

2.1 Data Visualisation

Actuarial Data Science Applications (ACTL4305/5305)

Fei Huang, UNSW



Table of contents

- Exploratory Data Analysis: an Introduction
- Data Visualization using `ggplot2`
- Declaring Data
- Aesthetic Mapping
- Facets
- Geometric Object



Reading List

- [The Art of Data Science](#), Chapter 4.1, 4.2.
- R for Data Science [Online Book](#), Chapters 2, 3



Learning Objectives of Exploratory Data Analysis

- Understand how to do exploratory data analysis with the tidyverse package in R
- Explore features of data using data visualisation.
- Explain common features of data, such as categorical variables, missing values, unreliable/non-validated data, outliers and high cardinality features, that may lead to problems.
- Apply appropriate methods to deal with common data problems.
- Apply a range of techniques to assess data quality.
- Use R to manipulate (eg filter, merge, sort, group by, summarise, etc.)
- Use R to import, export and tidy data.
- Apply the process to do exploratory data analysis with practical datasets
- Apply R Markdown for communication and reproducible analysis



Exploratory Data Analysis: an Introduction



Exploratory Data Analysis (EDA)

- A set of procedures to produce descriptive and graphical summaries of the data
- Explore the data as they are without making assumptions
- To examine your data and understand relationship among variables
- To determine if there are any problems with your dataset
- To determine whether the question you are asking can be answered by the data that you have
- To develop a sketch of the answer to your question



The process of exploratory data analysis

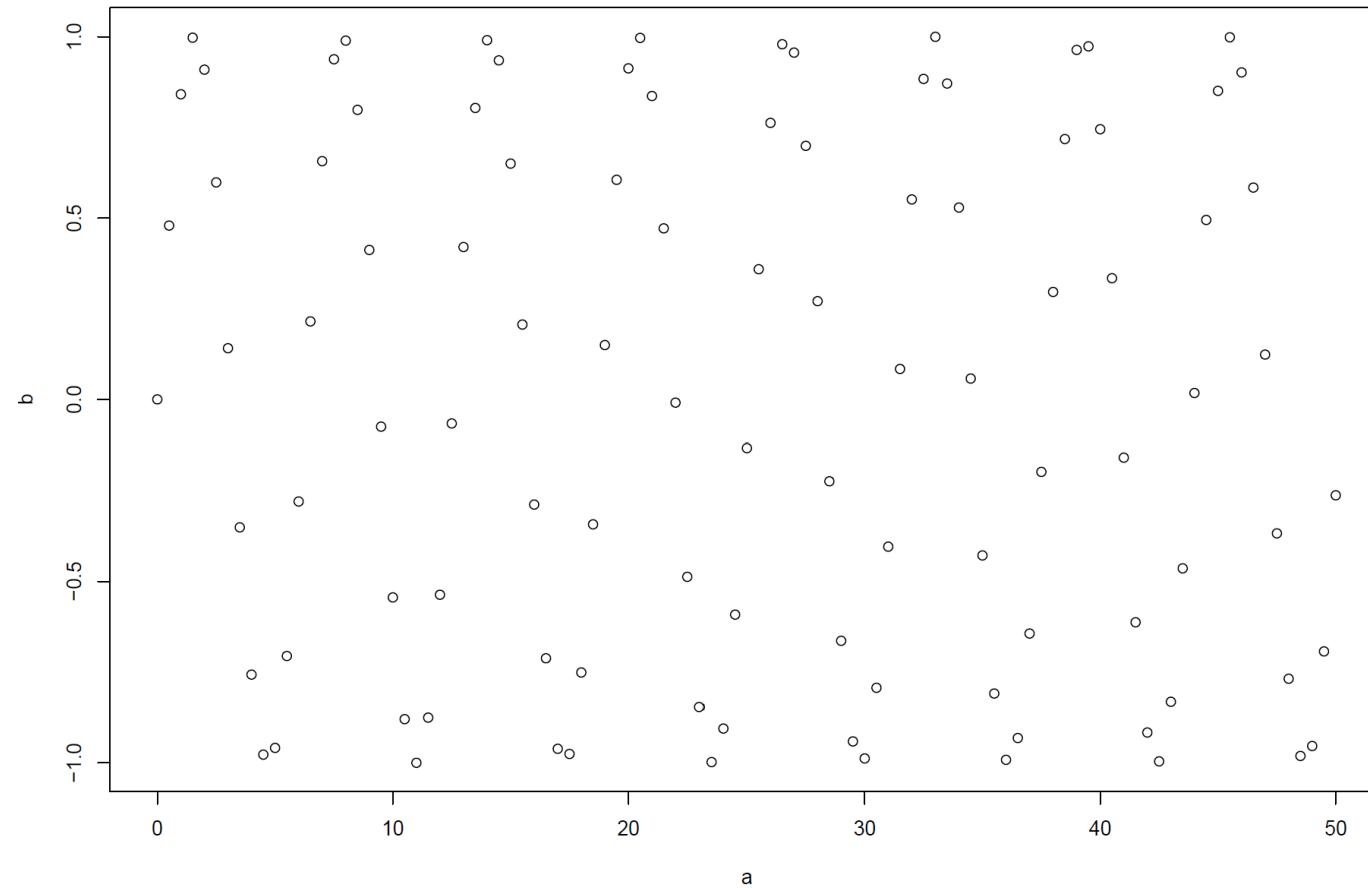
The EDA is an iterative cycle. You:

1. Formulate your question
2. Search for answers by
 1. Collect and Import data
 2. Check data quality and Cleansing data
 3. Manipulate and Transform data
 4. Visualise data
3. Use what you learn to refine your questions and/or generate new questions
 - Data Visualization is arguably the most important tool for EDA.



Example

- What pattern can you see from this plot?

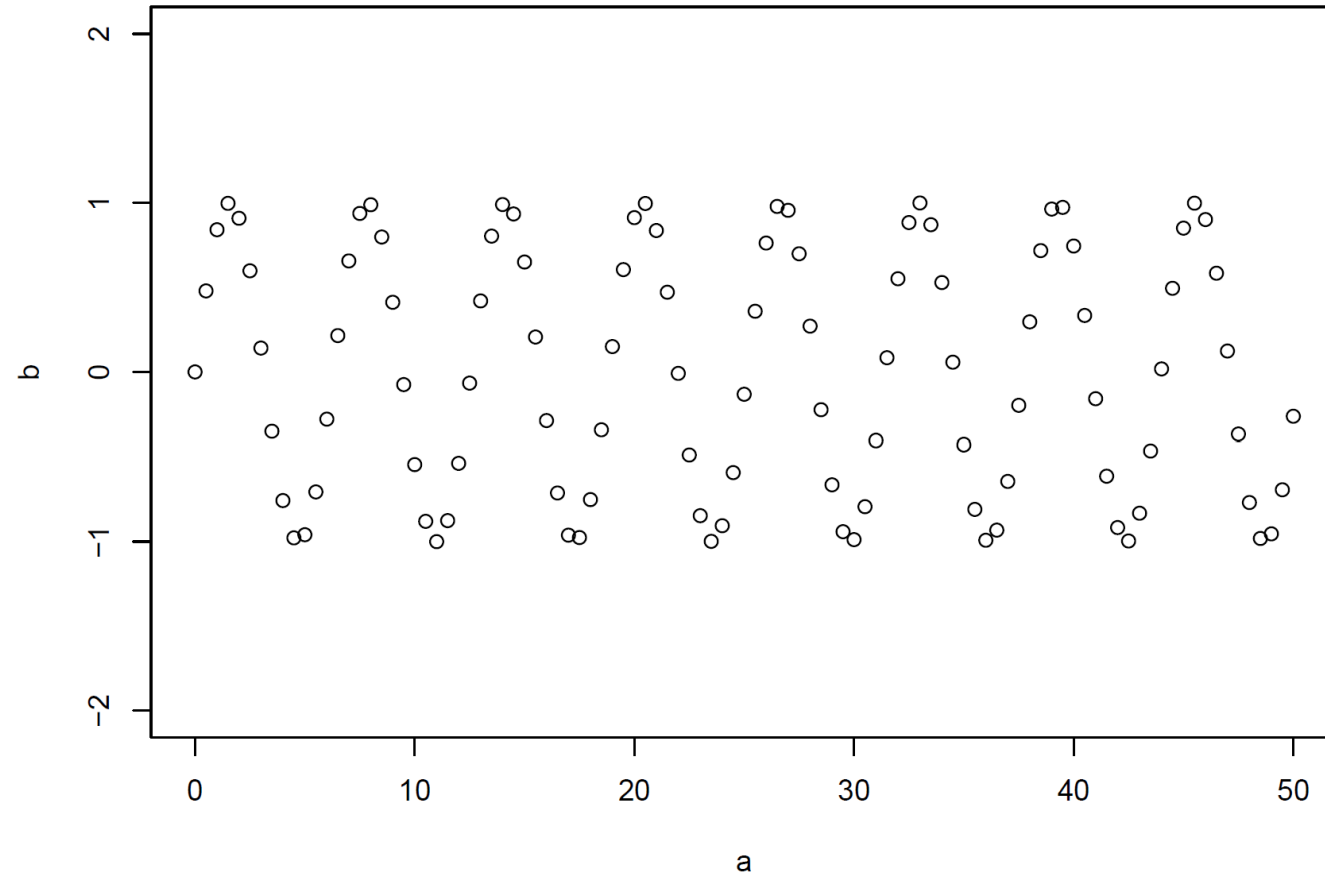


Plot 1



Example

- What pattern can you see from this plot?



