

R and Statistics

Learning R and Statistics.

Alberto Valdez

October 13, 2022

Contents

1	MPG Dataset	1
1.1	Plotting bars.	2
1.2	Formatting output	3
2	MPG Summary	4
2.1	Summary Table	4
2.2	Mpg dataset	5
3	Boxplot	8
4	Heatmaps	10
4.1	Class and Year Summary	10
4.2	Model and Year Summary	12
5	Layered Plots	14
5.1	Summary of Class	14

1 MPG Dataset

Learning R with Emacs.¹ Trying to follow Google's R style guide.²

```
library(ggplot2)
library(tidyverse)
```

¹<https://orgmode.org/worg/org-contrib/babel/languages/ob-doc-R.html>

²<https://web.stanford.edu/class/cs109l/unrestricted/resources/google-style.html>

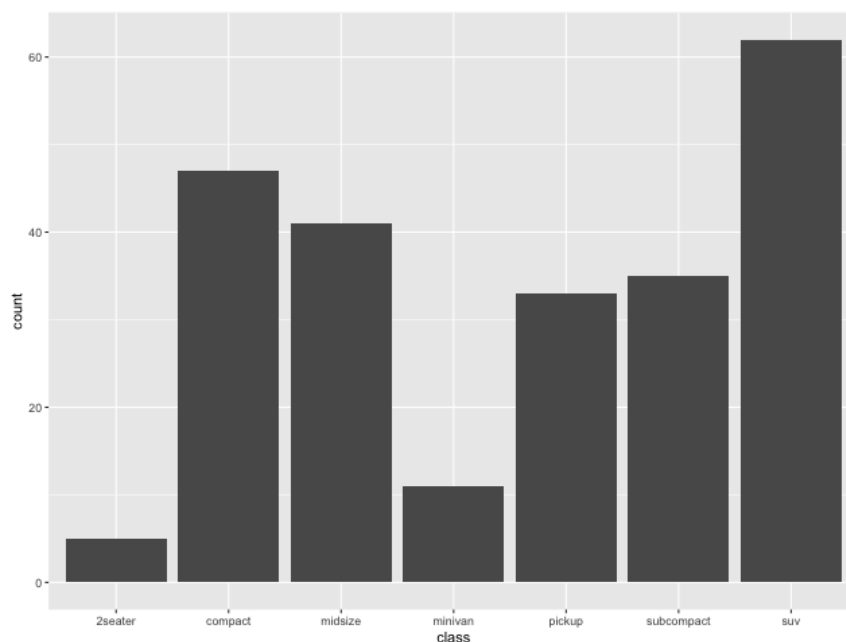
The mpg dataset contains fuel economy data from the EPA for vehicles manufactured between 1999 and 2008. The mpg dataset is built into R and is used throughout R documentation due to its availability, diversity of variables, and overall cleanliness of data. For our purposes, we'll use the mpg data to demonstrate how to implement each of our ggplot visualizations.

```
head(mpg)
```

