# Intro to ggplot2

## Hadley Wickham

Assistant Professor / Dobelman Family Junior Chair
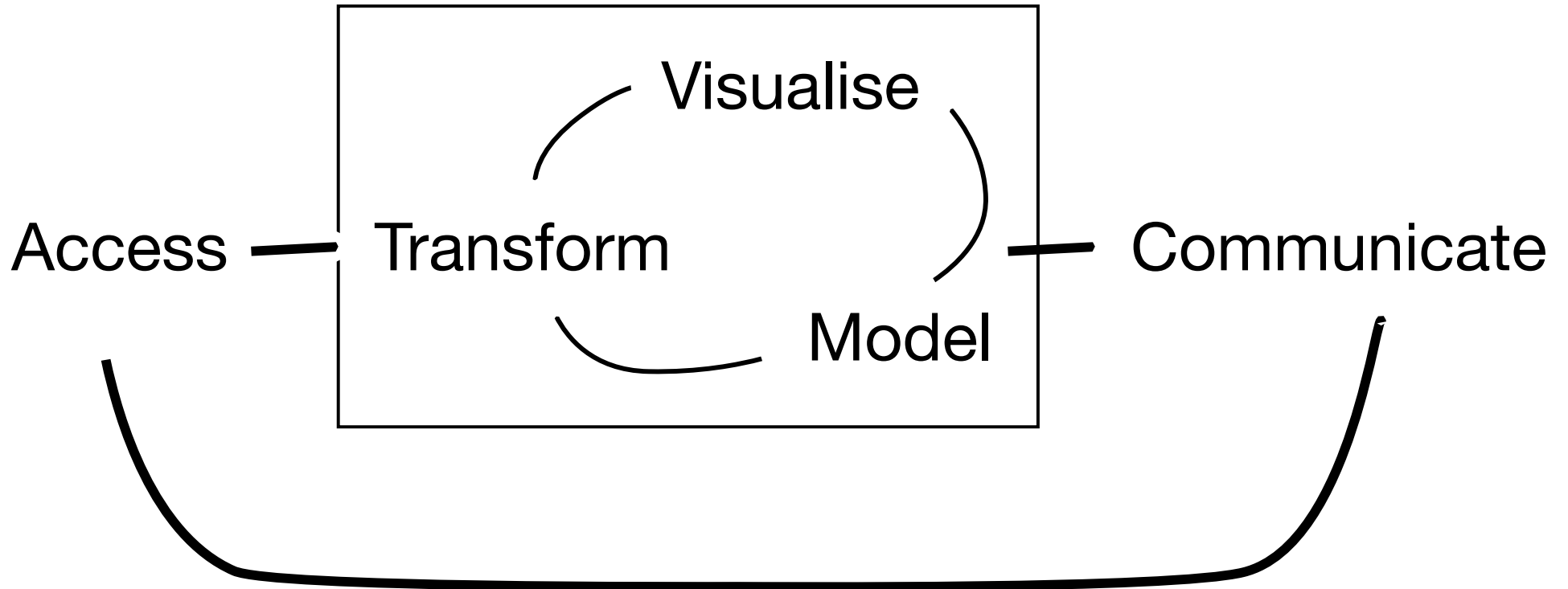Department of Statistics / Rice University
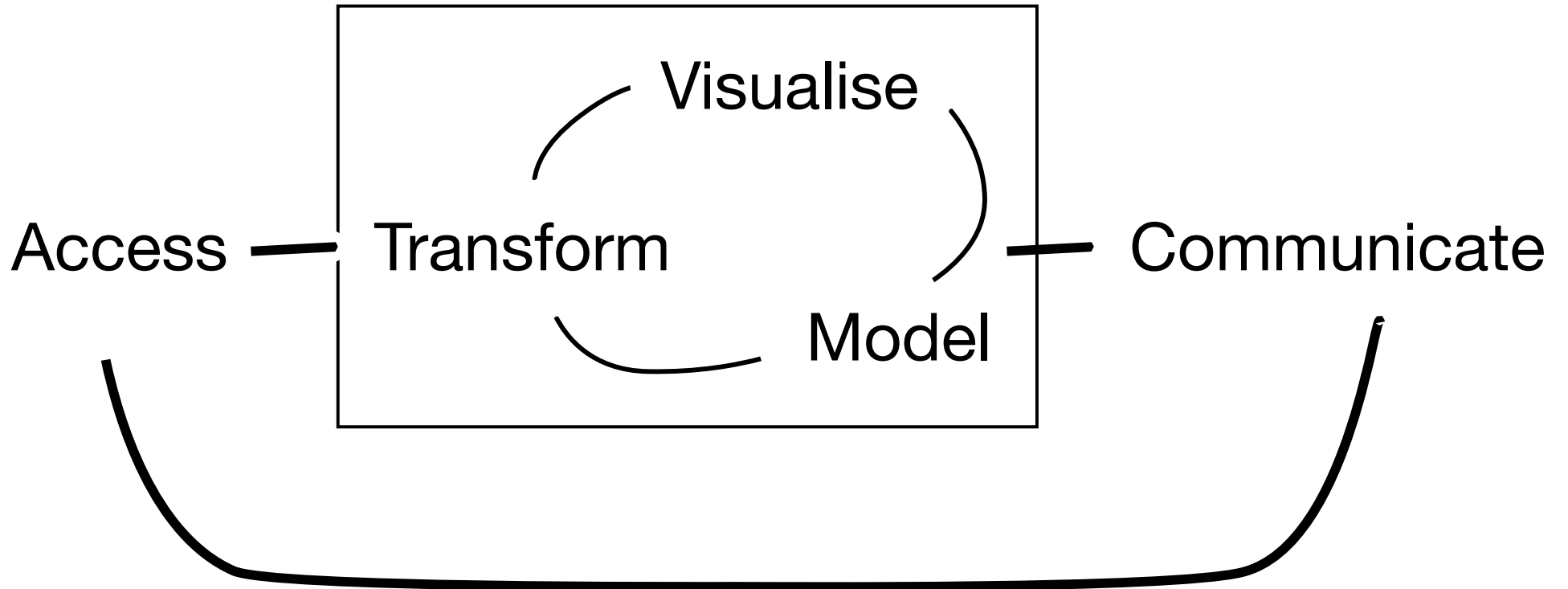
# HELLO

