# Introduction to R and the tidyverse

## – plotting part 1 –

Paolo Crosetto

# Why plot?

# Why do we plot

> *Why do we want to `plot` data?*

- we are human beings – we are `pattern recognizers`
- we can *see* things we are not able to grasp from data
- good to `explore` a dataset and look for regularities
- good to `convey` *a clear message*
- to have `fun`

# Why plot? Eyeballing

`Eyeballing` the data first is **always** a good idea

- data could look `similar` at a first glance

- and even have similar descriptive statistics

- but still be *very*`different` in practice

# An example

- data contains vars x and y, over 13 different conditions

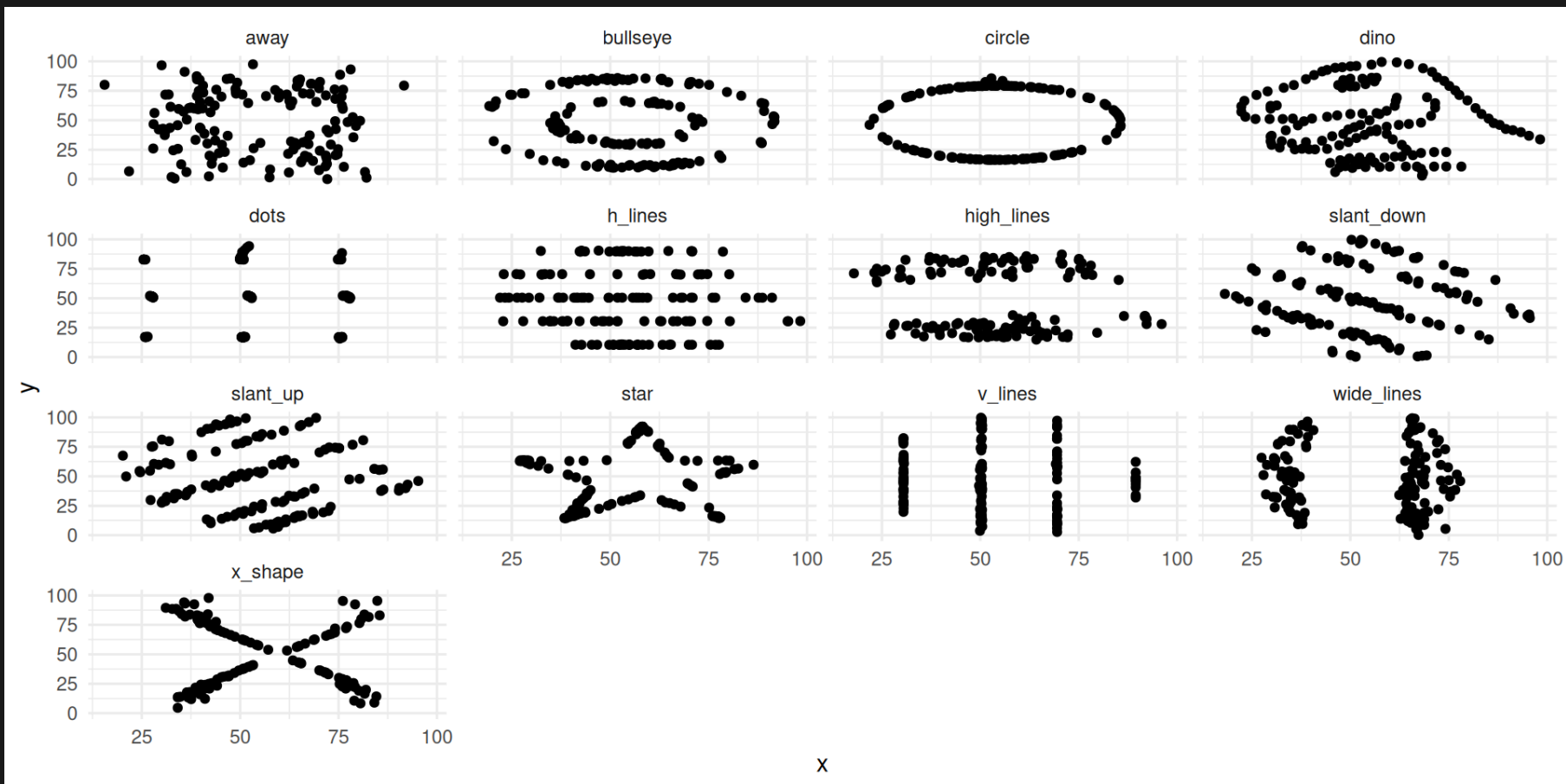- import `data/plotme.tsv`, compute $\mu$, $\sigma$ by `dataset`

# An example

- data contains vars `x` and `y`, over 13 different conditions
- import `data/plotme.tsv`, compute $\mu$, $\sigma$ by `dataset`

| dataset | mean_x | sd_x | mean_y | sd_y |
|---------|--------|------|--------|------|
| away | 54.27 | 16.77 | 47.83 | 26.94 |
| bullseye | 54.27 | 16.77 | 47.83 | 26.94 |
| circle | 54.27 | 16.76 | 47.84 | 26.93 |
| dino | 54.26 | 16.77 | 47.83 | 26.94 |
| dots | 54.26 | 16.77 | 47.84 | 26.93 |

# Now let's plot this!

But if you `plot` it, you'll see **stark** differences

# Why plot? Compact information

Plotting allows you to convey a `lot of info`

- humans are `pattern recognizers`

- several `geometric` objects can convey meaning

  - position (x,y,z)

  - color, size, shape

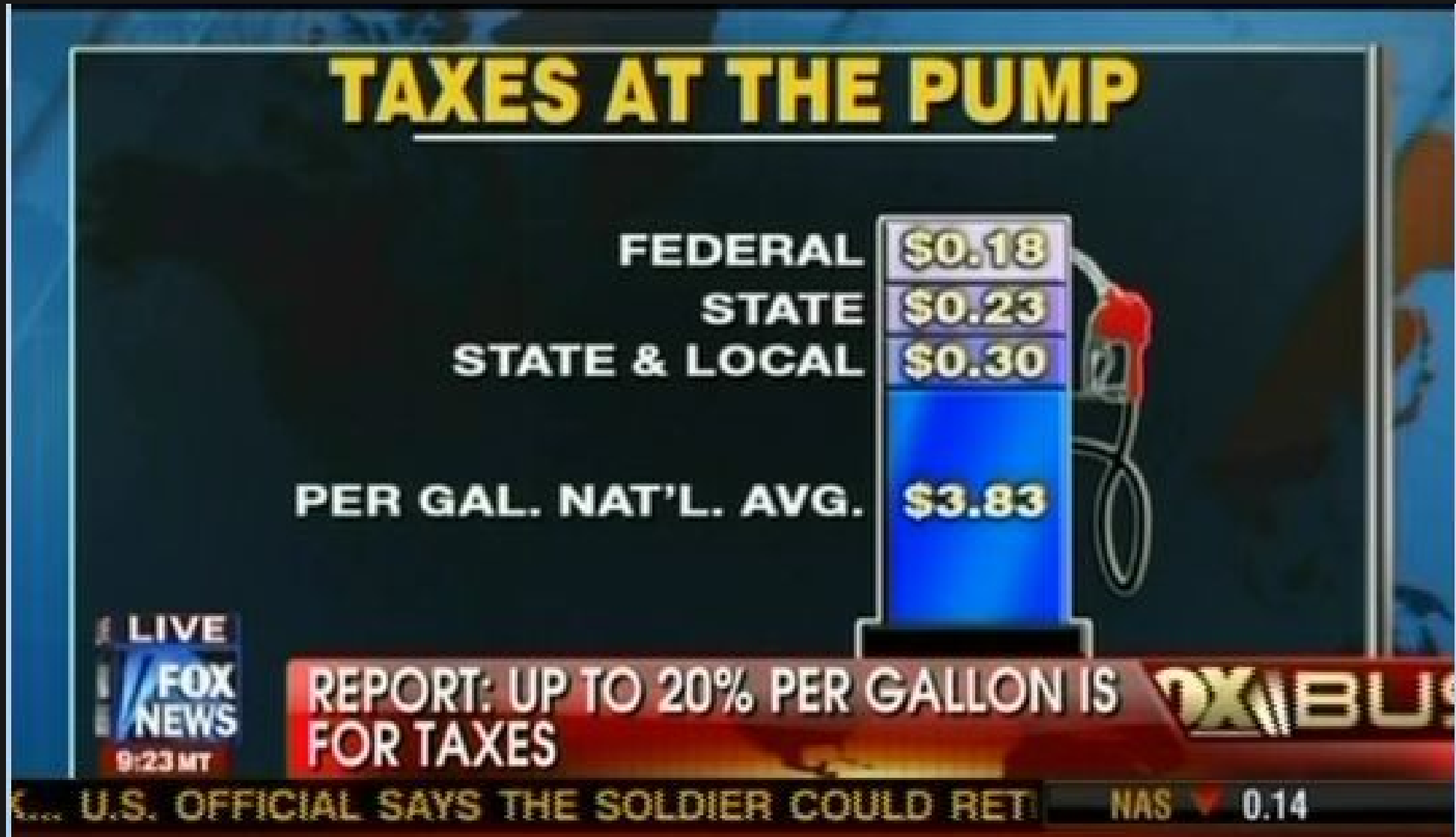- you can combine multiple plots into `infographics`

# Good and bad plots
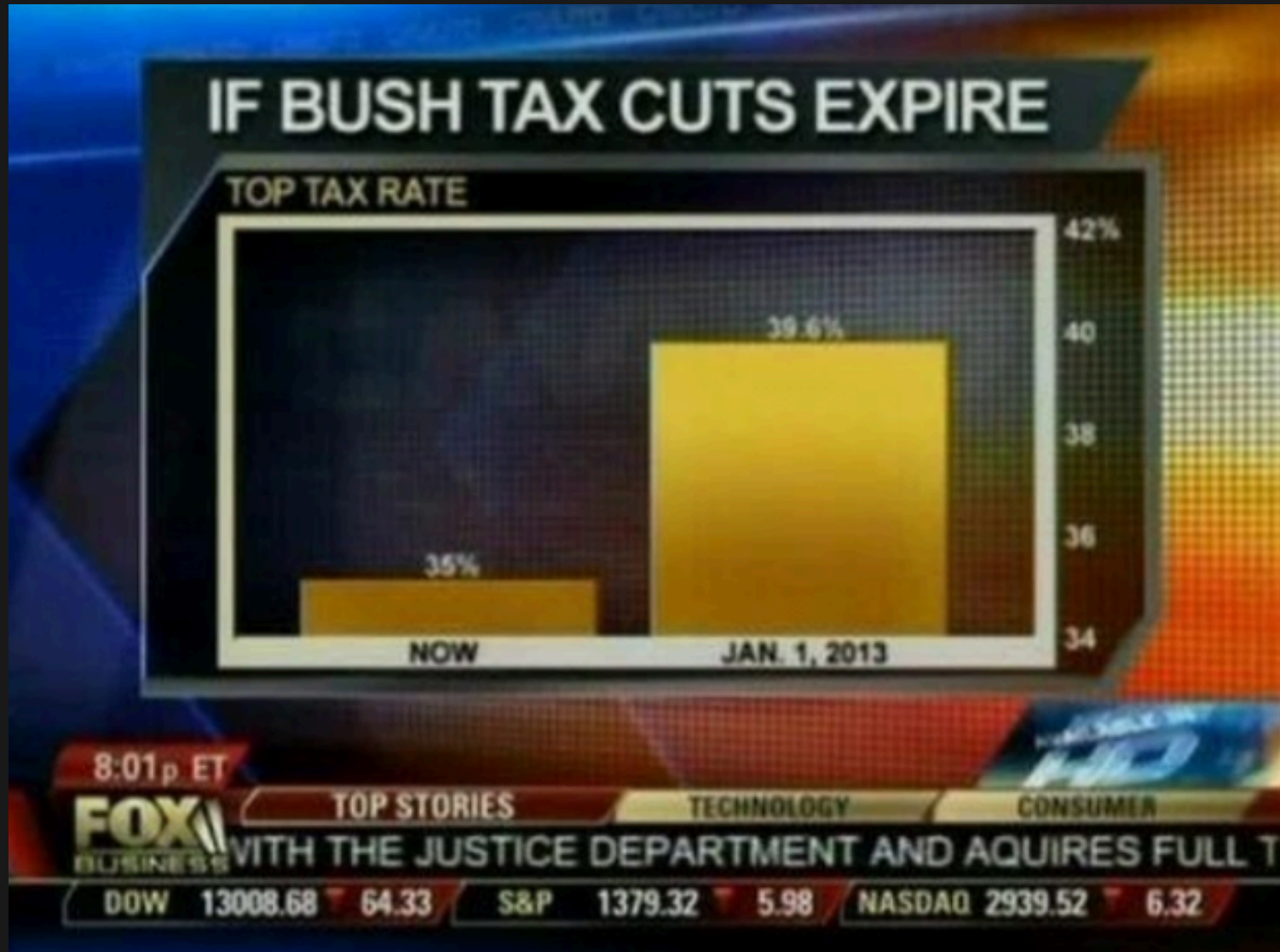
# Good plots, bad plots

- It is important to make *good* plots

- i.e., plots that *look good*…

- …and are *honest* to the data

- it is `very easy to hide` the message rather than *highlighting* it

- it is `very easy to mislead` with a plot

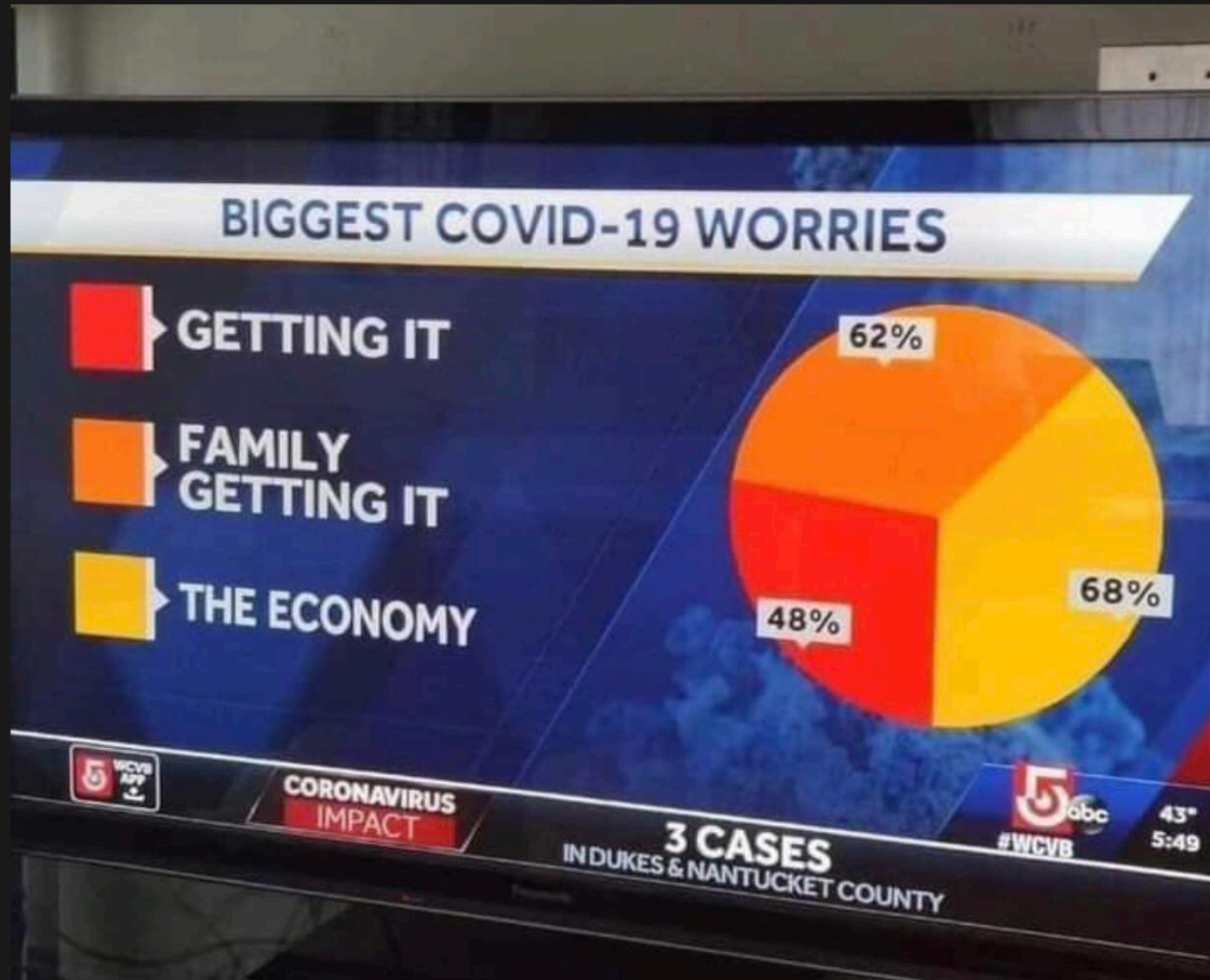let's start with **bad plots**. *Why* are they bad?
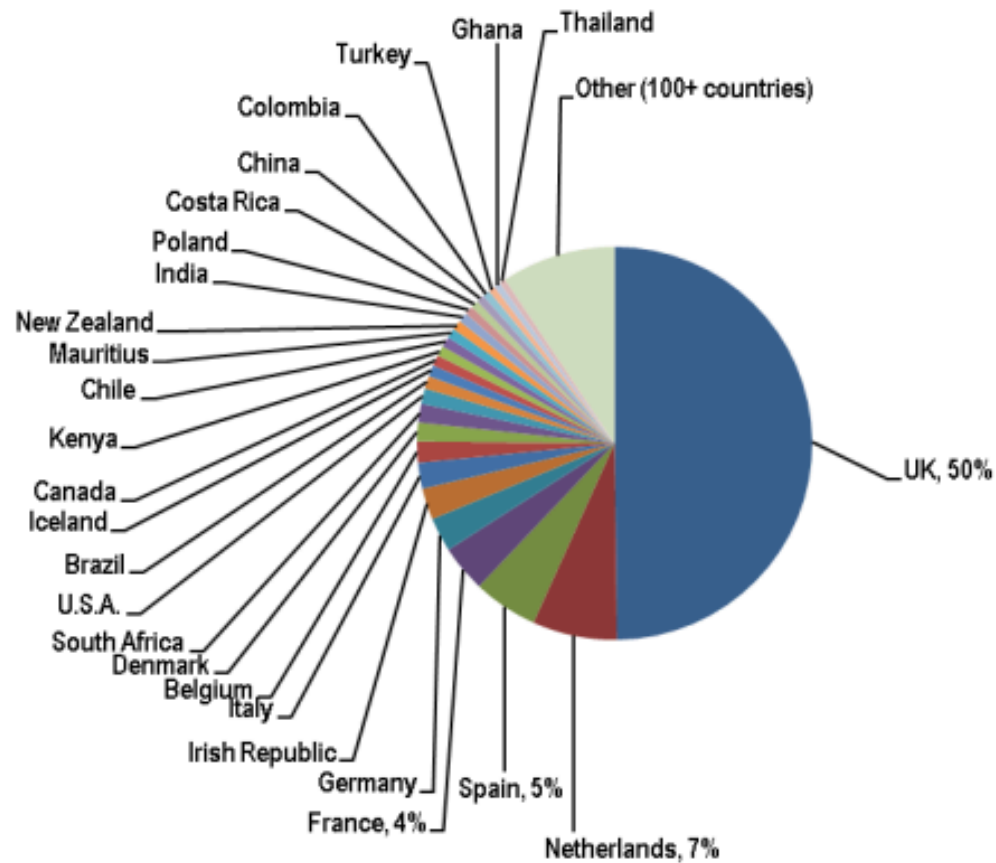
# Bad plotting 1

