

Stat405

Statistical computing & graphics

Hadley Wickham

1. Introductions

2. Syllabus

3. Introduction to linux

4. Introduction to R

5. Diving in

HELLO

my name is

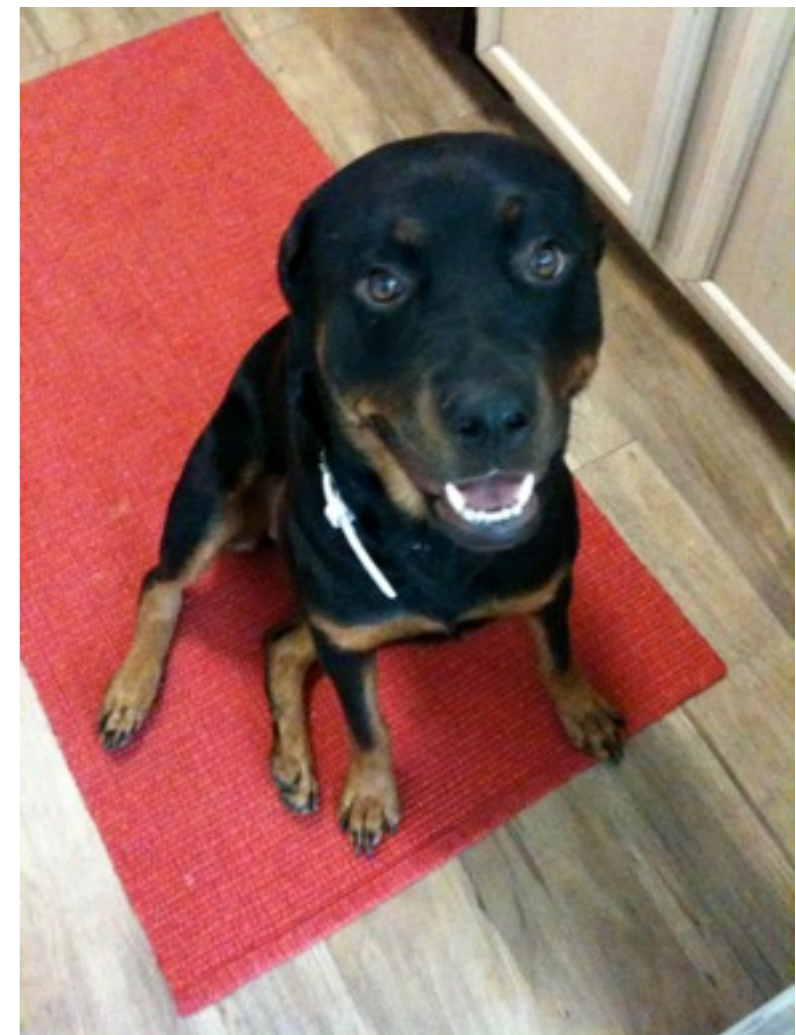
Hadley

At Rice, I'm a

- McMurry divisional advisor,
- major advisor for Statistics

Away from Rice, I

- love to cook
- have two dogs
- travel too much



<http://stat405.had.co.nz>
hadley@rice.edu

Syllabus

Tools

Computer: mac/windows/linux

Software: R, text editor, latex (= Rstudio)

Brain: scepticism, curiosity, organisation

Homework

Lowest grade dropped.

Honour code: you can discuss ideas with other class members, but you must present your own work. All code should be your own.

Late policy: 20% penalty if turned in by 9am Monday. Homeworks not accepted after that time.

All homeworks must be submitted in physical form. Electronic versions will only be accepted in exceptional circumstances.

Team projects

3 bigger team projects, culminating in poster presentation at the end of year.

Teams of 4-5 people, assigned by Hadley.

Will teach team work skills. Option to disband after first project. Firing and quitting.

Rstudio

Setup

Install R and Rstudio on your computer, following the instructions on the class website.

You can also use Rstudio online:

<https://www.clear.rice.edu/rstudio>

(note that the files are saved in your Rice computing account)

RStudio

Workspace History

Files Plots Packages Help

Zoom Export Clear All

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

1:30 (Top Level) R Script

Console ~/

or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

```
> library(ggplot2)  
Loading required package: reshape  
Loading required package: plyr  
  
Attaching package: 'reshape'  
  
The following object(s) are masked from 'package:plyr':  
  
  rename, round_any  
  
Loading required package: grid  
Loading required package: proto  
> qplot(displ, cyl, data = mpg)  
> qplot(displ, hwy, data = mpg)  
>
```

hwy

displ

The image shows the RStudio interface. The top-left pane is the Source Editor, containing the code `1 qplot(displ, hwy, data = mpg)`. The top-right pane is the Plots window, displaying a scatter plot of highway mileage (hwy) on the y-axis (ranging from 15 to 40) against engine displacement (displ) on the x-axis (ranging from 2 to 7). The plot shows a negative correlation between the two variables. The bottom-left pane is the Console, showing the output of running the code, including the loading of the `ggplot2` package and its dependencies (`reshape`, `plyr`, `grid`, `proto`).

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

```
Console ~/
or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> library(ggplot2)
Loading required package: reshape
Loading required package: plyr

Attaching package: 'reshape'

The following object(s) are masked from 'package:plyr':

  rename, round_any

Loading required package: grid
Loading required package: proto
> qplot(displ, cyl, data = mpg)
> qplot(displ, hwy, data = mpg)
>
```

Console – run code here

RStudio

Workspace History

Untitled1* x

Source on Save Run Source

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

1:30 (Top Level) R Script

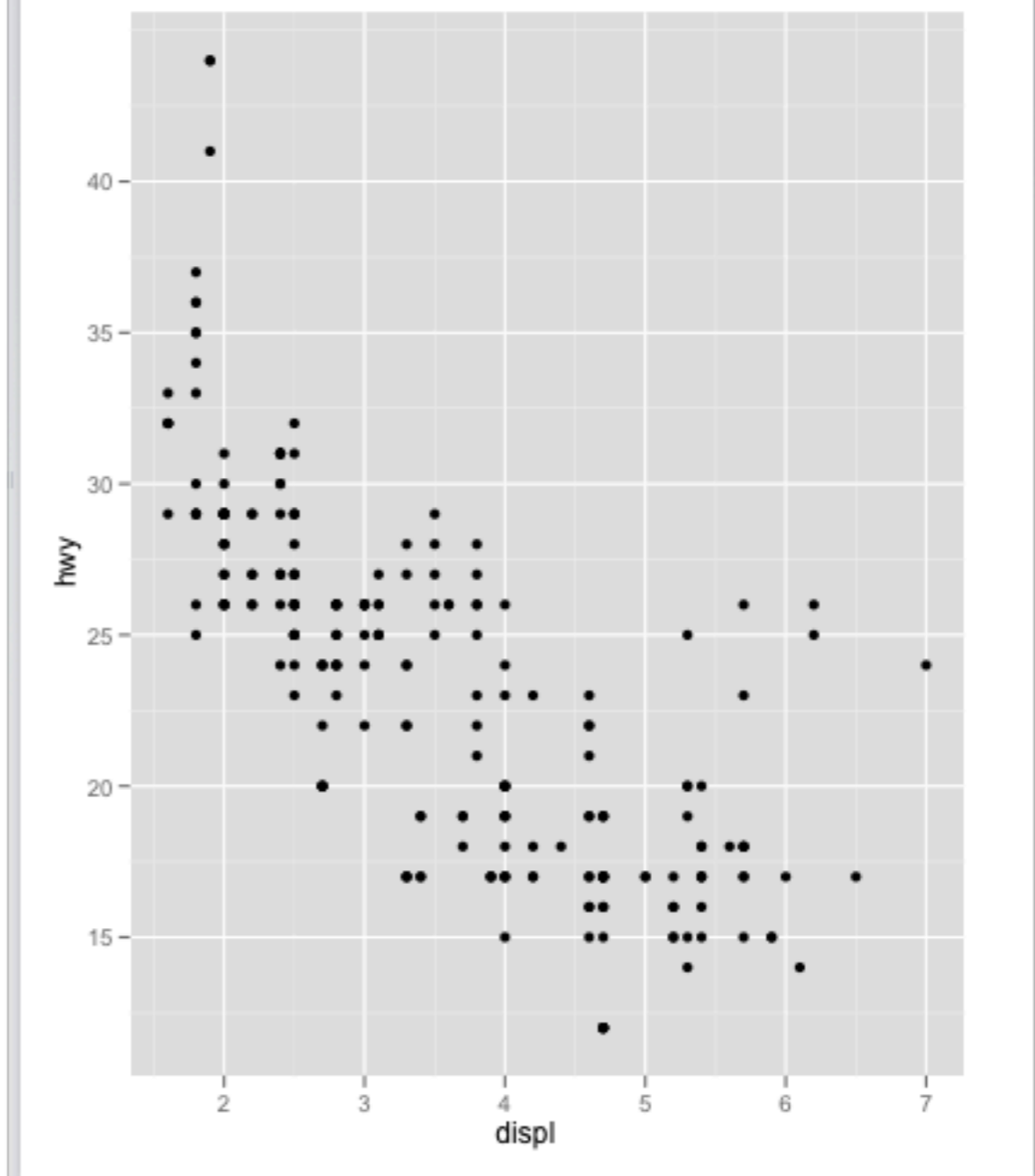
Console ~/

or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

```
> library(ggplot2)  
Loading required package: reshape  
Loading required package: plyr  
  
Attaching package: 'reshape'  
  
The following object(s) are masked from 'package:plyr':  
  
  rename, round_any  
  
Loading required package: grid  
Loading required package: proto  
> qplot(displ, cyl, data = mpg)  
> qplot(displ, hwy, data = mpg)  
>
```

