

Solutions to G. Grolemund & H. Wickhams's R for Data Science

Krista DeStasio

7/11/2017

Contents

Chapter 3, Data Visualisation	2
The mpg data frame	2
Making graphs with ggplot2	3
Exercises 3.2.4	4
1. Run ggplot(data = mpg). What do you see?	4
2. How many rows are in mpg? How many columns?	4
3. What does the drv variable describe? Read the help for ?mpg to find out.	4
4. Make a scatterplot of hwy vs cyl.	5
5. What happens if you make a scatterplot of class vs drv? Why is the plot not useful?	5
Aesthetic mappings	6
Exercises 3.3.1	11
1. What's gone wrong with this code? Why are the points not blue?	11
2. Which variables in mpg are categorical? Which variables are continuous? (Hint: type ?mpg to read the documentation for the dataset). How can you see this information when you run mpg?	13
3. Map a continuous variable to color, size, and shape. How do these aesthetics behave differently for categorical vs. continuous variables?	14
4. What happens if you map the same variable to multiple aesthetics?	15
5. What does the stroke aesthetic do? What shapes does it work with? (Hint: use ?geom_point)	16
6. What happens if you map an aesthetic to something other than a variable name, like aes(colour = displ < 5)?	18
Facets	19
Exercises 3.5.1	22
1. What happens if you facet on a continuous variable?	22
2. What do the empty cells in plot with facet_grid(drv ~ cyl) mean? How do they relate to this plot?	23
3. What plots does the following code make? What does . do?	25
4. Take the first faceted plot in this section. What are the advantages to using faceting instead of the colour aesthetic? What are the disadvantages? How might the balance change if you had a larger dataset?	27
5. Read ?facet_wrap. What does nrow do? What does ncol do? What other options control the layout of the individual panels? Why doesn't facet_grid() have nrow and ncol argument?	29
6. When using facet_grid() you should usually put the variable with more unique levels in the columns. Why?	29
Geometric objects	29
Exercises 3.6.1	39
1. What geom would you use to draw a line chart? A boxplot? A histogram? An area chart?	39
2. Run this code in your head and predict what the output will look like. Then, run the code in R and check your predictions.	39

3. What does `show.legend = FALSE` do? What happens if you remove it? Why do you think I used it earlier in the chapter? 40
4. What does the `se` argument to `geom_smooth()` do? 40
5. Will these two graphs look different? Why/why not? 40
6. Recreate the R code necessary to generate the following graphs. 42

Note: *This is a work in progress.*

A Brief Introduction to This File This R file walks through G. Grolemund & H. Wickham's online text, "R for Data Science." Much of the code is sourced directly from the book and credit belongs to the authors. Here, some sections of code are heavily commented so that the beginning R programmer can read through and understand what each line of code does and compare it to their own as they work through the text. Throughout, the book provides the primary and most thorough explanation. **For the greatest learning benefit, I suggest you attempt each exercise on your own before looking at the code or write-ups provided here.** Of course, there is more than one way to write code and you may find a more elegant solution that you prefer.

For those new to R and RStudio, it may be of additional benefit to knit the document and examine how the code in the Rmd file is visually expressed in the resultant knitted document. For example, see how the ["R for Data Science."](<http://r4ds.had.co.nz/index.html>) is expressed as a hyperlink in the preceeding paragraph where it was not surrounded by tick-marks and compare that to how the same text is expressed in this paragraph when surrounded by ticks. See also the difference in appearance when knitting to different document types (HTML, PDF, Word).

Tip: *If you are using RStudio, click the text next to the orange # box at the bottom of the editor window to easily navigate the code chunks.*

Tip: *Use the ? before any command to view the documentation on that function. Do this often. For example, type `?setwd` to see a description, usage, arguments, and more for the function `setwd()`.*

Tip: Find RStudio Cheatsheets at <https://www.rstudio.com/resources/cheatsheets/>

Chapter 3, Data Visualisation

To really understand ggplot2, I highly recommend reading "The Layered Grammar of Graphics" as suggested at the beginning of Chapter 3.

The mpg data frame

```
str(mpg) # Look at the structure of the mpg data frame

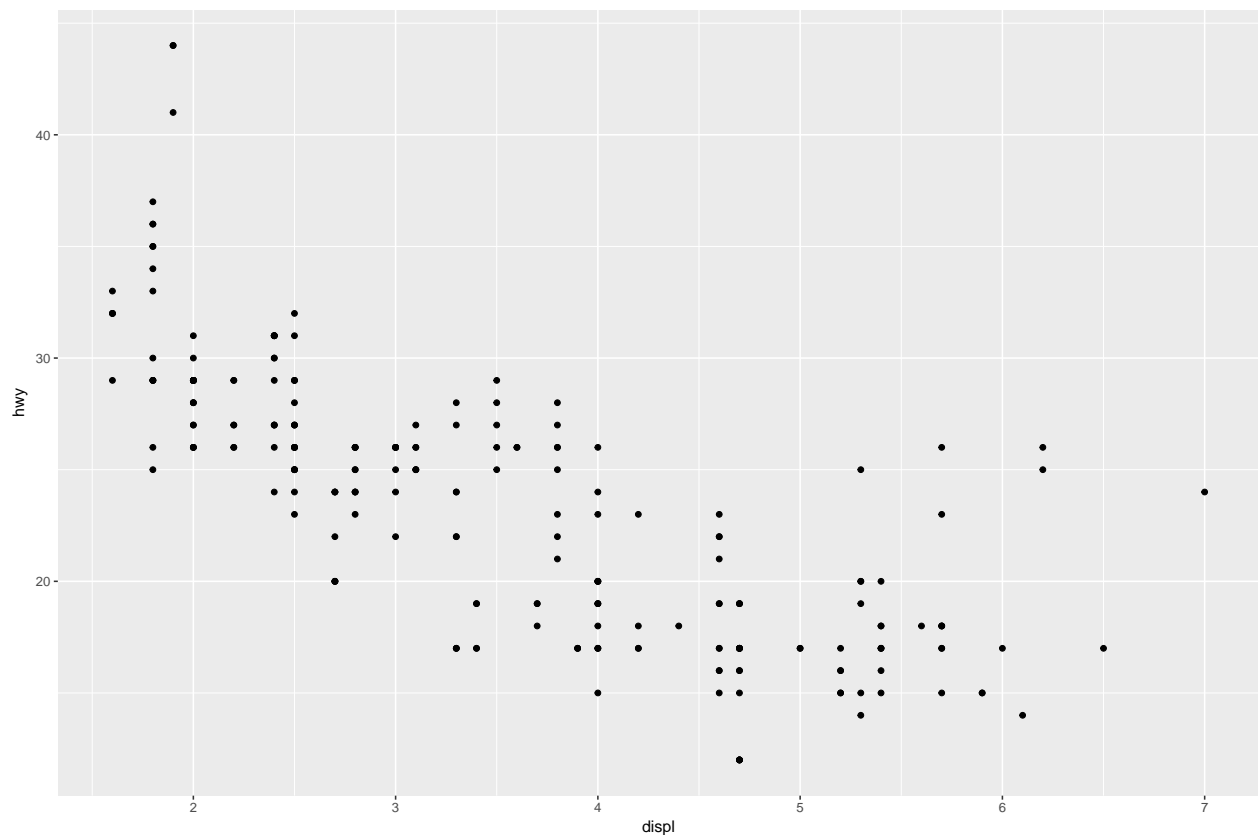
## Classes 'tbl_df', 'tbl' and 'data.frame':   234 obs. of  11 variables:
## $ manufacturer: chr  "audi" "audi" "audi" "audi" ...
## $ model       : chr  "a4" "a4" "a4" "a4" ...
## $ displ       : num  1.8 1.8 2 2 2.8 2.8 3.1 1.8 1.8 2 ...
## $ year        : int  1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 ...
## $ cyl         : int  4 4 4 4 6 6 6 4 4 4 ...
## $ trans       : chr  "auto(l5)" "manual(m5)" "manual(m6)" "auto(av)" ...
## $ drv         : chr  "f" "f" "f" "f" ...
## $ cty         : int  18 21 20 21 16 18 18 16 20 ...
## $ hwy         : int  29 29 31 30 26 26 27 26 25 28 ...
## $ fl         : chr  "p" "p" "p" "p" ...
## $ class       : chr  "compact" "compact" "compact" "compact" ...
```

```
mpg # Look at the first 10 rows of the mpg data frame
```

```
## # A tibble: 234 x 11
##   manufacturer    model displ  year  cyl    trans  drv   cty   hwy
##         <chr>      <chr> <dbl> <int> <int>    <chr> <chr> <int> <int>
## 1      audi        a4    1.8  1999    4  auto(l5)  f     18    29
## 2      audi        a4    1.8  1999    4 manual(m5)  f     21    29
## 3      audi        a4    2.0  2008    4 manual(m6)  f     20    31
## 4      audi        a4    2.0  2008    4  auto(av)   f     21    30
## 5      audi        a4    2.8  1999    6  auto(l5)  f     16    26
## 6      audi        a4    2.8  1999    6 manual(m5)  f     18    26
## 7      audi        a4    3.1  2008    6  auto(av)   f     18    27
## 8      audi  a4 quattro  1.8  1999    4 manual(m5)  f     18    26
## 9      audi  a4 quattro  1.8  1999    4  auto(l5)   f     16    25
## 10     audi  a4 quattro  2.0  2008    4 manual(m6)  f     20    28
## # ... with 224 more rows, and 2 more variables: fl <chr>, class <chr>
```

Hypothesis: There is a negative linear relationship between engine size and fuel efficiency, such that as engine size increases fuel efficiency decreases.

```
ggplot(data=mpg) + # specify data frame
  geom_point(mapping = aes(x = displ, y = hwy)) # specify that plot is a scatterplot with displ on the x-axis and hwy on the y-axis
```



The plot confirms the hypothesis that there is a negative relationship between engine size and fuel efficiency.

Making graphs with ggplot2

Template:

```
ggplot(data = <DATA>) +  
  <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))
```

Exercises 3.2.4

1. Run `ggplot(data = mpg)`. What do you see?

There are no visible results from the code below.

```
ggplot(data = mpg)
```

2. How many rows are in `mpg`? How many columns?

Based on the output from `str(mpg)`, we see that there are 234 rows and 11 columns in the `mpg` data frame.

```
# Alternative means of finding number of rows and columns  
nrow(mpg) # Print the number of rows
```

```
## [1] 234
```

```
ncol(mpg)
```

```
## [1] 11
```

There are 234 rows and 11 columns in the `mpg` data frame.

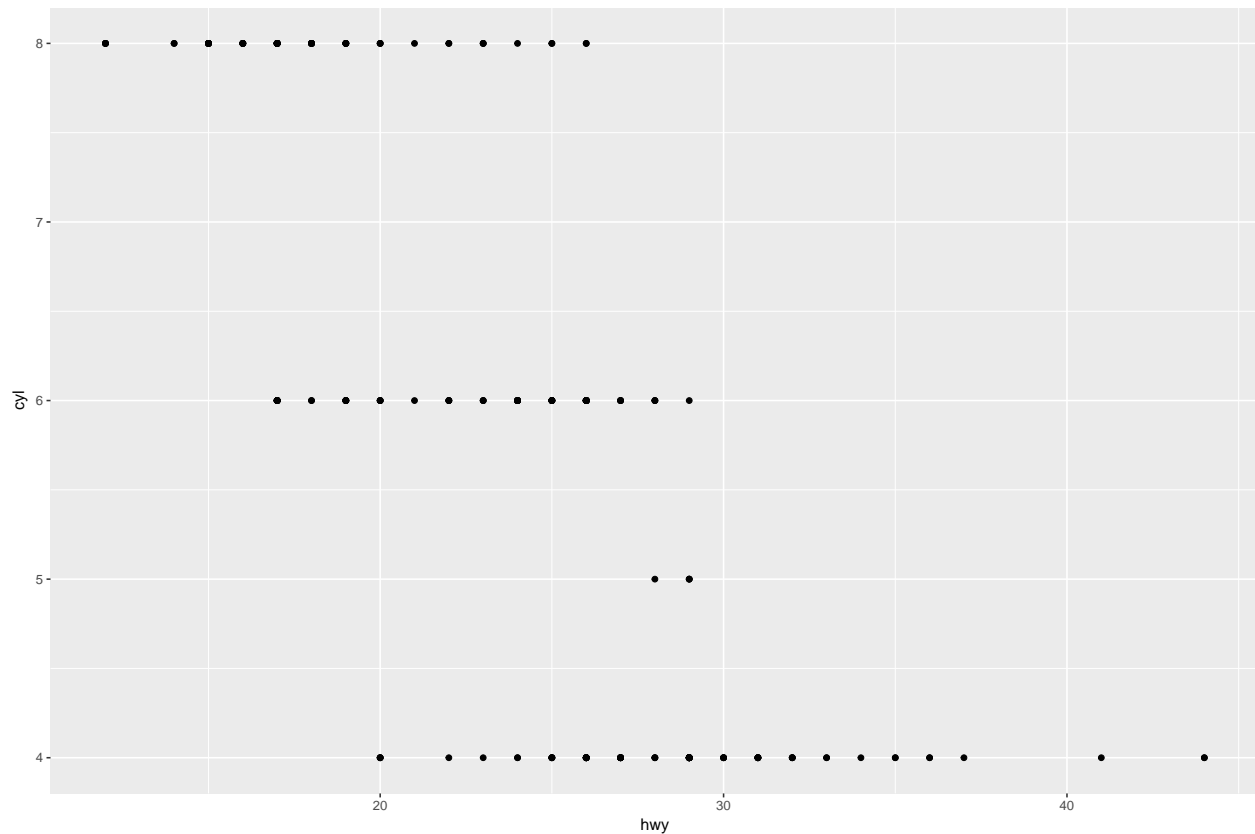
3. What does the `drv` variable describe? Read the help for `?mpg` to find out.

The `drv` variable describes whether the vehicle is front, rear, or 4-wheel drive.

```
?mpg
```

4. Make a scatterplot of `hwy` vs `cyl`.

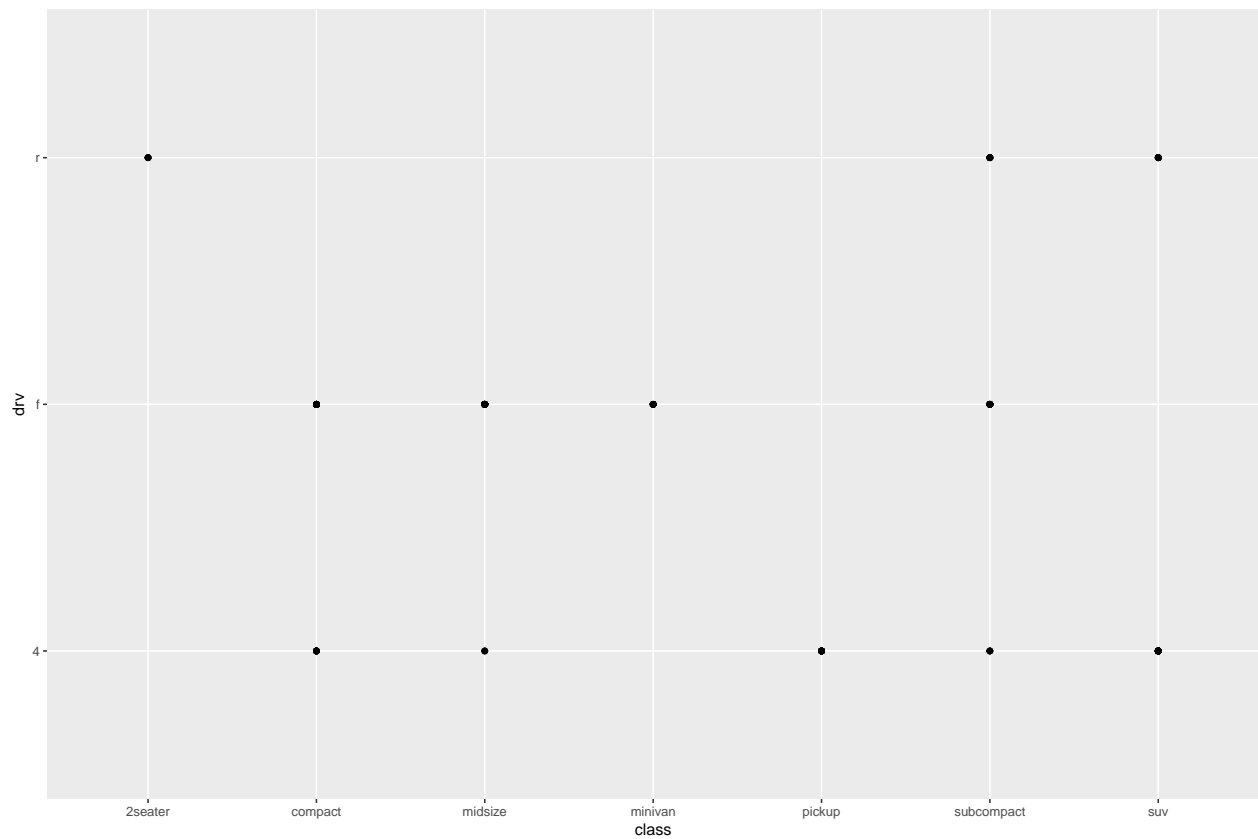
```
ggplot(data=mpg) +  
  geom_point(mapping = aes(x=hwy, y=cyl))
```



5. What happens if you make a scatterplot of `class` vs `drv`? Why is the plot not useful?

The plot is not useful because the variables are categorical and multiple points are plotted atop one another. We are unable to determine from this plot how many observations there are of each class-drive combination.

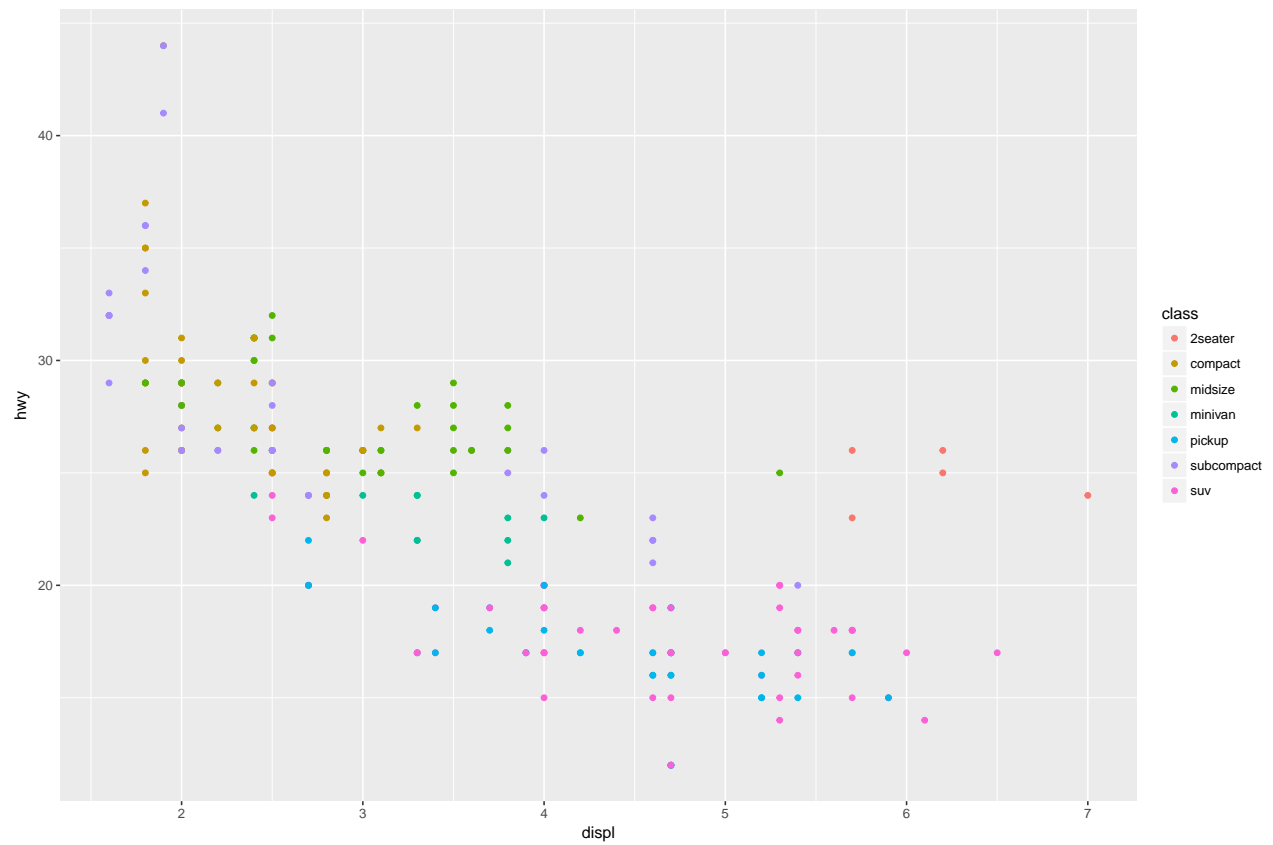
```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x=class, y=drv))
```



