

Southwestern University of Finance and Economics

Multivariate Statistical Analysis

Homework1

Student:
Wenjie Lan

Supervisor:
Prof Xu

An Assignment submitted:

Data Visualization With ggplot2

March 11, 2022

1 Exercise

PROBLEM 1.

SOLUTION.

This code creates an empty plot which plays as an "background".

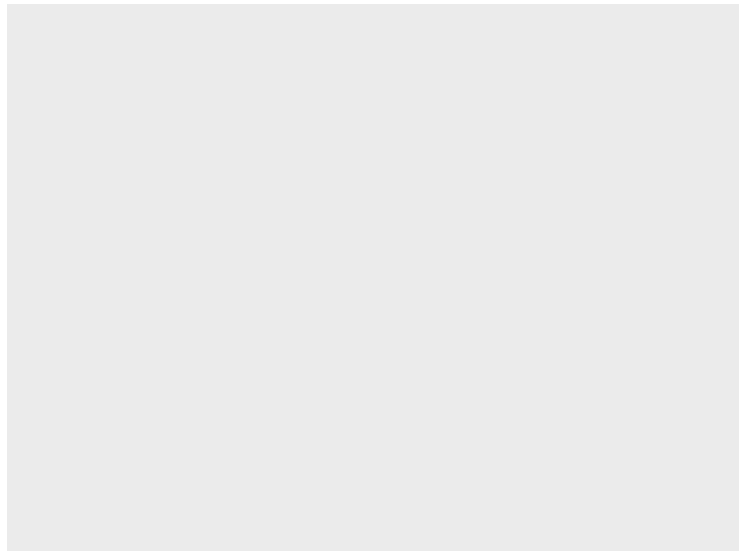


Figure 1: 1.1

PROBLEM 2.

SOLUTION.

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

PROBLEM 3.

SOLUTION.

The drv variable categorizes cars into front-wheels, rear-wheels, and four-wheel.

PROBLEM 4.

SOLUTION.

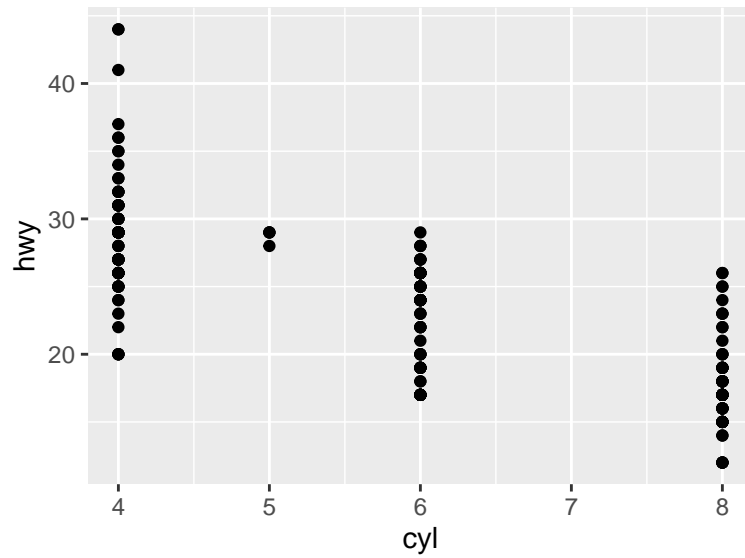


Figure 2: 1.4

PROBLEM 5.

SOLUTION.

Because drv and class are both categorical variables, and it only results in a few isolated points in the picture.

```
1 ggplot(mpg, aes(x = cyl, y = hwy)) +
2   geom_point()
```

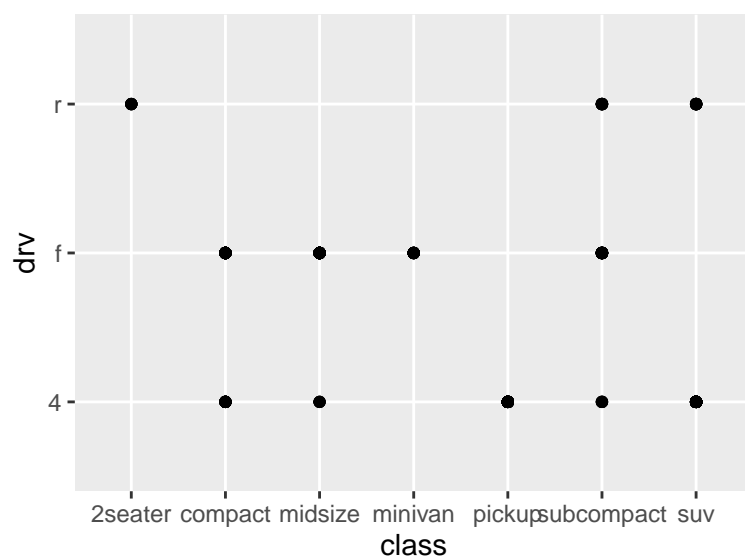


Figure 3: 1.5

2 Exercise

PROBLEM 6.

SOLUTION. The code "color = "blue"" should be outside of the aes argument.

Correct code:

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

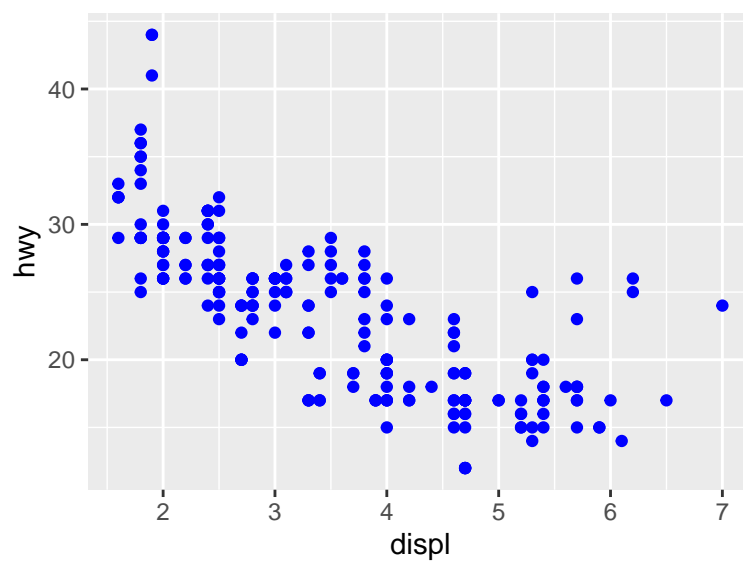


Figure 4: 2.1

PROBLEM 7.

SOLUTION.

Categorical variables:

manufacturer

model

trans

drv

fl

class

Continuous variables:

displ

year

cyl

cty

hwy

PROBLEM 8.

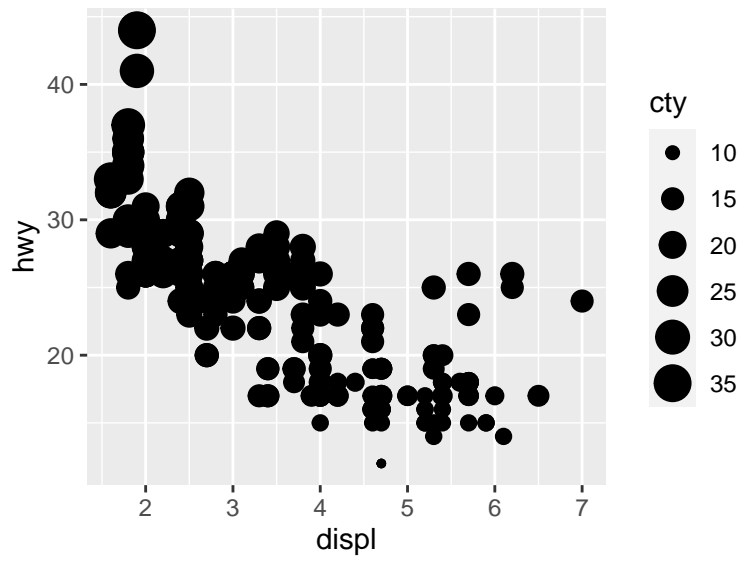
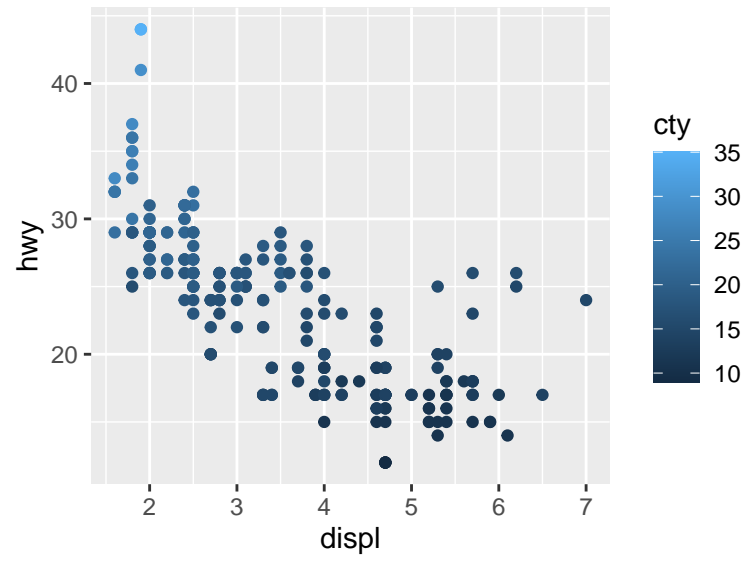
SOLUTION.