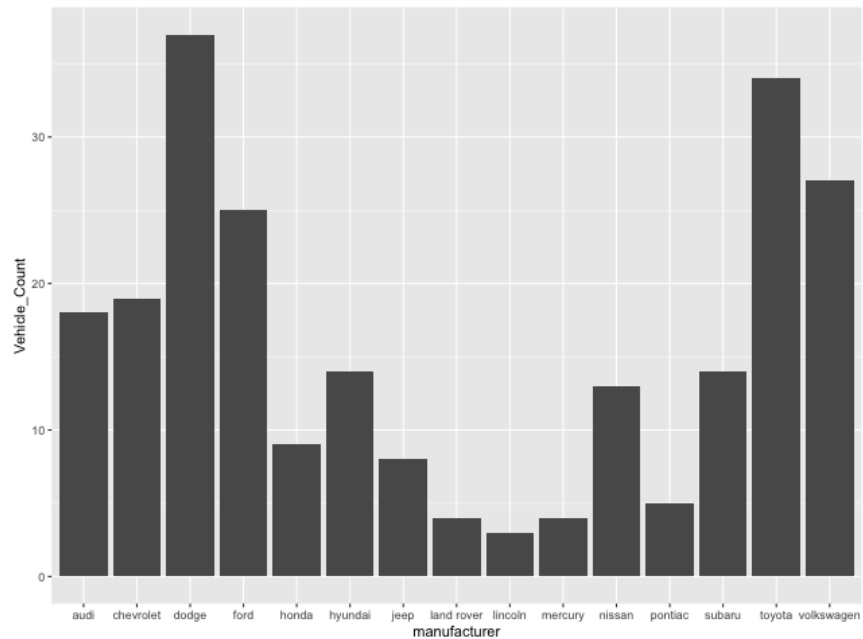
manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
audi	a4	1.8	1999	4	auto(l5)	f	18	29	p	compact
audi	a4	1.8	1999	4	manual(m5)	f	21	29	p	compact
audi	a4	2	2008	4	manual(m6)	f	20	31	p	compact
audi	a4	2	2008	4	auto(av)	f	21	30	p	compact
audi	a4	2.8	1999	6	auto(l5)	f	16	26	p	compact
audi	a4	2.8	1999	6	manual(m5)	f	18	26	p	compact

1.1 Plotting bars.

```
plt <- ggplot(mpg, aes(x=class))  
plt + geom_bar()
```



```
mpg_summary <- mpg %>%
  group_by(manufacturer) %>%
  summarize(Vehicle_Count=n(), .groups = 'keep') #create summary table
plt <- ggplot(
  mpg_summary,
  aes(x=manufacturer,y=Vehicle_Count)) #import dataset into ggplot2
plt + geom_col()
```



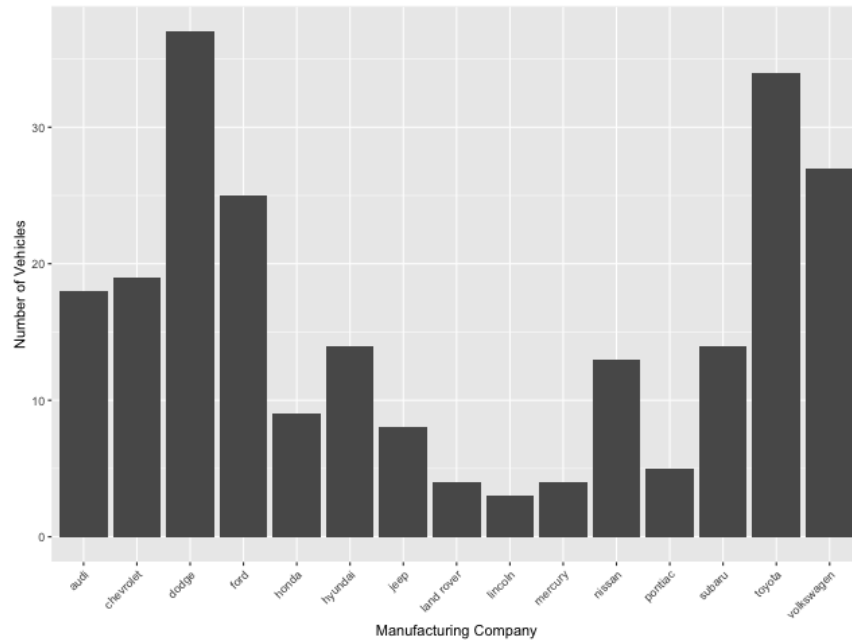
More info at the ggplot2 docs³.

