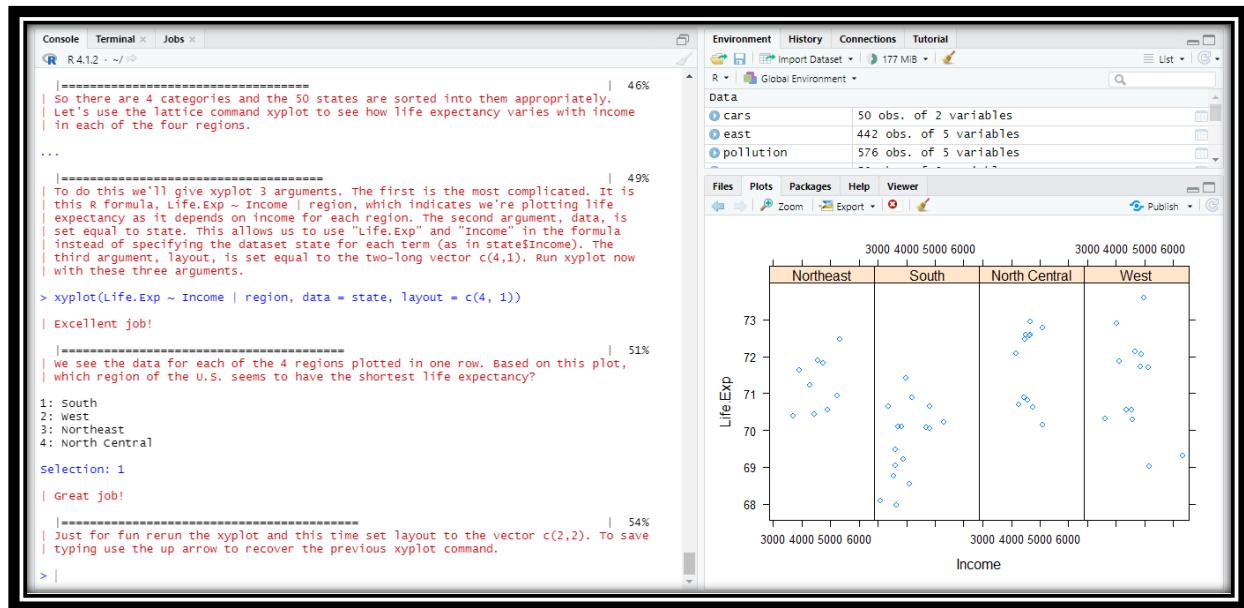
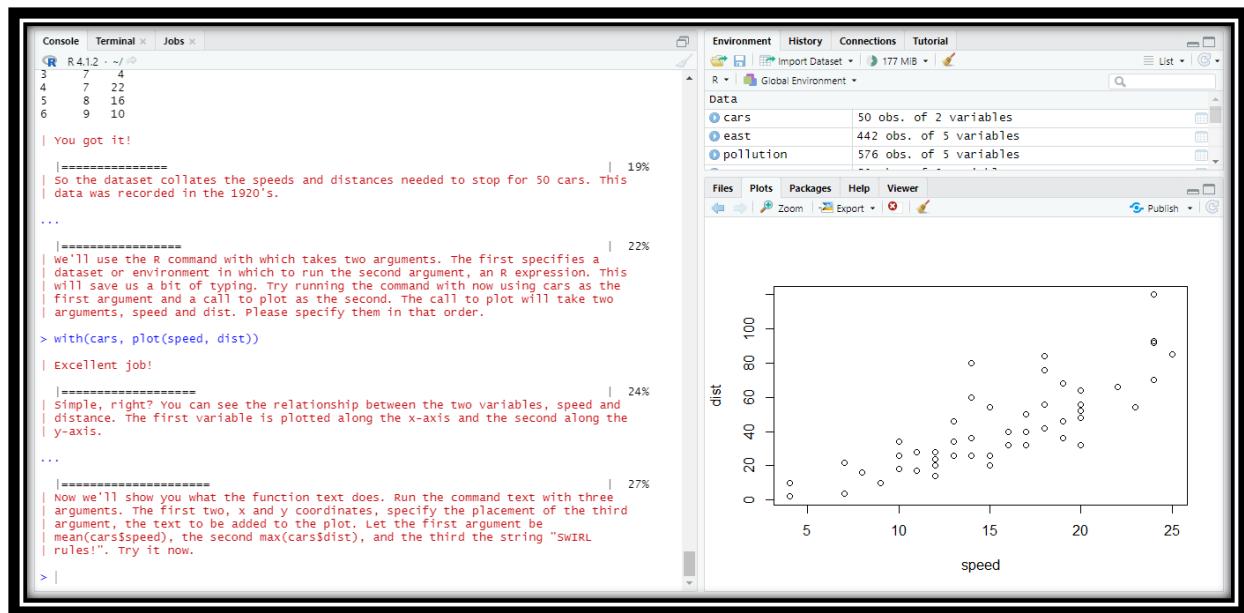
high	num [1:8]	16.2 15.8 18.4 16.7 15 ...
low	num [1:16]	3.49 4.19 4.92 4.5 4.79 ...
path_to_course	'C:/Users/bania/Documents/R/win-library/4.1/swirl...'	

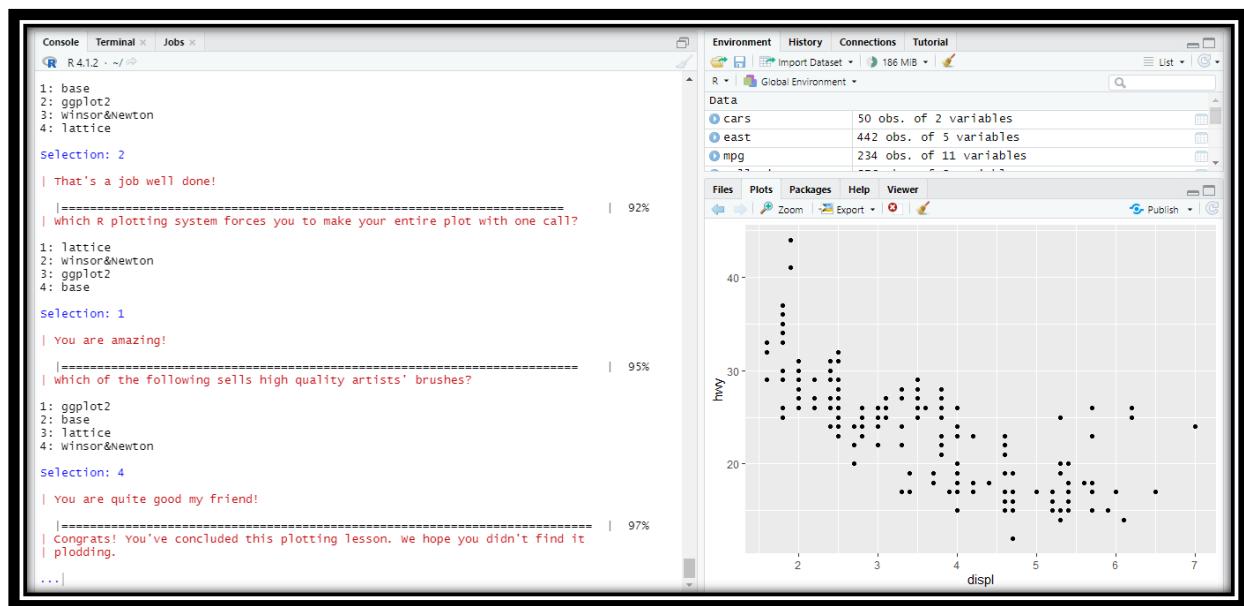
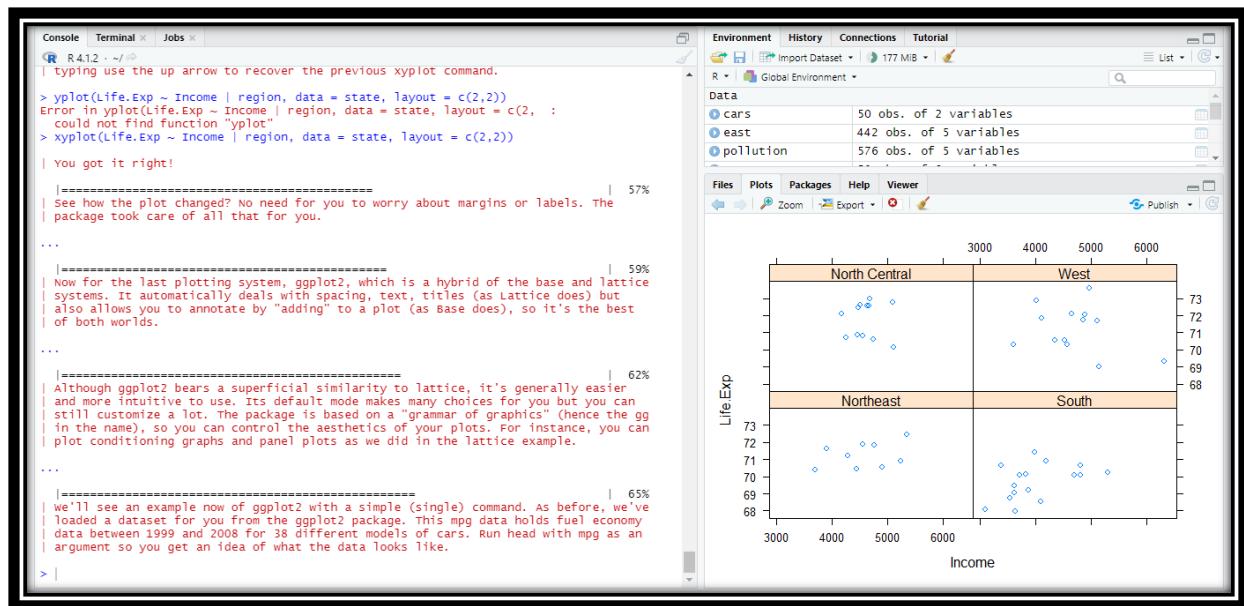
Files Plots Packages Help Viewer

Old Faithful Geyser data

The scatterplot shows the relationship between 'eruptions' (x-axis, ranging from 1.5 to 5.0) and 'waiting' (y-axis, ranging from 50 to 90). The data points form a clear positive correlation, with most points clustered between 1.5 and 4.5 on the x-axis and 50 and 80 on the y-axis.

4. Plotting Systems





5. Base Plotting Systems

```

|-----| 17%
| see the dataset contains 6 columns of data. Run the command range with two
| arguments. The first is the ozone column of airquality, specified by
| airquality$ozone, and the second is the boolean na.rm set equal to TRUE. If you
| don't specify this second argument, you won't get a meaningful result.
> range(airquality$ozone,na.rm=TRUE)
[1] 1 168
| You got it right!
|-----| 18%
| So the measurements range from 1 to 168. First we'll do a simple histogram of this
| ozone column to show the distribution of measurements. Use the R command hist with
| the argument airquality$ozone.
> hist(airquality$ozone)
| Your dedication is inspiring!
|-----| 20%
| Simple, right? R put a title on the histogram and labeled both axes for you. What is
| the most frequent count?
1: Between 60 and 75
2: Over 150
3: Over 100
4: Under 25
selection: 4
| You are doing so well!
|-----| 21%
| Next we'll do a boxplot. First, though, run the R command table with the argument
| airquality$Month.
>

```

A histogram titled "Histogram of airquality\$Ozone". The x-axis is labeled "airquality\$Ozone" and ranges from 0 to 150 with major ticks at 0, 50, 100, and 150. The y-axis is labeled "Frequency" and ranges from 0 to 30 with major ticks at 0, 10, 20, and 30. The distribution is skewed right, with the highest frequency in the 0-25 bin (~35), followed by the 25-50 bin (~28).

```

| You're the best!
|-----| 24%
| Note that boxplot, unlike hist, did NOT specify a title and axis labels for you
| automatically.
...
|-----| 26%
| Let's call boxplot again to specify labels. (Use the up arrow to recover the
| previous command and save yourself some typing.) We'll add more arguments to the
| call to specify labels for the 2 axes. Set xlab equal to "Month" and ylab equal to
| "ozone (ppb)". Specify col.axis equal to "blue" and col.lab equal to "red". Try this
| now.
> boxplot(Ozone~Month, airquality, xlab="Month", ylab="ozone (ppb)", col.axis="blue", col.lab="red")
| You are quite good my friend!
|-----| 27%
| Nice colors, but still no title. Let's add one with the R command title. Use the
| argument main set equal to the string "ozone and wind in New York City".
> title(main="ozone and wind in New York City")
| That's the answer I was looking for.
|-----| 29%
| Now we'll show you how to plot a simple two-dimensional scatterplot using the R
| function plot. We'll show the relationship between wind (x-axis) and ozone (y-axis).
| We'll use the function plot with those two arguments (wind and ozone, in that
| order). To save some typing, though, we'll call the R command with using 2
| arguments. The first argument of with will be airquality, the dataset containing
| wind and ozone; the second argument will be the call to plot. Doing this allows us
| to avoid using the longer notation, e.g., airquality$wind. Try this now.
>

```

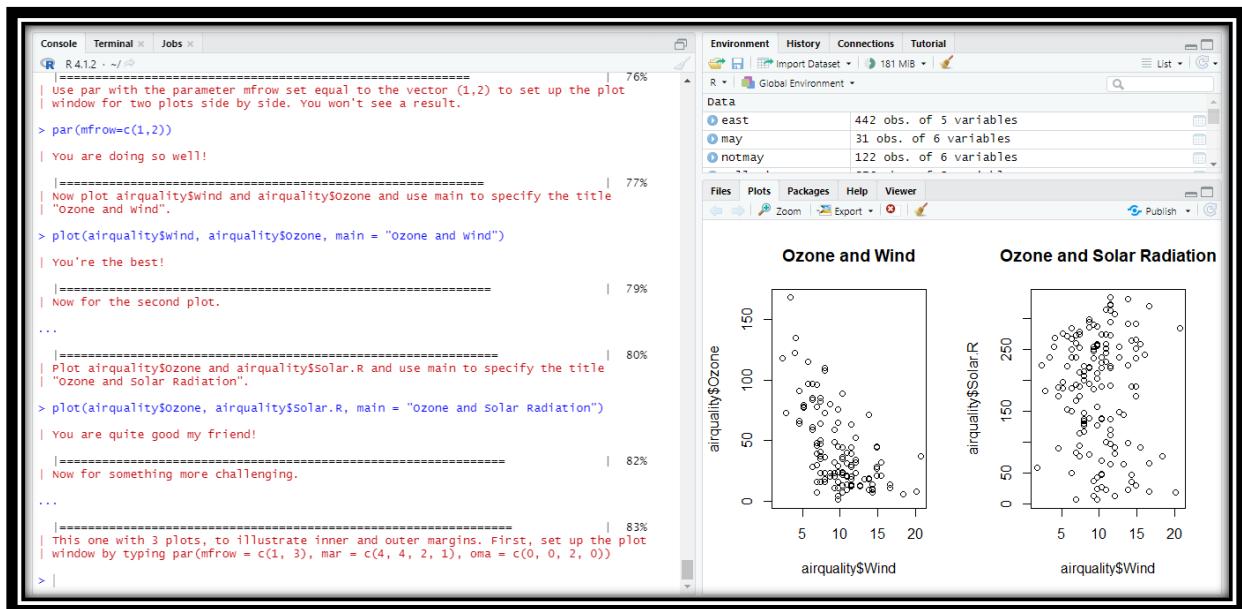
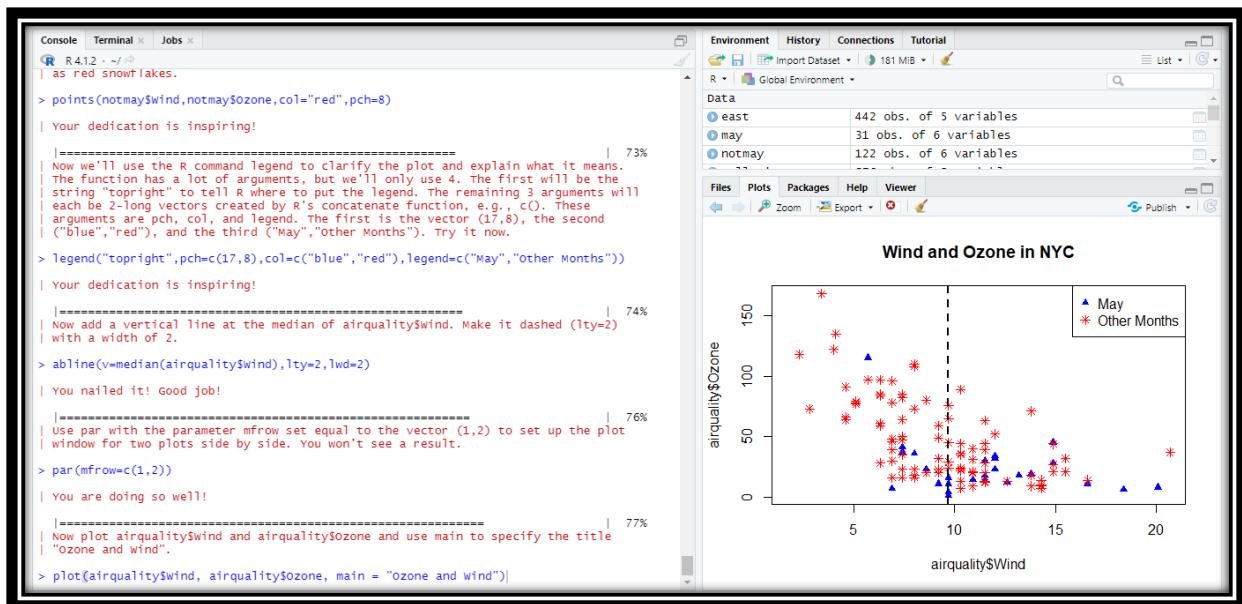
A boxplot titled "Ozone and Wind in New York City". The x-axis is labeled "Month" and has categories 5, 6, 7, 8, and 9. The y-axis is labeled "Ozone (ppb)" and ranges from 0 to 150 with major ticks every 50 units. The plot shows five boxplots corresponding to the months. The median ozone level increases from approximately 25 ppb in May to about 55 ppb in June, then rises sharply to around 80 ppb in July and August, before dropping back to about 30 ppb in September.

```

Console Terminal Jobs
R 4.1.2 ~/ ...
| Run the command par("lty") to see the default line type.
> par("lty")
[1] "solid"
| You are amazing!
|-----| 50%
| So the default line type is solid, but it can be dashed, dotted, etc. Once again,
| R's ?par documentation will tell you what other line types are available. The line
| width is a positive integer; the default value is 1.
...
|-----| 52%
| We've seen a lot of examples of col, the plotting color, specified as a number,
| string, or hex code; the colors() function gives you a vector of colors by name.
...
|-----| 53%
| What do you think the graphical parameters xlab and ylab control respectively?
1: labels for the x- and y- axes
2: labels for the y- and x- axes
selection: 1
| You're the best!
|-----| 55%
| The par() function is used to specify global graphics parameters that affect all
| plots in an R session. (Use dev.off or plot.new to reset to the defaults.) These
| parameters can be overridden when specified as arguments to specific plotting
| functions. These include las (the orientation of the axis labels on the plot),
| bg (background color), mar (margin size), oma (outer margin size), mfrow and mfcol
| (number of plots per row, column).
...

