```yaml
1     Course: Exploratory_Data_Analysis
2     Lesson: GGPlot2_Part1
3
4   - Class: text
5     Output: "GGPlot2_Part1. (Slides for this and other Data Science courses may be found
      at github 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/ggplot2.)"
6
7
8   - Class: text
9     Output:  In another lesson, we gave you an overview of the three plotting systems in
      R. In this lesson we'll focus on the third and newest plotting system in R, ggplot2.
      As we did with the other two systems, we'll focus on creating graphics on the screen
      device rather than another graphics device.
10
11  - Class: text
12    Output:  The ggplot2 package is an add-on package available from CRAN via
      install.packages(). (Don't worry, we've installed it for you already.) It is an
      implementation of The Grammar of Graphics, an abstract concept (as well as book)
      authored and invented by Leland Wilkinson and implemented by Hadley Wickham while he
      was a graduate student at Iowa State. The web site http://ggplot2.org provides
      complete documentation.
13
14  - Class: text
15    Output: A grammar of graphics represents an abstraction of graphics, that is,  a
      theory of graphics which conceptualizes basic pieces from which you can build new
      graphics and graphical objects. The  goal of the grammar is to "Shorten the distance
      from mind to page". From Hadley Wickham's book  we learn that
16
17  - Class: text
18    Output:  The ggplot2 package "is composed of a set of independent components that can
      be composed in many different ways. ... you can create new graphics that are
      precisely tailored for your problem." These components include aesthetics which are
      attributes such as colour, shape, and  size, and geometric objects or geoms such as
      points, lines, and bars.
19
20  - Class: text
21    Output: Before we delve into details, let's review the other 2 plotting systems.
22
23  - Class: mult_question
24    Output: Recall what you know about R's base plotting system. Which of the following
      does NOT apply to it?
25    AnswerChoices: Start with plot (or similar) function; Use annotation functions to
      add/modify (text, lines, points, axis); It is convenient and mirrors how we think of
      building plots and analyzing data; Can easily go back once the plot has started
      (e.g., to adjust margins or correct a typo)
26    CorrectAnswer: Can easily go back once the plot has started (e.g., to adjust margins
      or correct a typo)
27    AnswerTests: omnitest(correctVal='Can easily go back once the plot has started (e.g.,
      to adjust margins or correct a typo)')
28    Hint: Which choice is the only one which looks backward?
29
30
31  - Class: mult_question
32    Output: Recall what you know about R's lattice plotting system. Which of the
      following does NOT apply to it?
33    AnswerChoices:  Plots are created with a single function call (xyplot, bwplot, etc.);
      Most useful for conditioning types of plots and putting many panels on one plot;
      Margins and spacing are set automatically because entire plot is specified at once;
      Can always add to the plot once it is created
34    CorrectAnswer: Can always add to the plot once it is created
35    AnswerTests: omnitest(correctVal='Can always add to the plot once it is created')
36    Hint: Which choice is the only one which is inconsistent with the other three?
37
38  - Class: mult_question
39    Output: If we told you that ggplot2 combines the best of base and lattice, that would
      mean it ...?
40    AnswerChoices: Automatically deals with spacings, text, titles but also allows you to
```

```
            annotate; Like lattice it allows for multipanels but more easily and intuitively; Its
            default mode makes many choices for you (but you can customize!); All of the others
41          CorrectAnswer: All of the others
42          AnswerTests: omnitest(correctVal='All of the others')
43          Hint: Which choice is the only one that encompasses the other three?
44
45     - Class: text
46       Output: Yes, ggplot2 combines the best of base and lattice. It allows for multipanel
            (conditioning) plots (as lattice does) but also post facto annotation (as base does),
            so you can add titles and labels. It uses the low-level grid package (which comes
            with R) to draw the graphics. As part of its grammar philosophy, ggplot2 plots are
            composed of aesthetics (attributes such as size, shape, color) and geoms (points,
            lines, and bars), the geometric objects you see on the plot.
47
48     - Class: text
49       Output: The ggplot2 package has 2 workhorse functions. The more basic workhorse
            function is qplot, (think quick plot), which works like the plot function in the base
            graphics system. It can produce many types of plots (scatter, histograms, box and
            whisker) while hiding tedious details from the user. Similar to lattice functions, it
            looks for data in a data frame or parent environment.
50
51     - Class: text
52       Output: The more advanced workhorse function in the package is ggplot, which is  more
            flexible and can be customized for doing things qplot cannot do. In this lesson we'll
            focus on qplot.