1.2 Formatting output

Adding labels and themes.

```
plt + geom_col() +
  xlab("Manufacturing Company") +
  ylab("Number of Vehicles") +
  theme(axis.text.x=element_text(angle=45, hjust=1))
```

³<https://ggplot2.tidyverse.org/reference/index.html>



2 MPG Summary

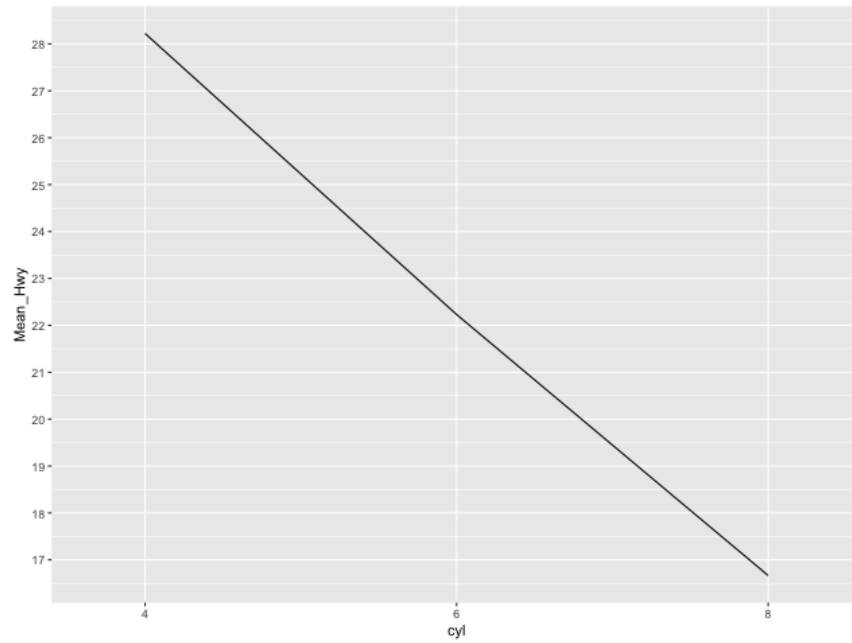
2.1 Summary Table

```
mpg_summary <-
  subset(mpg, manufacturer=="toyota") %>%
  group_by(cyl) %>%
  summarize(Mean_Hwy=mean(hwy), .groups="keep")
```

cyl	Mean _{Hwy}
4	28.222222222222
6	22.2307692307692
8	16.6666666666667

Import dataset into ggplot and plot the data and adjust the axis.

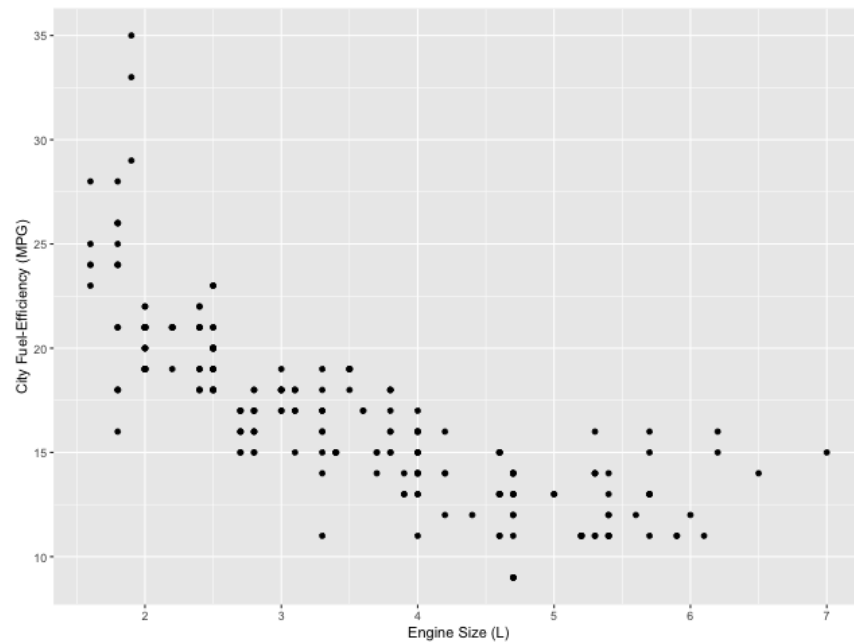
```
plt <- ggplot(mpg_summary, aes(x=cyl, y=Mean_Hwy))
plt + geom_line() +
  scale_x_discrete(limits=c(4, 6, 8)) +
  scale_y_continuous(breaks = c(15:30))
```



2.2 Mpg dataset

Import into ggplot and plot data with formatting.