Plot 2



Data Visualization

‘Visually attractive graphics also gather their power from content and interpretations beyond the immediate display of some numbers. The best graphics are about the useful and important, about life and death, about the universe. Beautiful graphics do not traffic with the trivial.’ — Edward Tufte



Data Visualization using `ggplot2`



Explore features of data using data visualisation

- A *statistical graphic* maps variables of
 1. a *dataset* to
 2. *aesthetic properties* of
 3. *geometric objects*.
- `ggplot2` is part of `tidyverse`
- Using `ggplot2` to visualise your data
- A [ggplot2 grammar guide](#)



ggplot2

```
1 #install.packages("tidyverse")  
2 library(tidyverse)
```



- reload the package everytime you start a new session.
- `package::function()`
 - `ggplot2::ggplot()`



First steps

- Question: Do cars with big engines use more fuel than cars with small engines?
- Data: The `mpg` data frame in `ggplot2` (`ggplot2::mpg`)
 - `mpg` contains observations collected by the US Environmental Protection Agency on 38 models of car, 1999-2008.
 - A data frame with 234 rows and 11 variables:
- Variables:
 - `manufacturer`: manufacturer name
 - `model`: model name
 - `displ`: engine displacement, in litres
 - `year`: year of manufacture
 - `cyl`: number of cylinders
 - `trans`: type of transmission
 - `drv`: the type of drive train, where f = front-wheel drive, r = rear wheel drive, 4 = 4wd



The mpg data

1 mpg



manufacturer <chr>	model <chr>	displ <dbl>	year <int>	cyl <int>	
audi	a4	1.8	1999	4	
audi	a4	1.8	1999	4	
audi	a4	2.0	2008	4	
audi	a4	2.0	2008	4	
audi	a4	2.8	1999	6	
audi	a4	2.8	1999	6	
audi	a4	3.1	2008	6	
audi	a4 quattro	1.8	1999	4	
audi	a4 quattro	1.8	1999	4	
audi	a4 quattro	2.0	2008	4	

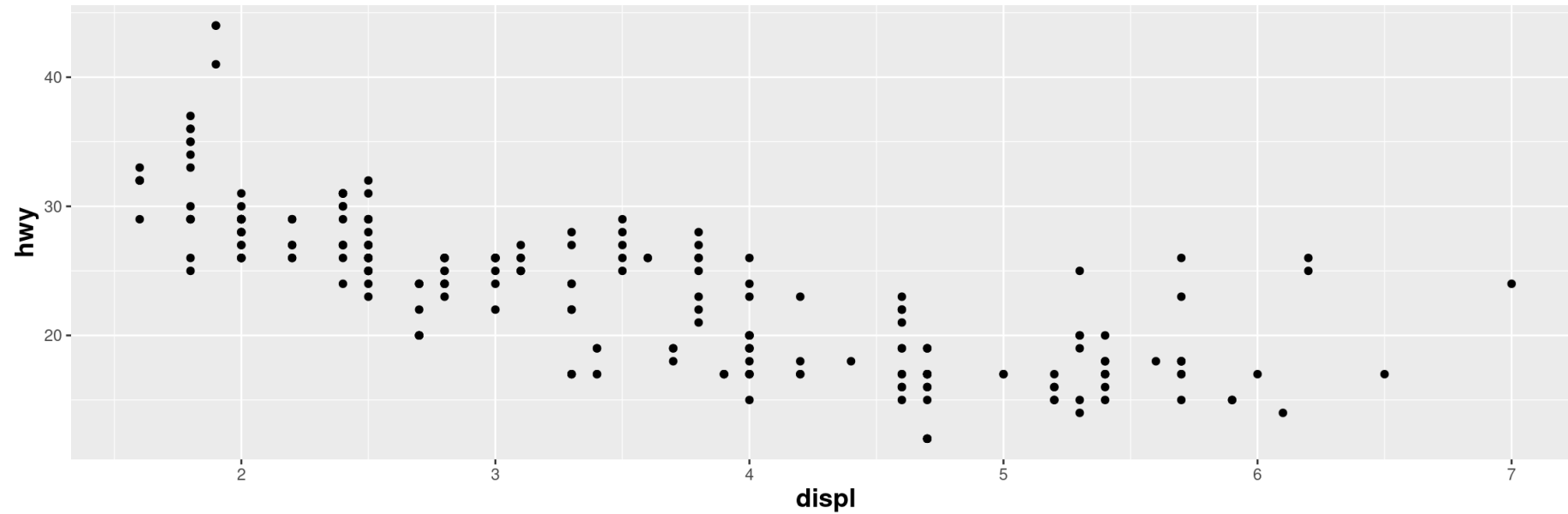
1-10 of 234 rows | 1-5 of 11 columns

Previous 1 2 3 4 5 6 ... 24 Next



Creating a plot

```
1 ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) +
4   # the y position
5   geom_point() + # the point geometric shape
6   # Adjust axis titles' front size
7   theme(axis.title=element_text(size=14, face="bold"))
```

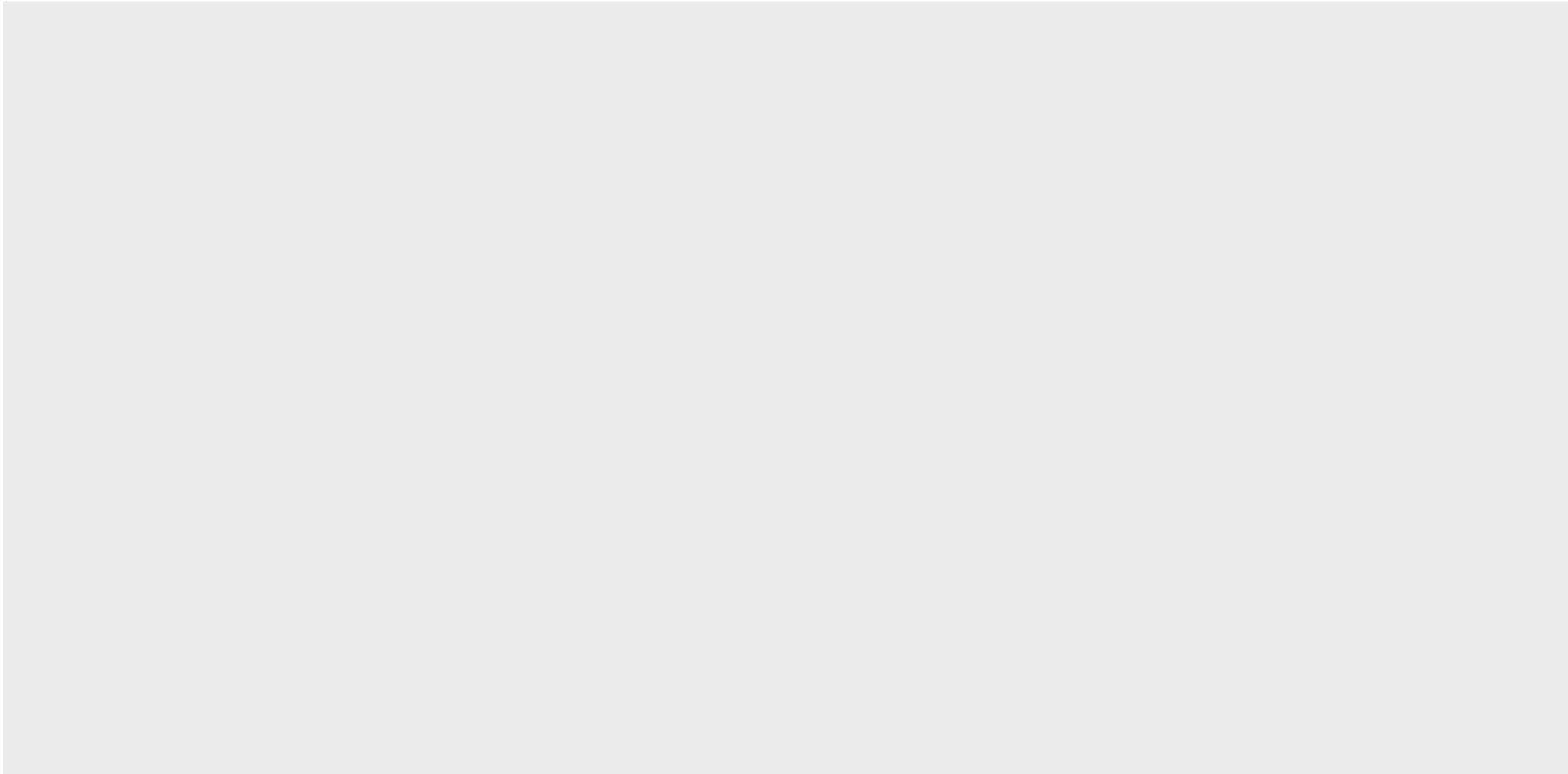


Declaring Data



Declaring data: method 1

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



Declaring data: method 2

- Pipe data into `ggplot()` using the pipe operator: `%>%`

```
1 mpg %>% # data piped into  
2   ggplot() # initiating plot
```



Exercises

1. Run `ggplot(data = mpg)`. What do you see?
2. How many rows are in `mpg`? How many columns?
3. What does the `drv` variable describe? Read the help for `?mpg` to find out.
4. Make a scatterplot of `hwy` vs `cyl`.
5. What happens if you make a scatterplot of `class` vs `drv`? Why is the plot not useful?