53
54     - Class: cmd_question
55       Output: We'll start by showing how easy and versatile qplot is. First, let's look at
            some data which comes with the ggplot2 package. The mpg data frame contains fuel
            economy data for 38 models of cars manufactured in 1999 and 2008. Run the R command
            str with the argument mpg. This will give you an idea of what mpg contains.
56       CorrectAnswer: str(mpg)
57       AnswerTests: omnitest(correctExpr='str(mpg)')
58       Hint: Type str(mpg) at the command prompt.
59
60     - Class: cmd_question
61       Output: We see that there are 234 points in the dataset concerning 11 different
            characteristics of the cars. Suppose we want to see if there's a correlation between
            engine displacement (displ) and highway miles per gallon (hwy). As we did with the
            plot function of the base system we could simply call  qplot with 3 arguments, the
            first two are the variables we want to examine and the third argument data is set
            equal to the name of the dataset which contains them (in this case, mpg). Try this now.
62       CorrectAnswer: qplot(displ, hwy, data = mpg)
63       AnswerTests: omnitest(correctExpr='qplot(displ, hwy, data = mpg)')
64       Hint: Type qplot(displ, hwy, data = mpg) at the command prompt.
65
66     - Class: cmd_question
67       Output: A nice scatterplot done simply, right? All the labels are provided. The first
            argument is shown along the x-axis and the second along the y-axis. The negative
            trend (increasing displacement and lower gas mileage) is pretty clear. Now suppose we
            want to do the same plot but this time use different colors to distinguish between
            the 3 factors (subsets) of different types of drive (drv) in  the data (front-wheel,
            rear-wheel, and 4-wheel). Again, qplot makes this very easy. We'll just add what
            ggplot2 calls an aesthetic, a fourth argument, color, and set it equal to drv. Try
            this now. (Use the up arrow key to save some typing.)
68       CorrectAnswer: qplot(displ, hwy, data = mpg, color = drv)
69       AnswerTests: omnitest(correctExpr='qplot(displ, hwy, data = mpg, color = drv)')
70       Hint: Type qplot(displ, hwy, data = mpg, color = drv) at the command prompt.
71
72     - Class: text
73       Output: Pretty cool, right? See the legend to the right which qplot helpfully
            supplied? The colors were automatically assigned by qplot so the legend decodes the
            colors for you. Notice that qplot automatically used dots or points to indicate the
            data. These points are geoms (geometric objects). We could have used a different
            aesthetic, for instance shape instead of color, to distinguish between the drive types.
74
75     - Class: cmd_question
76       Output: Now let's add a second geom to the default points. How about some smoothing
            function to produce trend lines, one for each color? Just add a fifth argument, geom,
```

and using the R function c(), set it equal to the concatenation of the two strings
"point" and "smooth". The first refers to the data points and second to the trend
lines we want plotted. Try this now.

77   **CorrectAnswer:** qplot(displ, hwy, data = mpg, color=drv, geom = c("point", "smooth"))
78   **AnswerTests:** omnitest(correctExpr='qplot(displ, hwy, data = mpg, color=drv, geom =
     c("point", "smooth"))')
79   **Hint:** Type qplot(displ, hwy, data = mpg, color=drv, geom = c("point", "smooth")) at
     the command prompt.
80
81   - **Class:** text
82     **Output:** Notice the gray areas surrounding each trend lines. These indicate the 95%
       confidence intervals for the lines.
83
84   - **Class:** cmd_question
85     **Output:** Before we leave qplot's scatterplotting ability, call qplot again, this time
       with 3 arguments. The first is y set equal to hwy, the second is data set equal to
       mpg, and the third is color set equal to drv. Try this now.
86     **CorrectAnswer:** qplot(y=hwy, data = mpg, color = drv)
87     **AnswerTests:** omnitest(correctExpr='qplot(y=hwy, data = mpg, color = drv)')
88     **Hint:** Type qplot(y=hwy, data = mpg, color = drv) at the command prompt.
89
90   - **Class:** cmd_question
91     **Output:** What's this plot showing? We see the x-axis ranges from 0 to 250 and we
       remember that we had 234 data points in our set, so we can infer that each point in
       the plot represents one of the hwy values (indicated by the y-axis). We've created
       the vector myhigh for you which contains the hwy data from the mpg dataset. Look at
       myhigh now.