# Bad plotting 2

# Bad plotting 3

# Bad plotting 4



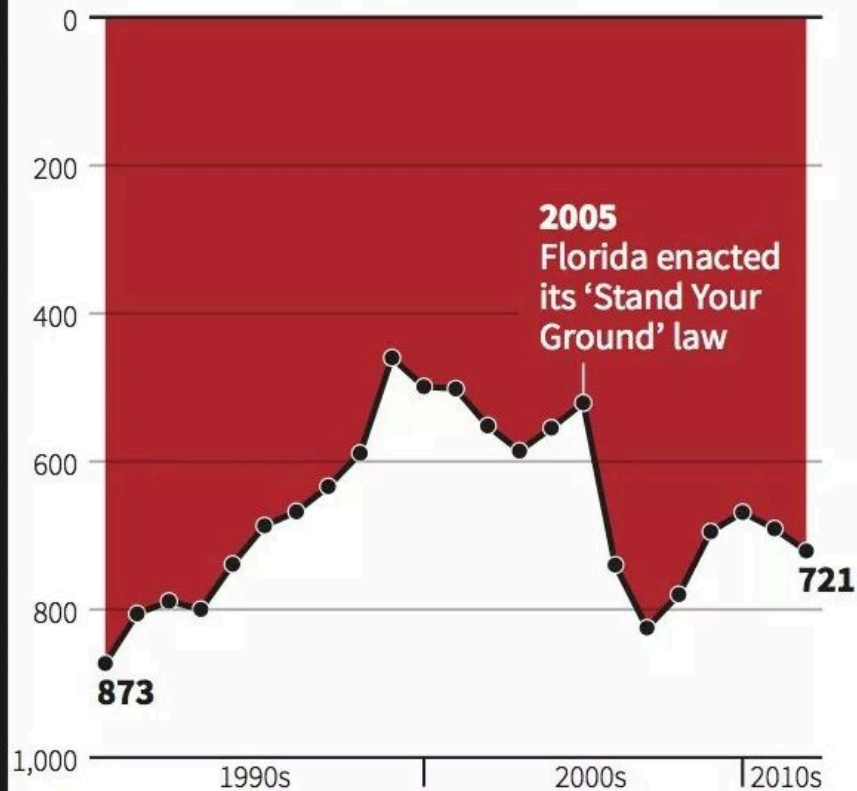Origins of food consumed in the UK by value: 2007

UK, 50%
Netherlands, 7%
Spain, 5%
France, 4%
Germany
Irish Republic
Italy
Belgium
Denmark
South Africa
U.S.A.
Brazil
Iceland
Canada
Kenya
Chile
Mauritius
New Zealand
India
Poland
Costa Rica
China
Colombia
Turkey
Ghana
Thailand
Other (100+ countries)

Based on the farm-gate value of unprocessed food

# Bad plotting 5



Gun deaths in Florida

Number of murders committed using firearms

2005
Florida enacted
its 'Stand Your
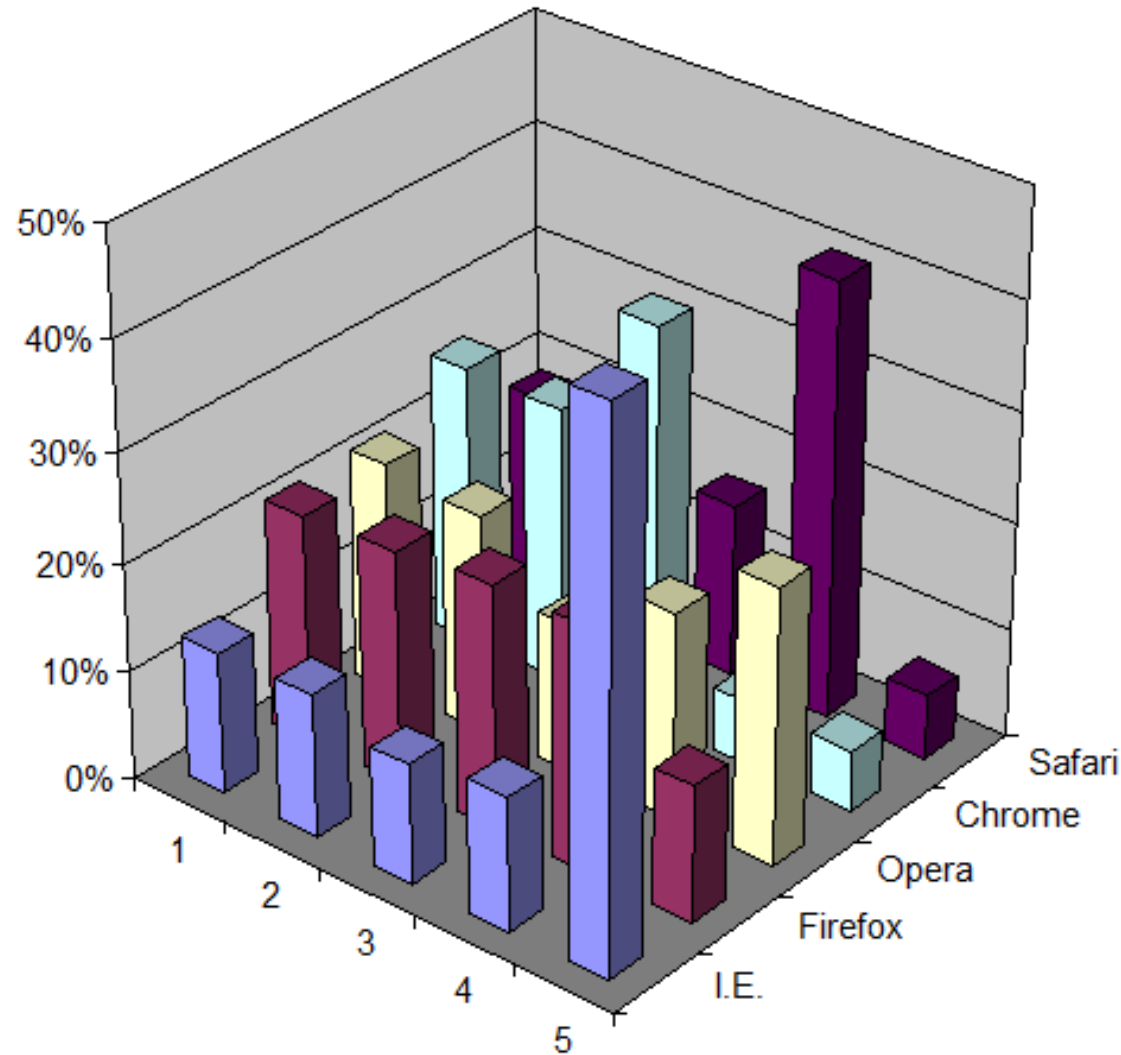Ground' law

873

721

1990s    2000s    2010s

Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

REUTERS

# Bad plotting 6



Number of singles sold

1995
1996
1997
1998

# Bad plotting 7 (really, NO 3D plots)

# The road to good plotting

- `know` your data

- `think` before you hit the enter button

- `sketch` on paper first

- be `honest`

- draw your `axis` first

- choose your `visualization` wisely

- a good plot: lots of precise information in a concise way.

# Good plots, 1

# Good plots, 2



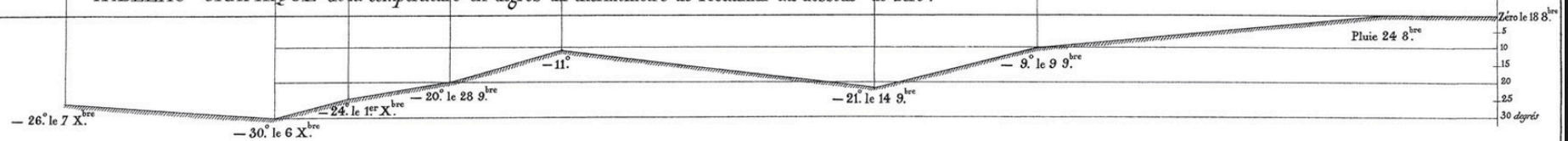Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812~1813.

# Good plots, 3



**Measles**

Vaccine introduced

0k 1k 2k 3k 4k
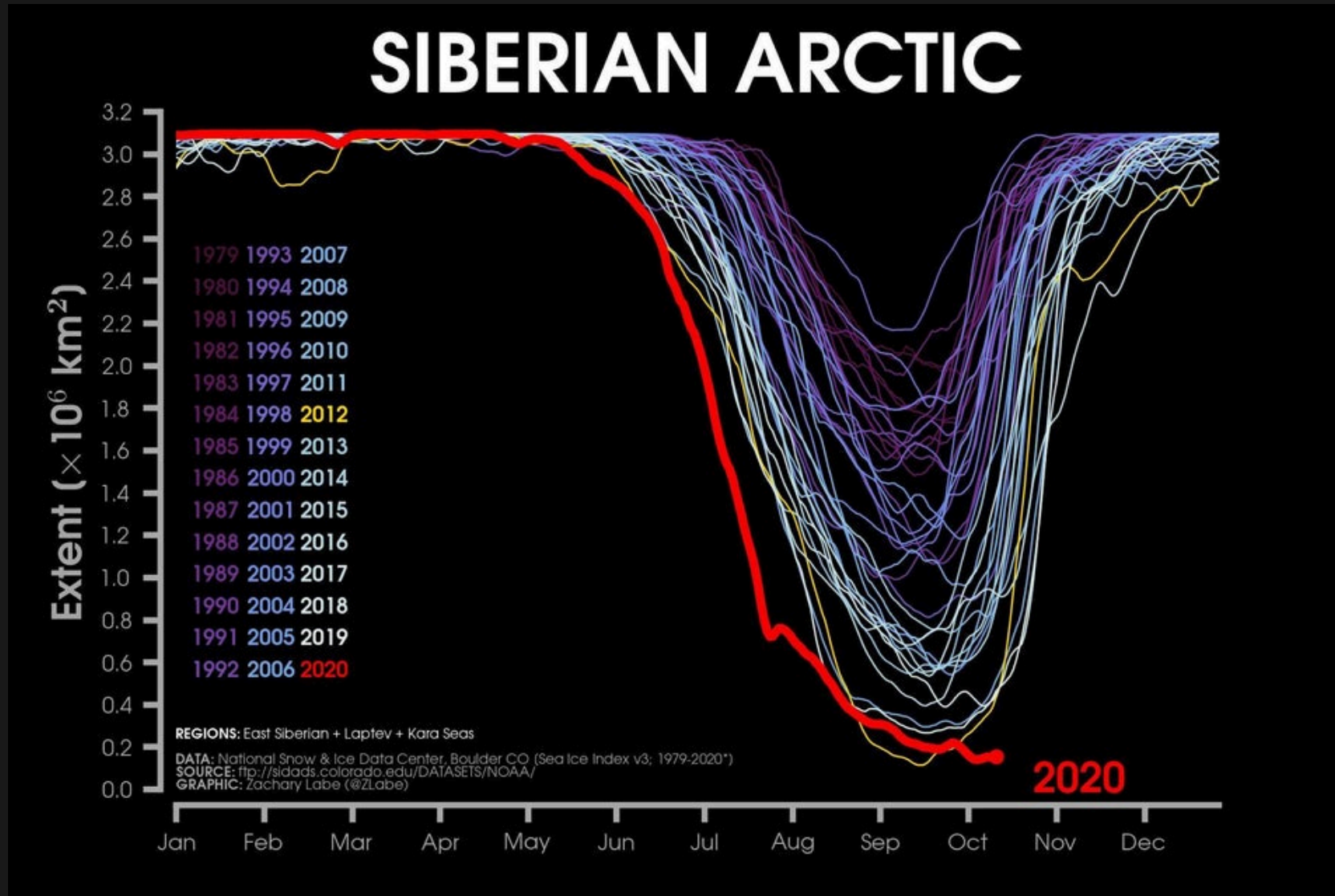
# Good plots, 4

# ggplot2: the basics

# Some data

We will start by using the *built-in dataset* **mpg**

```
1  mpg
```

```
# A tibble: 234 × 11
   manufacturer model    displ  year   cyl trans drv     cty   hwy fl
class
   <chr>        <chr>    <dbl> <int> <int> <chr> <chr> <int> <int> <chr>
<chr>
 1 audi         a4         1.8  1999     4 auto… f        18    29 p
comp…
 2 audi         a4         1.8  1999     4 manu… f        21    29 p
comp…
 3 audi         a4         2    2008     4 manu… f        20    31 p
comp…
 4 audi         a4         2    2008     4 auto… f        21    30 p
comp…
 5 audi         a4         2.8  1999     6 auto… f        16    26 p
```

# A look at the data

# A look at the data

```r
1  skimr::skim(mpg)
```

## Data summary

| Name | mpg |
|---|---|
| Number of rows | 234 |
| Number of columns | 11 |
| | |
| _____ | |
| Column type frequency: | |
| character | 6 |
| numeric | 5 |

| | Group variables | None |
|---|---|---|

## Variable type: character

| skim_variable | n_missing | complete_rate | min | max | e |
|---|---|---|---|---|---|
| manufacturer | 0 | 1 | 4 | 10 | |
| model | 0 | 1 | 2 | 22 | |
| trans | 0 | 1 | 8 | 10 | |
| drv | 0 | 1 | 1 | 1 | |
| fl | 0 | 1 | 1 | 1 | |
| class | 0 | 1 | 3 | 10 | |

## Variable type: numeric

| skim_variable | n_missing | complete_rate | mean | sd |
| --- | --- | --- | --- | --- |
| displ | 0 | 1 | 3.47 | 1.29 |
| year | 0 | 1 | 2003.50 | 4.51 |
| cyl | 0 | 1 | 5.89 | 1.61 |
| cty | 0 | 1 | 16.86 | 4.26 |
| hwy | 0 | 1 | 23.44 | 5.95 |

# Why `ggplot2`?

Advantages of `ggplot2`

- consistent underlying `grammar of graphics`

- plot specification at a high level of `abstraction`

- very flexible

- mature and complete graphics system

- `theme` system for polishing plot appearance

- many users, active, fast & competent `support`

- arguably the `best` plotting system on the planet

# Grammar of graphics

Independently specify plot `building blocks` & combine them to create *any* plot.
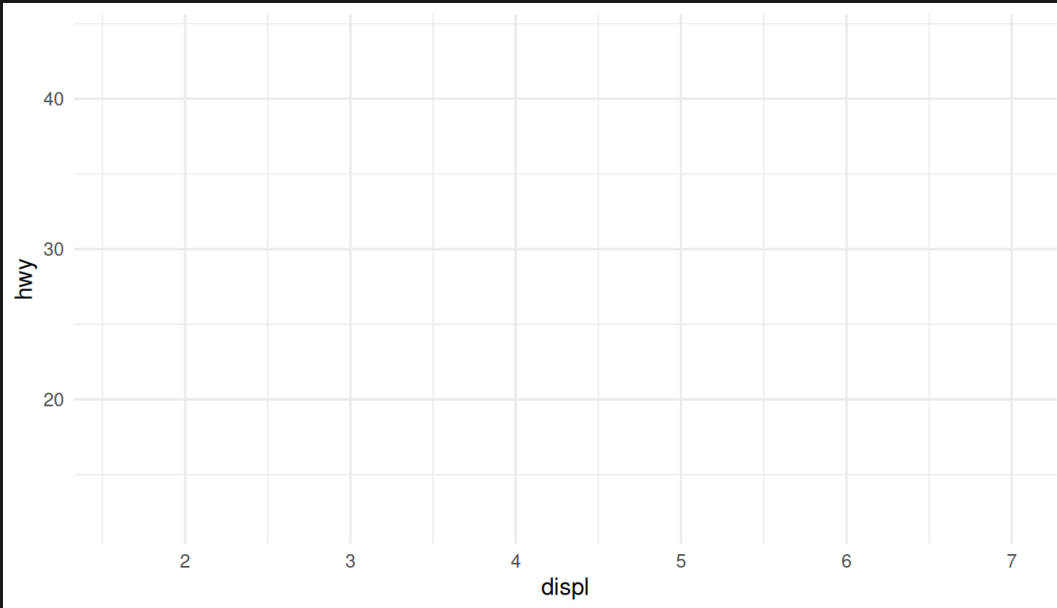
# Starting from the basics

As in a grammar the minimal sentence is a subject in a plot the minimal object is data

```r
1 p <- ggplot(mpg)
```