Aesthetic Mapping

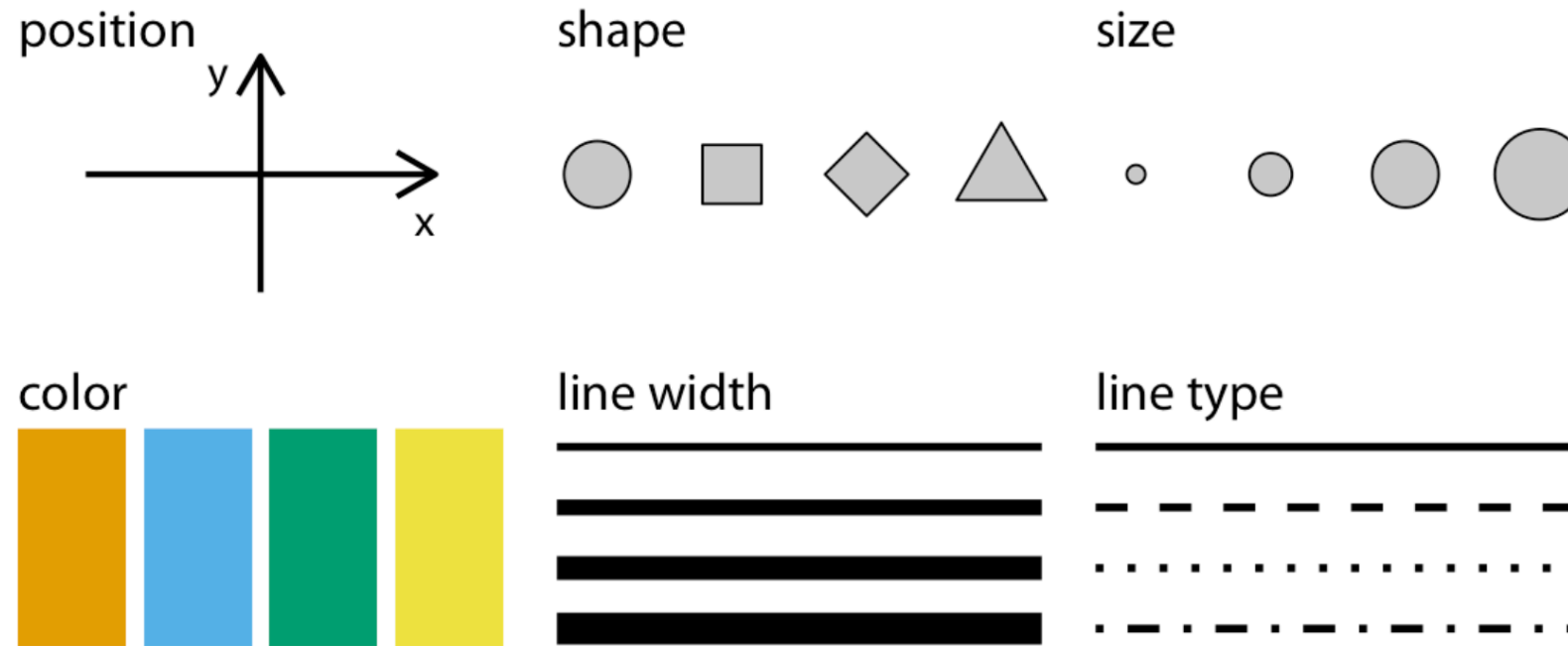


Aesthetic mappings

- An aesthetic is a visual property of the objects in your plot.
- Aesthetics include things like the *position*, *size*, the *shape*, or the *color* of your points.
- Mapping: variables are ‘mapped’ to (represented by) aesthetics.



A main pool of Aesthetics



A main pool of aesthetics

Note: This figure is from Wilke's *Fundamentals of data Visualization*.

`aes()` means 'Ask'

- `aes()`: What variables are we asking the aesthetic (color, position, shape, etc.) to represent?
- `aes(color=gender)`: 'Please represent the variable `gender` for me using different colors.'



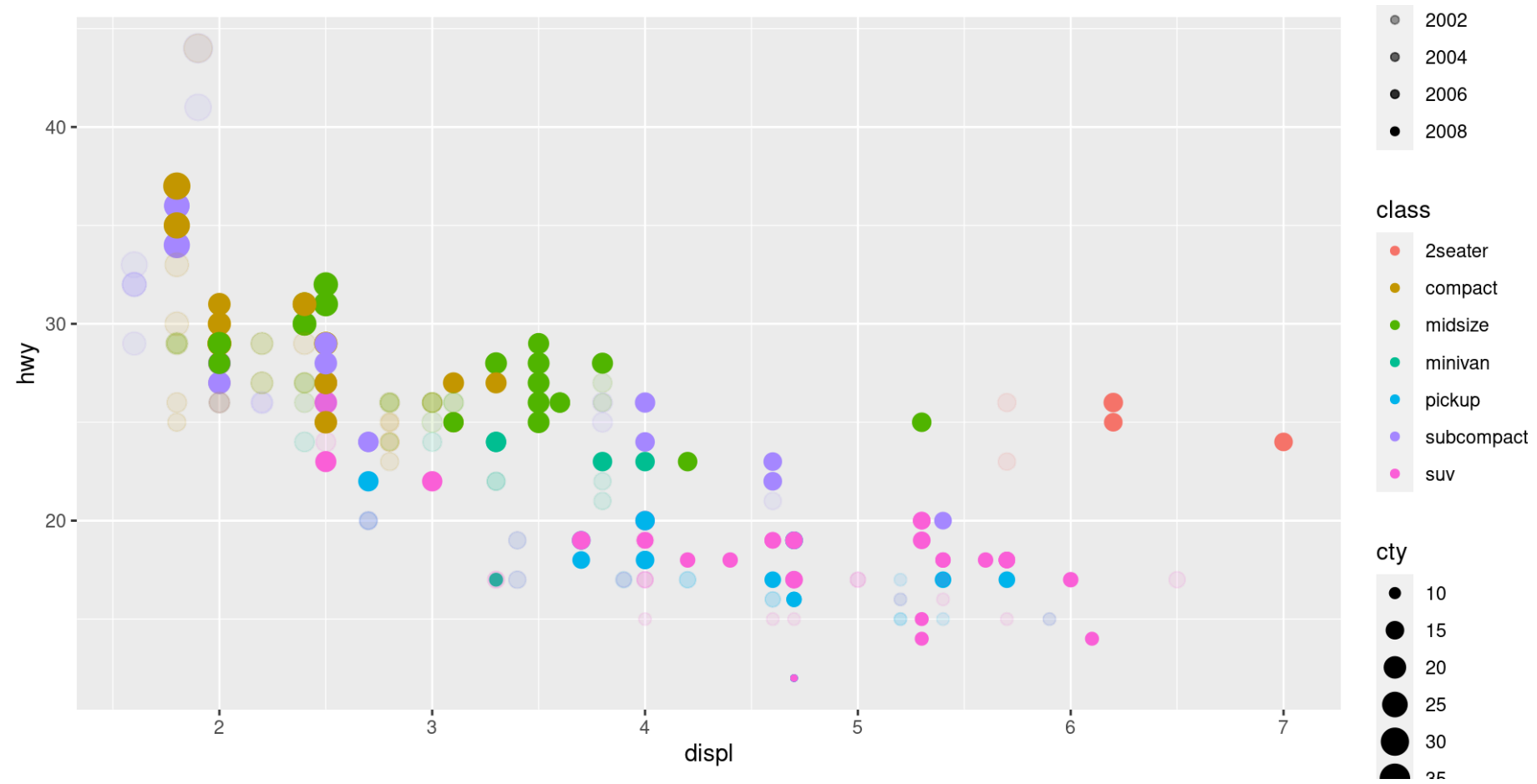
More `aes` mappings

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() +
5   #the point geometric shape, the above aes are required and
6   #the below are optional
7   #theme(axis.title=element_text(size=14,face="bold"))+
8   aes(color = class) + # Color for type of car
9   #aes(shape = class) +
10  #ggplot2 will only use six shapes at a time. By default,
11  #additional groups will go unplotted when using 'shape'.
12  aes(size = cty) + # Size for city miles per gallon
13  aes(alpha = year) # transparency for year of manufacture
```



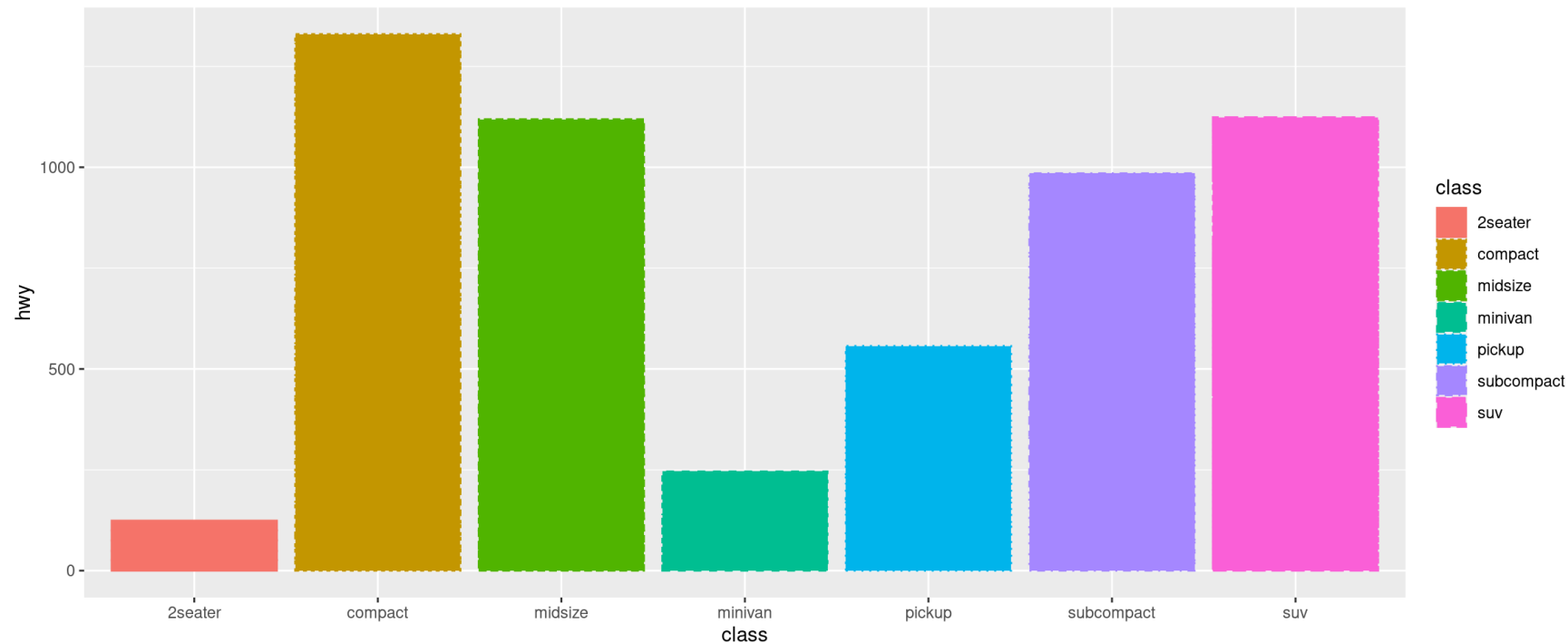
Plot