```

A scatterplot titled "Ozone and Wind in New York City". The x-axis is labeled "Wind" and ranges from 0 to 20 with major ticks every 5 units. The y-axis is labeled "Ozone" and ranges from 0 to 150 with major ticks every 50 units. The data points show a clear negative correlation, where ozone levels decrease as wind speeds increase. Most data points are clustered between wind speeds of 5 and 15, and ozone levels between 0 and 100 ppb.



R 4.1.2 - ~/

```

Console Terminal Jobs
...
|-----| 95%
| Since this is the main title, we specify it with the R command mtext. call mtext
| with the string "ozone and weather in New York City" and the argument outer set
| equal to TRUE.
> mtext("Ozone and Weather in New York City", outer = TRUE)
| Great job!
|-----| 97%
| Voila! Beautiful, right?
...
|-----| 98%
| Congrats! You've weathered this lesson nicely and passed out of the NoZone.
...
|-----| 100%
| would you like to receive credit for completing this course on Coursera.org?
1: Yes
2: No
Selection: 2
| Keep up the great work!
| You've reached the end of this lesson! Returning to the main menu...
| Please choose a course, or type 0 to exit swirl.
1: Exploratory Data Analysis
2: Take me to the swirl course repository!
Selection: |

```

Environment History Connections Tutorial

R - Global Environment

Data

- east 442 obs. of 5 variables
- may 31 obs. of 6 variables
- notmay 122 obs. of 6 variables

Files Plots Packages Help Viewer

Ozone and Weather in New York City

Ozone and Wind: airquality\$Ozone vs airquality\$Wind

Ozone and Solar Radiation: airquality\$Ozone vs airquality\$Solar.R

Ozone and Temperature: airquality\$Ozone vs airquality\$Temp

6. Lattice Plotting System

R 4.1.2 - ~/

```

Console Terminal Jobs
...
|-----| 22%
| That's a job well done!
|-----| 24%
| Red snowflakes are cool, right? Now that you've seen the basic xyplot() and some of
| its arguments, you might want to experiment more by yourself when you're done with
| this lesson. If you do, you'll see what other arguments and options are available. (If you can't
| wait to experiment, recall that swirl has play() and nxt() functions. At a command
| prompt, typing play() allows you to leave swirl temporarily so you can try different
| R commands at the console. Typing nxt() when you're done playing brings you back to
| swirl and you can resume your lesson.)
...
|-----| 25%
| Now you'll see how easy it is to generate a multipanel plot using a single lattice
| command.
...
|-----| 27%
| Run xyplot with the formula ozone~wind | as.factor(month) as the first argument and
| the second argument data set equal to airquality (use the up arrow to save typing).
| So far, not much is different, right? Add a third argument, layout, set equal to
| c(5,1).
> xyplot(ozone ~ wind | as.factor(month), data = airquality, layout=c(5,1))
| You are really on a roll!
|-----| 27%
| Note that the default color and plotting character are back. What did the
| as.factor(month) do?
1: Huh?
2: Randomly divided the data into 5 panels
3: Displayed and labeled each subplot with the month's integer
4: Displayed the data by individual months
Selection: |

```

Environment History Connections Tutorial

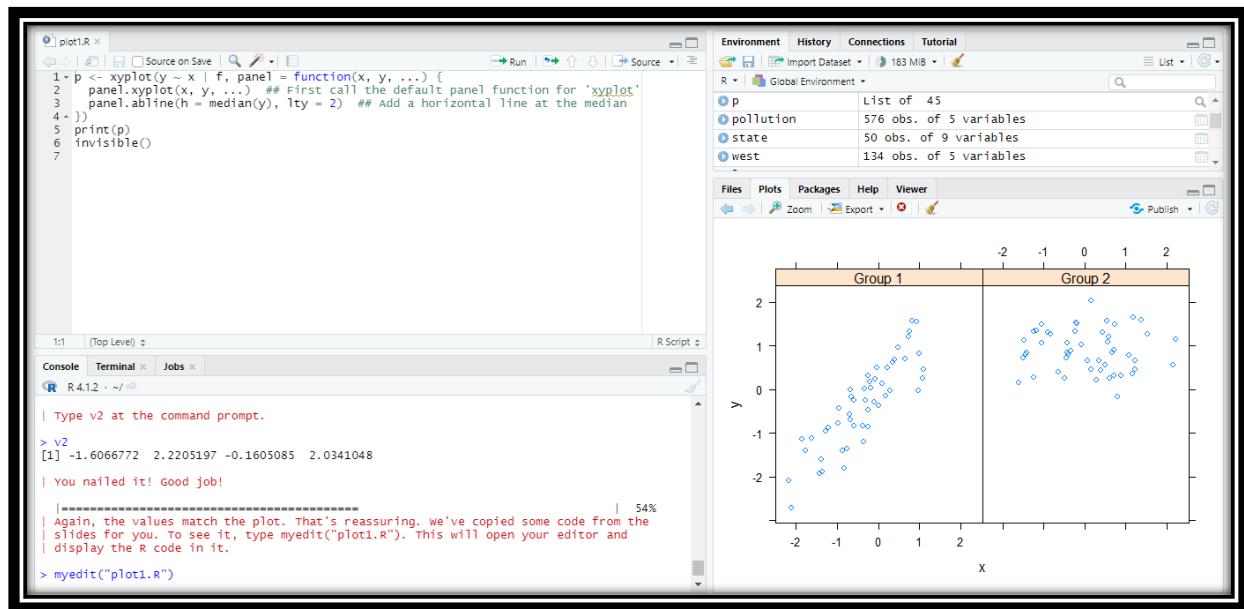
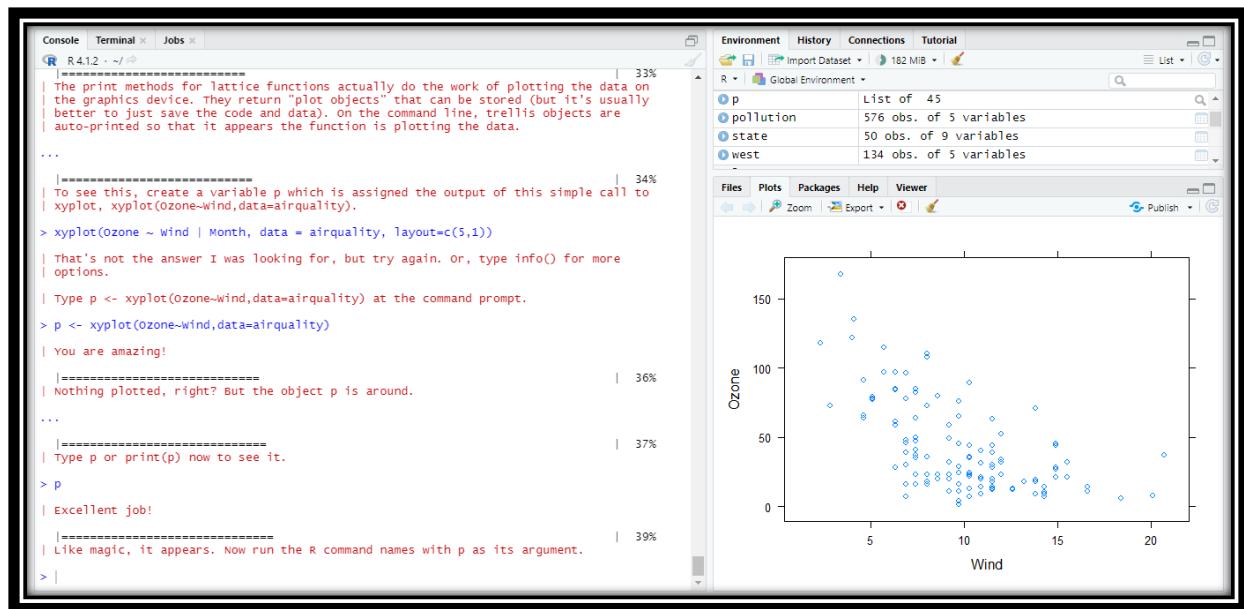
R - Global Environment

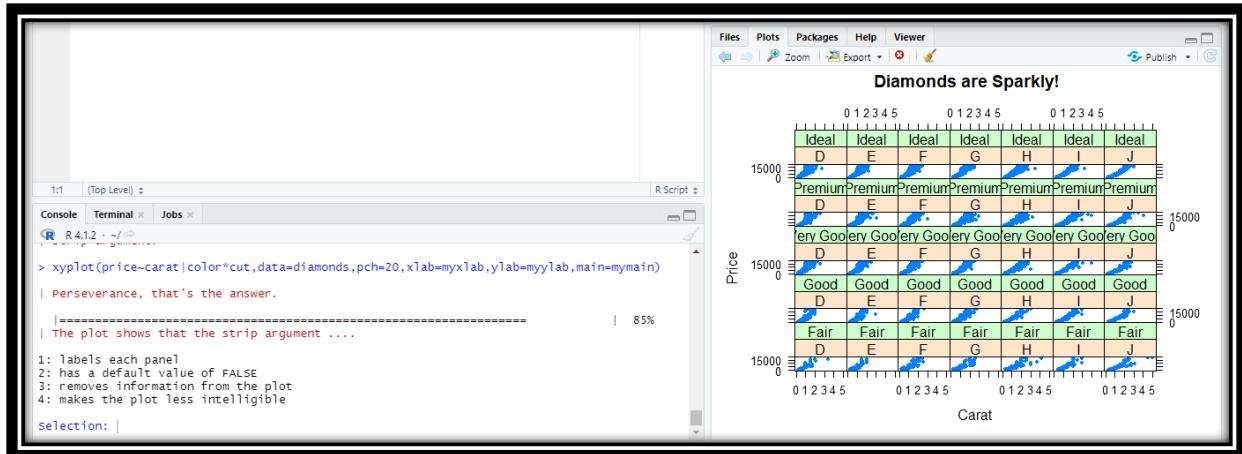
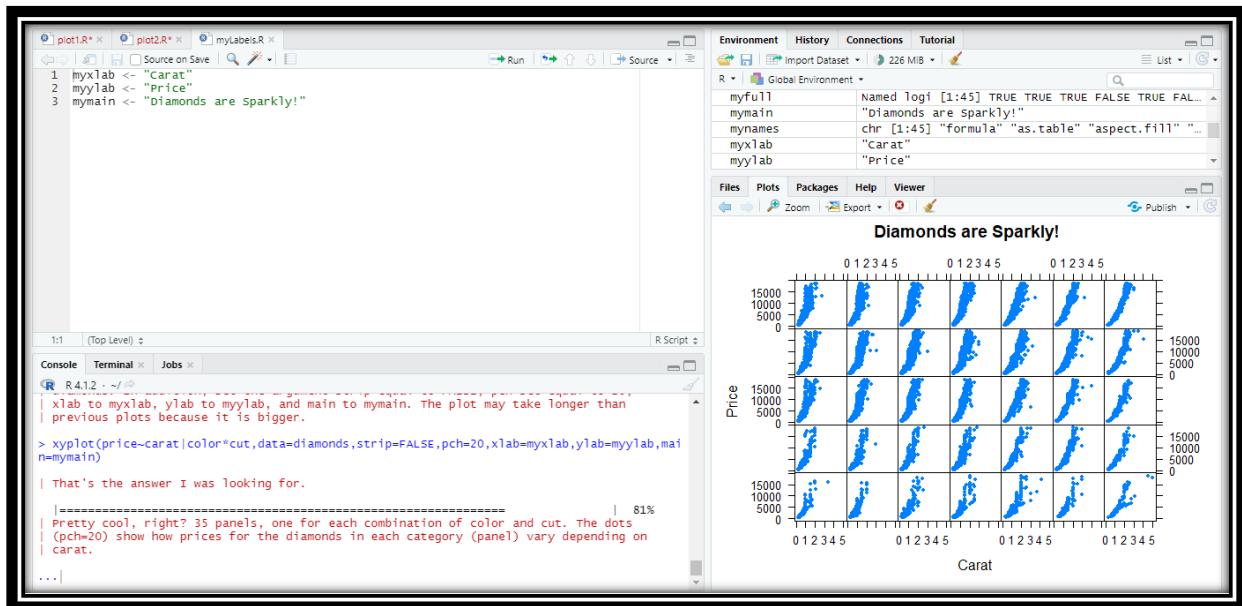
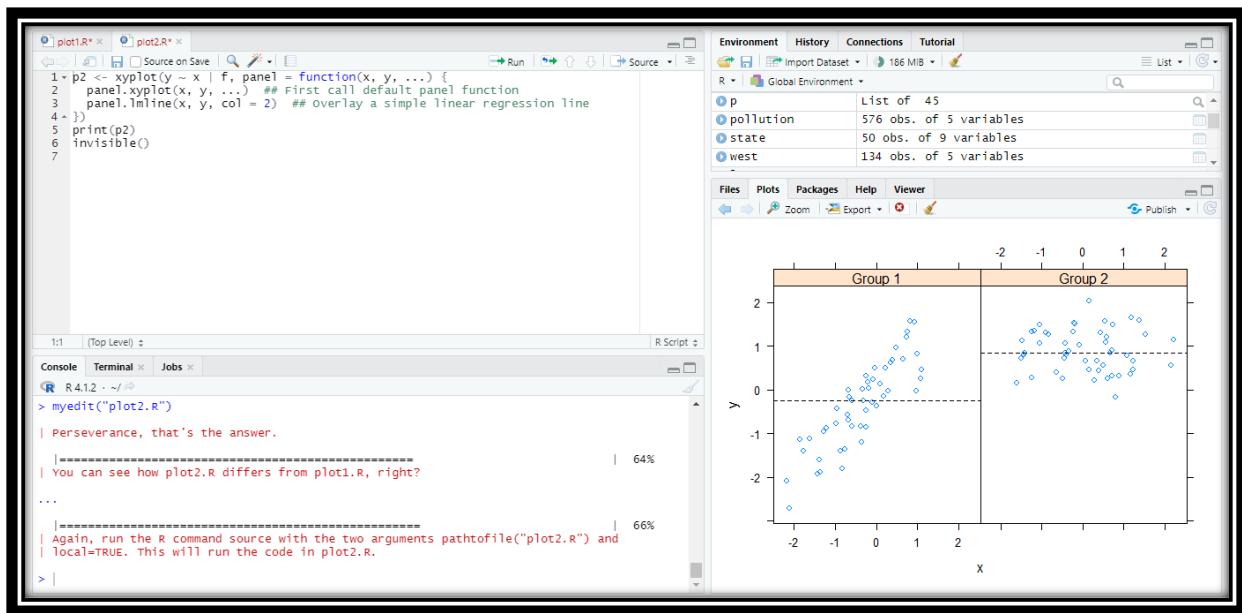
ry2	num [1:2] -0.161 2.034
v1	num [1:4] -2.19 1.1 -2.72 1.57
v2	num [1:4] -1.607 2.221 -0.161 2.034
x	num [1:100] 0.0187 -0.1843 -1.3713 -0.5992 0.29...
y	num [1:100] -0.3622 0.0254 -1.8913 -0.2434 -0.0...