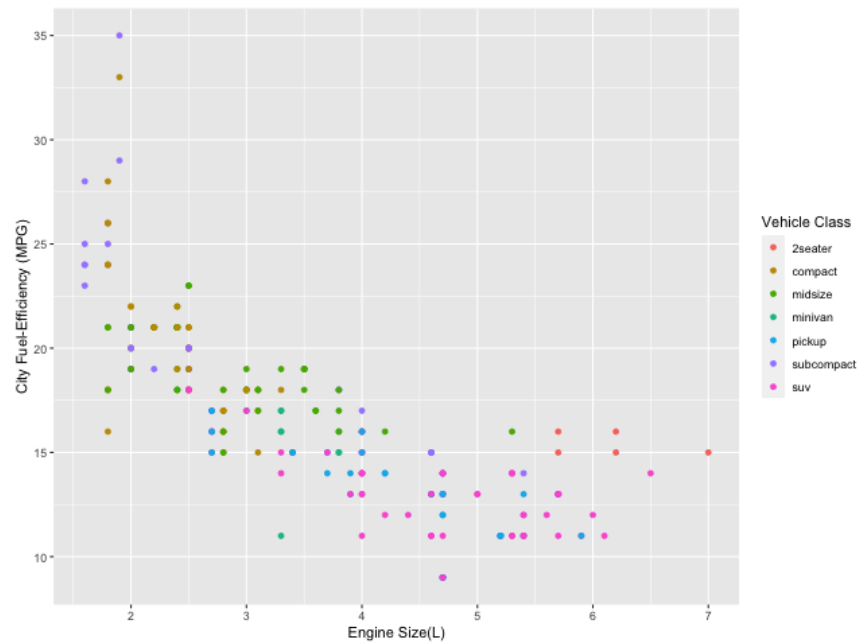
```
plt <- ggplot(mpg, aes(x=displ, y=cty))  
plt + geom_point() +  
  xlab("Engine Size (L)") +  
  ylab("City Fuel-Efficiency (MPG)")
```



Aesthetic changes.

- alpha changes the transparency of each data point
- color changes the color of each data point
- shape changes the shape of each data point
- size changes the size of each data point

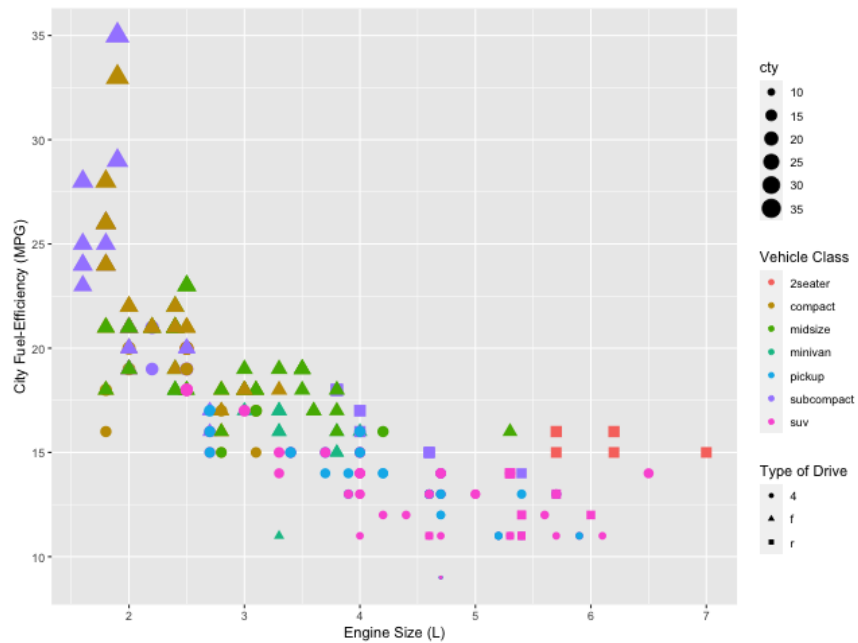
```
plt <- ggplot(mpg, aes(x=displ, y=cty, color=class))
plt + geom_point() +
  labs(
    x="Engine Size(L)",
    y="City Fuel-Efficiency (MPG)",
    color="Vehicle Class"
  )
```



