# Before we start...

Access the RStudio Cloud Project (Link in https://datafest.psu.edu/go/)

- Workshop - Introduction to Data Visualization in R with ggplot2

## OR

1. Download R
2. Download RStudio
3. Install packages `tidyverse` or `ggplot2` (check with library(ggplot2))

# Introduction to Data Visualization in R with ggplot2

Penn State, University Libraries, Research Informatics and Publishing

# Workshop Housekeeping

## Questions?
Use the Q&A, chat, or raise hand feature.

## Feedback Survey
After the workshop, please fill out the Qualtrics survey (link in chat).

# Credits

Images and content sourced/based on:

*ggplot2: Elegant Graphics for Data Analysis* by Hadley Wickham

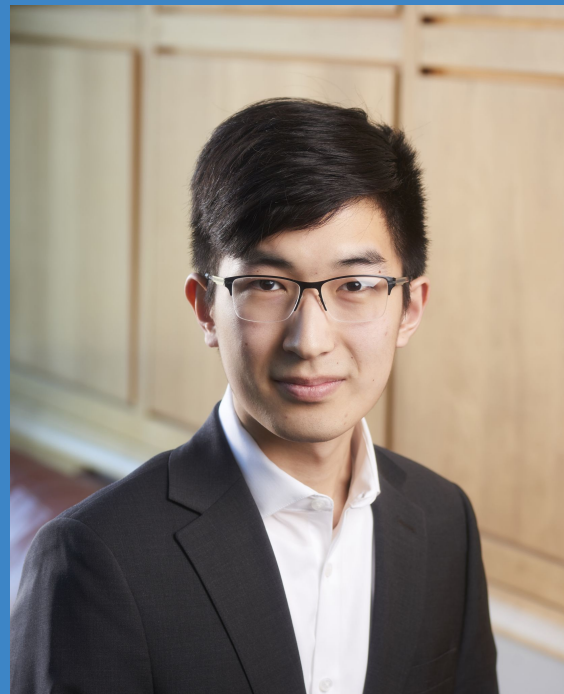*Introduction to data visualization with ggplot2*, Data Camp - Rick Scavetta

# About Me

- Master's in Applied Statistics
- Bachelor's in Computational Statistics

R Experience:

- Self-taught for research (2017-Present)
- Statistics Courses (STAT 184/380)
- Other Projects

Research Consultant, University Libraries

David Chen
dzc89@psu.edu
Research Consultant
github.com/TheDavidChen

# Introduce Yourself!

Name, Major, Year,
How do you pronounce the word: Data

## General Comments

If you see: > print("Hello World")

Run print("Hello World") in R Script. Do not include the `>`.

Use `#` to write comments - code after # is not run.

> # This is not run

# What is the purpose of Data Visualization?

# Agenda

- Base R vs ggplot2
- The logic of ggplot2
- Basic plots
- Data viz cheat sheet
- Faceting
- Themes
- Exporting plots
- The little things

# Base R vs ggplot2

Base R you plot based on
specific functions

- hist()
- boxplot()
- barplot()
- pie()
- plot()

# Base R vs ggplot2

ggplot2 is based on the Grammar of Graphics

Thinks about visualizations as layers/components

Greater flexibility, more personalization, more intricate plots, while simple plots are still simple

# ggplot2 layers

mpg dataset from ggplot2 - "Fuel economy data from 1999 to 2008 for 38 popular models of cars"

Variables

- displ - engine displacement (size)
- year - year of manufacture
- cty - city miles per gallon
- hwy - highway miles per gallon
- class - type of car (suv, pickup, etc.)

# ggplot2 layers

ggplot(mpg)

# ggplot2 layers

ggplot(mpg, aes(x = cty, y = hwy))

# ggplot2 layers

ggplot(mpg, aes(x = cty, y = hwy)) +

 geom_point()

# ggplot2 layers

ggplot(mpg, aes(x = cty, y = hwy)) +

 geom_point() +

 facet_wrap(~year)

# ggplot2 layers

ggplot(mpg, aes(x = cty, y = hwy)) +

  geom_point() +

  facet_wrap(~year) +

  geom_smooth(method = "lm")

# ggplot2 layers

ggplot(mpg, aes(x = cty, y = hwy)) +

geom_point() +

facet_wrap(~year) +

geom_smooth(method = "lm") +

theme_bw()

# ggplot2 layers

**Data**            **What dataset**

**Aesthetics**      **Select x and y-axis, colors, shapes, etc.**

**Geometries**      **Which geometric object (points, bars, lines)**

# ggplot2 layers

**Data**             **What dataset**

**Aesthetics**       **Select x and y-axis, colors, shapes, etc.**

**Geometries**       **Which geometric object (points, bars, lines)**

Facets               Plot subsets of the data separately

# ggplot2 layers

| | |
|---|---|
| **Data** | **What dataset** |
| **Aesthetics** | **Select x and y-axis, colors, shapes, etc.** |
| **Geometries** | **Which geometric object (points, bars, lines)** |
| Facets | Plot subsets of the data separately |
| Statistics | Statistical transformations |

# ggplot2 layers

| | |
|---|---|
| **Data** | **What dataset** |
| **Aesthetics** | **Select x and y-axis, colors, shapes, etc.** |
| **Geometries** | **Which geometric object (points, bars, lines)** |
| Facets | Plot subsets of the data separately |
| Statistics | Statistical transformations |
| Coordinates | Coordinate system (Cartesian, polar, map) |

# ggplot2 layers

| | |
|---|---|
| **Data** | **What dataset** |
| **Aesthetics** | **Select x and y-axis, colors, shapes, etc.** |
| **Geometries** | **Which geometric object (points, bars, lines)** |
| Facets | Plot subsets of the data separately |
| Statistics | Statistical transformations |
| Coordinates | Coordinate system (Cartesian, polar, map) |
| Themes | Font size, background color, etc. |

Breaking down ggplot2

# Before we start plotting

Make sure to load tidyverse (or just ggplot2)

`> library(tidyverse)`

We will load datasets like mpg by doing the following:

`> data(mpg)`

To learn about the dataset, type:

`> ?mpg`

# Before we start plotting

Plots will appear under the `Plots` tab in the bottom right section of RStudio

Click the arrow buttons to go between plots

ggplot(dataset, aes(<.........>))

Aesthetics

x and y-axis, color, etc.

# Aesthetics

| Aesthetic | Description |
|-----------|-------------|
| x | X axis |
| y | Y axis |
| fill | Fill color |
| color | Color of points or outlines of other geoms |
| size | Size of points, thickness of lines |
| alpha | Transparency |

# Geometries

ggplot(dataset, aes(xvar, yvar)) +

geom_*()

How to represent the data points

geom_point() = scatterplot

geom_histogram() = histogram

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point()

# ggplot(mpg, aes(x = displ, y = hwy)) +

# geom_point()



ggplot2: Elegant Graphics for Data Analysis,
by Hadley Wickham
https://ggplot2-book.org/getting-started.html

# ggplot(mpg, aes(x = displ, y = hwy)) + geom_point()

# ggplot(mpg, aes(x = displ, y = hwy)) +

## geom_point()



ggplot2: Elegant Graphics for Data Analysis, by Hadley Wickham
https://ggplot2-book.org/getting-started.html

# ggplot(mpg, aes(x = displ, y = hwy, color = class)) + geom_point()

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point(color = "blue")

# mpg Example

1. Load the mpg dataset (from the ggplot2 package)

> data(mpg)          > ?mpg # To learn about the variables

2. Experiment with different aesthetics (alpha, shape, size, color)

```
> ggplot(mpg, aes(x = displ, y = hwy)) +
```
```
> ggplot(mpg, aes(x = displ, y = hwy)) +
```
```
>   geom_point(alpha = 0.2)
```
```
>   geom_point(size = 5, color = "red")
```
```
> ggplot(mpg, aes(x = displ, y = hwy)) +
```
```
> ggplot(mpg, aes(x = displ, y = hwy)) +
```
```
>   geom_point(shape = 15)
```
```
>   geom_point()
```

ggplot(mpg, aes(x = drv)) +

geom_bar()

Some geometries only need one variable.

Bar Plots - One (Discrete) Variable

# ggplot(mpg, aes(x = drv, fill = year)) + geom_bar()



Arguments may require discrete (categorical) vs continuous variables.

Fill must be categorical, but the data has it as continuous.

# ggplot(mpg, aes(x = drv, fill = as.factor(year))) + geom_bar()

Setting the year (continuous) as a factor makes it categorical. We can now use it to fill the colors!

Note: For other plots, the color can be a continuous variable.

ggplot(mpg, aes(x = drv, fill = as.factor(year))) +

geom_bar(position = "dodge")

Different `geom_*` have
different arguments

ggplot(mpg, aes(y = drv, fill = as.factor(year))) +

geom_bar(position = "dodge")

# Data Viz Cheat Sheet

ggplot does not suggest what plot to make

How do you decide what type of plot to make?

- Number of variables
- Continuous vs Discrete
- Goal of the visualization

# Data Viz Cheat Sheet

# Economics Example

1. Load the `economics` dataset:

```
> data(economics)        > ?economics # To learn about the variables
```

2. Create a line plot to visualize the trend in personal savings rate (`psavert`) over time (`date`)

```
> ggplot(economics, aes(x = ____, y = _____) +

>   geom_____()
```

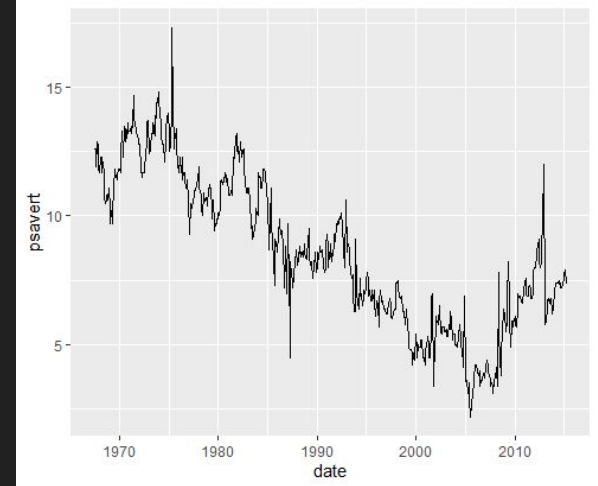3. Color the above plot based on the population size `pop`.