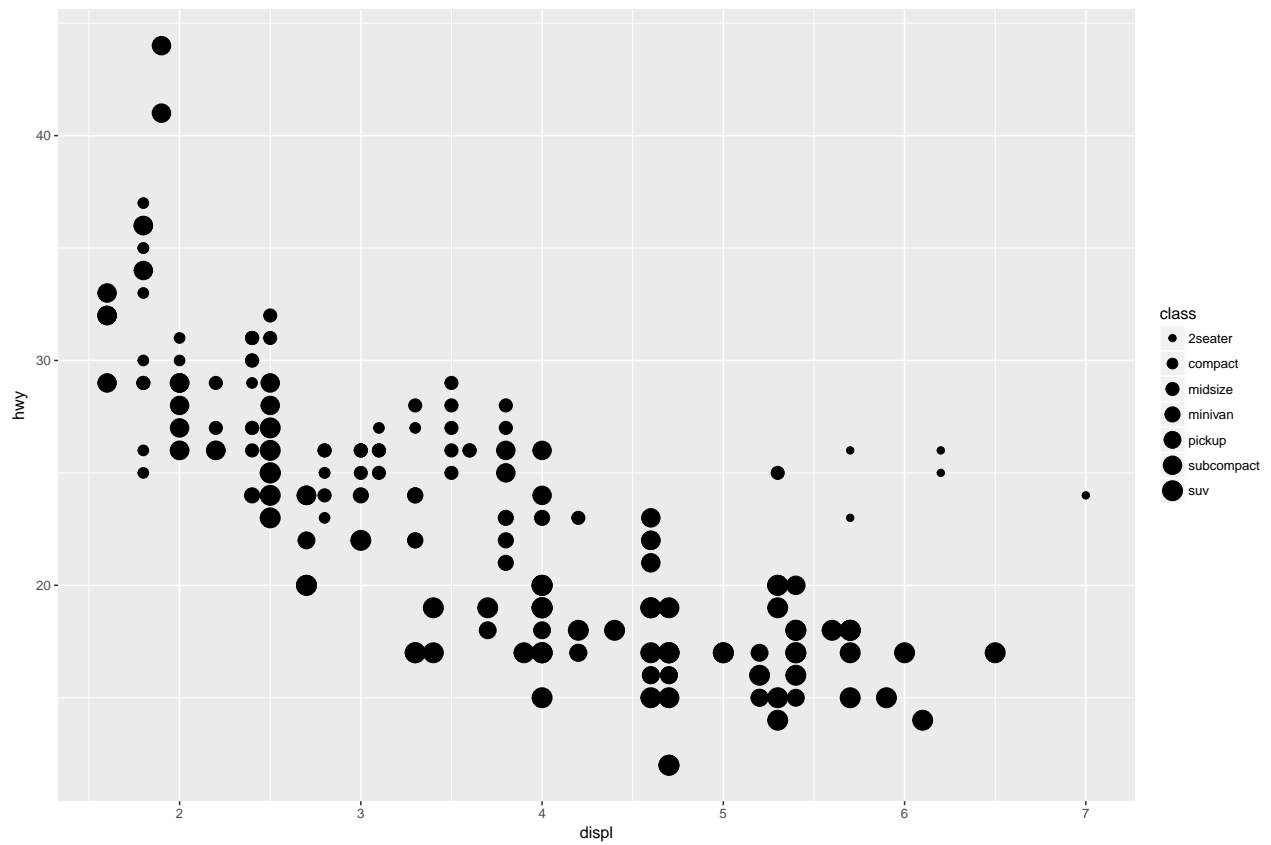
Aesthetic mappings

Test the hypothesis that the cars highlighted in red are hybrids by mapping car class to an aesthetic.

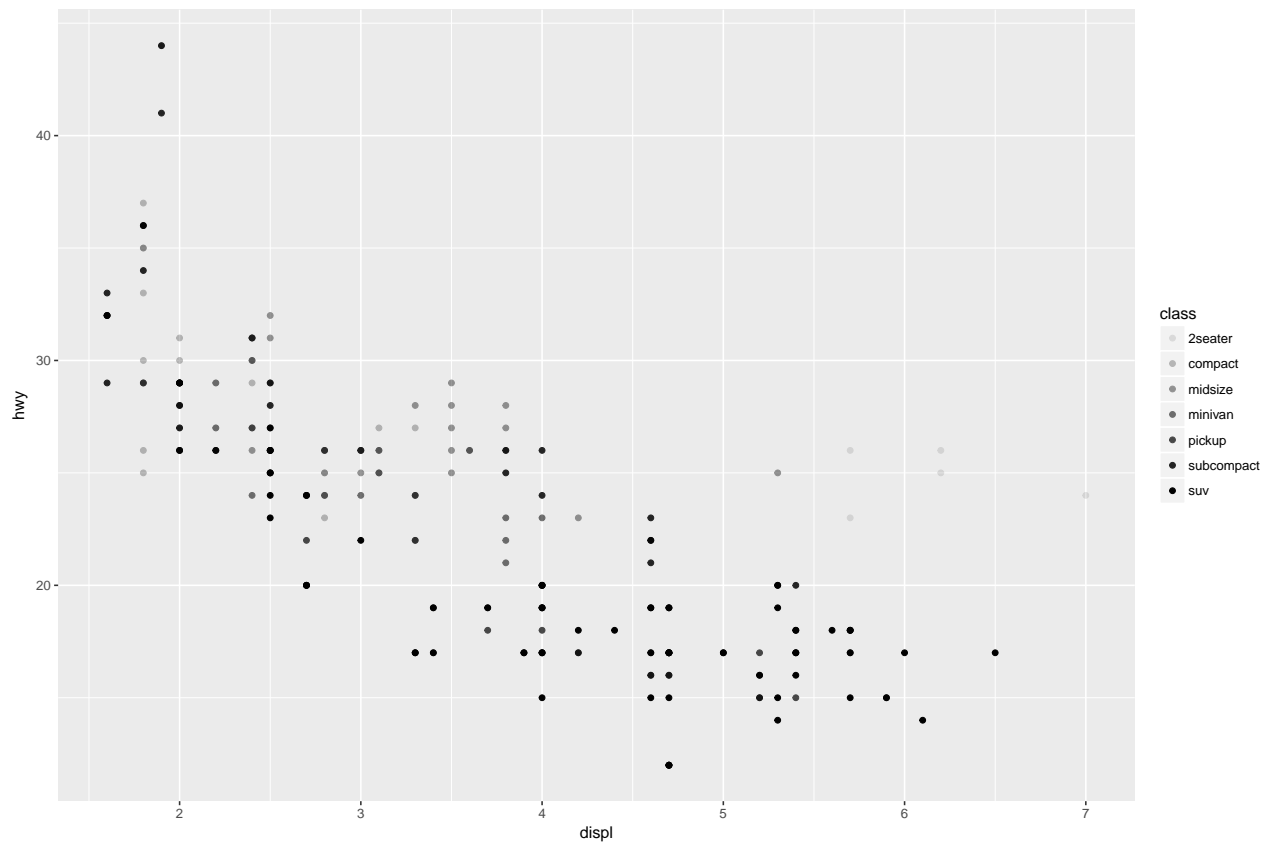
```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, color = class)) # map class to the color aesthetic so th
```



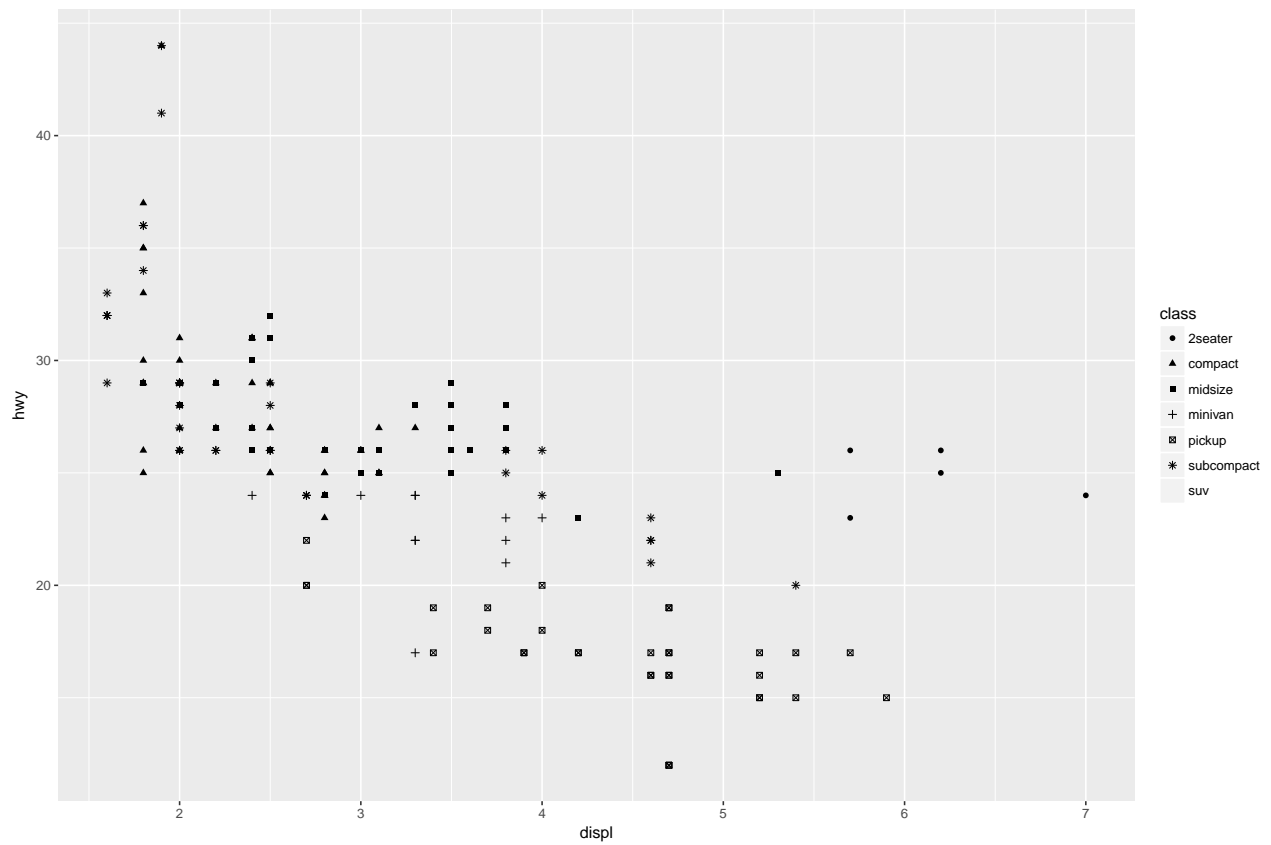
```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, size = class))
```



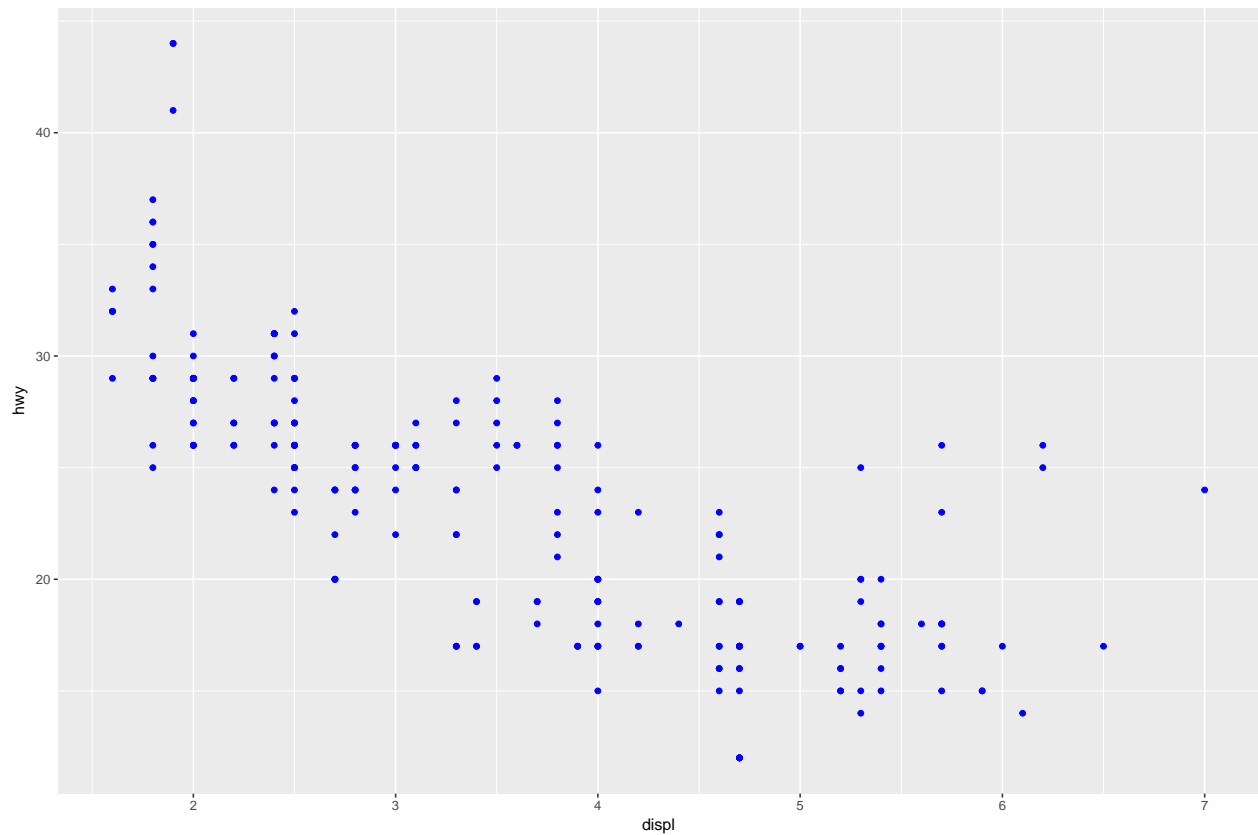
```
# Left  
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, alpha = class))
```

```
# Right
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, shape = class))
```



```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy), color = "blue") # Set the aesthetic outside of aes() to
```



Aesthetic shapes:

□ 0	✕ 4	⊕ 10	■ 15	■ 22
○ 1	▽ 6	⊗ 11	● 16	● 21
△ 2	⊠ 7	⊞ 12	▲ 17	▲ 24
◇ 5	✳ 8	⊗ 13	◆ 18	◆ 23
⊕ 3	⊞ 9	⊠ 14	● 19	● 20

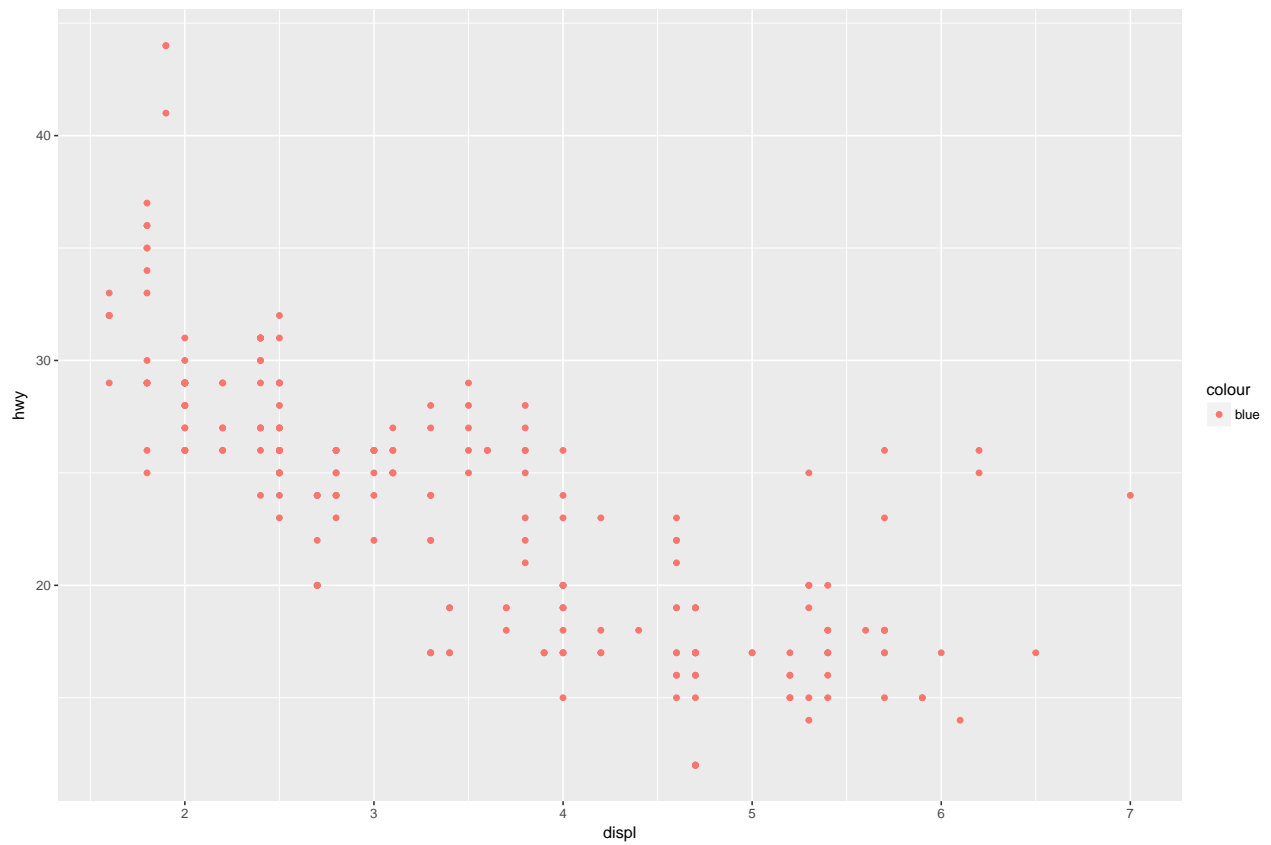
Figure 1:

Exercises 3.3.1

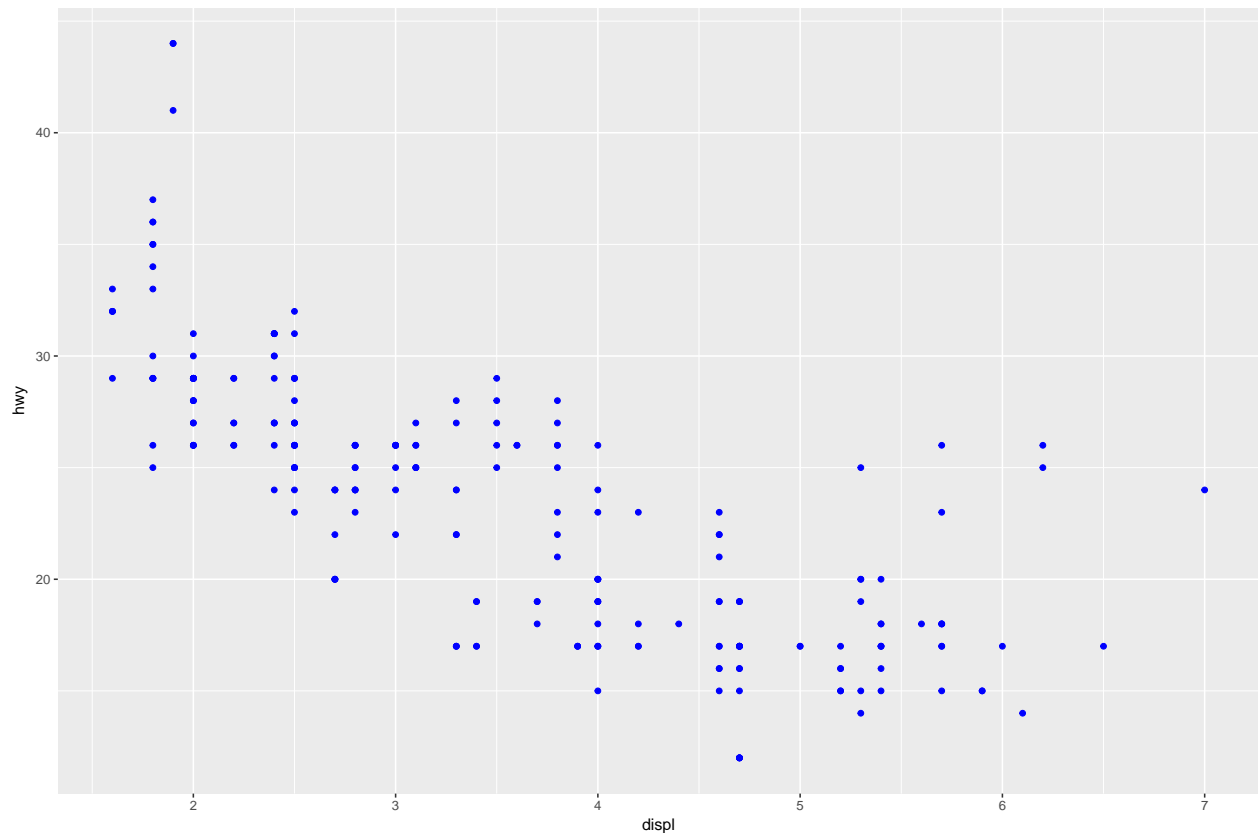
1. What's gone wrong with this code? Why are the points not blue?