# basics

## In a grammar, you need a verb. In plots, this is axis

```
1  p <- p + aes(x = displ, y =hwy)
2  p
```
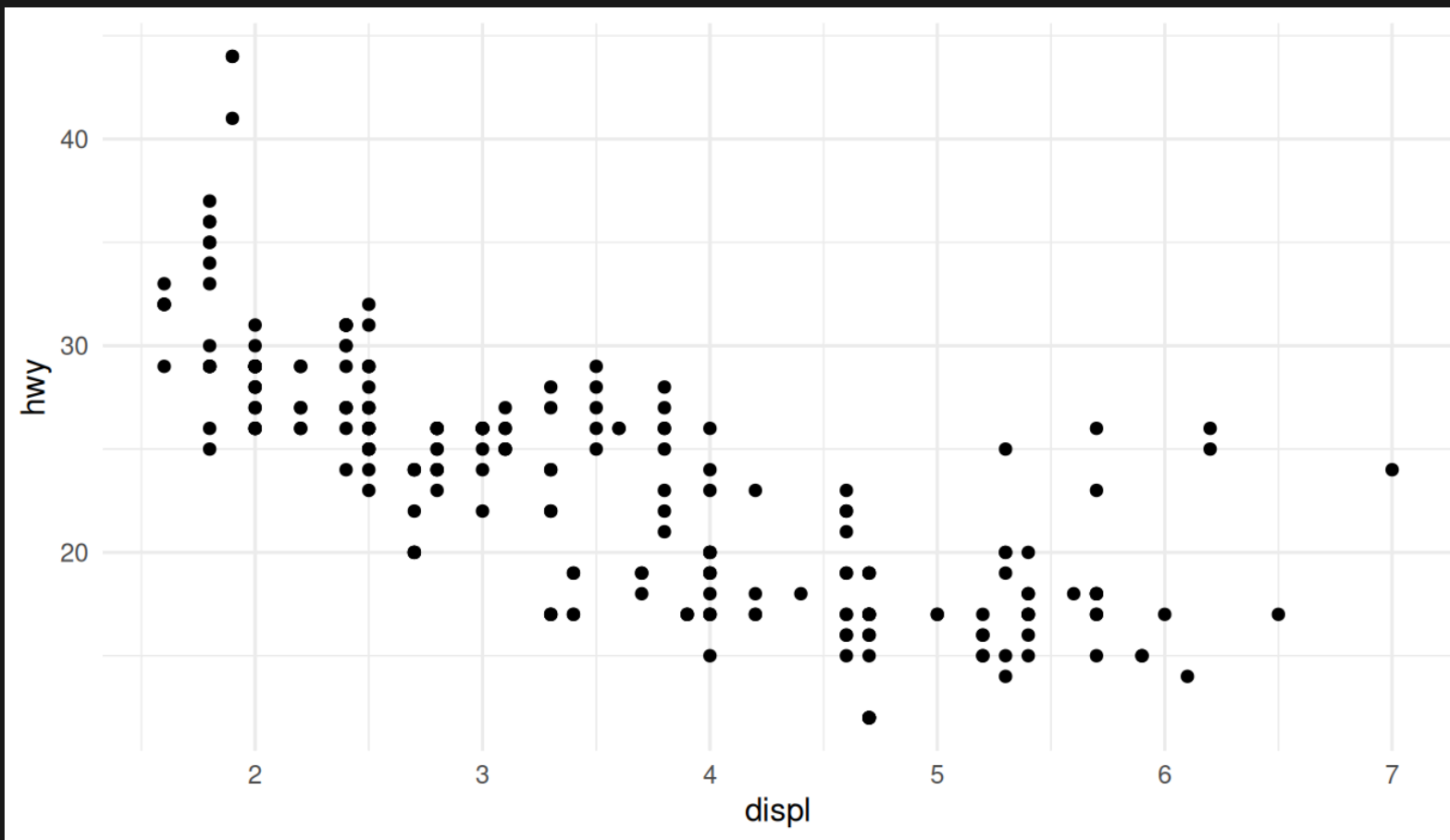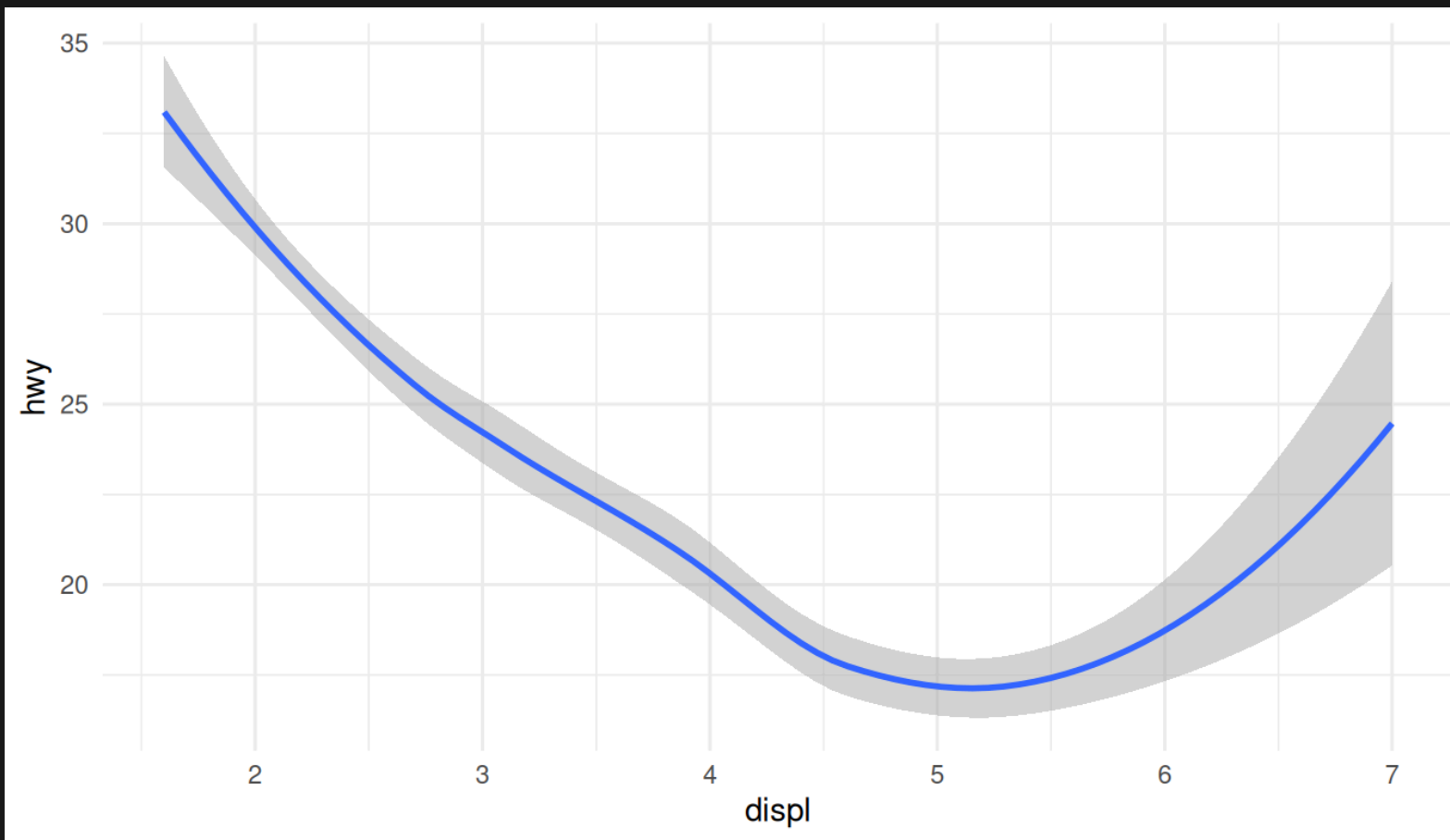


Still no real plot generated!

# Generating a plot

## But you also need an object. In ggplot, this is *geoms*

```
1  p + geom_point()
```

# Generating a plot, 2

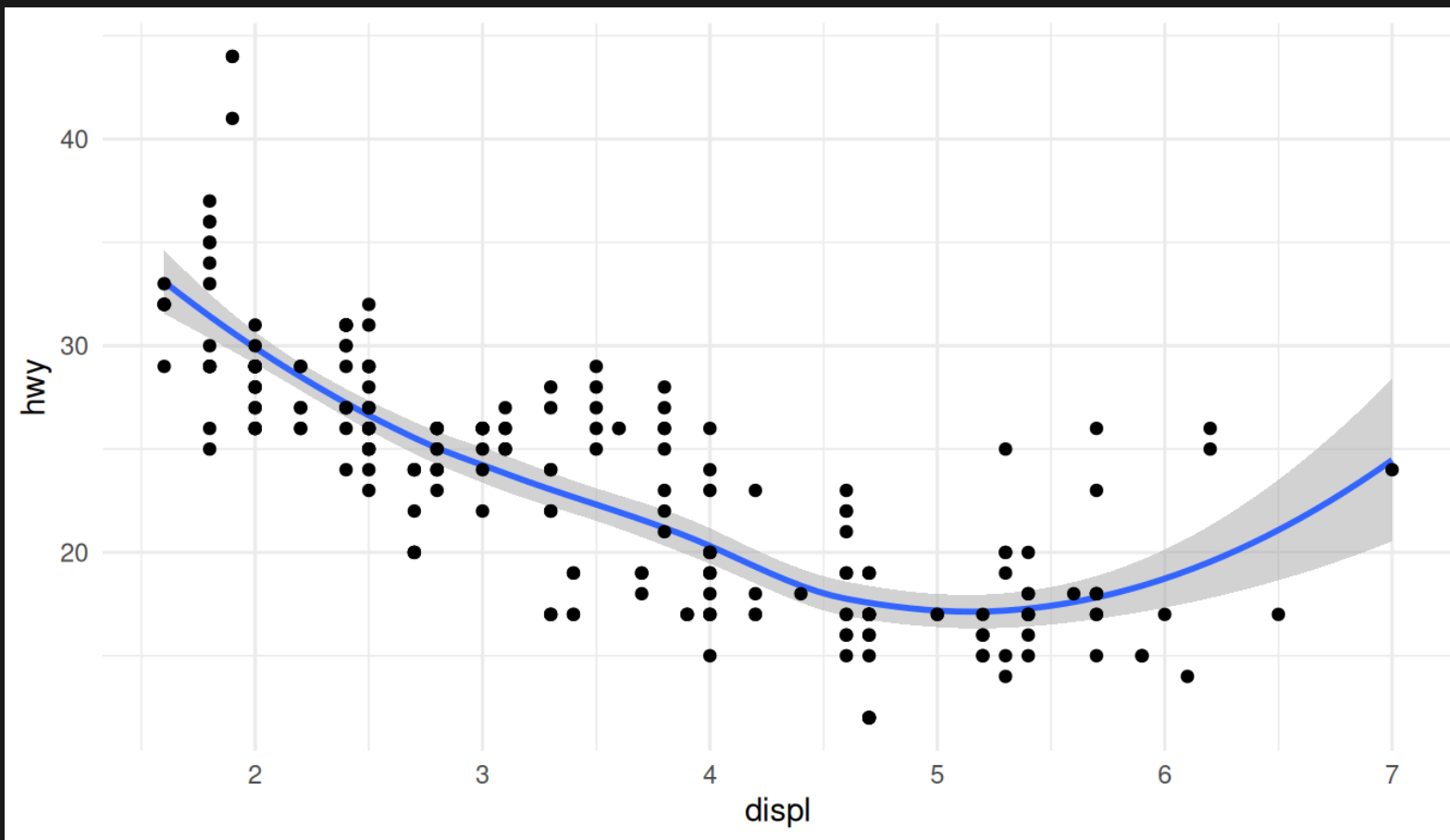## But you also need an object. In ggplot, this is *geoms*

```
1  p + geom_smooth()
```

# Generating a plot, 3

## You can add (+) as many *geoms* as you wish
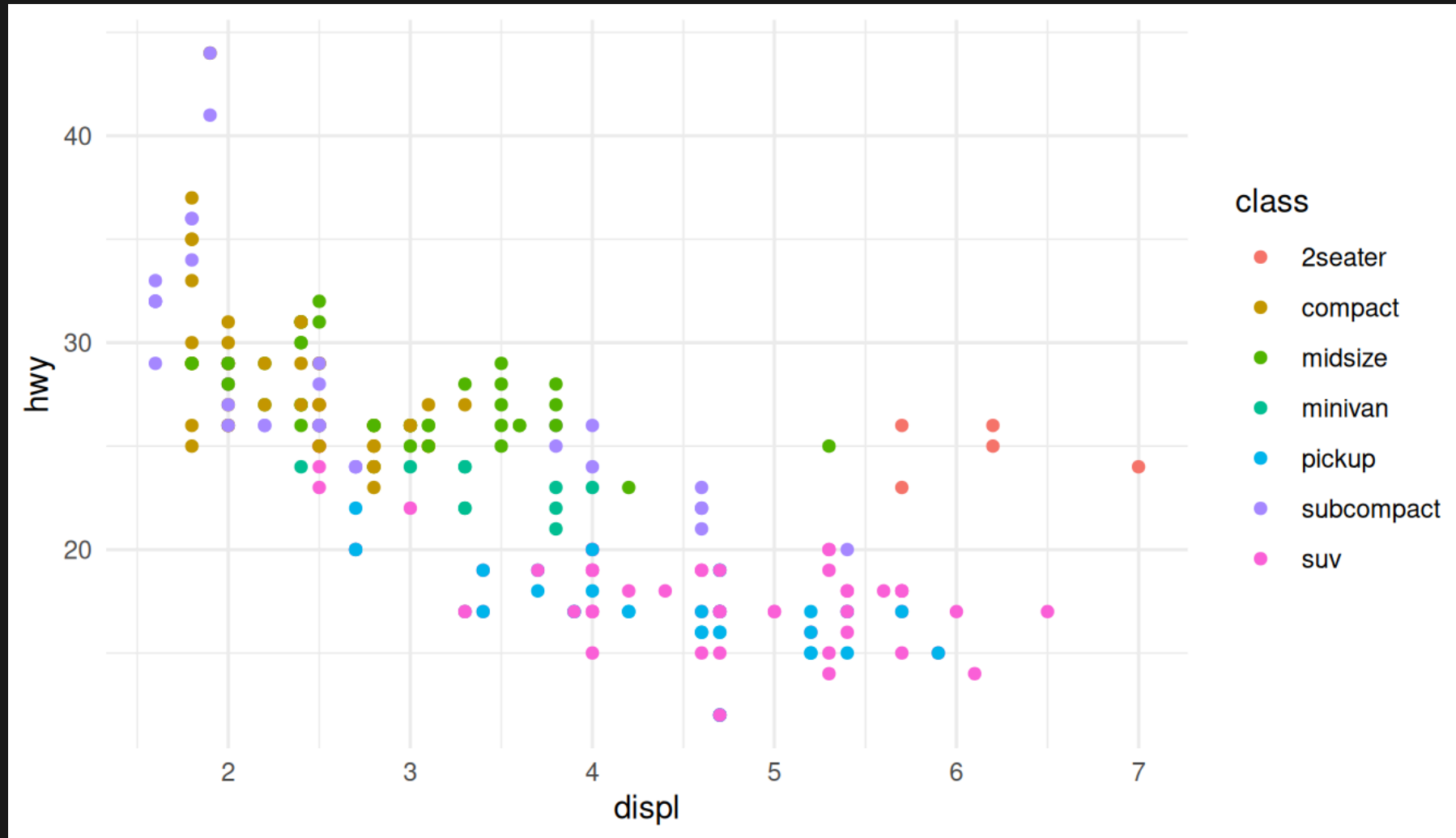
```
1  p + geom_smooth()+geom_point()
```

# The beauty of a grammar `metaphor`

- once you get the main idea, `adding` things is easy

- a plot is a **sentence made with data**

- you add layers with +

- as you would add words to a sentence

- as in grammar you use adjectives to give more nuanced meaning, in plots you could use + to add color, fill, size, shape, etc…
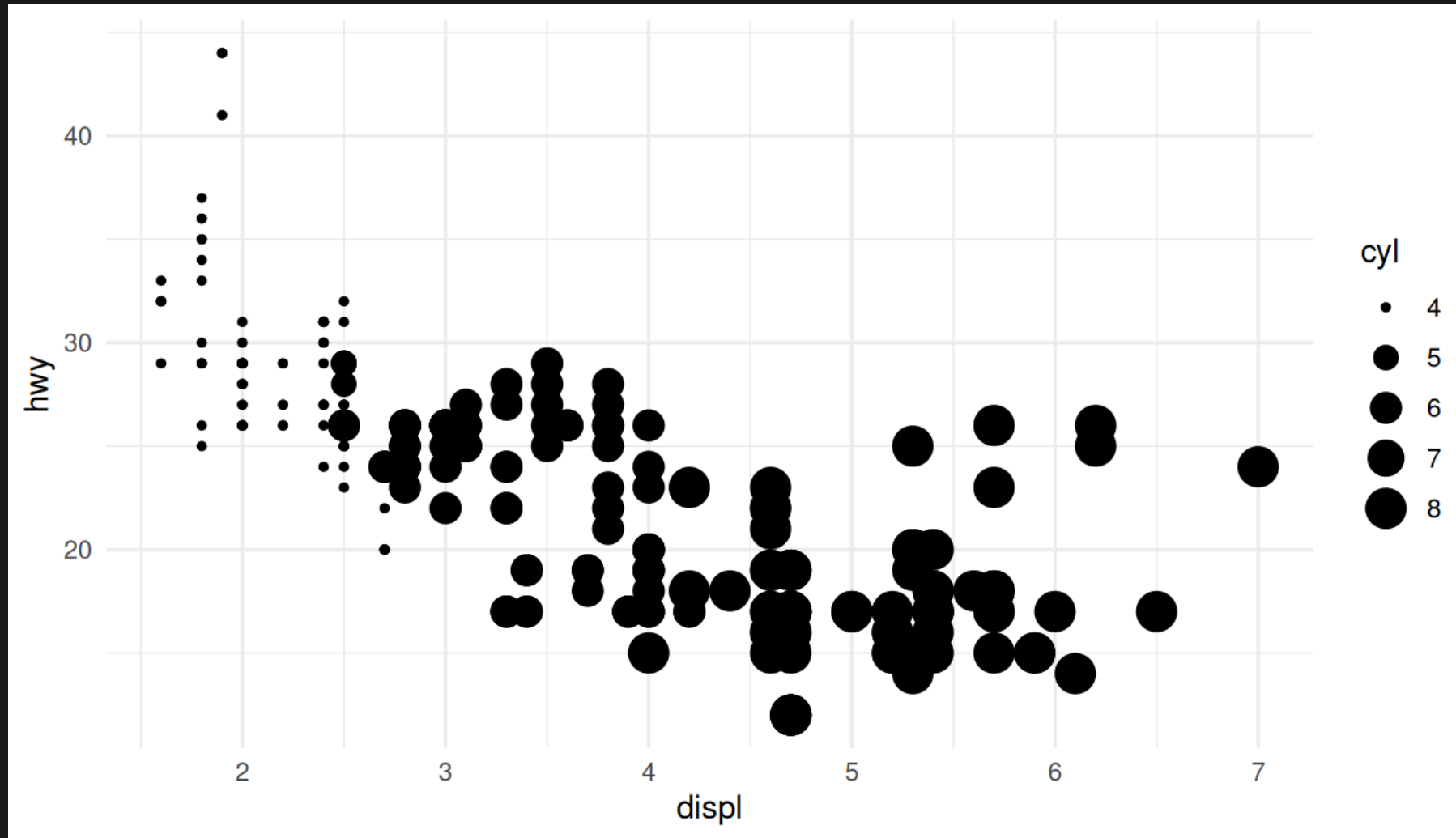
# Adding meaning: color
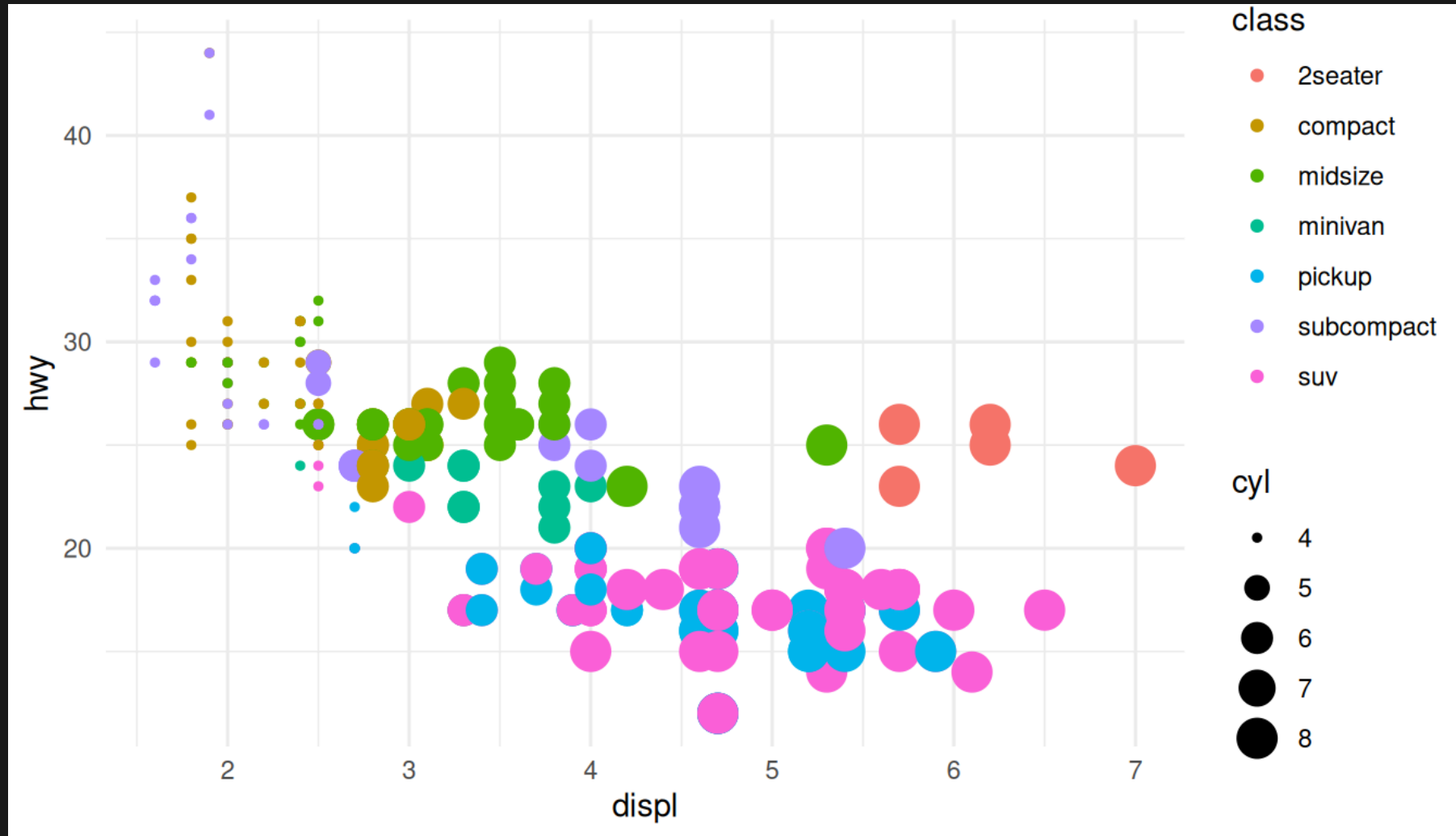
```
1 p + geom_point(aes(color=class))
```