4. Create as many different types of plots for (2) as possible.

# Economics Example - Solutions

2. Create a line plot to visualize the trend in personal savings rate (`psavert`) over time (`date`)

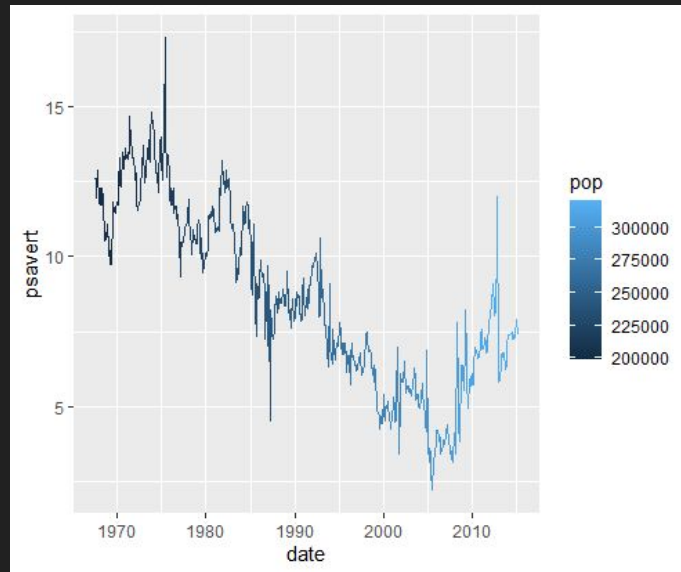> ggplot(economics, aes(x = date, y = psavert) +

>   geom_line()

# Economics Example - Solutions

3. Color the above plot based on the population size `pop`.

> ggplot(economics, aes(x = date, y = psavert, color = pop)) +

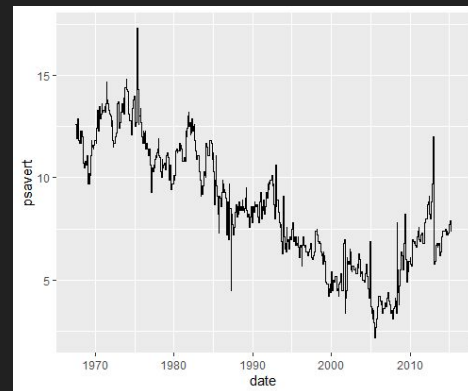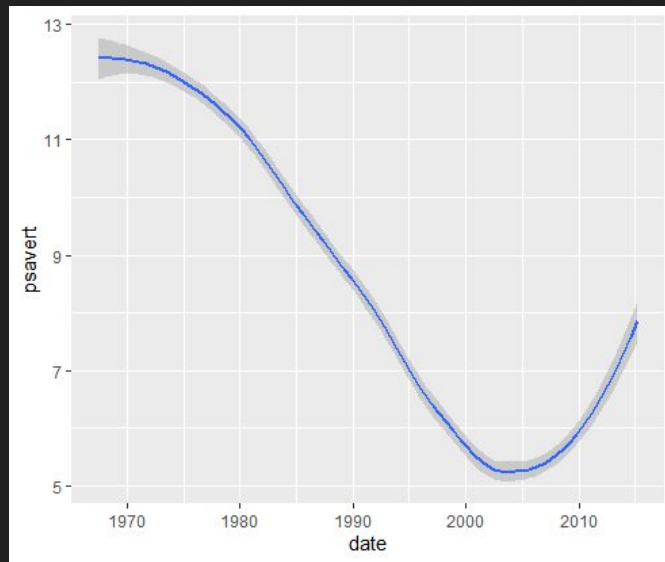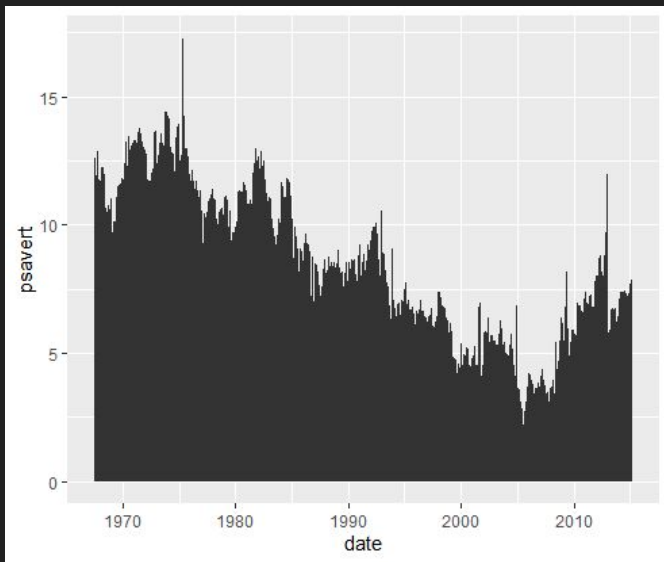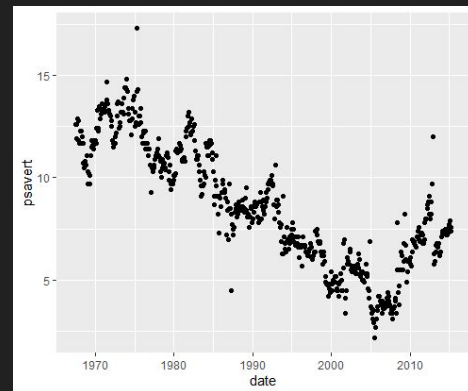>   geom_line()

# Economics Example - Solutions

4. Create as many different types of plots for (2) as possible.

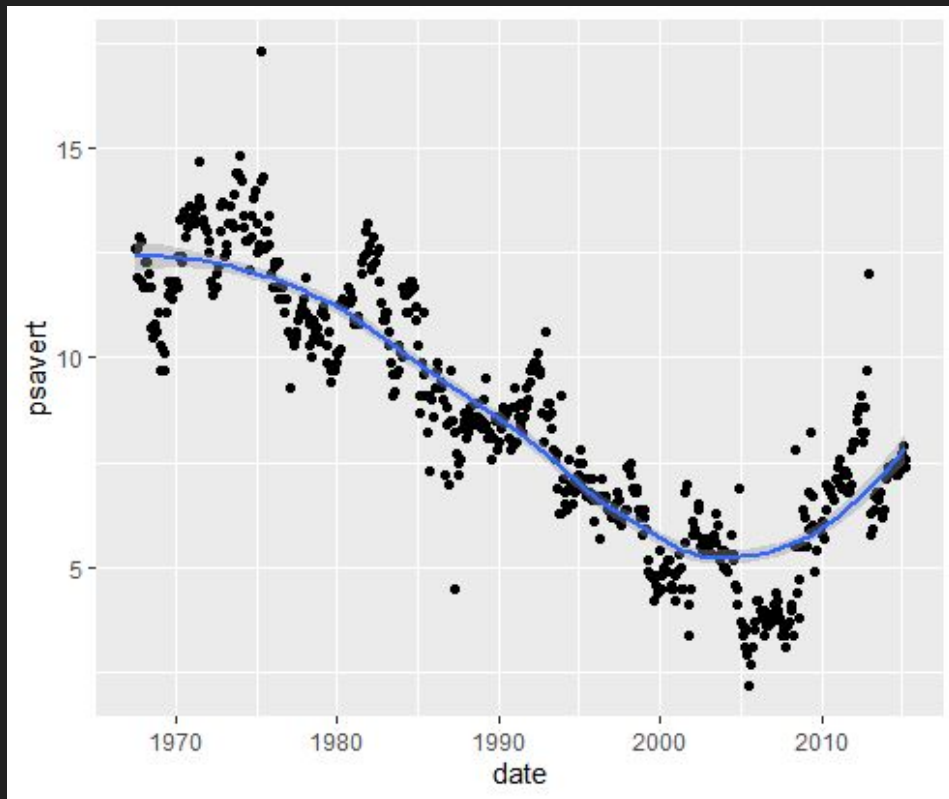Essentially all the geometries for 2 continuous variables work!