92     **CorrectAnswer:** myhigh
93     **AnswerTests:** omnitest(correctExpr='myhigh')
94     **Hint:** Type myhigh at the command prompt.
95
96   - **Class:** text
97     **Output:** Comparing the values of myhigh with the plot, we see the first entries in the
       vector (29, 29, 31, 30,...) correspond to the leftmost points in the the plot (in
       order), and the last entries in myhigh (28, 29, 26, 26, 26) correspond to the
       rightmost plotted points. So, specifying the y parameter only, without an x argument,
       plots the values of the y argument in the order in which they occur in the data.
98
99   - **Class:** cmd_question
100    **Output:** The all-purpose qplot can also create box and whisker plots.  Call qplot now
       with 4 arguments. First specify the variable by which you'll split the data, in this
       case drv, then specify the variable which you want to examine, in this case hwy. The
       third argument is data (set equal to mpg), and the fourth, the geom, set equal to the
       string "boxplot"
101    **CorrectAnswer:** qplot(drv,hwy,data=mpg,geom="boxplot")
102    **AnswerTests:** omnitest(correctExpr='qplot(drv,hwy,data=mpg,geom="boxplot")')
103    **Hint:** Type qplot(drv,hwy,data=mpg,geom="boxplot") at the command prompt.
104
105  - **Class:** cmd_question
106    **Output:** We see 3 boxes, one for each drive. Now to impress you, call qplot with 5
       arguments. The first 4 are just as you used previously, (drv, hwy, data set equal to
       mpg, and geom set equal to the string "boxplot"). Now add a fifth argument, color,
       equal to manufacturer.
107    **CorrectAnswer:** qplot(drv,hwy,data=mpg,geom="boxplot",color=manufacturer)
108    **AnswerTests:**
       omnitest(correctExpr='qplot(drv,hwy,data=mpg,geom="boxplot",color=manufacturer)')
109    **Hint:** Type qplot(drv,hwy,data=mpg,geom="boxplot",color=manufacturer) at the command
       prompt.
110
111  - **Class:** text
112    **Output:** It's a little squished but we just wanted to illustrate qplot's capabilities.
       Notice that there are still 3 regions of the plot (determined by the factor drv).
       Each is subdivided into several boxes depicting different manufacturers.
113
114  - **Class:** cmd_question
115    **Output:** Now, on to histograms. These display frequency counts for a single variable.
       Let's start with an easy one. Call qplot with 3 arguments. First specify the variable
       for which you want the frequency count, in this case hwy, then specify the data (set
       equal to mpg), and finally, the aesthetic, fill, set equal to drv. Instead of a plain

old histogram, this will again use colors to distinguish the 3 different drive factors.
116    **CorrectAnswer**: qplot(hwy, data = mpg, fill = drv)
117    **AnswerTests**: omnitest(correctExpr='qplot(hwy, data = mpg, fill = drv)')
118    **Hint**: Type qplot(hwy, data = mpg, fill = drv) at the command prompt.
119
120    - **Class**: text
121    **Output**: See how qplot consistently uses the colors. Red (if 4-wheel drv is in the
       bin) is at the bottom of the bin, then green on top of it (if present), followed by
       blue (rear wheel drv). The color lets us see right away that 4-wheel drive vehicles
       in this dataset don't have gas mileages exceeding 30 miles per gallon.
122
123    - **Class**: text
124    **Output**: It's cool that qplot can do this  so easily, but some people may find this
       multi-color histogram hard to interpret. Instead of using colors to distinguish
       between the drive factors let's use facets or panels. (That's what lattice called
       them.) This just means we'll split the data into 3 subsets (according to drive) and
       make 3 smaller individual plots of each subset in one plot (and with one call to
       qplot).
125
126    - **Class**: text
127    **Output**: Remember that with base plot we had to do each subplot individually. The
       lattice system made plotting  conditioning plots  easier. Let's see how easy it is
       with qplot.
128
129    - **Class**: cmd_question
130    **Output**:  We'll do two plots, a scatterplot and then a histogram, each with 3 facets.
       For the scatterplot, call qplot with 4 arguments. The first two are displ and hwy and
       the third is the argument data set equal to mpg. The fourth is the argument facets
       which will be set equal to the expression . ~ drv which is ggplot2's shorthand for
       number of rows (to the left of the ~) and number of columns (to the right of the ~).
       Here the . indicates a single row and drv implies 3, since there are 3 distinct drive
       factors. Try this now.