The points are not blue because the color aesthetic is set inside `aes()`.

```
# Problematic code
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, color = "blue"))
```



```
# Corrected code  
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy), color = "blue")
```



2. Which variables in mpg are categorical? Which variables are continuous? (Hint: type `?mpg` to read the documentation for the dataset). How can you see this information when you run `mpg`?

To determine what the categorical and continuous variables are, one can either view the tibble by typing `mpg` or by viewing the documentation `?mpg`. One may decide whether a variable is categorical or continuous by checking whether it is stored as a character, integer, or double (floating point integer) value. However, this can lead to miscategorization in some cases. For example, while `year` is an integer, it is typically considered a whole number, a discrete variable without a meaningful 0 value anchor, and therefore not continuous.

The categorical variables are:

- `model`
- `year` (discrete, rather than categorical)
- `trans`
- `drv`
- `fl`
- `class`

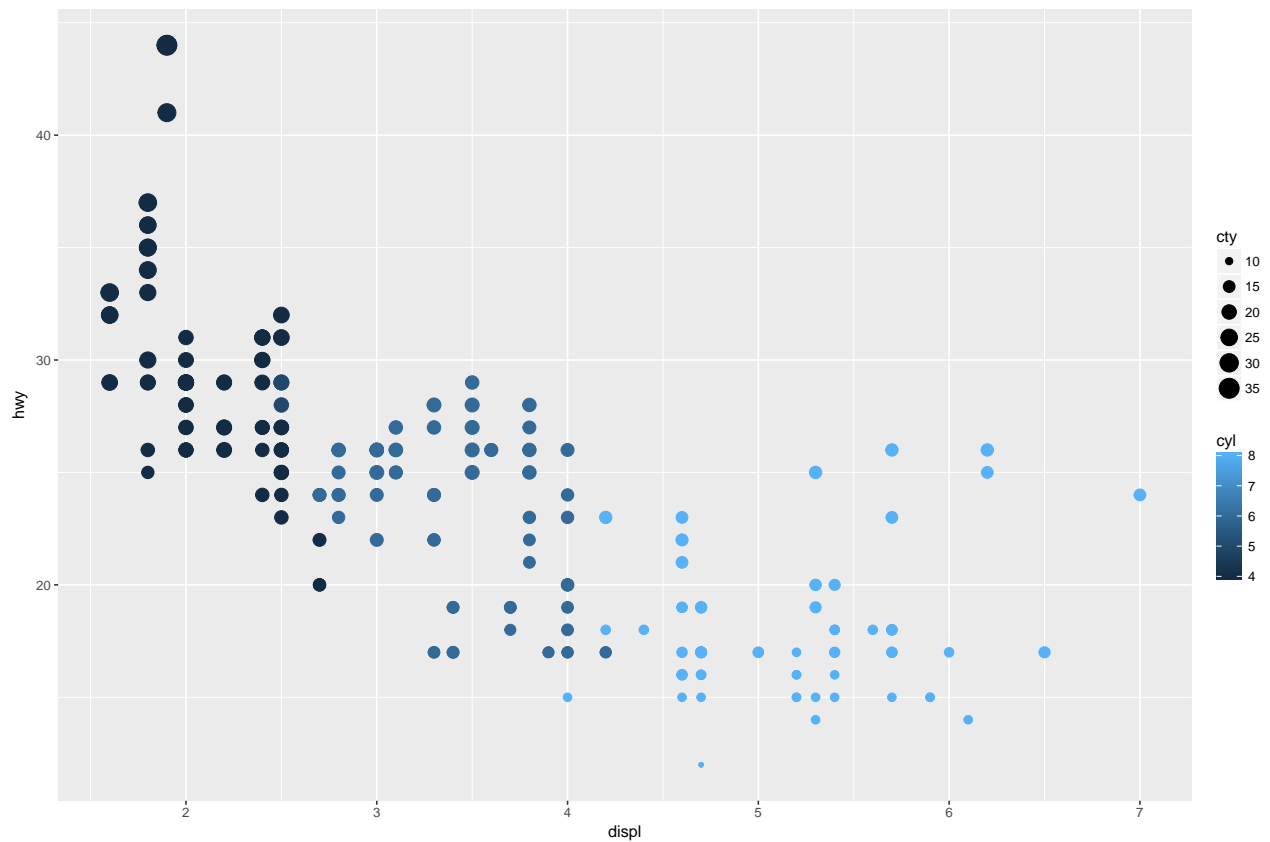
The continuous variables are:

- `displ`
- `cyl`
- `cty`
- `hwy`
- `year` (in this data set, `year` is treated as an integer variable. Better to consider this “quantitative”, rather than “continuous”)

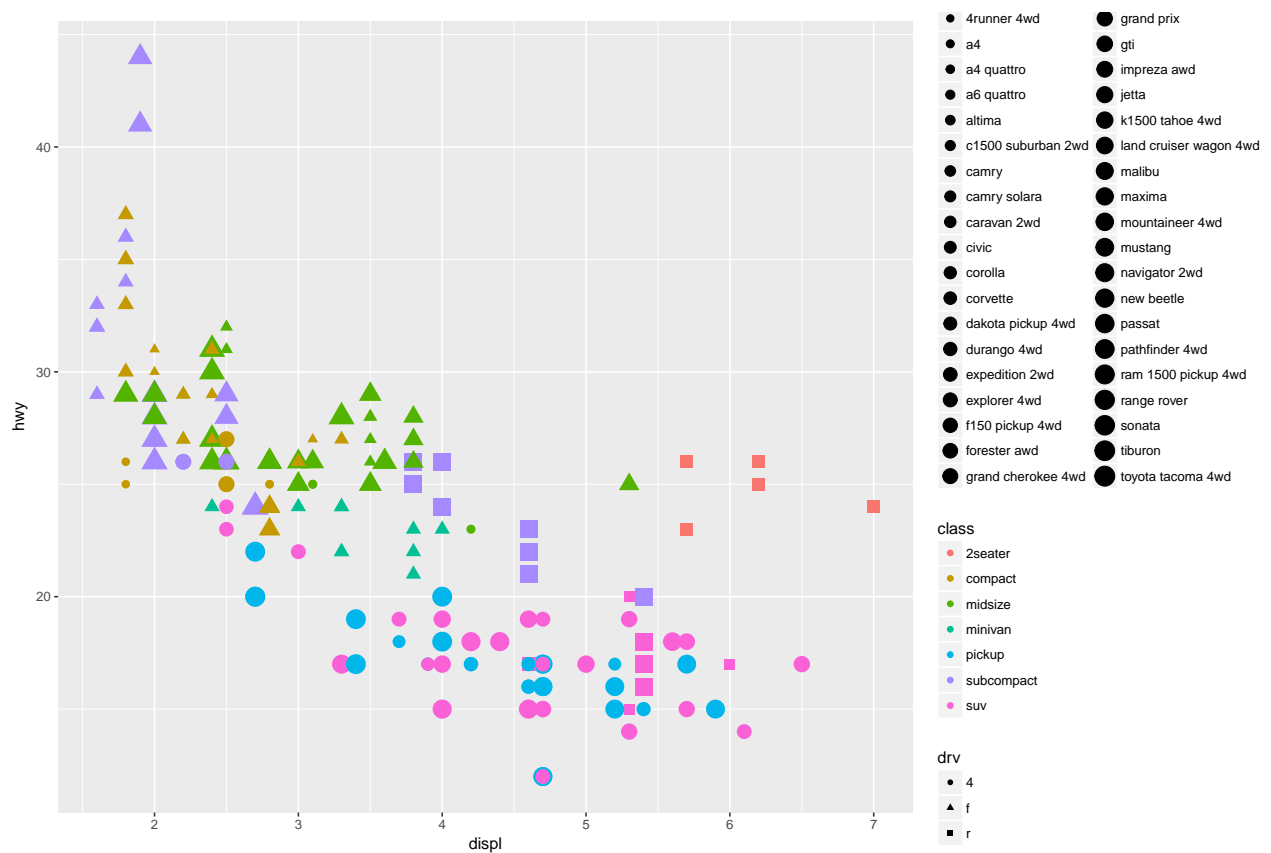
3. Map a continuous variable to color, size, and shape. How do these aesthetics behave differently for categorical vs. continuous variables?

A continuous variable cannot be mapped to shape. When mapped to size or color, the continuous variable is binned by equal intervals (in this case, intervals of 5 mpg). When mapped to the size aesthetic, points scale by the intervals. Continuous variables when mapped to a color aesthetic are mapped along a gradient scale.

```
# Mapping a continuous variable to the shape aesthetic  
ggplot(data=mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, shape = cty))  
  
# Mapping continuous variables to the color and size aesthetics  
ggplot(data=mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, color = cyl, size = cty))
```



```
# Mapping categorical variables to size, color, and shape  
ggplot(data=mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, size = model, color = class, shape = drv))
```



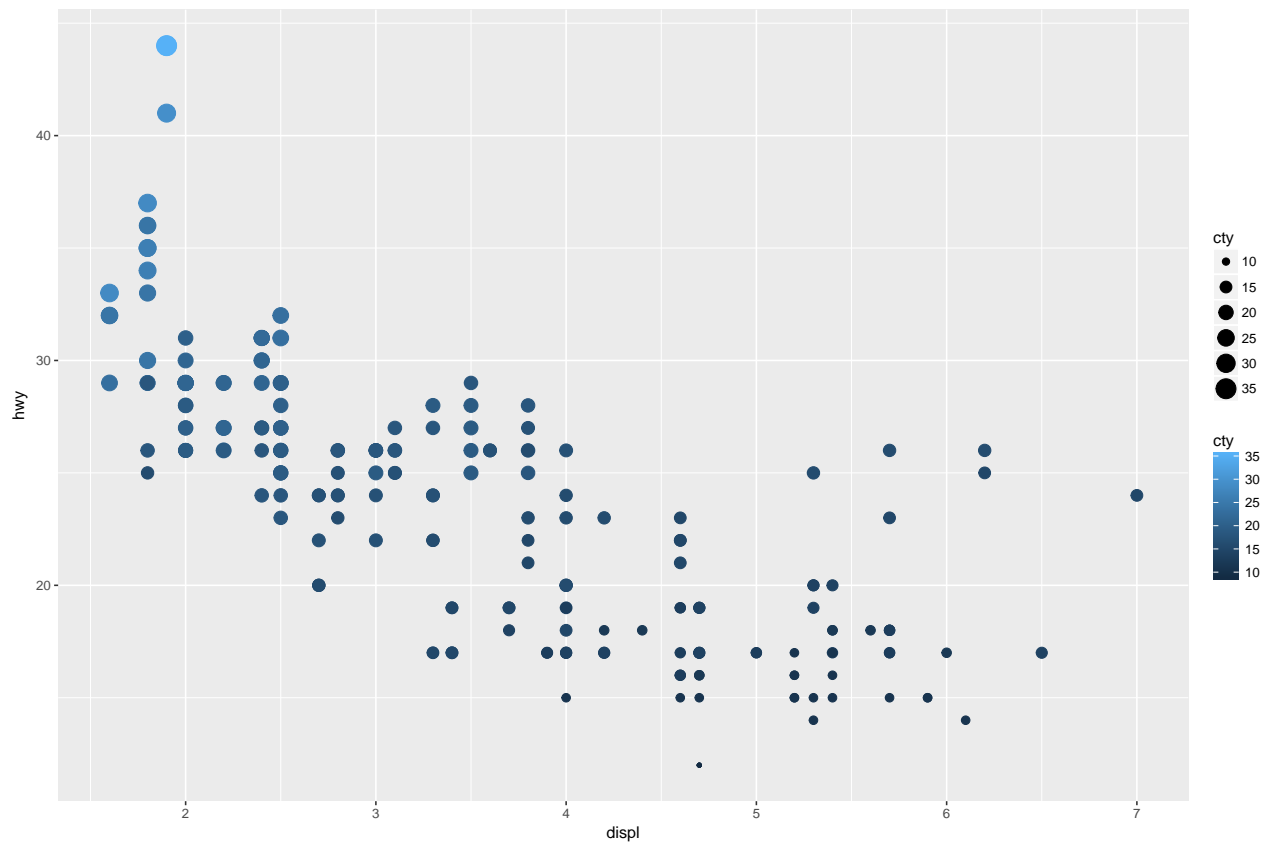
4. What happens if you map the same variable to multiple aesthetics?

When the same variable is mapped to multiple aesthetics, it is represented by those aesthetics.

Mapping the same variable to multiple aesthetics

`ggplot(data=mpg) +`

`geom_point(mapping = aes(x = displ, y = hwy, color = cty, size = cty))` *# Here, city is mapped to th*



5. What does the **stroke** aesthetic do? What shapes does it work with? (Hint: use `?geom_point`)

According to the R documentation:

“For shapes that have a border (like 21), you can colour the inside and outside separately. Use the **stroke** aesthetic to modify the width of the border.”

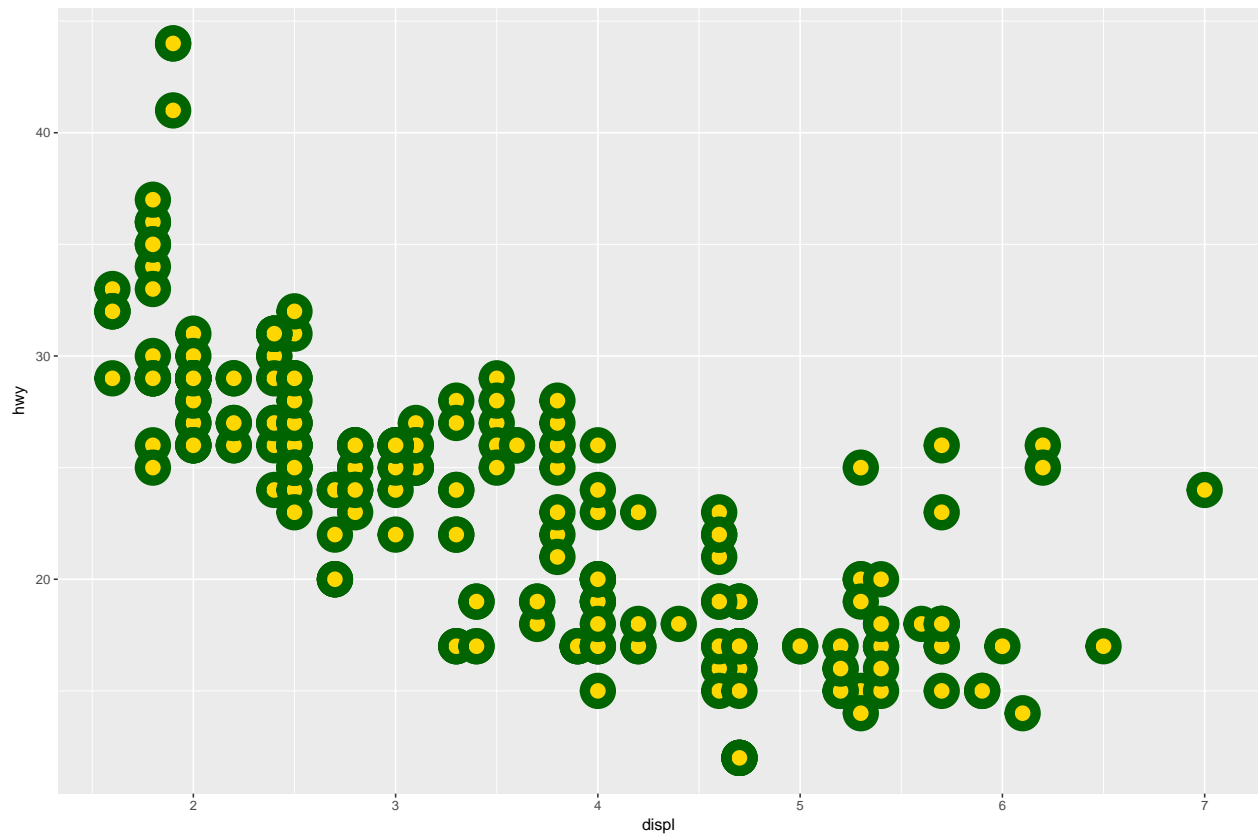
Tip: You can find documentation of available colors [here](#).

`?geom_point`

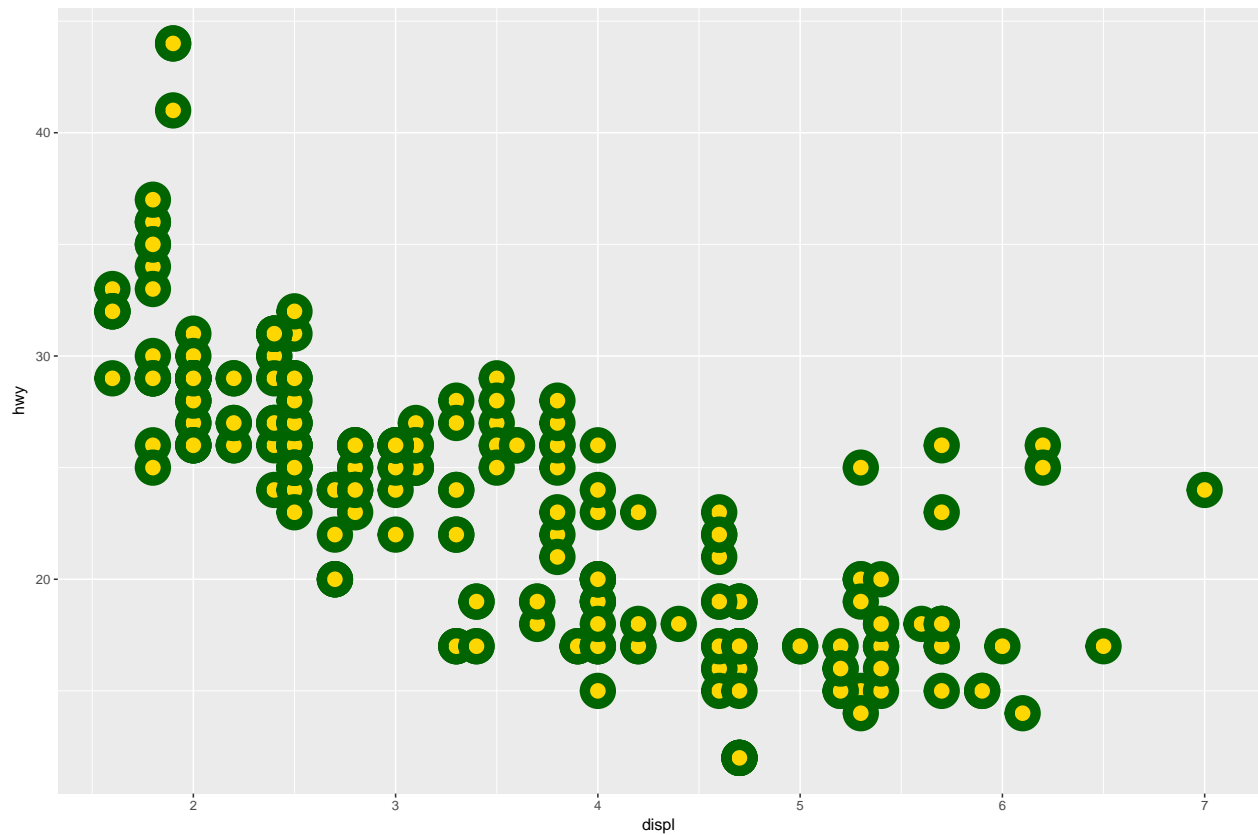
Example using `stroke`

`ggplot(data=mpg)+`

`geom_point(mapping = aes(x=displ, y=hwy), shape = 21, colour = "darkgreen", fill = "gold", size = 5`