# Economics Example - Solutions

Note: You can overlap geoms!
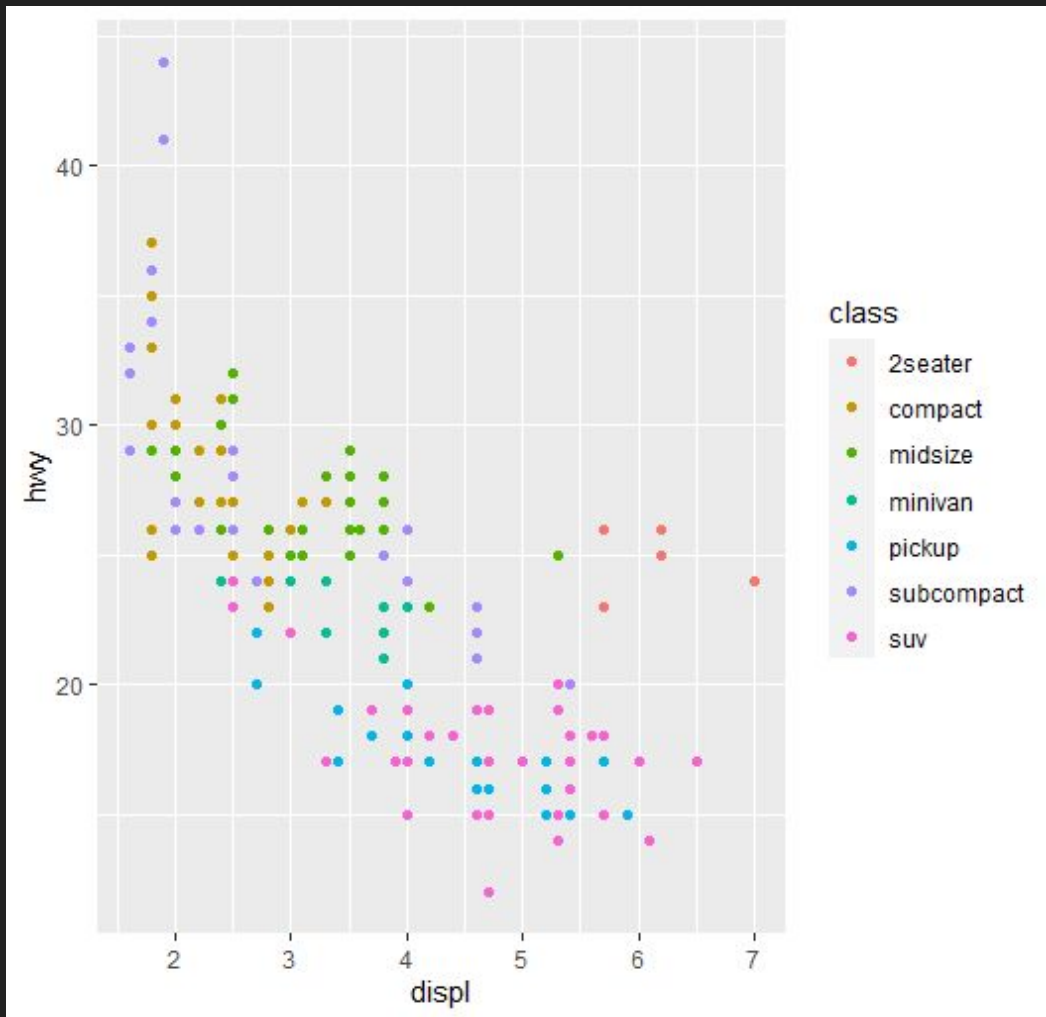
ggplot(economics, aes(x = date, y = psavert)) +
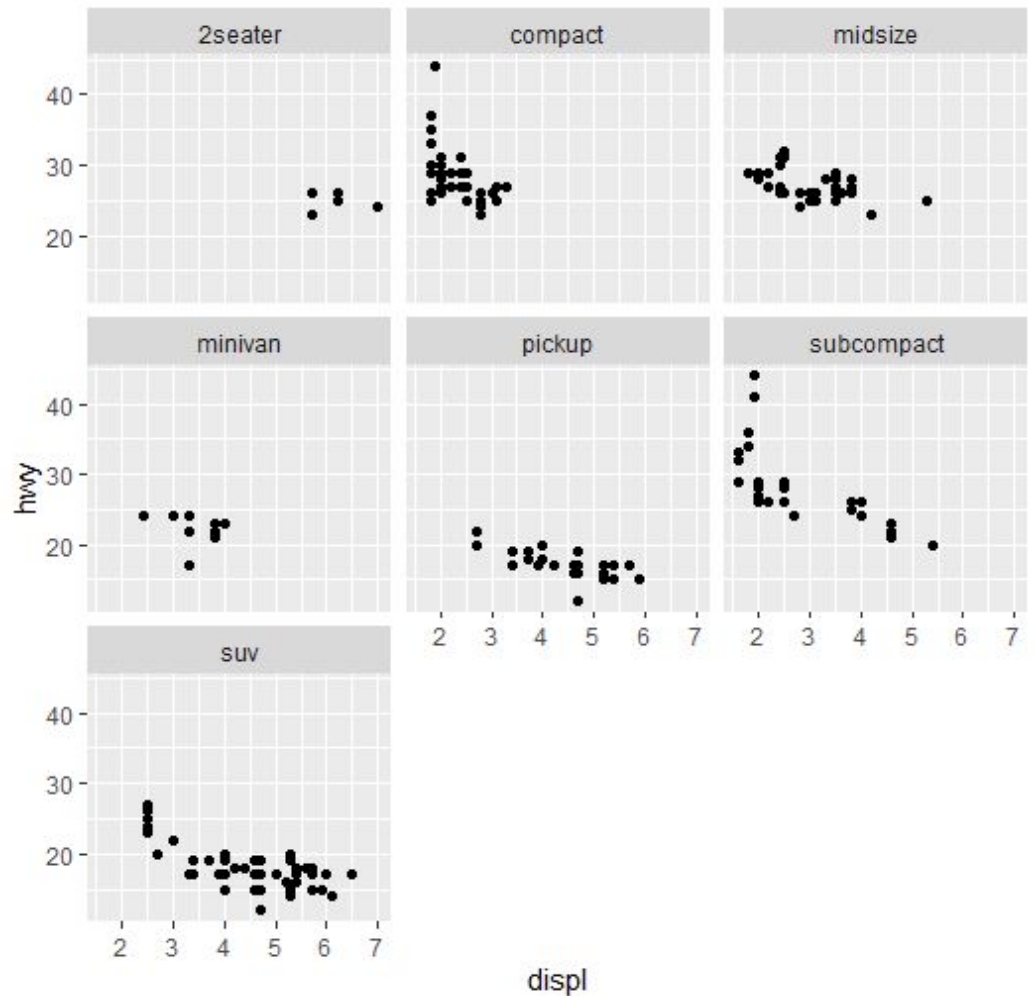
  geom_point() +

  geom_smooth()

# Facets

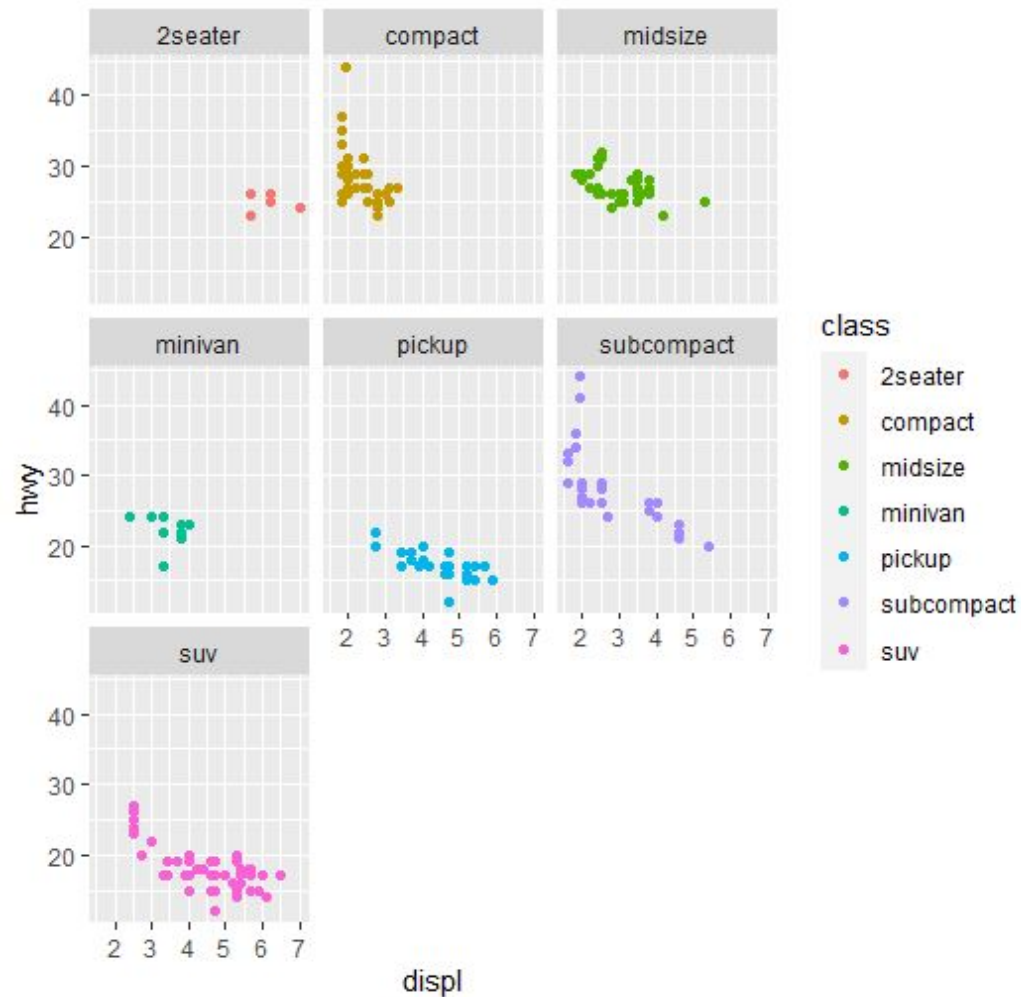Grouping by color may be hard to differentiate
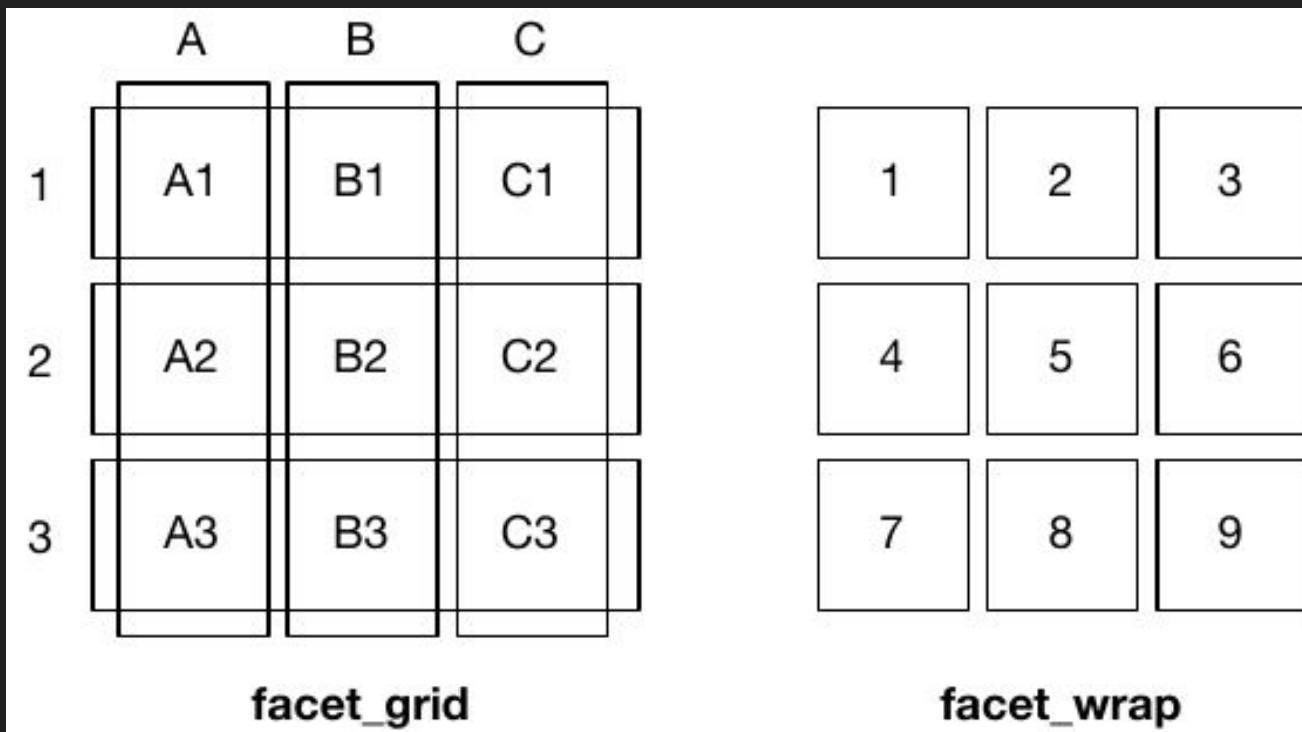
# Facets