# Adding meaning: size
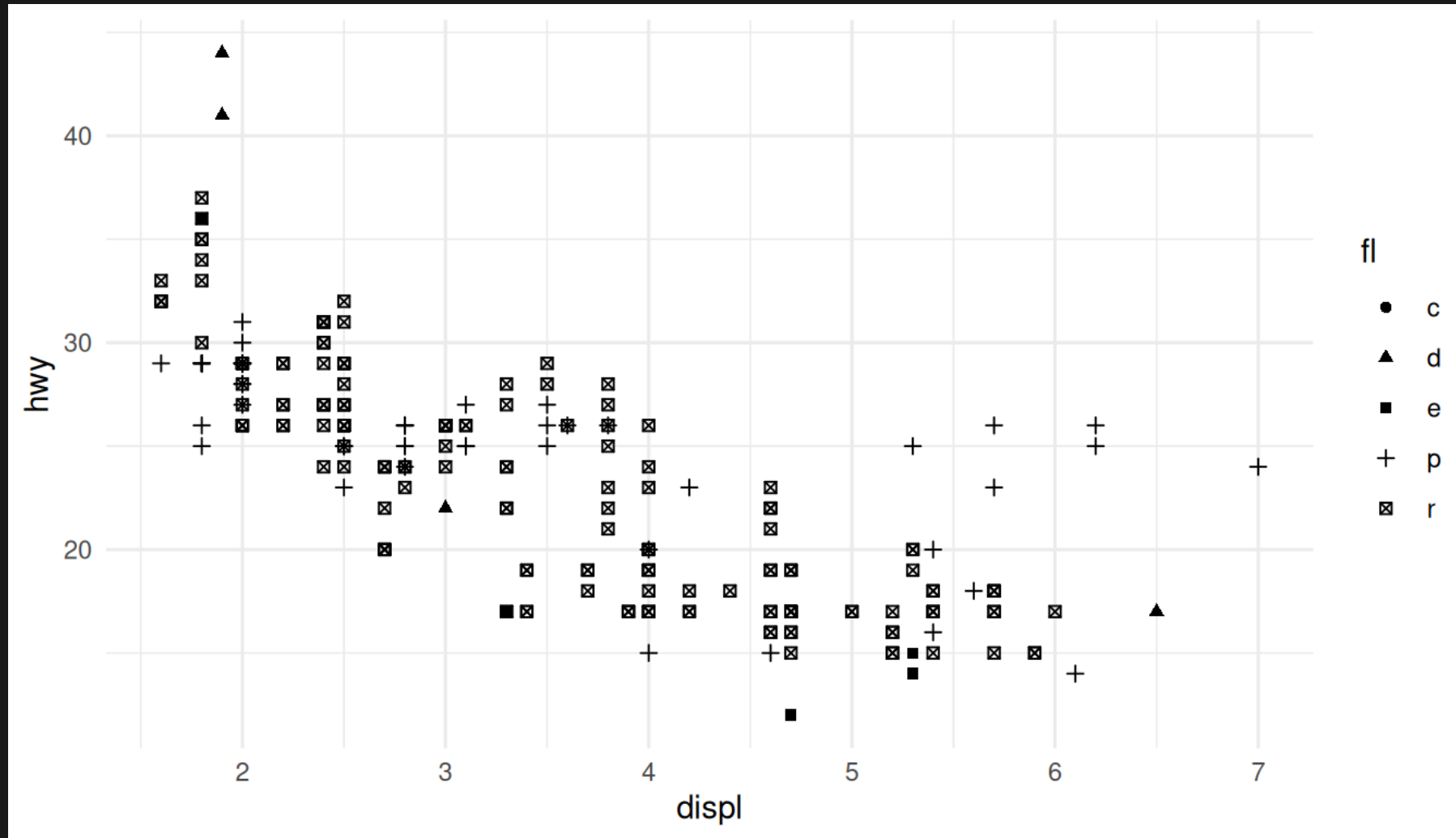
```
1  p + geom_point(aes(size=cyl))
```

# Adding meaning: color AND size

```
1  p + geom_point(aes(size = cyl, color=class))
```
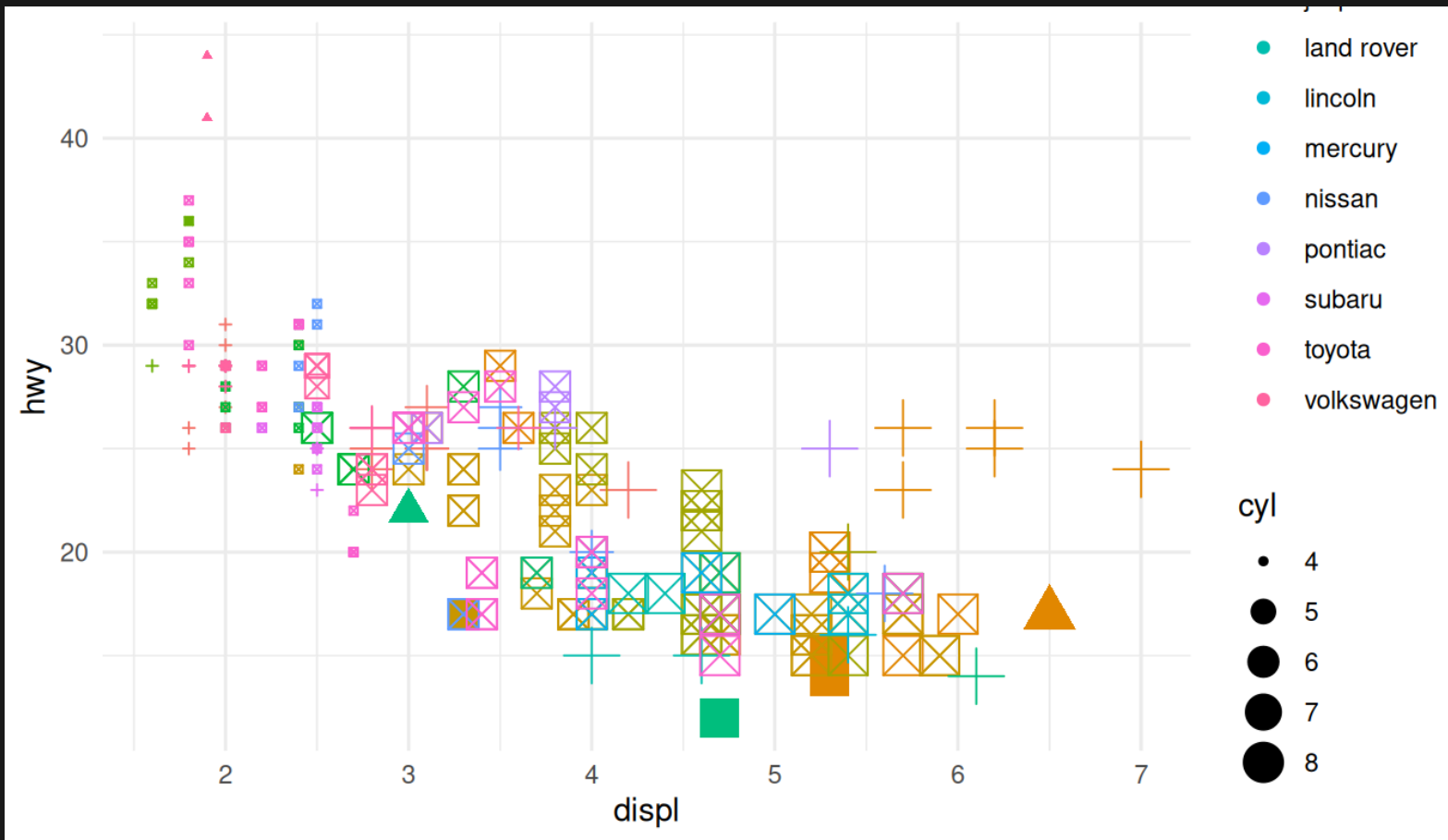
# Adding meaning: shape

```
1  p + geom_point(aes(shape=fl))
```

# Adding meaning: all together (...)

## Possibly not a good idea though

```
1  p + geom_point(aes(color=manufacturer, shape =fl, size = cyl))
```
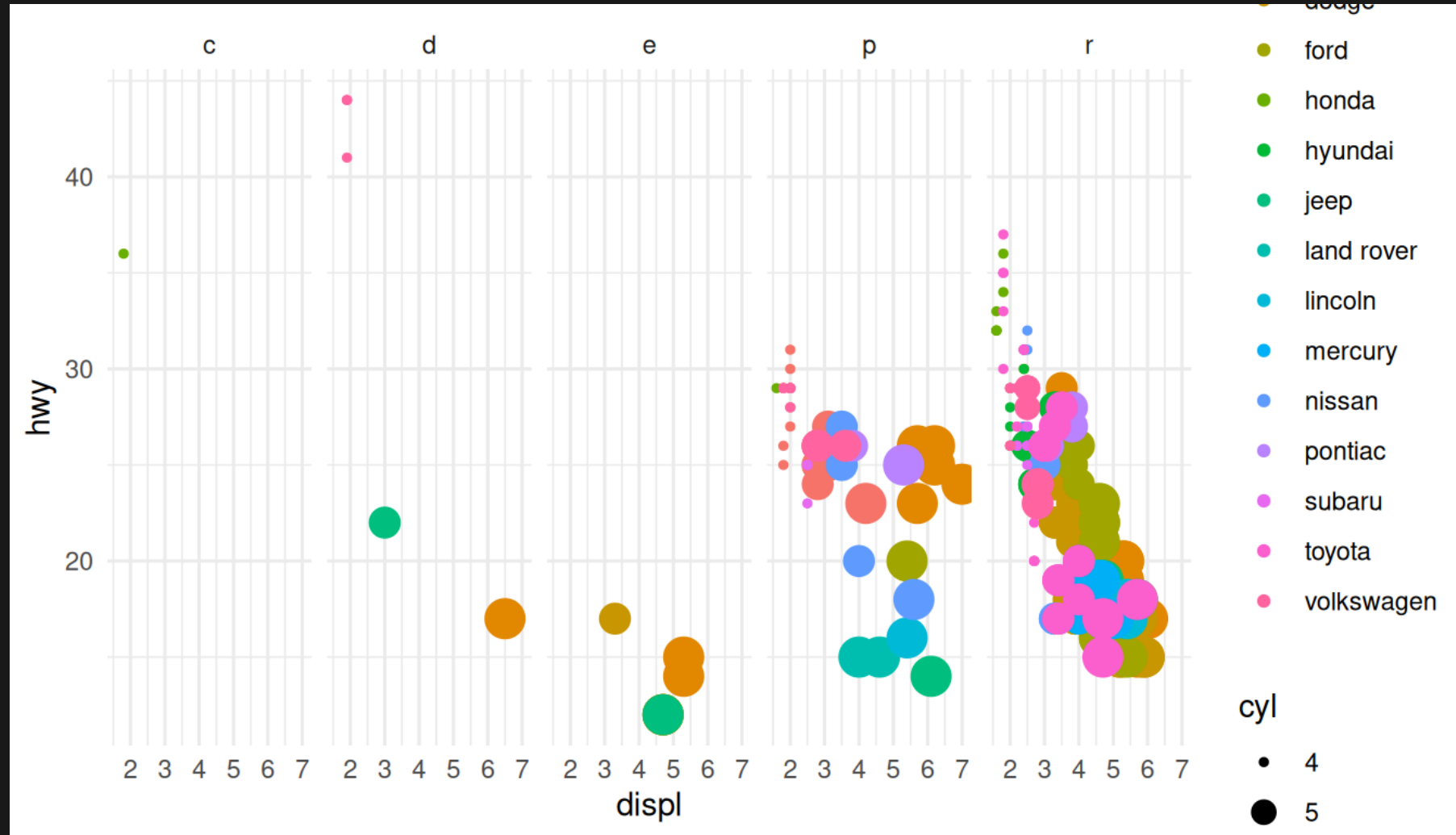
# Recap so far

- ggplot works like a grammar

- start with `ggplot()`

- first argument: *data*: `ggplot(df, ...)`

- then map variables to `aes`thetics (x, y, color, fill, …)

- `ggplot(df, aes(dimension = variable))`

- then add meaning with geometric objects: `geom_*`

- notes:

    - `geom`s inherit the `aes` of the plot if not specified

    - all variables mapped to `aes` vary with the data

# Facets

- sometimes sentences become too long

- it is useful to `split` them up in shorter sentences

- you could first talk about a car, *then* another one

- in plots, you can split up the plot along a `variable`

- one plot is drawn `for each level` of a given variable

# Facets

```
1  p + geom_point(aes(color=manufacturer, size = cyl))+facet_grid(.~fl)
```

# More details on the grammar

Once your main plot is done, you can tweak it

- coordinate functions (changing the axis)

- scale functions (changing how geoms look)

- theme functions (changing how the plot looks)