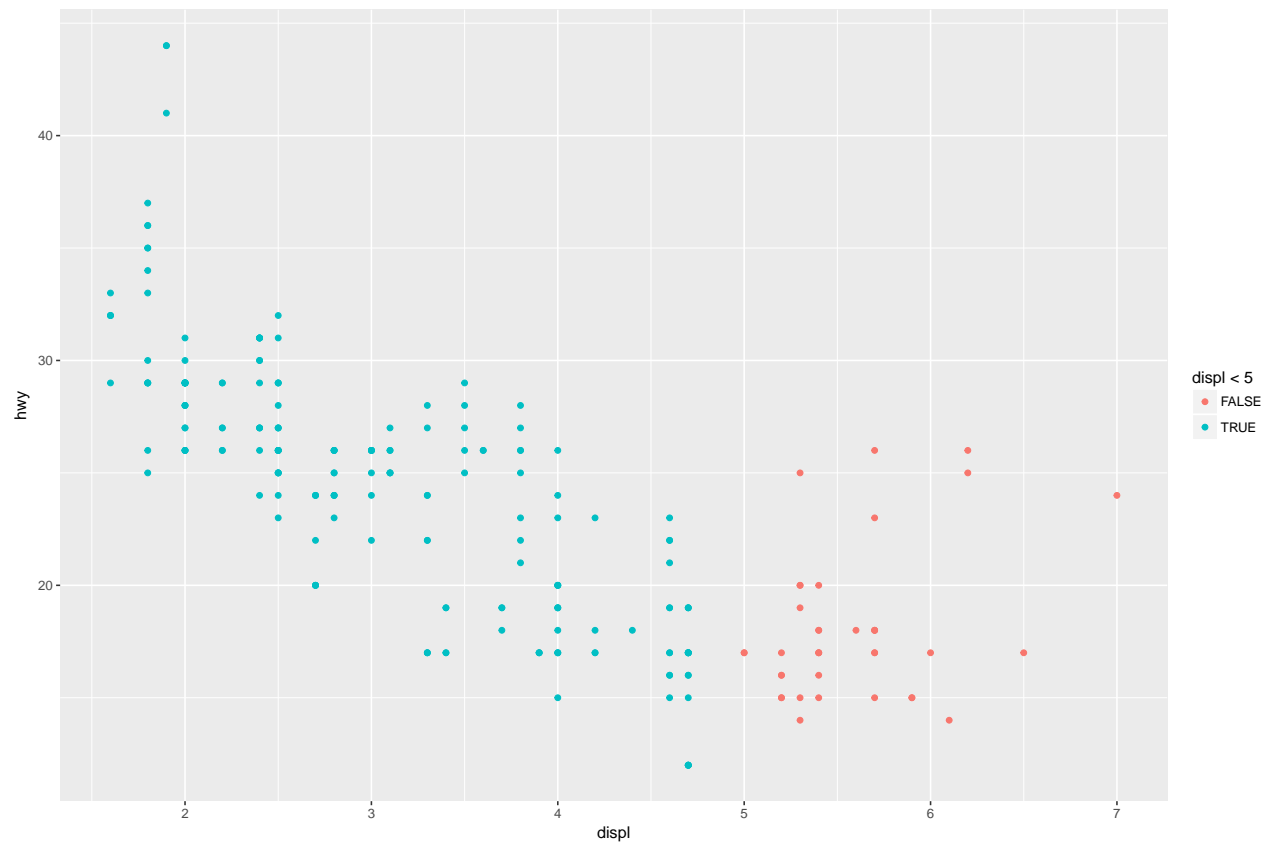
```
# Just for fun, let's write short-hand code make the same plot  
ggplot(mpg, aes(displ, hwy)) +  
  geom_point(shape = 21, colour = "darkgreen", fill = "gold", size = 5, stroke = 5)
```



6. What happens if you map an aesthetic to something other than a variable name, like `aes(colour = displ < 5)`?

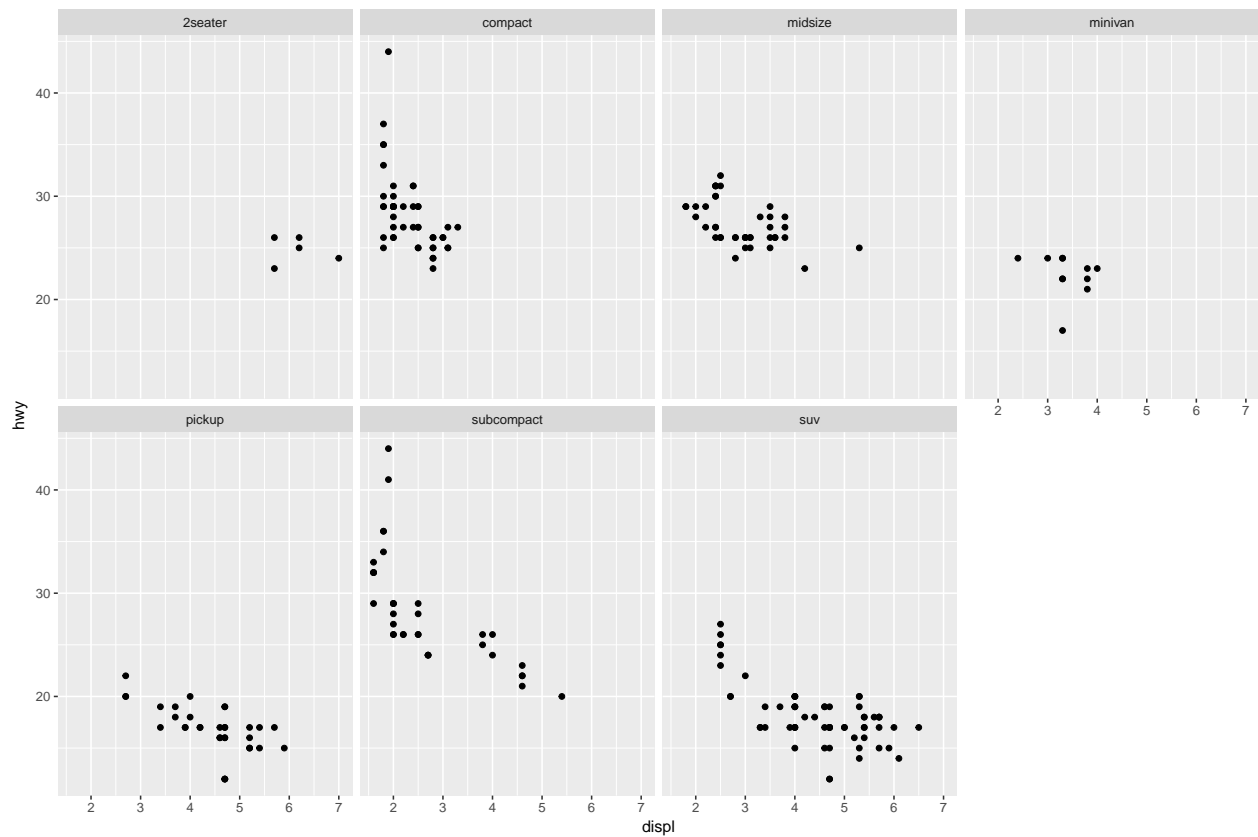
Setting the color aesthetic to `displ < 5` will assign one color to all x-axis (hwy) values < 5 and a different color to x-axis values ≥ 5 . Since the color palette is not specified, default colors are used.

```
ggplot(data=mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, color = displ < 5))
```

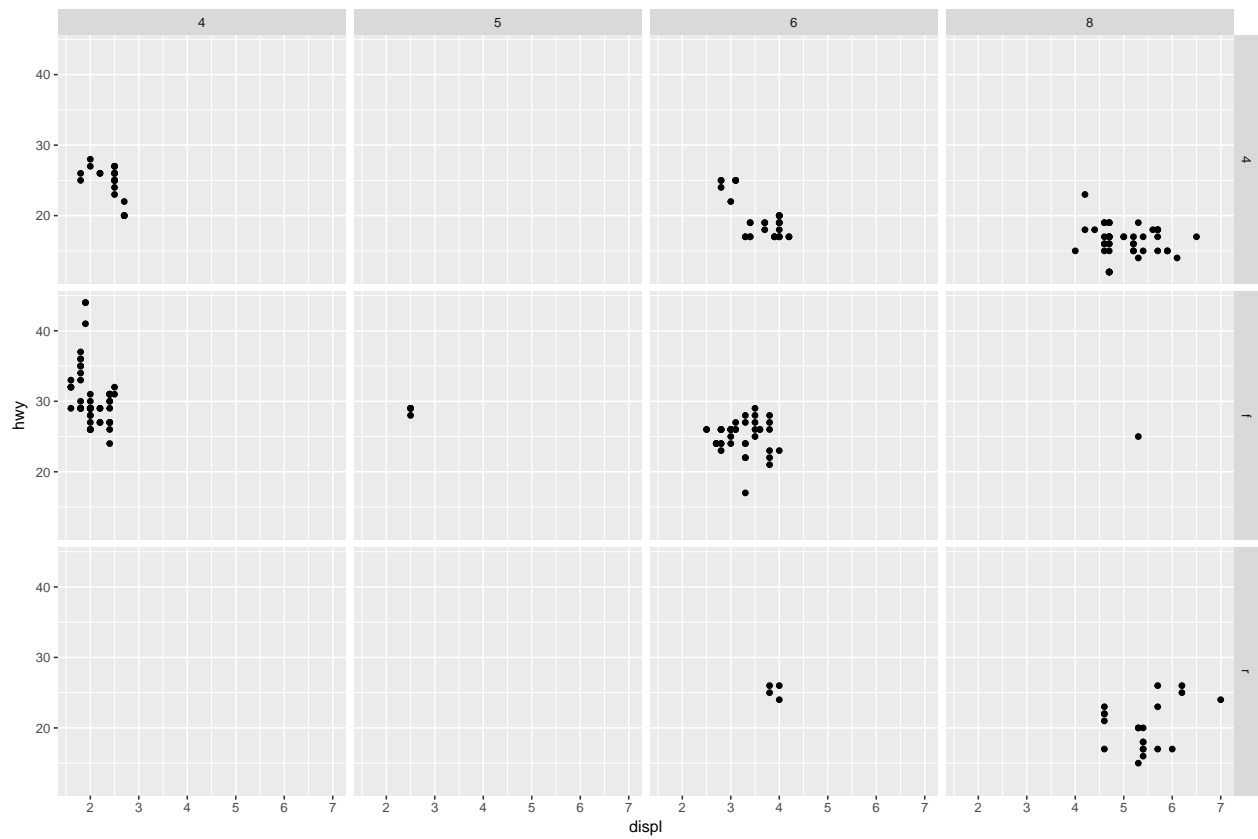


Facets

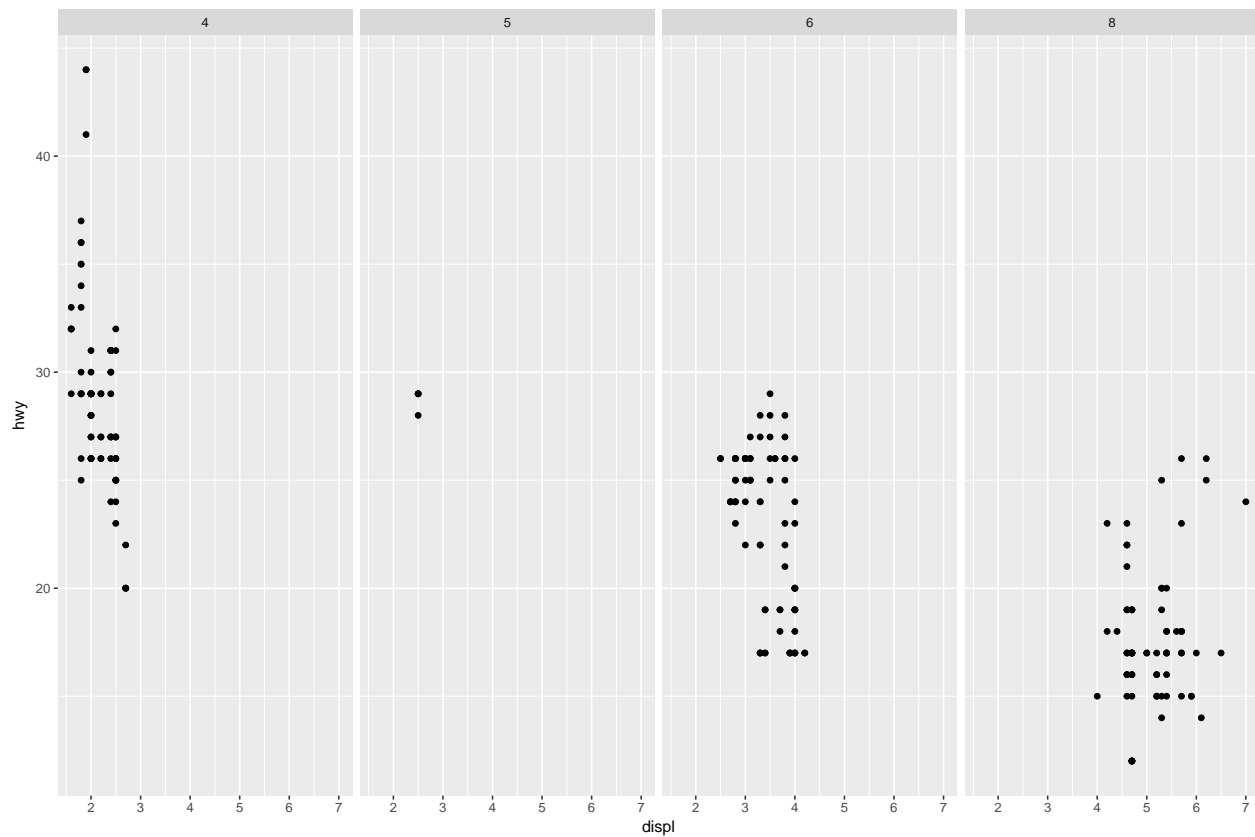
```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_wrap(~ class, nrow = 2) # This will create a separate plot for each class of vehicle and will f
```



```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(drv ~ cyl) # This will create a grid of plots with one plot for each combination of drv and cyl
```



```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(. ~ cyl) # Use the . to create plots for each level of cylinder (cyl) in the columns dim
```

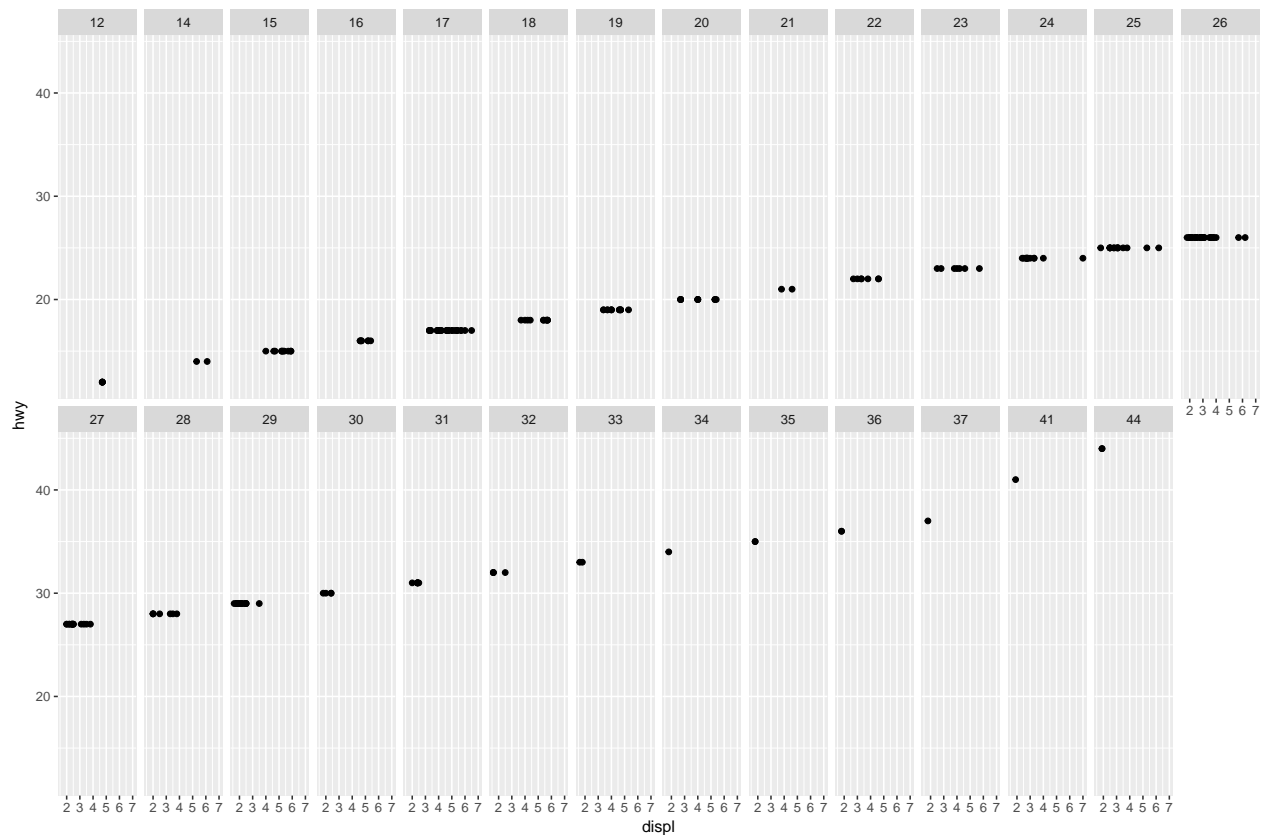


Exercises 3.5.1

1. What happens if you facet on a continuous variable?

If faceting is done with a continuous variable, a plot is created for each value for which there is at least one observation.

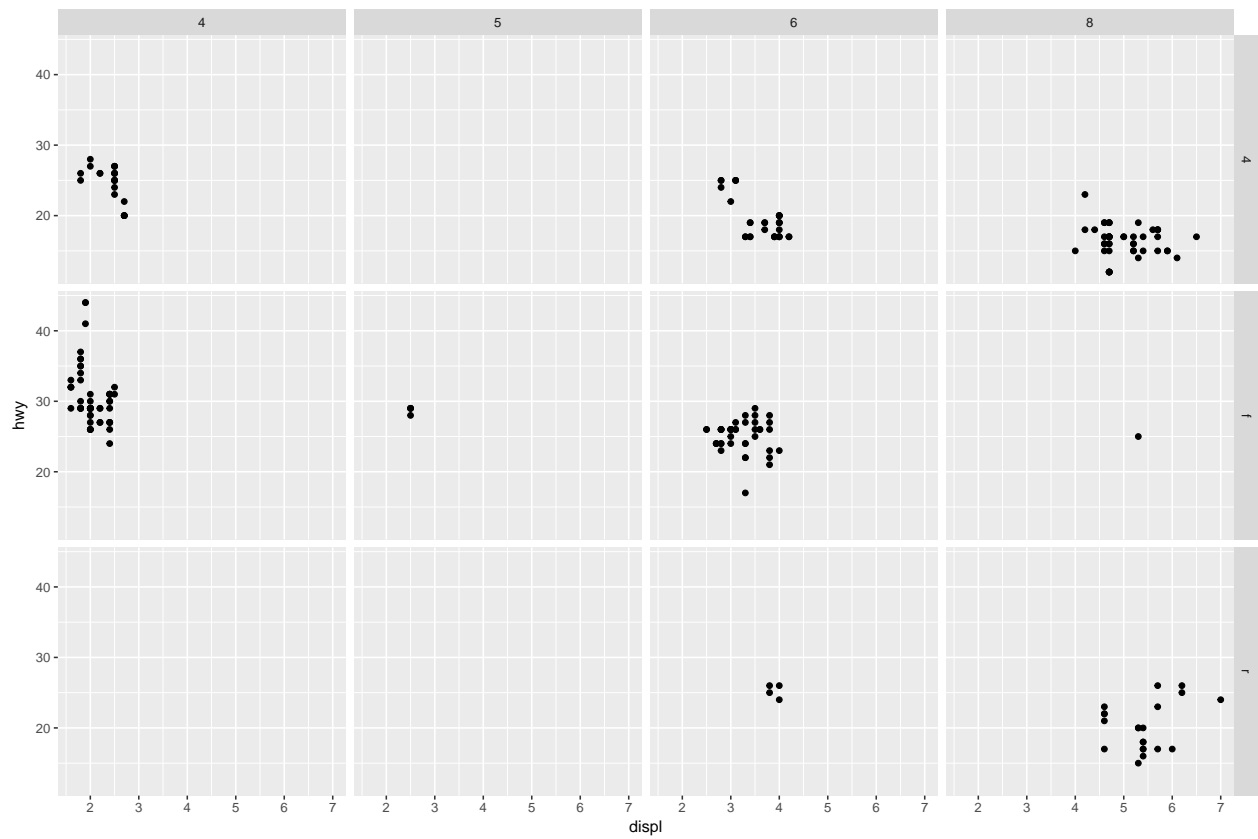
```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_wrap(~ hwy, nrow = 2)
```



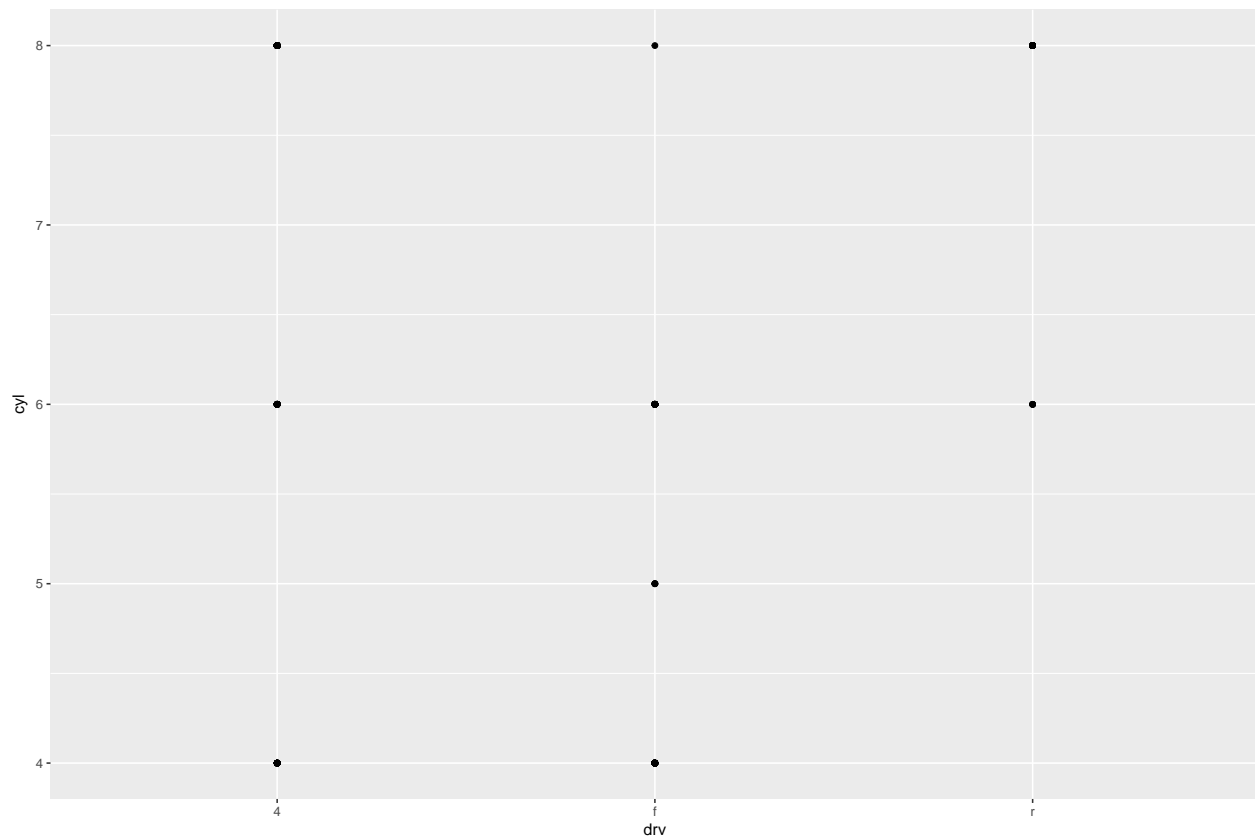
2. What do the empty cells in plot with `facet_grid(drv ~ cyl)` mean? How do they relate to this plot?