Different shapes.⁴

```
plt <- ggplot(
  mpg,
  aes(x=displ, y=cty, color=class, shape=drv, size=cty)
)
plt + geom_point() +
  labs(
    x="Engine Size (L)",
    y="City Fuel-Efficiency (MPG)",
    color="Vehicle Class",
    shape="Type of Drive"
  )
```

⁴https://ggplot2.tidyverse.org/reference/geom_point.html#aesthetics



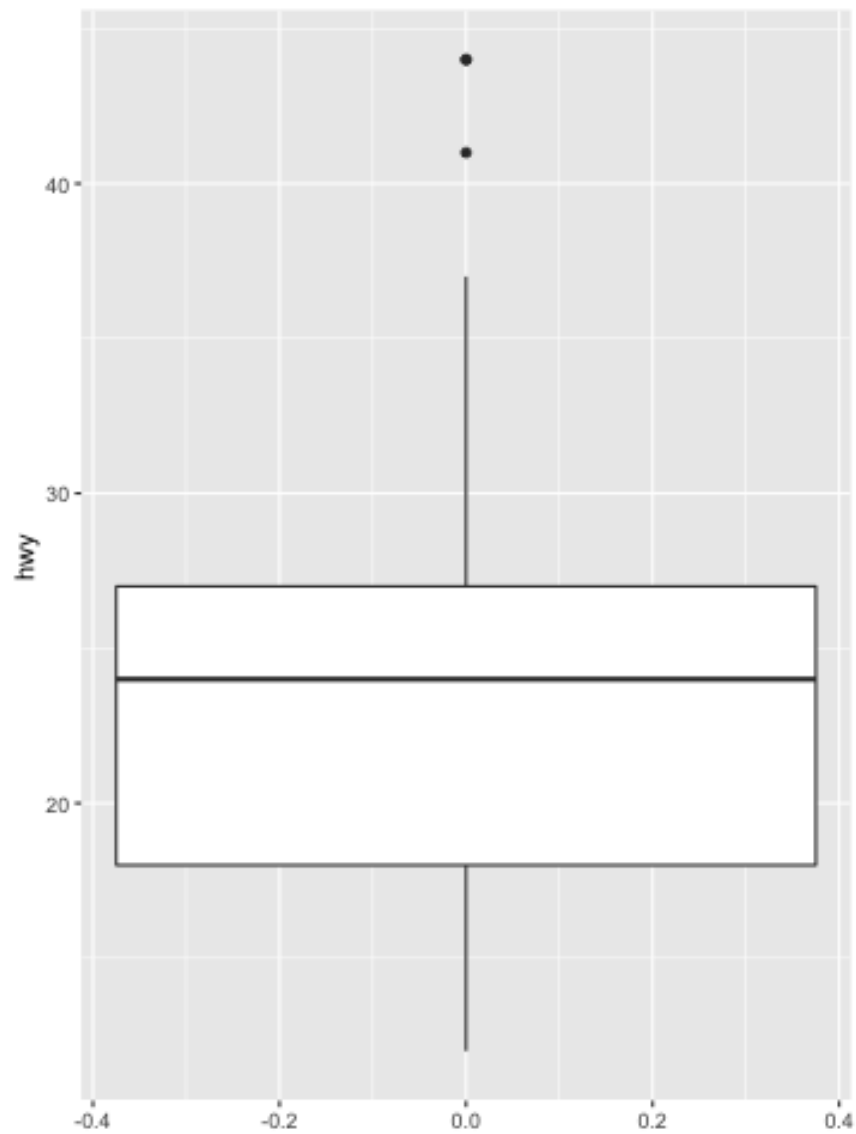
Although there is no technical limit to the number of variables we can add to a ggplot figure, there are diminishing returns. A good rule of thumb is to limit the number of variables displayed in a single figure to a maximum of 3 or 4.