Files Plots Packages Help

Zoom Export Clear All



displ	hwy
1.8	44
1.8	42
1.8	37
1.8	36
1.8	35
1.8	34
1.8	33
1.8	32
1.8	31
1.8	30
1.8	29
1.8	28
1.8	27
1.8	26
1.8	25
1.8	24
1.8	23
1.8	22
1.8	21
1.8	20
1.8	19
1.8	18
1.8	17
1.8	16
1.8	15
2.0	31
2.0	30
2.0	29
2.0	28
2.0	27
2.0	26
2.0	25
2.0	24
2.0	23
2.0	22
2.0	21
2.0	20
2.0	19
2.0	18
2.0	17
2.0	16
2.0	15
2.4	31
2.4	30
2.4	29
2.4	28
2.4	27
2.4	26
2.4	25
2.4	24
2.4	23
2.4	22
2.4	21
2.4	20
2.4	19
2.4	18
2.4	17
2.4	16
2.4	15
2.8	29
2.8	28
2.8	27
2.8	26
2.8	25
2.8	24
2.8	23
2.8	22
2.8	21
2.8	20
2.8	19
2.8	18
2.8	17
2.8	16
2.8	15
3.2	29
3.2	28
3.2	27
3.2	26
3.2	25
3.2	24
3.2	23
3.2	22
3.2	21
3.2	20
3.2	19
3.2	18
3.2	17
3.2	16
3.2	15
3.6	29
3.6	28
3.6	27
3.6	26
3.6	25
3.6	24
3.6	23
3.6	22
3.6	21
3.6	20
3.6	19
3.6	18
3.6	17
3.6	16
3.6	15
4.0	29
4.0	28
4.0	27
4.0	26
4.0	25
4.0	24
4.0	23
4.0	22
4.0	21
4.0	20
4.0	19
4.0	18
4.0	17
4.0	16
4.0	15
4.4	29
4.4	28
4.4	27
4.4	26
4.4	25
4.4	24
4.4	23
4.4	22
4.4	21
4.4	20
4.4	19
4.4	18
4.4	17
4.4	16
4.4	15
4.8	29
4.8	28
4.8	27
4.8	26
4.8	25
4.8	24
4.8	23
4.8	22
4.8	21
4.8	20
4.8	19
4.8	18
4.8	17
4.8	16
4.8	15
5.2	29
5.2	28
5.2	27
5.2	26
5.2	25
5.2	24
5.2	23
5.2	22
5.2	21
5.2	20
5.2	19
5.2	18
5.2	17
5.2	16
5.2	15
5.6	29
5.6	28
5.6	27
5.6	26
5.6	25
5.6	24
5.6	23
5.6	22
5.6	21
5.6	20
5.6	19
5.6	18
5.6	17
5.6	16
5.6	15
6.0	29
6.0	28
6.0	27
6.0	26
6.0	25
6.0	24
6.0	23
6.0	22
6.0	21
6.0	20
6.0	19
6.0	18
6.0	17
6.0	16
6.0	15
6.4	29
6.4	28
6.4	27
6.4	26
6.4	25
6.4	24
6.4	23
6.4	22
6.4	21
6.4	20
6.4	19
6.4	18
6.4	17
6.4	16
6.4	15
6.8	29
6.8	28
6.8	27
6.8	26
6.8	25
6.8	24
6.8	23
6.8	22
6.8	21
6.8	20
6.8	19
6.8	18
6.8	17
6.8	16
6.8	15
7.2	29
7.2	28
7.2	27
7.2	26
7.2	25
7.2	24
7.2	23
7.2	22
7.2	21
7.2	20
7.2	19
7.2	18
7.2	17
7.2	16
7.2	15

Output – plots and help

The image shows the RStudio interface. The top-left pane is the R Script editor, titled 'Untitled1*'. It contains the following code:

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

The bottom-left pane is the Console, showing the following output:

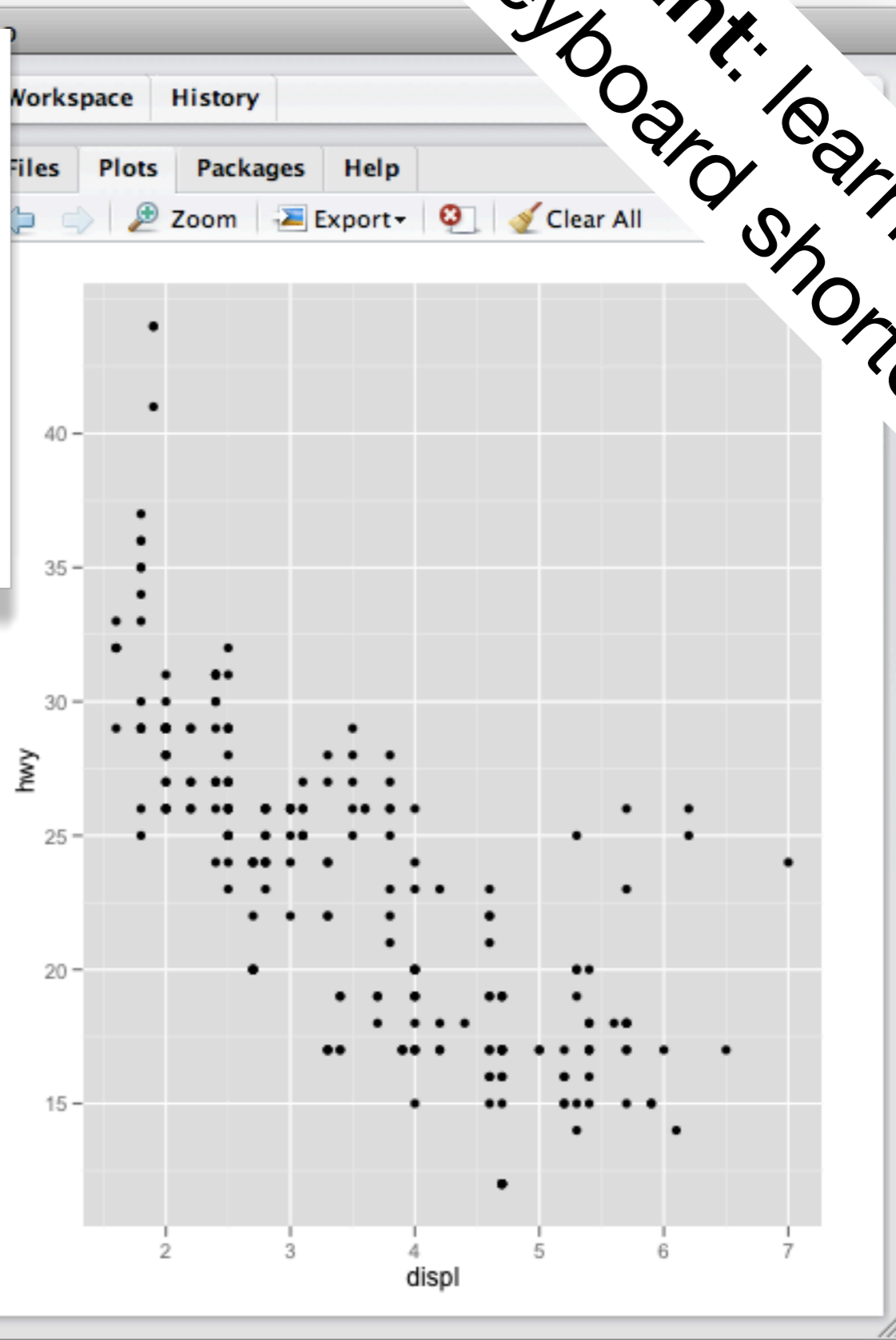
```
or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
> library(ggplot2)  
Loading required package: reshape  
Loading required package: plyr  
  
Attaching package: 'reshape'  
  
The following object(s) are masked from 'package:plyr':  
  
  rename, round_any  
  
Loading required package: grid  
Loading required package: proto  
> qplot(displ, cyl, data = mpg)  
> qplot(displ, hwy, data = mpg)  
>
```

The right pane shows the resulting scatter plot. The x-axis is labeled 'displ' and ranges from approximately 1.8 to 7.0. The y-axis is labeled 'hwy' and ranges from approximately 12 to 45. The plot shows a negative correlation between engine displacement and highway mileage.

