

# Visualization in R

**Hadley Wickham**

Assistant Professor / Dobelman Family Junior Cha  
Department of Statistics / Rice University

1. Scatterplot basics
2. Adding extra variables
3. Jittering and boxplots
4. Histograms and barcharts
5. Learn more



Learning a new  
language is hard!

# Scatterplot basics

```
install.packages("ggplot2")  
library(ggplot2)  
mpg  
head(mpg)  
str(mpg)  
summary(mpg)  
ggplot(displ, hwy, data = mpg)
```

# Scatterplot basics

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

```
library(ggplot2)
```

```
mpg
```

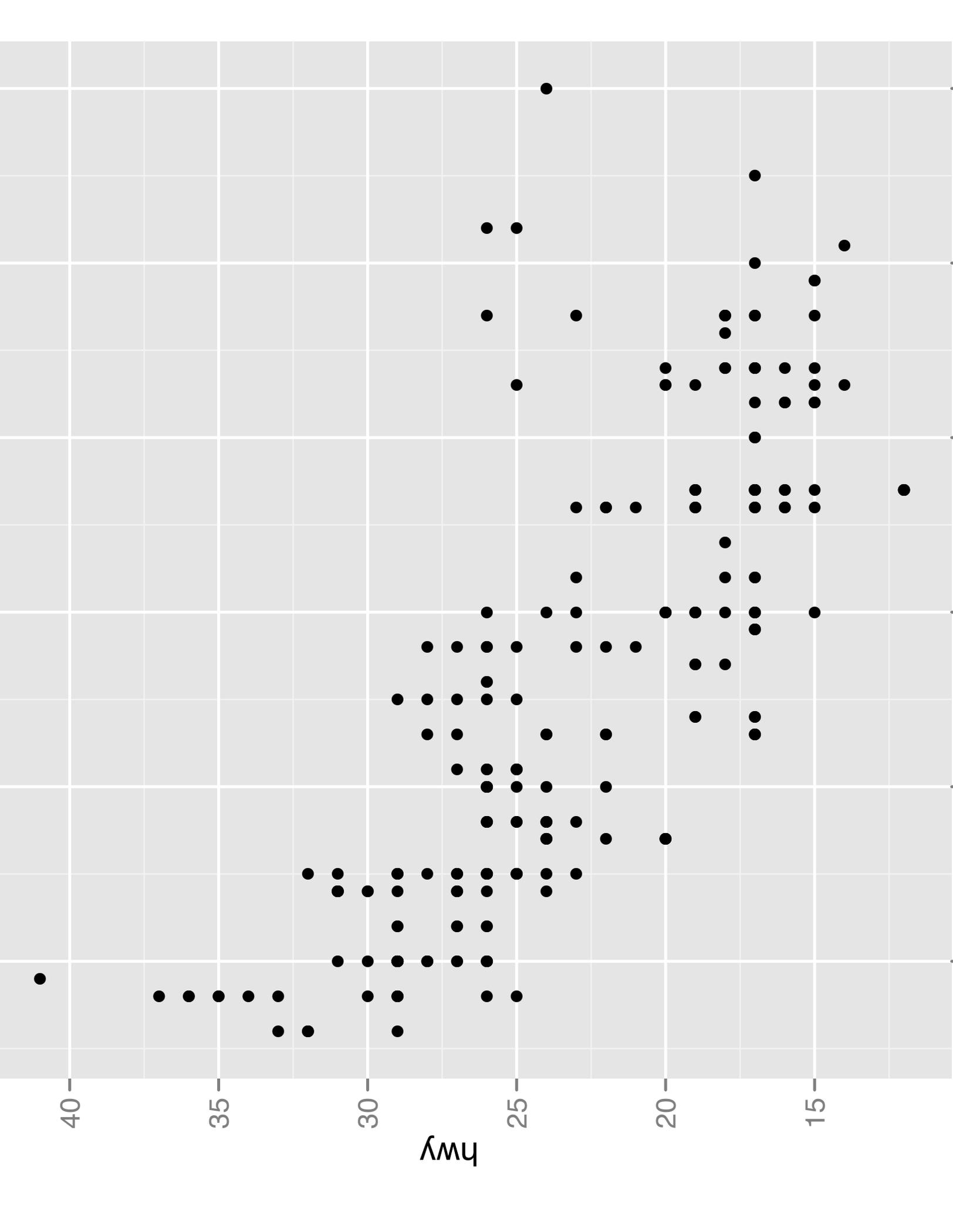
```
head(mpg)
```

```
str(mpg)
```

```
summary(mpg)
```

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

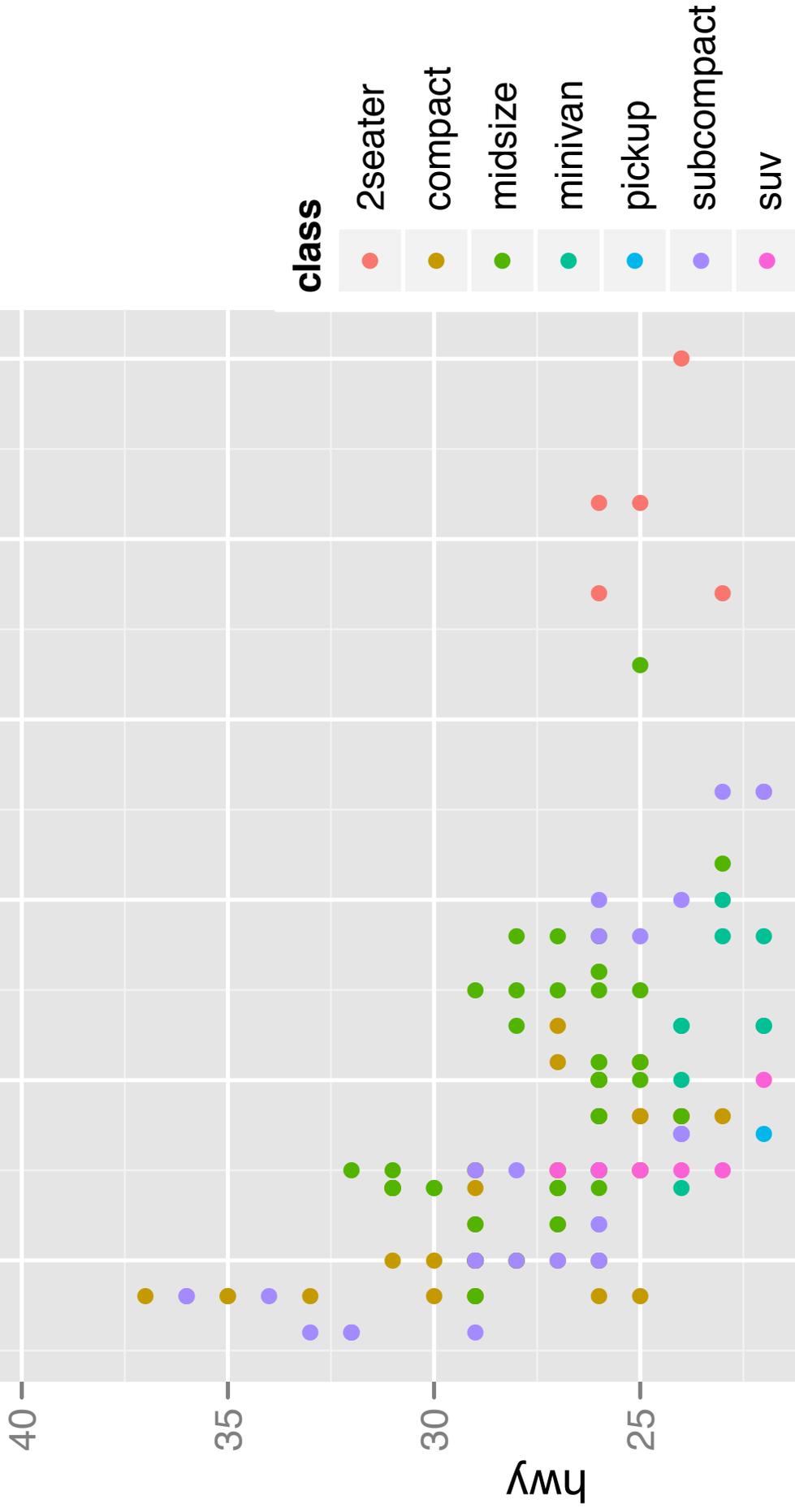
In ggplot2, we  
always explicitly  
specify the data



# Additional variables

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





Legend chosen displayed automa

# Your turn

Experiment with colour, size, and shape aesthetics.