Color:appears as gradient color.

Size:appear as different size continuously.

shape:error since cty is a continuous variable.

```
1 ggplot(mpg, aes(x = displ, y = hwy, color = cty)) +  
2   geom_point()  
3 ggplot(mpg, aes(x = displ, y = hwy, size = cty)) +  
4   geom_point()  
5 ggplot(mpg, aes(x = displ, y = hwy, shape = cty)) +  
6   geom_point()
```



```
## Error in 'scale_f()':## ! A continuous variable can not be  
mapped to shape
```

PROBLEM 9.

SOLUTION.

```
1 ggplot(mpg, aes(x = displ, y = hwy, colour = hwy,  
2   size = displ)) +  
   geom_point()
```

3.

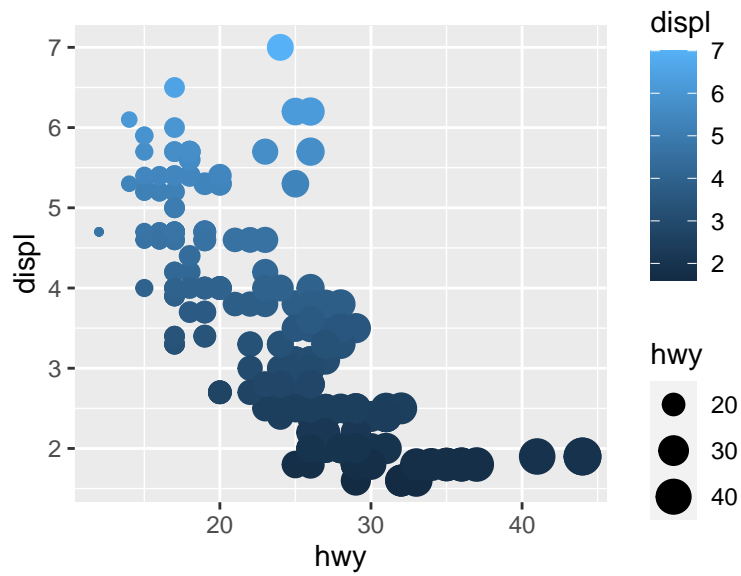


Figure 5: 2.4

PROBLEM 10.

SOLUTION.

It can change points' fill color and border.

3 Exercise

PROBLEM 11.

SOLUTION.

The continuous variable appears as categorical form.

```
1 ggplot(mpg, aes(x = displ , y = hwy)) +
2   geom_point() +
3   facet_grid(. ~ cty)
```

3.

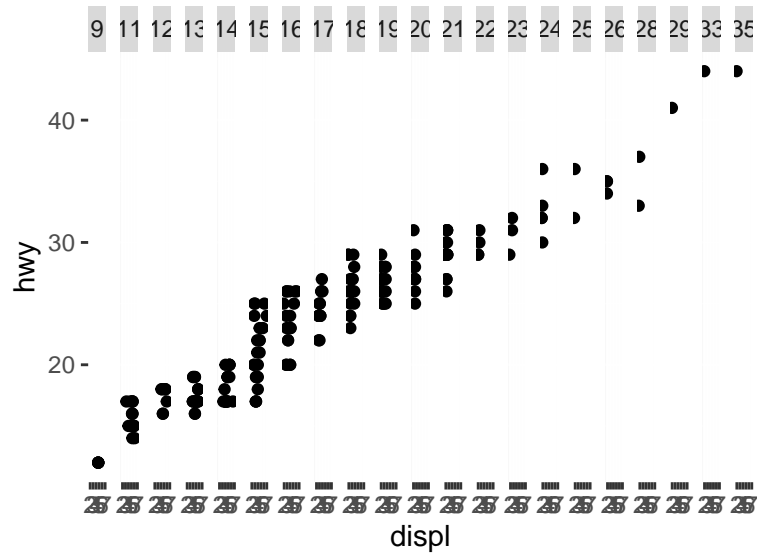


Figure 6: 3.1

PROBLEM 12.

SOLUTION.

```

1 ggplot(data = mpg) +
2   geom_point(mapping = aes(x = displ, y = hwy))
3 ggplot(data = mpg) +
4   geom_point(mapping = aes(x = displ, y = hwy)) +
5   facet_grid(drv ~ cyl)

```

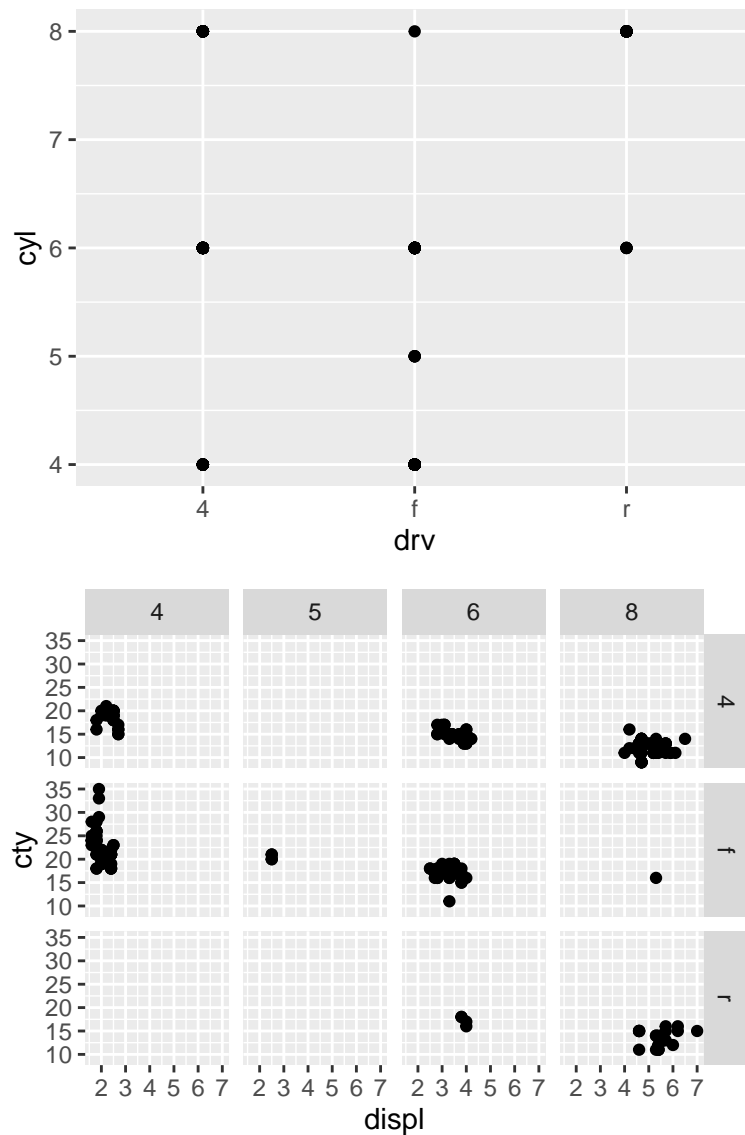


Figure 7: 3.2

PROBLEM 13.

SOLUTION. The "." ignores dimensions and facet by values on x or y axis according to where it is placed in the code.

```
1 ggplot(data = mpg) +
2   geom_point(mapping = aes(x = displ, y = hwy)) +
3   facet_grid(drv ~ .)
4 ggplot(data = mpg) +
5   geom_point(mapping = aes(x = displ, y = hwy)) +
6   facet_grid(. ~ cyl)
```

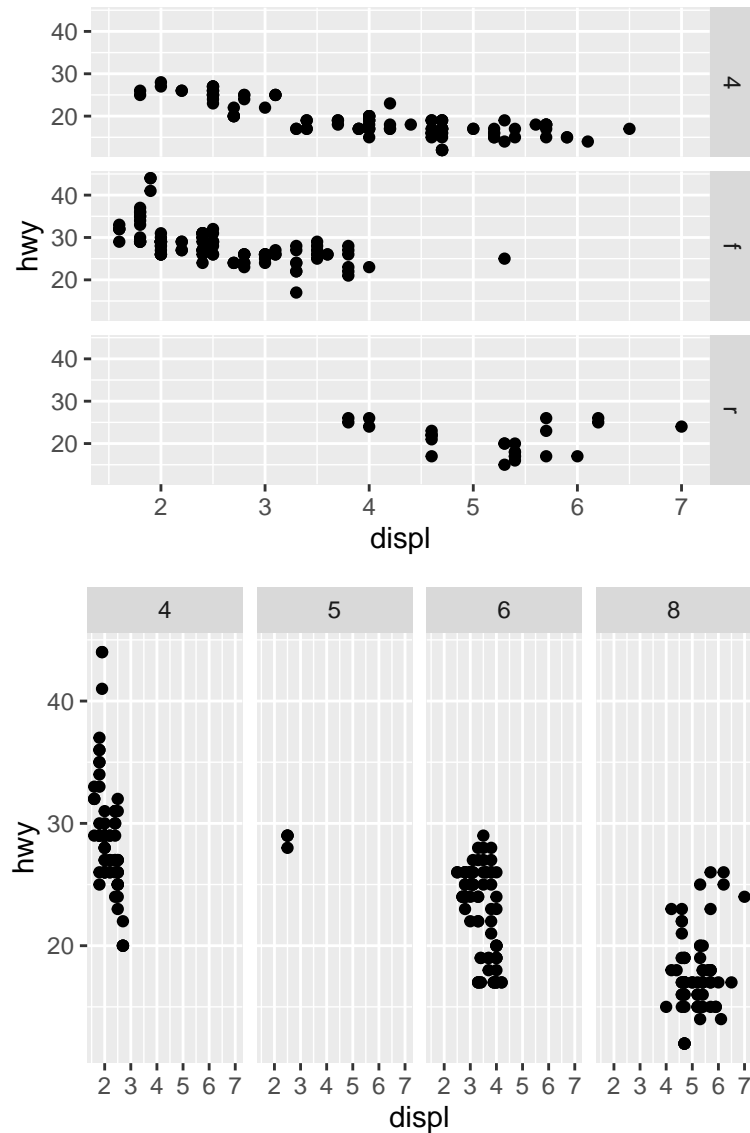


Figure 8: 3.3

PROBLEM 14.