Editor – save code here

Hint: learn the keyboard shortcuts

```
Untitled1* x  
Source on Save | Search | Run | Source  
1 qplot(displ, hwy, data = mpg)|  
1:30 (Top Level) R Script
```



```
Console ~/  
or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
> library(ggplot2)  
Loading required package: reshape  
Loading required package: plyr  
  
Attaching package: 'reshape'  
  
The following object(s) are masked from 'package:plyr':  
  
  rename, round_any  
  
Loading required package: grid  
Loading required package: proto  
> qplot(displ, cyl, data = mpg)  
> qplot(displ, hwy, data = mpg)  
>
```

Editor – save code here

Short cuts

In editor:

Command/ctrl + enter: send code to console

Ctrl + 2: move cursor to console

In console:

Up arrow: retrieve previous command

Ctrl + up arrow: search commands

Ctrl + 1: move cursor to editor

Introduction to R



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)
```

Scatterplot basics

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

```
?mpg
```

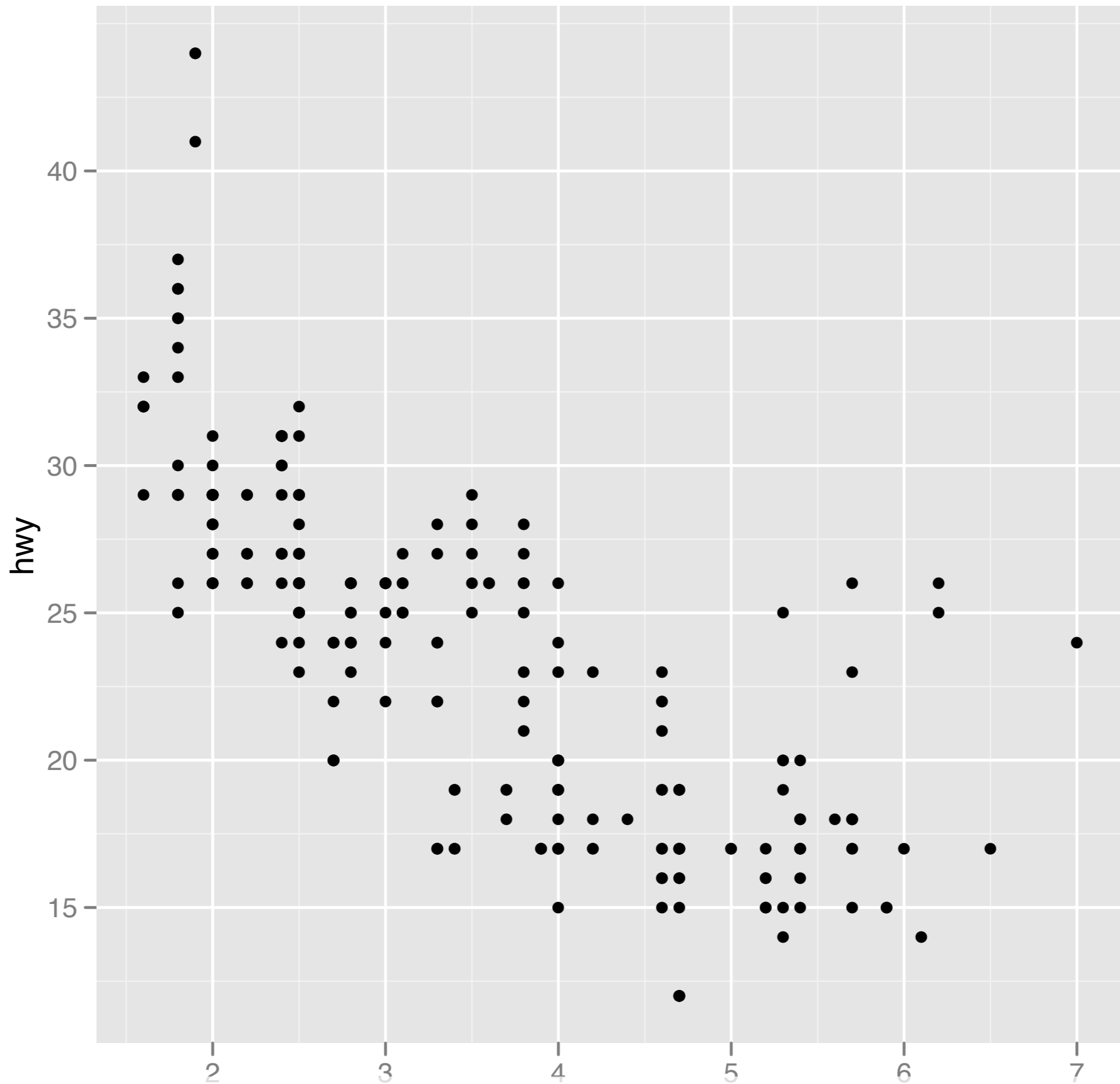
```
head(mpg)
```

```
str(mpg)
```

```
summary(mpg)
```

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

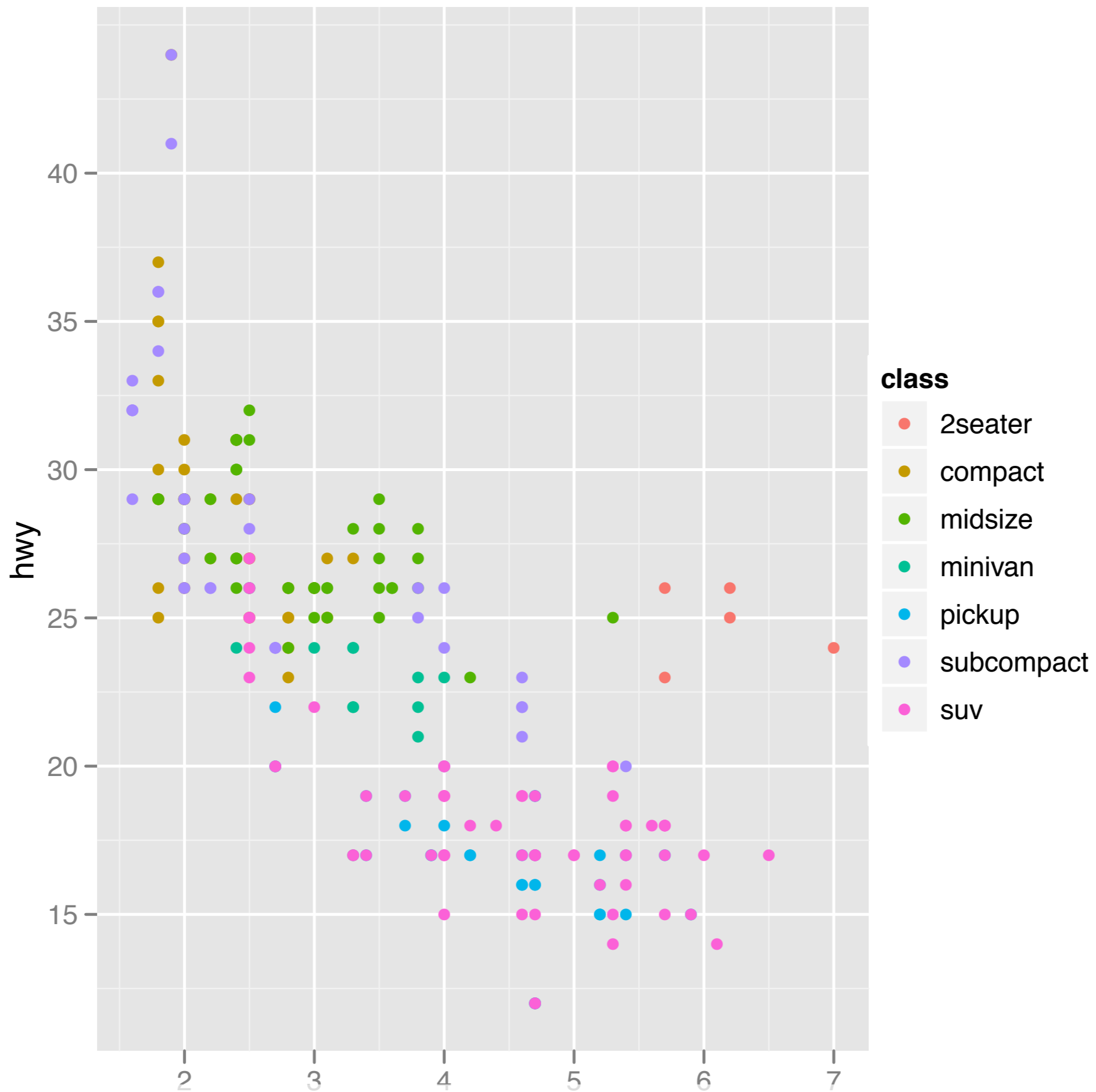
Always explicitly
specify the data



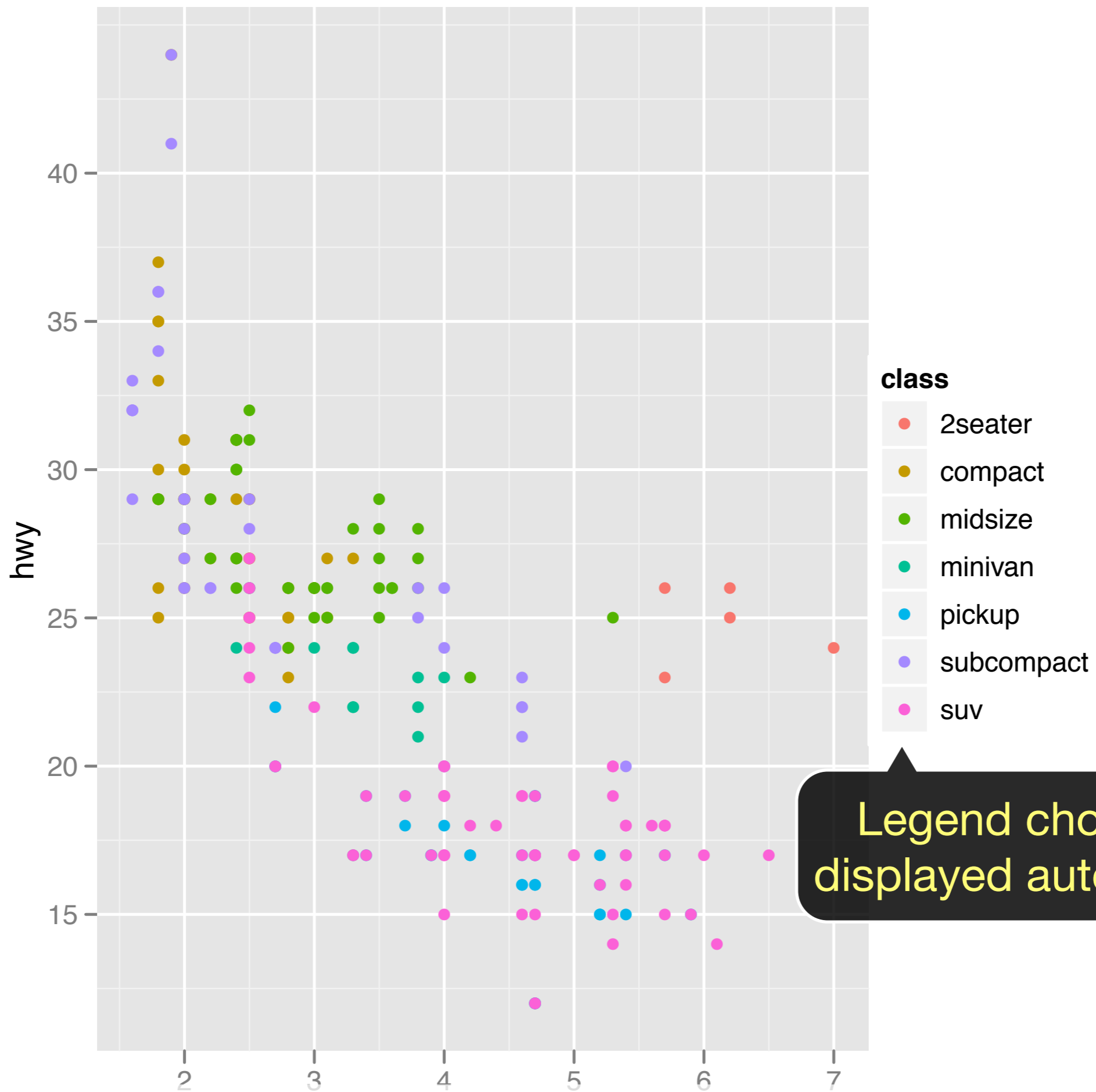
```
qplot(displ, hwy, data = mpg)
```

Additional variables

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



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

Legend chosen and displayed automatically.

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