my name is

# Hadley

had.co.nz/courses/
10-tokyo

# Outline

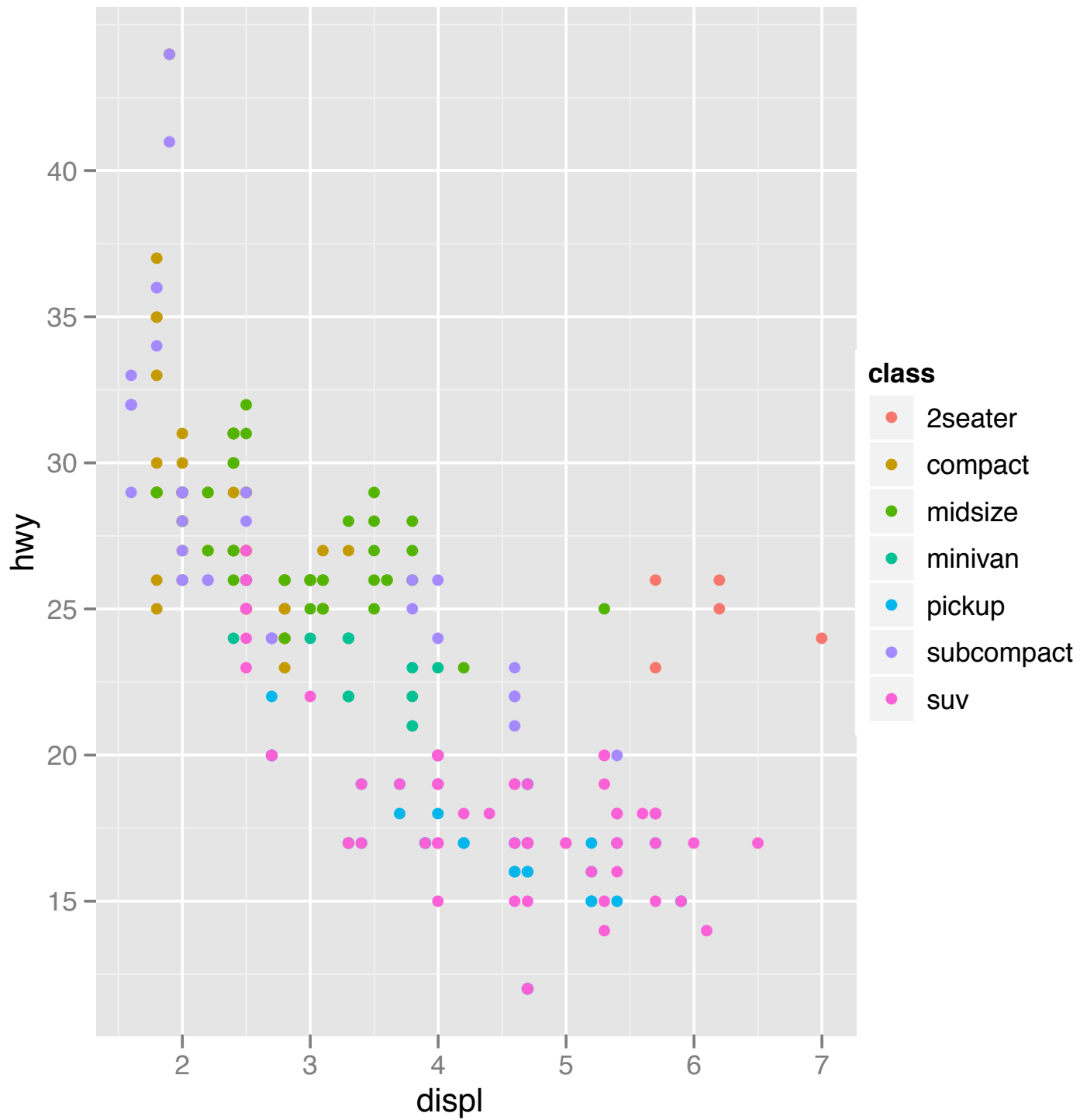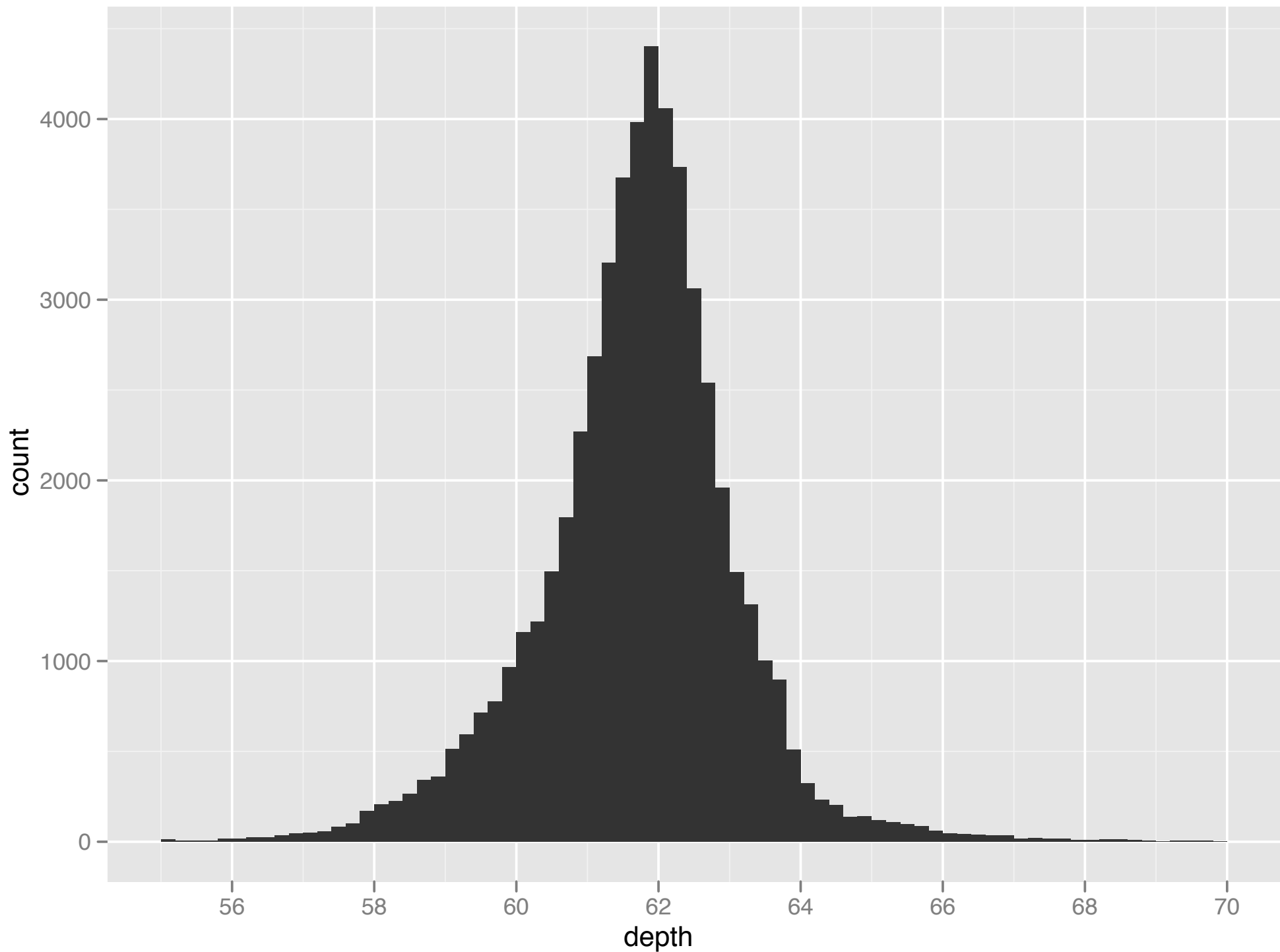# Data analysis is the process by which data becomes understanding, knowledge and insight