SOLUTION.

Advantage: It is convenient to observe the condition within each catalog.

Disadvantage: It is inconvenient to compare different catalogs.

I would choose faceting instead of color when having a large dataset since colors can be difficult when distinguishing.

```

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

```
color = class))
```

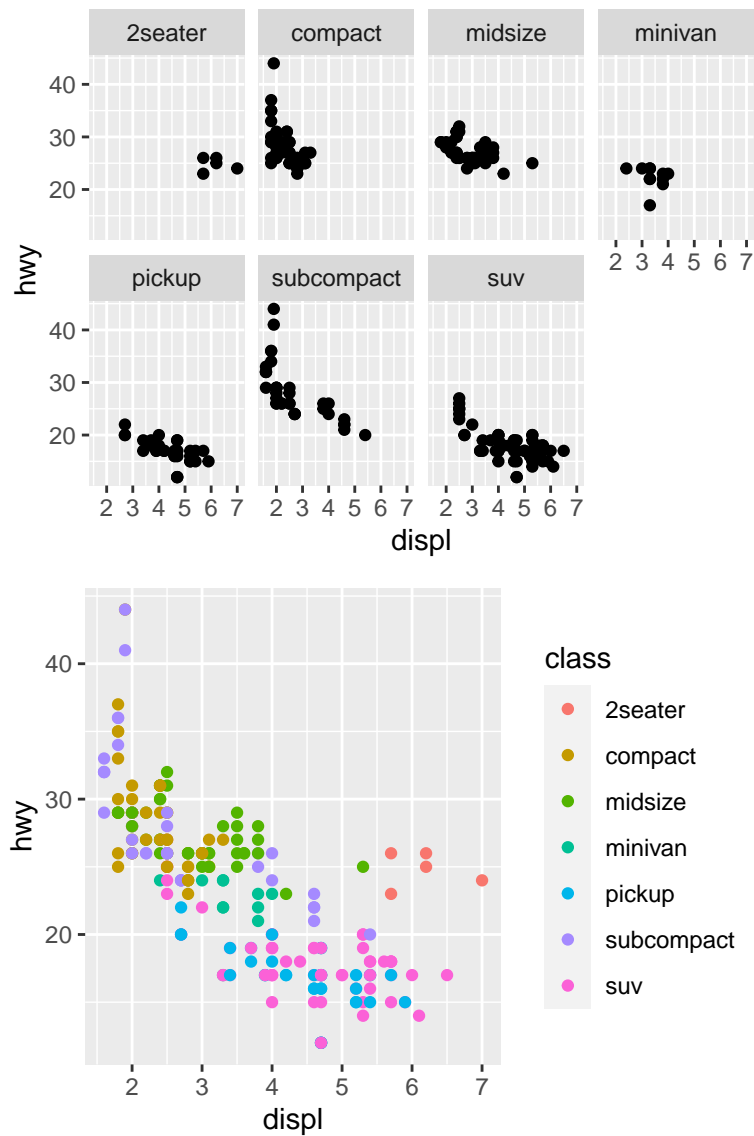


Figure 9: 3.4

PROBLEM 15.

SOLUTION.

"nrow (ncol)" determines the number of rows (columns) of facets we create. Because the number of rows and columns represent the number of variables.

PROBLEM 16.

SOLUTION.

When we use `facet_wrap` we usually put the variable with more unique levels in the

columns to create less columns, since it is not easy to display too many columns and inconvenient for we to observe data.

4 Exercise

PROBLEM 17.

SOLUTION.

Line chart: `geom_line()`

Boxplot: `geom_boxplot()`

Histogram: `geom_histogram()`

Area chart: `geom_area()`

PROBLEM 18.

SOLUTION.

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

'geom_smooth()' using method = 'loess' and formula 'y ~ x'

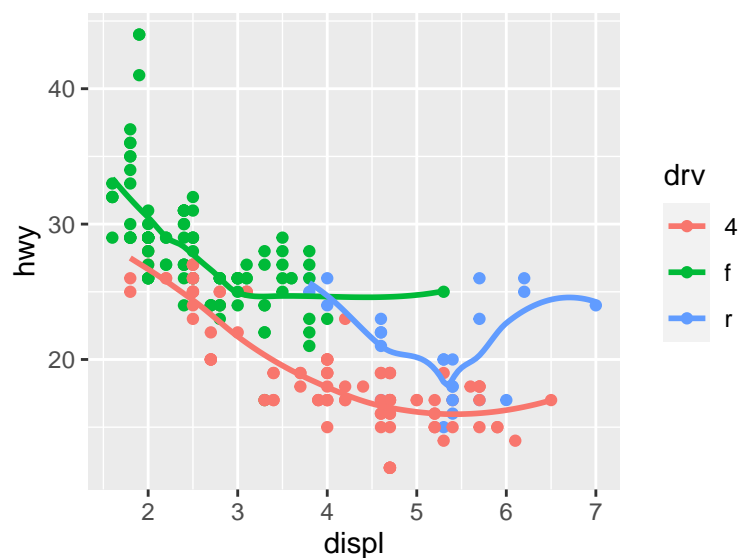


Figure 10: 4.2

PROBLEM 19.

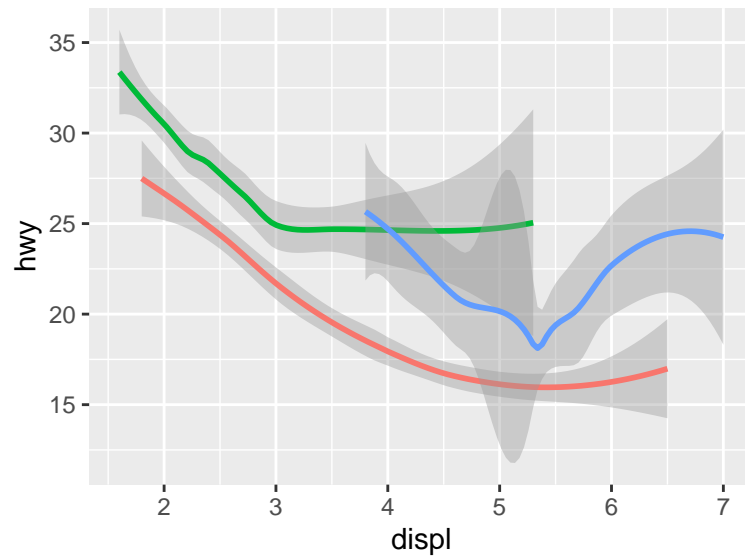
SOLUTION.

"show.legend = FALSE" help to remove the legend.

The author use it

```
1 ggplot(data = mpg) +  
2   geom_smooth(  
3     mapping = aes(x = displ, y = hwy, colour = drv)  
4     ,  
5     show.legend = FALSE  
6   )  
7 ggplot(data = mpg) +  
8   geom_smooth(  
9     mapping = aes(x = displ, y = hwy, colour = drv)  
10  )
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

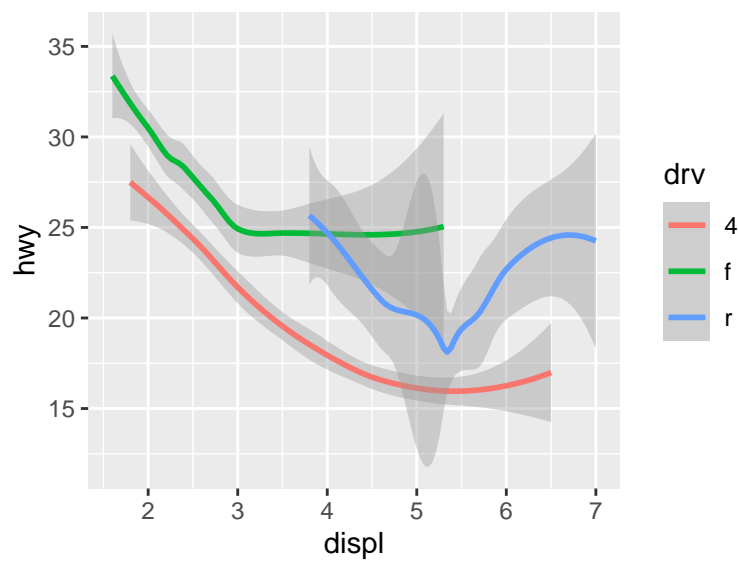


Figure 11: 4.3

PROBLEM 20.