Your turn

Experiment with colour, size, and shape aesthetics.

What's the difference between discrete or continuous variables?

What happens when you combine multiple aesthetics?

	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	Shouldn'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

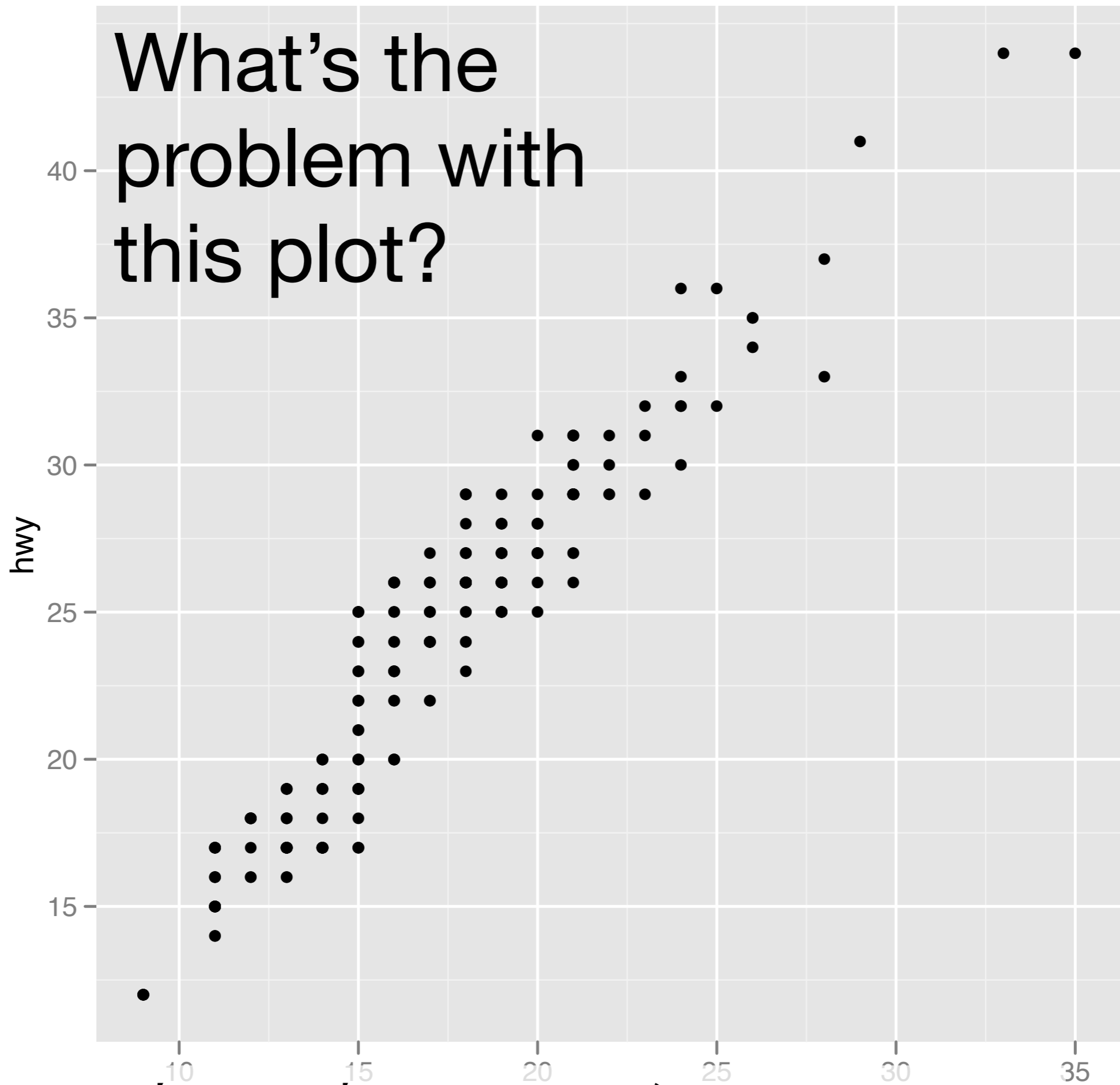
`facet_wrap()`: 1d ribbon wrapped into 2d

Aside: workflow

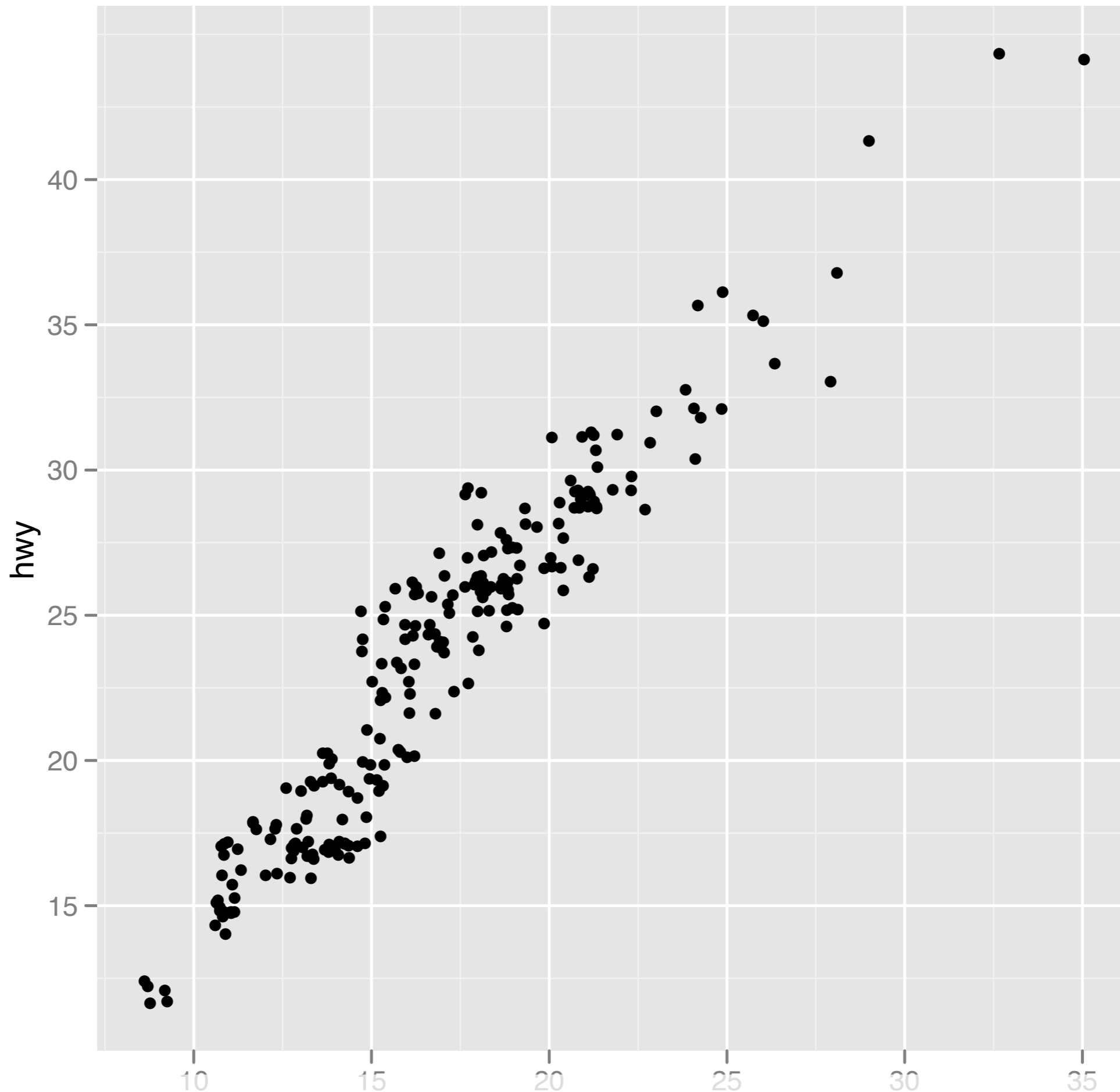
Keep a copy of the slides open so that you can copy and paste the code.

For complicated commands, write them in the editing area and then copy and paste.

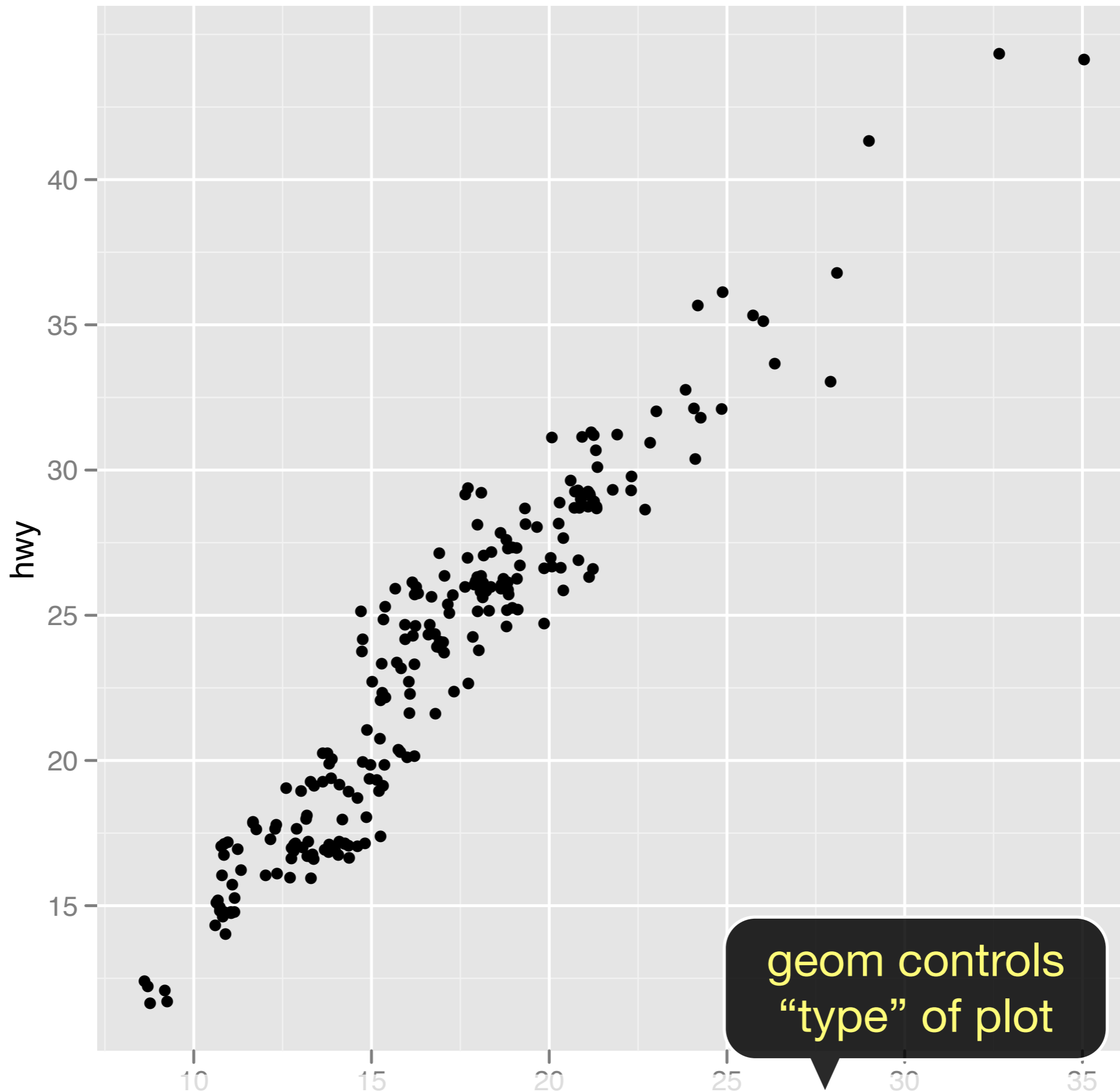
What's the
problem with
this plot?



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

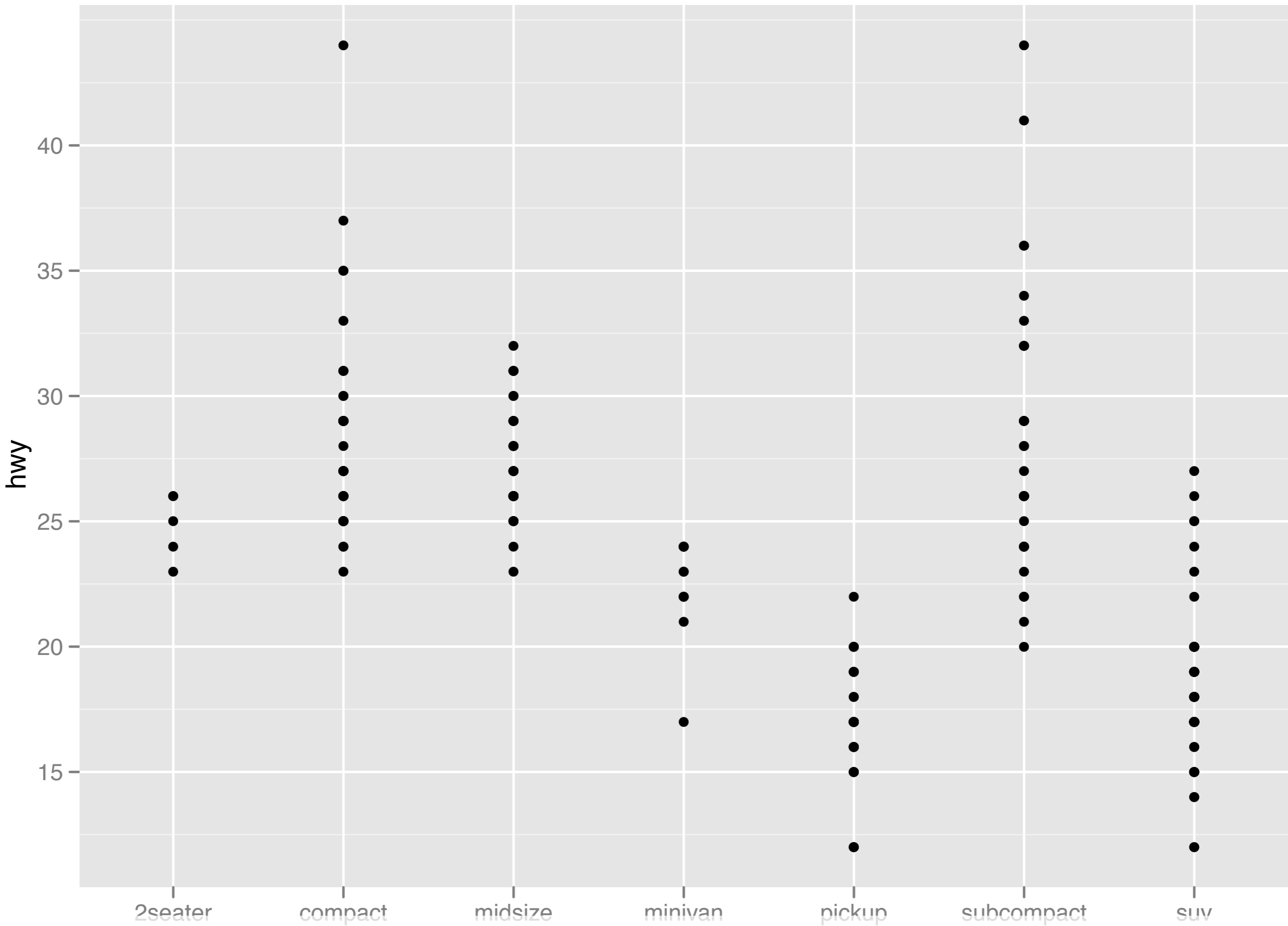