**Understand**

Access — Transform — Visualise / Model — Communicate

**Understand**

Access — Transform — Communicate
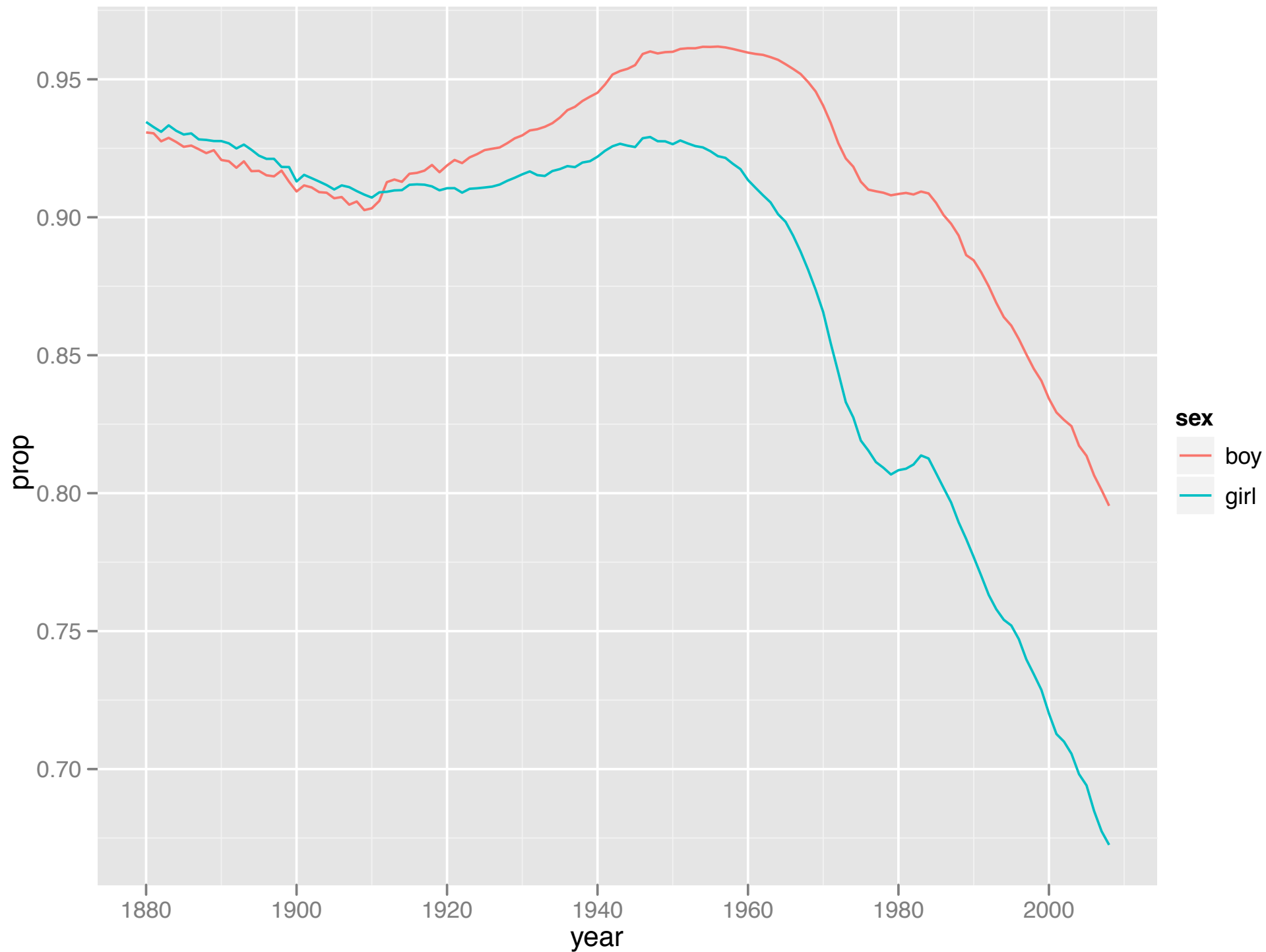
Visualise

Model

# Plotting basics

Learning a new language is hard!

# Scatterplot basics

```
install.packages("ggplot2")
library(ggplot2)

?mpg
head(mpg)
str(mpg)
summary(mpg)

qplot(displ, hwy, data = mpg)
```

Always explicitly specify the data

```
qplot(displ, hwy, data = mpg)
```

# Additional variables

Can display additional variables with **aesthetics** (like shape, colour, size) or **facetting** (small multiples displaying different subsets)