The empty cells in the plot with `facet_grid(drv ~ cyl)` indicate that there are no cars with at the intersection of that number of cylinders and that type of drivetrain (e.g. no cars with 5 cylinders and 4-wheel drive). The absence of vehicles corresponding to specific cylinder-drive combinations is also evident in the second plot. Those intersections in the second plot without a point correspond to the empty cells in the first plot (see again cars with 5 cylinders on the y-axis and 4-wheel drive on the x-axis).

```
# First plot, with drivetrain and cylinder faceted
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(drv ~ cyl)
```



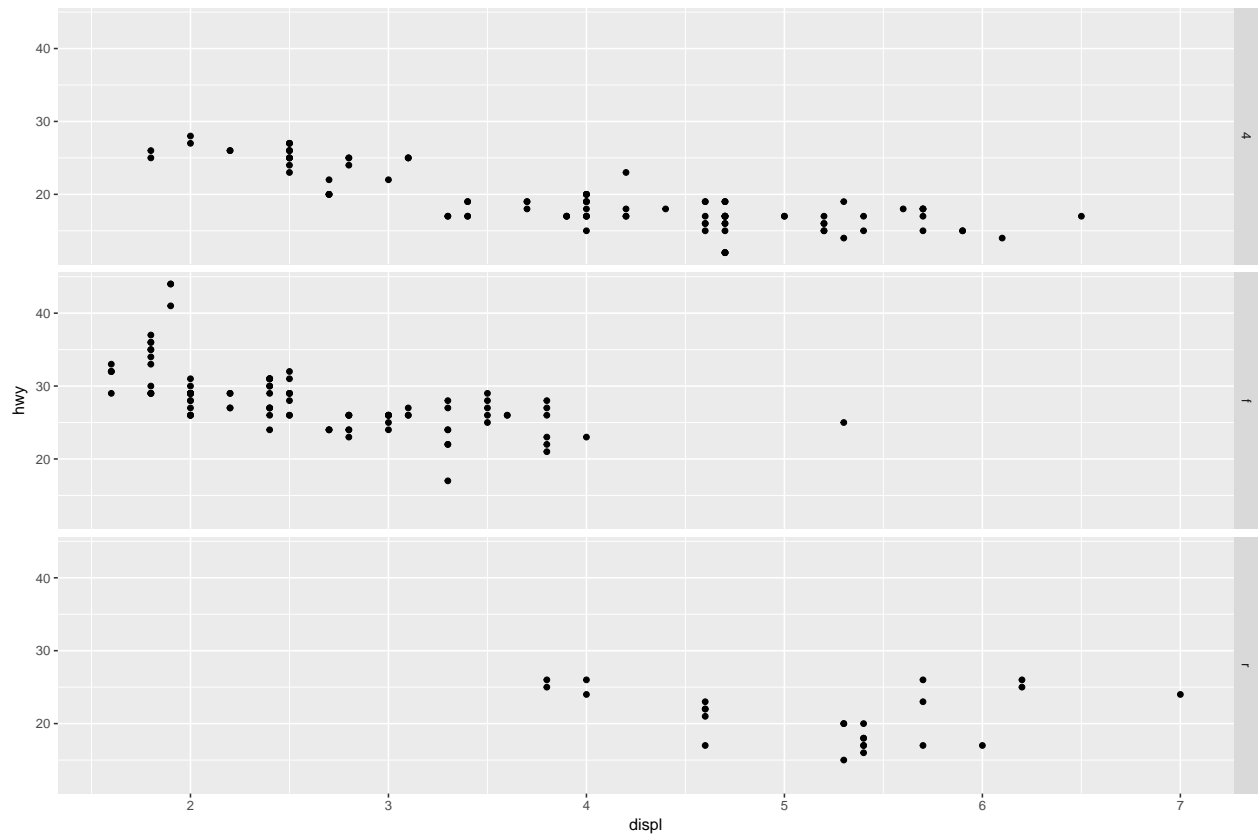
```
# Second plot, with drivetrain and cylinder represented on the axes of a single plot
ggplot(data = mpg) +
  geom_point(mapping = aes(x = drv, y = cyl))
```

3. What plots does the following code make? What does . do?

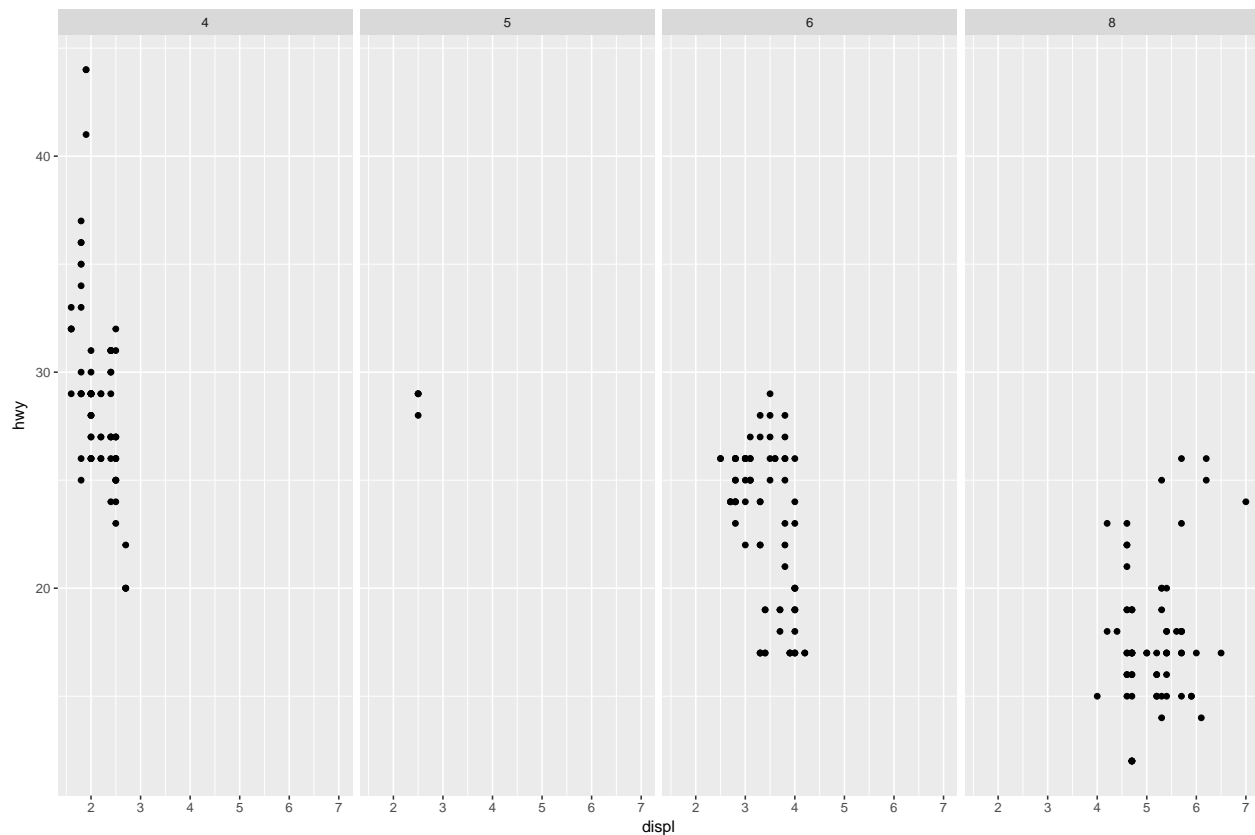
The first plot shows highway miles per gallon and engine displacement faceted by drivetrain type. The . in the second position specifies that drivetrain type should be displayed in rows. The second plot shows highway miles per gallon and engine displacement faceted by number of cylinders. The . in the first position specifies that number of cylinders should be displayed in columns.

```
# Plot of highway mpg and engine displacement faceted by drivetrain type  
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_grid(drv ~ .)
```



```
# The above is the same as the following except that the drivetrain labels shift from right to top align
#ggplot(data = mpg) +
#  geom_point(mapping = aes(x = displ, y = hwy)) +
#  facet_wrap(~ drv, nrow = 3)

# Plot of highway mpg and engine displacement faceted by number of cylinders
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_grid(. ~ cyl)
```

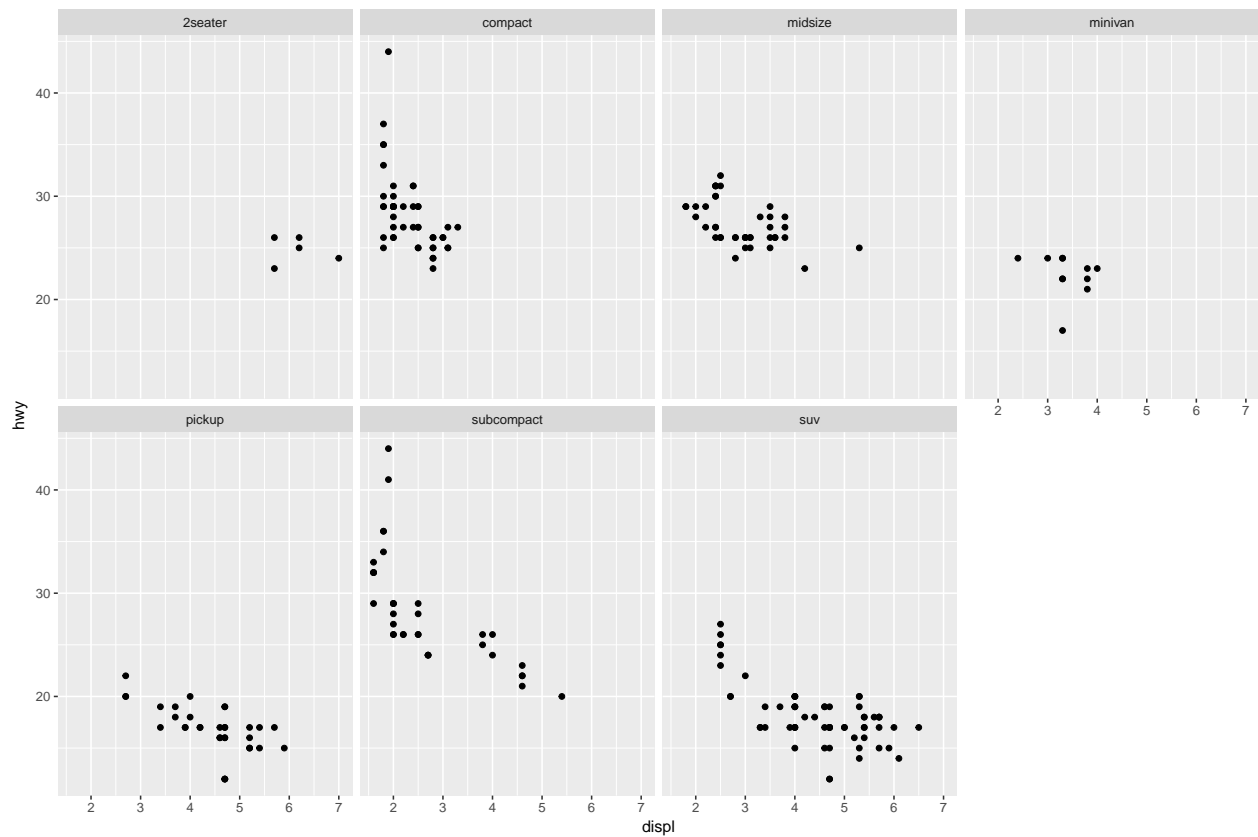


```
# The above is the same as the following. Uncomment and run the code to see.
#ggplot(data = mpg) +
#  geom_point(mapping = aes(x = displ, y = hwy)) +
#  facet_wrap(~ cyl, nrow = 1)
```

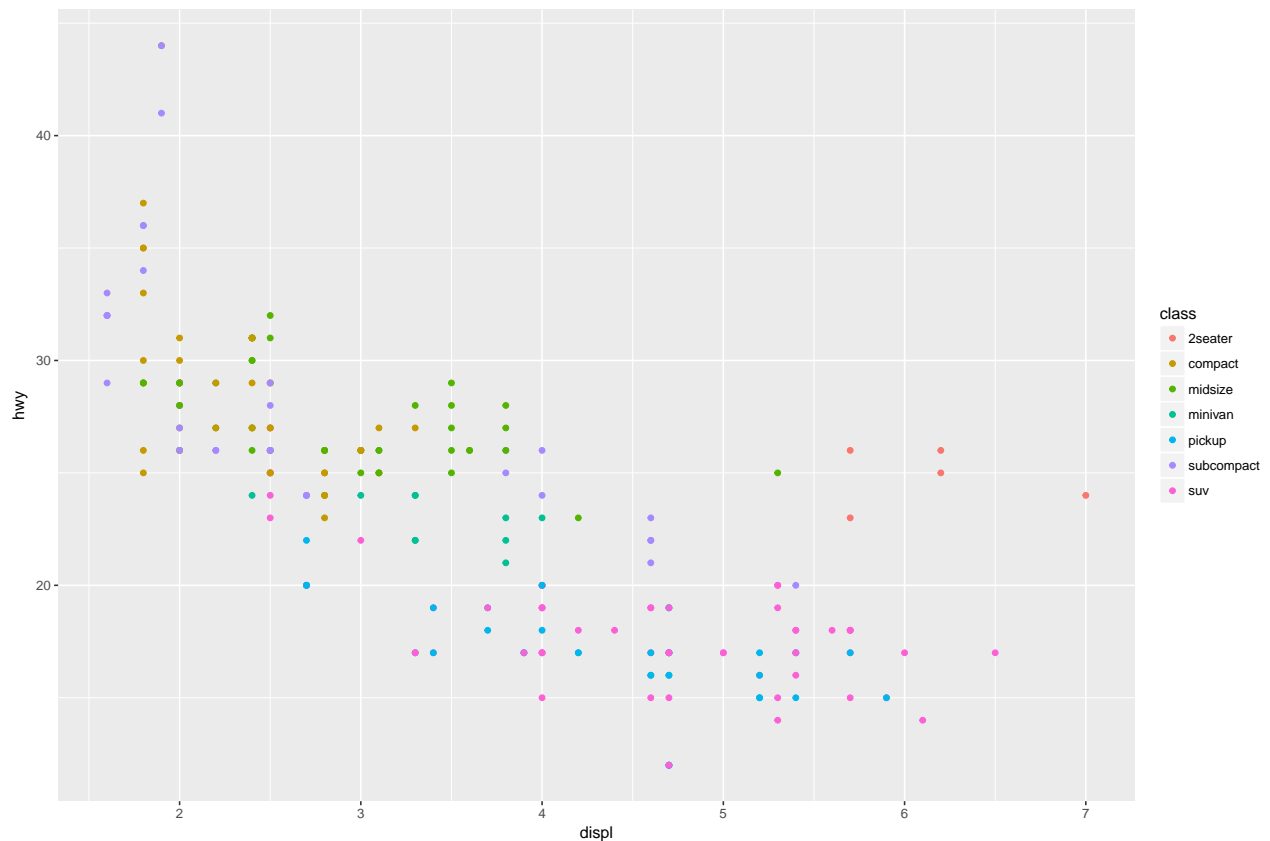
4. Take the first faceted plot in this section. What are the advantages to using faceting instead of the colour aesthetic? What are the disadvantages? How might the balance change if you had a larger dataset?

The advantage of using faceting rather than the color aesthetic is that with separate plots it is easier to see the shape and spread of the data points for each level of the variable. A disadvantage is that it's difficult to see the overall shape and spread of the observations across levels of the faceted variable. While using the color aesthetic works well with the mpg dataset, with a larger dataset, the likelihood of overlapping data points increases and with enough overlapping observations jittering may be insufficient. It may therefore be preferable to use faceting with large datasets.

```
# Plot with facets
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_wrap(~ class, nrow = 2)
```



```
# Plot with color aesthetic
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, color = class))
```



5. Read `?facet_wrap`. What does `nrow` do? What does `ncol` do? What other options control the layout of the individual panels? Why doesn't `facet_grid()` have `nrow` and `ncol` argument?

`nrow` - specifies the number of rows into which the faceted plots are fitted.

`ncol` - specifies the number of columns into which the faceted plots are fitted.

`facet_grid()` does not have `nrow` or `ncol` arguments because the number of rows and columns is determined by the number of levels of the row and column facetting variables.