What's the difference between discrete and continuous variables?

What happens when you combine multiple aesthetics?

Discrete

Continuous

Rainbow of  
colours

Gradient from  
red to blue

Discrete size  
steps

Linear mapp  
between rad  
and value

Different shape  
for each

Doesn't wo

colour

size

shape

# 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)
```

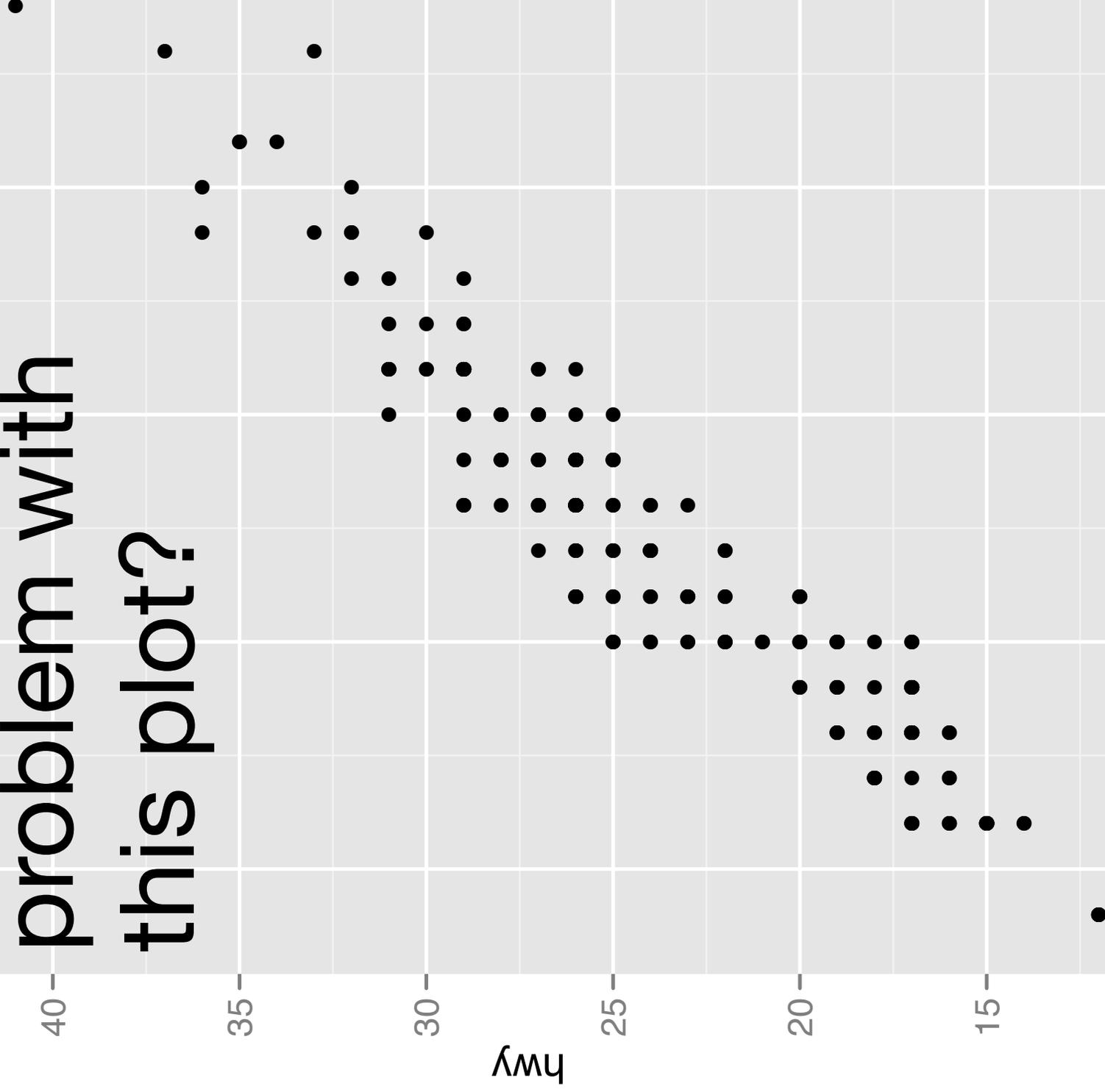
# Summary

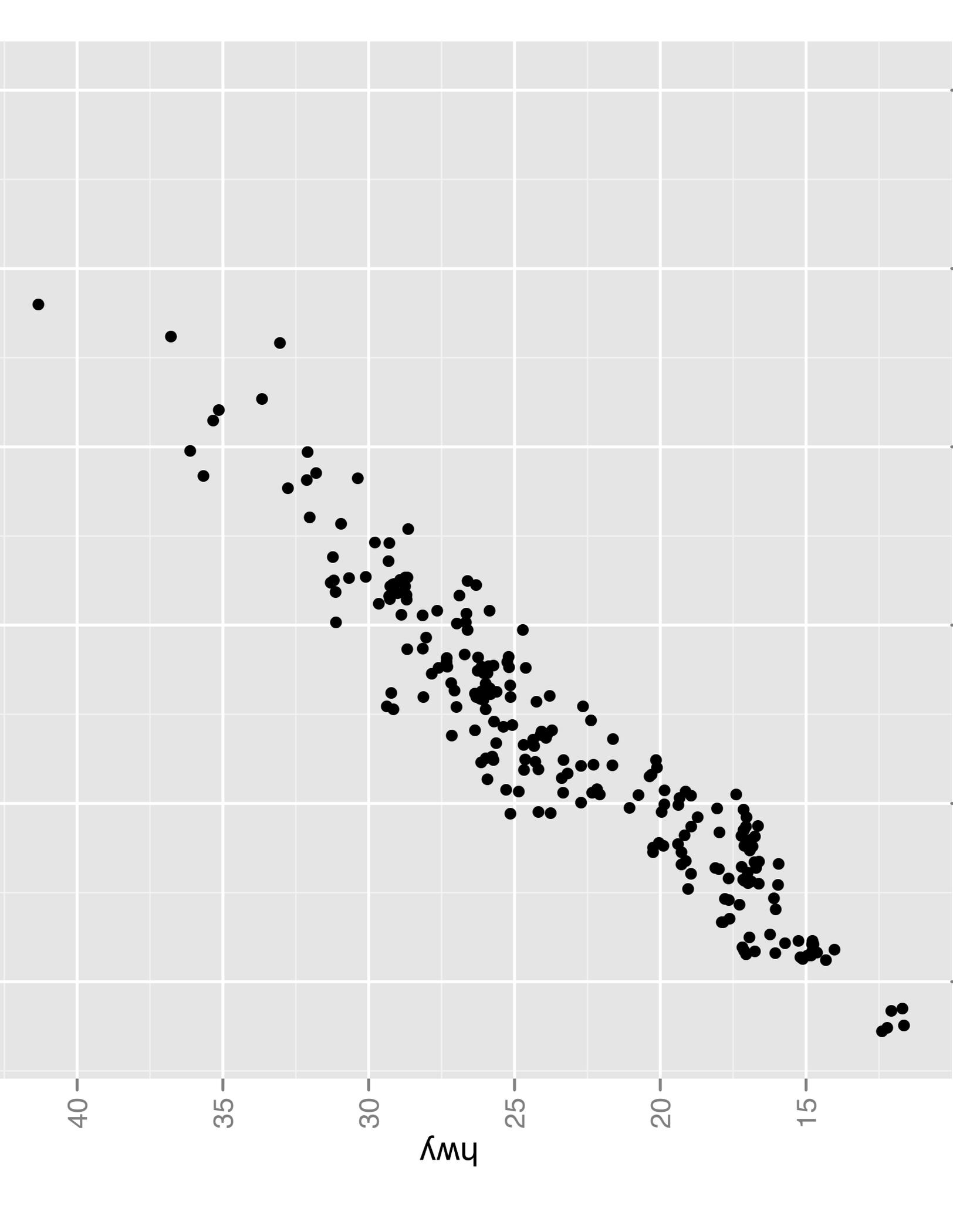