Legend chosen and displayed automatically.

```
qplot(displ, hwy, colour = class, data = mpg)
```
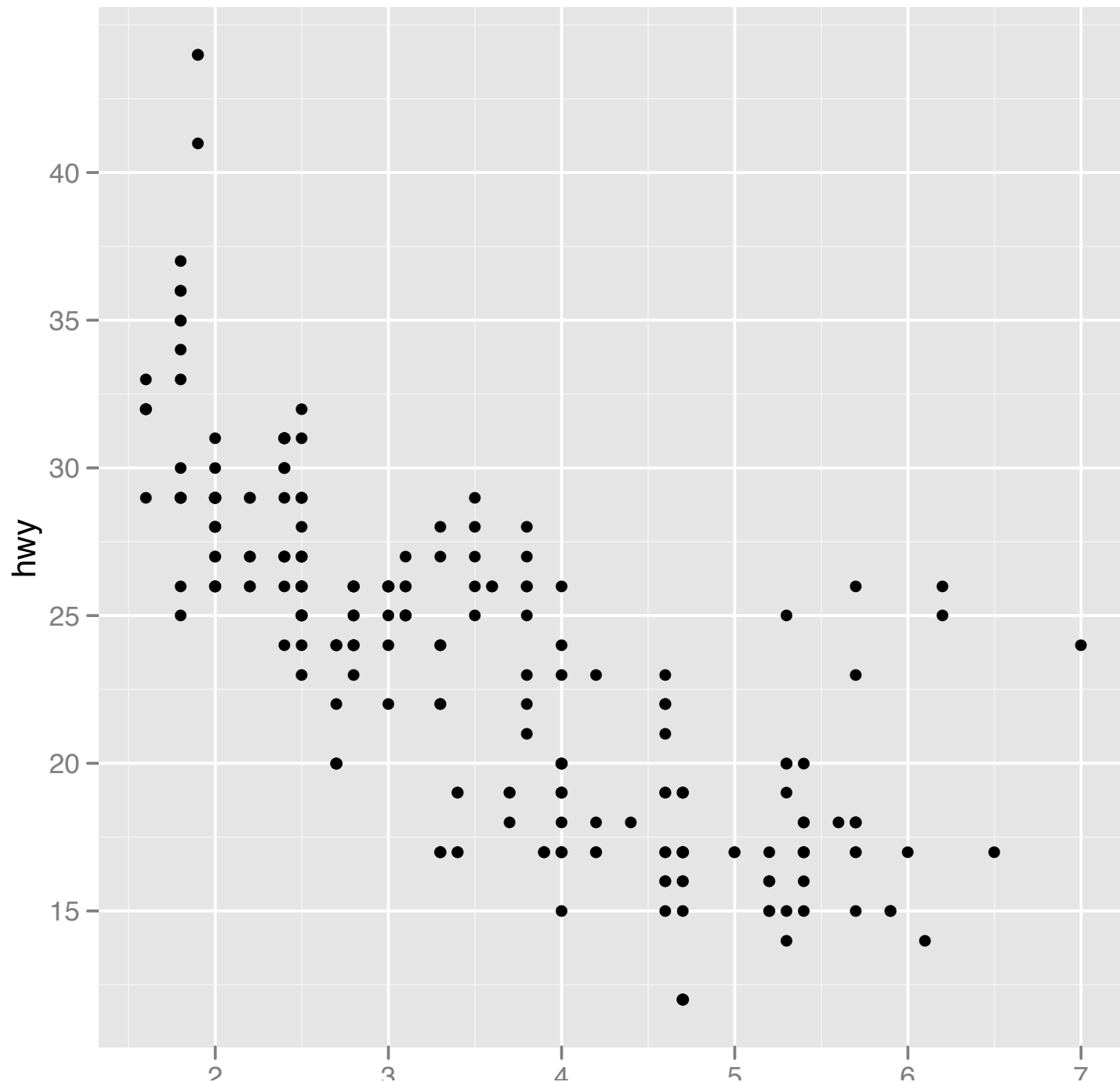
# Your turn

Try mapping different variables to the colour, size, and shape aesthetics.  Is there a difference between discrete and continuous variables? What happens when you use multiple aesthetics?

http://had.co.nz/courses/10-tokyo

# Aside: workflow

Keep a copy of the slides open so that you can copy and paste the code.

For complicated commands, write them in the script editor and then copy and paste.

|        | Discrete                | Continuous                                  |
| ------ | ----------------------- | ------------------------------------------- |
| Colour | Rainbow of colours      | Gradient from red to blue                   |
| Size   | Discrete size steps     | Linear mapping between radius and value     |
| Shape  | Different shape for each | Doesn't work                                |

# Faceting

Small multiples displaying different subsets of the data.

Useful for exploring conditional relationships.  Useful for large data.