**We will do this in Lecture 5 – advanced plotting**
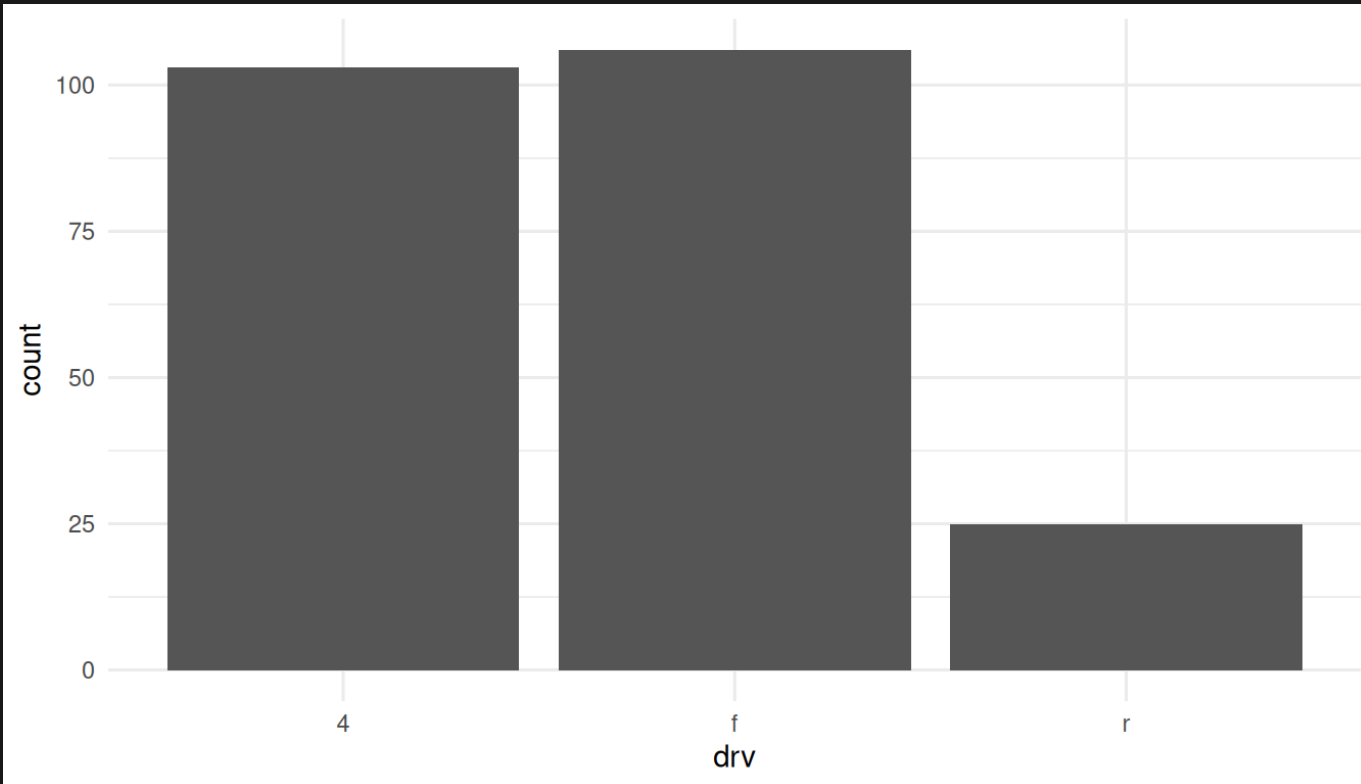
# ggplot2: gallery

# Exploring data: one variable

*Plot types depend on the variable type*

- *one-variable plots, discrete variable*: `barplot`

- *one-variable plots, continuous variable*: `distribution`, `density`

# Barplots

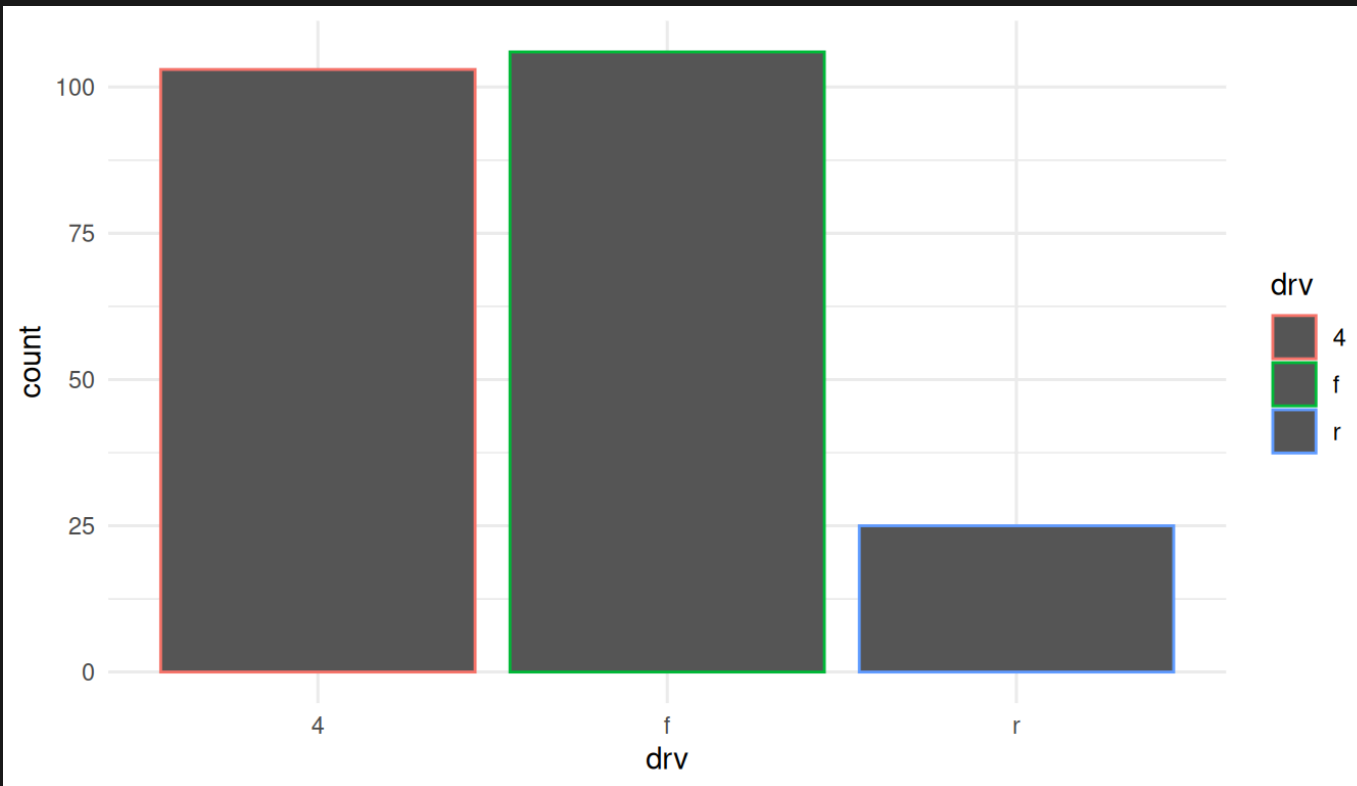- let's look at the drive type of the cars: front, rear, or 4wd

```
1  p <- ggplot(mpg, aes(drv))
2  p + geom_bar()
```
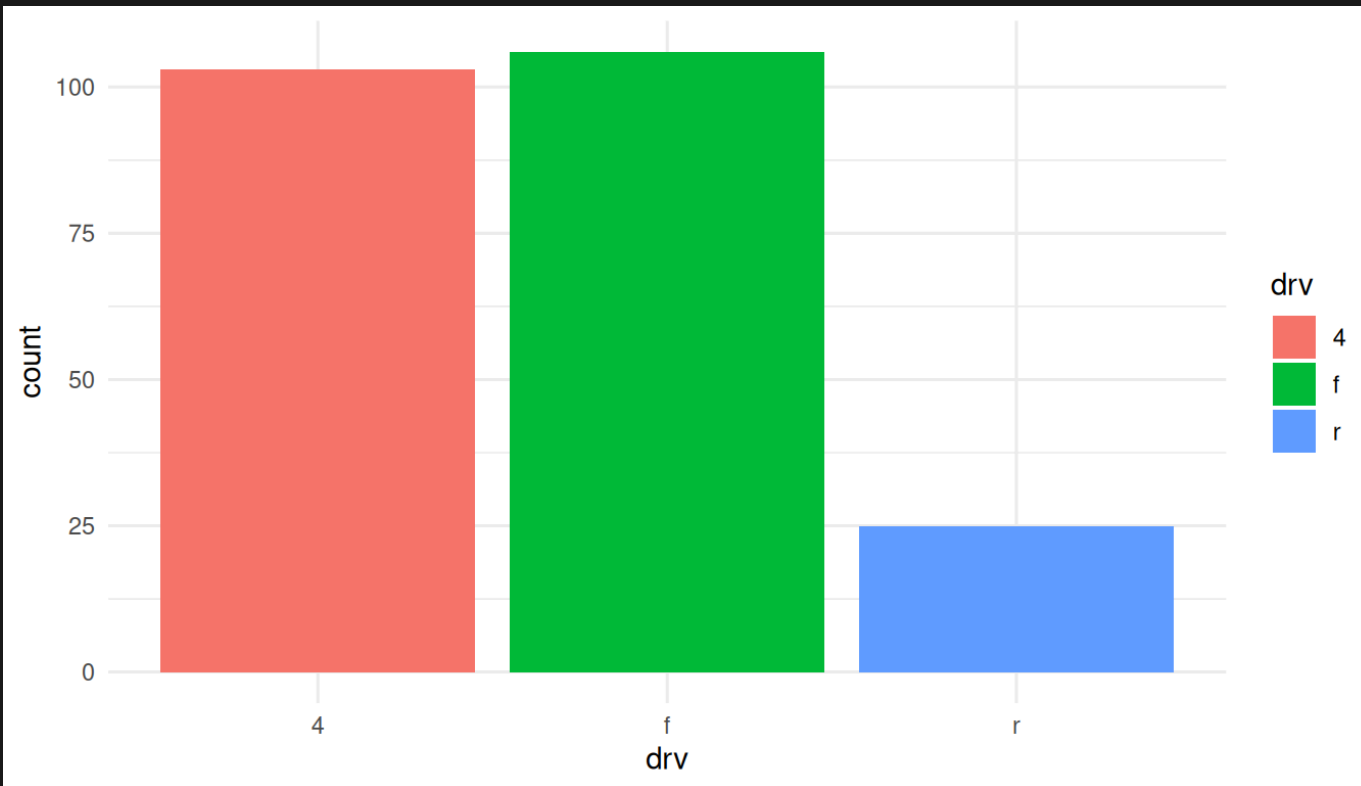
# Barplots

- not so fancy. should we add color?

```r
1  p <- ggplot(mpg, aes(drv))
2  p + geom_bar(aes(color=drv))
```
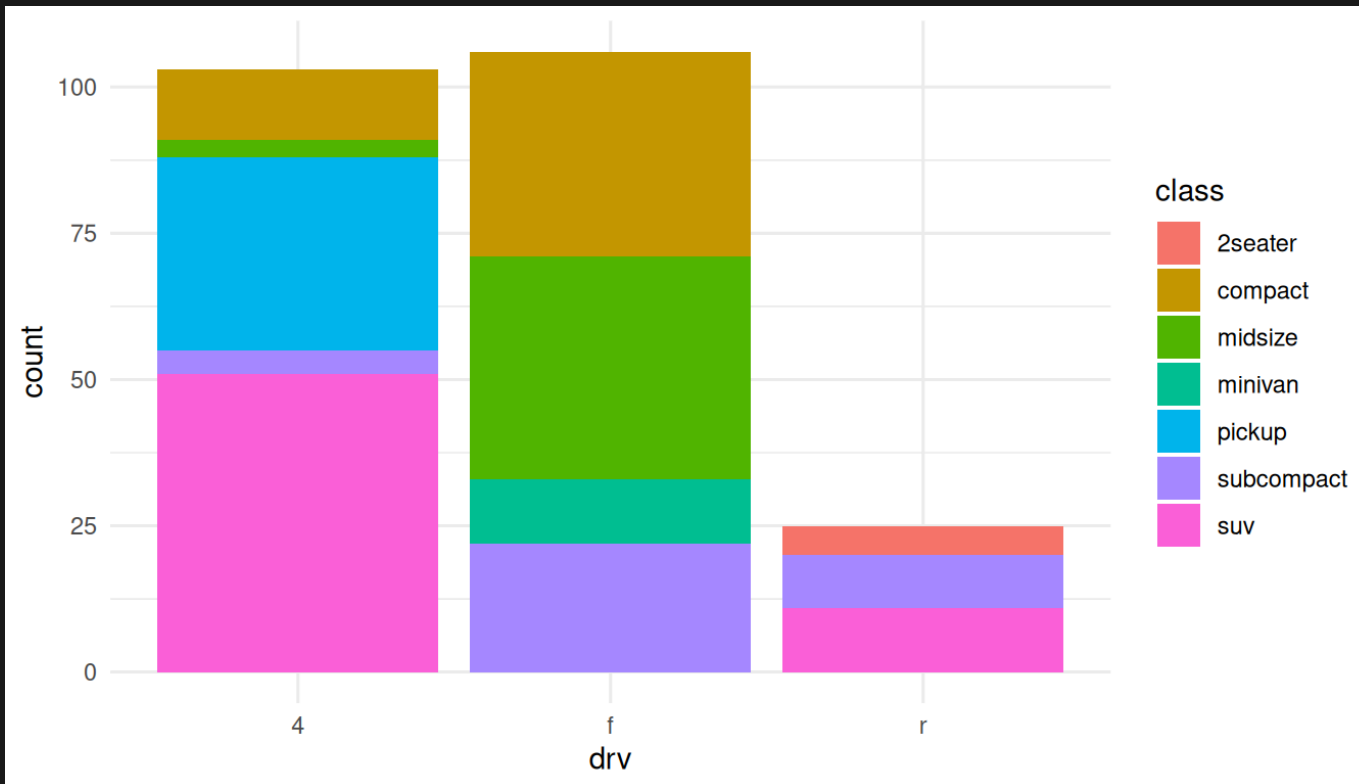
# Barplots

- ups. Maybe we meant fill?

```
1  p <- ggplot(mpg, aes(drv))
2  p + geom_bar(aes(fill=drv))
```

# Barplots

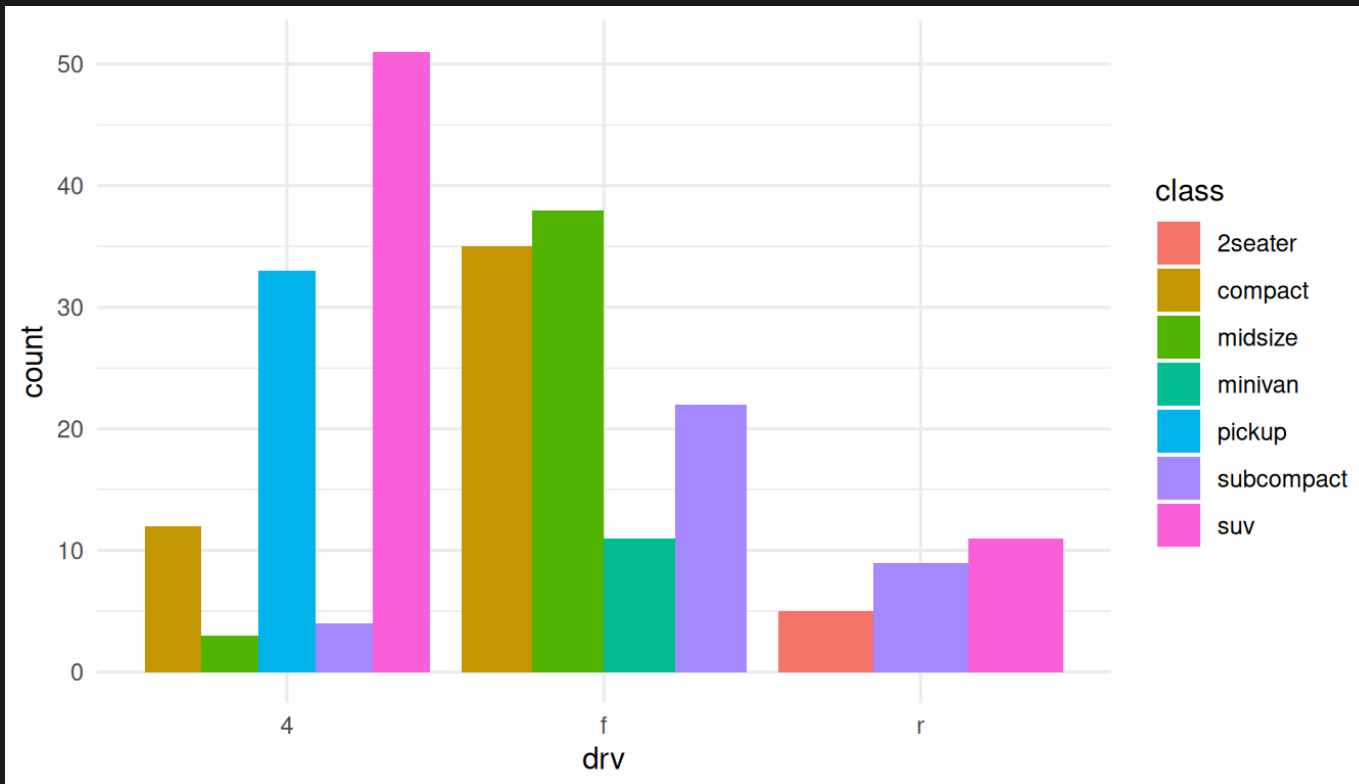- what if we cross it with another variable?

```r
1 p <- ggplot(mpg, aes(drv))
2 p + geom_bar(aes(fill=class))
```

# Barplots

- By default stacked. How to unstack?

```
1  p <- ggplot(mpg, aes(drv))
2  p + geom_bar(aes(fill=class), position = position_dodge())
```
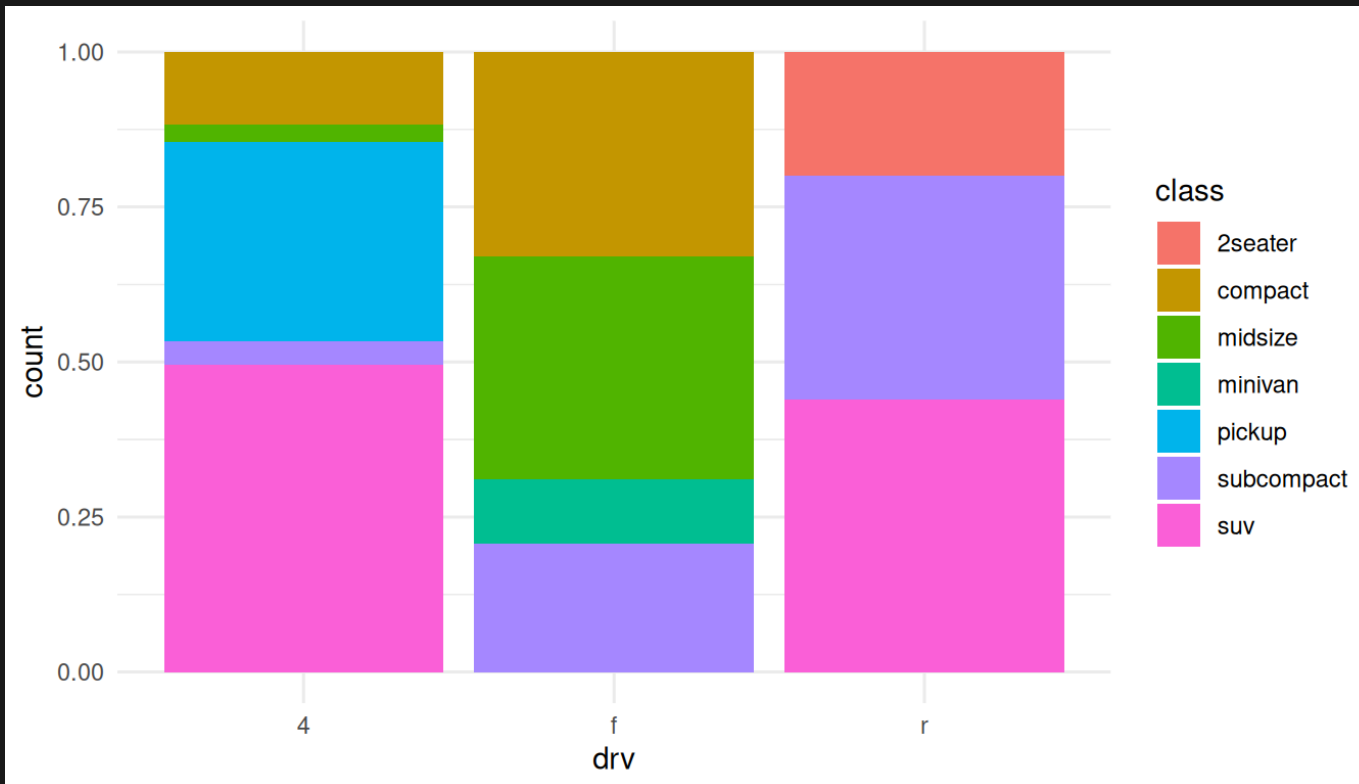
# Barplots