`facet_grid(): 2d grid, rows ~ cols, . for  
no split`

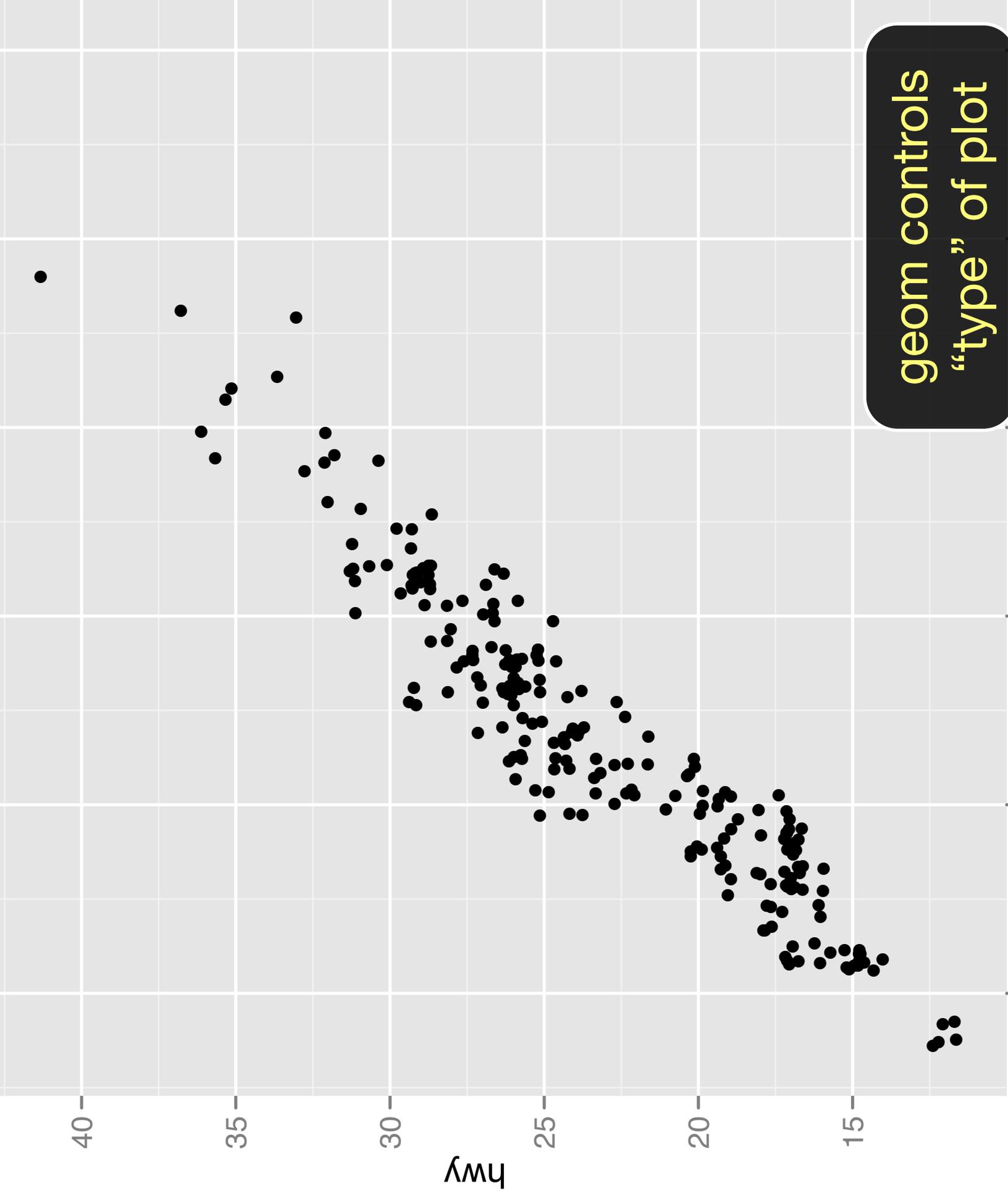
`facet_wrap(): 1d ribbon wrapped into 2d`

Scales argument controls whether  
position scales are fixed or free.

problem with  
this plot?





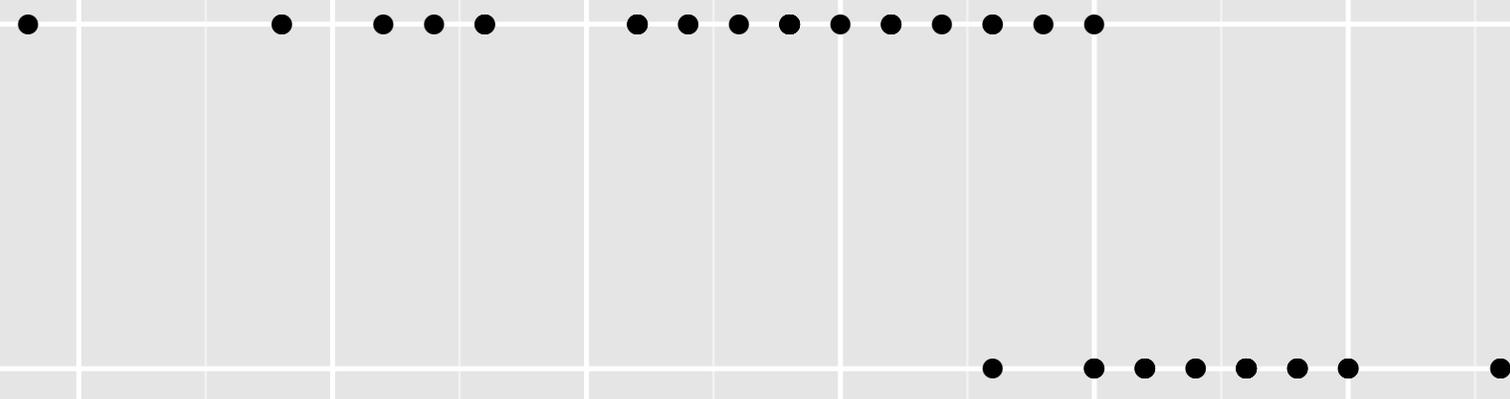


geom controls  
"type" of plot

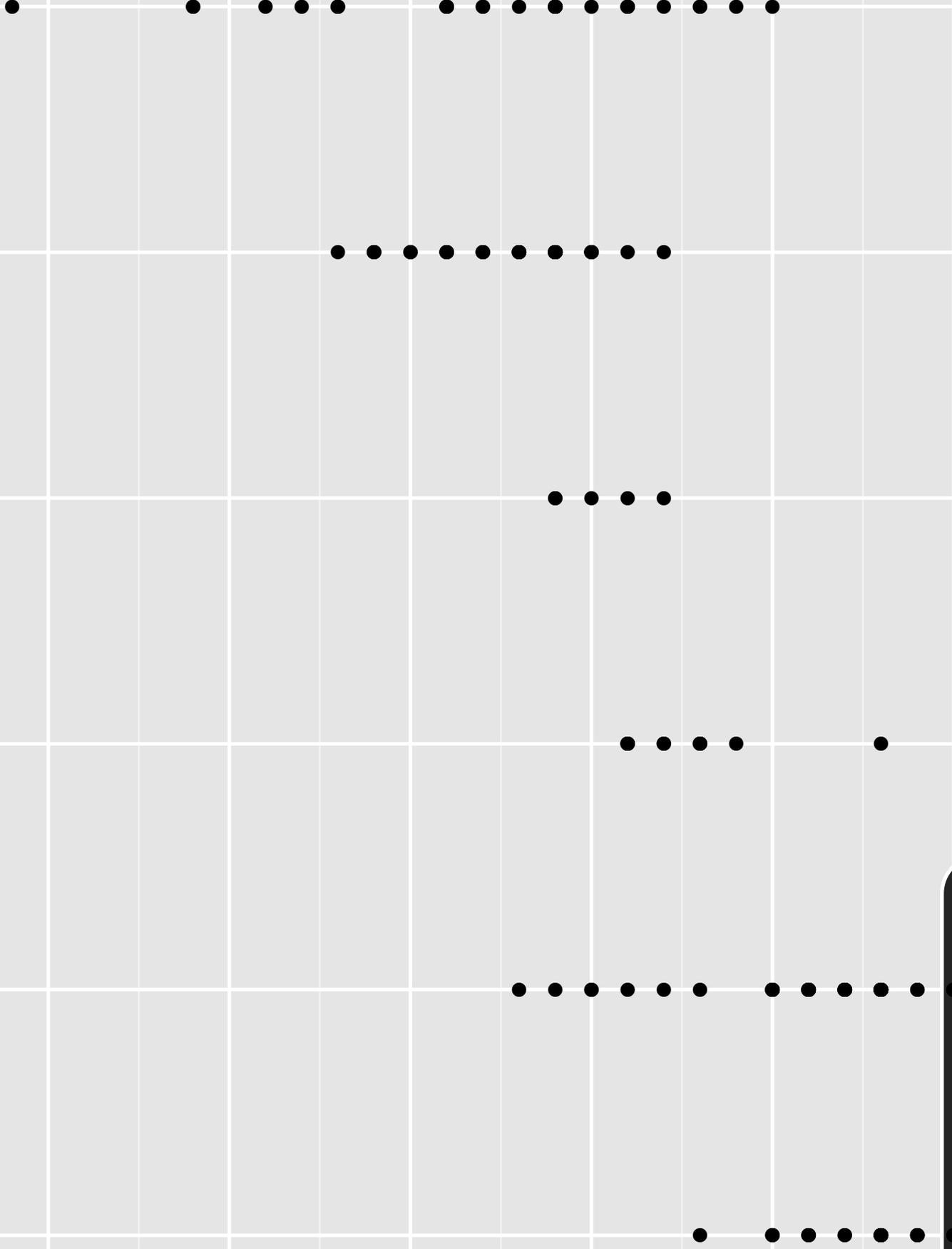


How can we improve this plot?

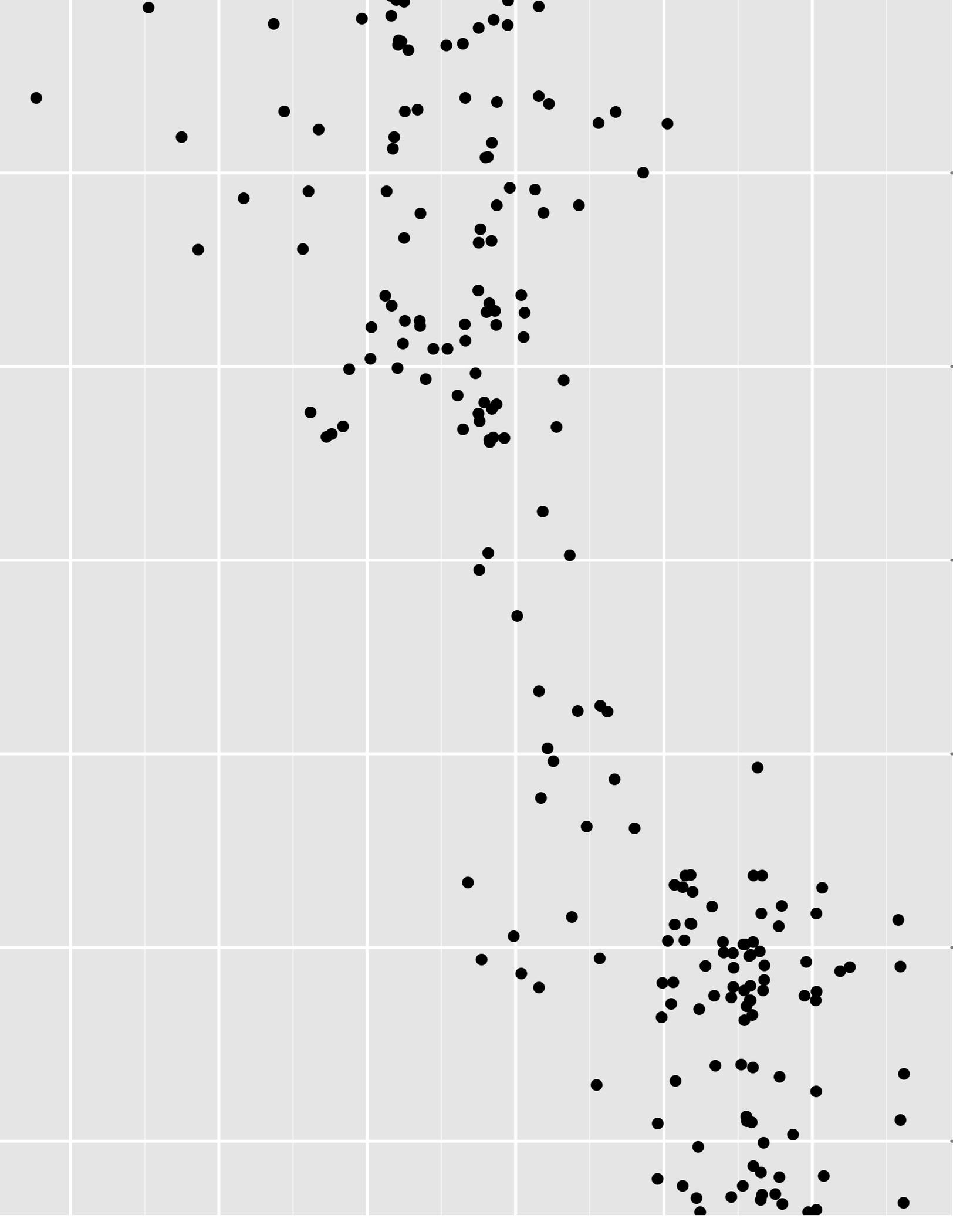