```
1 print(mpg_plot)
```



Other `aes` for other geometric objects

```
1 mpg %>% # data piped into
2   ggplot() + # initiating plot
3   aes(x = class) + #categorical variable
4   aes(y = hwy) +
5   geom_col() + #Use `geom_col` to creat a column geometry
6   aes(color = class) +
7   aes(fill = class) + # new aes 'fill'
8   aes(linetype = class) #new aes 'linetype'
```



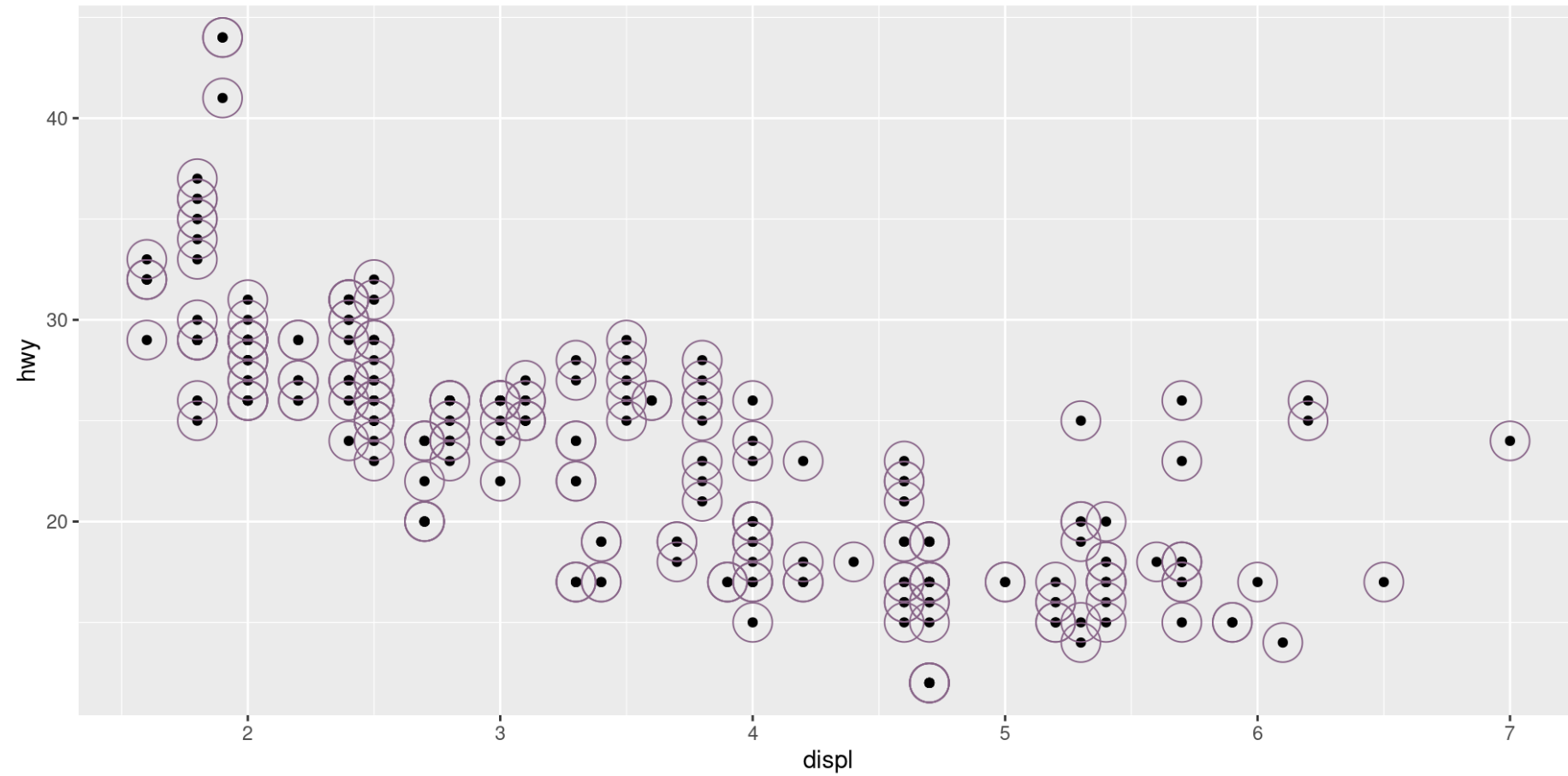
Unmapped aesthetics

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() + # the point geometric shape
5   #Another geom layer with aesthetics
6   #that don't do representation
7   geom_point(
8     color="plum4",
9     size=8,
10    shape=21
11  )
```



Plot

```
1 print(mpg_plot)
```

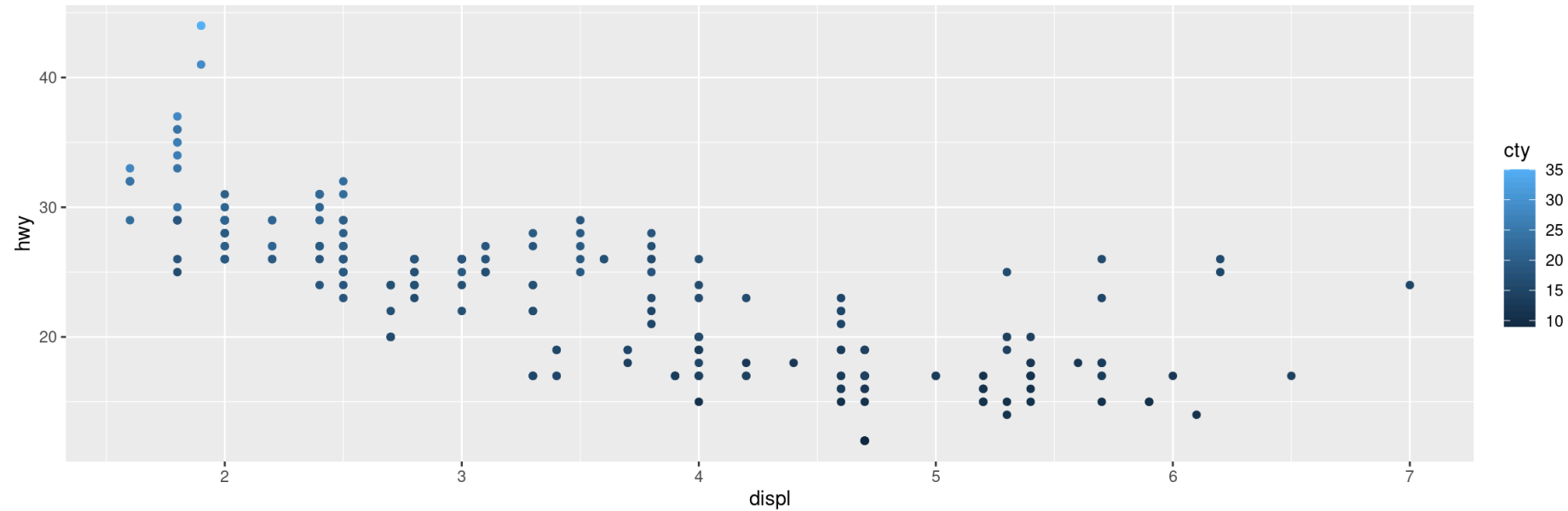


Exercises:

1. Look at the `help` for `geom_text` (`?geom_text`). What are the required aesthetics?
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`?
3. Map a continuous variable to color, size, and shape. How do these aesthetics behave differently for categorical vs. continuous variables?



```
1 ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) +
4   aes(color = cty) + #color, size and shape
5   # the y position
6   geom_point() # the point geometric shape
```



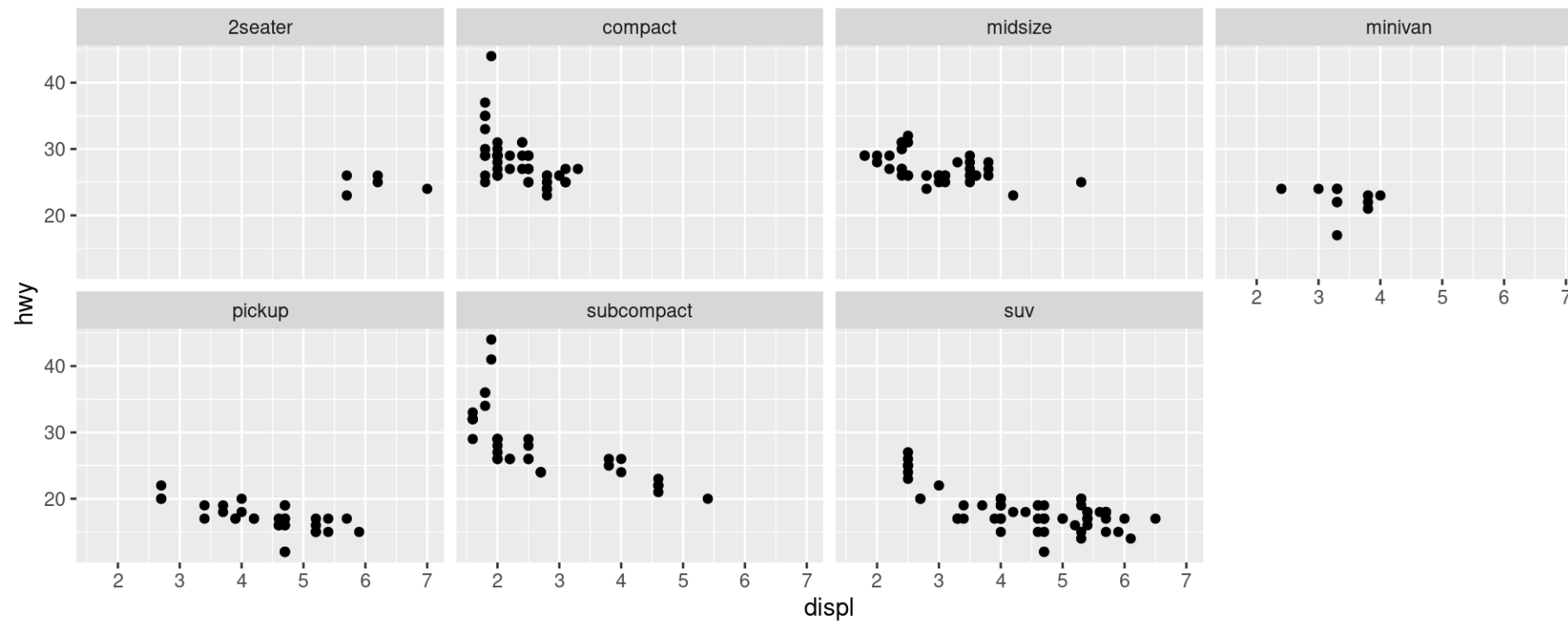
Facets



Facets: facet_wrap

- facet your plot by a single variable

```
1 ggplot(data = mpg) +  
2   geom_point(mapping = aes(x = displ, y = hwy)) +  
3   # ~ followed by a discrete variable  
4   facet_wrap(~ class, nrow = 2)
```