Files Plots Packages Help Viewer

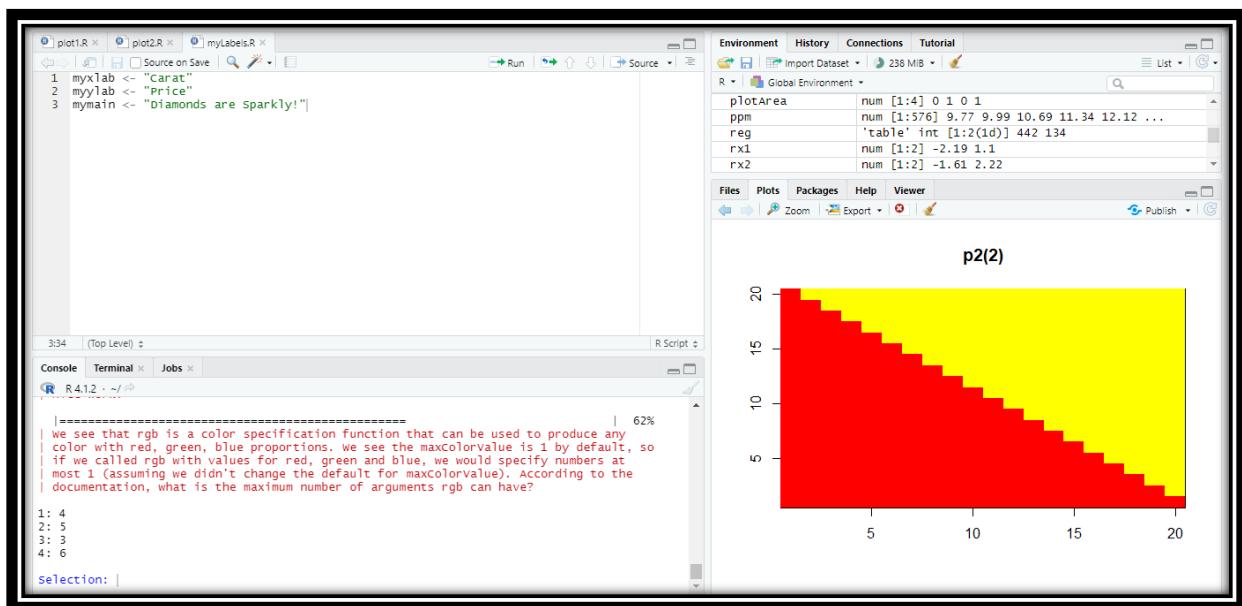
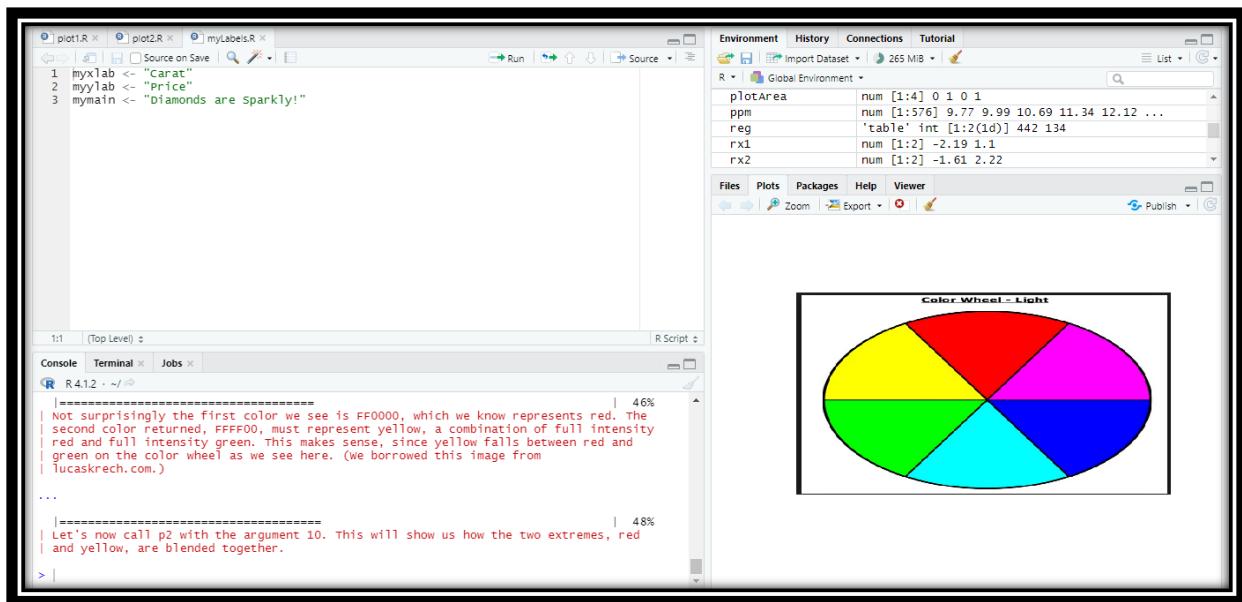
Ozone

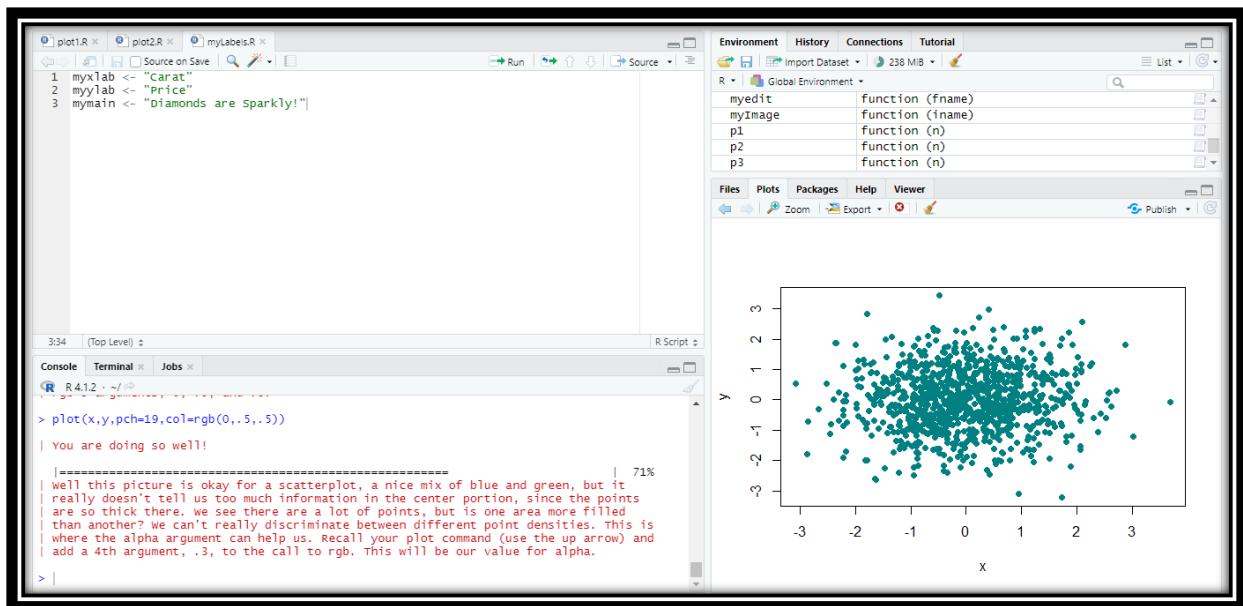
Wind



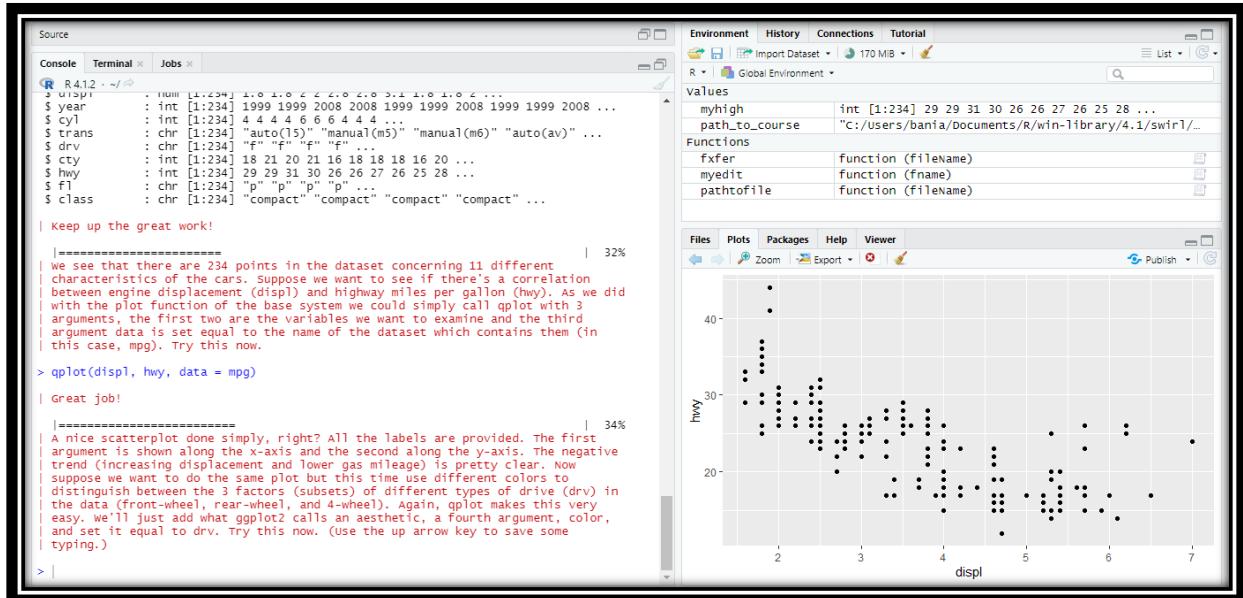


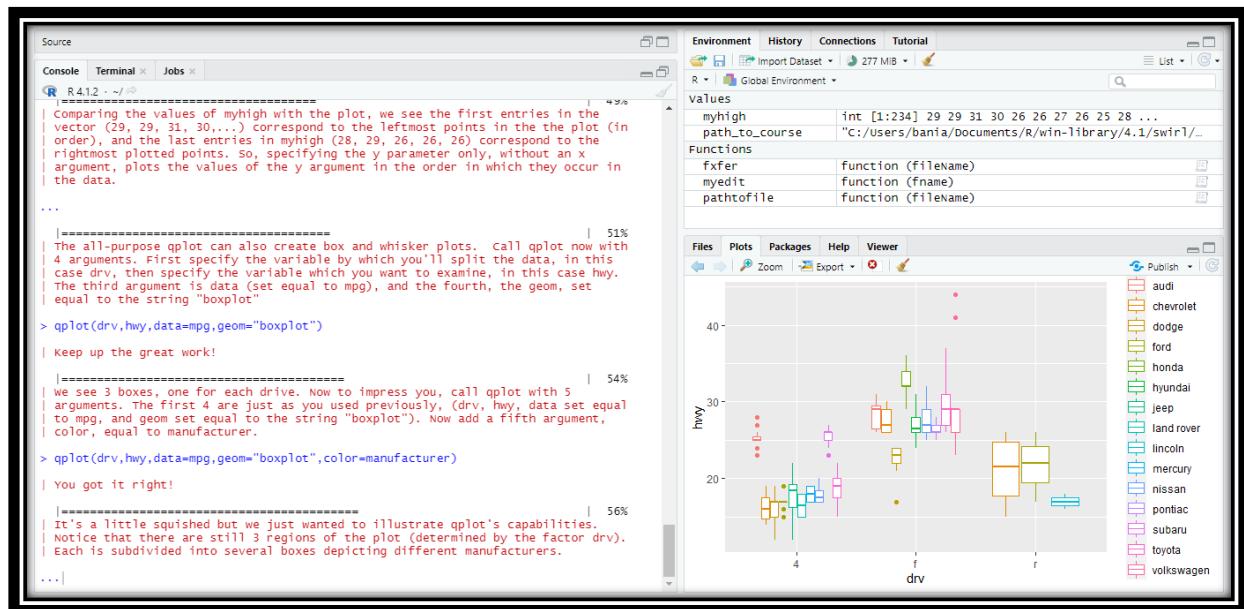
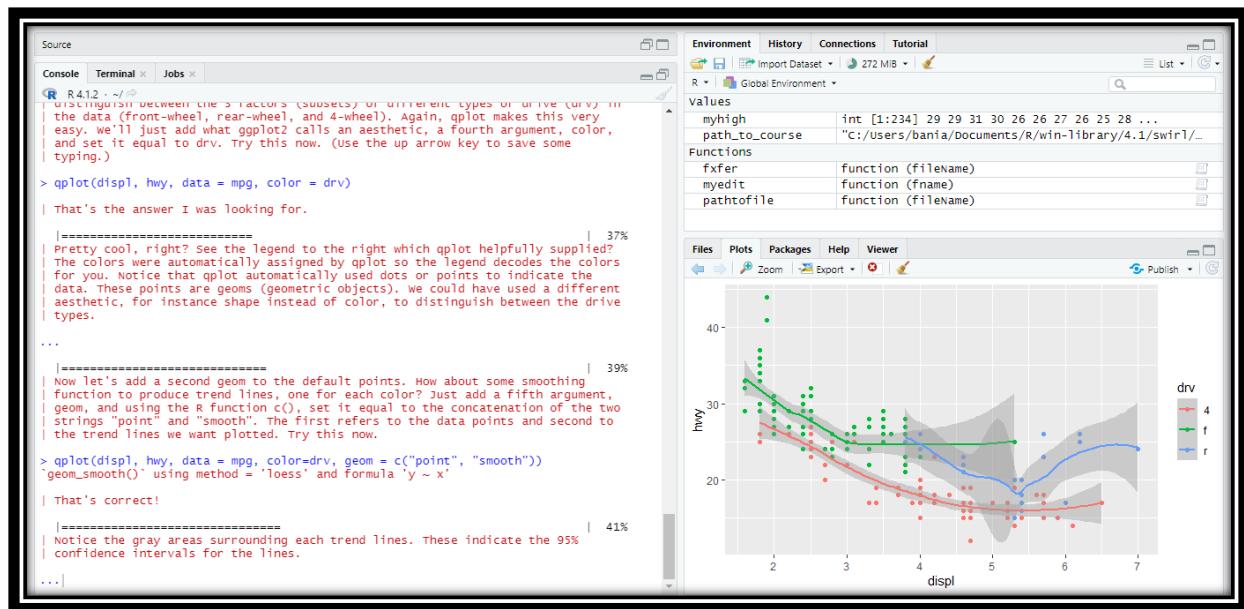
7. Working with Colors