Brainstorm for 1 minute.

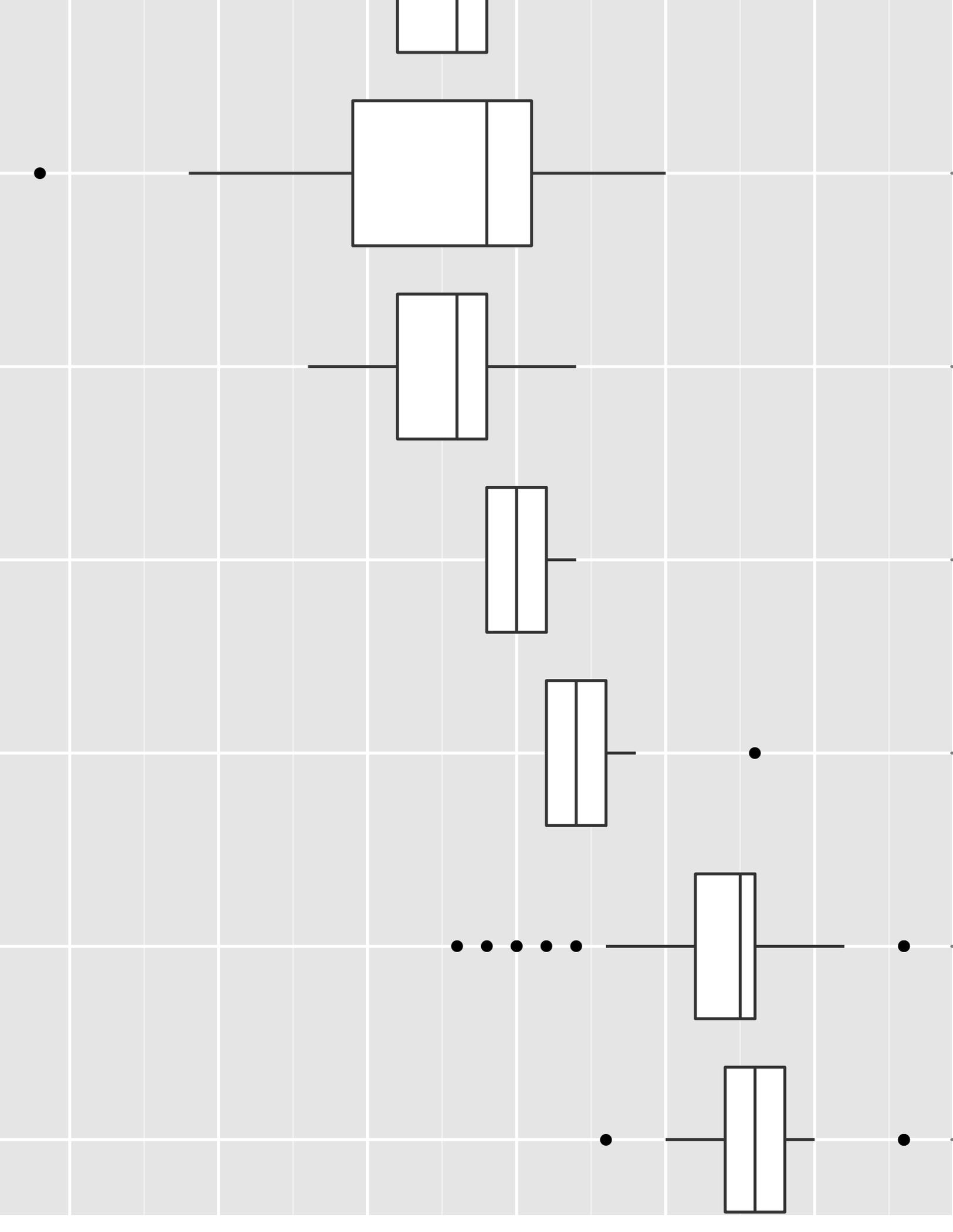


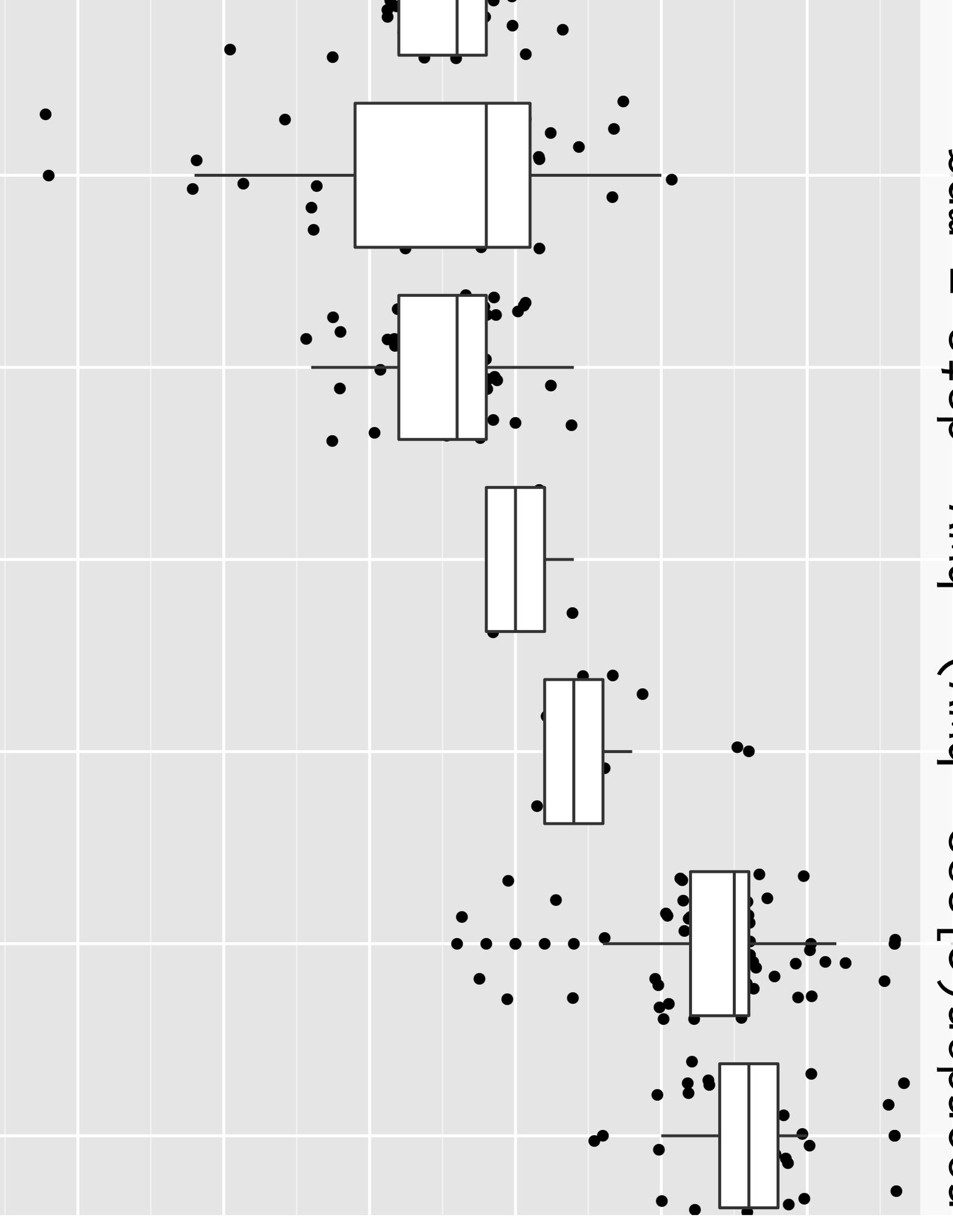




**Incredibly useful**  
• technique!







# Your turn

Read the help for `reorder`. Redraw the previously 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. Save everything you did with `savehistory()` then remove the unimportant bits.
2. Build up the important bits as you go. (this is how I work)

# Diamonds data

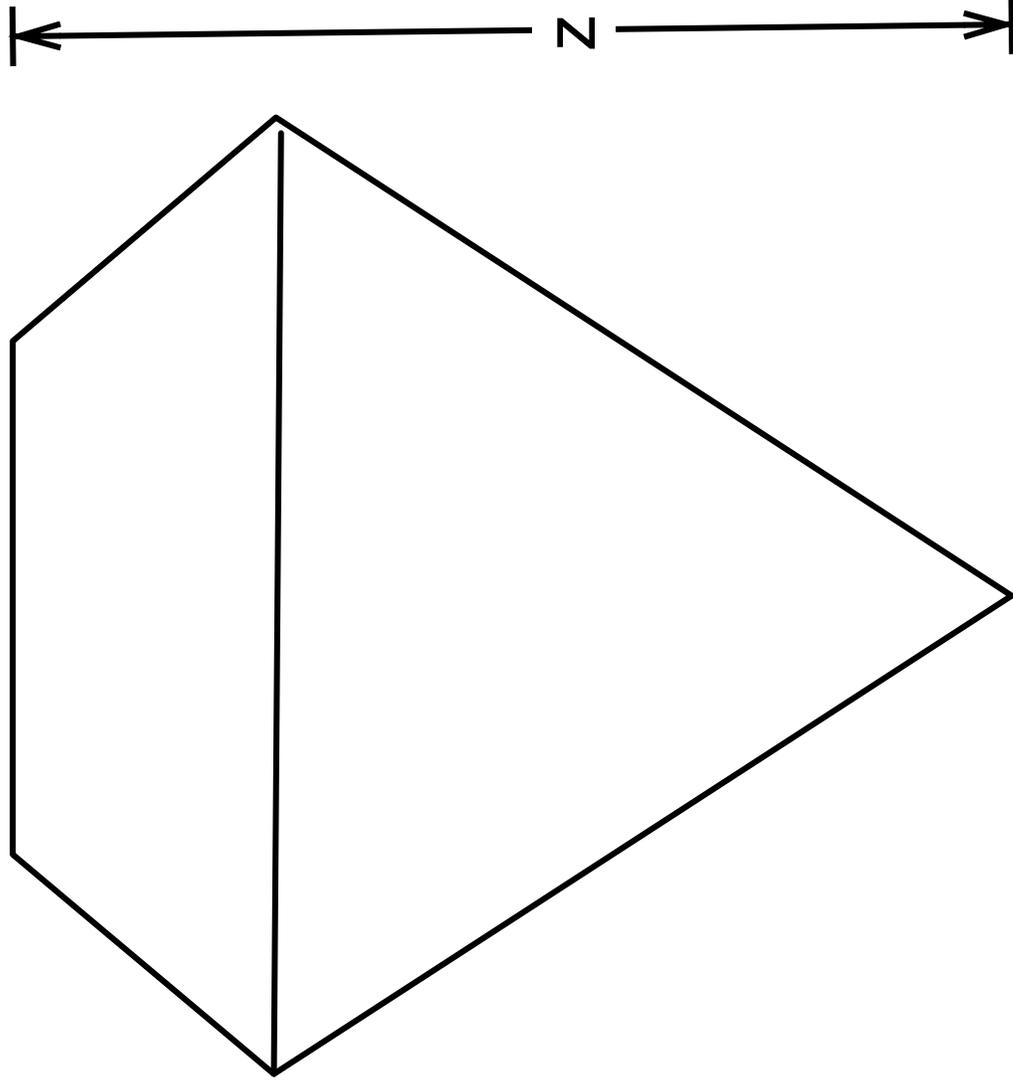
~54,000 round diamonds from  
<http://www.diamondse.info/>

Carat, colour, clarity, cut  
Total depth, table, depth,  
width, height

Price



↔ table width →



depth = z / diameter  
table = table width / y \* 100

# Histogram & bar charts

Histograms and

# bar charts

Used to display the distribution of a variable

Categorical variable → bar chart

Continuous variable → histogram

Always

xperiment with

the bin width

# Examples