- By default stacked. How to show relative weight?

```
1  p <- ggplot(mpg, aes(drv))
2  p + geom_bar(aes(fill=class), position = position_fill())
```
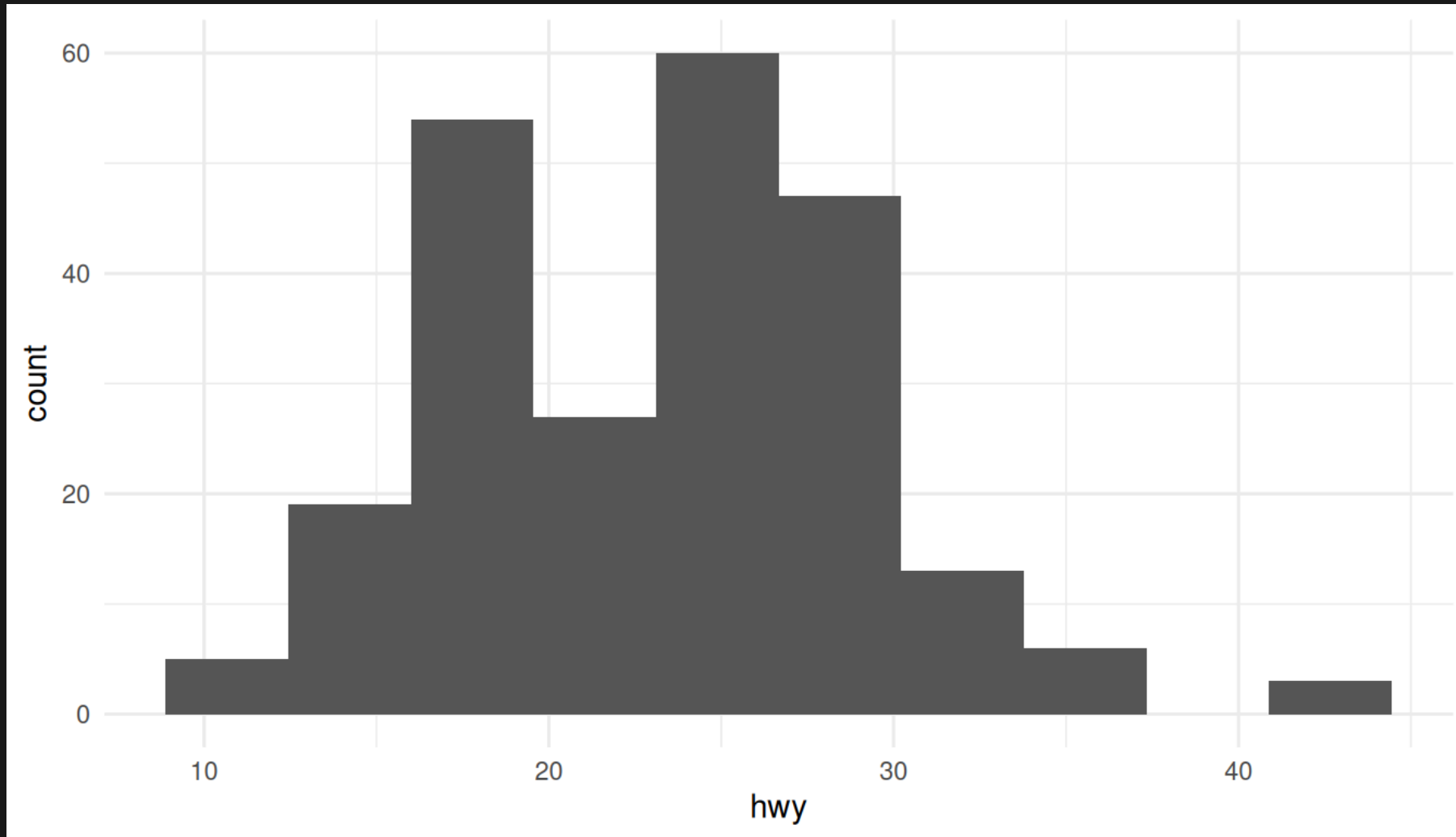
# One variable, continuous

- If var is continuous, it makes more sense to show distributions

```
1  p <- ggplot(mpg, aes(hwy))
2  p + geom_histogram()
```
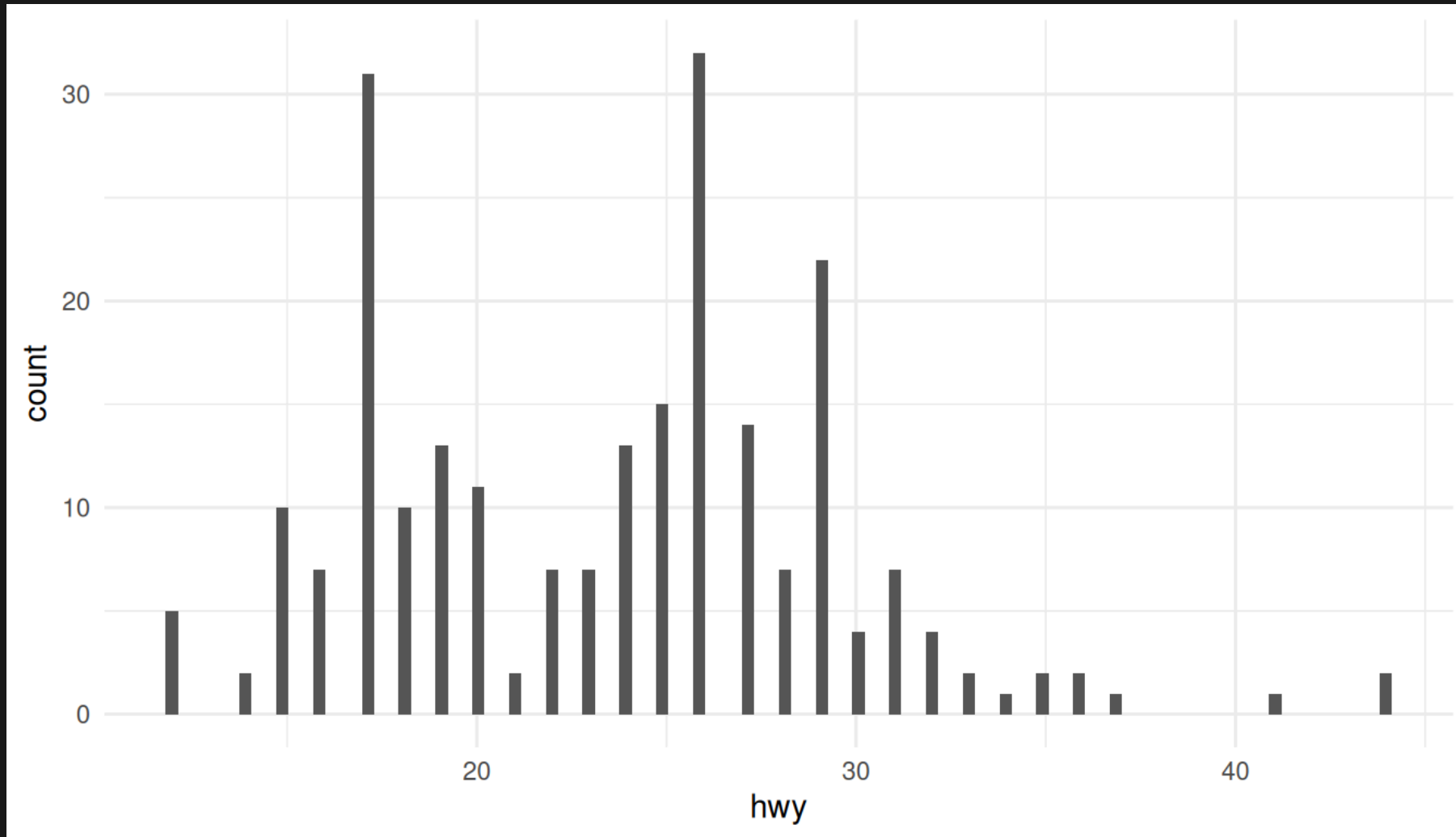
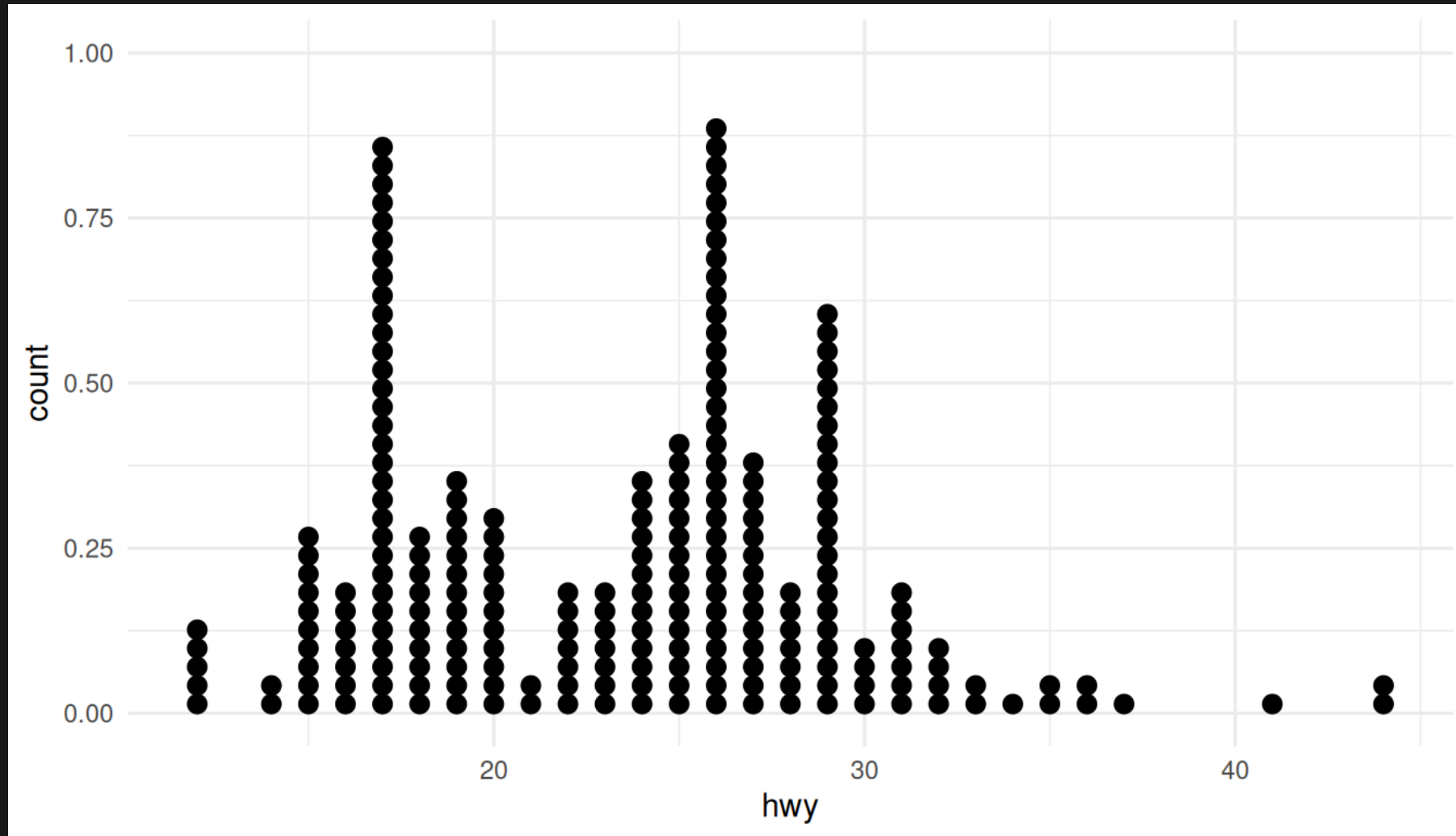# Histograms: binwidth

```
1  p + geom_histogram(bins = 10)
```

# Histograms: binwidth

```
1 p + geom_histogram(bins = 100)
```

# An alternative to histogram: dotplot
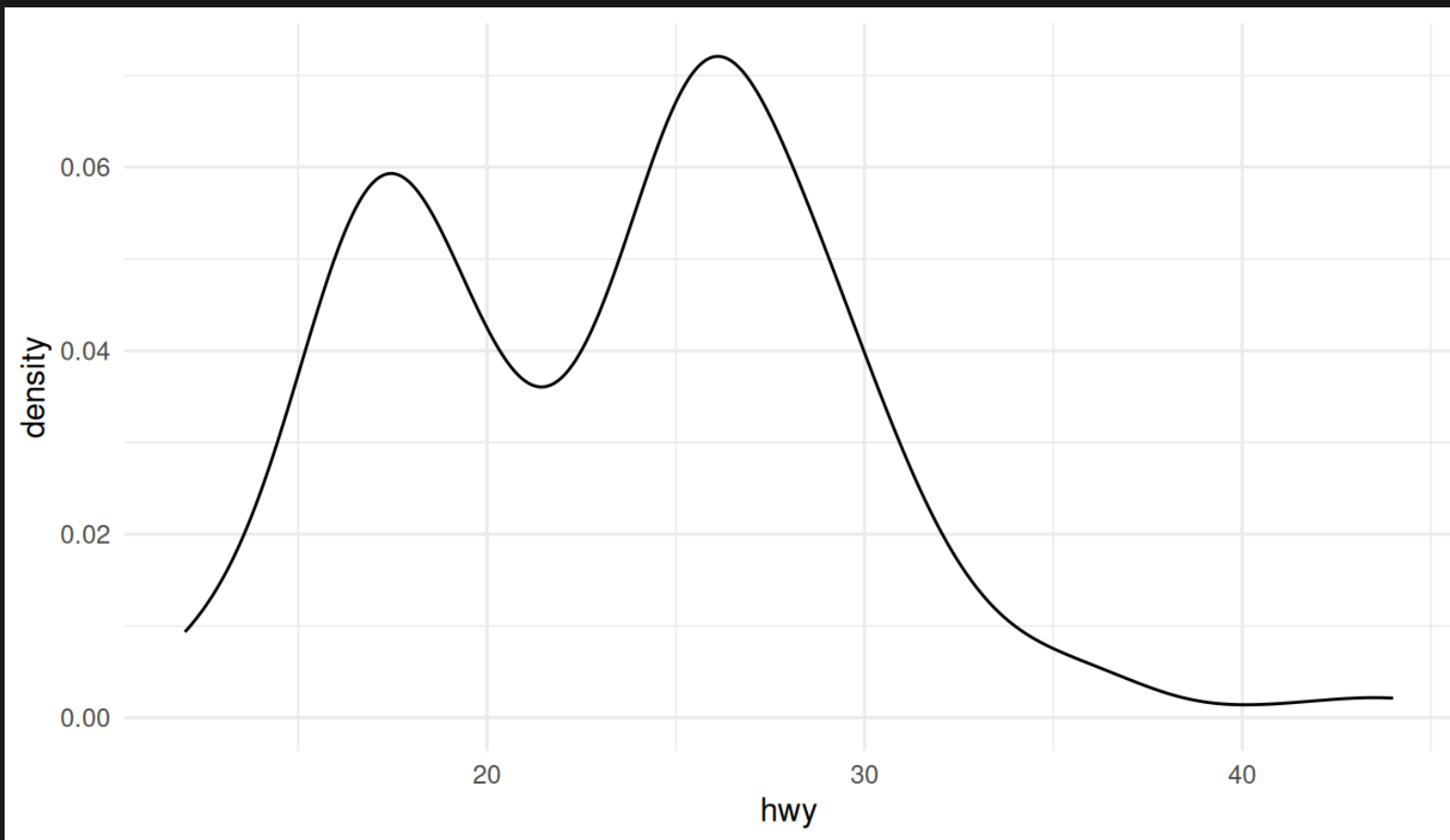
```
1  p + geom_dotplot(binwidth = 0.5)
```

# Continuous distributions
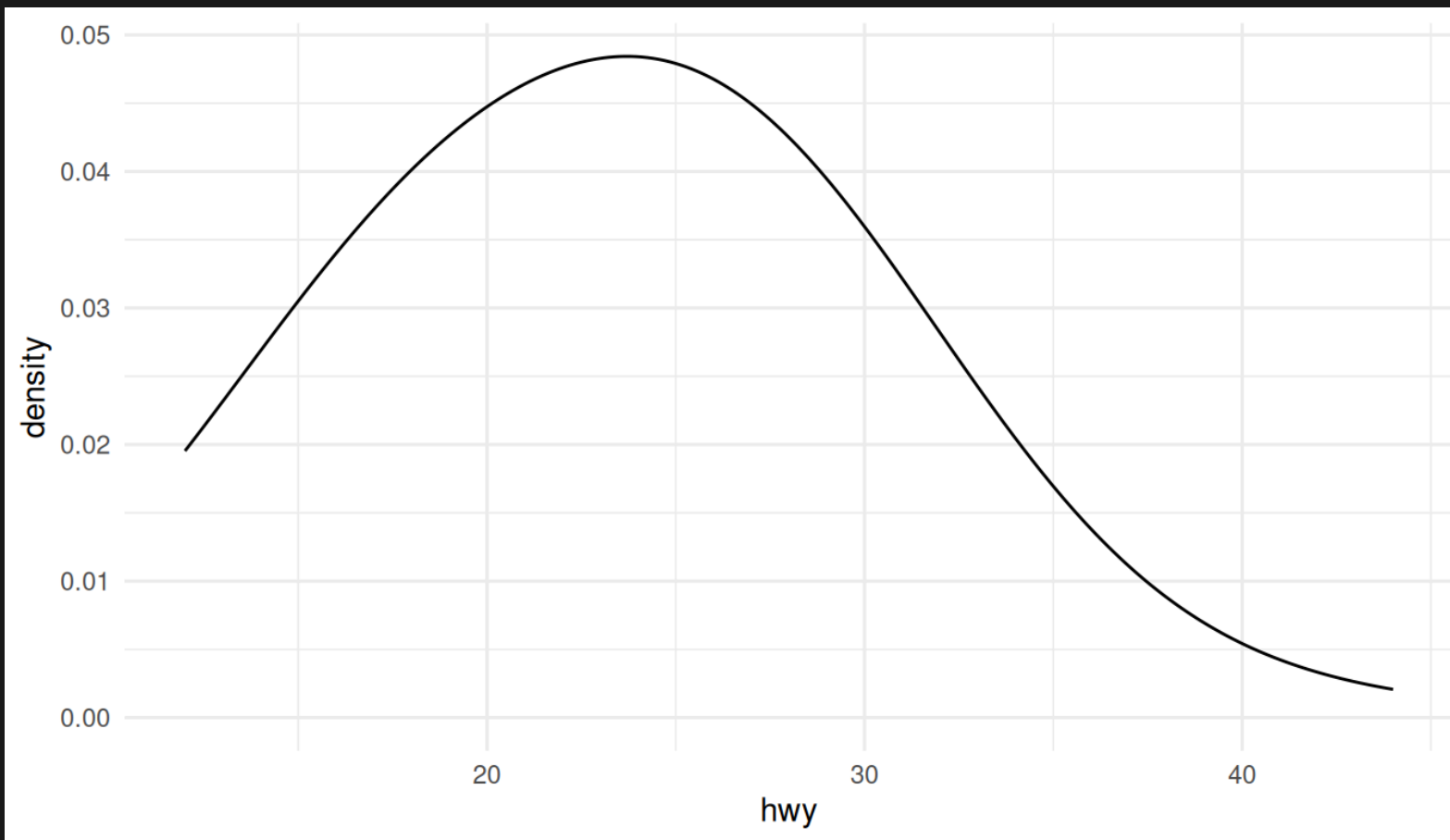
In this case use kernel density estimation

```
1  p + geom_density()
```