Facets clearly show groupings

# Facets

Facets clearly show groupings

# Facets

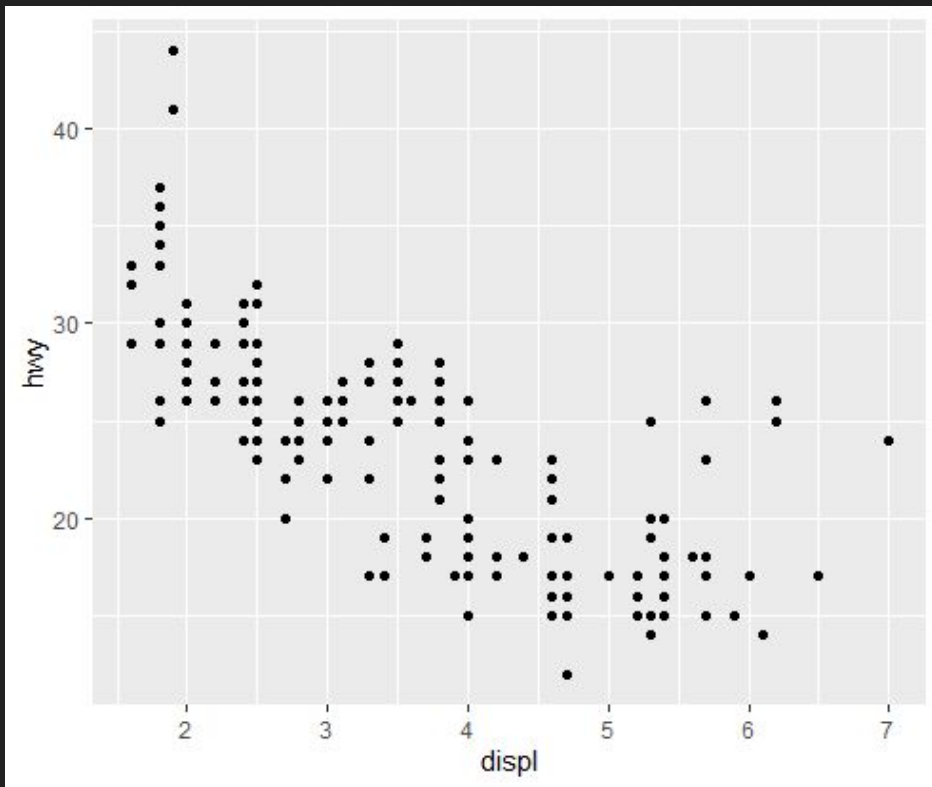**facet_wrap**
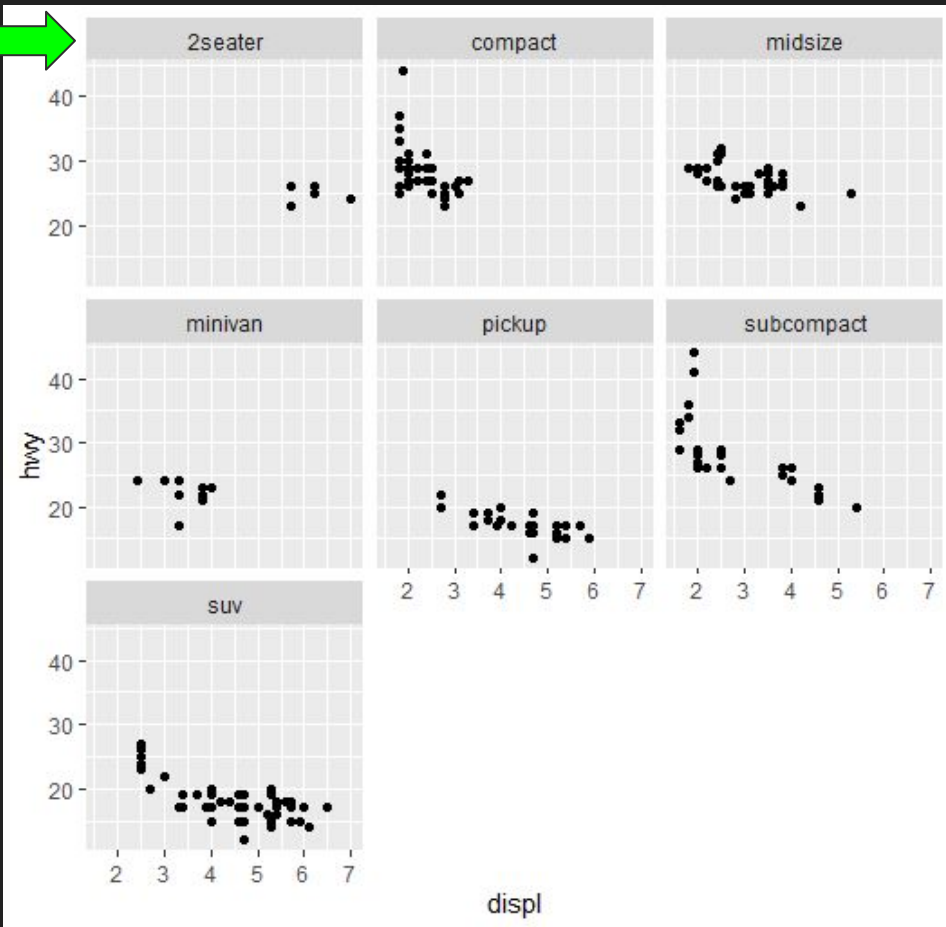
... + facet_wrap(~x)

Split the plots by
variable x

# facet_wrap

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point()

# facet_wrap

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point() +

facet_wrap(~class)

... + facet_grid(y~x)
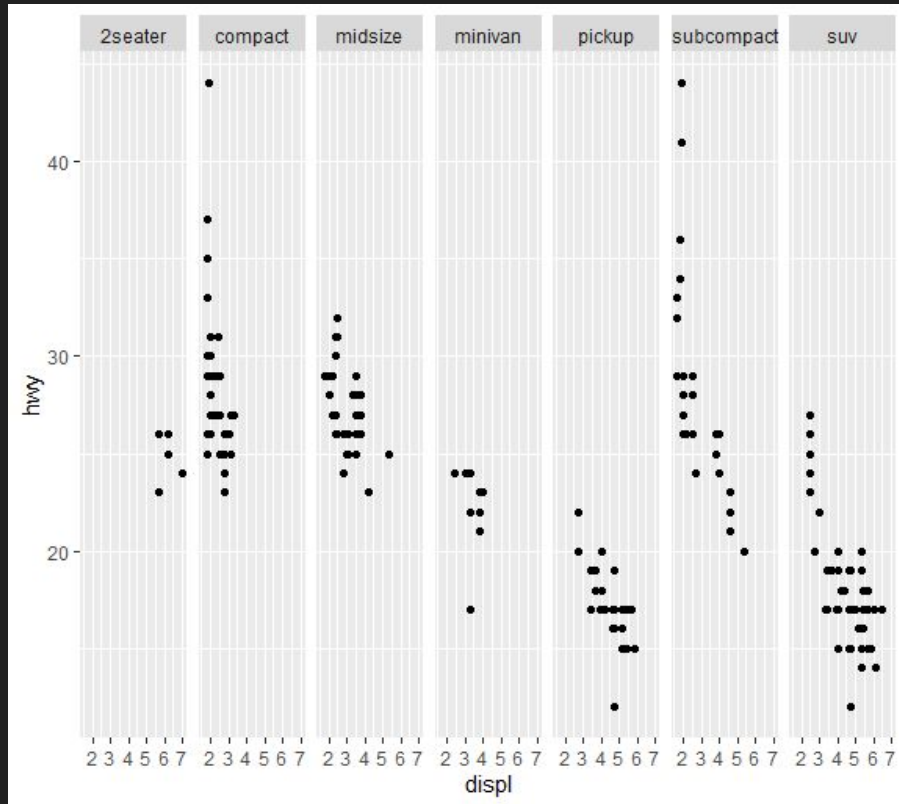
y variable on columns