```
qplot(cty, hwy, data = mpg, geom = "jitter")
```



geom controls
"type" of plot

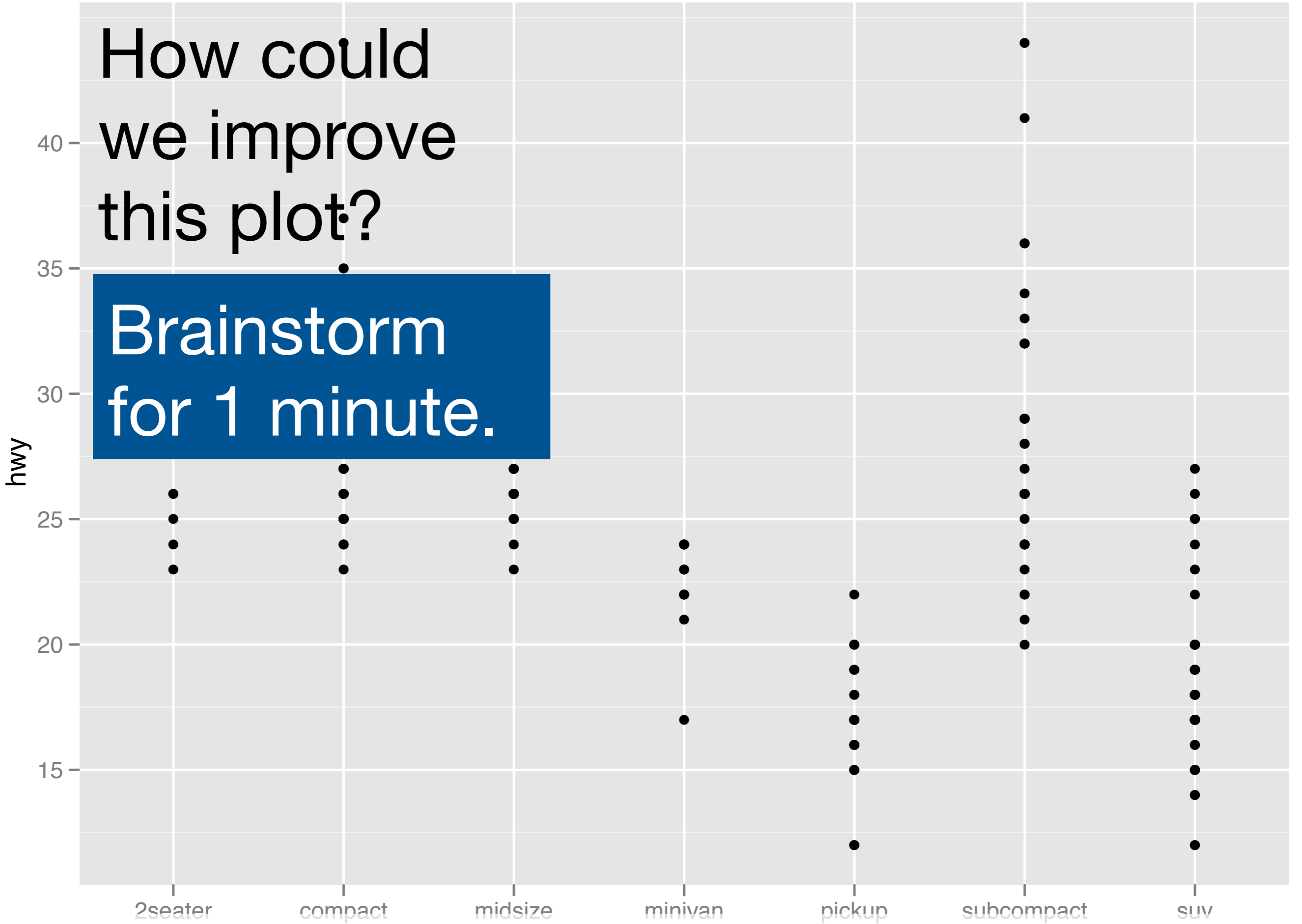
```
qplot(cty, hwy, data = mpg, geom = "jitter")
```



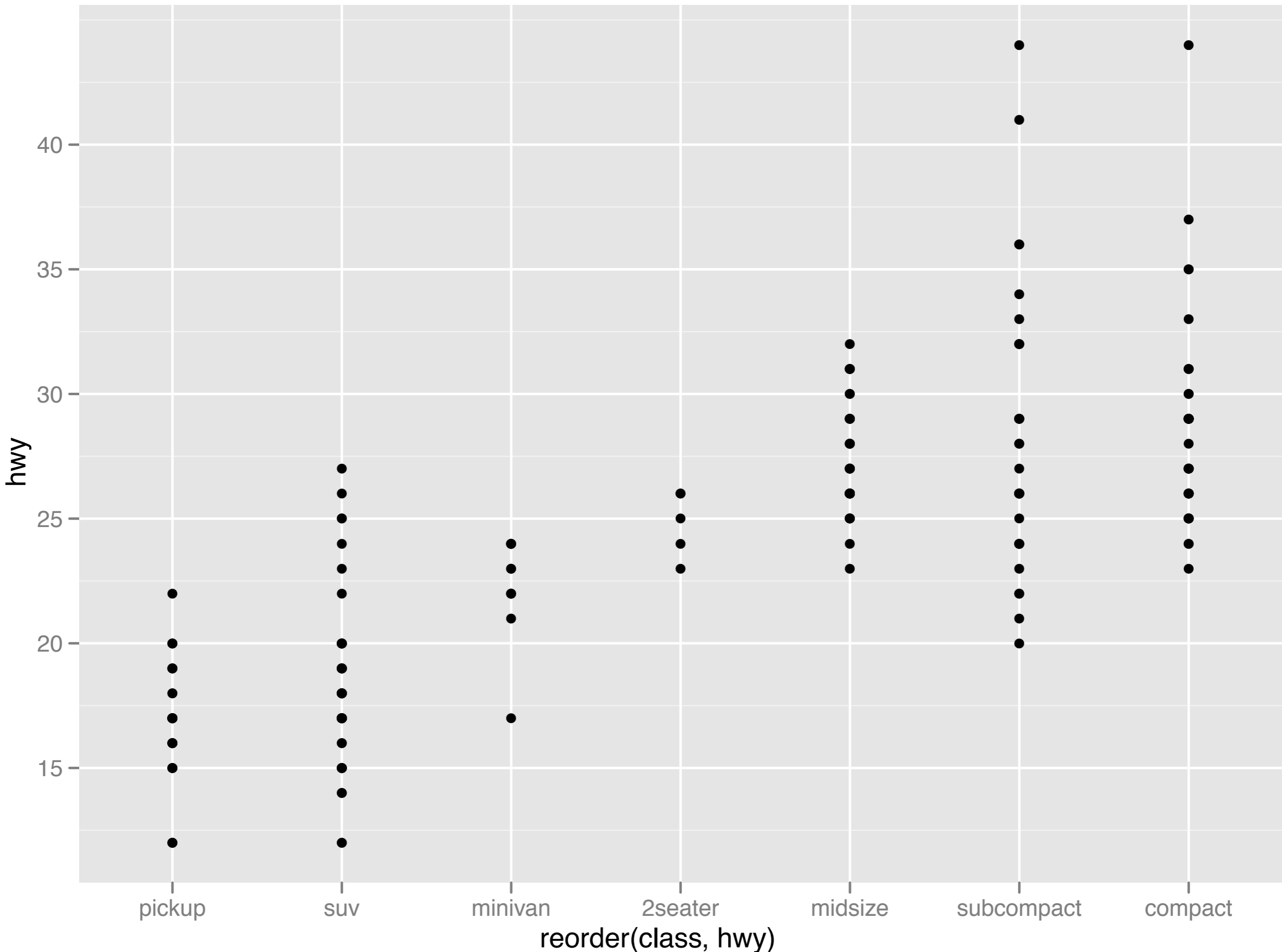
`qplot(class, hwy, data = mpg)`

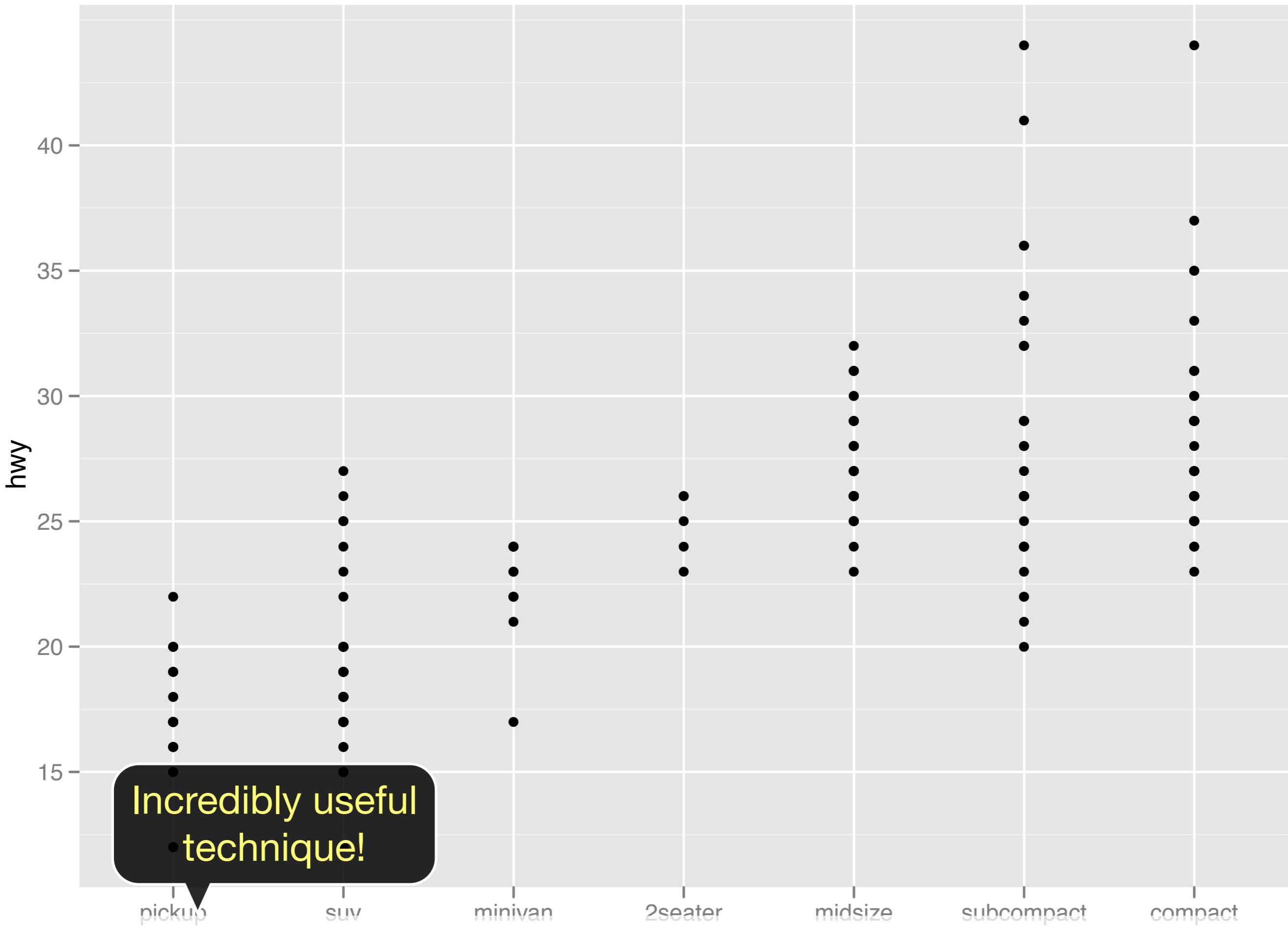