```
# With only one variable, qplot guesses that  
# you want a bar chart or histogram
```

```
qplot(cut, data = diamonds)
```

```
qplot(carat, data = diamonds)
```

```
qplot(carat, data = diamonds, binwidth = 1)
```

```
qplot(carat, data = diamonds, binwidth = 0.1)
```

```
qplot(carat, data = diamonds, binwidth = 0.01)
```

```
resolution(diamonds$carat)
```

```
last_plot() + xlim(0, 3)
```

# Examples

```
# With only one variable, qplot guesses that  
# you want a bar chart or histogram  
qplot(cut, data = diamonds)
```

```
qplot(carat, data = diamonds)  
qplot(carat, data = diamonds, binwidth = 1)  
qplot(carat, data = diamonds, binwidth = 0.1)  
qplot(carat, data = diamonds, binwidth = 0.01)  
resolution
```

Common ggplot2  
technique: adding  
together plot  
components

```
geom_histogram() + xlim(0, 3)
```

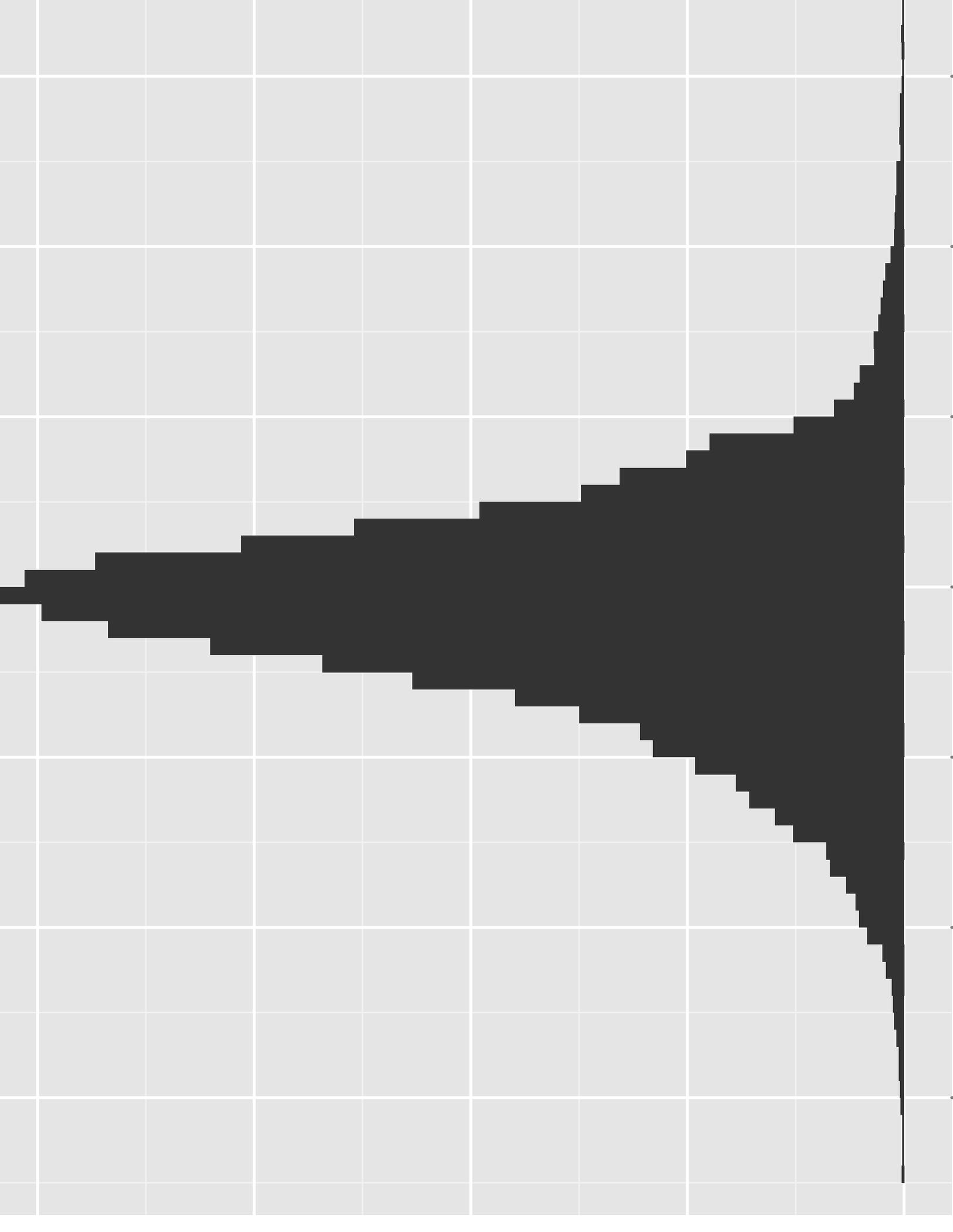
```
t(table, data = diamonds, binwidth = 1)
zoom in on a plot region use xlim() and ylim()
t(table, data = diamonds, binwidth = 1) +
  xlim(50, 70)
t(table, data = diamonds, binwidth = 0.1) +
  xlim(50, 70)
t(table, data = diamonds, binwidth = 0.1) +
  xlim(50, 70) + ylim(0, 50)
```

te that this type of zooming discards data  
ide of the plot regions  
e coord\_cartesian() for an alternative

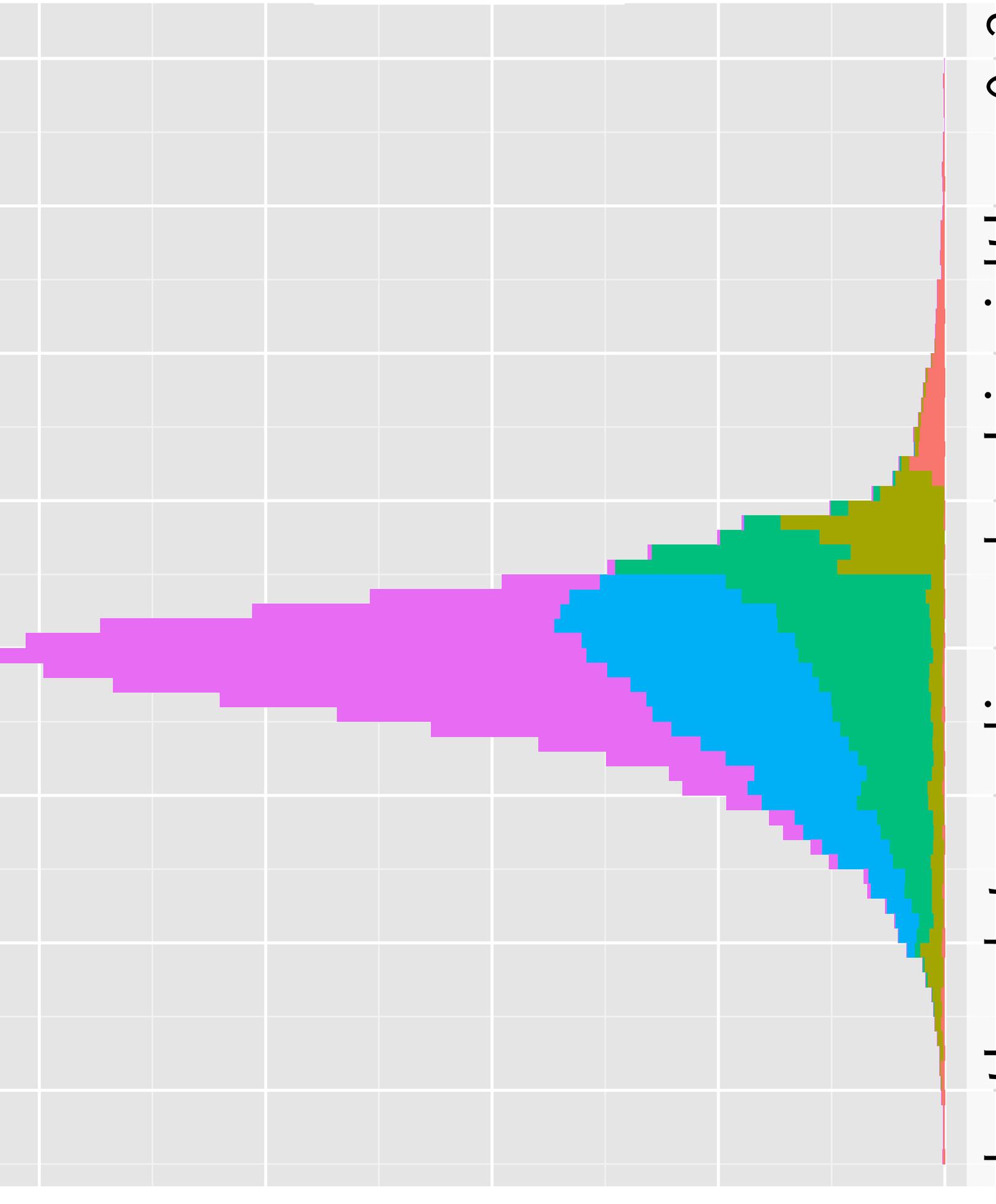
# Additional variables

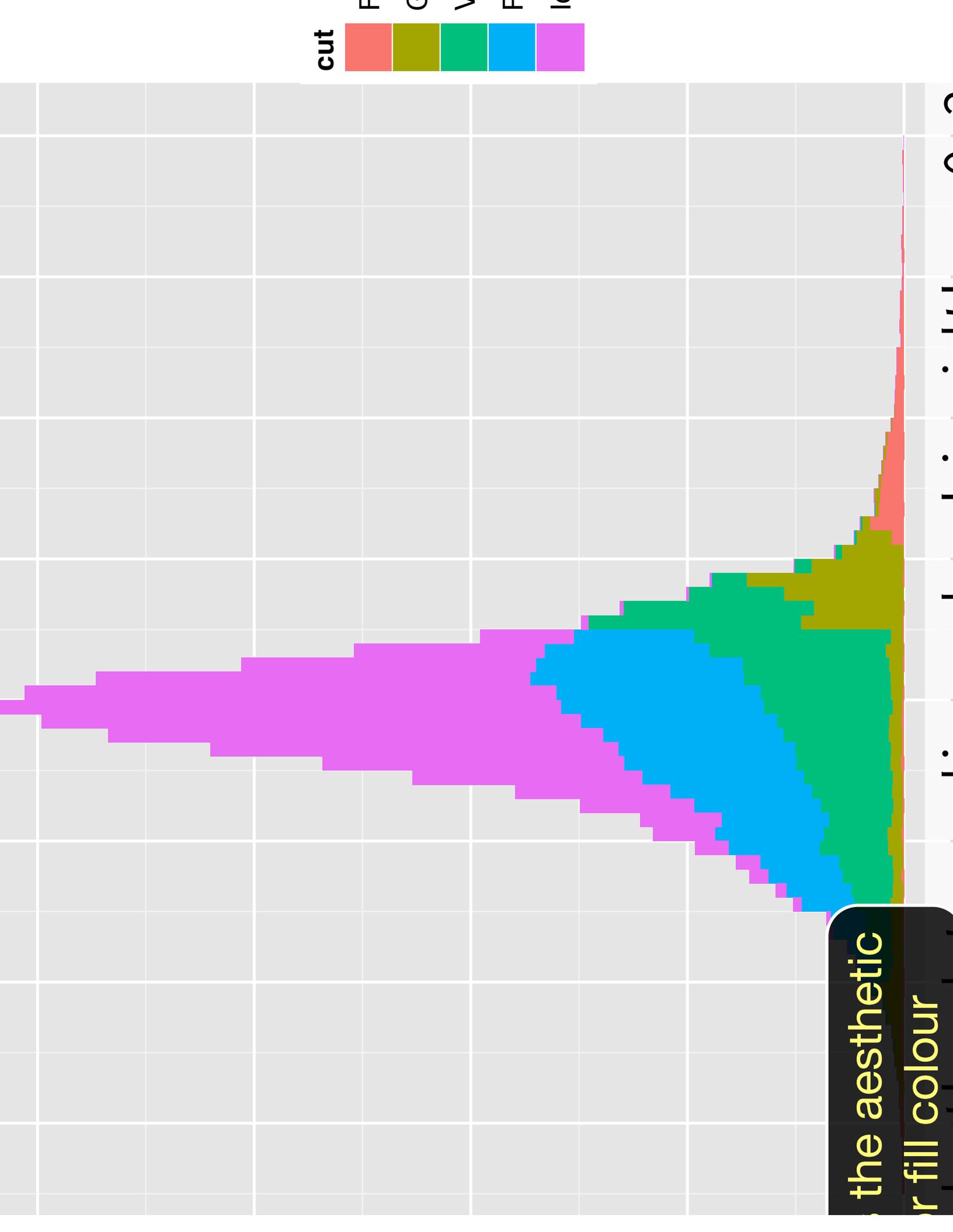
As with scatterplots can use aesthetics or **faceting**. Using aesthetics creates pretty, but ineffective, plots.

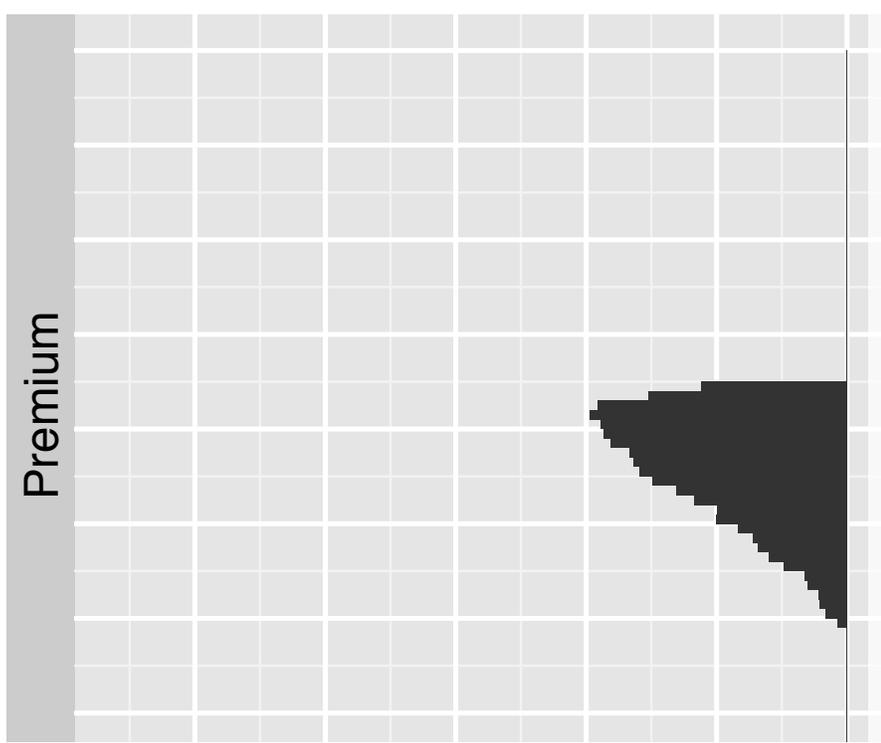