x variable on rows

# facet_grid

ggplot(mpg, aes(x = displ, y = hwy)) +
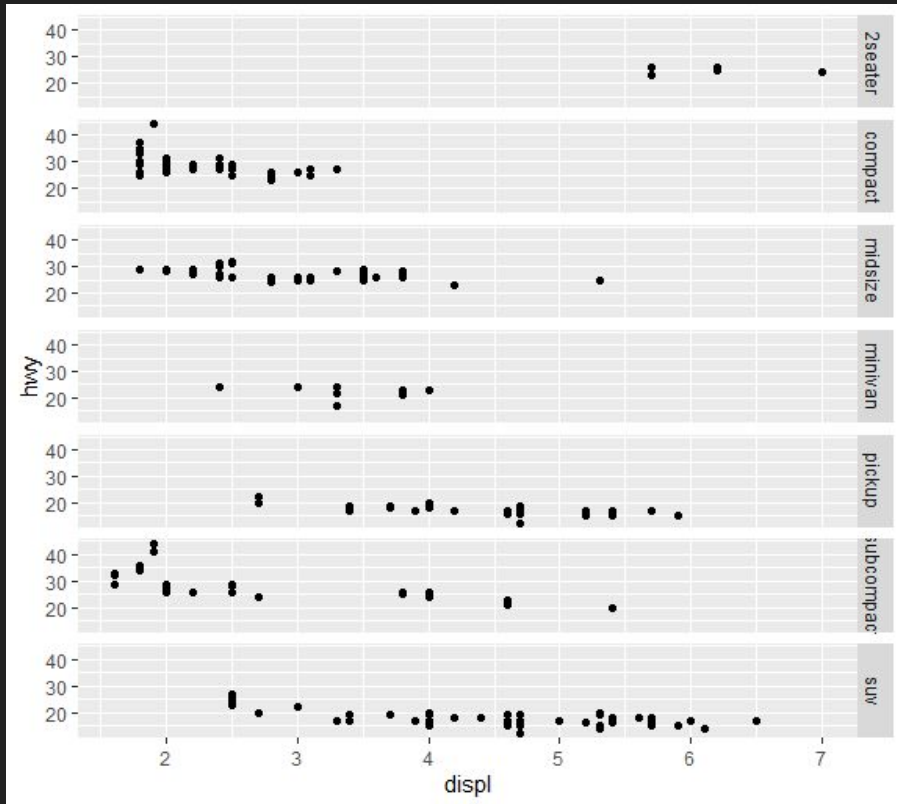
geom_point() +

facet_grid(class~.)

`class` as columns

# facet_grid



ggplot(mpg, aes(x = displ, y = hwy)) +
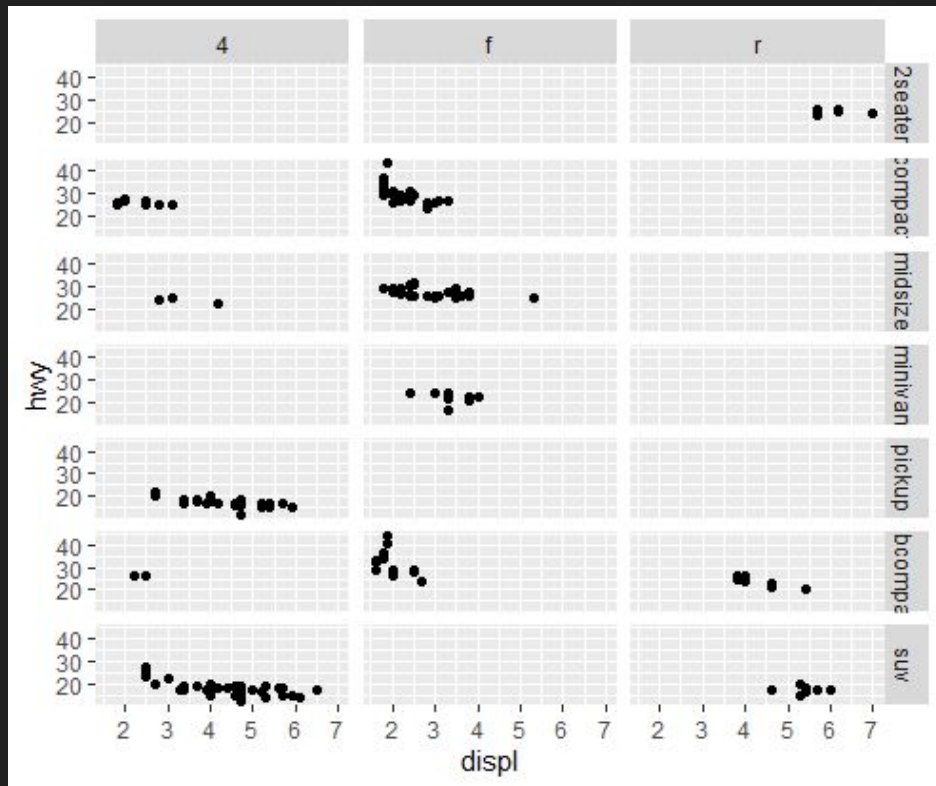
geom_point() +

facet_grid(.~class)

`class` as rows

# facet_grid

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point() +

facet_grid(drv ~ class)

`drv` as columns

`class` as rows

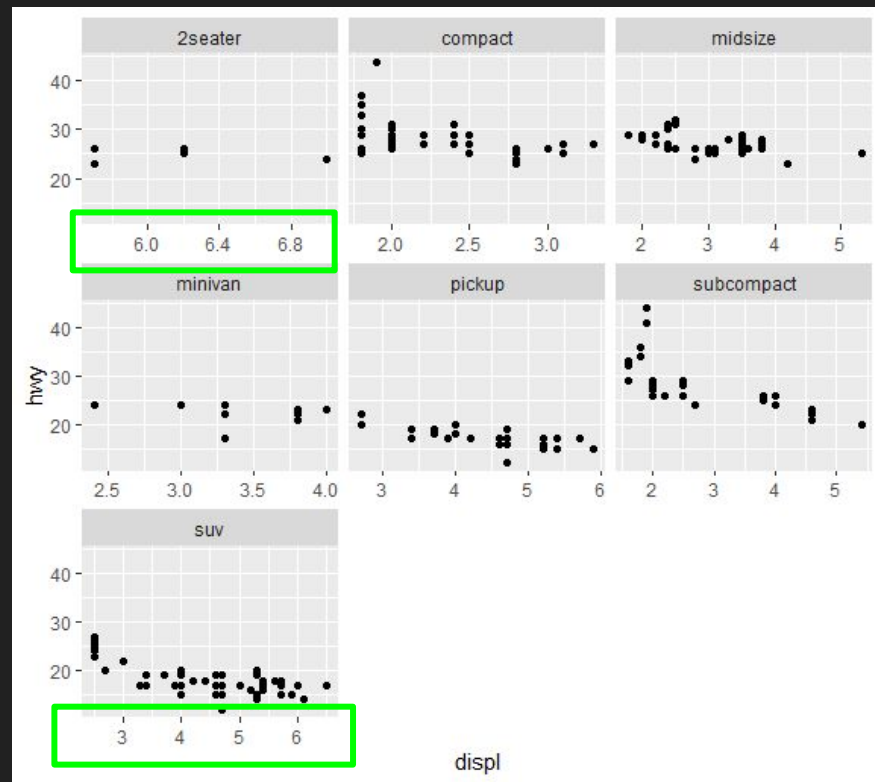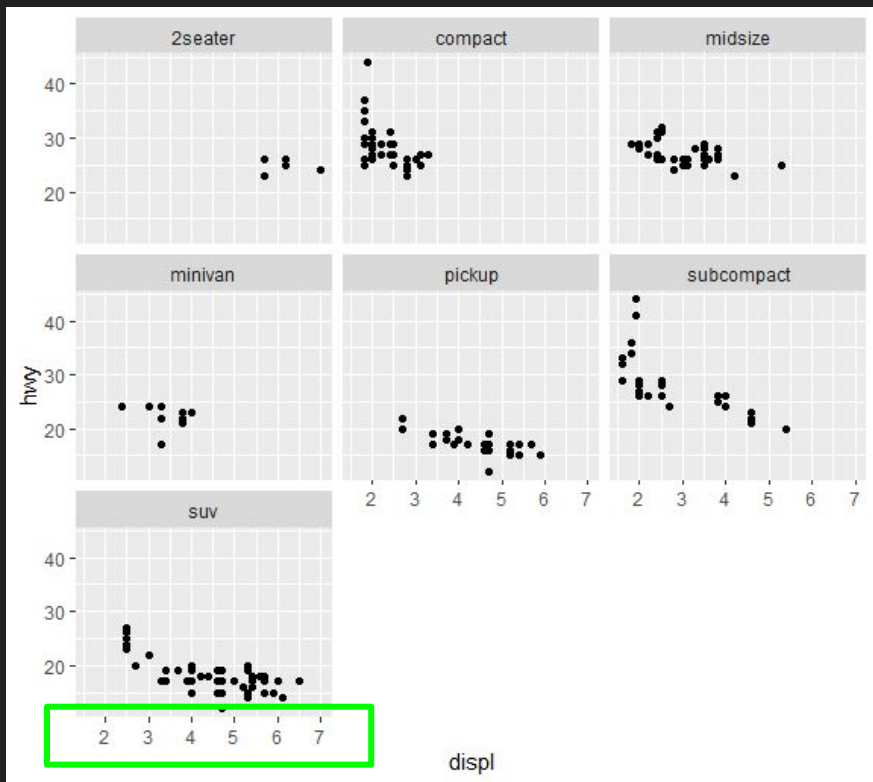## facet_grid and facet_wrap

... + facet_grid(y~x, scales = "_____")