3 Boxplot

Unlike the previous ggplot objects, `geom_boxplot()` expects a numeric vector assigned to the y-value⁵.

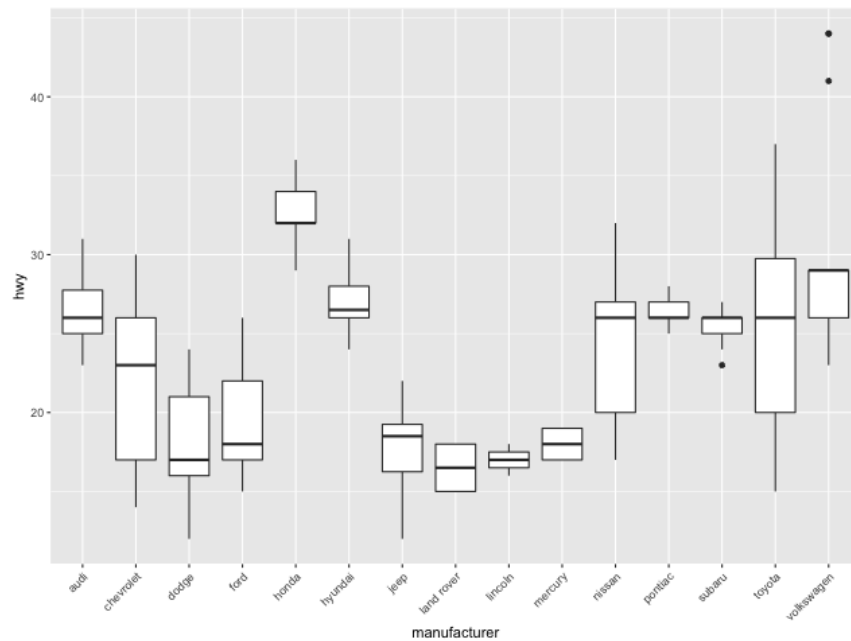
```
plt <- ggplot(mpg, aes(y=hwy))
plt + geom_boxplot()
```

⁵https://ggplot2.tidyverse.org/reference/geom_boxplot.html#aesthetics



Creating multiple boxes.

```
plt <- ggplot(mpg, aes(x=manufacturer, y=hwy))  
plt +  
  geom_boxplot() +  
  theme(axis.text.x=element_text(angle=45, hjust=1))
```



4 Heatmaps

Heatmap plots help visualize the relationship between one continuous numerical variable and two other variables (categorical or numerical).