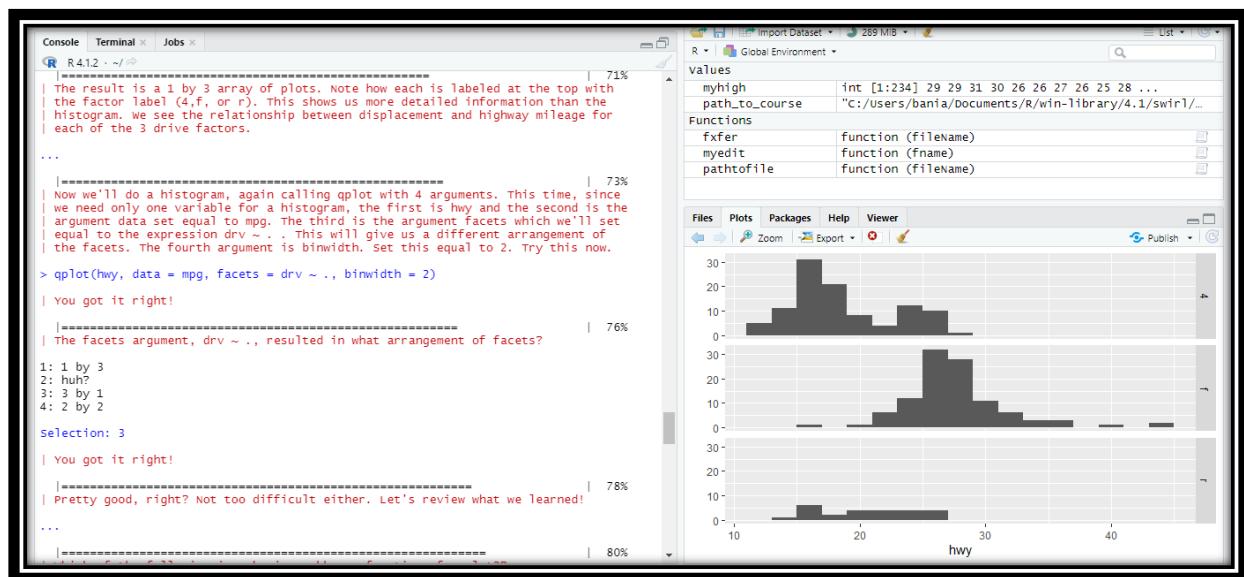
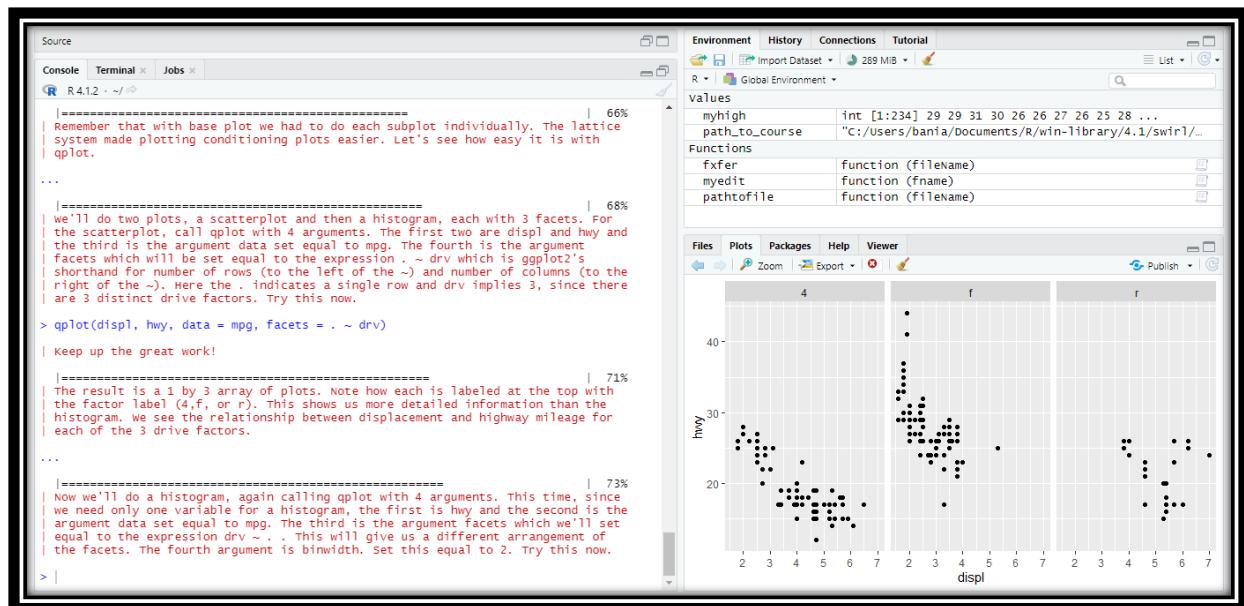




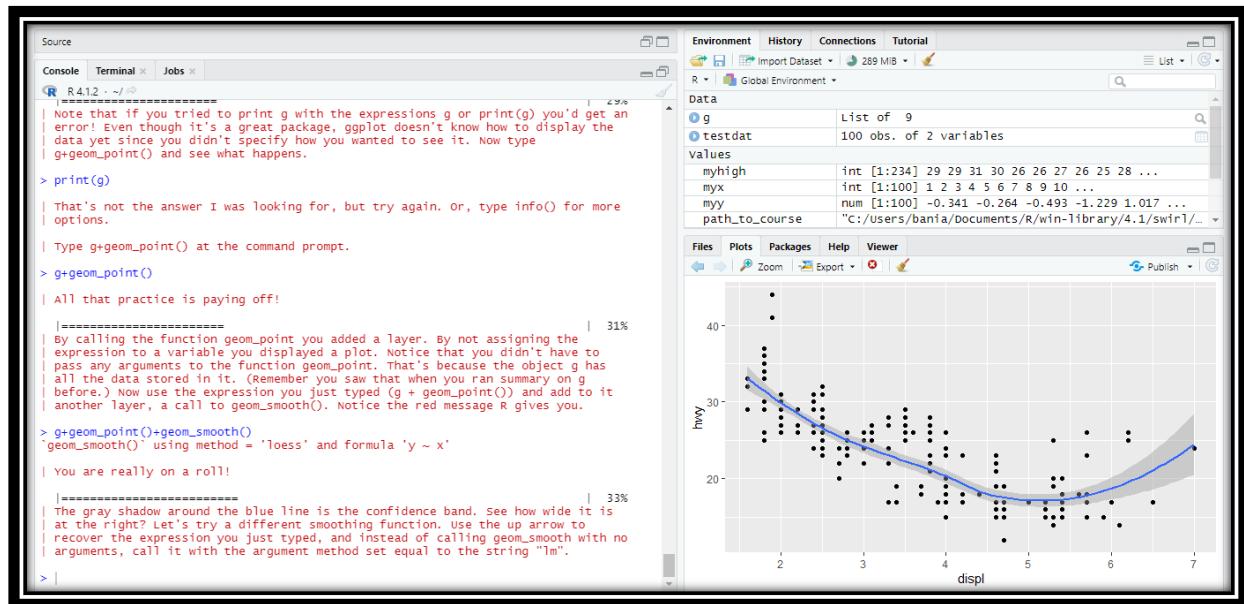
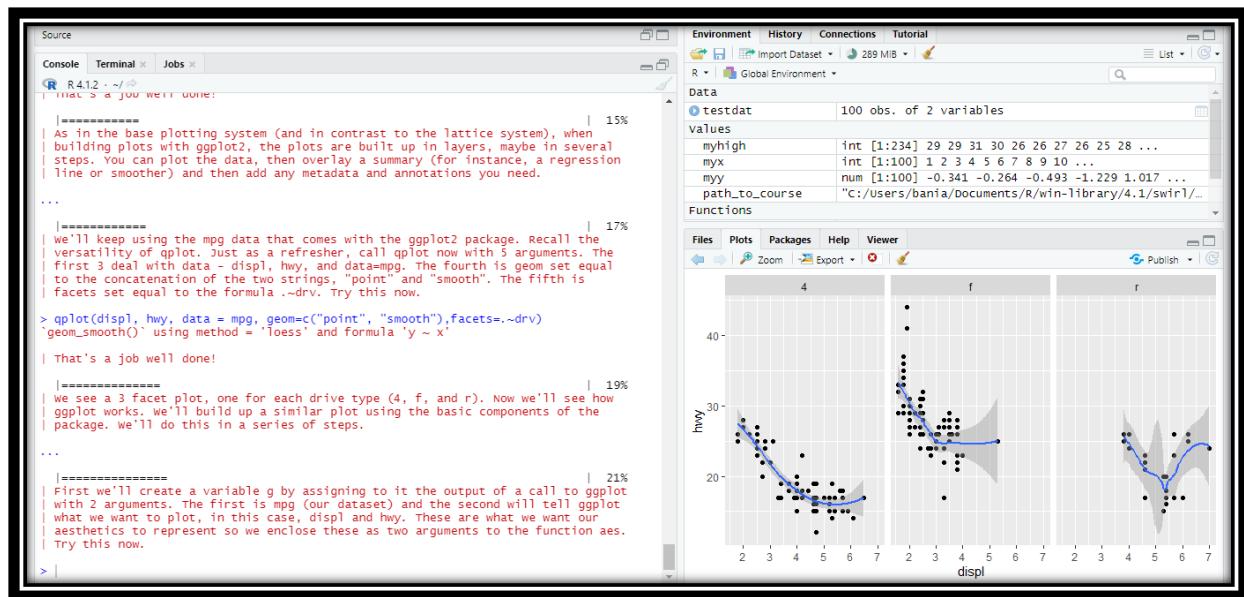
8. GGPLOT2 Part 1

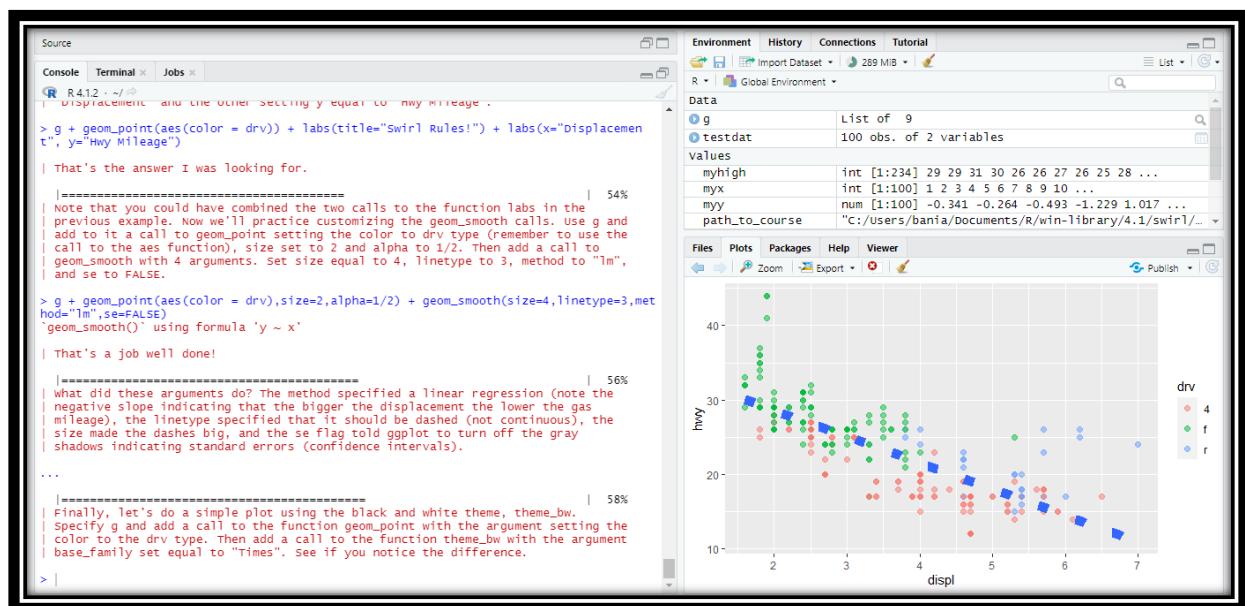
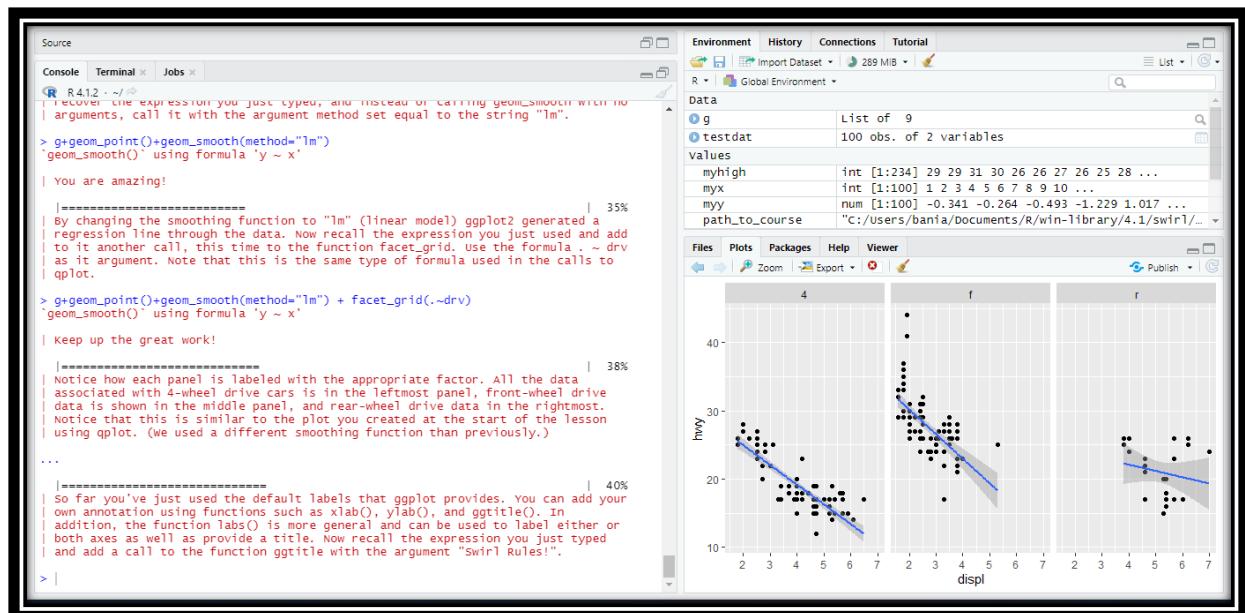






9. GGPLOT2 Part 2





Source

Console Terminal Jobs

```
R 4.1.2 : ~
> plot(myx, myy, type = "l", ylim = c(-3,3))
| All that practice is paying off!
| ====== | 67%
| Notice how plot plotted the points in the (-3,3) range for y-values. The outlier
| at (50,100) is NOT shown on the line plot. Now we'll plot the same data with
| ggplot. Recall that the name of the data frame is testdat. Create the graphical
| object g with a call to ggplot with 2 arguments, testdat (the data) and a call to
| aes with 2 arguments, x set equal to myx, and y set equal to myy.
> g <- ggplot(testdat, aes(x = myx, y = myy))
| That's a job well done!
| ====== | 69%
| Now add a call to geom_line with 0 arguments to g.
> g + geom_line()
| All that hard work is paying off!
| ====== | 71%
| Notice how ggplot DID display the outlier point at (50,100). As a result, the rest
| of the data is smashed down so you don't get to see what the bulk of it looks
| like. The single outlier probably isn't important enough to dominate the graph.
| How do we get ggplot to behave more like plot in a situation like this?
...
| ====== | 73%
| Let's take a guess that in addition to adding geom_line() to g we also just have
| to add ylim(-3,3) to it as we did with the call to plot. Try this now to see what
| happens.
> |
```

Environment History Connections Tutorial

Data

- g List of 9
- testdat 100 obs. of 2 variables

Values

myhigh	int [1:234] 29 29 31 30 26 26 27 26 25 28 ...
myx	int [1:100] 1 2 3 4 5 6 7 8 9 10 ...
myy	num [1:100] -0.341 -0.264 -0.493 -1.229 1.017 ...
path_to_course	"C:/Users/bania/Documents/R/win-library/4.1/swirl/..."