"fixed" = x and y fixed

"free_x" = x scale free

"free_y" = y scale free

"free: = x and y free

# ... + facet_wrap(~class, scales = "free_x")

# Diamonds Example

1. Load the `diamonds` dataset:

`> data(diamonds)`          `> ?diamonds # To learn about the variables`

2. What is the most common diamond color (`color`)?

`> ggplot(_____, aes(x = _____) +`

`>   geom_____()`

3. What is the most common diamond color (`color`) for each type of cut `cut`? Use facet_wrap() or facet_grid().
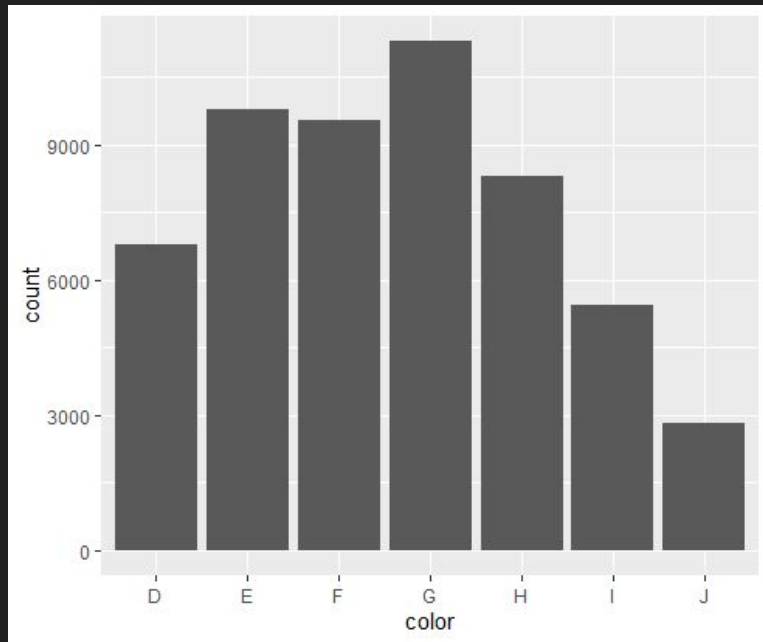
4. Set the facet scales to be "free_y", and color by `cut`

# Diamonds Example - Solutions

2. What is the most common diamond color (`color`)? Answer with a plot.
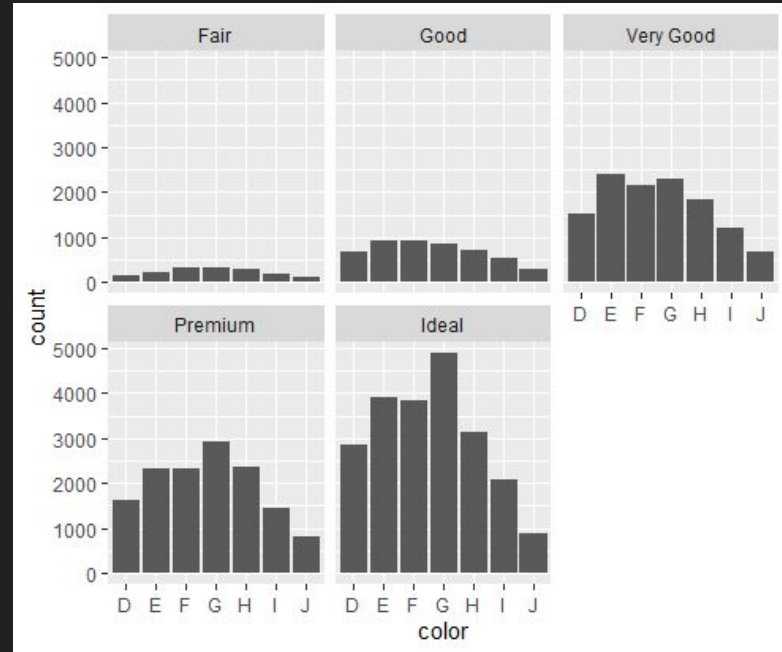
> ggplot(diamonds, aes(x = color) +

>   geom_bar()

# Diamonds Example - Solutions

3. What is the most common diamond color (`color`) for each type of cut `cut`? Use facet_wrap() or facet_grid().

> ggplot(diamonds, aes(x = color)) +

>   geom_bar() +

>   facet_wrap(~cut)

# Diamonds Example - Solutions

4. Set the facet scales to be "free_y", and color by `cut`

```
> ggplot(diamonds, aes(x = color, fill = cut)) +

>   geom_bar() +

>   facet_wrap(~cut, scales = "free_y")

> ggplot(diamonds, aes(x = color, color = cut)) +

>   geom_bar() +

>   facet_wrap(~cut, scales = "free_y")
```
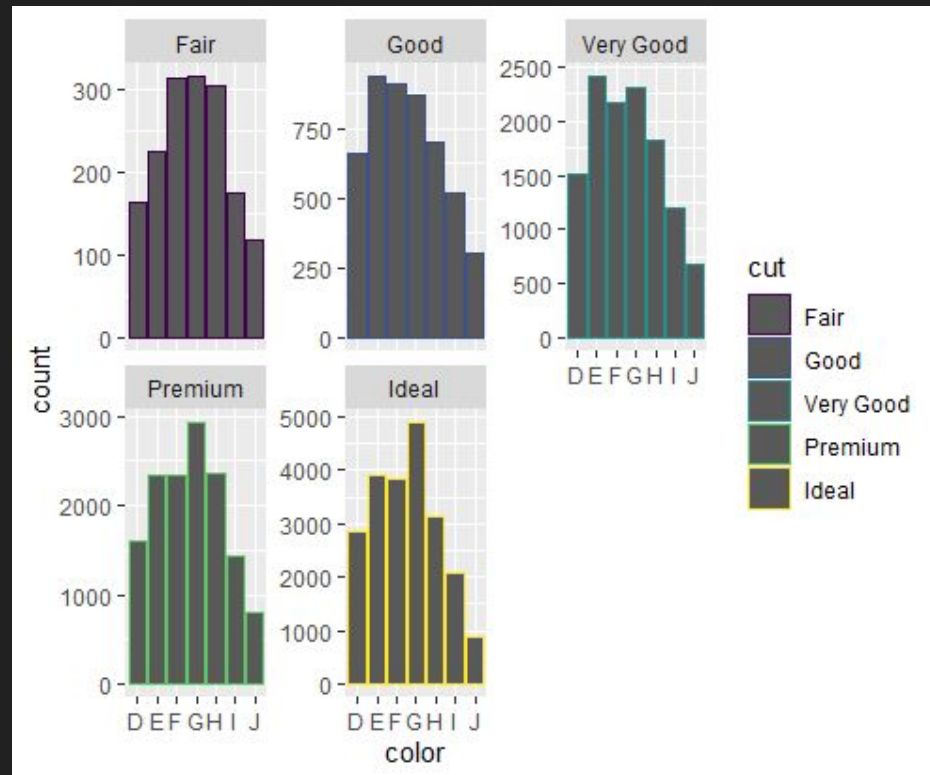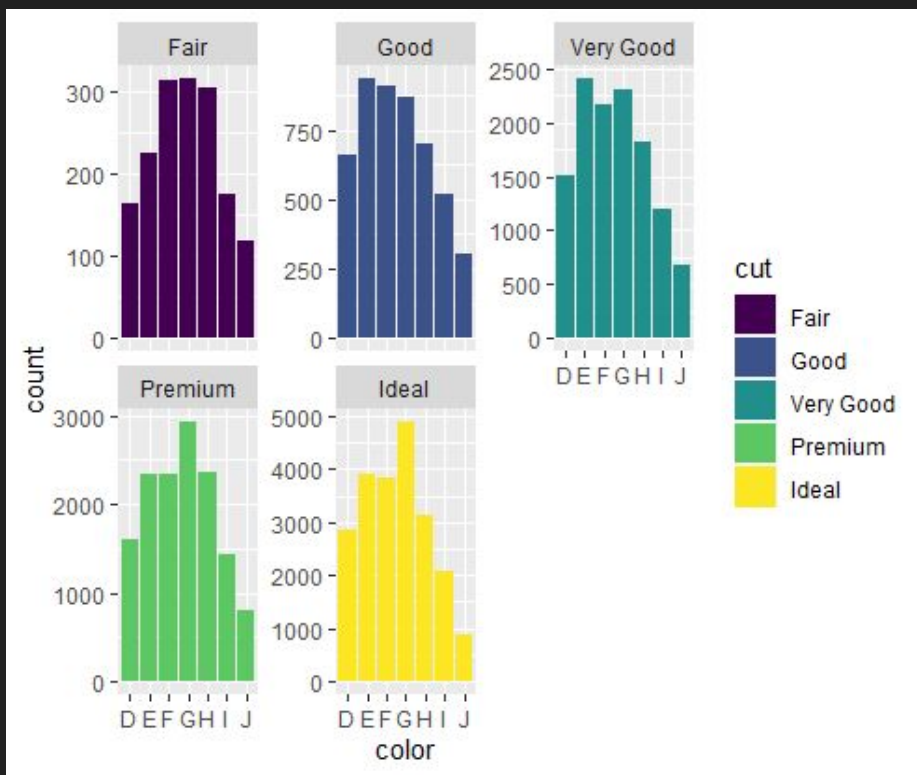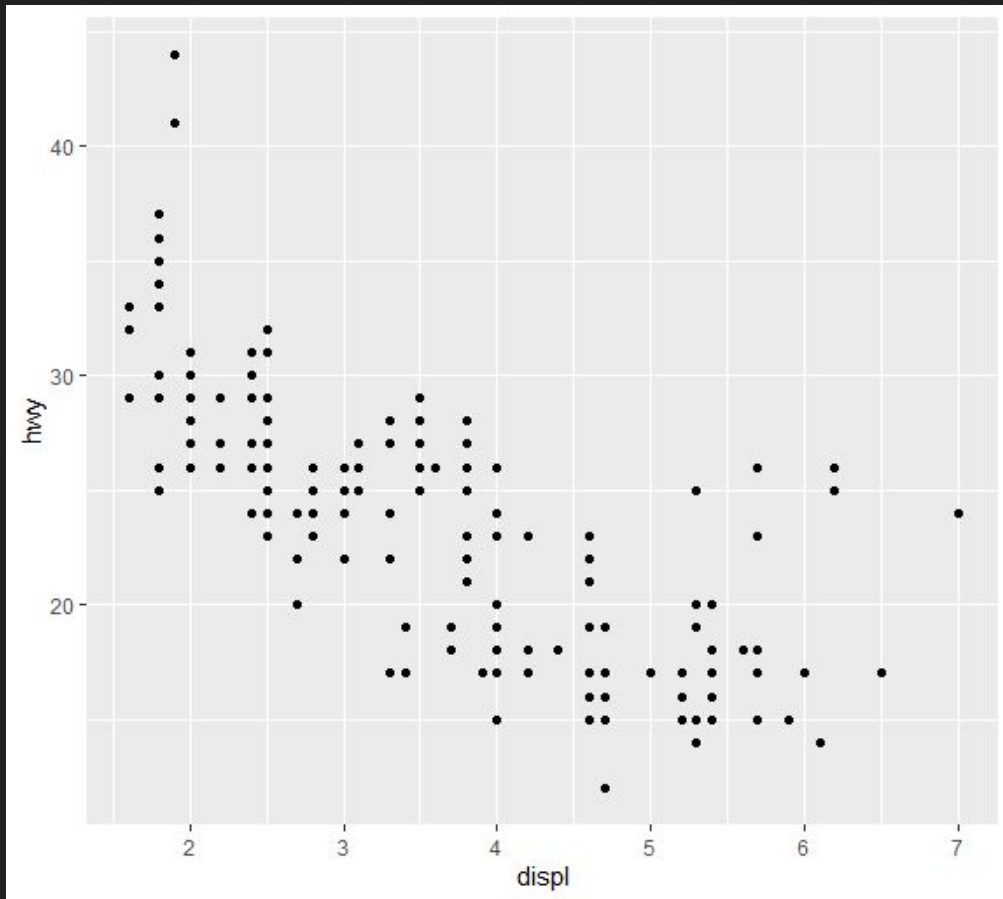
# Diamonds Example - Solutions