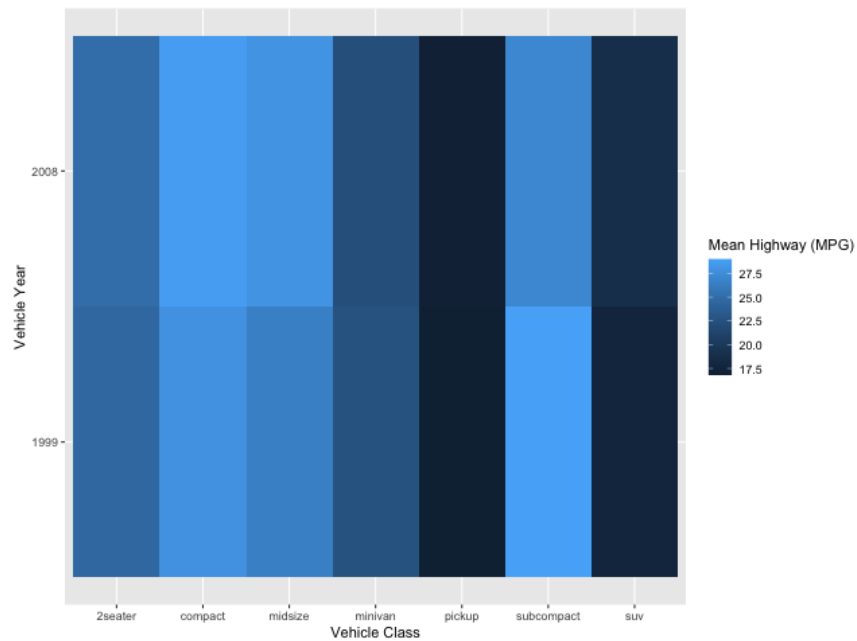
4.1 Class and Year Summary

```
mpg_summary <- mpg %>%
  group_by(class, year) %>%
  summarize(Mean_Hwy=mean(hwy), .groups='keep')
```

class	year	Mean _{Hwy}
2seater	1999	24.5
2seater	2008	25
compact	1999	27.92
compact	2008	28.7272727272727
midsize	1999	26.5
midsize	2008	28.047619047619
minivan	1999	22.5
minivan	2008	22.2
pickup	1999	16.8125
pickup	2008	16.9411764705882
subcompact	1999	29
subcompact	2008	27.125
suv	1999	17.551724137931
suv	2008	18.6363636363636

Plotting heatmap.

```
plt <- ggplot(
  mpg_summary,
  aes(x=class, y=factor(year), fill=Mean_Hwy))
plt + geom_tile() +
  labs(
    x="Vehicle Class",
    y="Vehicle Year",
    fill="Mean Highway (MPG)")
```



4.2 Model and Year Summary

```
mpg_summary <- mpg %>%
  group_by(model, year) %>%
  summarize(Mean_Hwy=mean(hwy), .groups='keep')
mpg_summary %>% head
```

model	year	Mean _{Hwy}
4runner 4wd	1999	19
4runner 4wd	2008	18.5
a4	1999	27.5
a4	2008	29.3333333333333
a4 quattro	1999	25.25
a4 quattro	2008	26.25

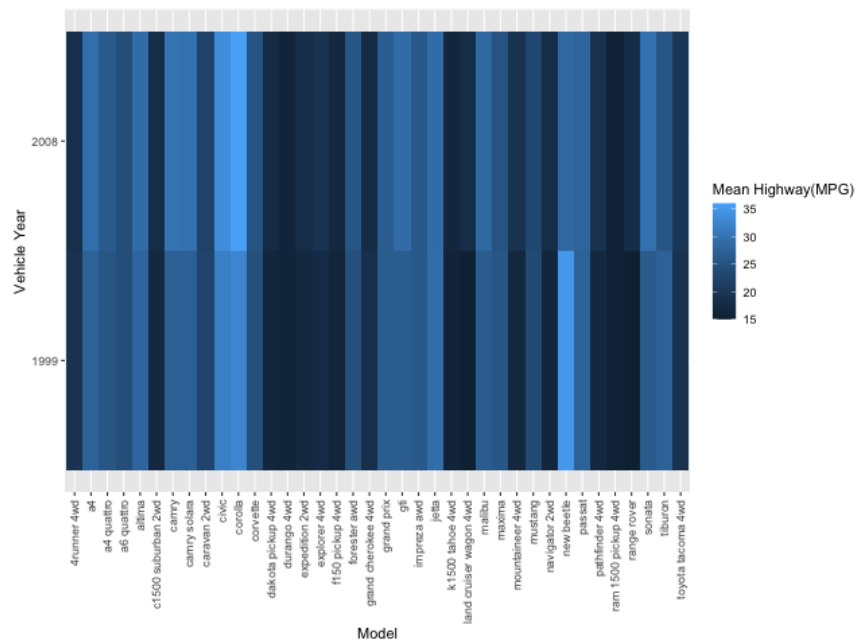
Adding labels to heatmap.

```
plt <- ggplot(
  mpg_summary,
  aes(
    x=model,
    y=factor(year),
```

```

    fill=Mean_Hwy)
  )
plt + geom_tile() +
  labs(
    x="Model",
    y="Vehicle Year",
    fill="Mean Highway(MPG)"
  ) +
  theme(
    axis.text.x = element_text(
      angle=90,
      hjust=1,
      vjust=0.5
    )
  )
)