Files Plots Packages Help Viewer

Source

Console Terminal Jobs

```
R 4.1.2 : ~
| A simple, yet comfortingly familiar scatterplot appears. Let's make our display a
| 2 dimensional multi-panel plot. Recall your last command (with the up arrow) and
| add to it a call to the function facet_grid. Give it 2 arguments. The first is the
| formuladrv~cyl, and the second is the argument margins set equal to TRUE. Try
| this now.
> g + geom_point() + facet_grid(drv~cyl,margins=TRUE)
| Great job!
| ====== | 92%
| A 4 by 5 plot, huh? The margins argument tells ggplot to display the marginal
| totals over each row and column, so instead of seeing 3 rows (the number of cyl
| factors) and 4 columns (the number of cyl factors) we see a 4 by 5 display. Note
| that the panel in position (4,5) is a tiny version of the scatterplot of the
| entire dataset.
...
| ====== | 94%
| Now add to your last command (or retype it if you like to type) a call to
| geom_smooth with 4 arguments. These are method set to "lm", se set to FALSE, size
| set to 2, and color set to "black".
> g + geom_point() + facet_grid(drv~cyl,margins=TRUE)+geom_smooth(method="lm",size=2,se
|=FALSE,color="black")
| 'geom_smooth()' using formula 'y ~ x'
| Excellent work!
| ====== | 96%
| Angry Birds? Finally, add to your last command (or retype it if you like to type)
| a call to the function labs with 3 arguments. These are x set to "displacement",
| y set to "Highway Mileage", and title set to "Swirl Rules!".
> |
```

Environment History Connections Tutorial

Data

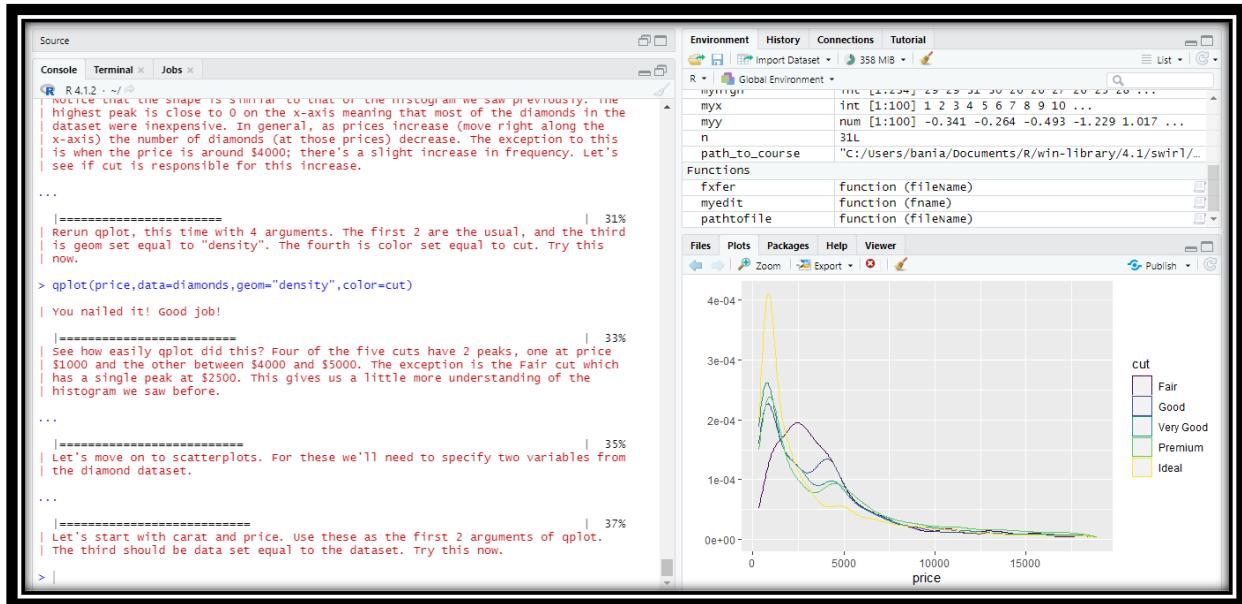
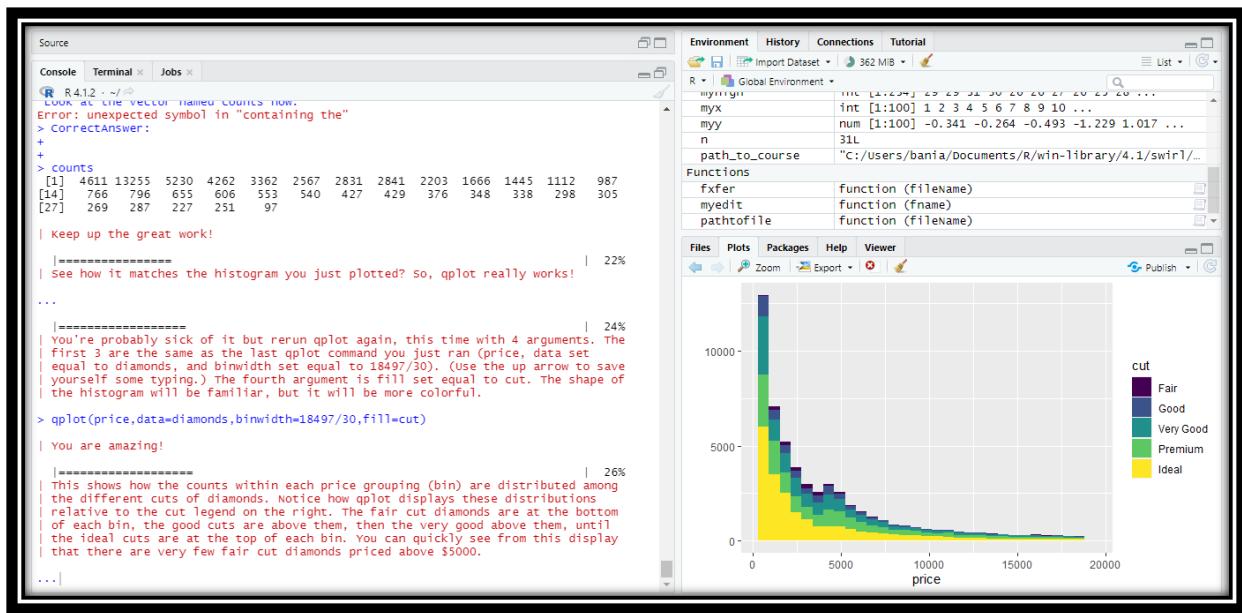
- g List of 9
- testdat 100 obs. of 2 variables

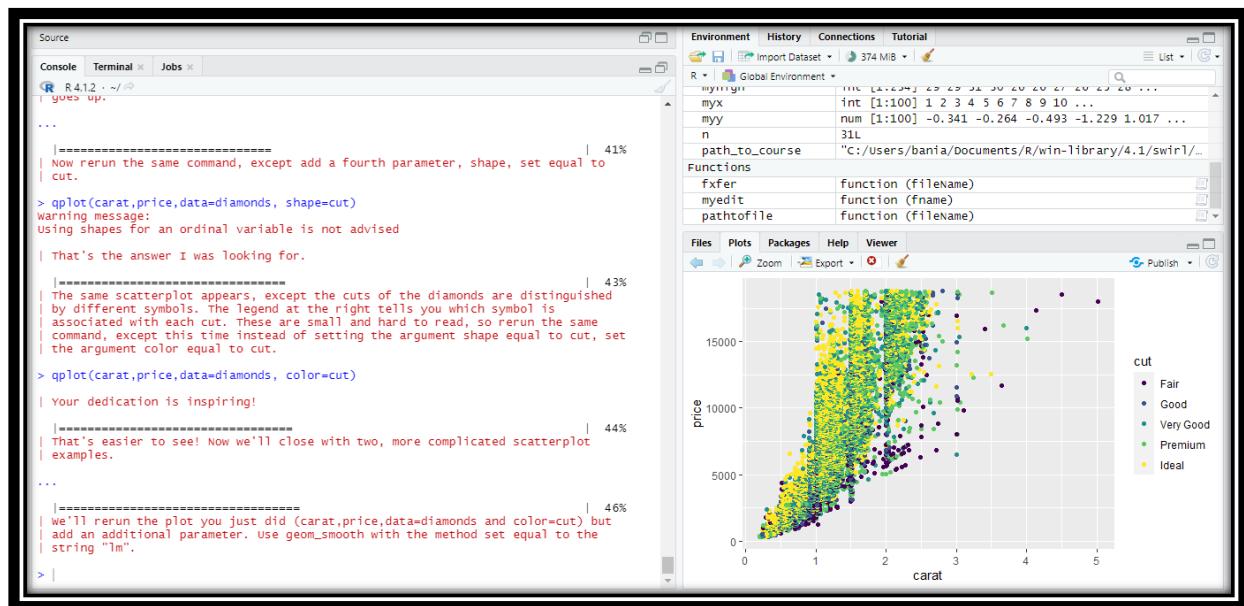
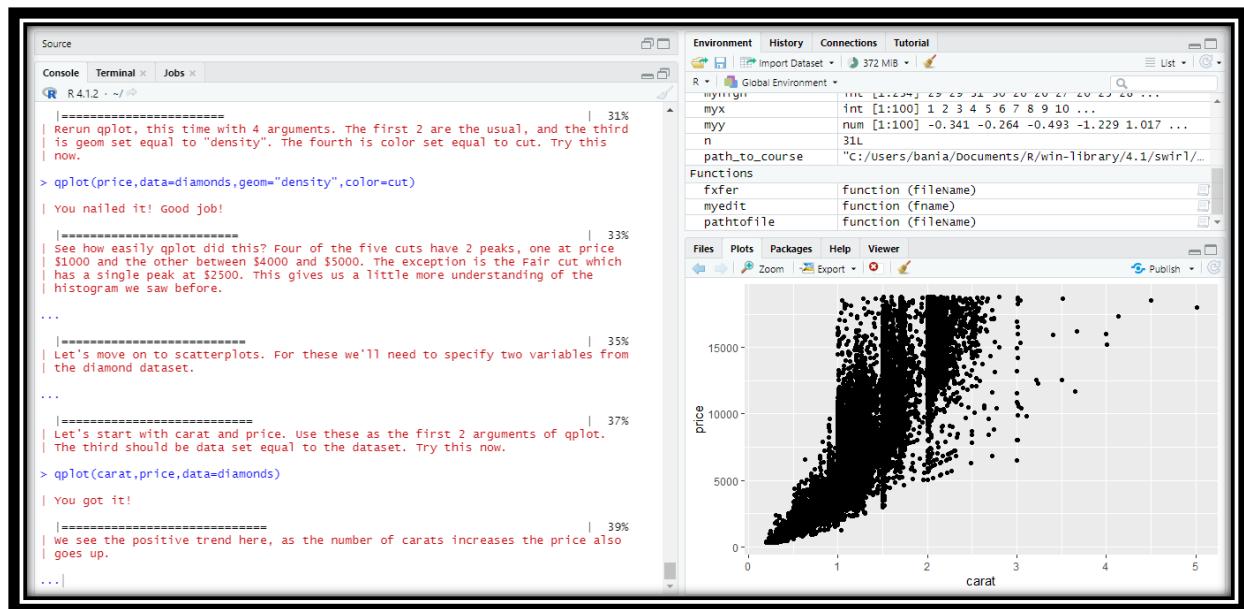
Values

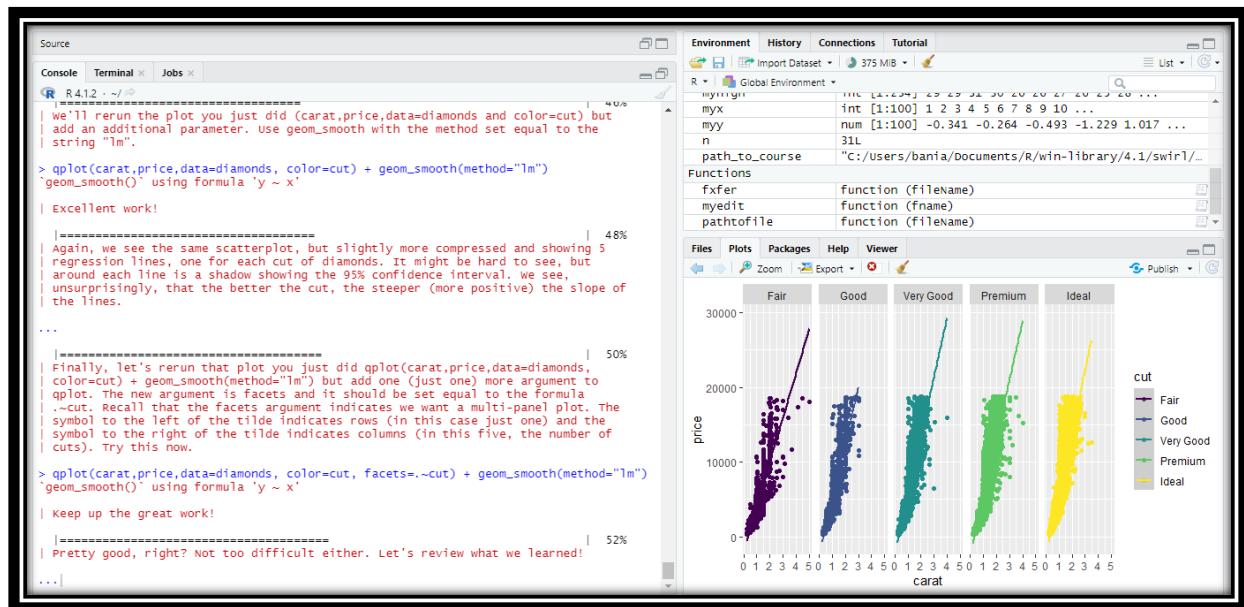
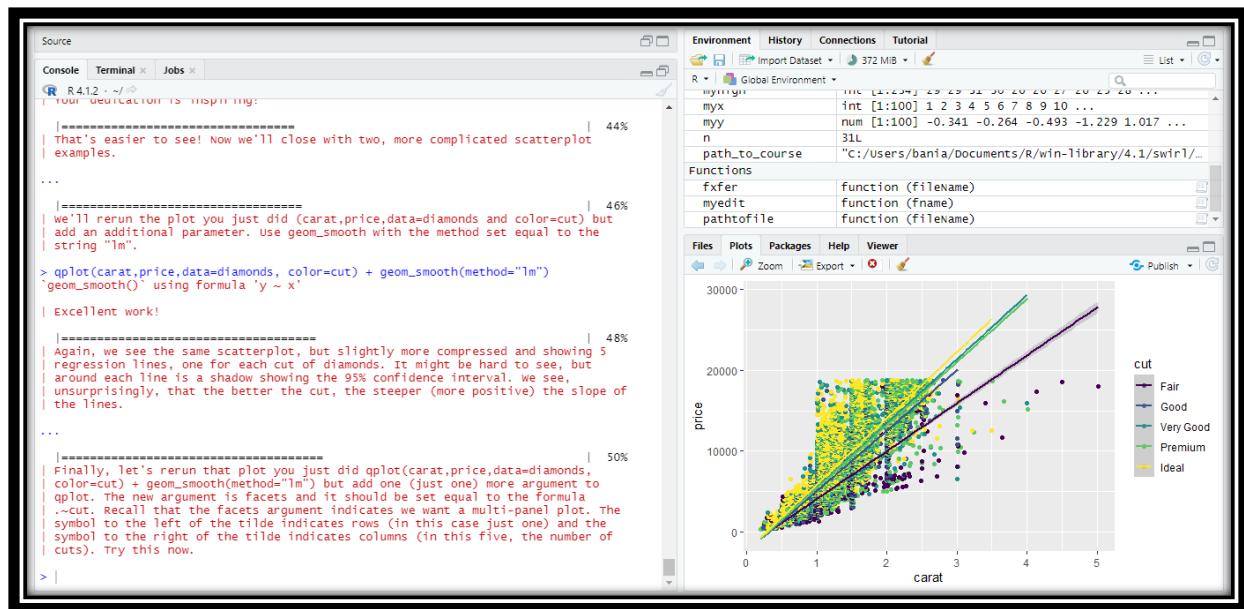
myhigh	int [1:234] 29 29 31 30 26 26 27 26 25 28 ...
myx	int [1:100] 1 2 3 4 5 6 7 8 9 10 ...
myy	num [1:100] -0.341 -0.264 -0.493 -1.229 1.017 ...
path_to_course	"C:/Users/bania/Documents/R/win-library/4.1/swirl/..."