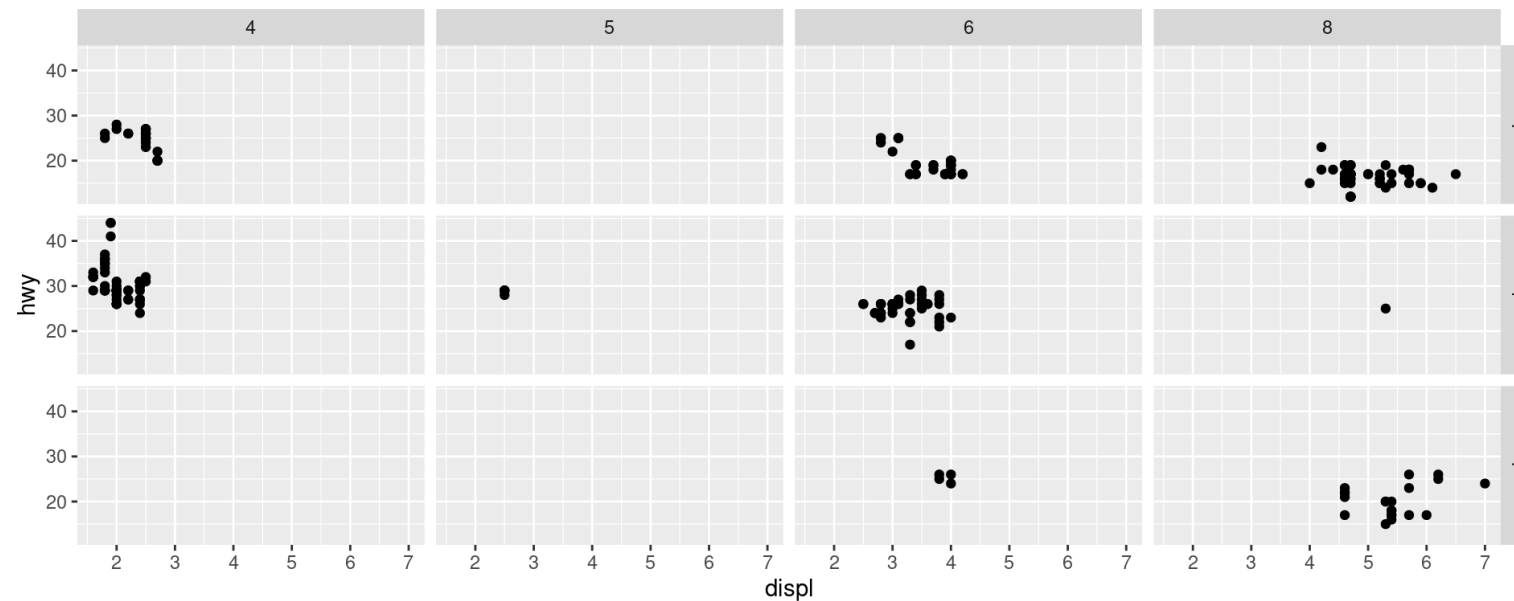


Facets: facet_grid

- facet your plot on the combination of two variables

```
1 ggplot(data = mpg) +
2   geom_point(mapping = aes(x = displ, y = hwy)) +
3   # two variable names separated by a ~
4   facet_grid(drv ~ cyl)
```

```
1 #facet_grid(. ~ cyl) #not facet in the rows
2 #facet_grid(drv~.) #not facet in the rows
```



Geometric Object



A complete sentence of `ggplot`

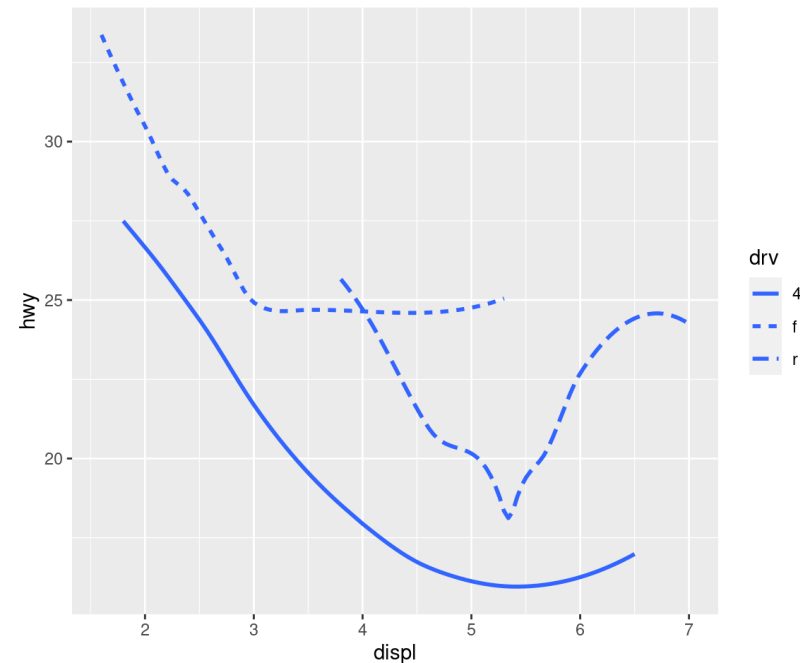
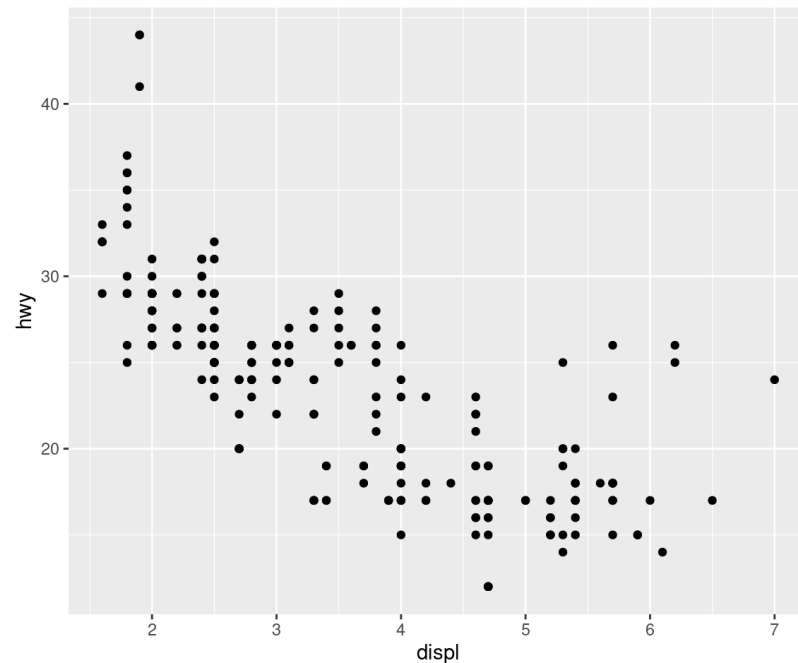
- `data` + `aes` + `geom`
- plots
- Nouns: geometric objects
 - `geom_point()`
 - `geom_col()`
 - `geom_line()`
 - `geom_text()`
 - `geom_segment()`
 - `geom_smooth()`
 - `geom_bar()`
 - etc.
- The conditional mood: `geom` specific data and aesthetic mapping