```



We can always refer to the ggplot cheatsheet⁶.

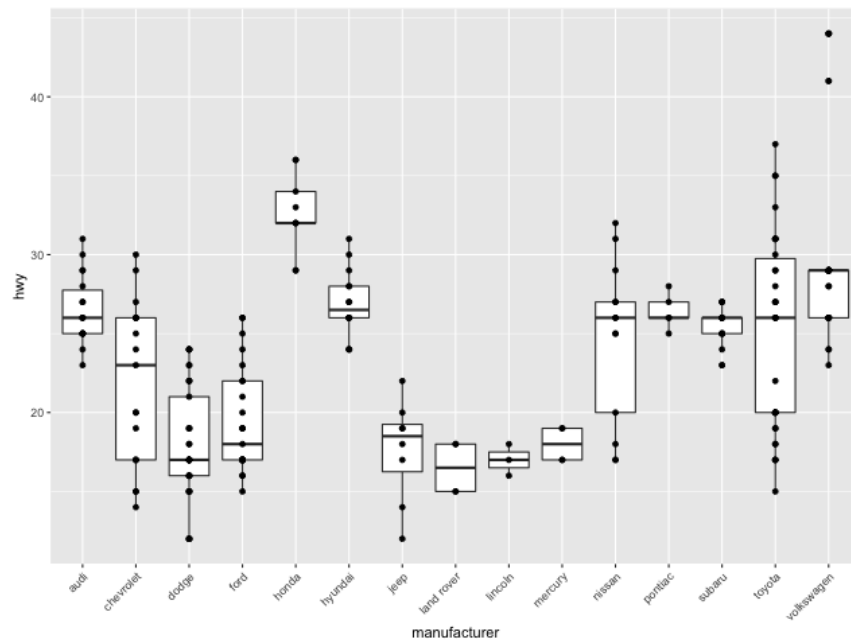
⁶<https://github.com/rstudio/cheatsheets/blob/main/data-visualization.pdf>

5 Layered Plots

There are two types of plot layers:

1. Layering additional plots that use the same variables and input data as the original plot
2. Layering of additional plots that use different but complementary data to the original plot

```
plt <- ggplot(mpg, aes(x=manufacturer,y=hwy))  
plt + geom_boxplot() +  
  theme(axis.text.x=element_text(angle=45,hjust=1)) +  
  geom_point()
```



By layering our data points on top of our boxplot, we can see the general distribution of values within each box as well as the number of data points.

5.1 Summary of Class

```
mpg_summary <- mpg %>%  
  group_by(class) %>%  
  summarize(Mean_Engine=mean(displ), .groups='keep')
```

class	Mean _{Engine}
2seater	6.16
compact	2.32553191489362
midsize	2.9219512195122
minivan	3.39090909090909
pickup	4.41818181818182
subcompact	2.66
suv	4.45645161290323

Plotting scatter plot.

```
plt <-
  ggplot(mpg_summary, aes(x=class,y=Mean_Engine))
plt +
  geom_point(size=4) +
  labs(x="Vehicle Class",y="Mean Engine Size")
```