SOLUTION.

It shows standard error bands by adding shadow to the lines.

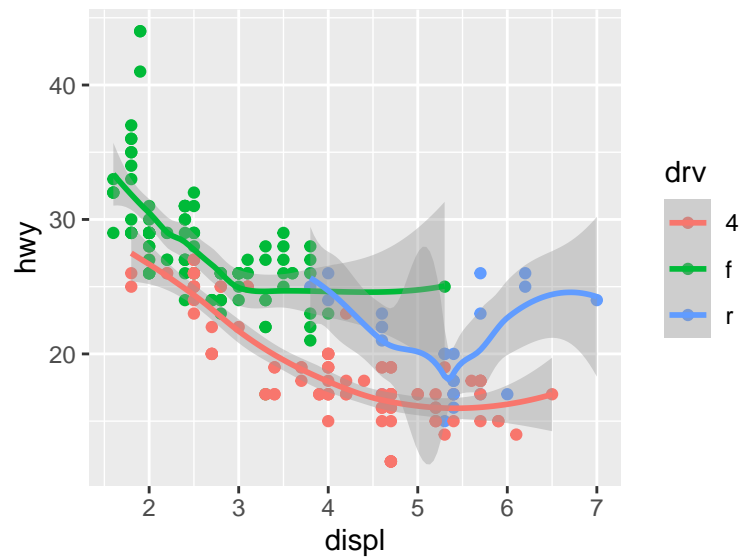
```
1 ggplot(data = mpg, mapping = aes(x = displ, y = hwy
2   , colour = drv)) +
  geom_point() +
```

```

3   geom_smooth(se = TRUE)
4   ggplot(data = mpg, mapping = aes(x = displ, y = hwy
5     , colour = drv)) +
6     geom_point() +
    geom_smooth(se = FALSE)

```

'geom_smooth()' using method = 'loess' and formula 'y ~ x'



'geom_smooth()' using method = 'loess' and formula 'y ~ x'

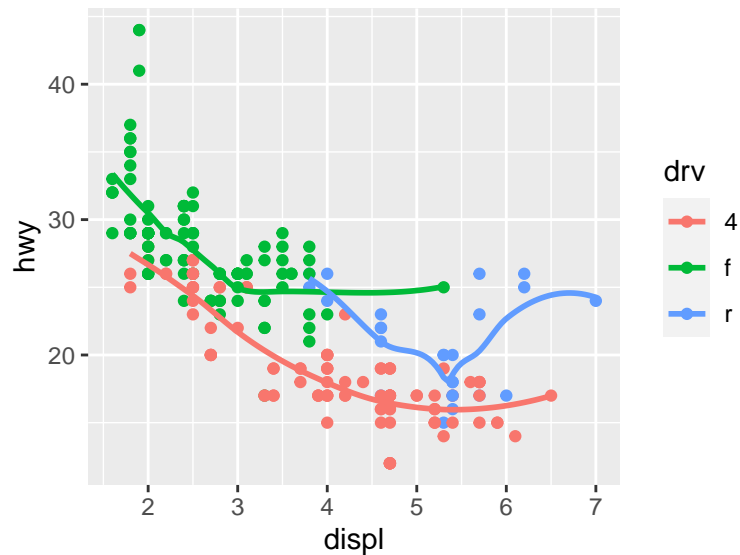


Figure 12: 4.4

PROBLEM 21.

SOLUTION.

The two graphs look the same. Since they display the same data by same method.

```
1 ggplot(data = mpg, mapping = aes(x = displ, y = hwy  
  )) +  
2   geom_point() +  
3   geom_smooth()  
4  
5 ggplot() +  
6   geom_point(data = mpg, mapping = aes(x = displ, y  
  = hwy)) +  
7   geom_smooth(data = mpg, mapping = aes(x = displ,  
    y = hwy))
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

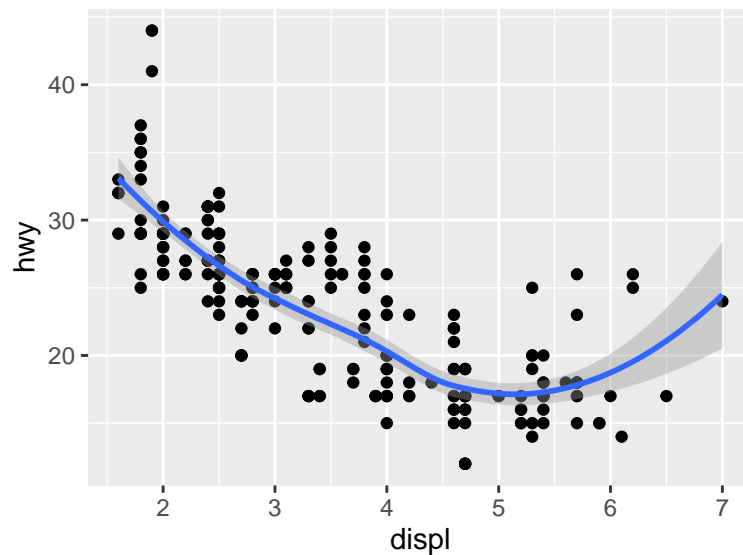


Figure 13: 4.5

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

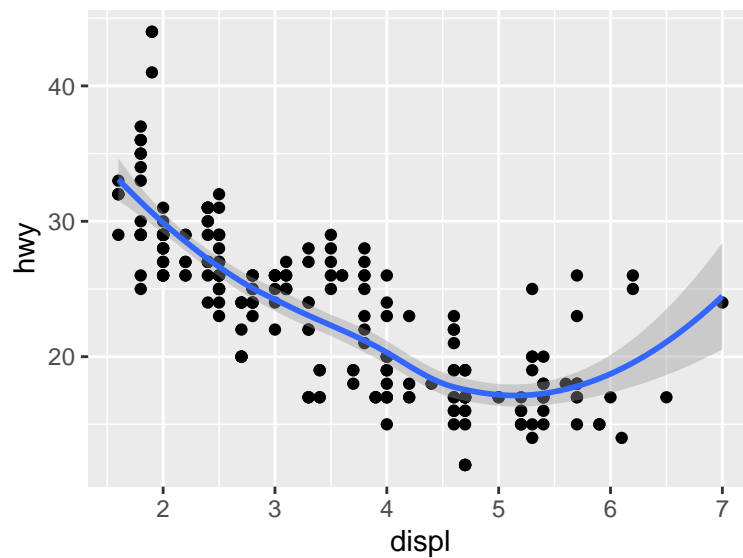


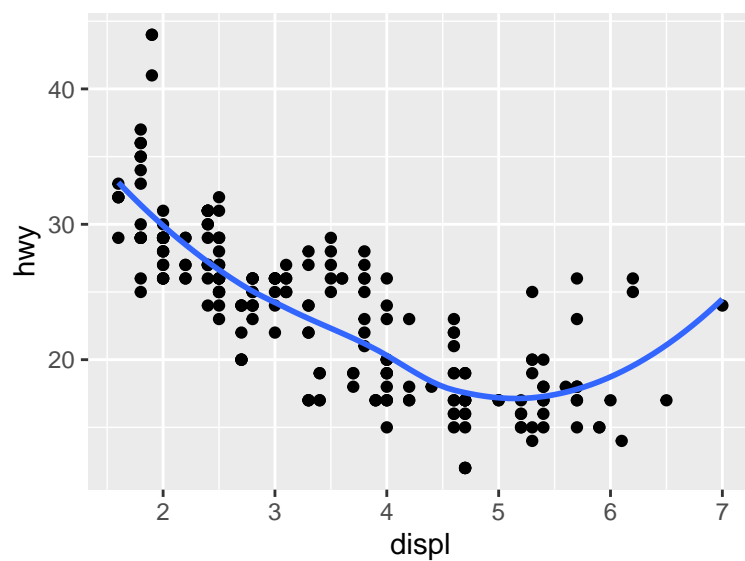
Figure 14: 4.5

PROBLEM 22.

SOLUTION.