# Continuous distribution
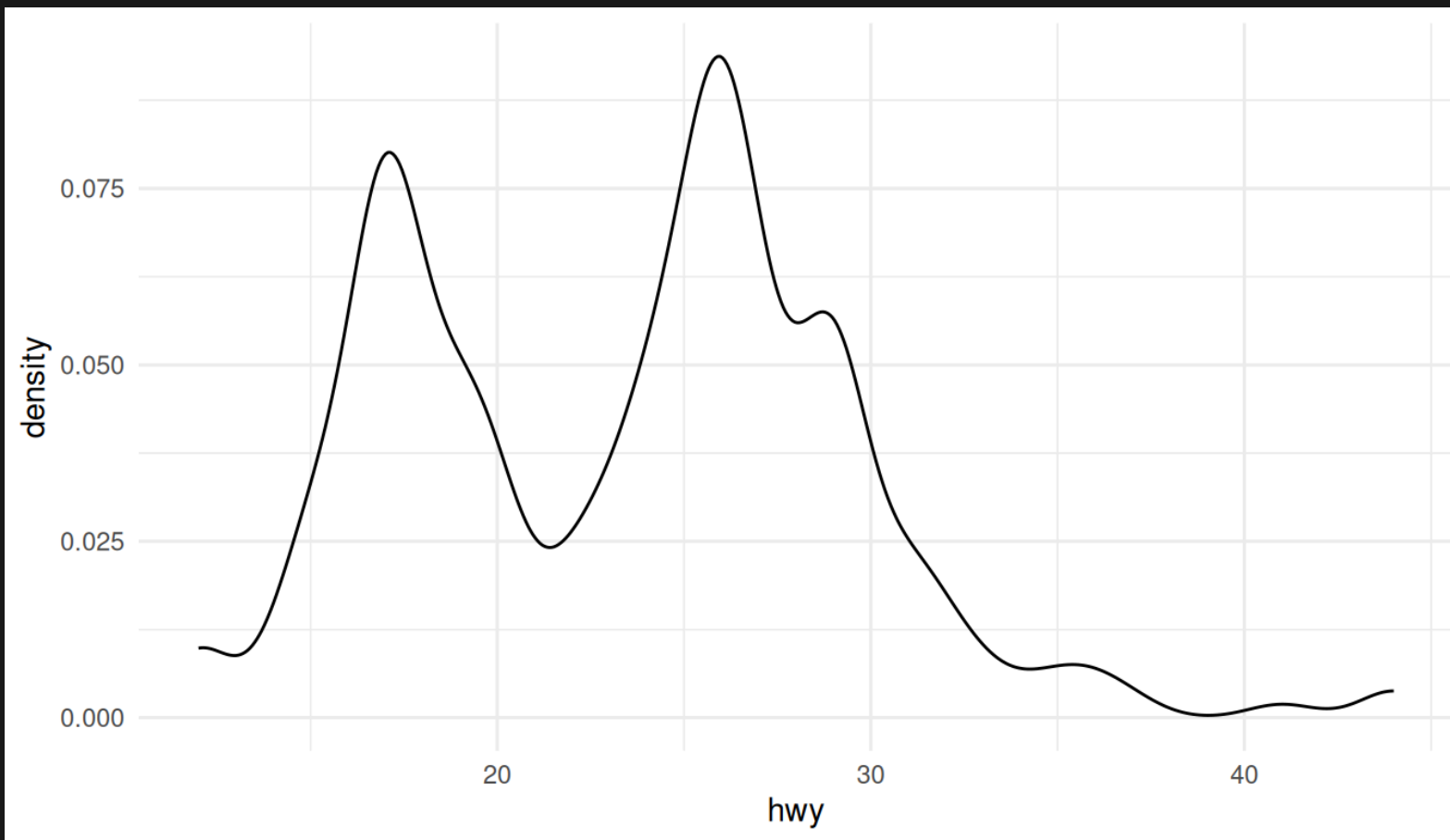
In this case use kernel density estimation

```
1  p + geom_density(adjust = 3)
```

# Continuous distribution

In this case use kernel density estimation

```
1  p + geom_density(adjust = 0.5)
```
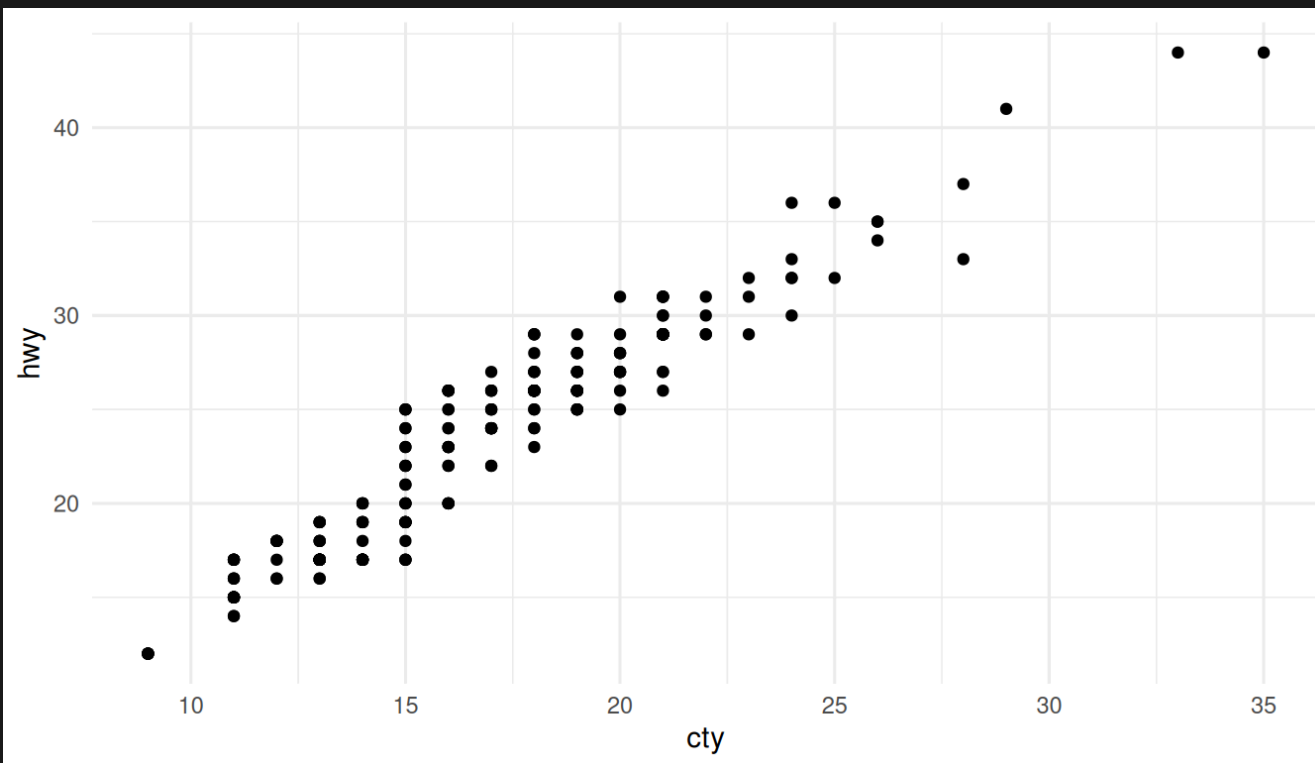
# Exploring data: **two** variables

*Plot types depend on the variable type*

- *both vars continuous*: scatter, smooth

- *one continuous, one discrete*: columns, boxplot, violins

- *both discrete*: count

# Scatter

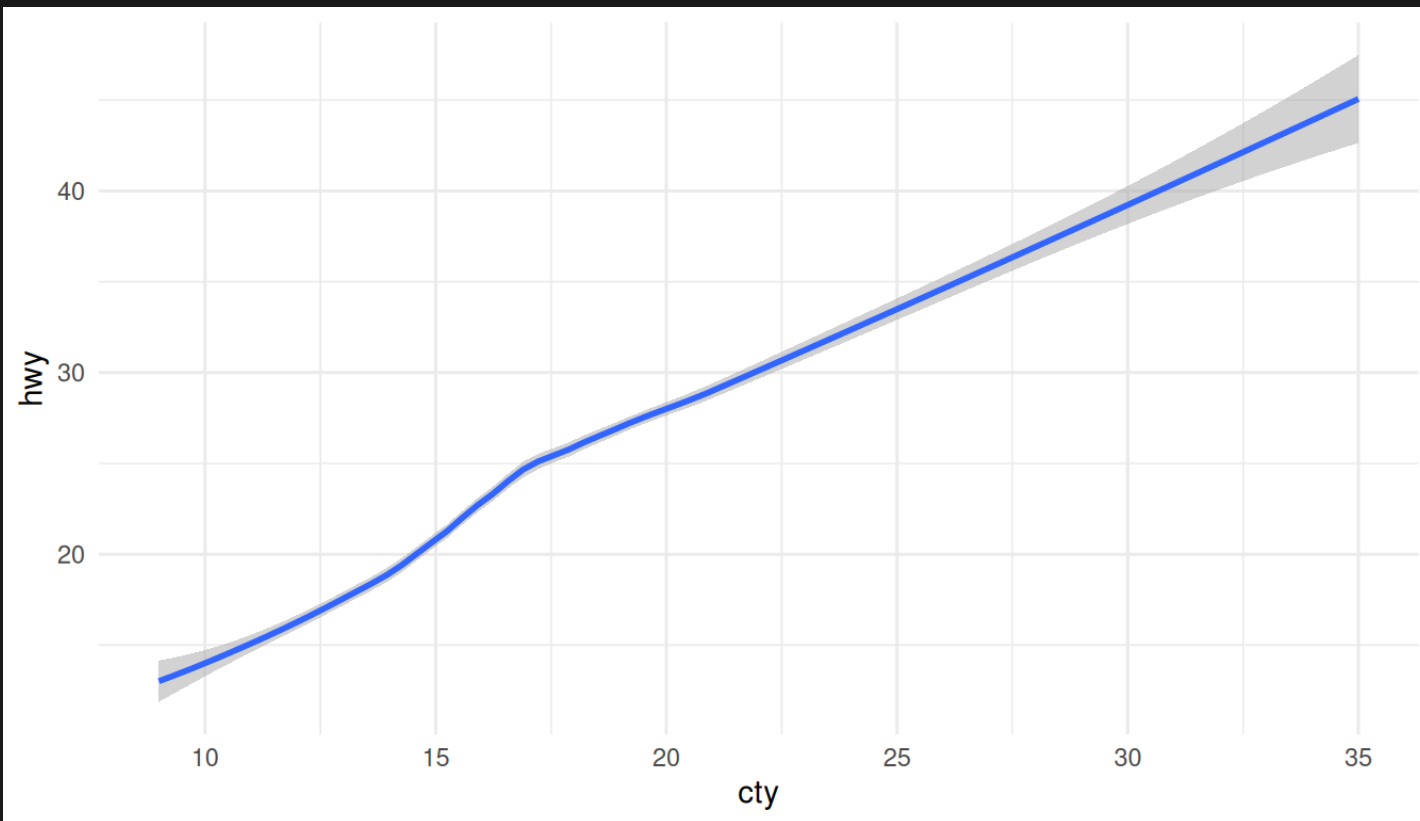if two variables are continuous, your choice is `scatter`

```r
1  p <- ggplot(mpg, aes(x = cty, y = hwy))
2  p + geom_point()
```

# Smooth

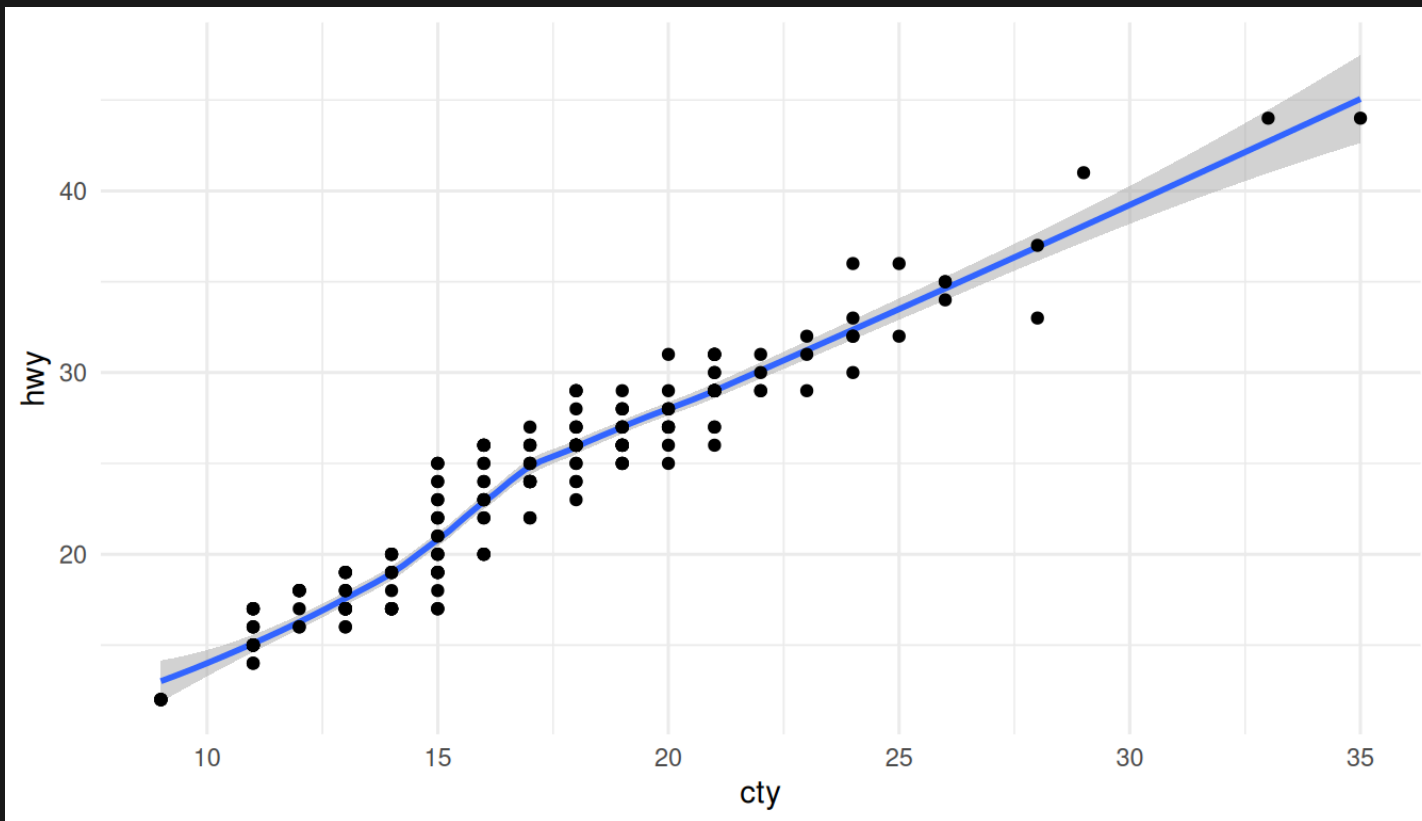still, you might just want to show the general tendency

```
1  p + geom_smooth()
```

# Scatter + smooth

> or both

```
1  p + geom_smooth() + geom_point()
```
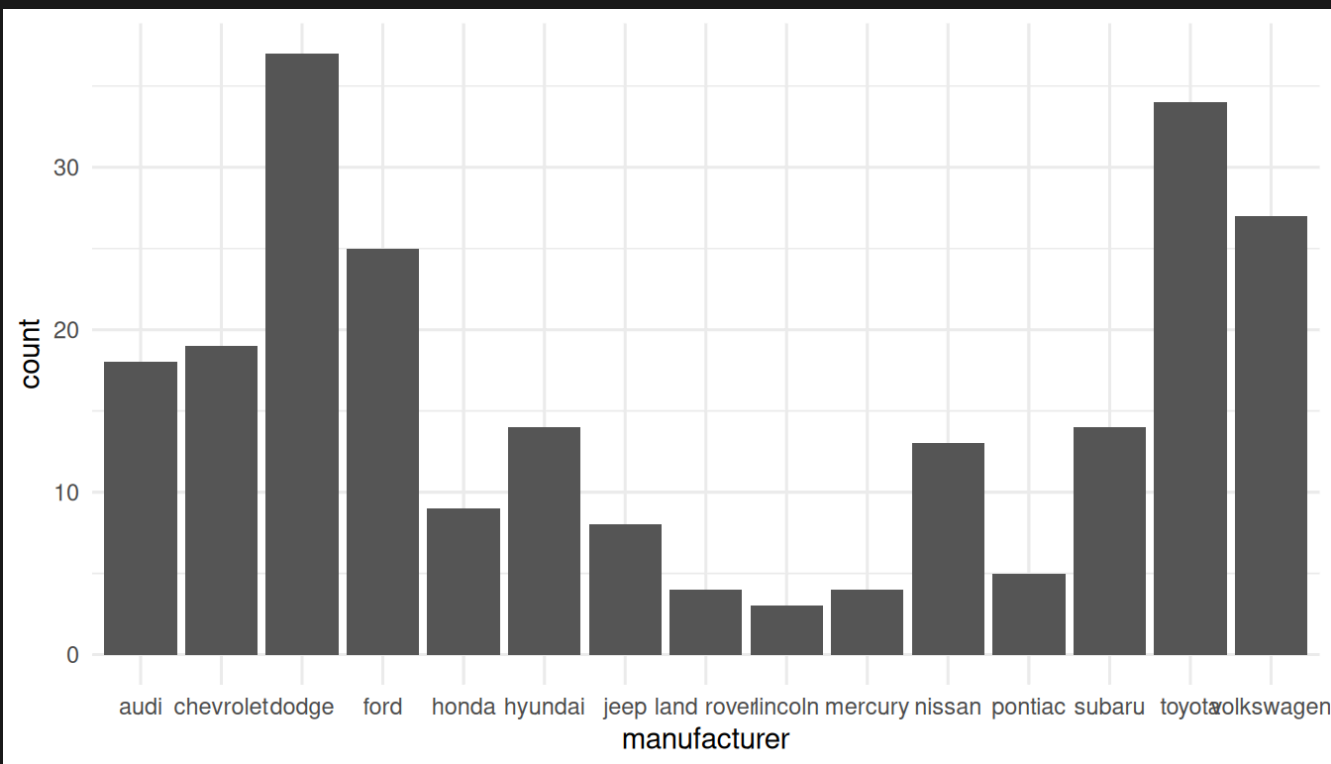
# Columns: a special type of bars

one variable discrete, one continuous (needs
`summarise()`!)

```
1  mpg %>% group_by(manufacturer) %>% summarise(n = n()) %>%
2  ggplot(aes(manufacturer, n))+
3    geom_col()
```

# Columns: why bother?

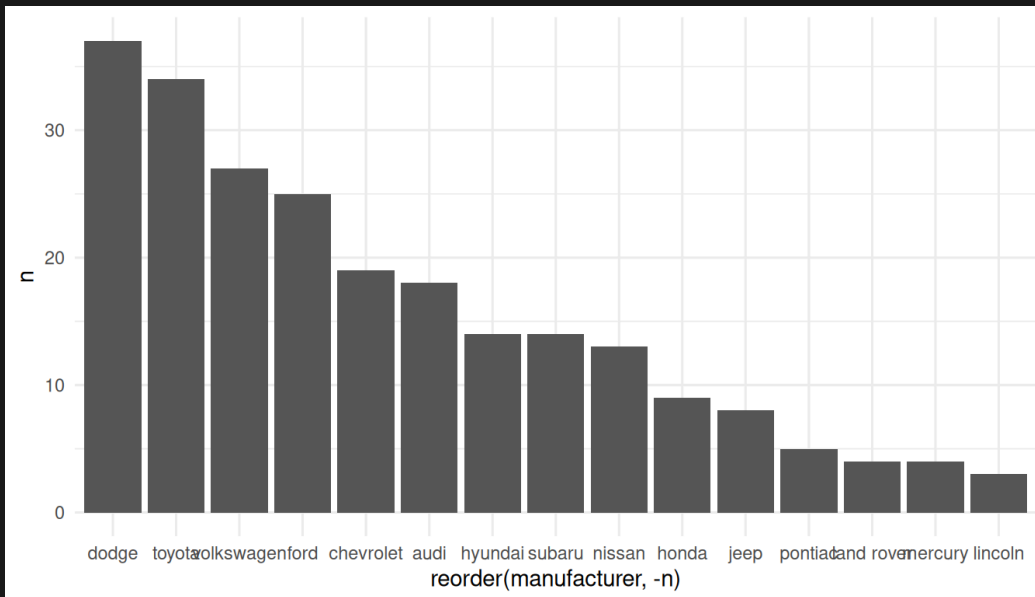we could have used `geom_bar` (that counts for us)

```
1  mpg %>% ggplot(aes(manufacturer))+
2    geom_bar()
```