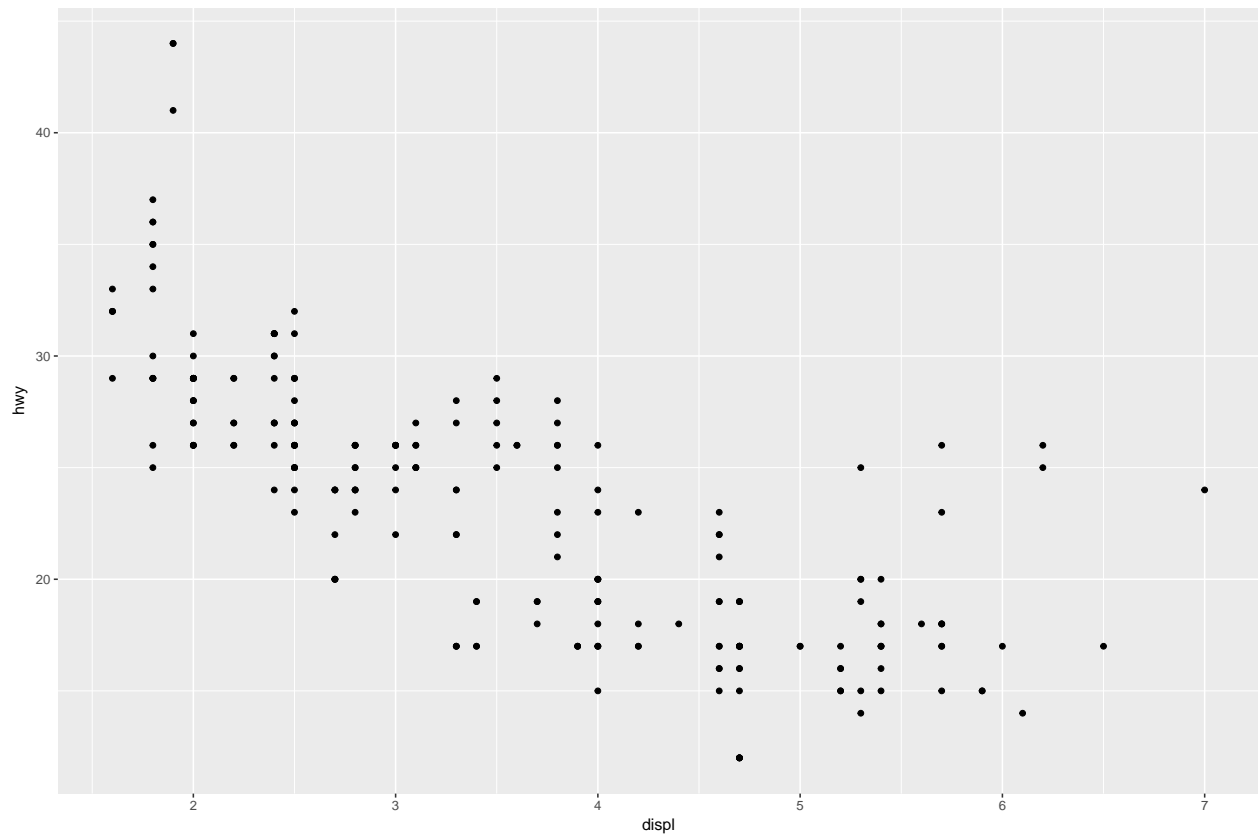
```
?facet_wrap
```

6. When using `facet_grid()` you should usually put the variable with more unique levels in the columns. Why?

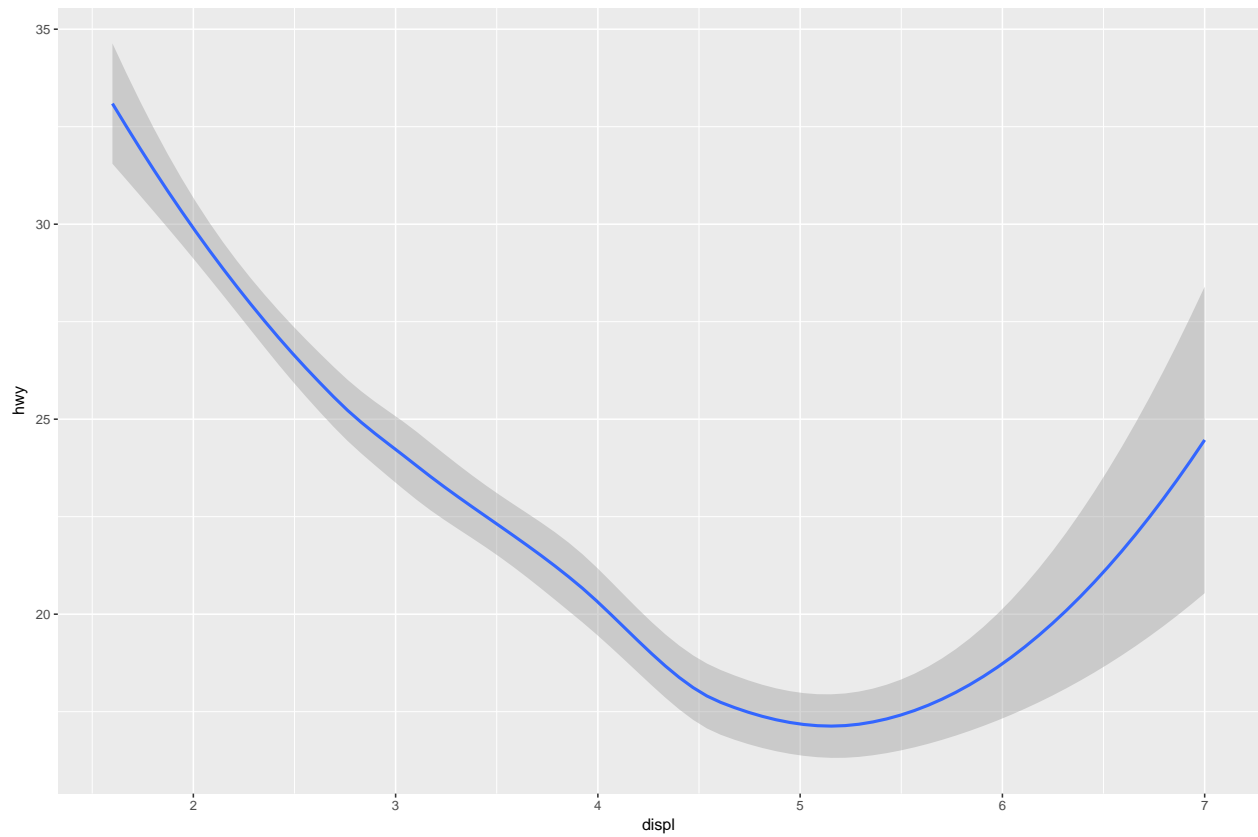
One should put the variable with more unique levels in the columns so the plots can extend vertically where there is more space. The horizontal space is limited by the page width and adding more plots compresses them, making them difficult to read.

Geometric objects

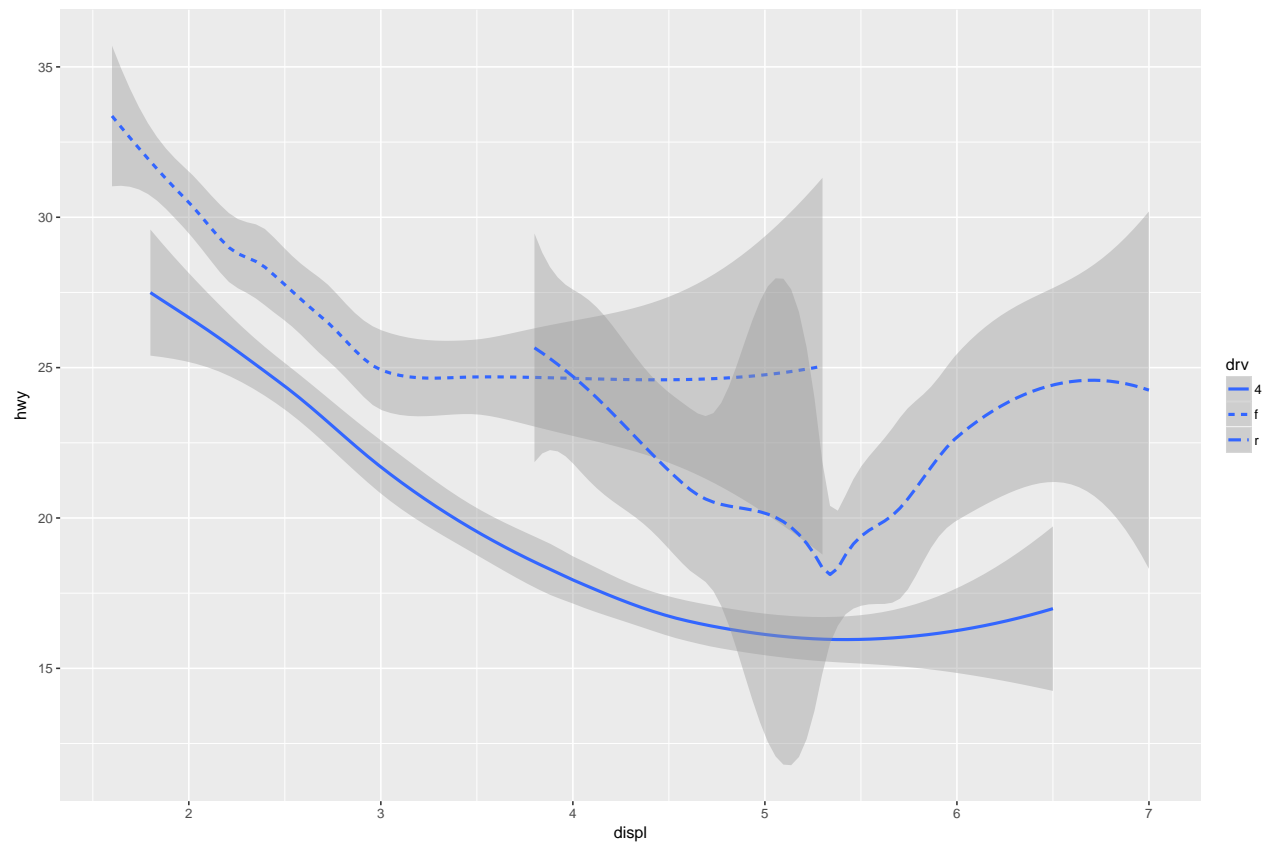
```
# Scatterplot
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy))
```



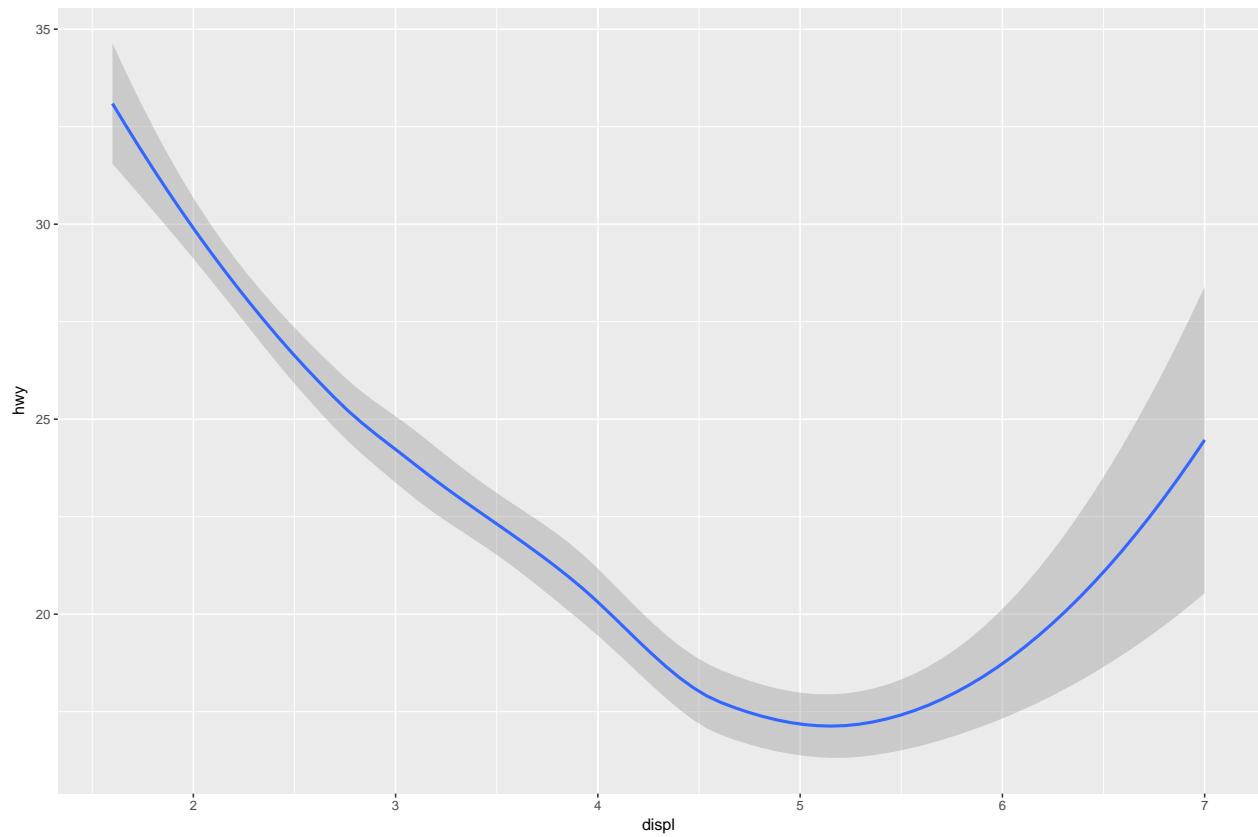
```
# Smooth geom plot  
ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy))
```



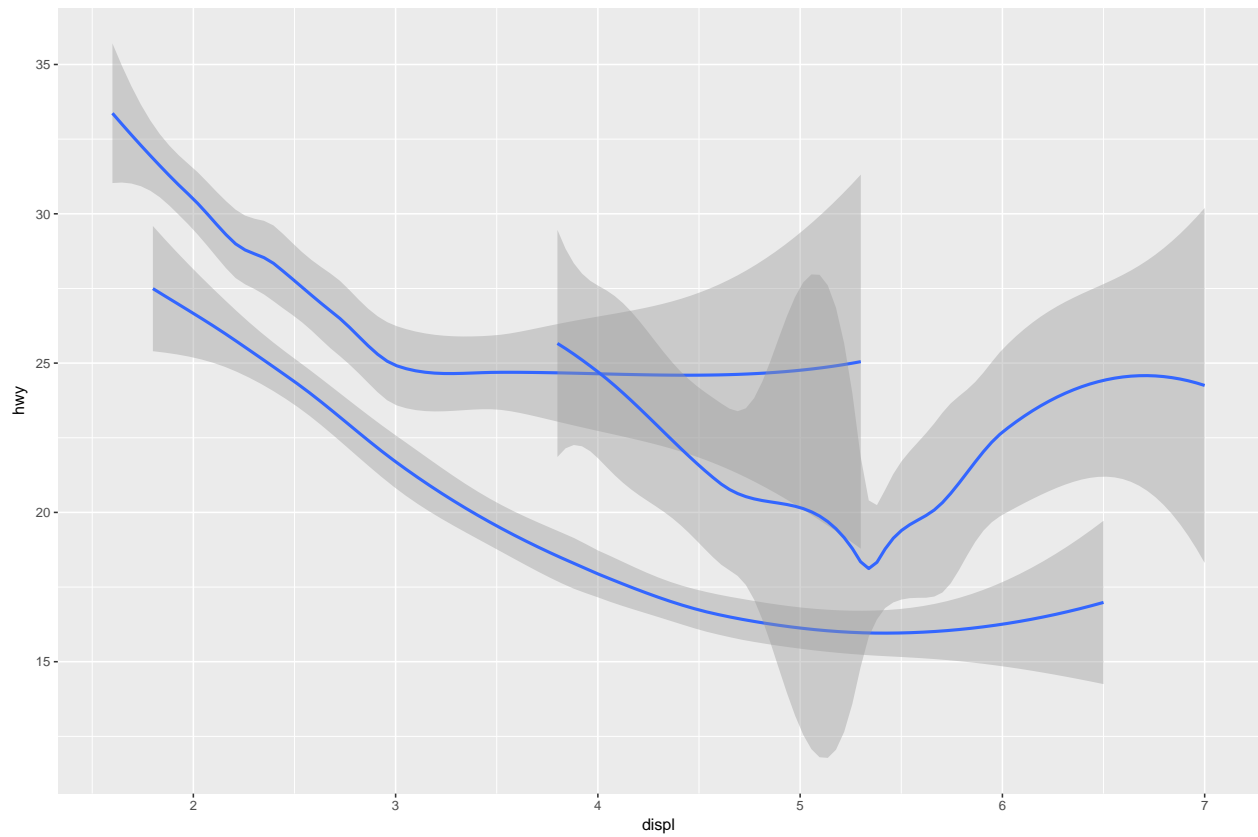
```
# Use a different linetype for each unique value of drv  
ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy, linetype = drv))
```



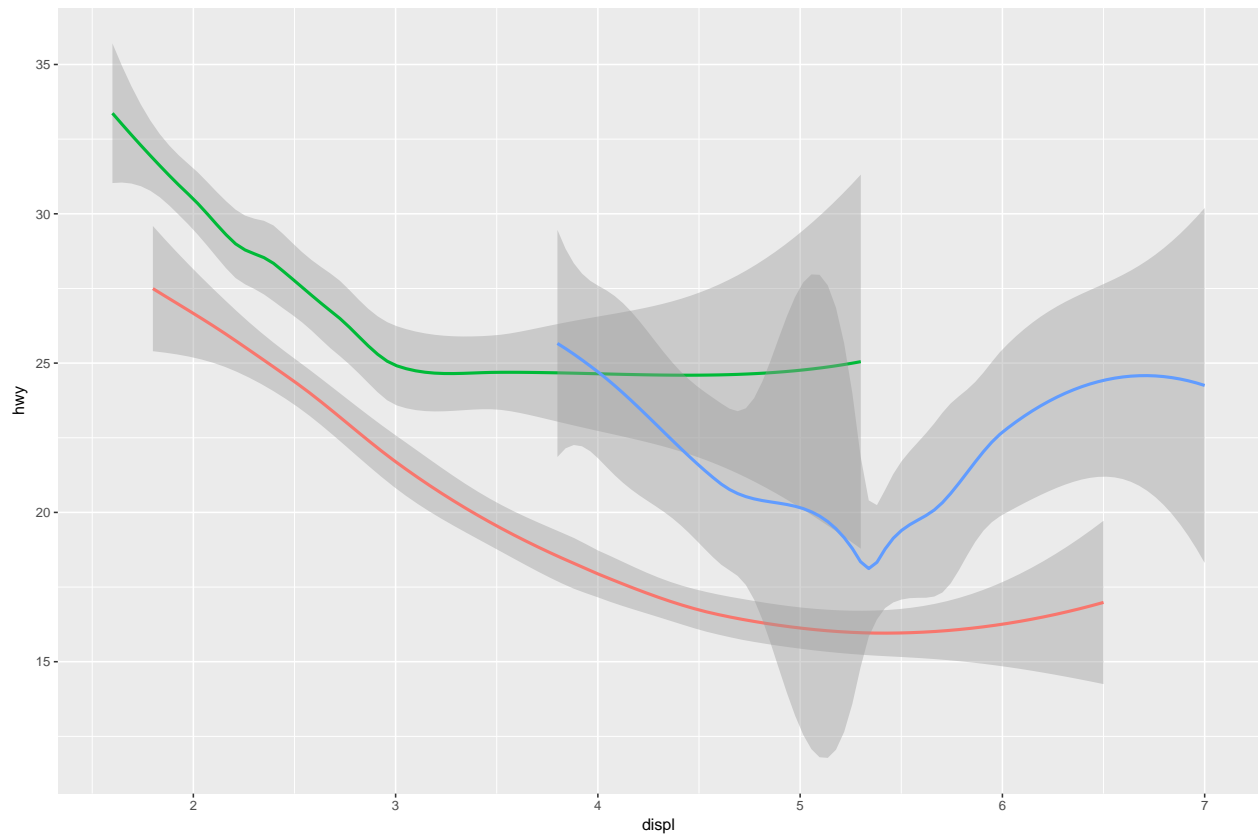
```
# Plot a single geom to display the data  
ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy))
```

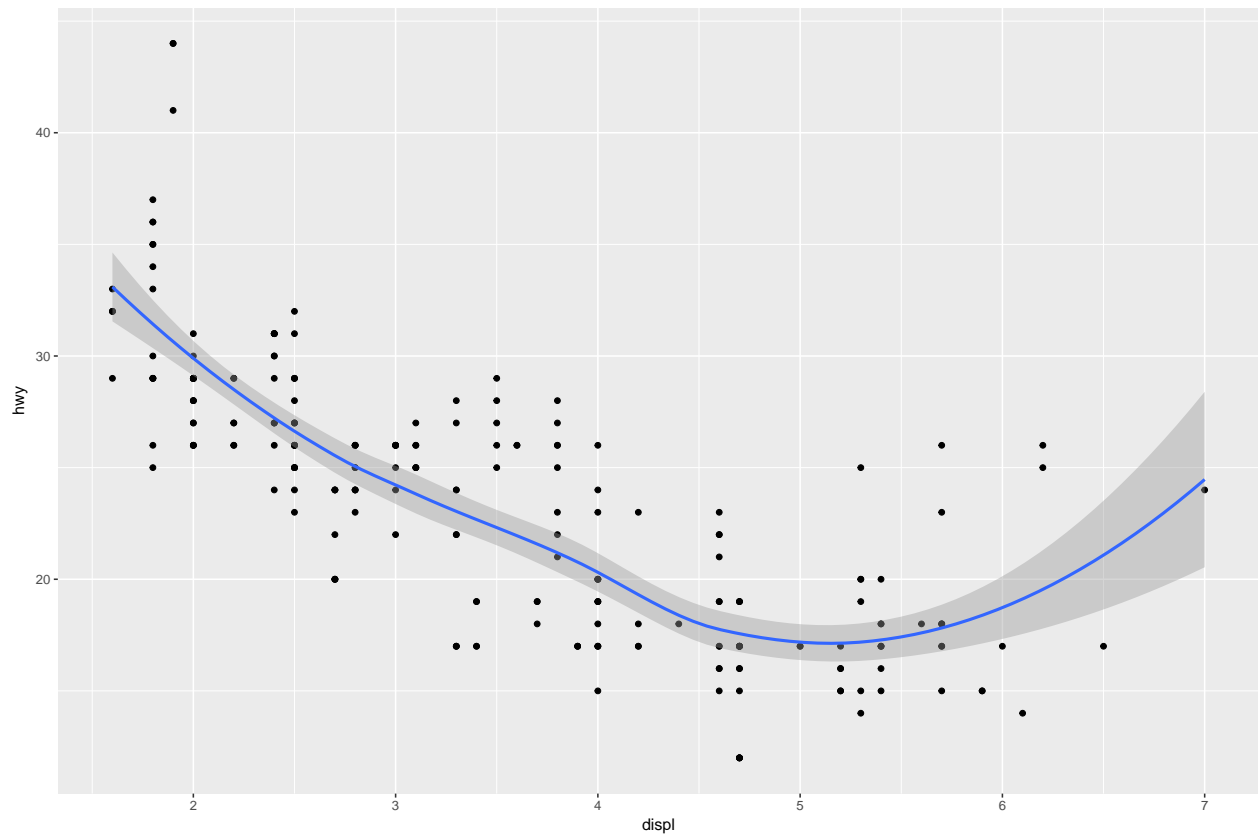
```
# Set the `group` aesthetic to `drv` to draw separate geoms for each unique value of the variable  
ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy, group = drv))
```