How could we improve this plot?

Brainstorm for 1 minute.



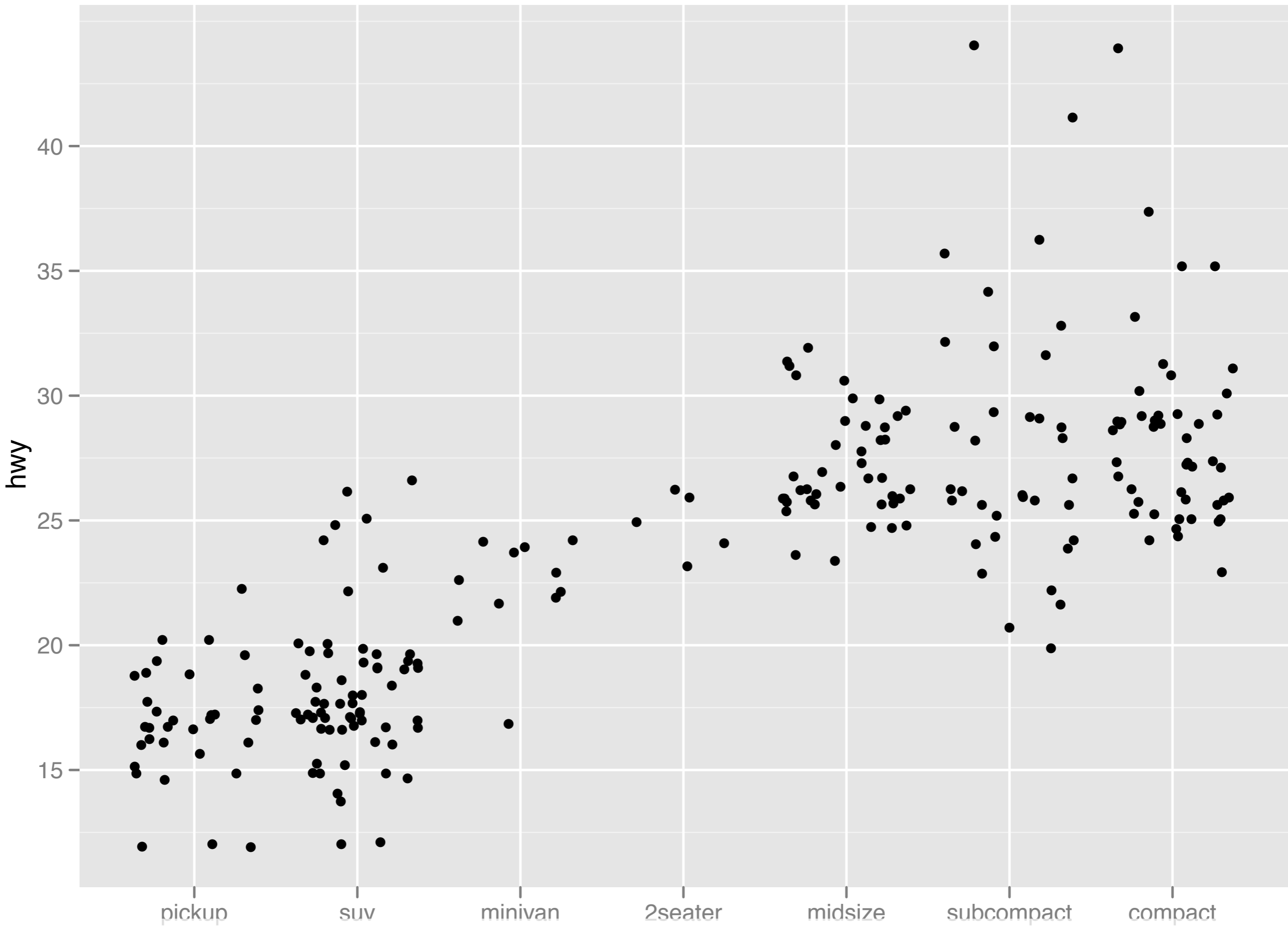
```
qplot(class, hwy, data = mpg)
```



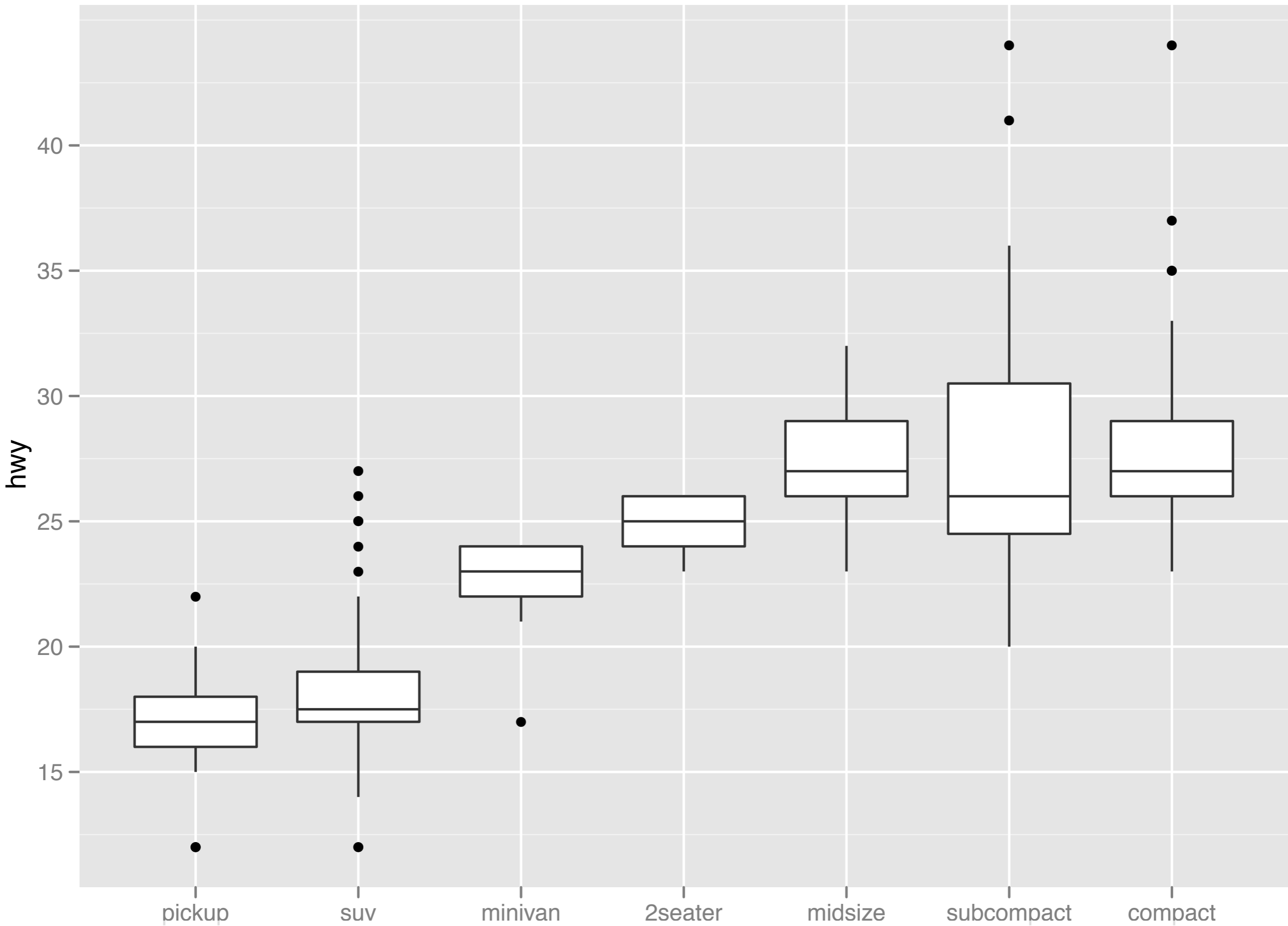


Incredibly useful
• technique!

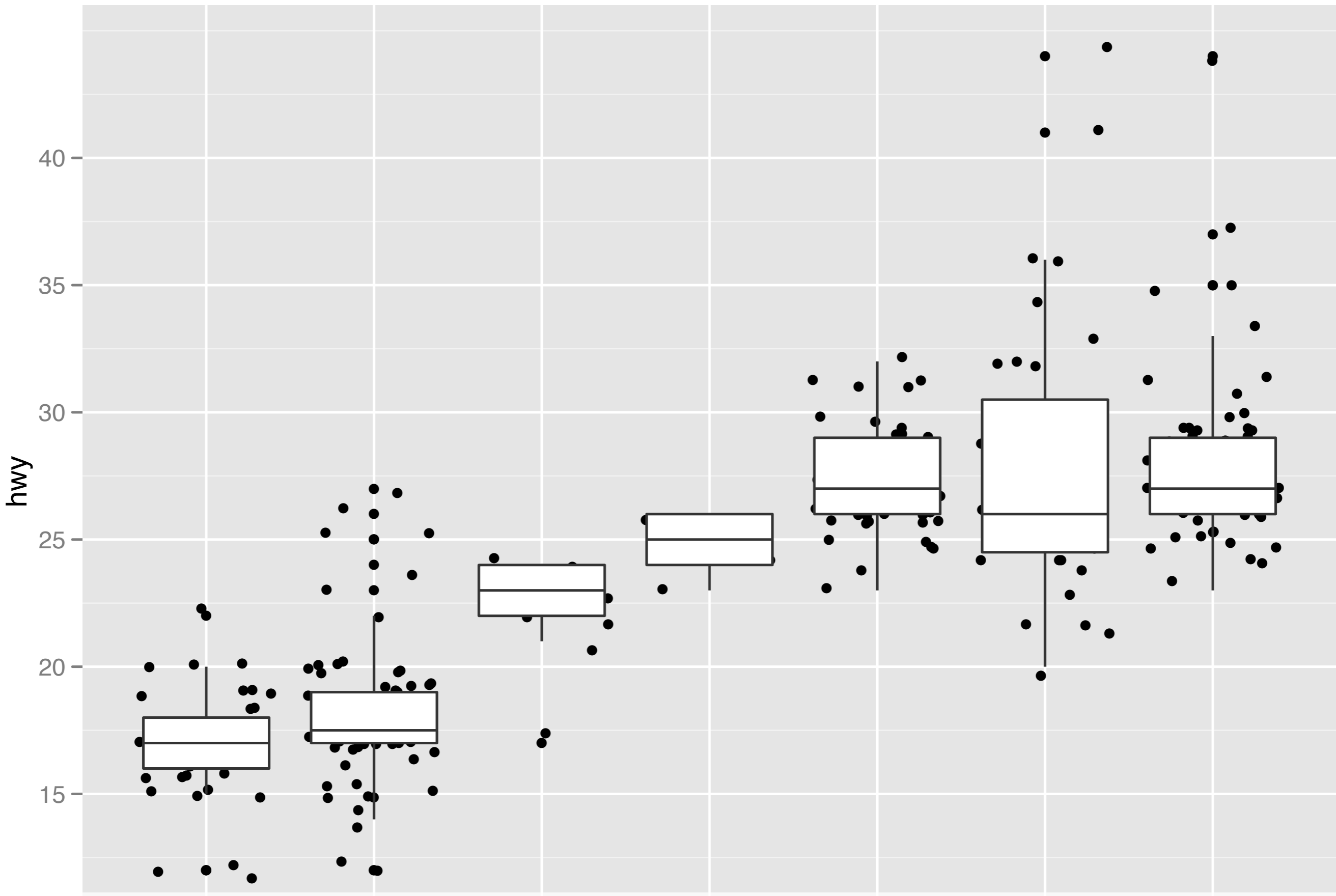
```
qplot(reorder(class, hwy), hwy, data = mpg)
```



```
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?

Aside: coding strategy

At the end of each interactive session, you want a summary of everything you did. Two options:

1. Copy from the history panel.
2. Build up the important bits as you go.
(recommended)