# labs – Adding labels

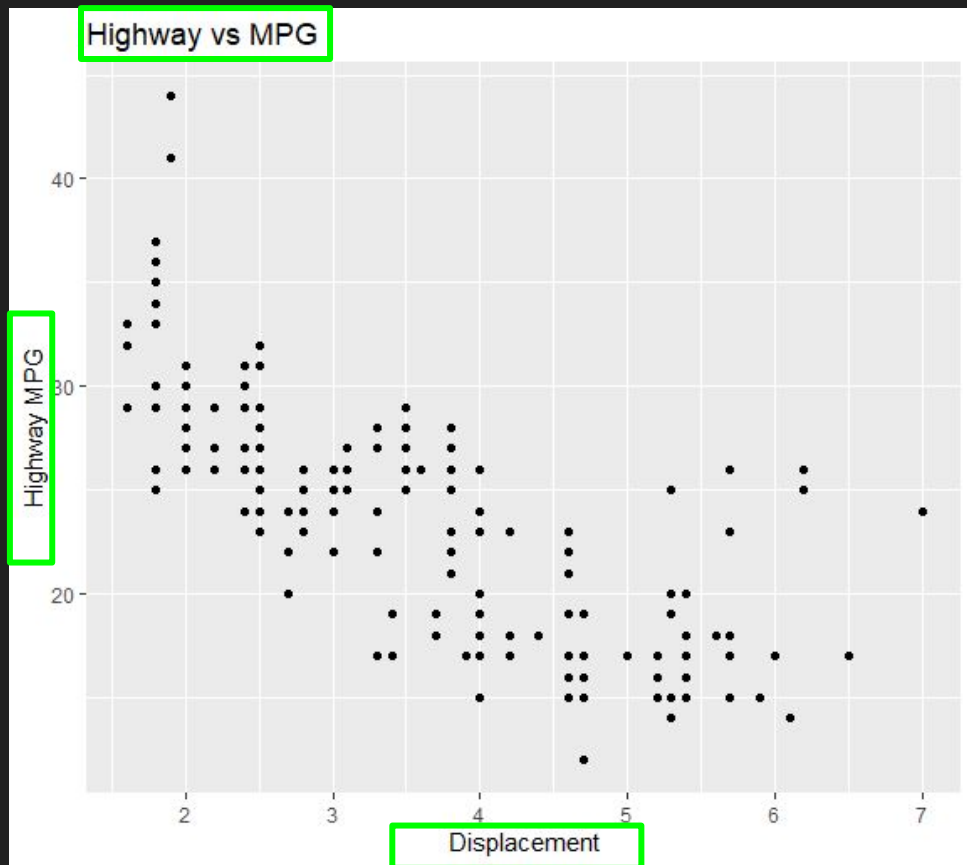ggplot(mpg, aes(x = displ, y = hwy)) +

 geom_point()

# labs - Adding labels

ggplot(mpg, aes(x = displ, y = hwy)) +

geom_point() +

labs(x = "Displacement",
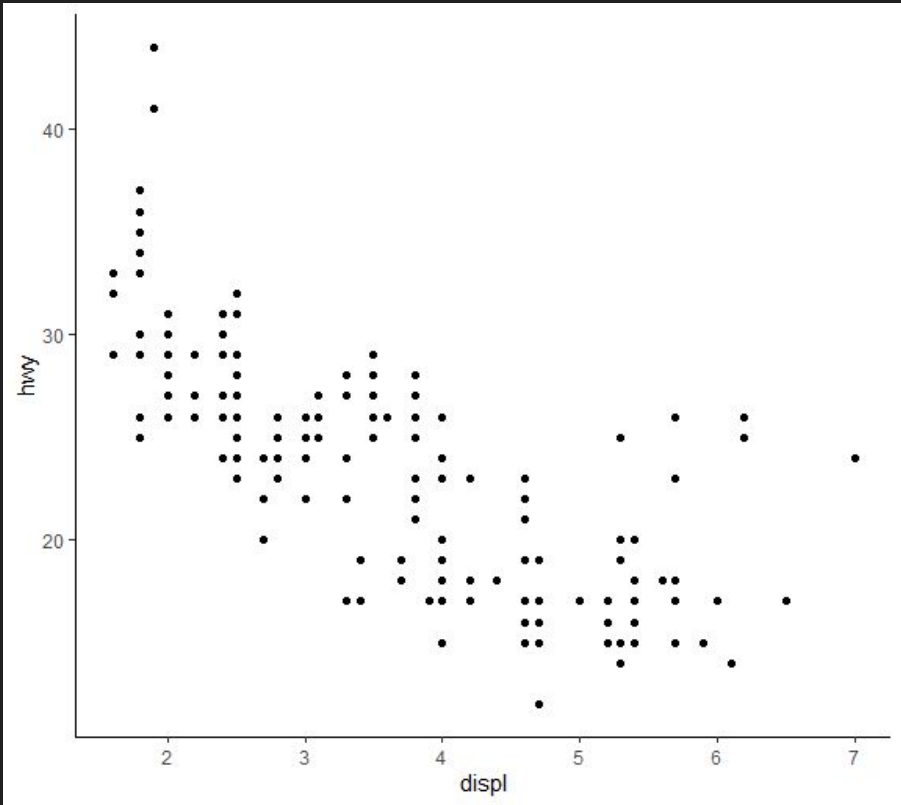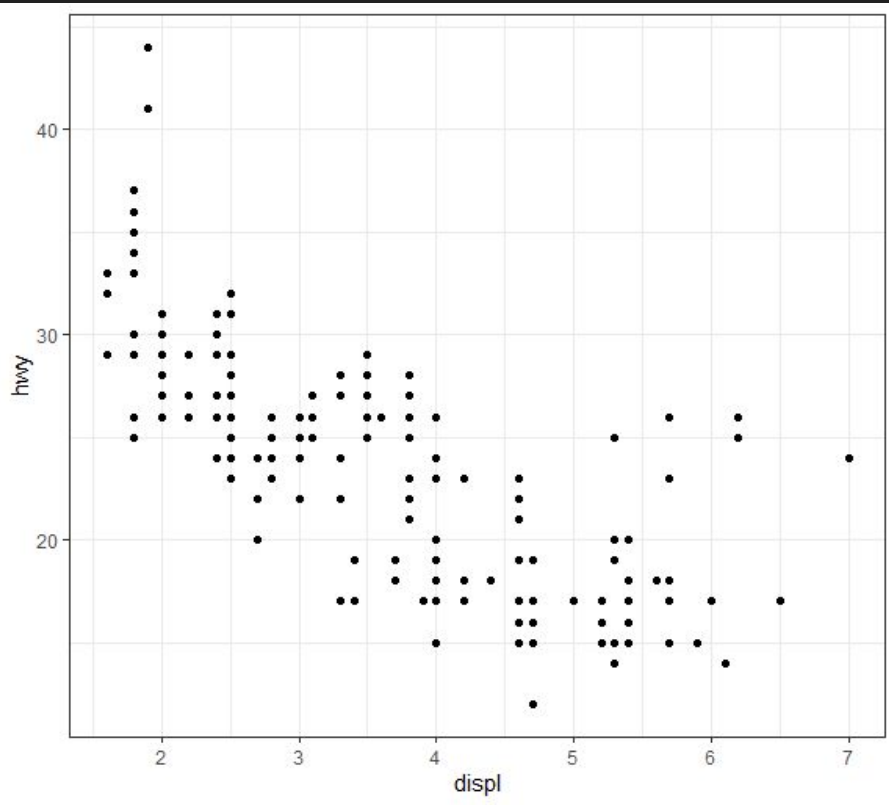
    y = "Highway MPG",

    title = "Highway vs MPG")