# Columns: a special type of bars

but `geom_col` gives more options, since now you condition on a proper variable (n). For instance: order by n
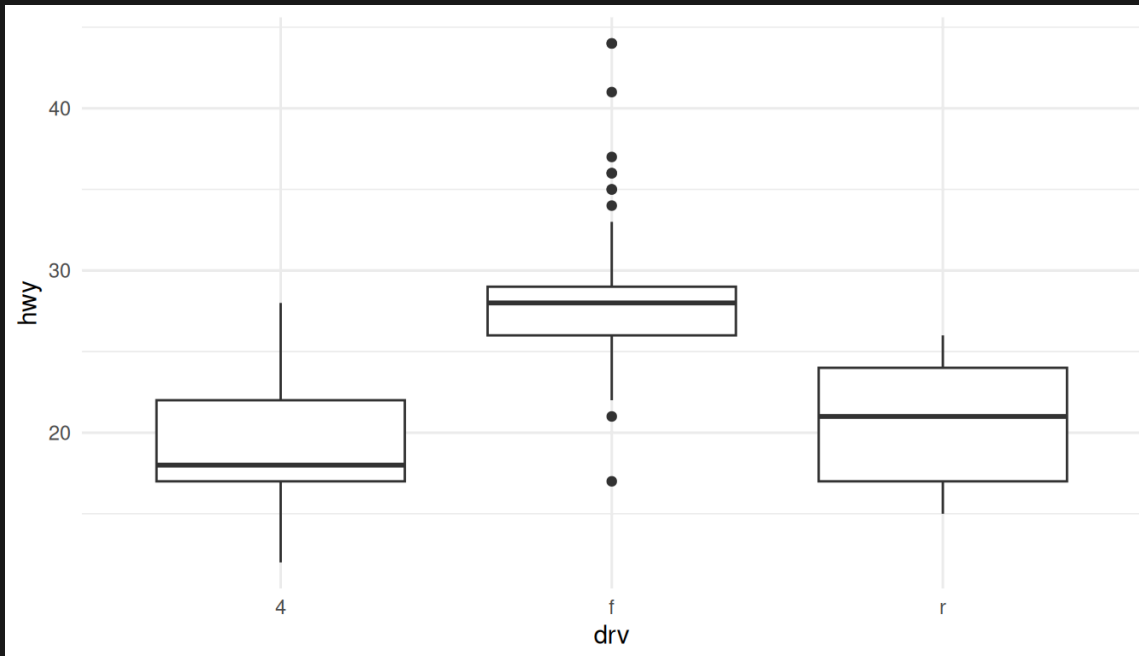
```
1  mpg %>% group_by(manufacturer) %>% summarise(n = n()) %>%
2  ggplot(aes(reorder(manufacturer, -n), n))+
3    geom_col()
```

# Boxplots

boxplots show a distribution but can do so over different levels of a categorical var

```
1  mpg %>% ggplot(aes(drv, hwy))+
2      geom_boxplot()
```
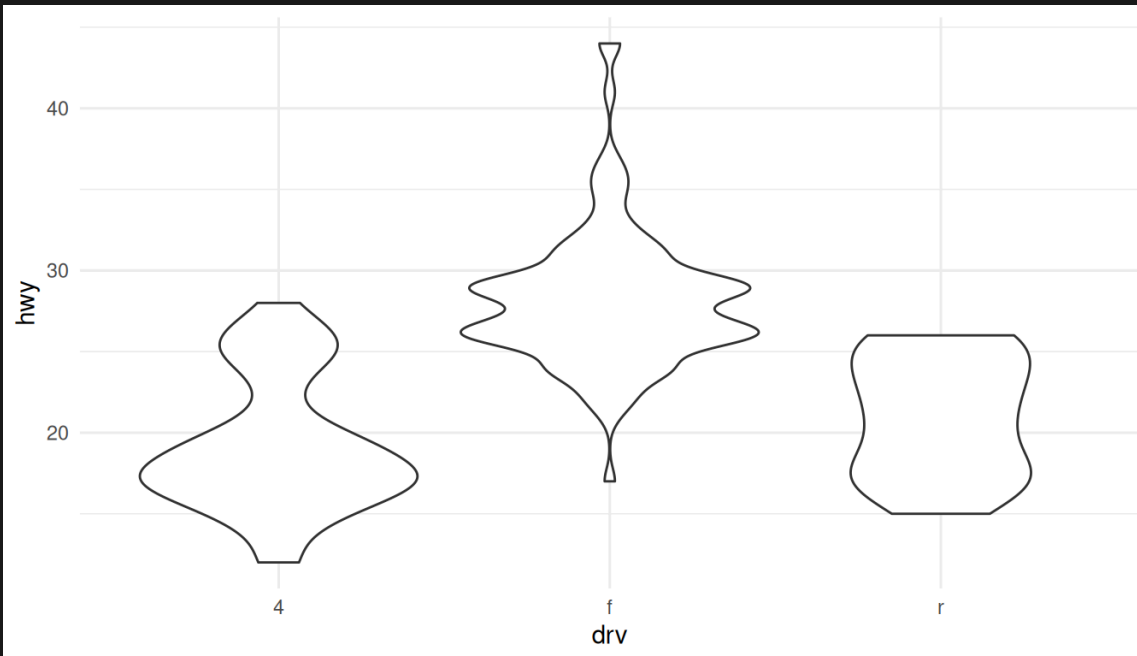
# An alternative to boxplot: violin

boxplots are bulky and only show relevant info. Want full distribution? Use violins
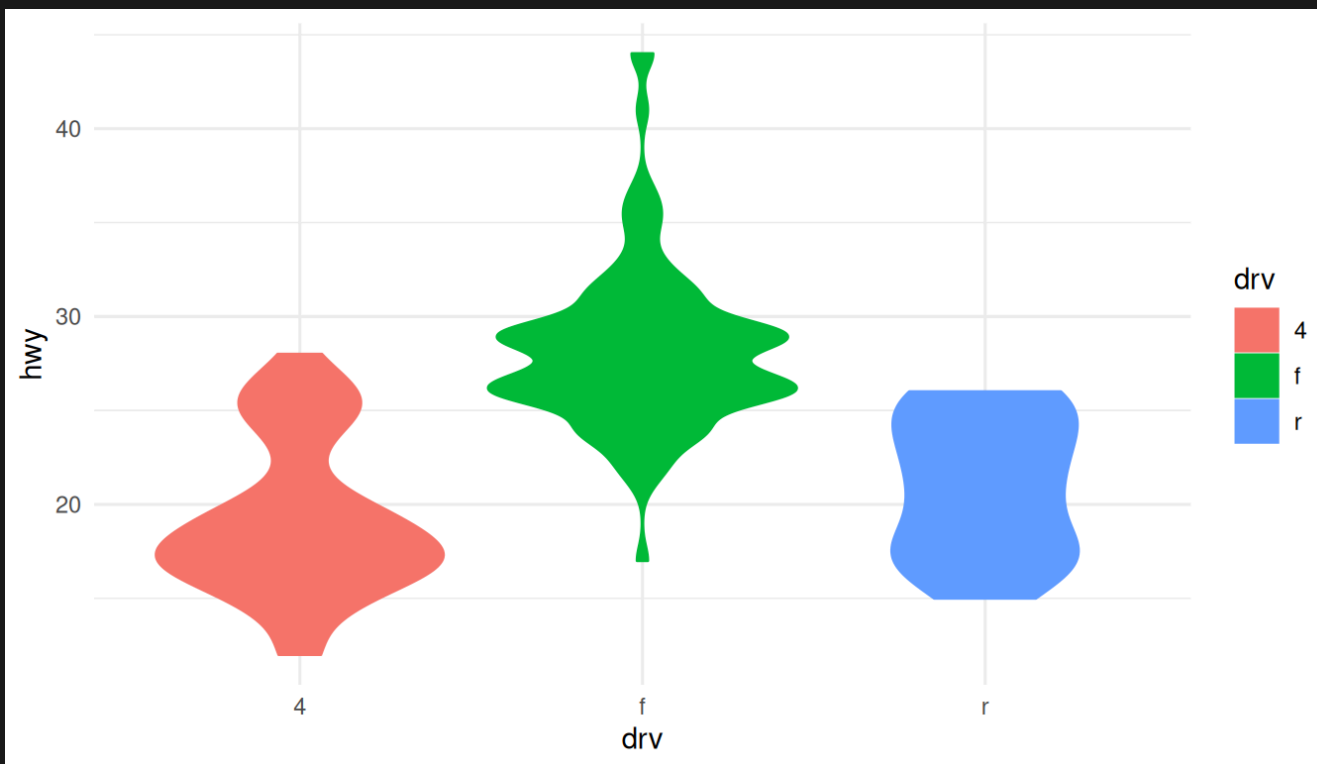
```
1  mpg %>% ggplot(aes(drv, hwy))+
2     geom_violin()
```

# An alternative to boxplot: violin

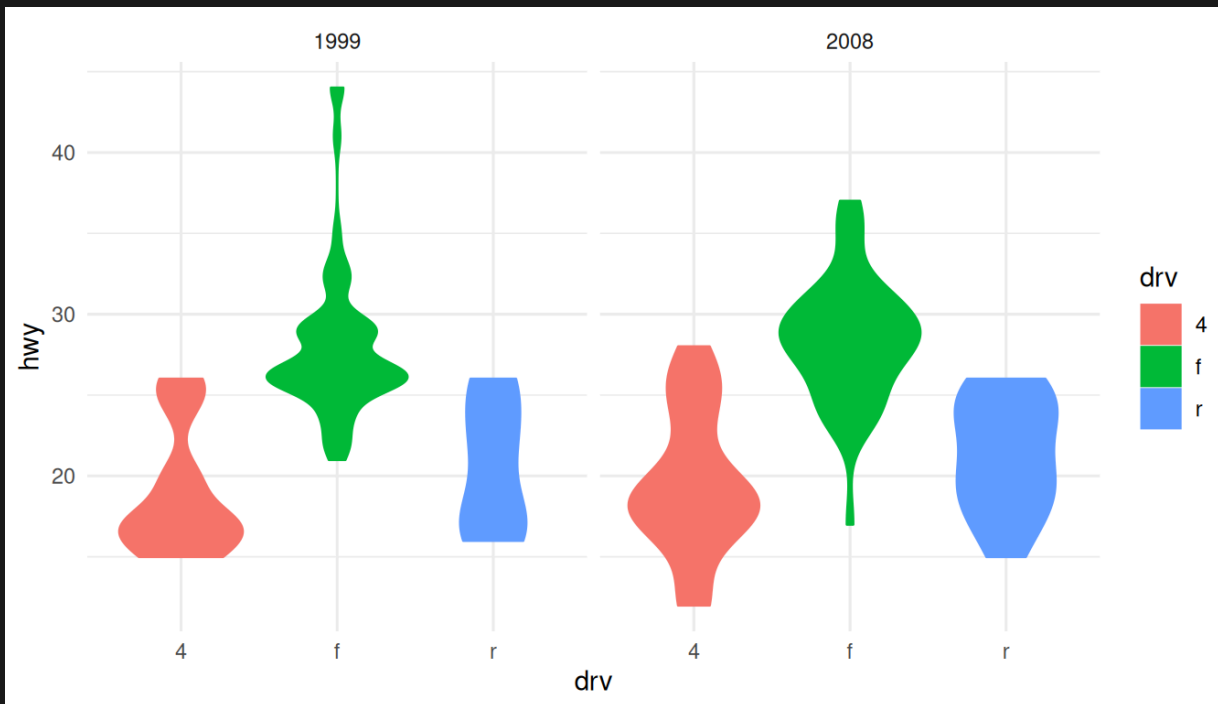remember: all is modular. We can always add color, fill...

```
1  mpg %>% ggplot(aes(drv, hwy, color = drv, fill = drv))+
2      geom_violin()
```

# An alternative to boxplot: violin
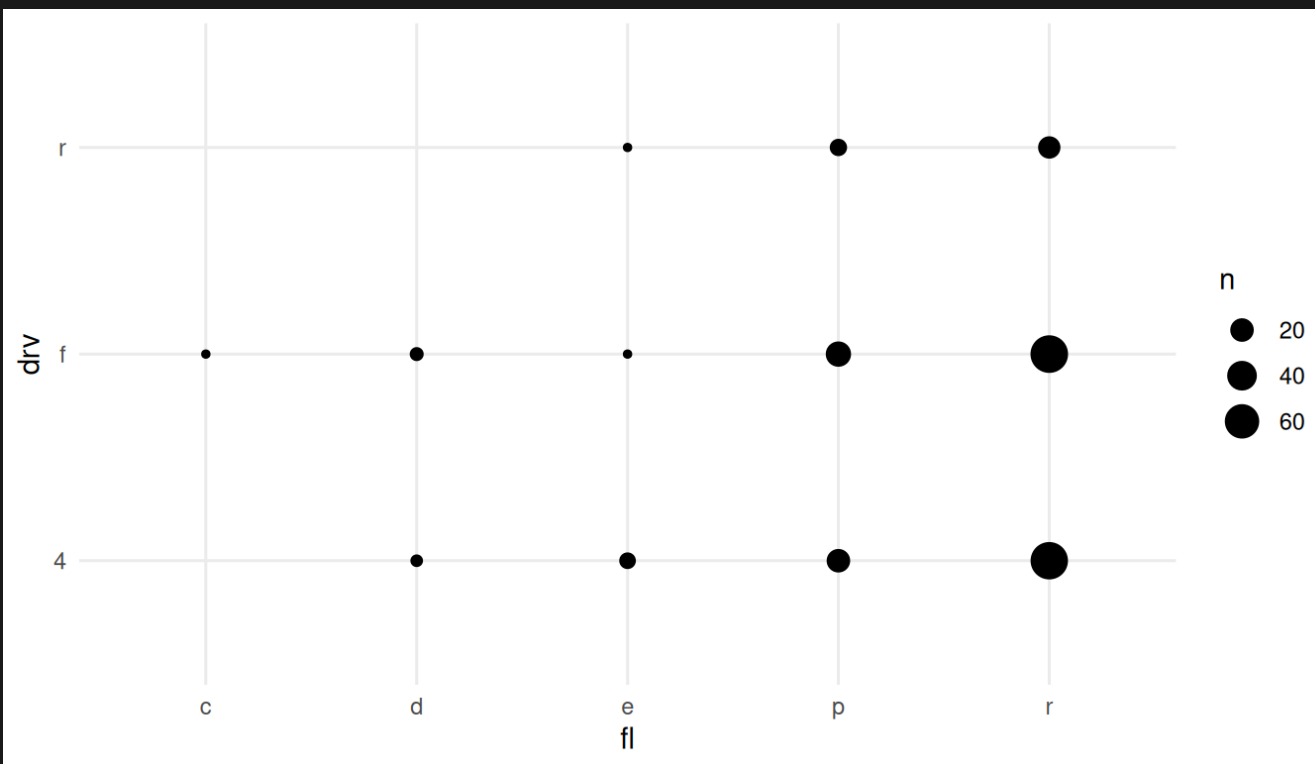
remember: all is modular. ...facets

```
1  mpg %>% ggplot(aes(drv, hwy, color = drv, fill = drv))+
2    geom_violin()+
3    facet_grid(.~year)
```

# Counts

two categorical variables: count their cross-tabulation

```
1  mpg %>% ggplot(aes(fl, drv))+
2     geom_count()
```

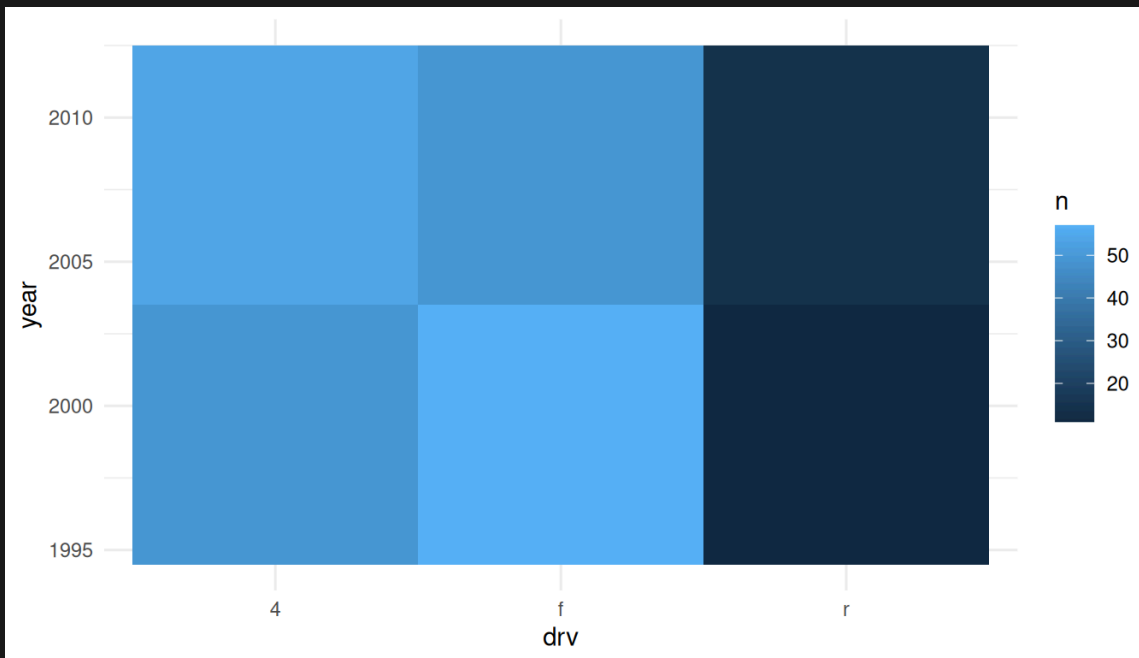# Exploring data: **three** variables

*Plot types depend on the variable type*

- *all continuous*: contour plot (think: elevation in maps)

- *some discrete*: tile

# Tile

two variables for the x,y grid. A third the color of the cell.
(needs `summarise()`!)

```
1  mpg %>% group_by(year, drv) %>% summarise(n = n()) %>%
2    ggplot(aes(x = drv, y = year, fill = n)) +  geom_tile()
```

# Additional resources

- the ggplot `cheatsheet` is your friend (Help -> cheatsheets)

- `stack overflow` helps out for trickier questions

- `chatGPT` is your friend, too (but beware)

- not feeling inspired?

  - 50 cool visualisations

  - a complete list of possibilities in R and python