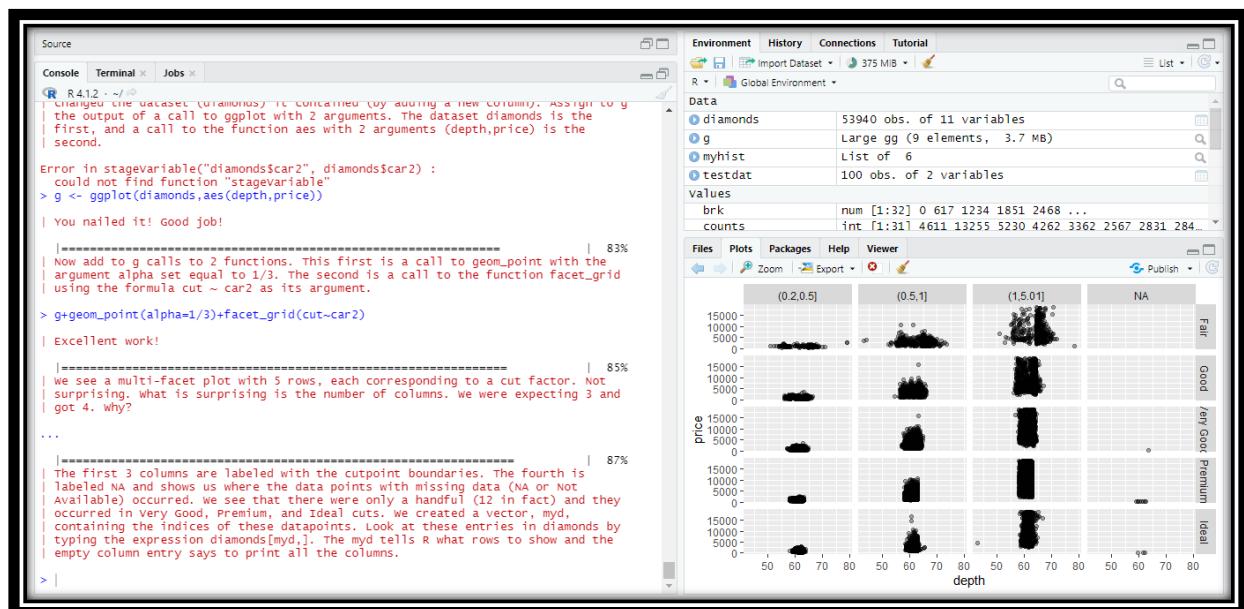
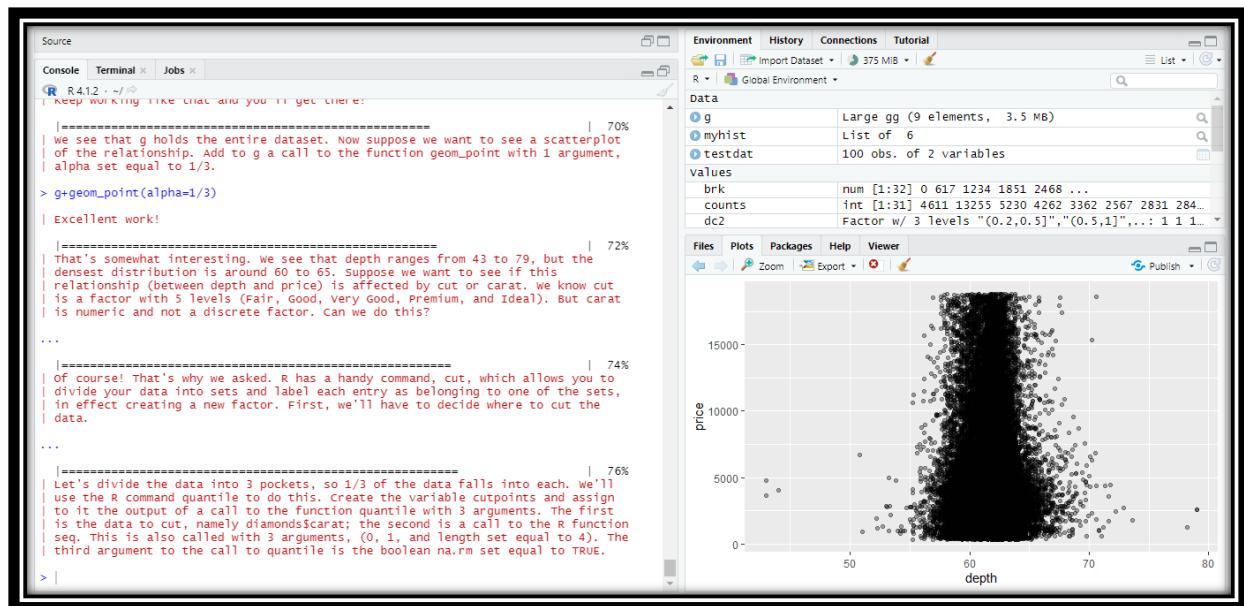
Files Plots Packages Help Viewer

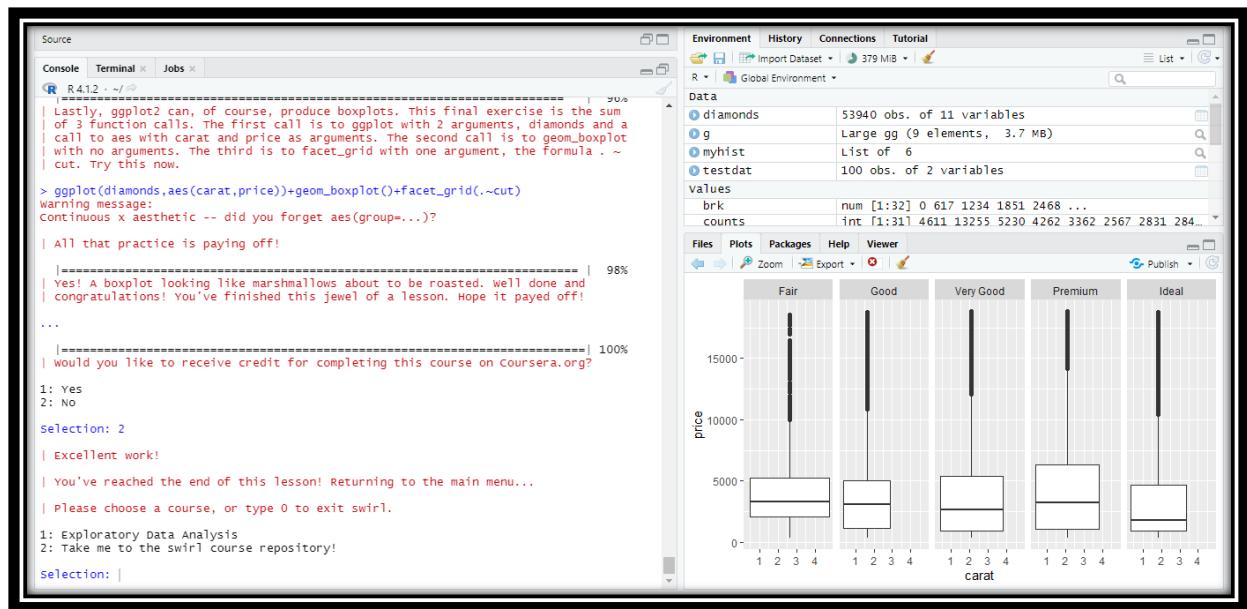
10. GGPLOT2 Extra



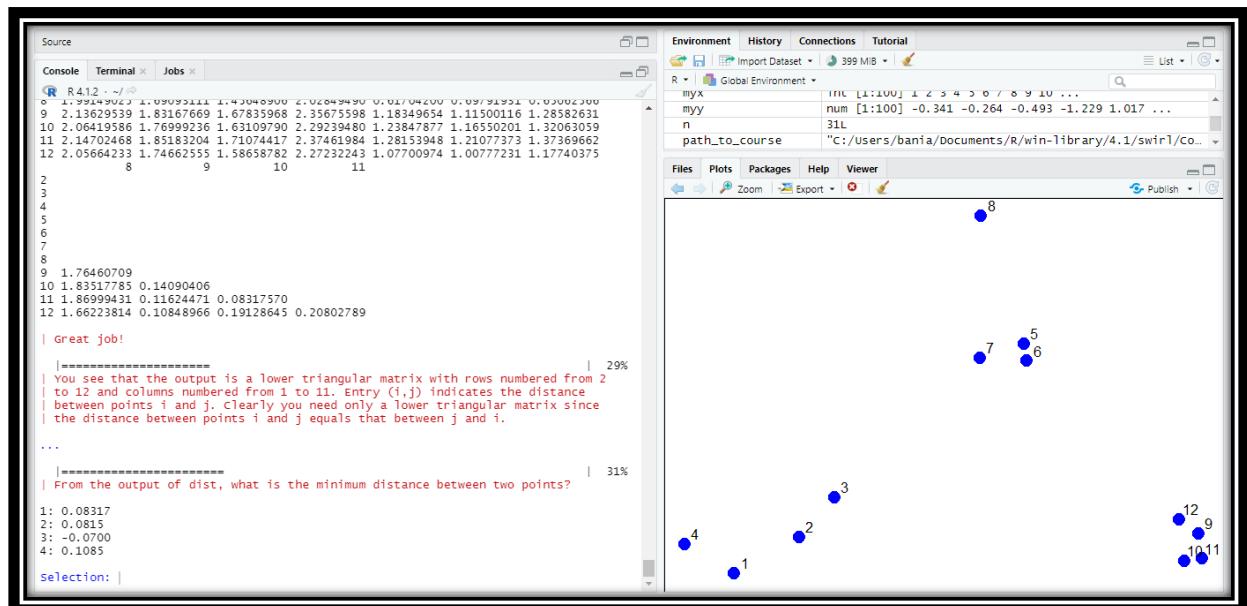


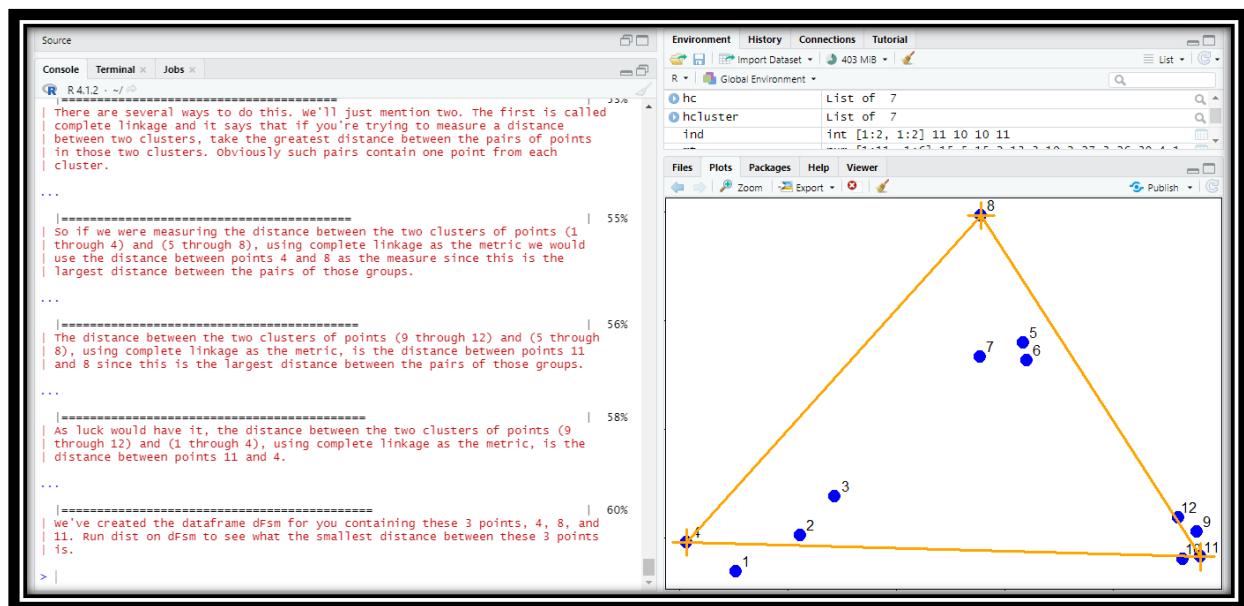
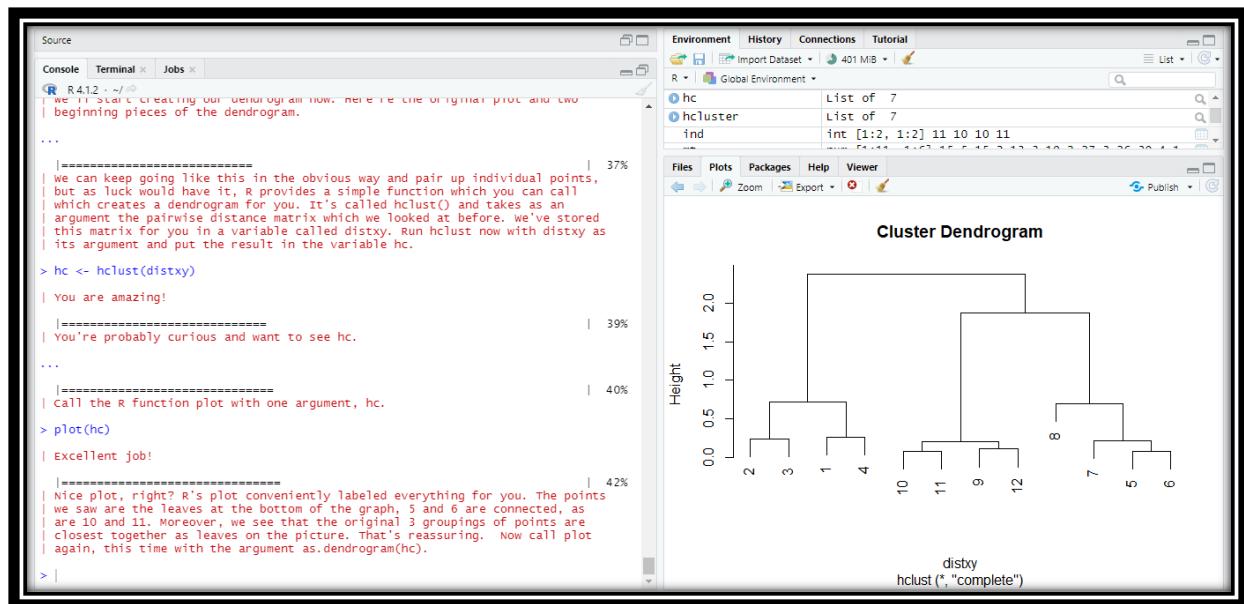






11. Hierarchical Clustering





The last method of visualizing data we'll mention in this lesson concerns heat maps. Wikipedia (http://en.wikipedia.org/wiki/Heat_map) tells us a heat map is a graphical representation of data where the individual values contained in a matrix are represented as colors. ... Heat maps originated in 2D displays of the values in a data matrix. Larger values were represented by small dark gray or black squares (pixels) and smaller values by lighter squares."

You've probably seen many examples of heat maps, for instance weather radar and displays of ocean salinity. From Wikipedia (http://en.wikipedia.org/wiki/Heat_map) we learn that heat maps are often used in molecular biology "to represent the level of expression of many genes across a number of comparable samples (e.g. cells in different states, samples from different patients) as they are obtained from DNA microarrays."

we won't say too much on this topic, but a very nice concise tutorial on creating heatmaps in R exists at http://sebastiansraschka.com/Articles/heatmaps_in_r.html#clustering. Here's an image from the tutorial to start you thinking about the topic. It shows a sample heat map with a dendrogram on the left edge mapping the relationship between the rows. The legend at the top shows how colors relate to values.

R provides a handy function to produce heat maps. It's called heatmap. We've put the point data we've been using throughout this lesson in a matrix. Call heatmap now with 2 arguments. The first is dataMatrix and the second is col set equal to cm.colors(25). This last is optional, but we like the colors better than the default ones.

12. K Means Clustering

gravity) of the cluster or points assigned to them, we have to do the x and y coordinates separately. We'll do the x coordinate first. Recall that the x vectors x and y hold the respective coordinates of our 12 data points.

We can use the R function tapply which applies a function over a ragged array'. This means that every element of the array is assigned a factor and the function is applied to subsets of the array (identified by the factor vector). This allows us to take advantage of the factor vector newClust we calculated. Call tapply now with 3 arguments, x (the data), newClust (the factor array), and mean (the function to apply).

```
> tapply(x,newClust,mean)
 1 2 3
1.210767 1.010320 2.498011
```