Different Geoms

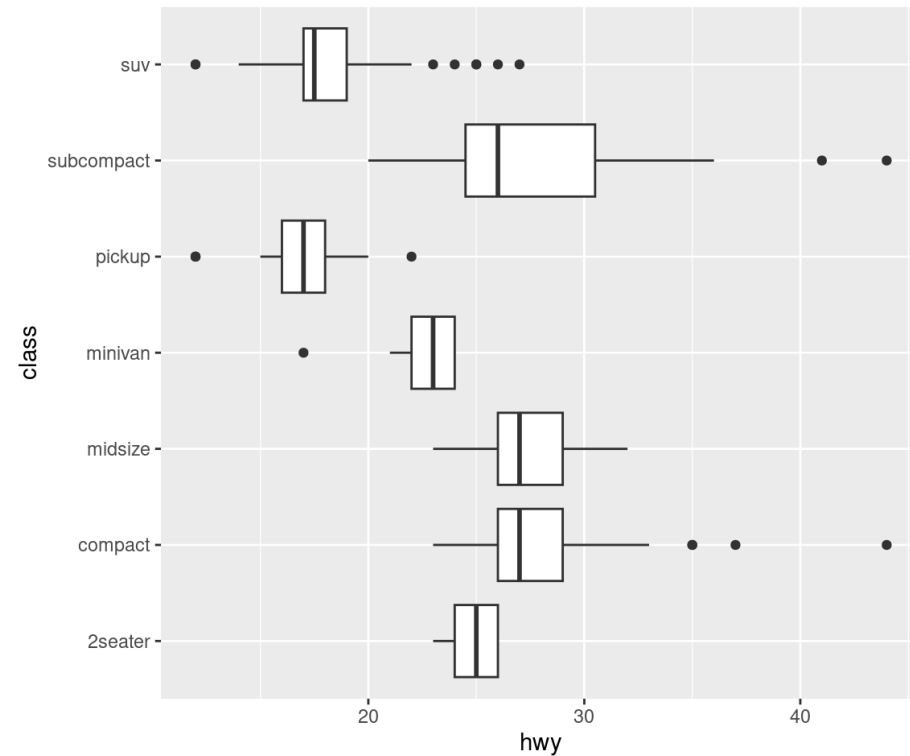
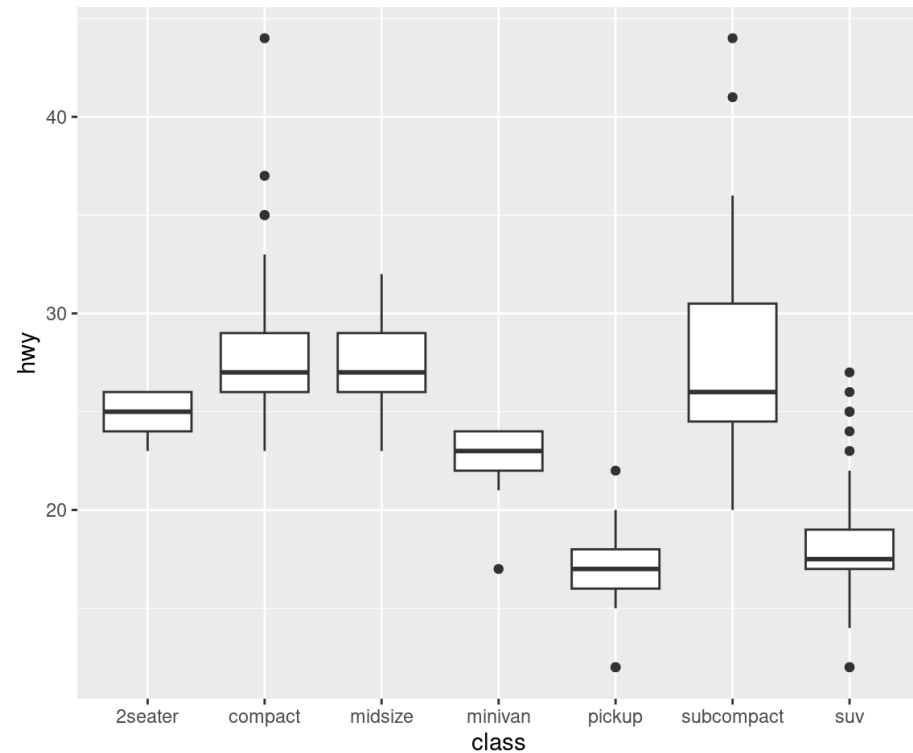
Compare with

```
1 library(grid)
2 library(gridExtra)
3 p1=ggplot(data = mpg) +
4   geom_point(mapping = aes(x = displ, y = hwy)) # geom_point
5 p2=ggplot(data = mpg) +
6   geom_smooth(mapping = aes(x = displ, y = hwy, linetype=drv), se=FALSE) # geom_smooth
7 grid.arrange(p1, p2, ncol = 2)
```



Another Example: boxplot

```
1 p1=ggplot(data = mpg, mapping = aes(x = class, y = hwy)) +
2   geom_boxplot()
3 p2=ggplot(data = mpg, mapping = aes(x = class, y = hwy)) +
4   geom_boxplot() +
5   coord_flip() #switches the x and y axes
6 grid.arrange(p1, p2, ncol = 2)
```



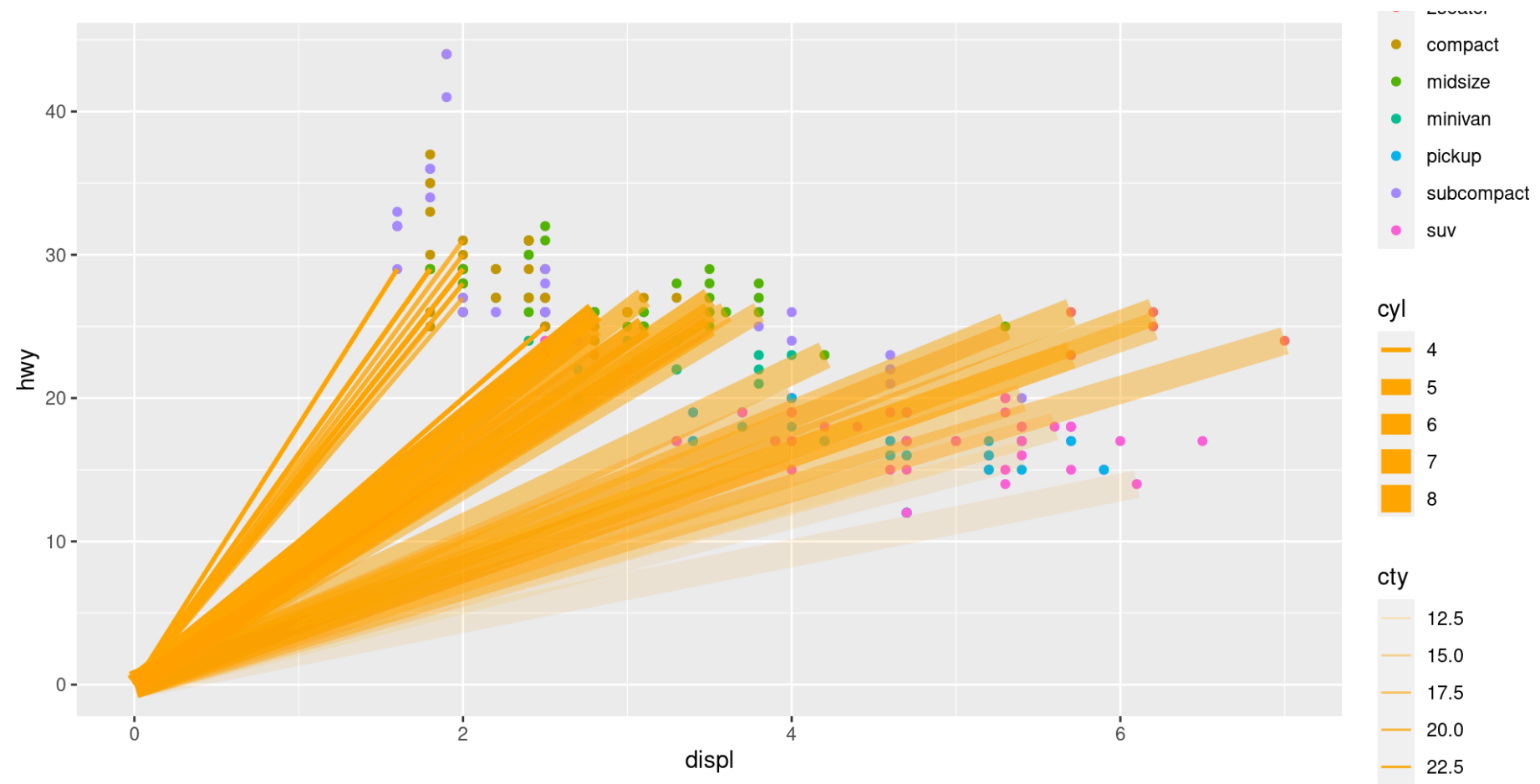
Going local with data and aesthetics

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ, y = hwy) +
3   geom_point() +
4   aes(color = class) +
5   #xend and yend are required for geom_segment
6   #like creating a column with a single value
7   aes(xend = 0) + aes(yend = 0) +
8   #geom_segment() draws a straight line
9   #between points (x, y) and (xend, yend)
10  geom_segment(
11    #geom specific data, using 'subset' to select data
12    data = subset(mpg, fl=="p"),
13    # geom specific (local) aesthetics
14    aes(size = cyl, alpha = cty),
15    color = "orange"
16  )
```



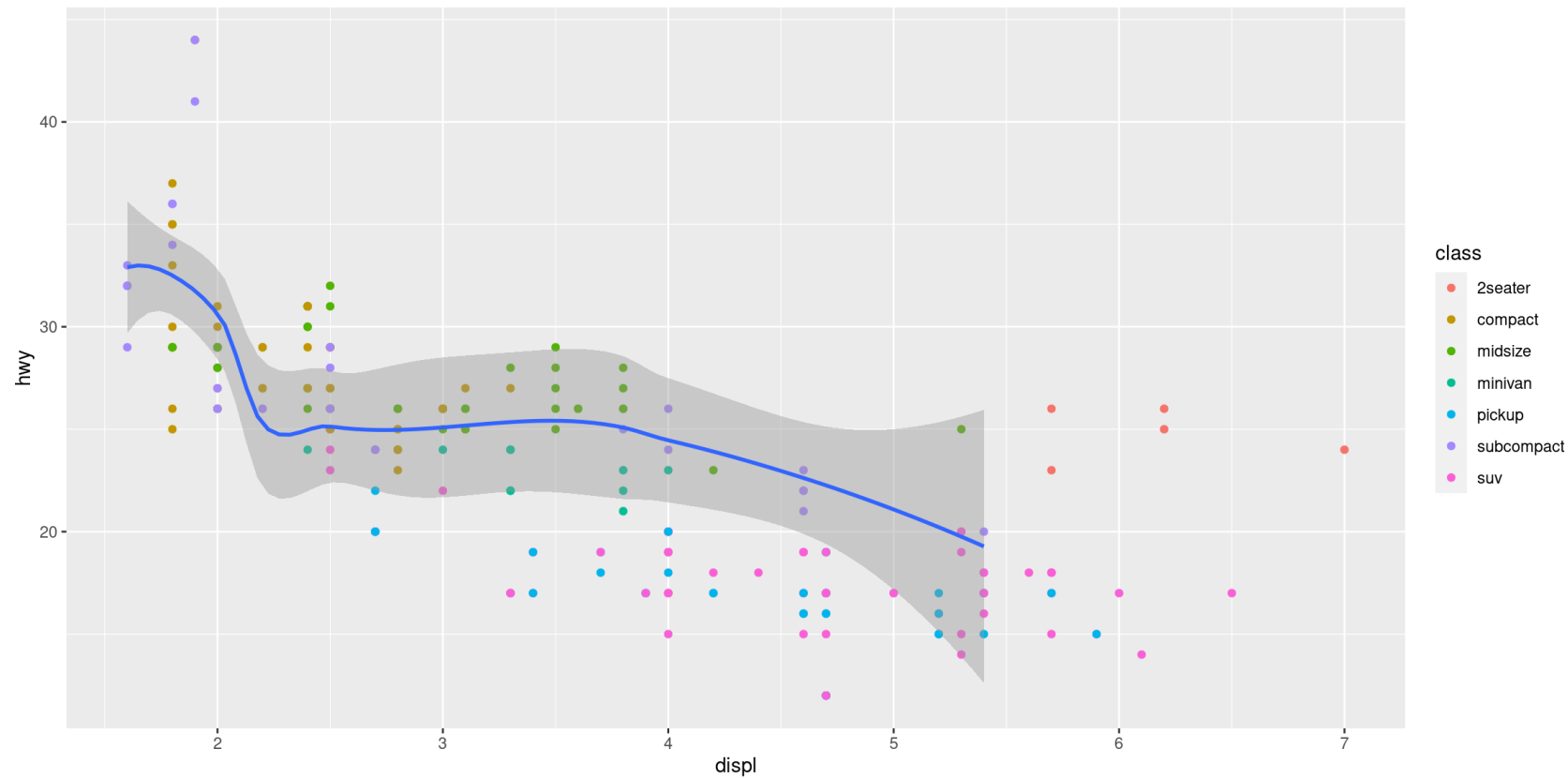
Plot