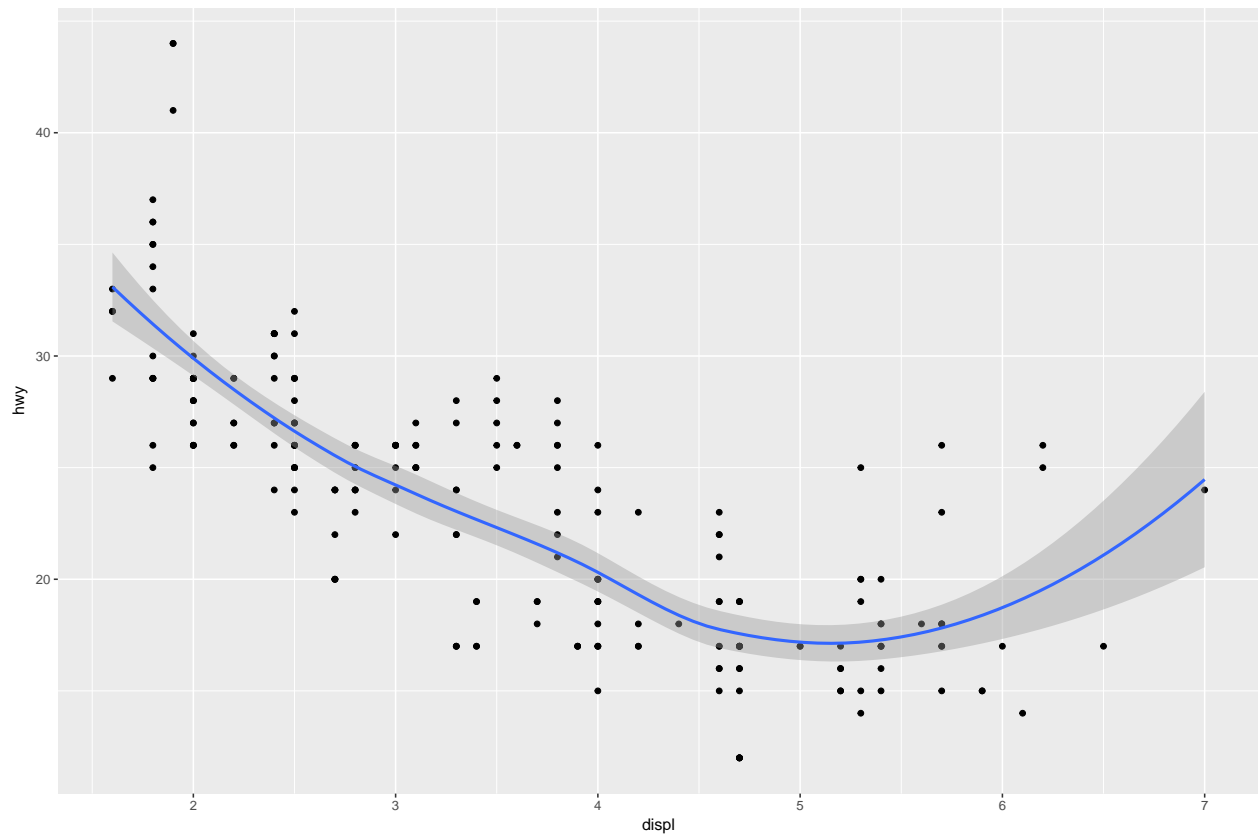
```
# Set the color aesthetic to `drv` to automatically group the data by drivetrain and distinguish them by color
ggplot(data = mpg) +
  geom_smooth(
    mapping = aes(x = displ, y = hwy, color = drv),
    show.legend = FALSE
  )
```



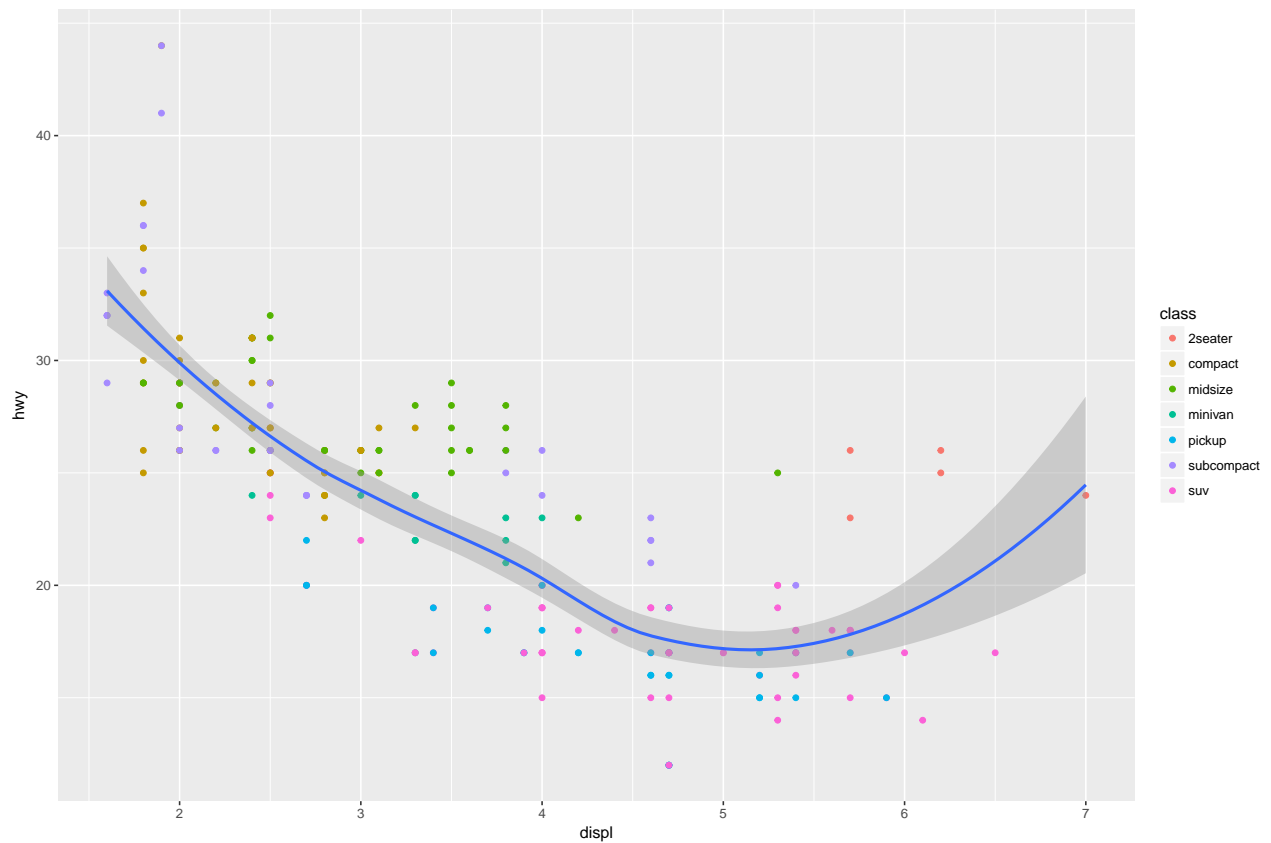
```
# Plot a smooth geom over a scatterplot of the data, the verbose way  
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  geom_smooth(mapping = aes(x = displ, y = hwy))
```



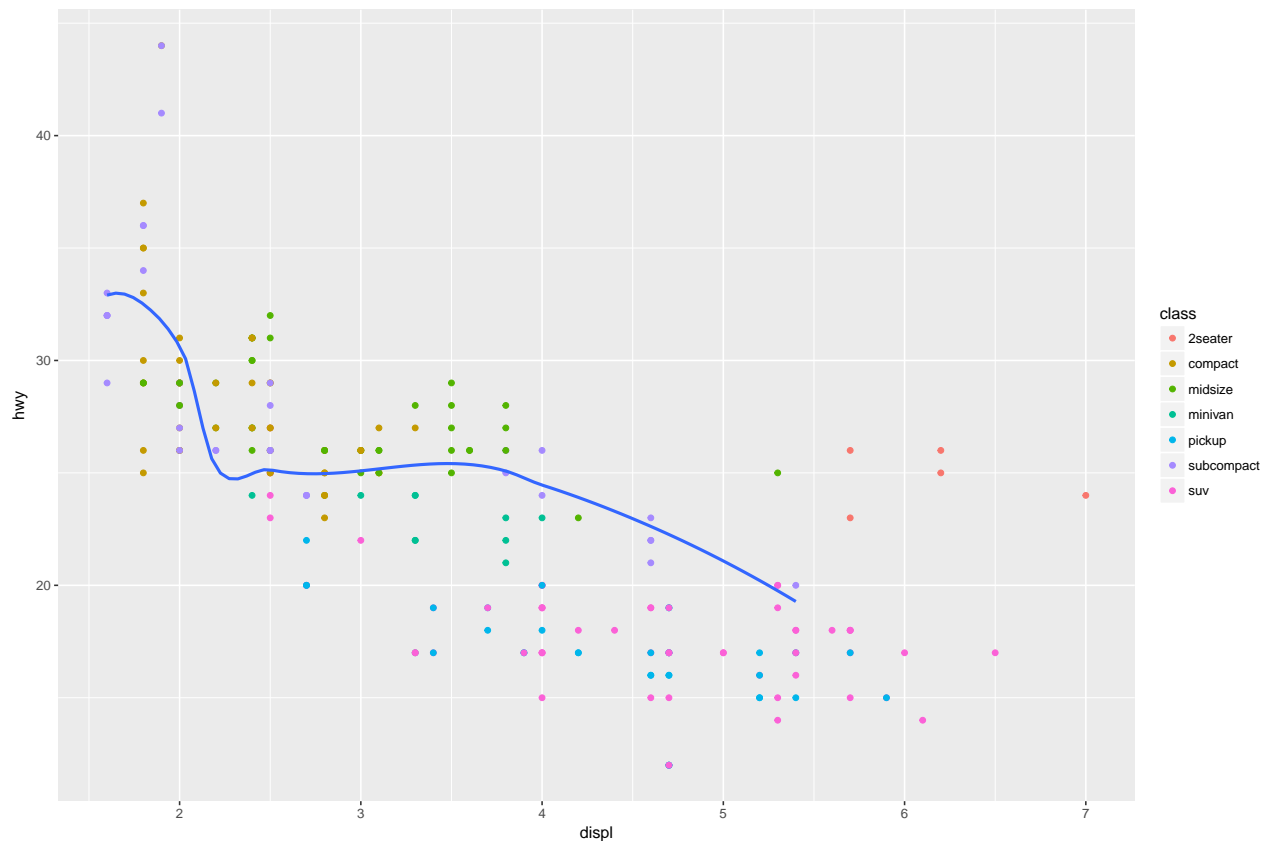
```
# The same plot with mappings passed to `ggplot()`  
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point() +  
  geom_smooth()
```



```
# Color to the points by car class  
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point(mapping = aes(color = class)) +  
  geom_smooth()
```



```
# Plot all classes of car, but draw a smooth line geom for only cars of the subcompact class
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
  geom_point(mapping = aes(color = class)) +
  geom_smooth(data = filter(mpg, class == "subcompact"), se = FALSE)
```



Exercises 3.6.1

1. What geom would you use to draw a line chart? A boxplot? A histogram? An area chart?

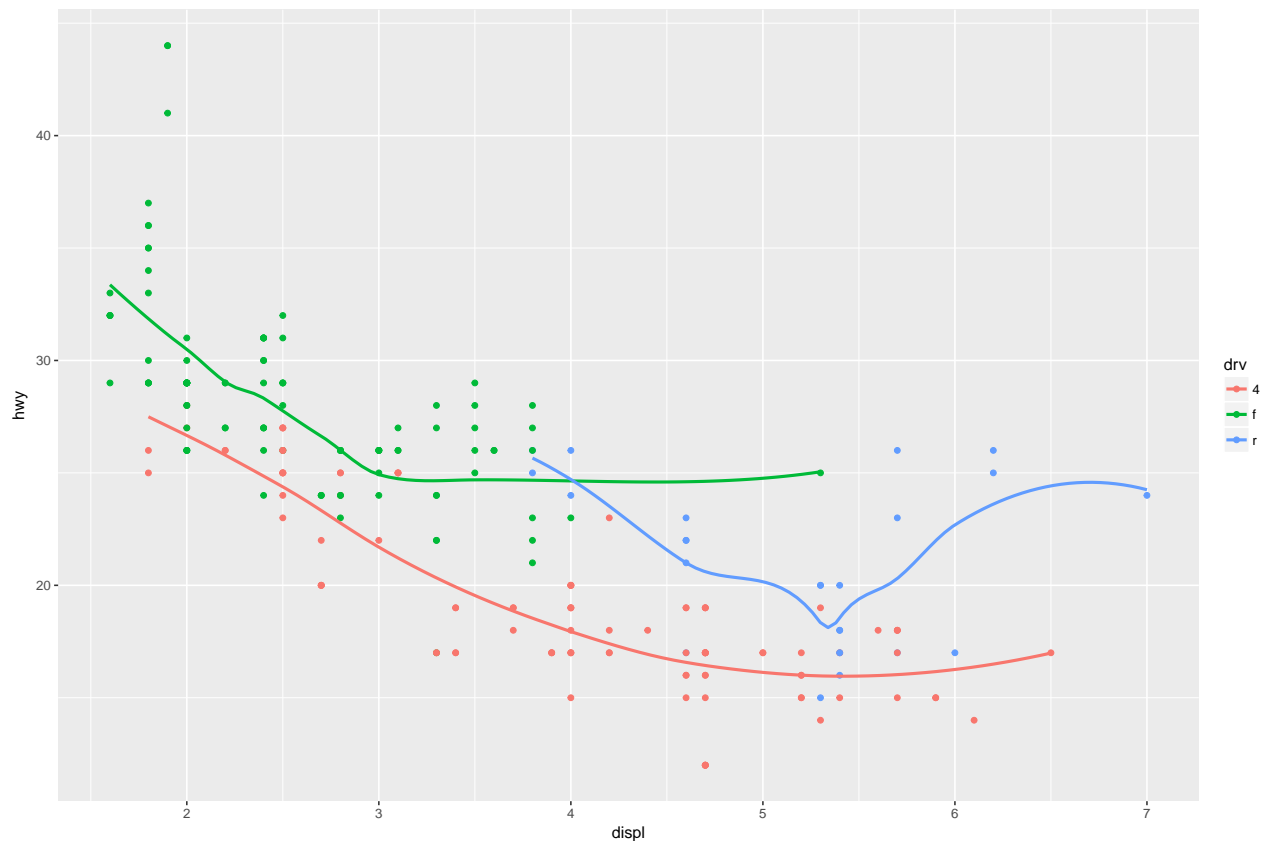
What geom would you use to draw a(n)

- line chart: `geom_line()` - boxplot: `geom_boxplot()` - histogram: `geom_histogram()` - area chart: `geom_area()`

2. Run this code in your head and predict what the output will look like. Then, run the code in R and check your predictions.

The output will be a scatterplot with engine displacement mapped to the x-axis, miles per gallon highway on the y-axis, and the points colored by type of drivetrain. The scatterplot will be overlaid by one solid smooth line for each drivetrain type and will not expand to indicate the confidence interval. The color of the line for a given type of drivetrain will match the color of the points for the same.

```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth(se = FALSE)
```



3. What does `show.legend = FALSE` do? What happens if you remove it? Why do you think I used it earlier in the chapter?

`show.legend = FALSE` indicates that a legend should not be included in the plot. If it is removed, the default is to include a legend if any aesthetics are mapped. It may have been used earlier in the chapter to introduce us to the option, to make the final graph in the set of three match the first two, which did not have legends, or for another unknown reason.

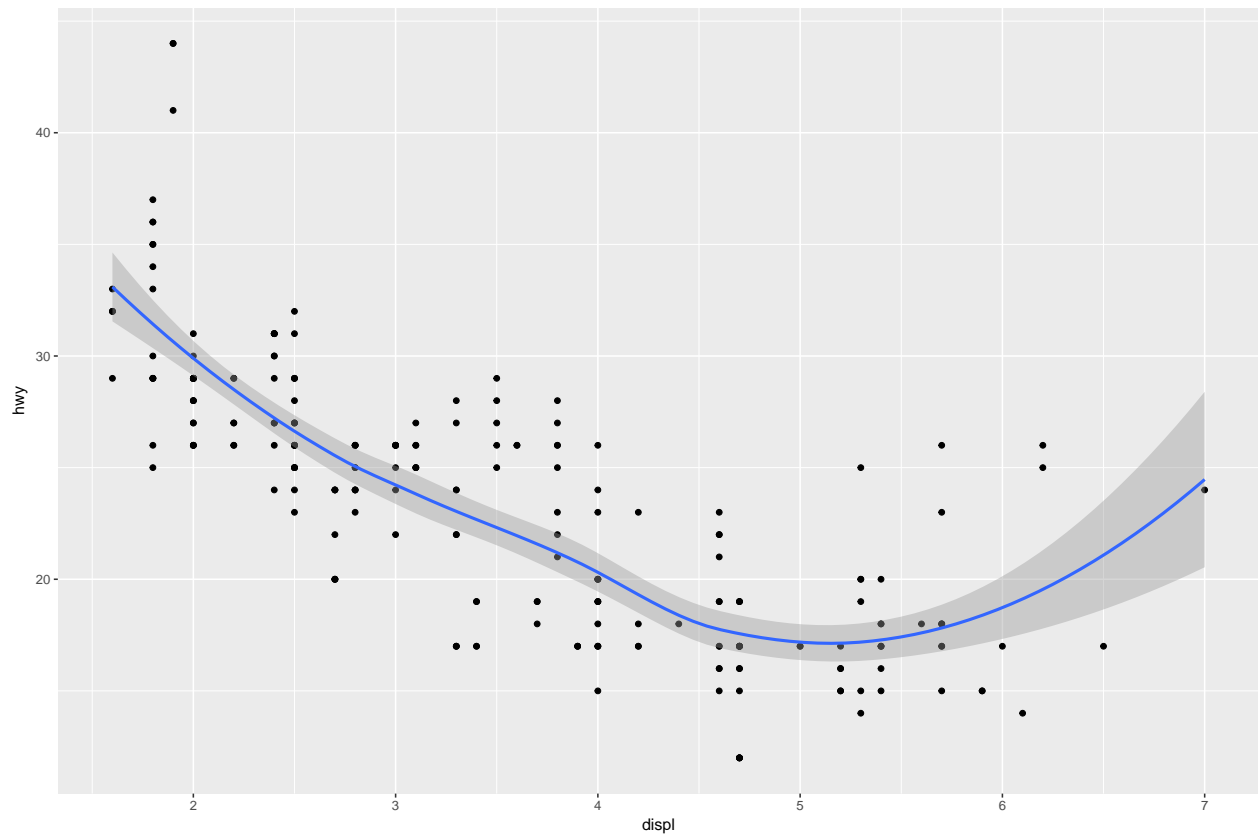
4. What does the `se` argument to `geom_smooth()` do?

The `se` argument to `geom_smooth()` indicates whether to include confidence intervals around smooth. The default is `TRUE`.