```
1 ggplot(mpg, aes(x = displ, y = hwy)) +
2   geom_point()+
3   geom_smooth(se = FALSE)
4 ggplot(mpg, aes(x = displ, y = hwy)) +
5   geom_point()+
6   geom_smooth(mapping = aes(group = drv), se = FALSE
7 )
8 ggplot(mpg, aes(x = displ, y = hwy)) +
9   geom_point(mapping = aes(color = drv))+
10  geom_smooth(mapping = aes(color = drv), se = FALSE
11 )
12 ggplot(mpg, aes(x = displ, y = hwy)) +
13   geom_point(mapping = aes(color = drv))+
14   geom_smooth(se = FALSE)
15 ggplot(mpg, aes(x = displ, y = hwy)) +
16   geom_point(mapping = aes(color = drv))+
17   geom_smooth(mapping = aes(linetype = drv), se =
18     FALSE)
19 ggplot(mpg, aes(x = displ, y = hwy)) +
20   geom_point(size = 3, color = "white") +
21   geom_point(aes(colour = drv))
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

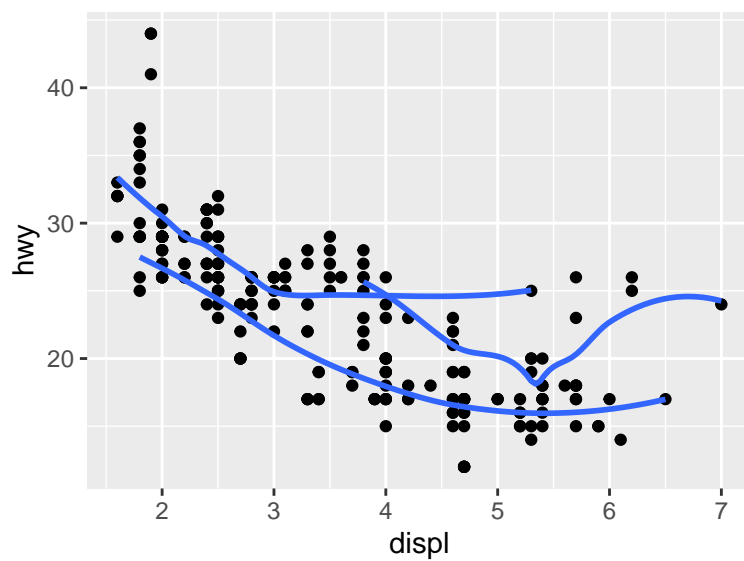
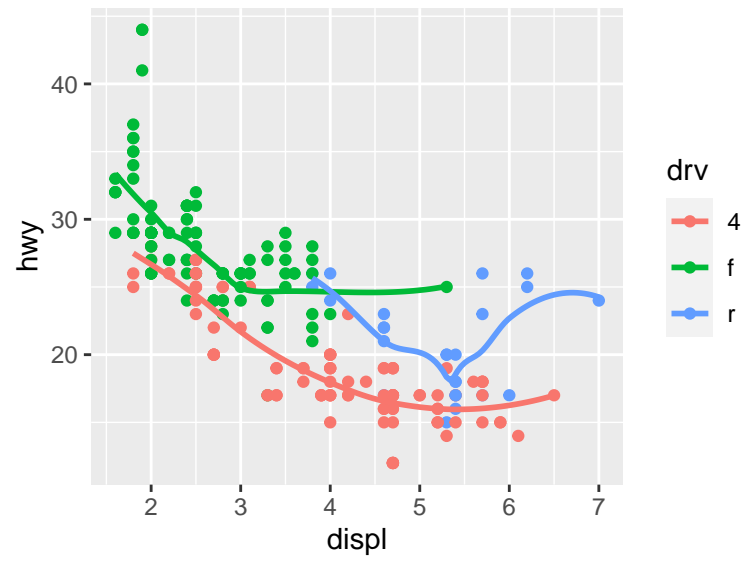


Figure 15: 4.6

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

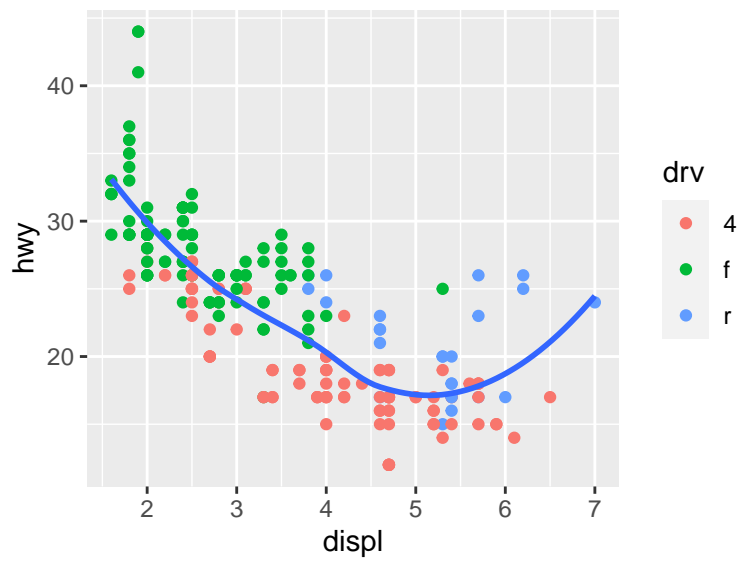


Figure 16: 4.6

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```

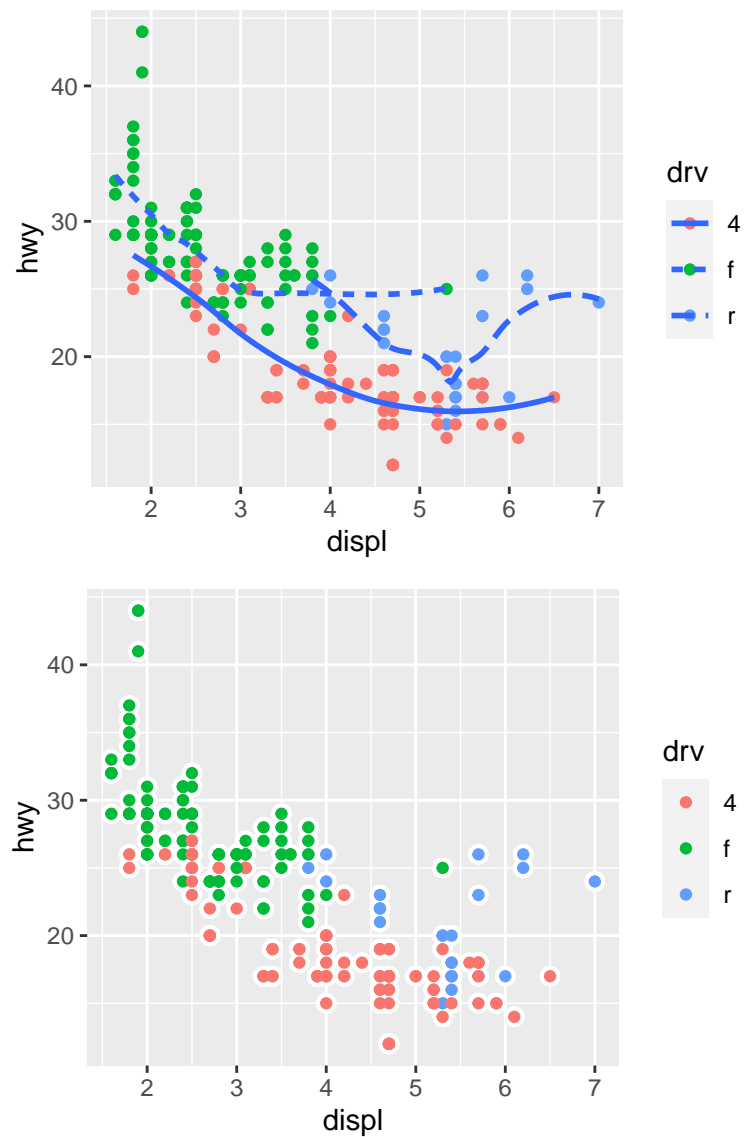


Figure 17: 4.6

5 Exercise

PROBLEM 23.

SOLUTION.

```
1 ggplot(data = diamonds) +
2   stat_summary(
3     mapping = aes(x = cut, y = depth),
```

```

4     fun.min = min,
5     fun.max = max,
6     fun = median
7   )
8   ggplot(data = diamonds) +
9     geom_pointrange(
10      mapping = aes(x = cut, y = depth),
11      stat = "summary",
12      fun.min = min,
13      fun.max = max,
14      fun = median
15    )

```

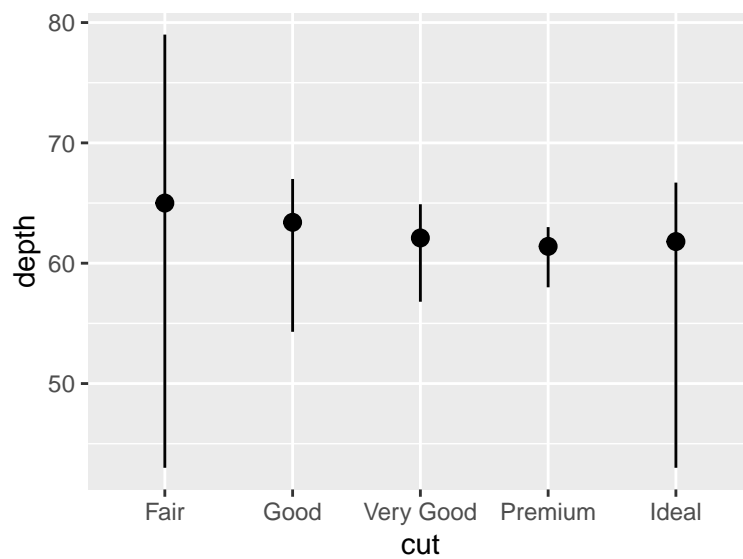


Figure 18: 5.1

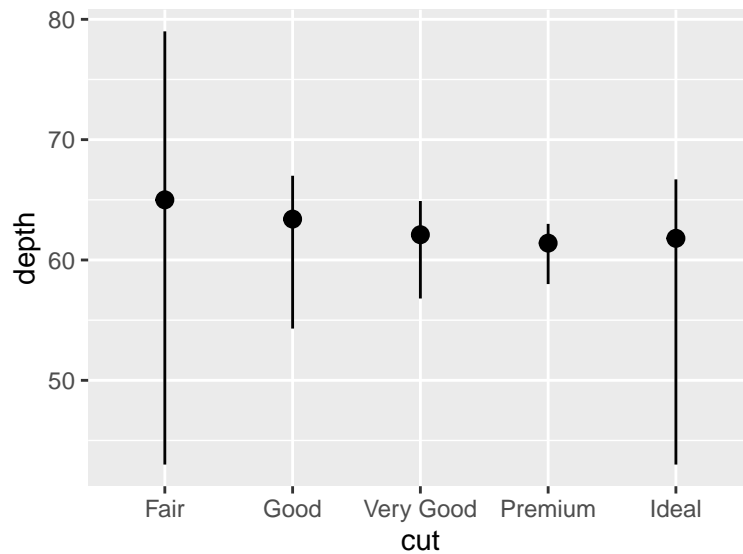


Figure 19: 5.1

PROBLEM 24.