```
1 print(mpg_plot)
```



Another Example

```
1 ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) + #global aes
2   geom_point(mapping = aes(color = class)) + #local aes
3   #local data and aes
4   geom_smooth(data = filter(mpg, class == "subcompact"),
5               se = TRUE) #se: standard error
```



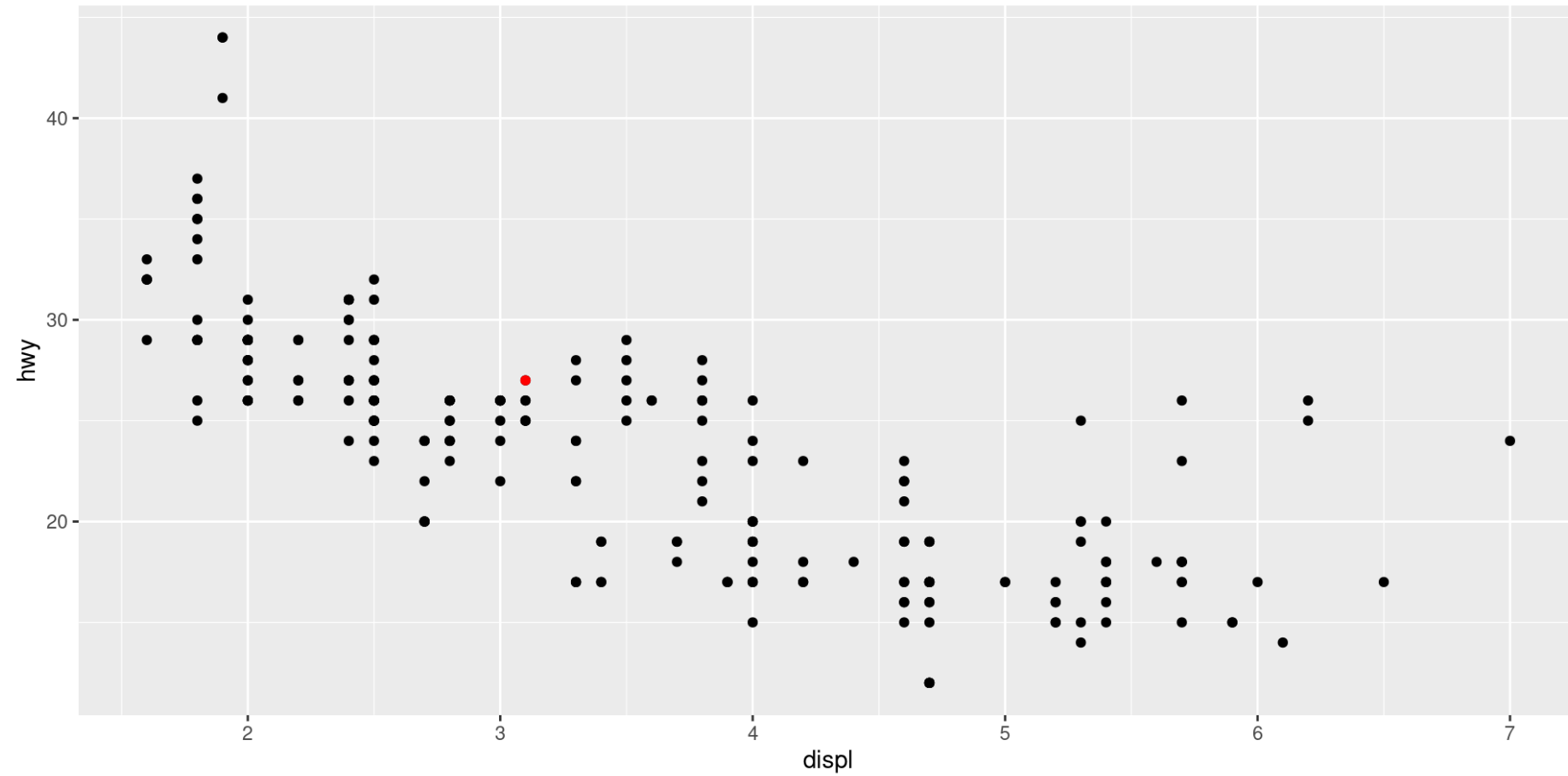
Annotation 1

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() + # the point geometric shape
5   annotate(geom="point",
6           x=3.1,
7           y=27,
8           color="red")
```



Plot 1

```
1 print(mpg_plot)
```



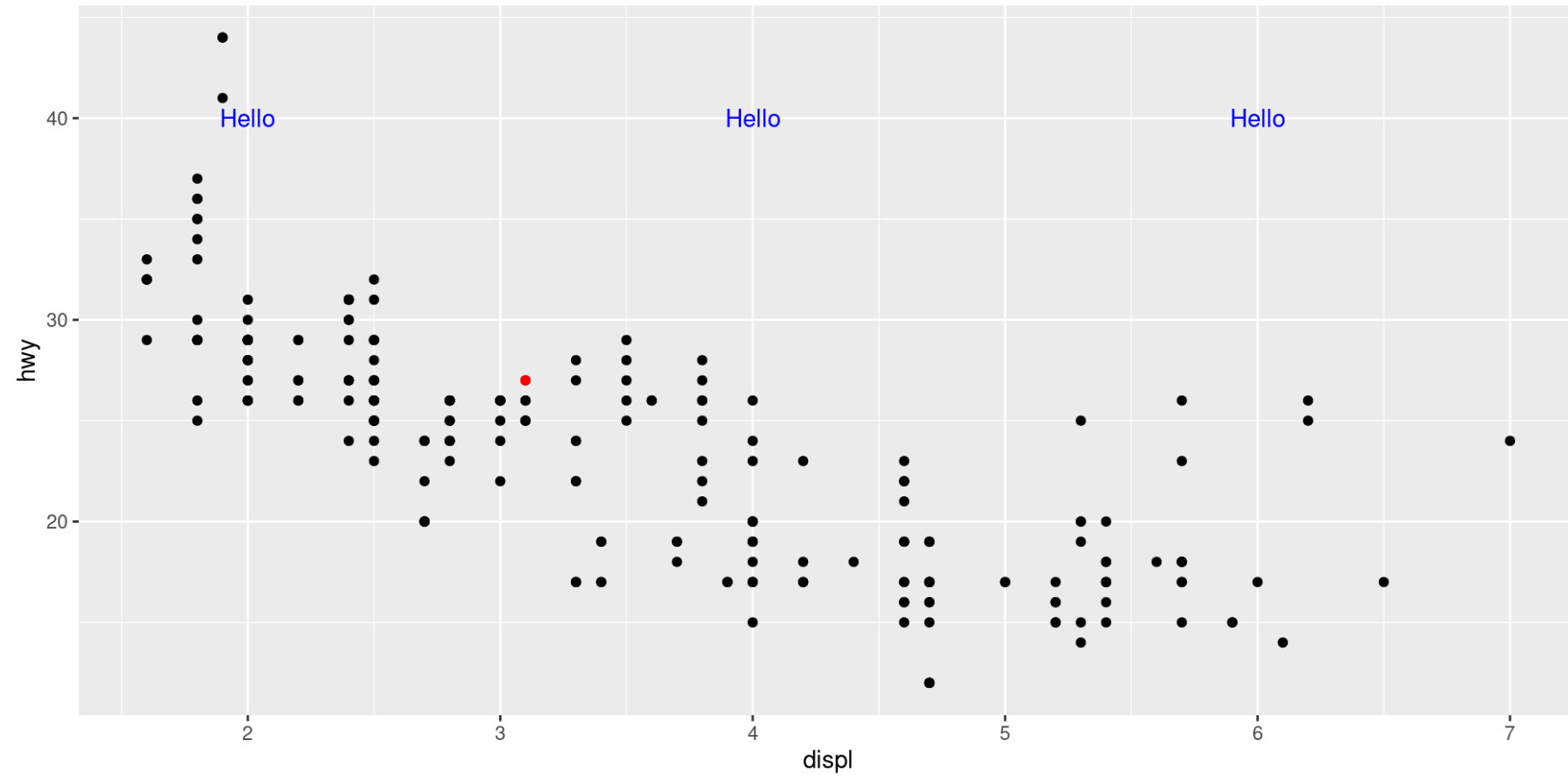
Annotation 2

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() + # the point geometric shape
5   annotate(geom="point",
6           x=3.1,
7           y=27,
8           color="red")+
9   annotate(geom="text",
10          x=c(2, 4, 6),
11          y=40,
12          label="Hello",
13          color="blue")
```



Plot 2