131    **CorrectAnswer**: qplot(displ, hwy, data = mpg, facets = . ~ drv)
132    **AnswerTests**: omnitest(correctExpr='qplot(displ, hwy, data = mpg, facets = . ~ drv)')
133    **Hint**: Type qplot(displ, hwy, data = mpg, facets = . ~ drv) at the command prompt.
134
135    - **Class**: text
136    **Output**: The result is a 1 by 3 array of plots. Note how each is labeled at the top
       with the factor label (4,f, or r). This shows us more detailed information than the
       histogram. We see the relationship between displacement and highway mileage for each
       of the 3 drive factors.
137
138    - **Class**: cmd_question
139    **Output**:  Now we'll do a histogram, again calling qplot with 4 arguments. This time,
       since we need only one variable for a histogram, the first is hwy and the second is
       the argument data set equal to mpg. The third is the argument facets which we'll set
       equal to the expression drv ~ . . This will give us a different arrangement of the
       facets. The fourth argument is binwidth. Set this equal to 2. Try this now.
140    **CorrectAnswer**: qplot(hwy, data = mpg, facets = drv ~ ., binwidth = 2)
141    **AnswerTests**: omnitest(correctExpr='qplot(hwy, data = mpg, facets = drv ~ ., binwidth
       = 2)')
142    **Hint**: Type qplot(hwy, data = mpg, facets = drv ~ ., binwidth = 2) at the command
       prompt.
143
144    - **Class**: mult_question
145    **Output**: The facets argument, drv ~ ., resulted in what arrangement of facets?
146    **AnswerChoices**:  1 by 3; 3 by 1; 2 by 2; huh?
147    **CorrectAnswer**:  3 by 1
148    **AnswerTests**: omnitest(correctVal='3 by 1')
149    **Hint**: How many row? How many columns?
150
151
152    - **Class**: text
153    **Output**:  Pretty good, right? Not too difficult either. Let's review what we learned!
154
155    - **Class**: mult_question
156    **Output**: Which of the following is a basic workhorse function of ggplot2?
157    **AnswerChoices**:  hist; xyplot; scatterplot; gplot; qplot
158    **CorrectAnswer**:  qplot

```
159      AnswerTests: omnitest(correctVal='qplot')
160      Hint: Which function did we invoke the most in this lesson?
161
162    - Class: mult_question
163      Output:  Which types of plot does qplot plot?
164      AnswerChoices:  histograms; scatterplots; box and whisker plots; all of the others
165      CorrectAnswer: all of the others
166      AnswerTests: omnitest(correctVal='all of the others')
167      Hint: That qplot is amazing! It seems to do everything!
168
169    - Class: mult_question
170      Output:  What does the gg in ggplot2 stand for?
171      AnswerChoices:  good grief; grammar of graphics; goto graphics; good graphics
172      CorrectAnswer: grammar of graphics
173      AnswerTests: omnitest(correctVal='grammar of graphics')
174      Hint: Think of building blocks and components.
175
176    - Class: mult_question
177      Output:  True or False? The geom argument takes a string for a value.
178      AnswerChoices:  True; False
179      CorrectAnswer: True
180      AnswerTests: omnitest(correctVal='True')
181      Hint: Recall our examples, for instance, geom="density".
182
183    - Class: mult_question
184      Output:  True or False? The data argument takes a string for a value.
185      AnswerChoices:  True; False
186      CorrectAnswer: True
187      AnswerTests: omnitest(correctVal='False')
188      Hint: Recall our examples. Did we ever put the dataset name in quotation marks?
189
190    - Class: mult_question
191      Output:  True or False? The binwidth argument takes a string for a value.
192      AnswerChoices:  True; False
193      CorrectAnswer: False
194      AnswerTests: omnitest(correctVal='False')
195      Hint: Recall our examples, for instance, binwidth=18497/30.
196
197    - Class: mult_question
198      Output:  True or False? The user must specify x- and y-axis labels when using qplot.
199      AnswerChoices:  True; False
200      CorrectAnswer: False
201      AnswerTests: omnitest(correctVal='False')
202      Hint: Recall our examples when we saw labels that we didn't specify.
203
204    - Class: text
205      Output: Congrats! You've finished plot 1 of  ggplot2. In the next lesson the plot
         thickens.
206
207    - Class: mult_question
208      Output: "Would you like to receive credit for completing this course on
209        Coursera.org?"
210      CorrectAnswer: NULL
211      AnswerChoices: Yes;No
212      AnswerTests: coursera_on_demand()
213      Hint: ""
214
```