# Your turn

```
qplot(displ, hwy, data = mpg) +
facet_grid(. ~ cyl)

qplot(displ, hwy, data = mpg) +
facet_grid(drv ~ .)

qplot(displ, hwy, data = mpg) +
facet_grid(drv ~ cyl)

qplot(displ, hwy, data = mpg) +
facet_wrap(~ class)
```
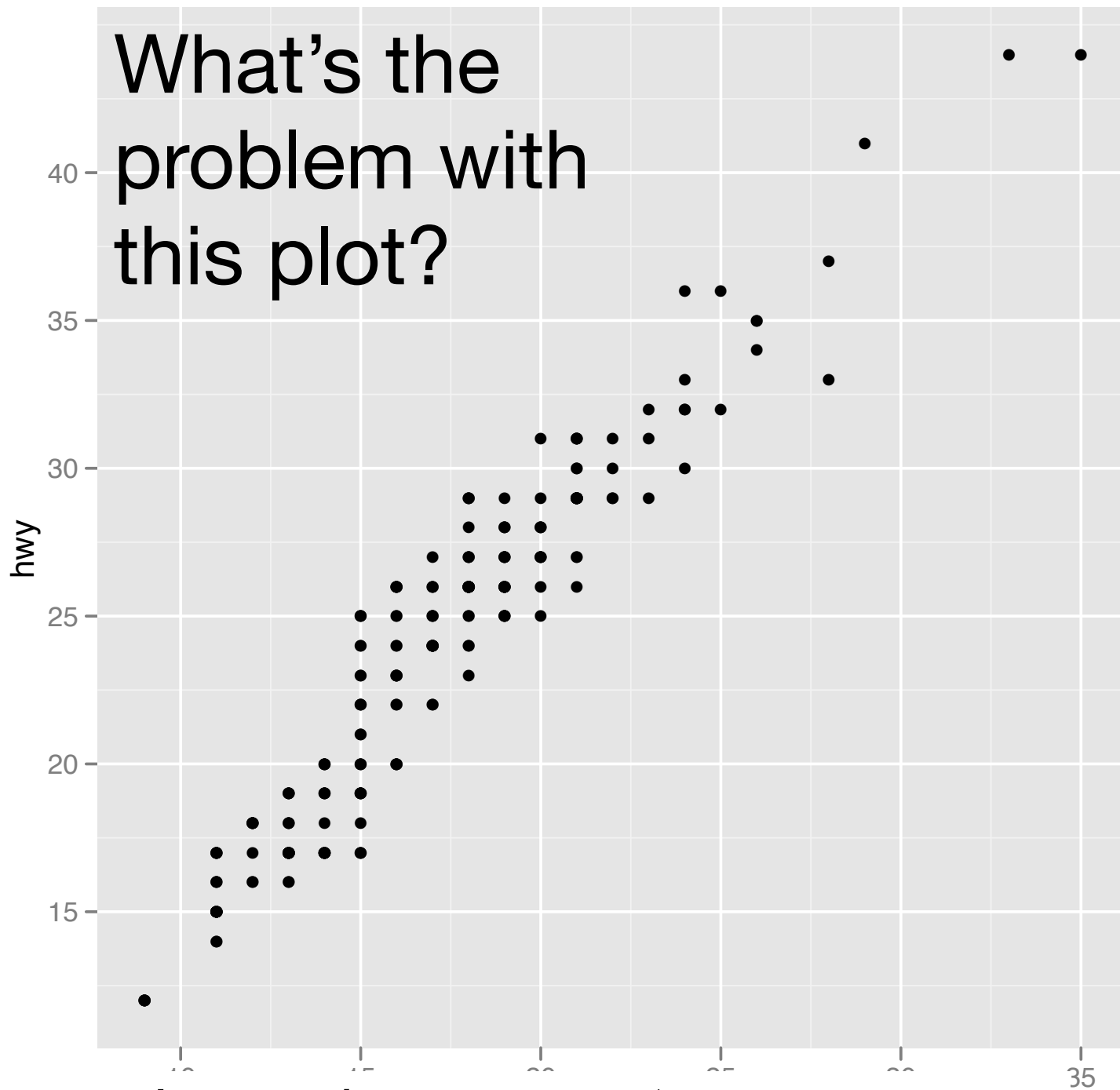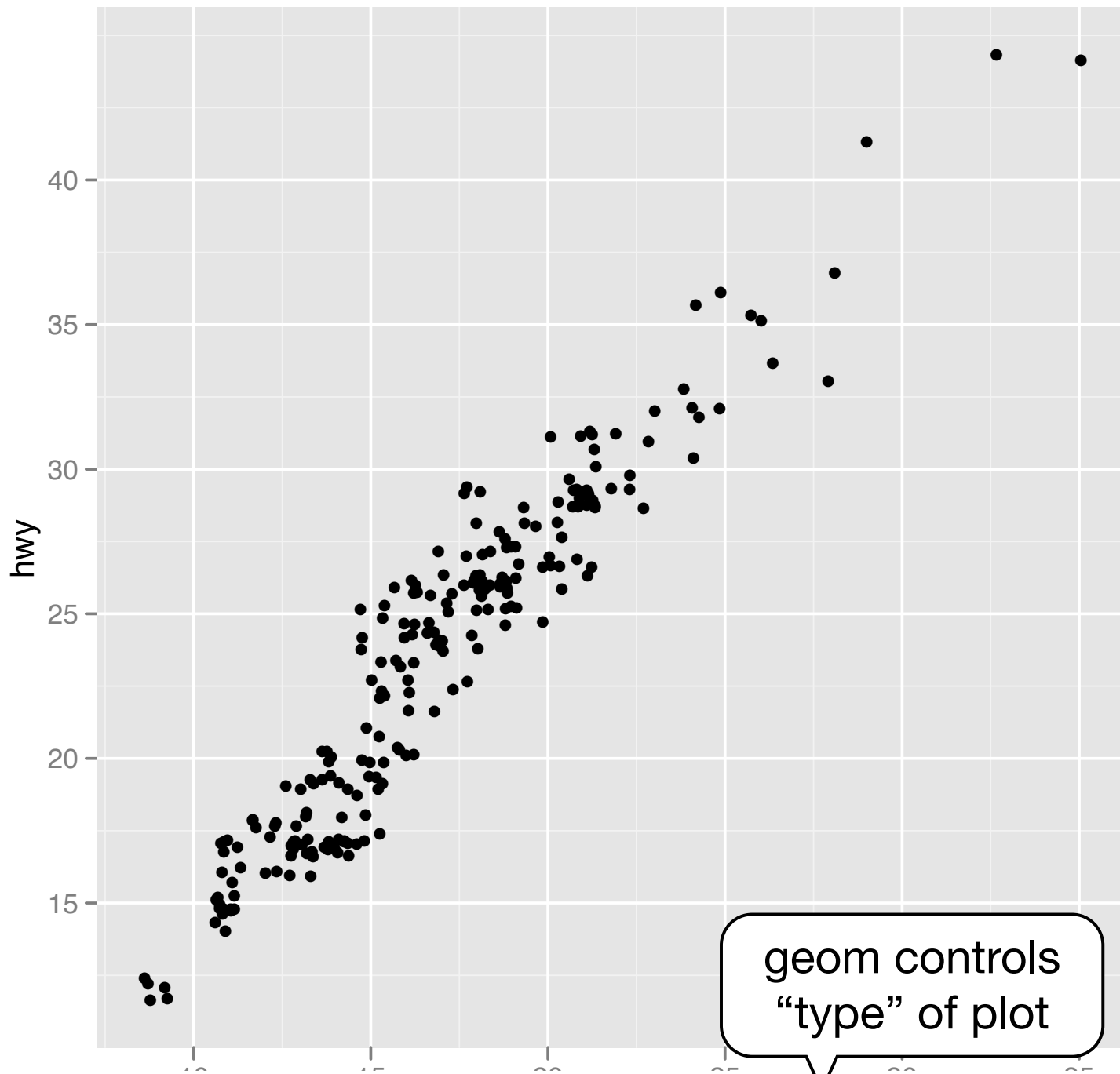
# Summary