```
1 print(mpg_plot)
```



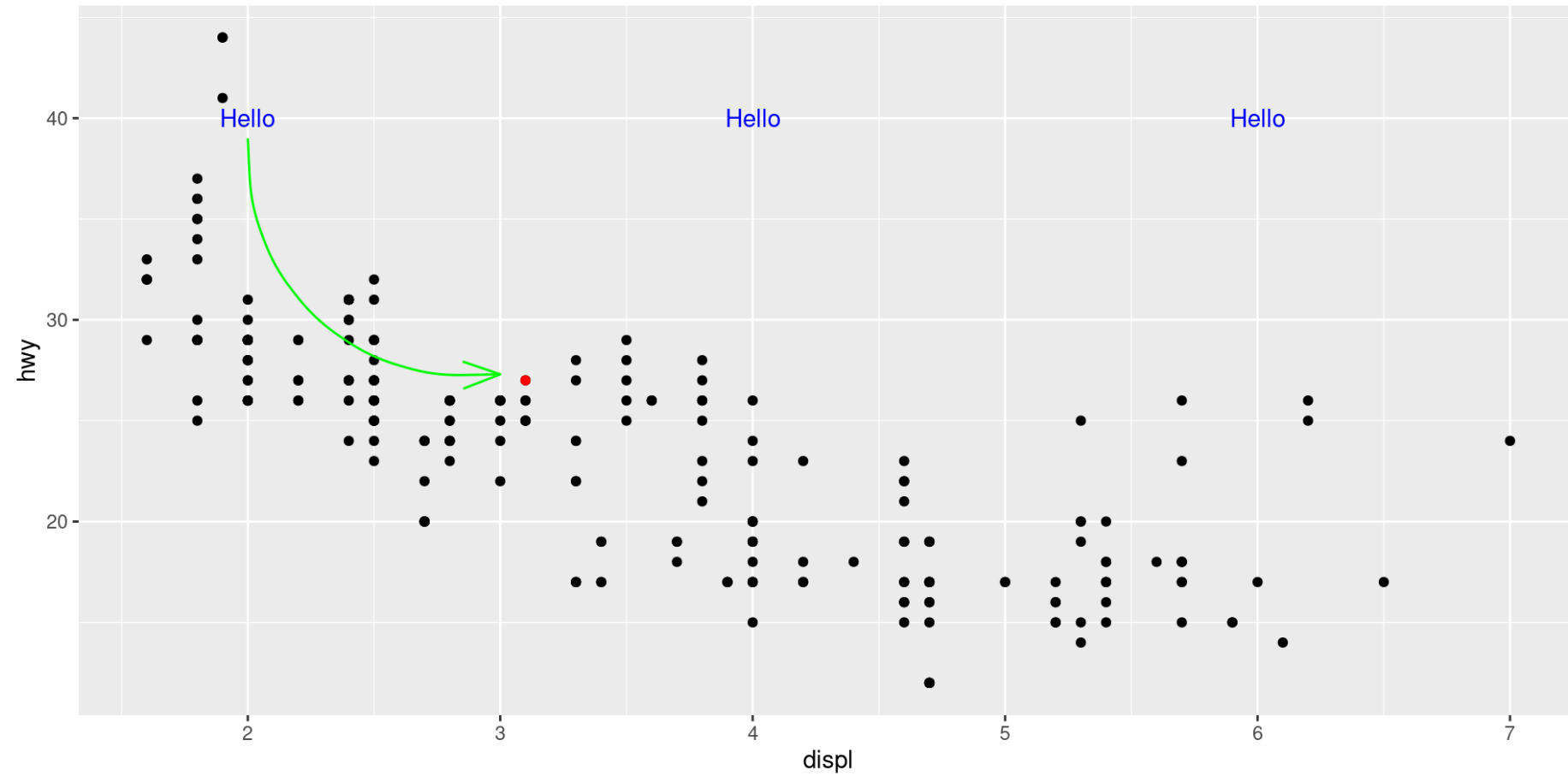
Annotation 3

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() + # the point geometric shape
5   annotate(geom="point",
6           x=3.1,
7           y=27,
8           color="red")+
9   annotate(geom="text",
10          x=c(2, 4, 6),
11          y=40,
12          label="Hello",
13          color="blue")+
14   annotate(geom="curve",
15          x=2,
16          y=39,
17          xend=3,
18          yend=27.3,
19          color="green",
20          arrow=arrow(angle=20))
```



Plot 3

```
1 print(mpg_plot)
```



Annotation 4

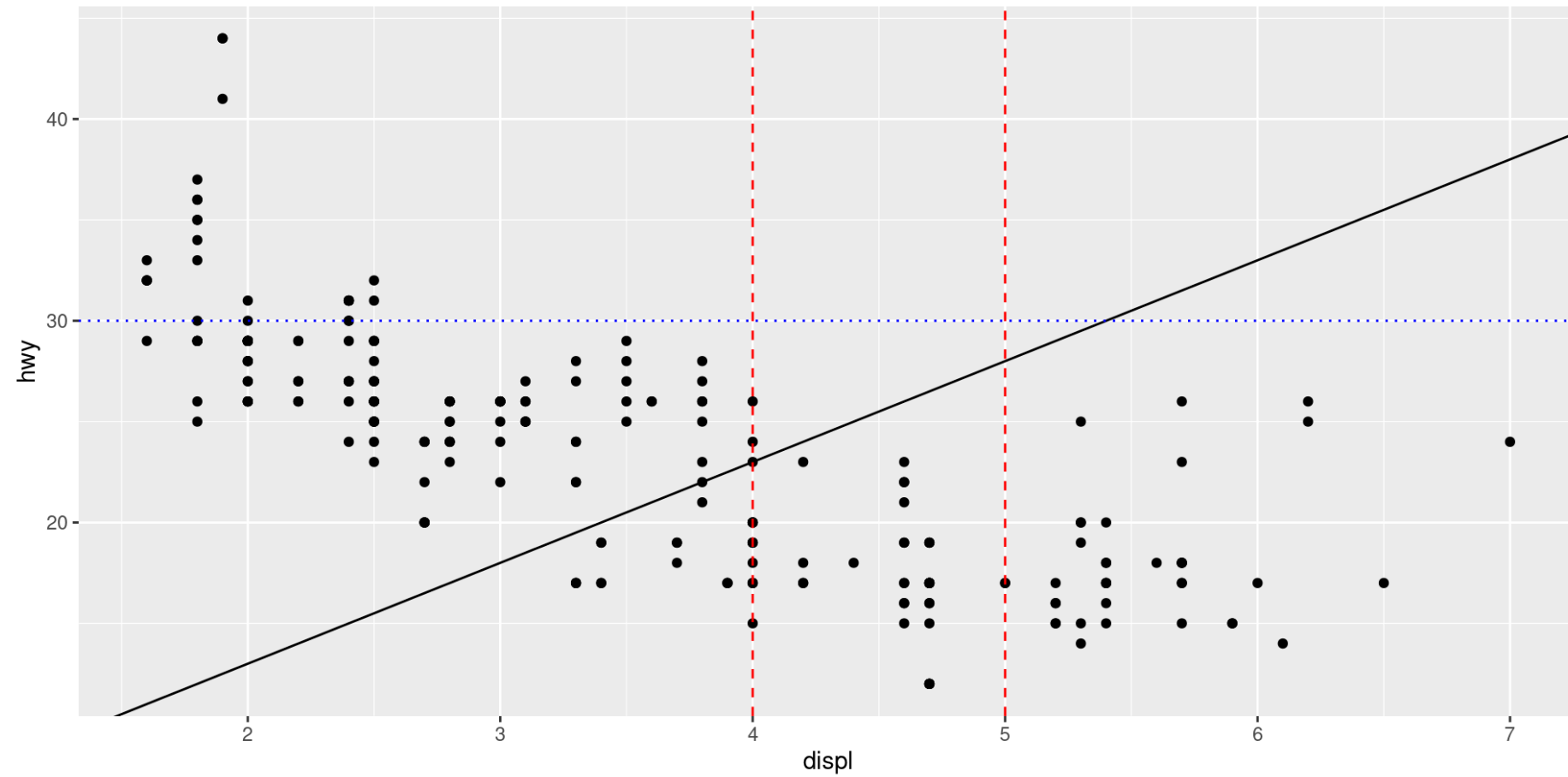
- use `geom_abline`, `geom_hline`, and `geom_vline`

```
1 mpg_plot= ggplot(data = mpg) + # the dataset
2   aes(x = displ) + # the x position
3   aes(y = hwy) + # the y position
4   geom_point() + # the point geometric shape
5   geom_abline(slope=5, intercept=3) +
6   geom_hline(yintercept= 30,
7             linetype="dotted", color="blue")+
8   geom_vline(xintercept=c(4,5),
9             linetype="dashed", color="red")
```



Plot 4

```
1 print(mpg_plot)
```



Exercises

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

```
1 ggplot(data = mpg) +  
2   geom_point(mapping = aes(x = displ,  
3                             y = hwy, color = "blue"))
```

2. What happens if you map an aesthetic to something other than a variable name, like `aes(colour = displ < 5)`? Note, you'll also need to specify x and y.



Interactive data visualisation (optional)

- R package: *Shiny*
- can host standalone apps on a webpage
 - Example: **Life Expectancy** using data from the
- can embed them in R Markdown documents or build dashboards.