SOLUTION.

Geom_col and geom_bar are used to create different types of bar Charts.

Official explanation are as follows:

Geom_bar() makes the height of the bar proportional to the number of cases in each group (or if the weight aesthetic is supplied, the sum of the weights). If you want the heights of the bars to represent values in the data, use geom_col() instead. geom_bar() uses stat_count() by default: it counts the number of cases at each x position. geom_col() uses stat_identity(): it leaves the data as is.

```

1 geom_bar (
2   mapping = NULL,
3   data = NULL,
4   stat = "count",
5   position = "stack",
6   ... ,
7   width = NULL,
8   na.rm = FALSE,
9   orientation = NA,
10  show.legend = NA,
11  inherit.aes = TRUE
12 )
13
14 geom_col (
15   mapping = NULL,
16   data = NULL,
17   position = "stack",
18   ... ,

```

```

19 width = NULL,
20 na.rm = FALSE,
21 show.legend = NA,
22 inherit.aes = TRUE
23 )

```

PROBLEM 25.

SOLUTION.

Some geom have stat as the default element while some stat also have geom as the default element.

geom	stat
geom_bar()	stat_count()
geom_bin2d()	stat_bin_2d()
geom_boxplot()	stat_boxplot()
geom_contour_filled()	stat_contour_filled()
geom_contour()	stat_contour()
geom_count()	stat_sum()
geom_density()	stat_density()
geom_dotplot()	stat_bindot()
geom_function()	stat_function()
geom_sf()	stat_sf()
geom_smooth()	stat_smooth()
geom_violin()	stat_ydensity()
geom_hex()	stat_bin_hex()
geom_qq()	stat_qq()
geom_qq_line()	stat_qq_line()
geom_quantile()	stat_quantile()

list pairs of geom and stat

PROBLEM 26.

SOLUTION.

"stat_smooth" is used to fit a continuous curve to discrete data points.

It calculates the following variables:

Predicted value(y) and its confidence interval

Example is as follow:

```

1 stat_smooth(
2   mapping = NULL,
3   data = NULL,
4   geom = "smooth",
5   position = "identity",
6   ... ,

```



```

7   method = NULL,
8   formula = NULL,
9   se = TRUE,
10  n = 80,
11  span = 0.75,
12  fullrange = FALSE,
13  level = 0.95,
14  method.args = list(),
15  na.rm = FALSE,
16  orientation = NA,
17  show.legend = NA,
18  inherit.aes = TRUE
19 )

```

```

1  ggplot(mtcars, aes(qsec, wt))+
2    stat_smooth()+
3    stat_smooth() + geom_point()

```

```

## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'##
'geom_smooth()' using method = 'loess' and formula 'y ~ x'

```

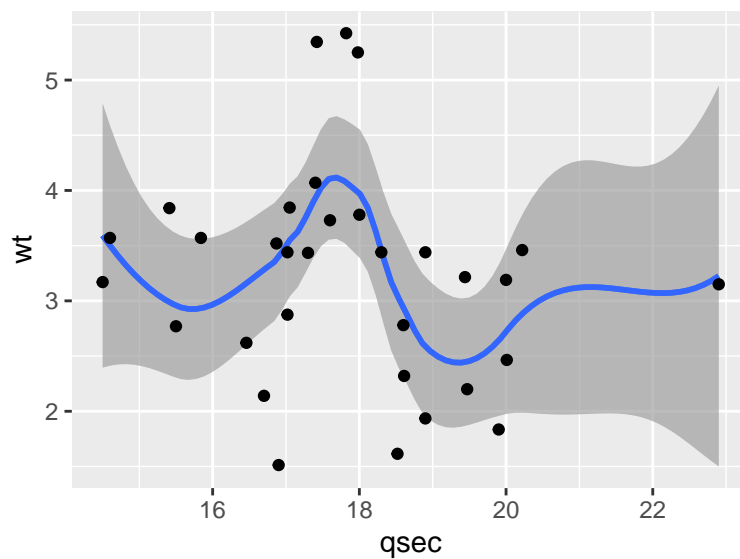


Figure 20: 5.4

PROBLEM 27.

SOLUTION.

The problem with these two graphs is that the proportions are calculated within every each type which results in 1 for every x.

```

1 ggplot(data = diamonds) +
2   geom_bar(mapping = aes(x = cut, y = ..prop..))
3
4 ggplot(data = diamonds) +
5   geom_bar(mapping = aes(x = cut, fill = color, y =
6     ..prop..))

```

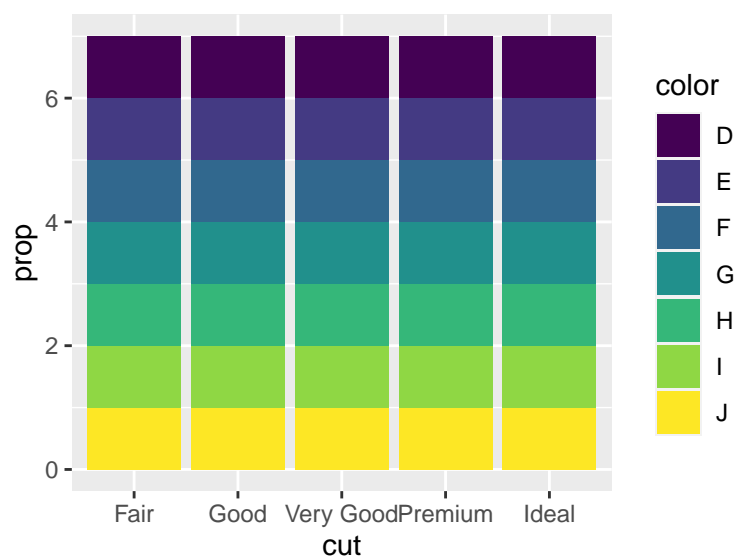
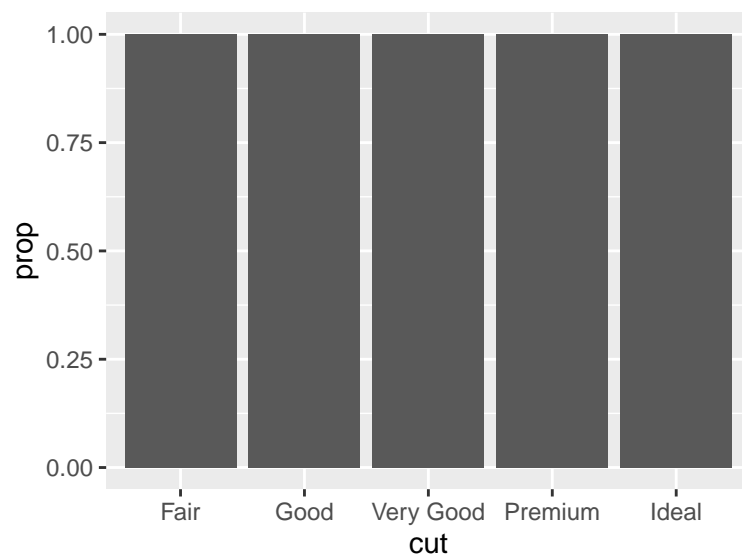


Figure 21: 5.5

6 Exercise

PROBLEM 28.

SOLUTION.