`facet_grid()`: 2d grid, rows ~ cols, . for no split

`facet_wrap()`: 1d ribbon wrapped into 2d

What's the problem with this plot?
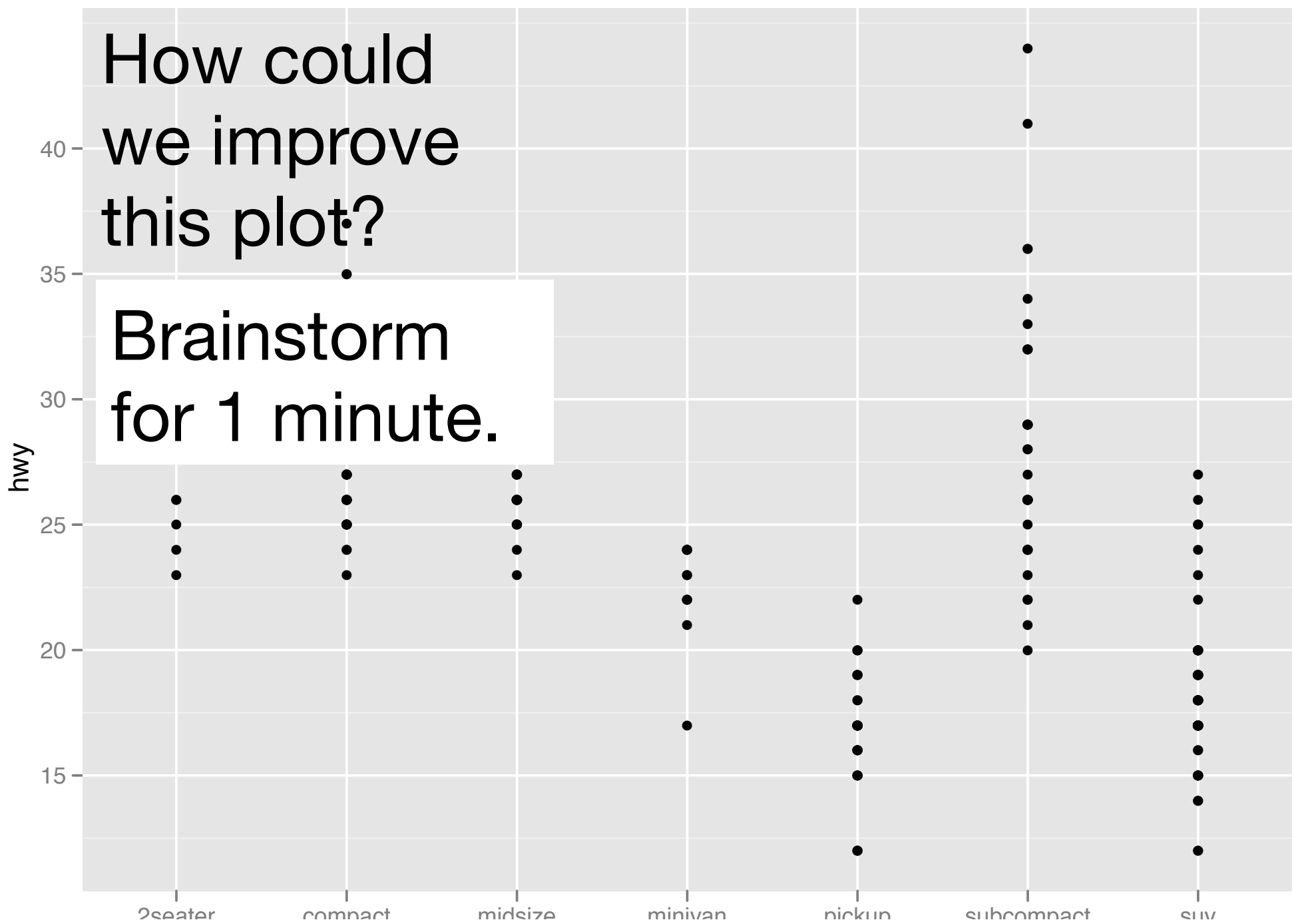
```
qplot(cty, hwy, data = mpg)