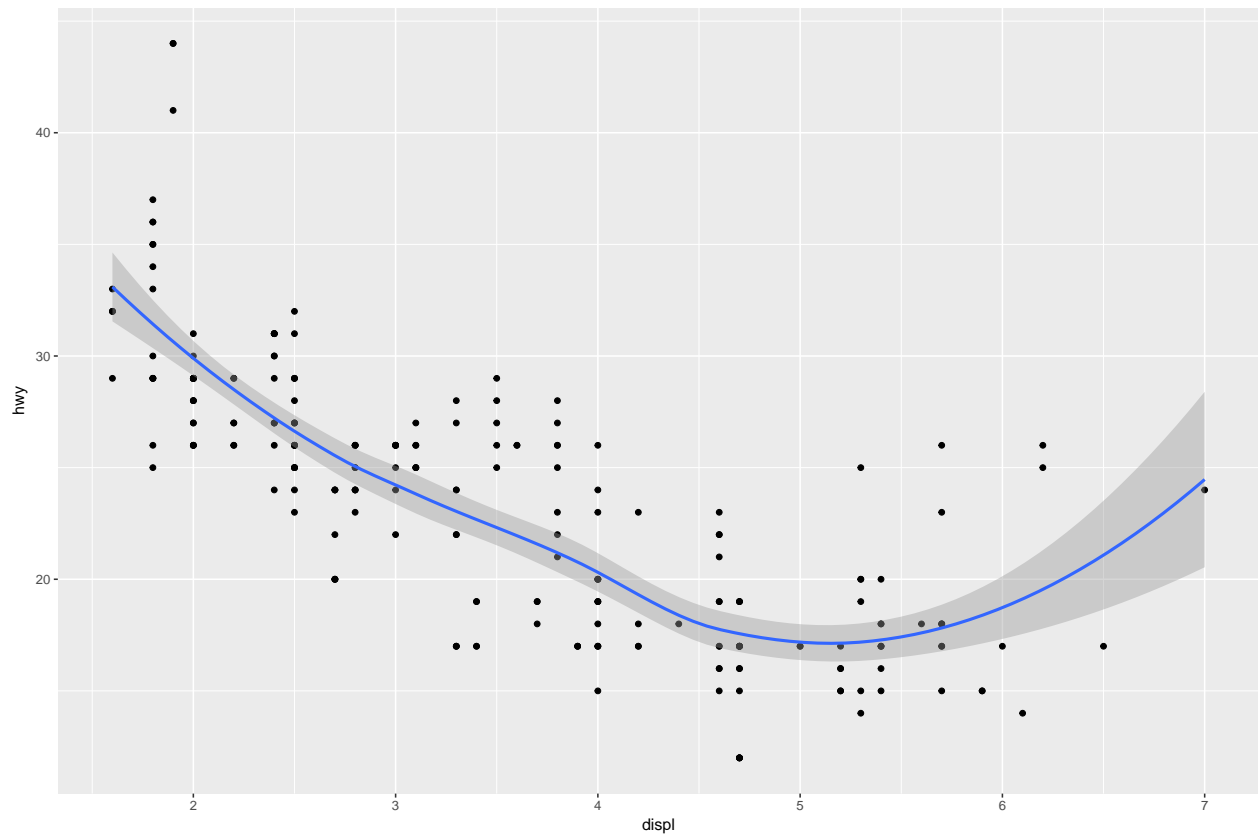
5. Will these two graphs look different? Why/why not?

The two plots will be identical. In the first plot, the selection of `mpg` as the data frame from which to draw the variables and the mapping of `displ` to the x-axis and `hwy` to the y-axis is specified in the global mappings, within `ggplot()`. These mappings extend to the following layers, in this case `geom_point()` and `geom_smooth()` unless those layers explicitly overwrite the global mappings, which they don't here. In the second plot, the same mappings are specified as in the global settings of the first plot. Therefore, in both plots, the mappings specified for `geom_point()` and `geom_smooth()` are identical, though they are explicitly laid out for each layer in plot 2, whereas they are carried over from the global mappings in plot 1.

```
# Plot using global mappings
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +
  geom_point() +
  geom_smooth()
```

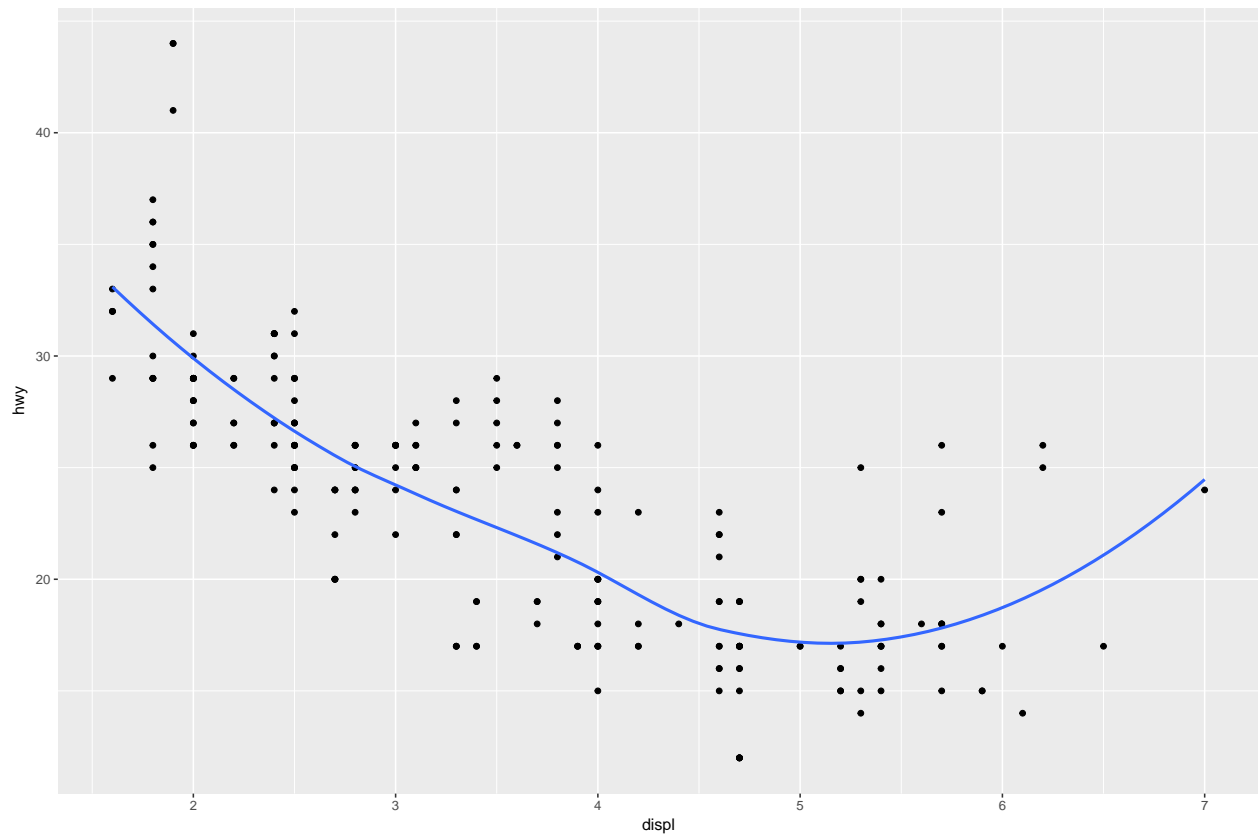



```
# Plot specifying mappings in each geom layer  
ggplot() +  
  geom_point(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_smooth(data = mpg, mapping = aes(x = displ, y = hwy))
```

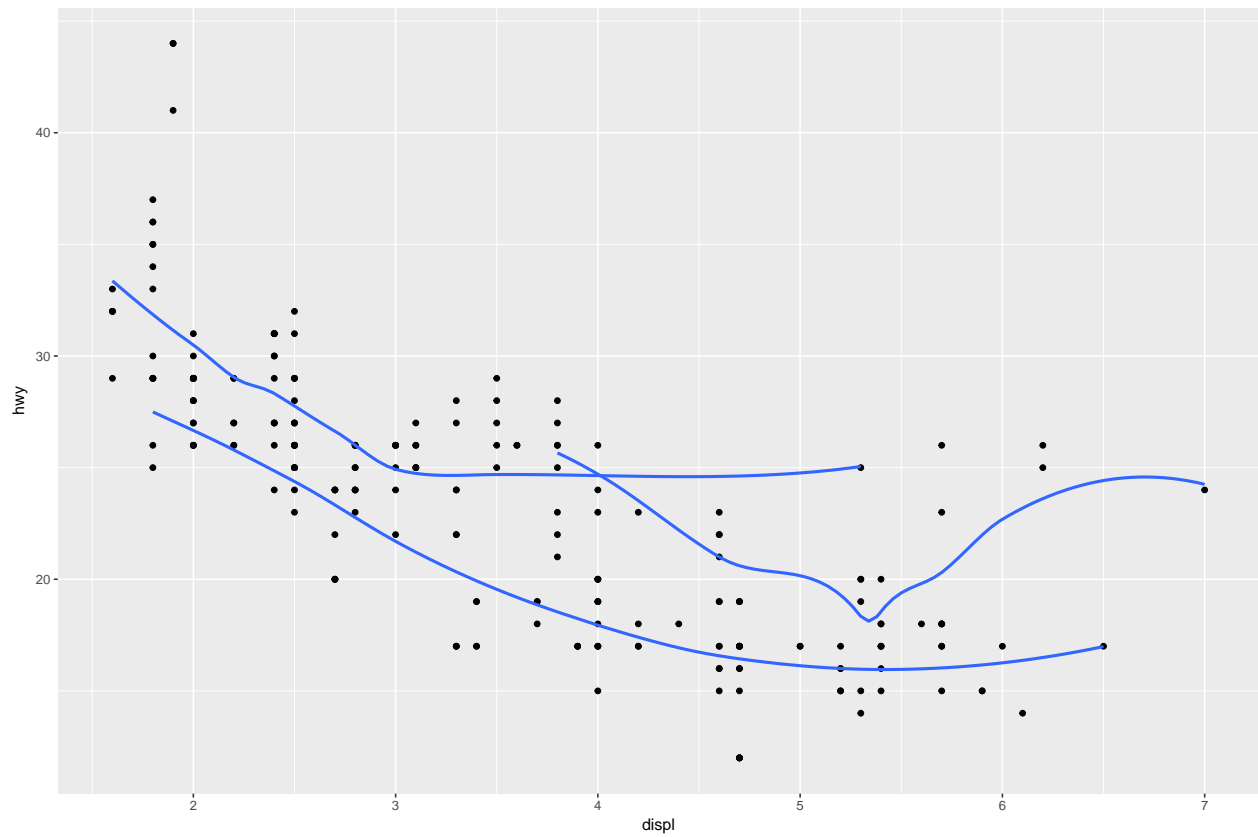


6. Recreate the R code necessary to generate the following graphs.

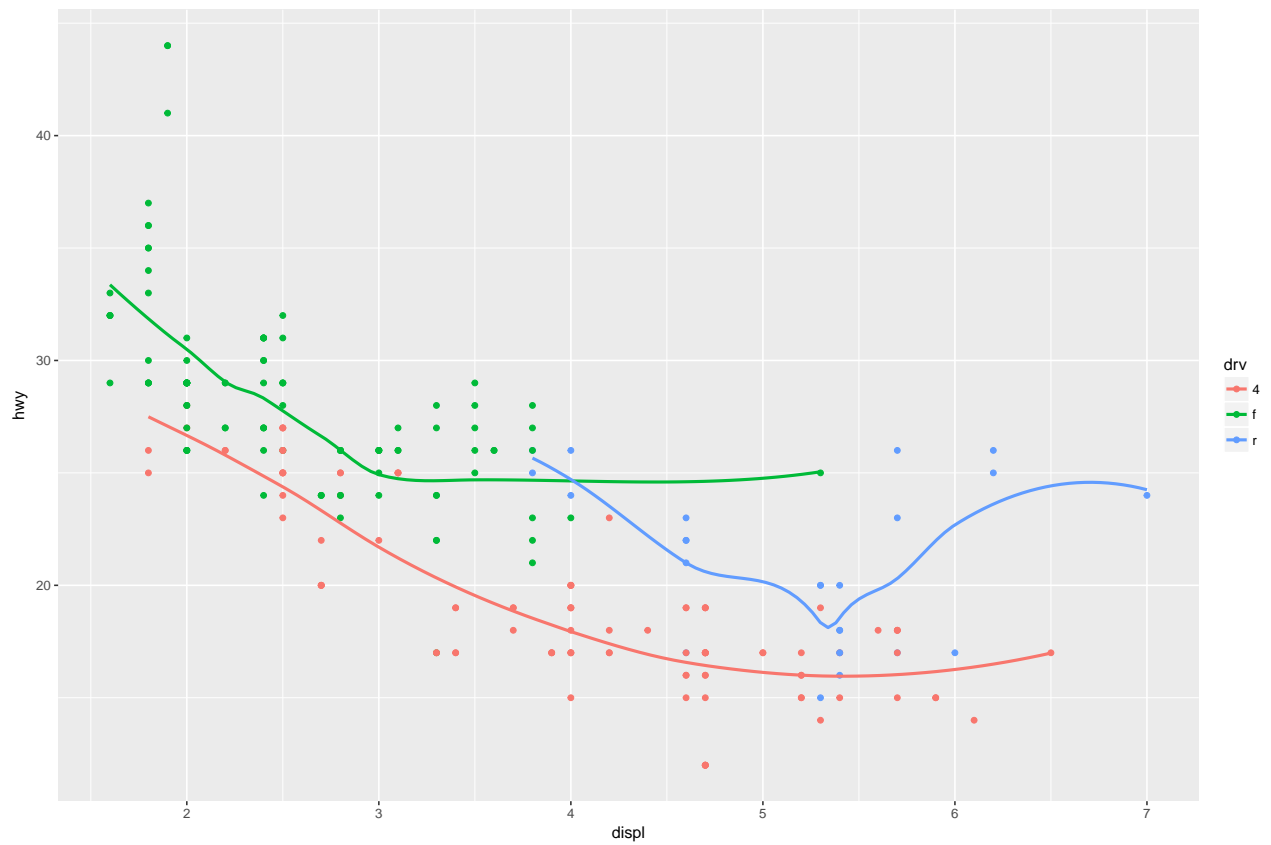
```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point() +  
  geom_smooth(se=FALSE)
```



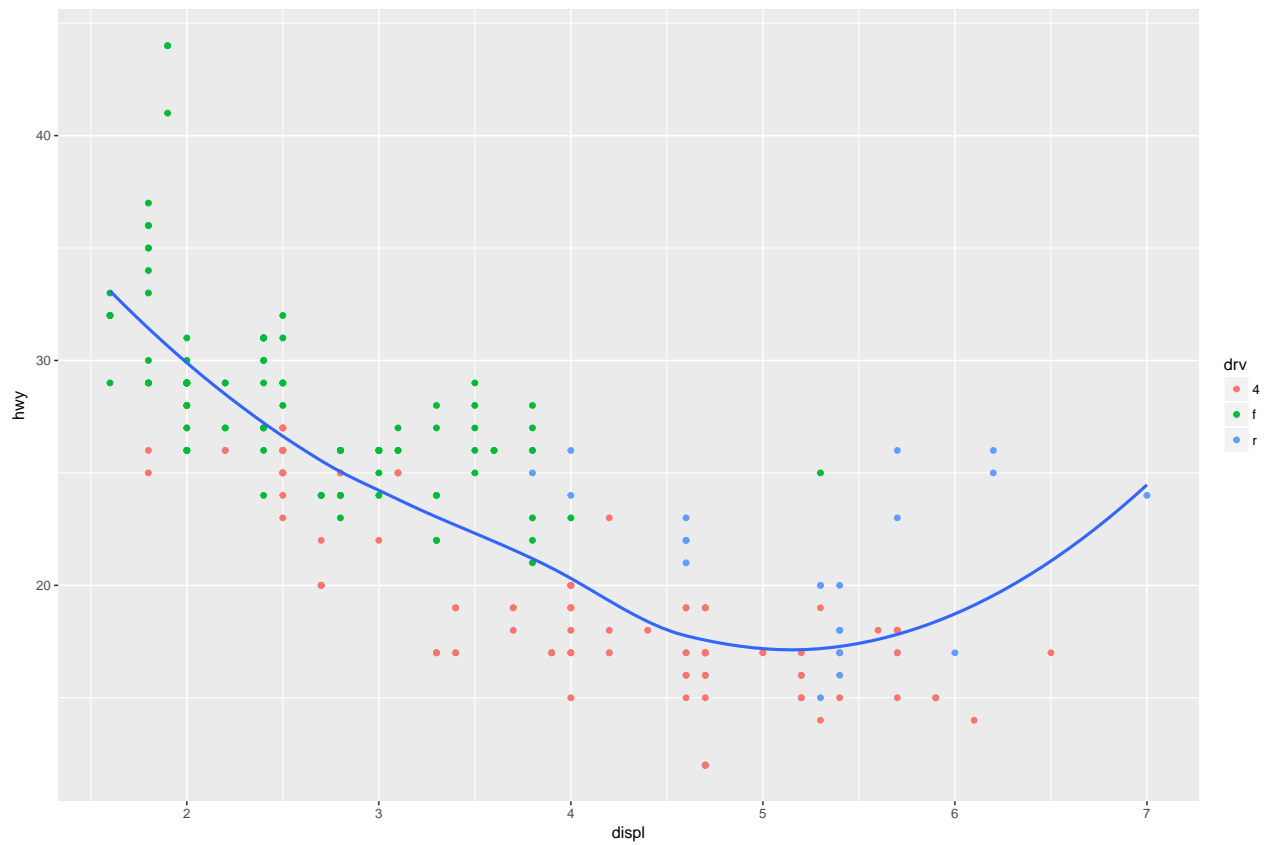
```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point() +  
  geom_smooth(aes(group = drv), se=FALSE)
```



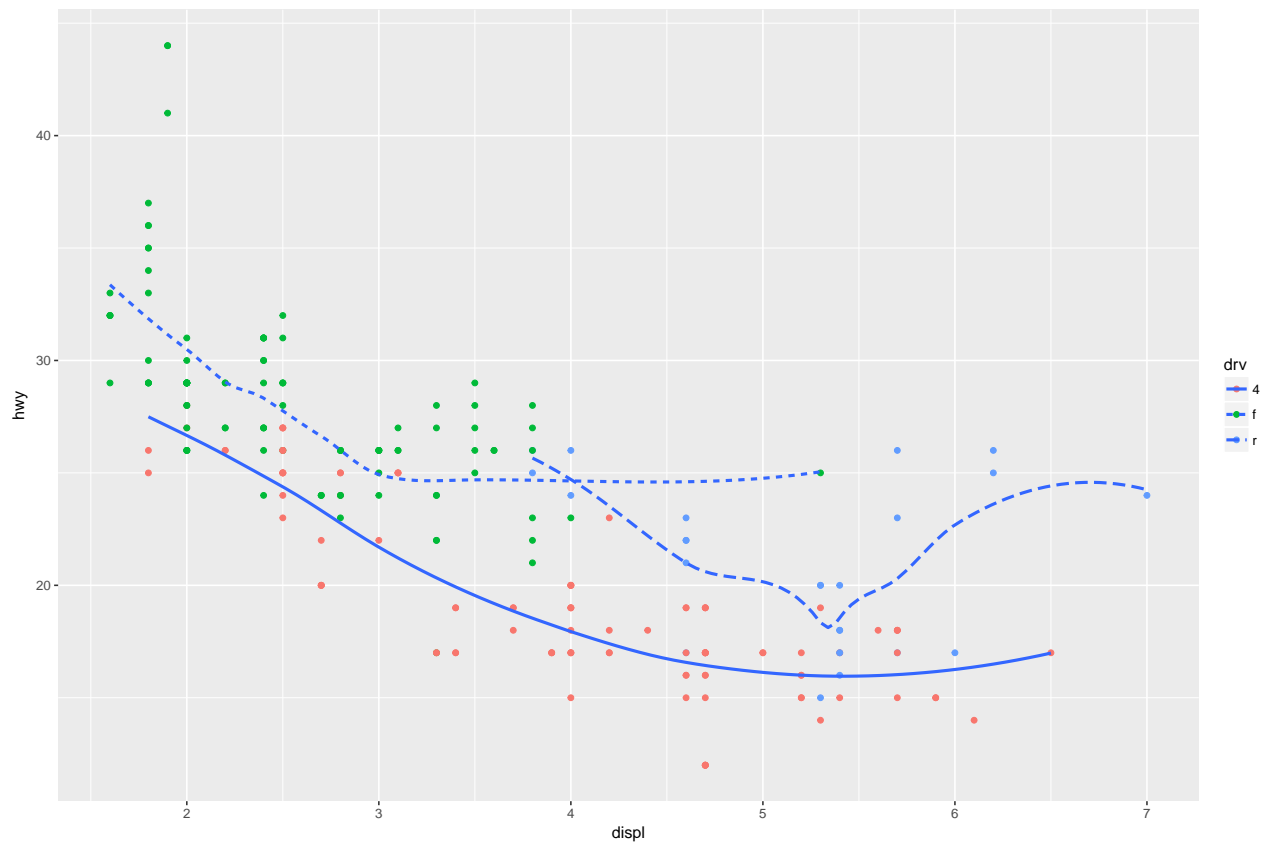
```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy, color = drv)) +  
  geom_point() +  
  geom_smooth(se=FALSE)
```



```
ggplot(data = mpg, mapping = aes( x = displ, y = hwy)) +  
  geom_point(aes(color = drv)) +  
  geom_smooth(se=FALSE)
```



```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point(aes(color=drv)) +  
  geom_smooth(aes(linetype=drv),se=FALSE)
```



```
ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_point(size = 4, colour = "white") +  
  geom_point(aes(colour = drv))
```