The problem with the plot is that many points overlap together. Thus, we have difficulty in seeing where the mass of data is. To solve the problem, we can add some random noise to the points.

```
1 %original codes
2 ggplot(data=mpg, mapping = aes(x=cty ,y=hwy)) +
3   geom_point()
4 %improved codes
5 ggplot(data=mpg, mapping = aes(x=cty ,y=hwy)) +
6   geom_jitter()
```

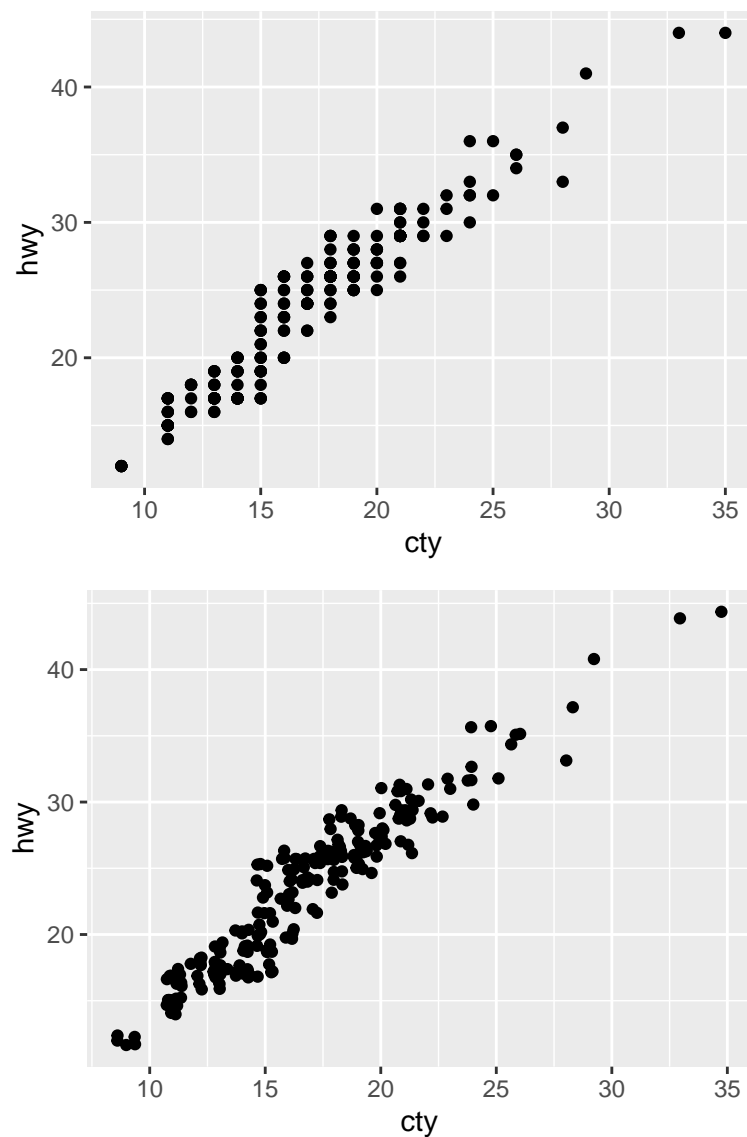


Figure 22: 6.1

PROBLEM 29.

SOLUTION.

The argument "width" and "height" are used to control the amount of jittering. Specifically, The jitter is added in both positive and negative directions, so the total spread is twice the value specified here.defaults to 0.4.

And for categorical data which is aligned on the integers, it not possible to see the distinction.So a width or height of 0.5 will spread the categorical data.

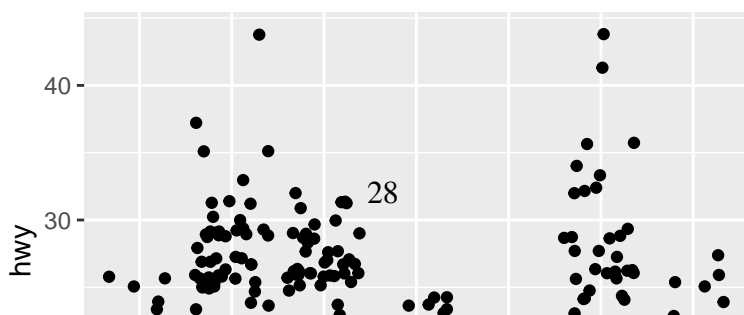
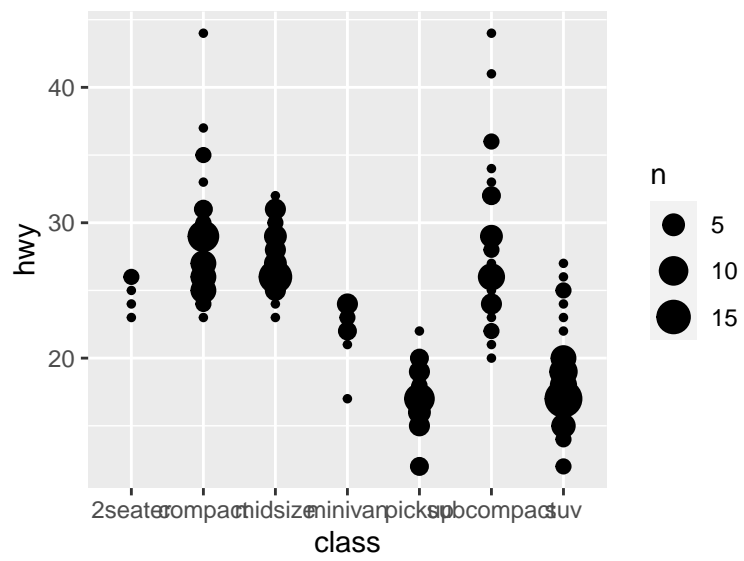
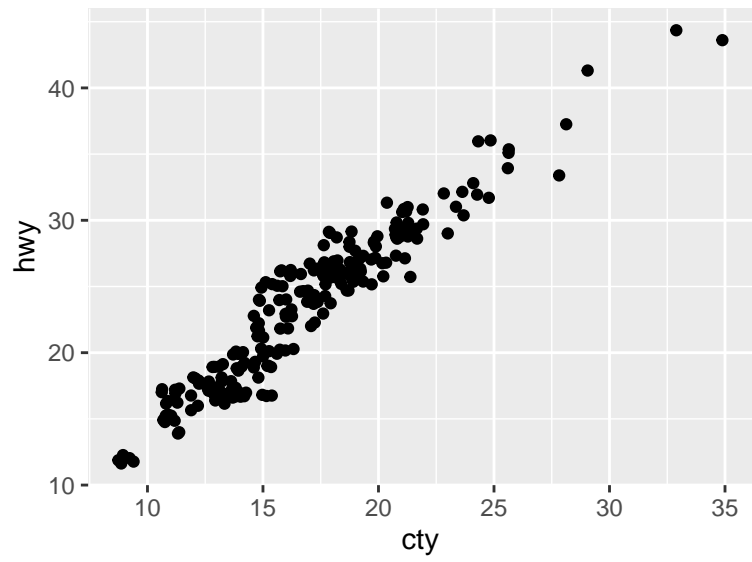
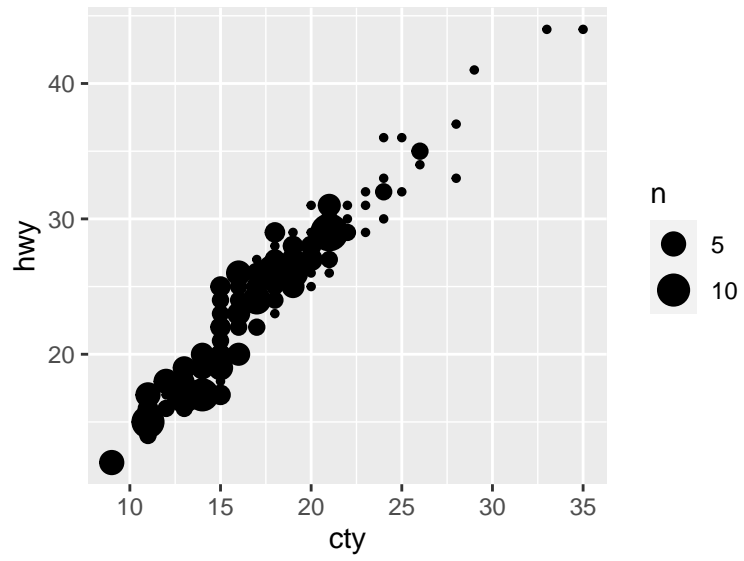
```
1 geom_jitter(  
2   mapping = NULL,  
3   data = NULL,  
4   stat = "identity",  
5   position = "jitter",  
6   ... ,  
7   width = NULL,  
8   height = NULL,  
9   na.rm = FALSE,  
10  show.legend = NA,  
11  inherit.aes = TRUE  
12 )
```

PROBLEM 30.

SOLUTION.

Both of geom_count and geom_jitter are used to solve overlapping. But, geom_count is used for categorical variables while geom_jitter is suitable for continuous variables.

```
1 ggplot(data=mpg, mapping = aes(x=cty ,y=hwy)) +  
2   geom_count()  
3 ggplot(data=mpg, mapping = aes(x=cty ,y=hwy)) +  
4   geom_jitter()  
5 ggplot(data=mpg, mapping = aes(x=class ,y=hwy)) +  
6   geom_count()  
7 ggplot(data=mpg, mapping = aes(x=class ,y=hwy)) +  
8   geom_jitter()
```



PROBLEM 31.

SOLUTION.