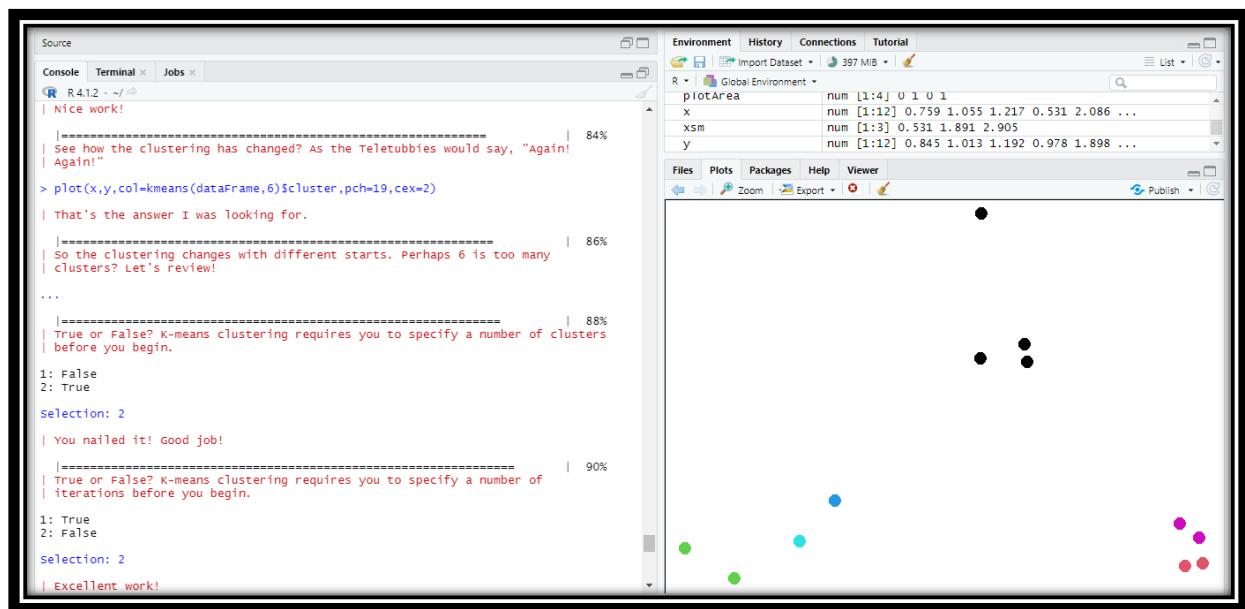
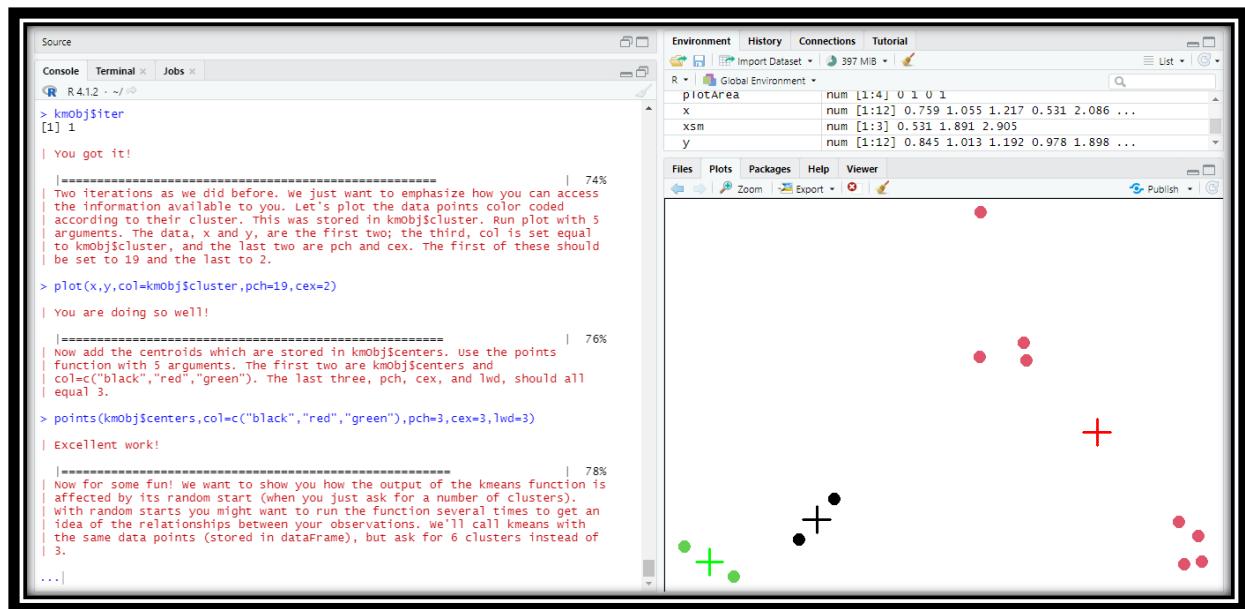
Perseverance, that's the answer.

Repeat the call, except now apply it to the vector y instead of x.

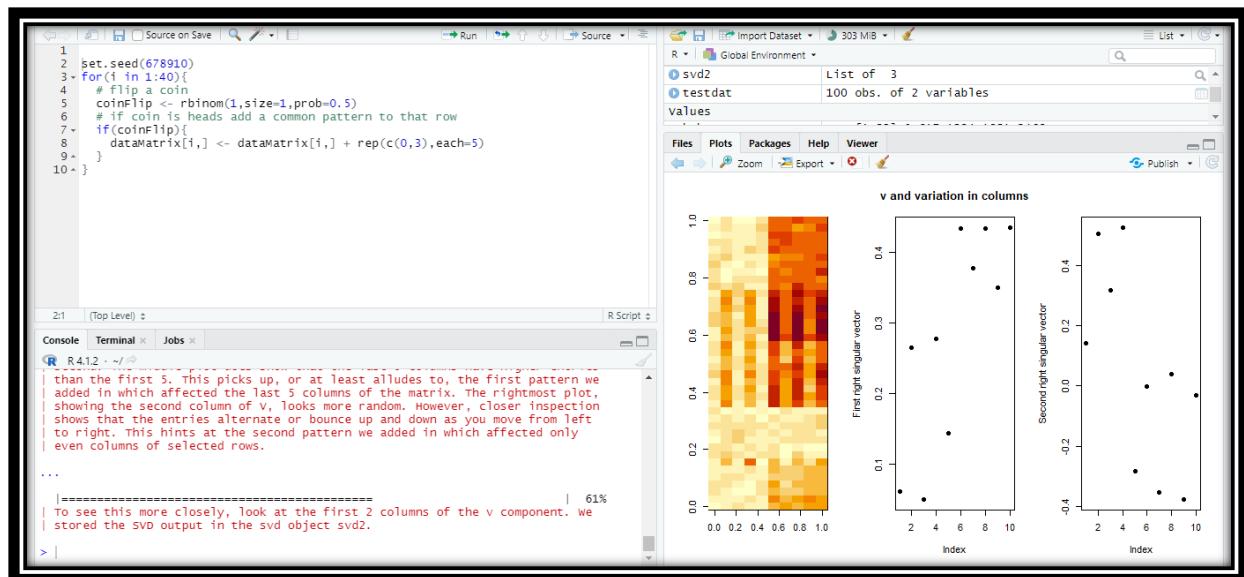
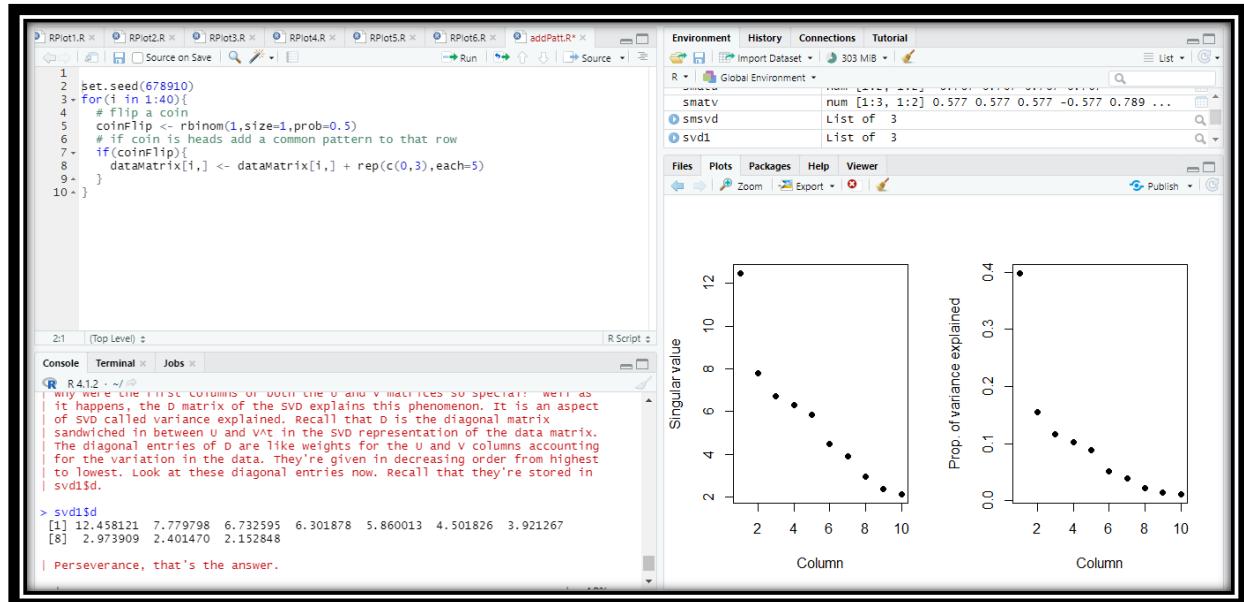
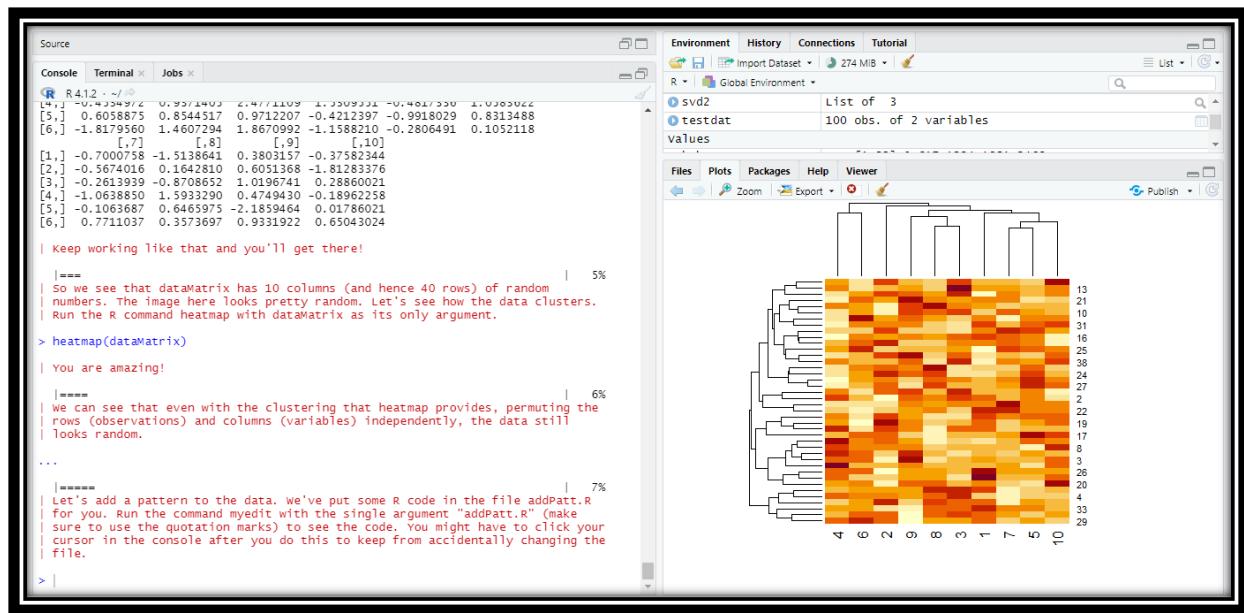
```
> tapply(y,newClust,mean)
 1 2 3
1.730555 1.016513 1.354373
```

Your dedication is inspiring!

Now that we have new x and new y coordinates for the 3 centroids we can plot them. We've stored off the coordinates for you in variables newCx and newCy. Use the R command points with these as the first 2 arguments. In addition, use the arguments col set equal to cols1, pch equal to 8, cex equal to 2 and lwd also equal to 2.



13. Dimension Reduction



R Script

```

1 set.seed(678910)
2 for(i in 1:40){
3   # flip a coin
4   coinFlip <- rbinom(1,size=1,prob=0.5)
5   # if coin is heads add a common pattern to that row
6   if(coinFlip){
7     dataMatrix[i,] <- dataMatrix[i,] + rep(c(0,3),each=5)
8   }
9 }
10

```

2:1 (Top Level) R Script

Console Terminal Jobs

R 4.1.2 ~/

...
|=====| 72%
| consider this low resolution image file showing a face. We'll use SVD and see
| how the first several components contain most of the information in the file
| so that storing a huge matrix might not be necessary.
...
|=====| 73%
| The image data is stored in the matrix faceData. Run the R command dim on
| faceData to see how big it is.
> |

R Script

```

1 set.seed(678910)
2 for(i in 1:40){
3   # flip a coin
4   coinFlip <- rbinom(1,size=1,prob=0.5)
5   # if coin is heads add a common pattern to that row
6   if(coinFlip){
7     dataMatrix[i,] <- dataMatrix[i,] + rep(c(0,3),each=5)
8   }
9 }
10

```

2:1 (Top Level) R Script

Console Terminal Jobs

R 4.1.2 ~/

> a2 <- svd1\$u[,1:2] %*% diag(svd1\$d[1:2]) %*% t(svd1\$v[,1:2])
| Great job!
|=====| 87%
| Use myImage again to see how a2 displays.
> myImage(svd1\$u[,1:5] %*% diag(svd1\$d[1:5]) %*% t(svd1\$v[,1:5]))
| One more time. You can do it! Or, type info() for more options.
| Type myImage(a2) at the command prompt.
> myImage(svd1\$u[,1:10] %*% diag(svd1\$d[1:10]) %*% t(svd1\$v[,1:10]))

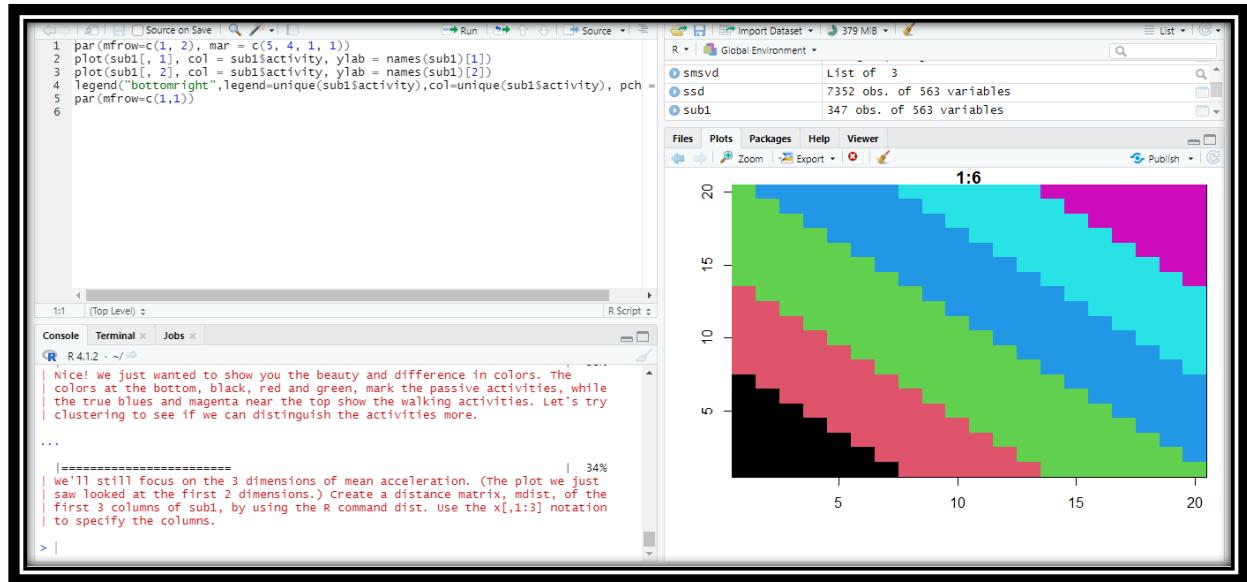
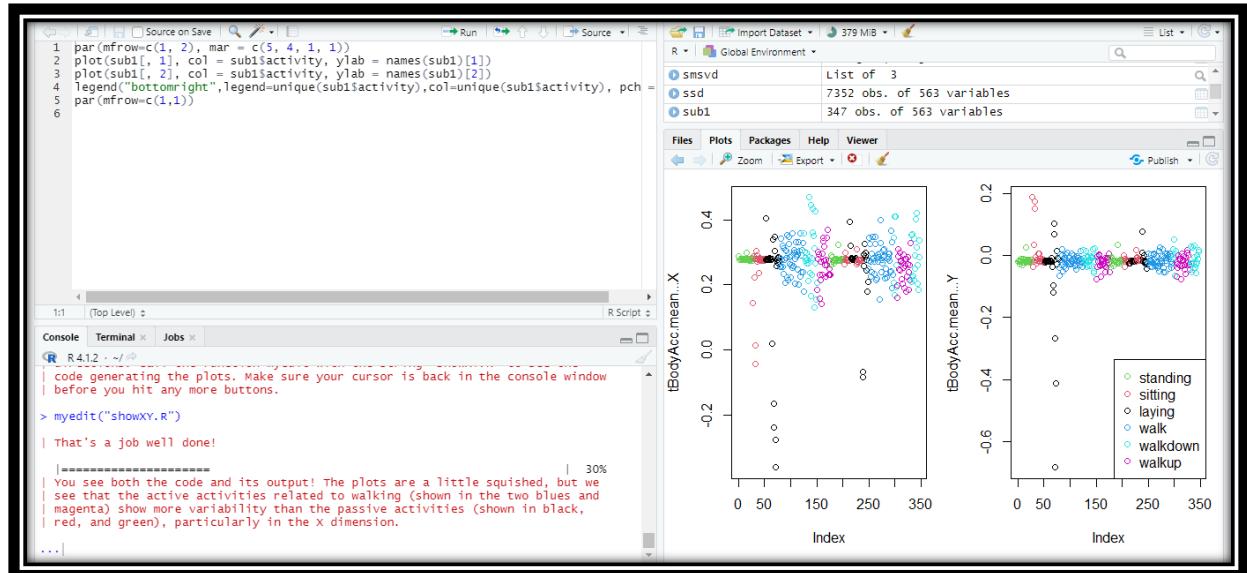
R Script