```
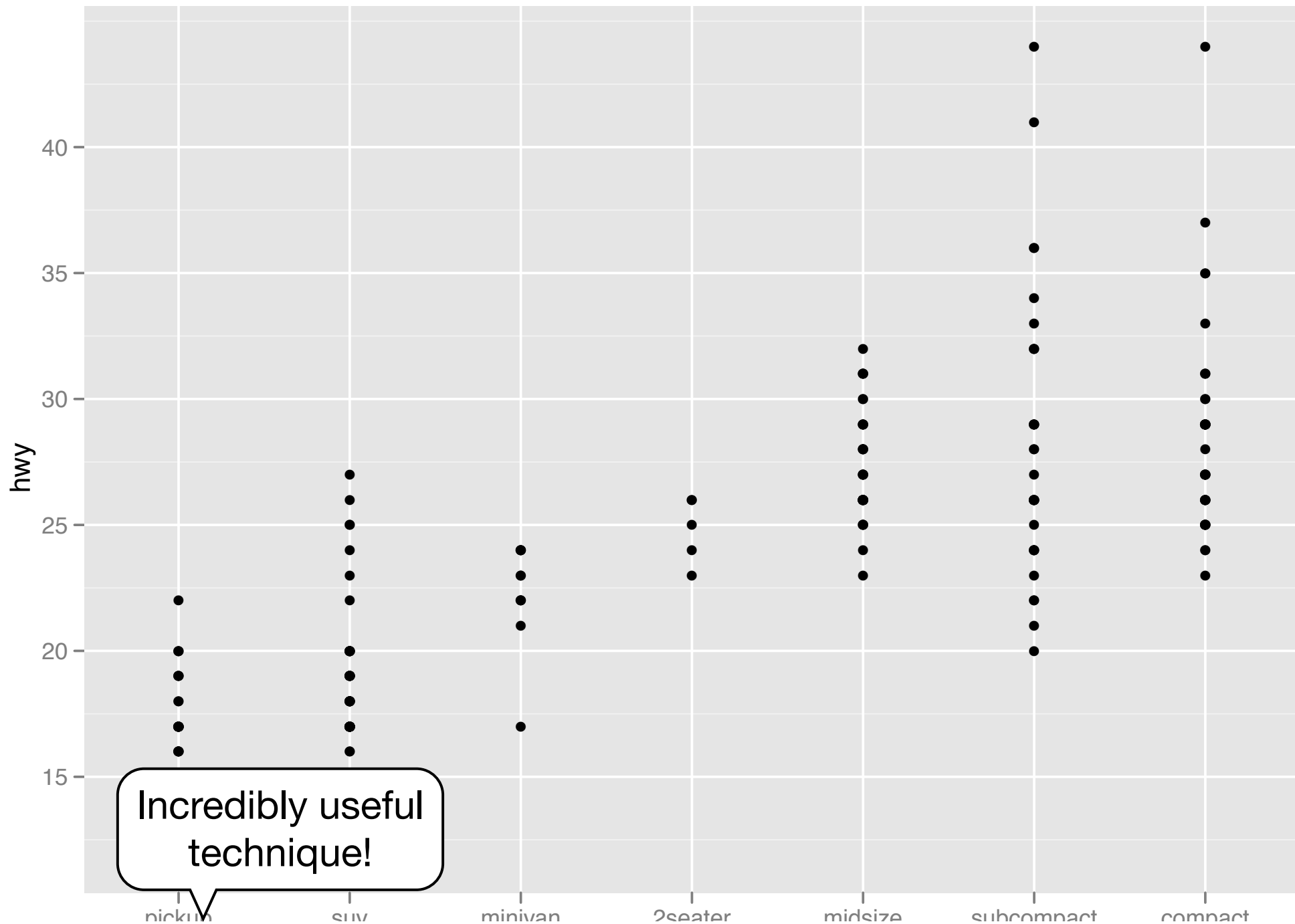
geom controls
"type" of plot

```
qplot(cty, hwy, data = mpg, geom = "jitter")
```
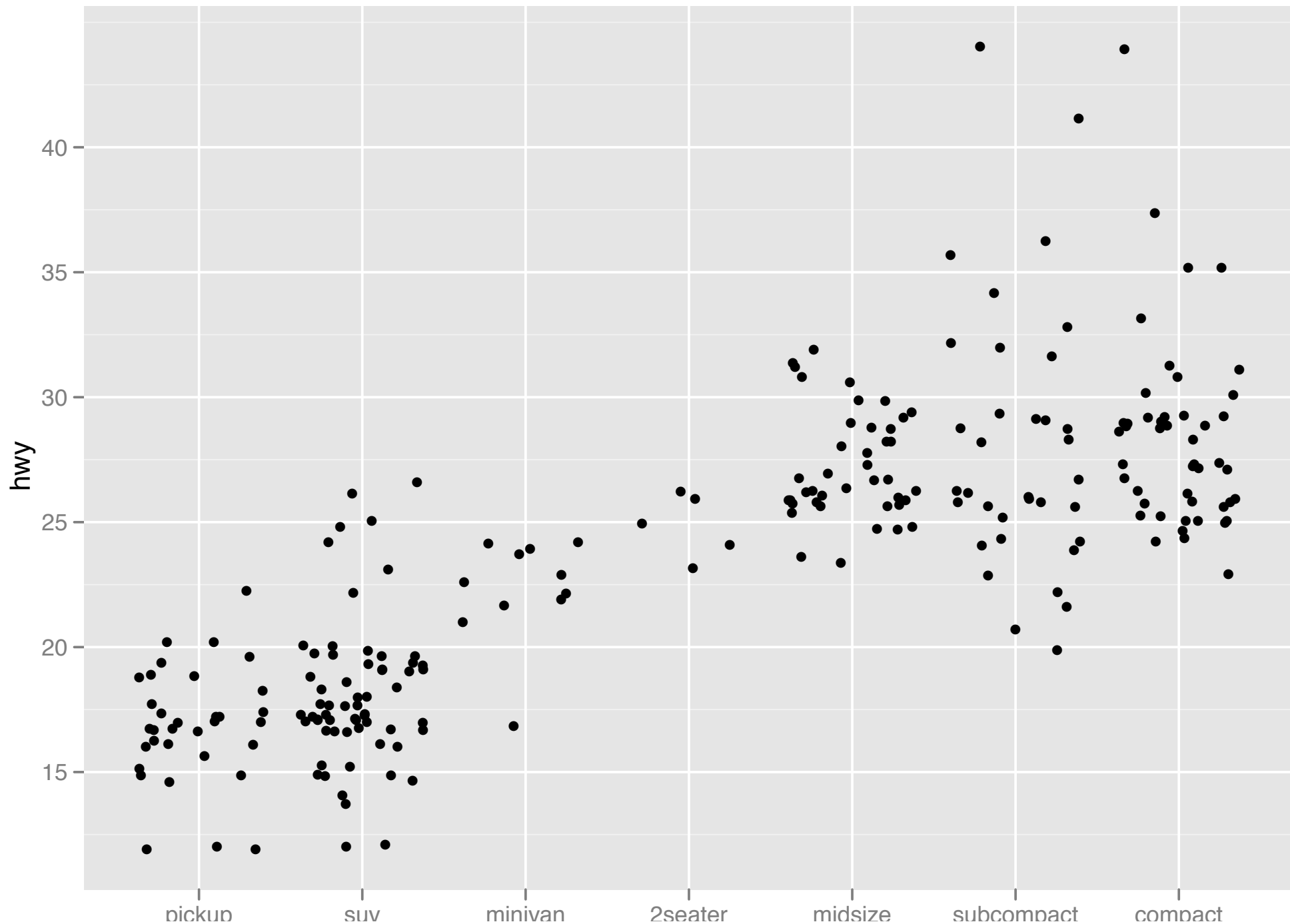
How could
we improve
this plot?