The boxplot compactly displays the distribution of a continuous variable. It visualises five summary statistics (the median, two hinges and two whiskers), and all "outlying" points individually.

```
1 ggplot(data=mpg, mapping = aes(x=class, y=hwy)) +  
2   geom_boxplot()
```

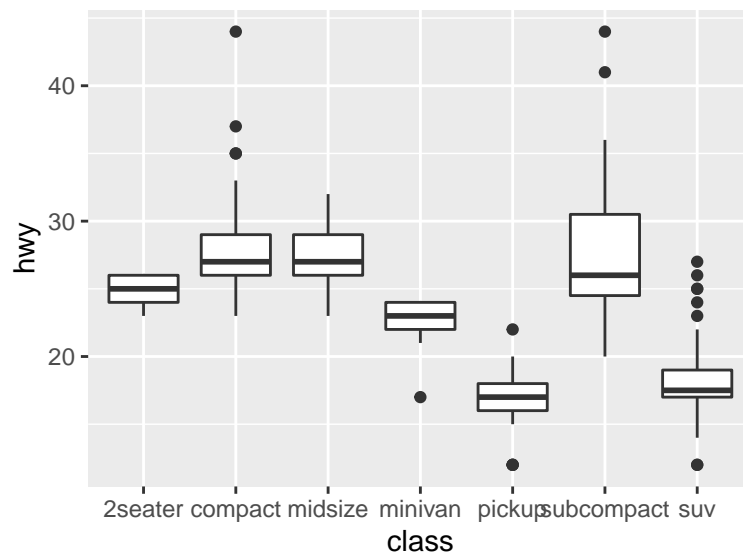


Figure 24: 6.4

7 Exercise

PROBLEM 32.

Turn a stacked bar chart to a pie chart.

SOLUTION.

```
1 bar = ggplot(data = diamonds) +  
2   geom_bar(mapping = aes(x = cut, fill = cut), show.  
3     legend = FALSE, width = 1)  
4 bar + coord_polar()
```

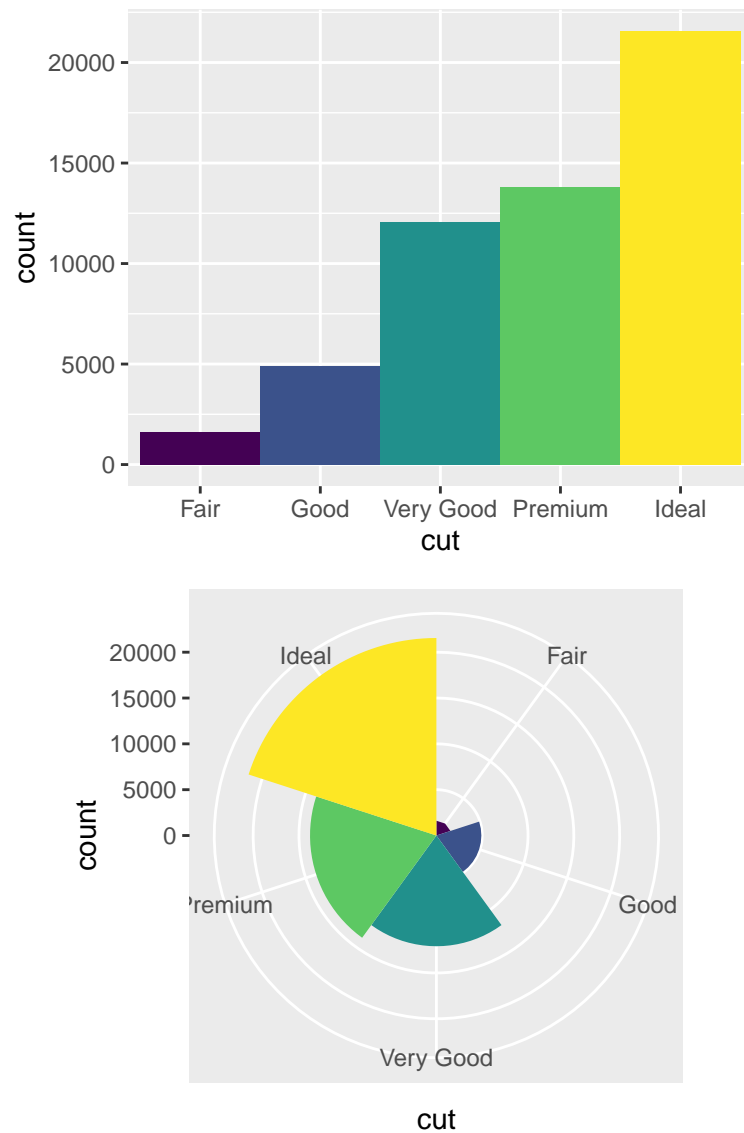


Figure 25: 7.1

PROBLEM 33.

SOLUTION.

"labs" modify axis, legend, and plot labels.

```

1 labs (
2   ...,
3   title = waiver() ,
4   subtitle = waiver() ,
5   caption = waiver() ,
6   tag = waiver() ,
7   %tag can be used for adding identification tags
   to differentiate between multiple plots.

```

```

8   alt = waiver(),
9   alt_insight = waiver()
10  )
11
12  xlab(label)
13
14  ylab(label)
15
16  ggtitle(label, subtitle = waiver())

```

PROBLEM 34.

SOLUTION.

They are projected differently. "Quick_mapping" ignores the curvature of the earth, so it is faster when calculating.

```

1  library (maps)
2  nz = map_data ("nz")
3  map = ggplot (nz, aes (long, lat, group=group)) +
4    geom_polygon (fill="white", color="black")
5  map + coord_map()
6  map + coord_quickmap()

```



```
#### 'maps' ## The following object is masked from
'package:purrr':#### map
```

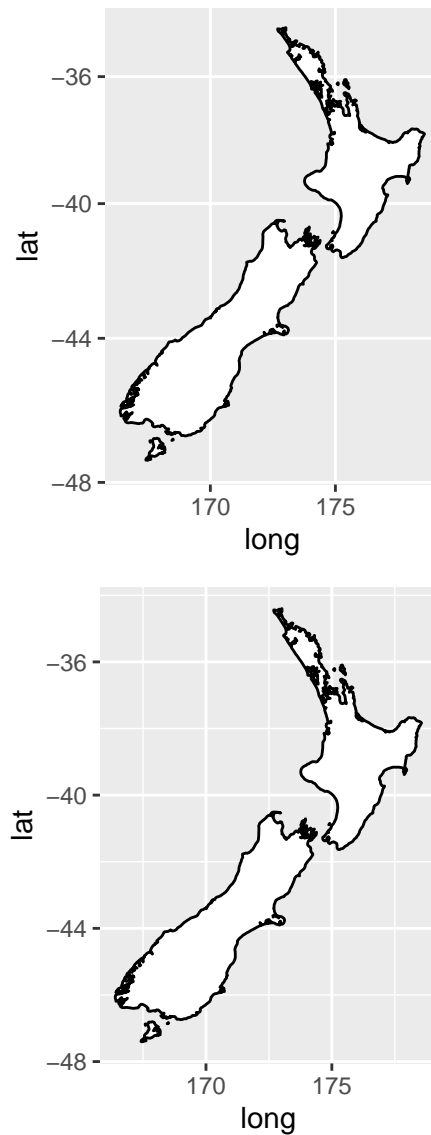


Figure 26: 7.3

PROBLEM 35.

SOLUTION.

"geom_abline" add line $y=x$ to the plot, helping us to explore their size and relationship.
 "coord_fixed" helps to remain the fixed ratio of the graph.

```
1 a <- ggplot(data = mpg, mapping = aes(x = cty, y =
  hwy)) +
2   geom_point() +
```

```

3 | geom_abline()
4 | a
5 | a + coord_fixed()

```

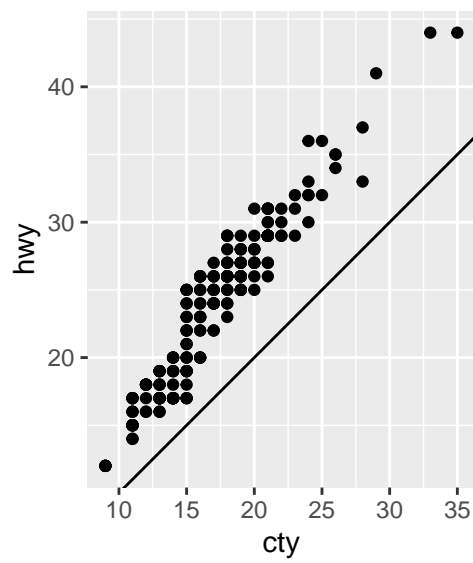
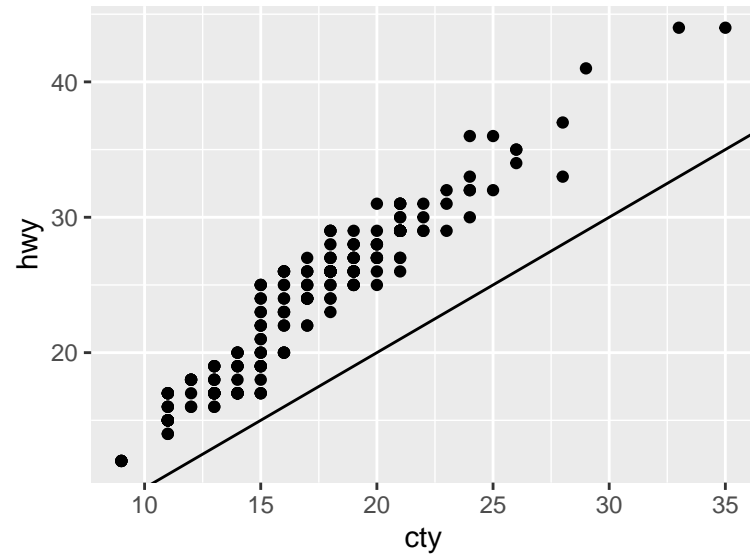


Figure 27: 7.4