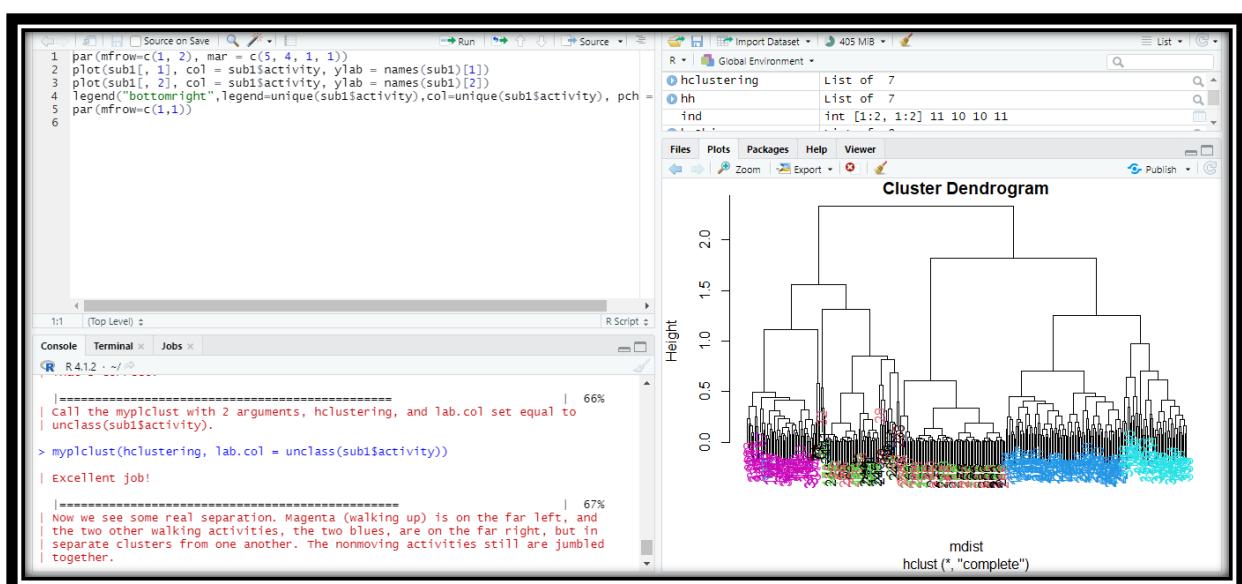
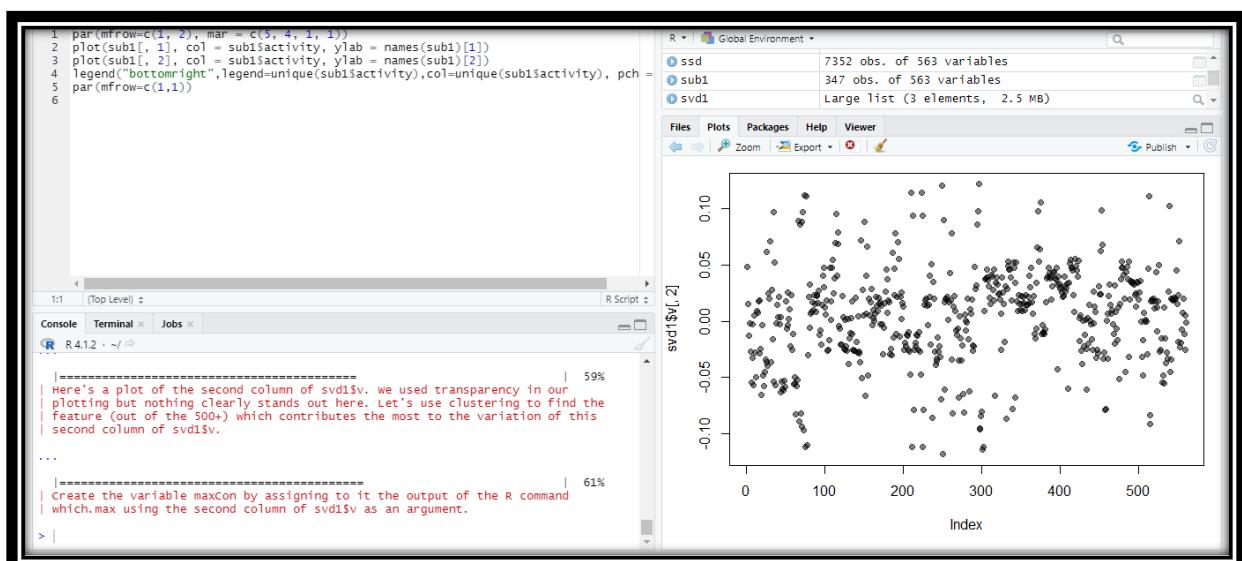
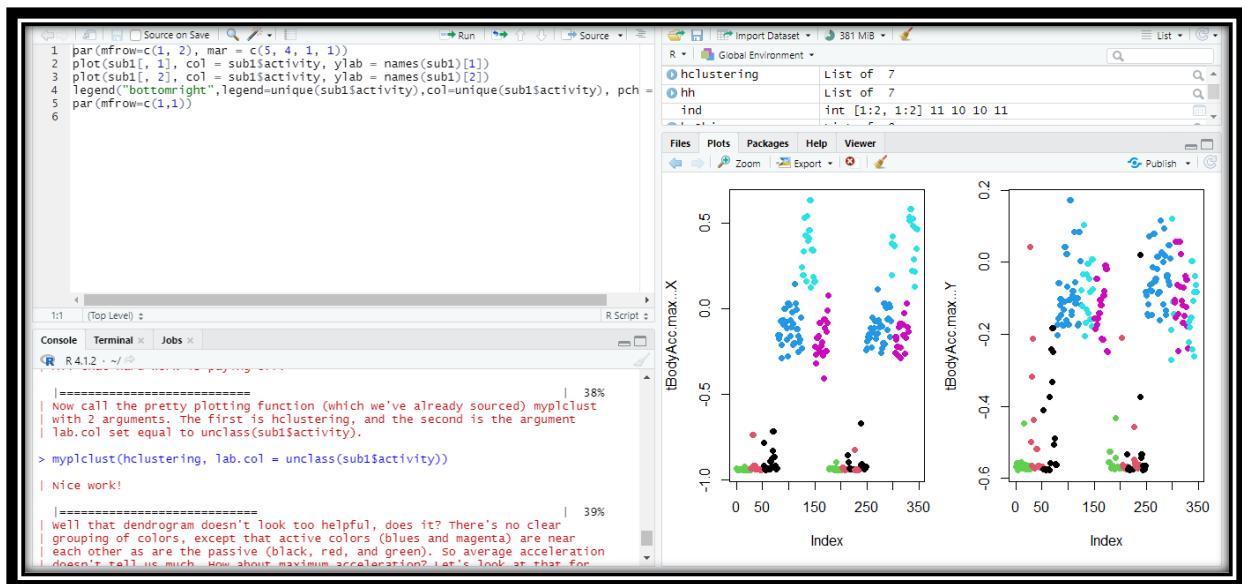
```

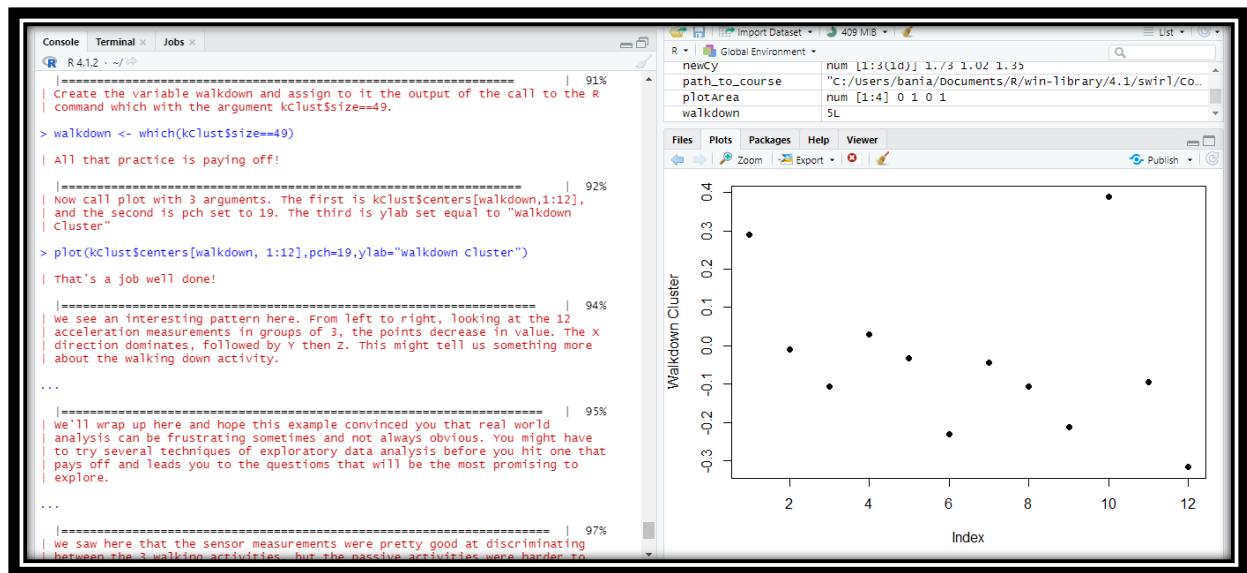
...  
|=====| 92%  
| we'll close now with a few comments. First, when reducing dimensions you have  
| to pay attention to the scales on which different variables are measured and  
| make sure that all your data is in consistent units. In other words, scales of  
| your data matter. Second, principal components and singular values may mix  
| real patterns, as we saw in our simple 2-pattern example, so finding and  
| separating out the real patterns require some detective work. Let's do a quick  
| review now.  
...  
|=====| 93%  
| Which of the following cliches LEAST captures the essence of dimension  
| reduction?  
1: see the forest through the trees  
2: a face that could launch a 1000 ships  
3: find the needle in the haystack  
4: separate the wheat from the chaff  
selection: 2  
| You nailed it! Good job!  
|=====| 94%  
| A matrix X has the singular value decomposition  $UDV^T$ . The principal  
| components of X are ?  
1: the rows of U  
2: the rows of V  
3: the columns of V  
4: the columns of U  
selection: |

```

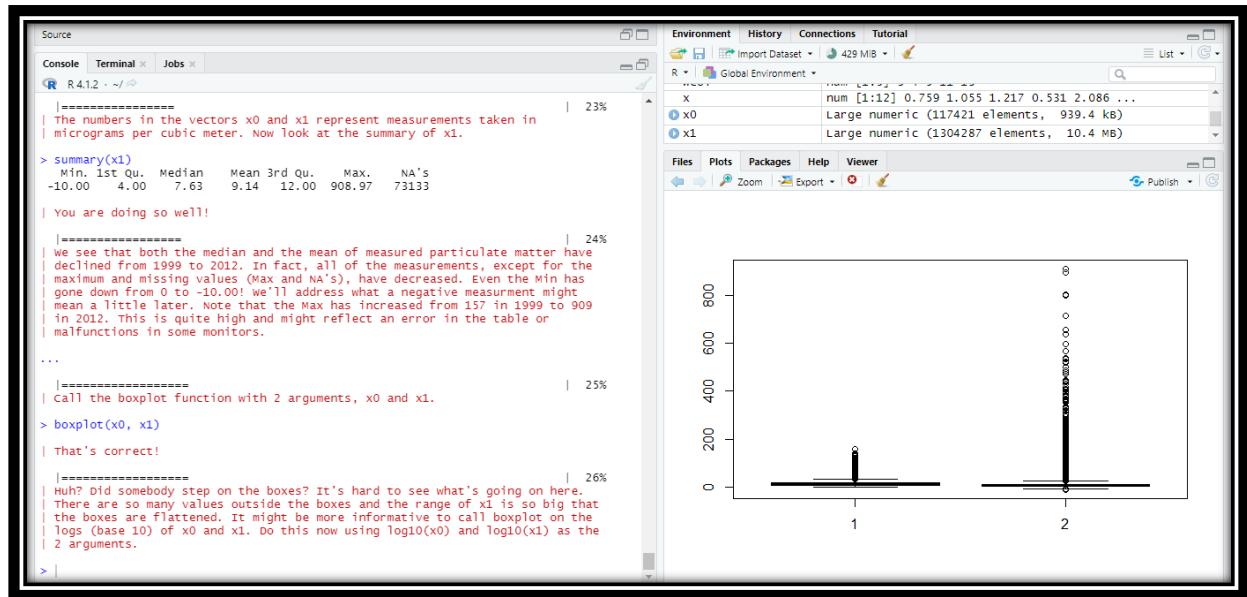
14. Clustering Example







15. CaseStudy



Source

Console Terminal Jobs

R 4.1.2 - ~

```
| Huh? Did somebody step on the boxes? It's hard to see what's going on here.
| There are so many values outside the boxes and the range of x1 is so big that
| the boxes are flattened. It might be more informative to call boxplot on the
| Logs (base 10) of x0 and x1. Do this now using log10(x0) and log10(x1) as the
| 2 arguments.

> boxplot(log10(x0), log10(x1))
Warning messages:
1: In boxplot.default(log10(x0), log10(x1)) : Nans produced
2: In bplot(at[i], wid = width[i], stats = z$stats[, i], out = z$out[z$group == :
outlier <-Inf] in boxplot.1 is not drawn
3: In bplot(at[i], wid = width[i], stats = z$stats[, i], out = z$out[z$group == :
outlier <-Inf] in boxplot.2 is not drawn

| You are really on a roll!

| -----
| A bonus! Not only do we get a better looking boxplot we also get some warnings
| from R in Red. These let us know that some values in x0 and x1 were
| "unlogable", no doubt the 0 (Min) we saw in the summary of x0 and the
| negative values we saw in the Min of the summary of x1.

...
| -----
| From the boxplot (x0 on the left and x1 on the right), what can you say about
| the data?

1: The median of x1 is less than the median of x0
2: The range of x0 is greater than the range of x1
3: The mean of x1 is less than the mean of x0
4: The boxes are too small to interpret

selection: 1

| All that practice is paying off!
```

Environment History Connections Tutorial

Import Dataset 501 MB

Global Environment

dates

dc2

distmt

distxv

Files Plots Packages Help Viewer

Source

Console Terminal Jobs

R 4.1.2 - ~

```
| Day: Reassuring dates the output of a call to as.Date with the z arguments
| as.character(dates) as the first argument and the string "%Y%m%d" as the
| second.

> dates <- as.Date(as.character(dates), "%Y%m%d")
| That's correct!

| -----
| Now when you run head on dates you'll see the dates in a nicer format. Try
| this now.

> head(dates)
[1] "2012-01-01" "2012-01-04" "2012-01-07" "2012-01-10" "2012-01-13" "2012-01-16"
| Nice work!

| -----
| Let's plot a histogram of the months when the particulate matter measurements
| are negative. Run hist with 2 arguments. The first is dates[negative] and the
| second is the string "month".

> hist(dates[negative], "month")
| All that practice is paying off!

| -----
| we see the bulk of the negative measurements were taken in the winter months,
| with a spike in May. Not many of these negative measurements occurred in
| summer months. we can take a guess that because particulate measures tend to
| be low in winter and high in summer, coupled with the fact that higher
| densities are easier to measure, that measurement errors occurred when the
| values were low. For now we'll attribute these negative measurements to
| errors. Also, since they account for only 2% of the 2012 data, we'll ignore
| them.

...
```

Environment History Connections Tutorial

Import Dataset 820 MB

Global Environment

dates

dc2

distmt

distxv

Files Plots Packages Help Viewer

Histogram of dates[negative]

Source

Console Terminal Jobs

R 4.1.2 - ~

```
| Do the same for the 2012 data. Specifically, create dates1 using pm25subdate
| as your input.

> dates1 <- as.Date(as.character(pm25subDate), "%Y%m%d")
| You are doing so well!

| -----
| Now we'll plot these 2 time series in the same panel using the base plotting
| system. call par with 2 arguments. The first is mfrow set equal to c(1,2).
| This will tell the system we're plotting 2 graphs in 1 row and 2 columns. The
| second argument will adjust the panel's margins. It is mar set to c(4,4,2,1).

> par(mfrow = c(1, 2), mar = c(4, 4, 2, 1))
| Great job!

| -----
| call plot with the 3 arguments dates0, x0sub, and pch set to 20. The first two
| arguments are the x and y coordinates. This will show the pm25 values as
| functions of time.

> plot(dates0, x0sub, pch = 20)
| Keep up the great work!

| -----
| Now we'll mark the median.

...
| -----
| use abline to add a horizontal line at the median of the pm25 values. Make the
| line width 2 (lwd is the argument), and when you call median with x0sub,
| specify the argument na.rm to be TRUE.

> |
```

Environment History Connections Tutorial

Import Dataset 846 MB

Global Environment

dates1

dc2

distmt

distxv

Files Plots Packages Help Viewer