Brainstorm
for 1 minute.

hwy

40
35
30
25
20
15

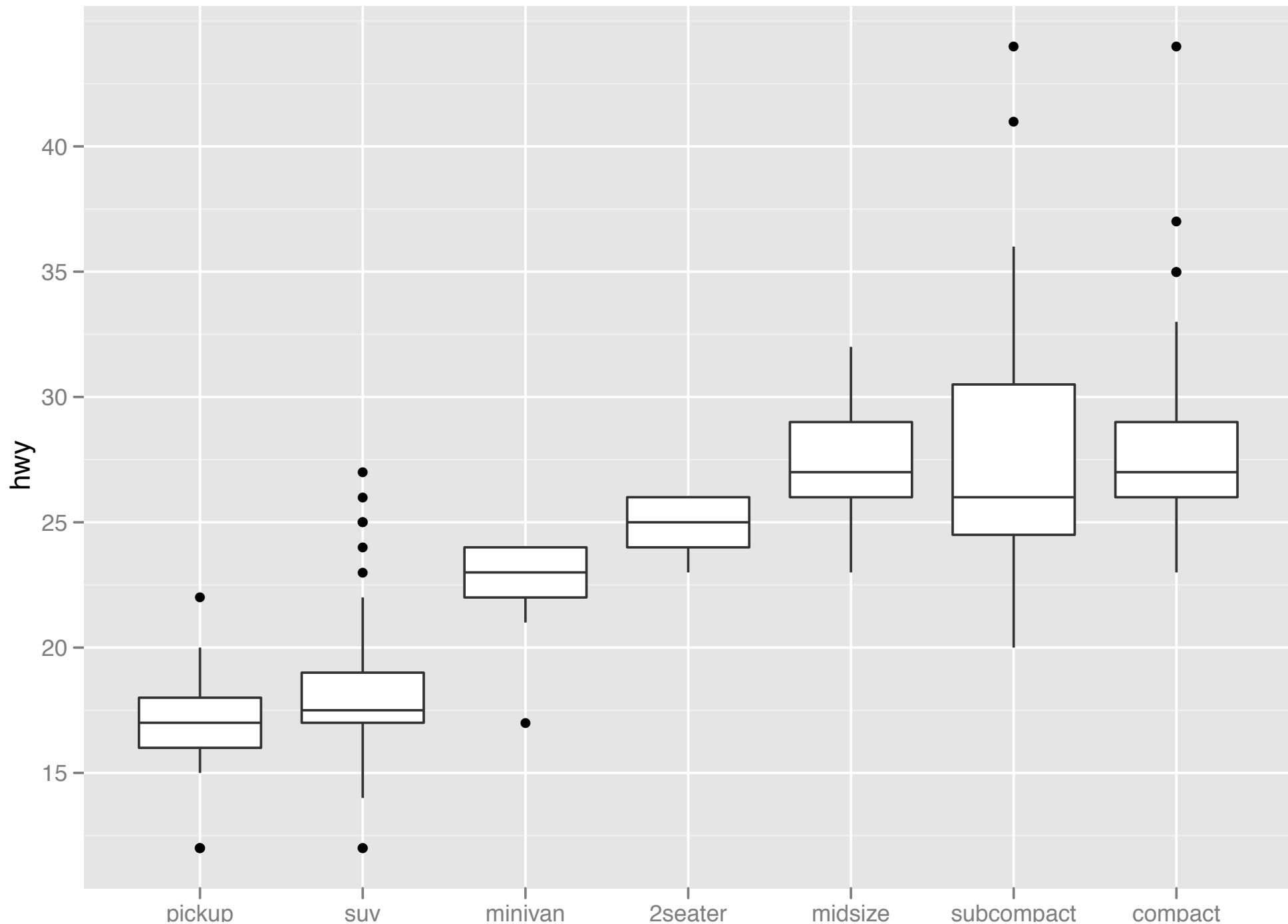2seater    compact    midsize    minivan    pickup    subcompact    suv
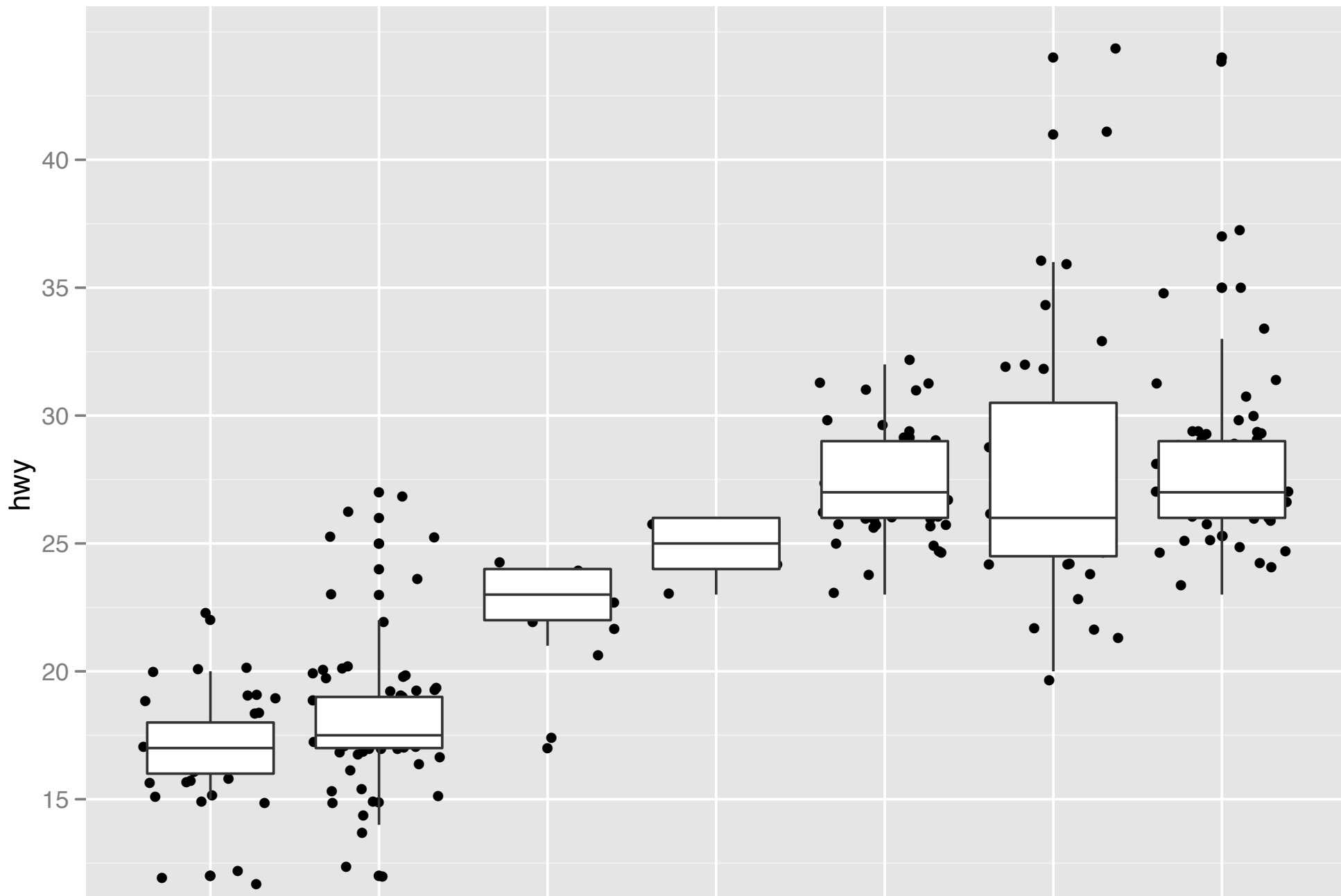
```
qplot(class, hwy, data = mpg)
```

```
qplot(reorder(class, hwy), hwy, data = mpg)
```

```
qplot(reorder(class, hwy), hwy, data = mpg, geom = "jitter")
```

```
qplot(reorder(class, hwy), hwy, data = mpg, geom = "boxplot")
```

```
qplot(reorder(class, hwy), hwy, data = mpg,
  geom = c("jitter", "boxplot"))
```

# Your turn

Read the help for reorder.  Redraw the previous plots with class ordered by median hwy.

How would you put the jittered points on top of the boxplots?