# Themes
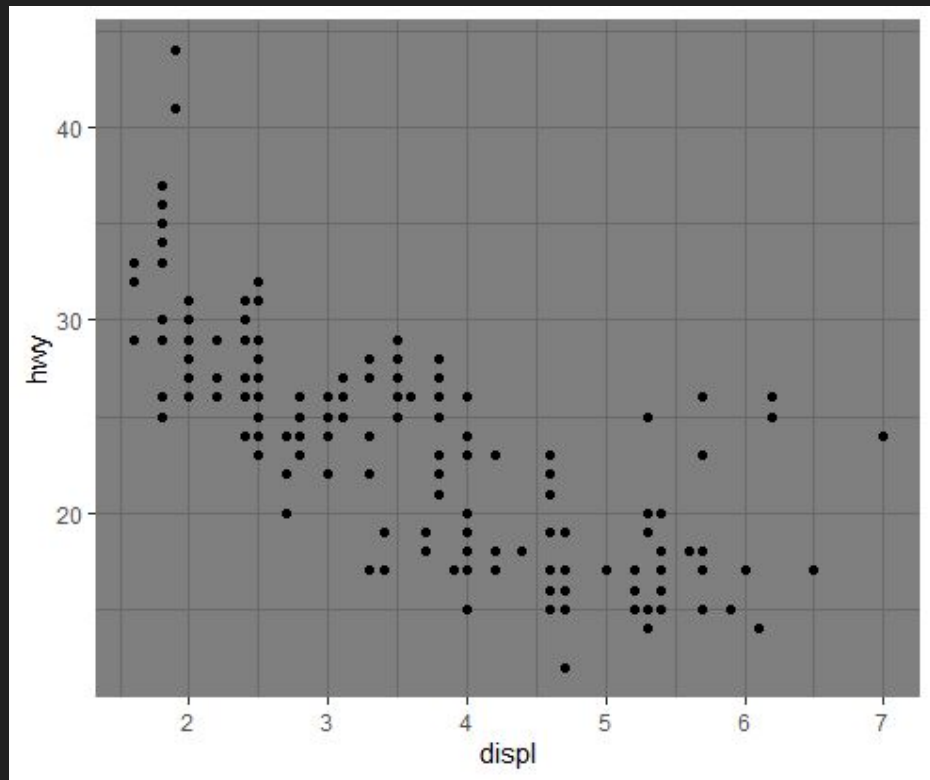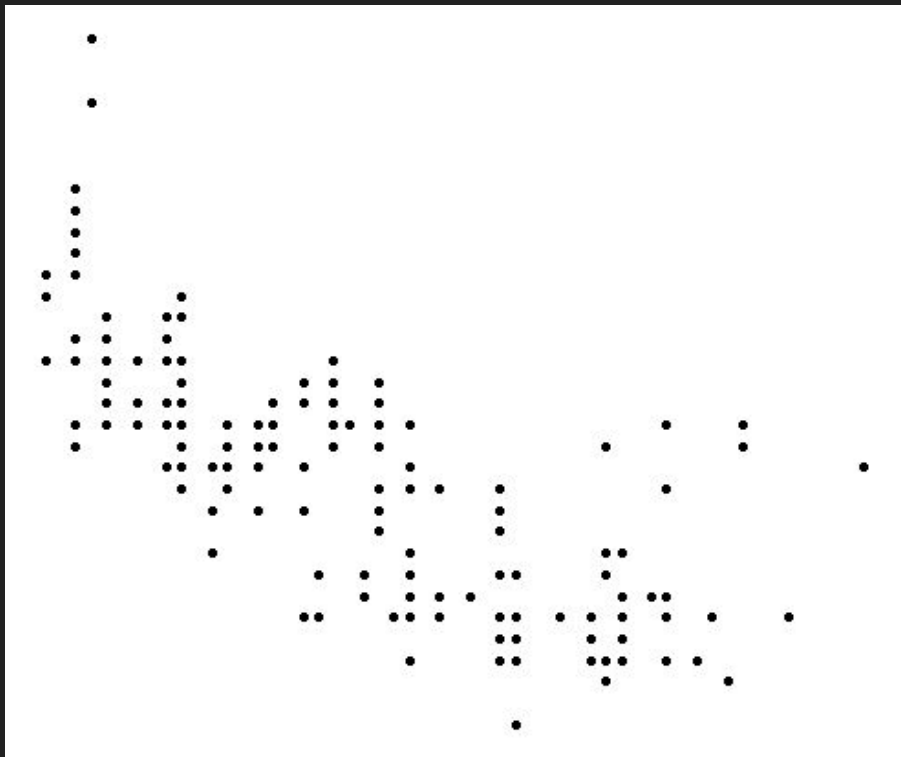
- theme_bw()
- theme_linedraw()
- theme_light()
- theme_dark()
- theme_minimal()
- theme_classic()
- theme_void()

Note: These are complete themes, you can manually adjust colors/scales/fonts too

# Themes
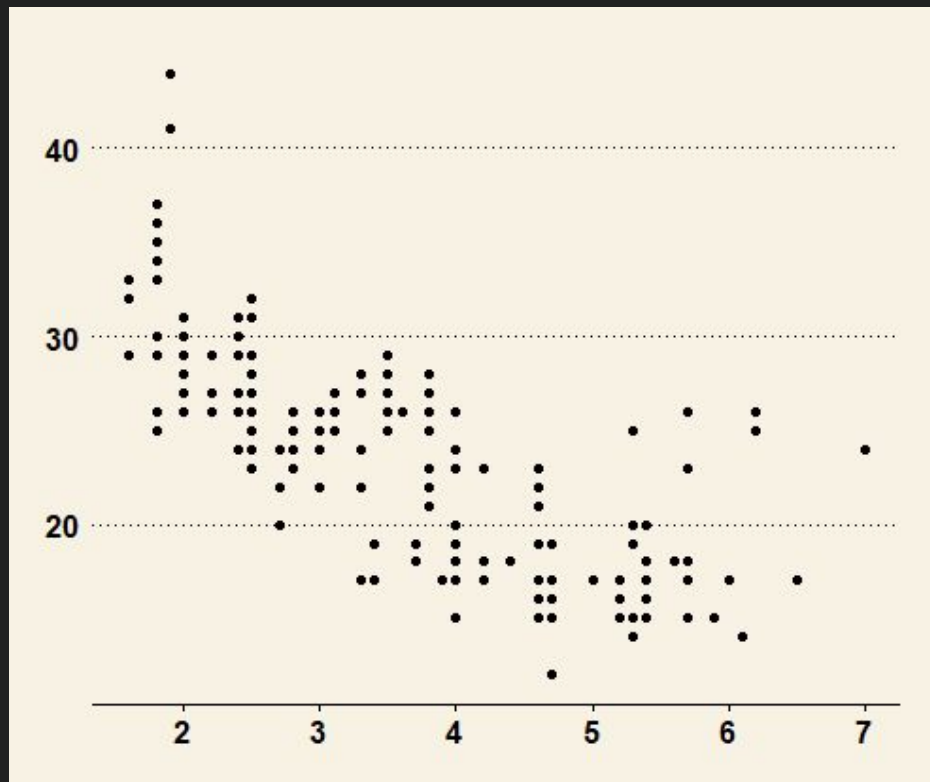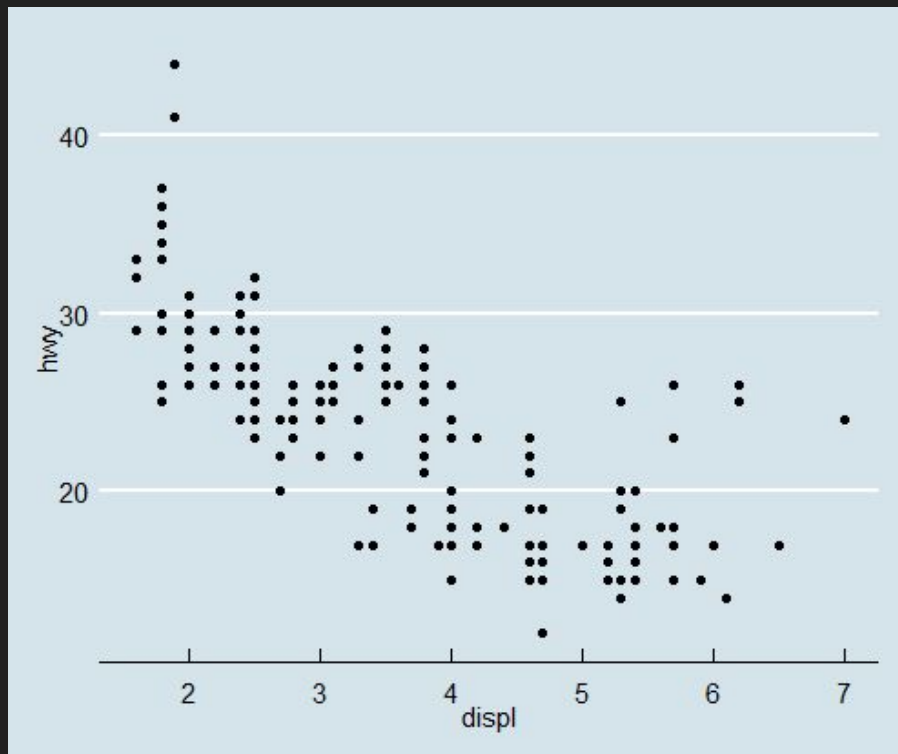
# Themes

# Themes - Extra themes from package `ggthemes`

# Exporting ggplots

ggsave("plotname.png")

Saves most recently displayed ggplot

Can save as a png, jpg, pdf, svg, etc.

Can specify height and width of output

# The Little Things

1. One plot = one primary message.
2. Do not use the default ggplot2 theme
3. Adjust axes as necessary
   a. Extreme outliers breaking scales - `scale_x_log10()`
   b. Manually select tick marks - `scale_<x, y>_manual()`
4. Rotate x labels if they are overlapping
   a. `+ theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))`
5. Alternative: Flip x and y axes if x has a lot of words
   a. `coord_flip()`
6. Want animated plots? Use the `**gganimate**` package

# Useful Resources

*Fundamentals of Data Visualization* by Claus O. Wilke

- https://clauswilke.com/dataviz/

*ggplot2: Elegant Graphics for Data Analysis* by Hadley Wickham

- https://ggplot2-book.org/index.html

# Good luck!

Fill out the feedback survey!