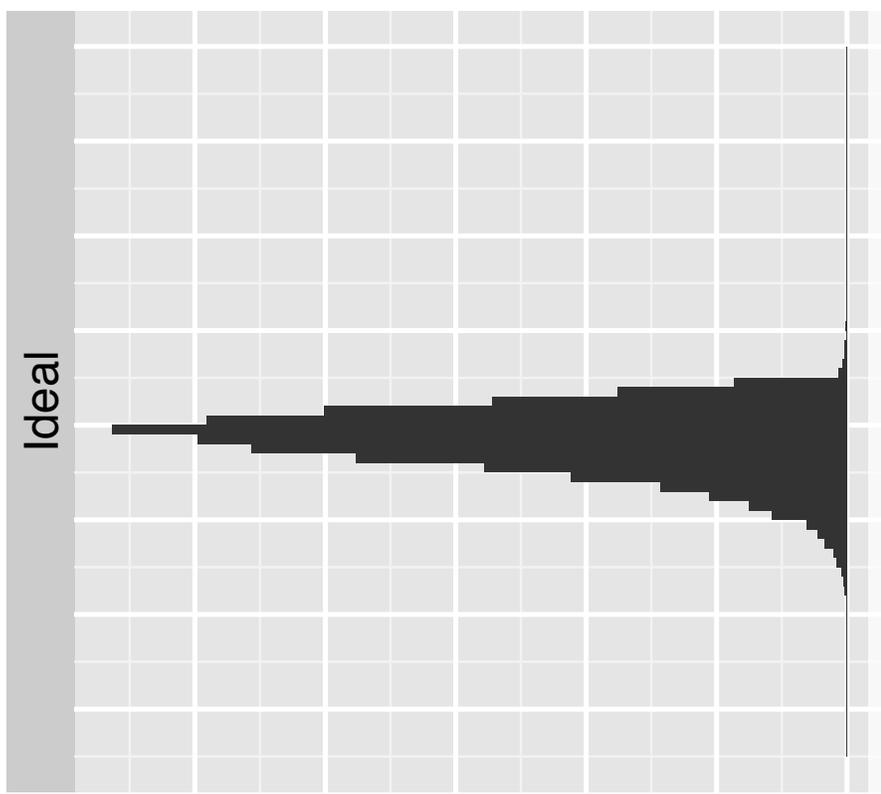
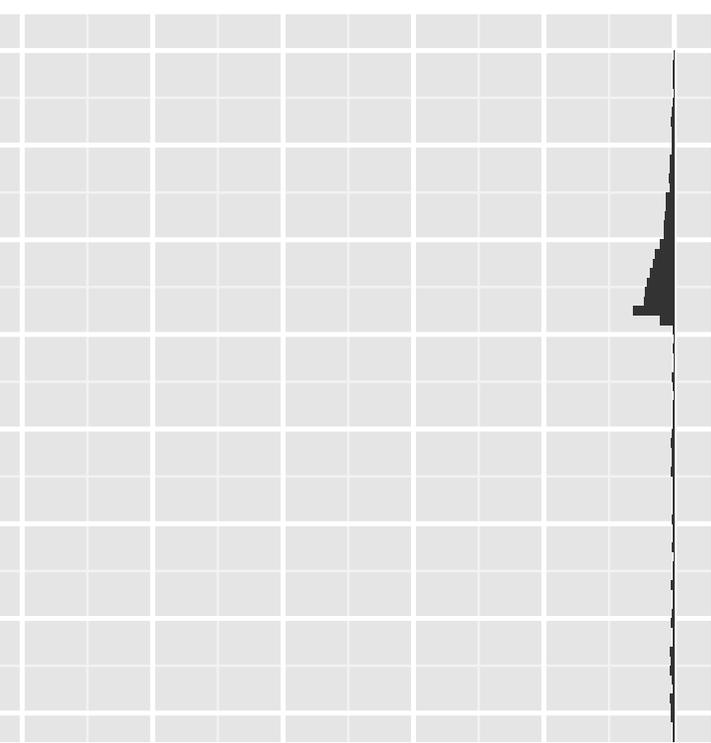
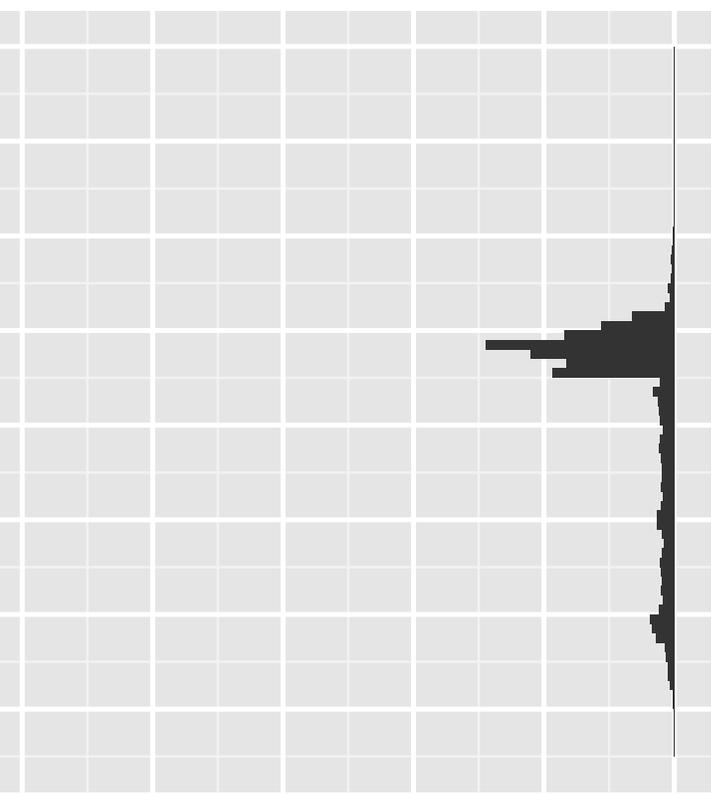
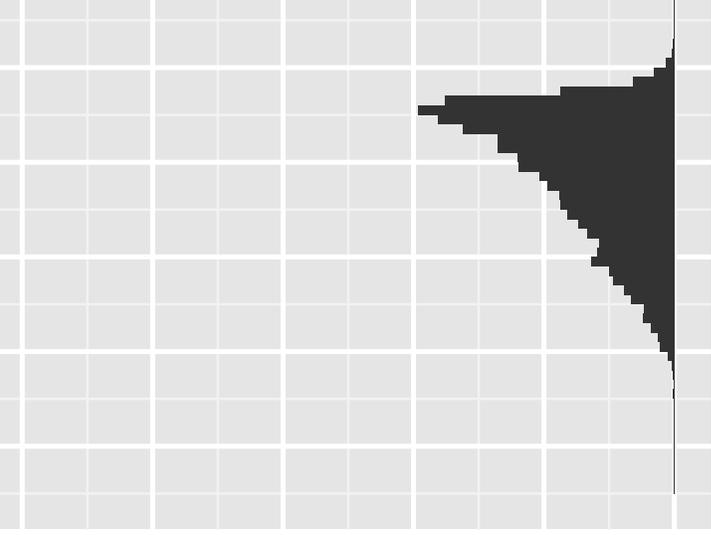
The following examples show the difference, when investigation the relationship between cut and depth.



cut F C V F K







# Your turn

Explore the distribution of price.

How does it vary with colour, or cut, and clarity?

Learn more

# About ggplot2

Graphical grammar (domain specific language), based on “The Grammar of Graphics” by Leland Wilkinson.

Specify what you want, not how to create it.  
Many fiddly details taken care of.

“Instead of spending time making your graph look pretty, you can focus on creating a graph that bests reveals the messages in your data.”

# Useful resources

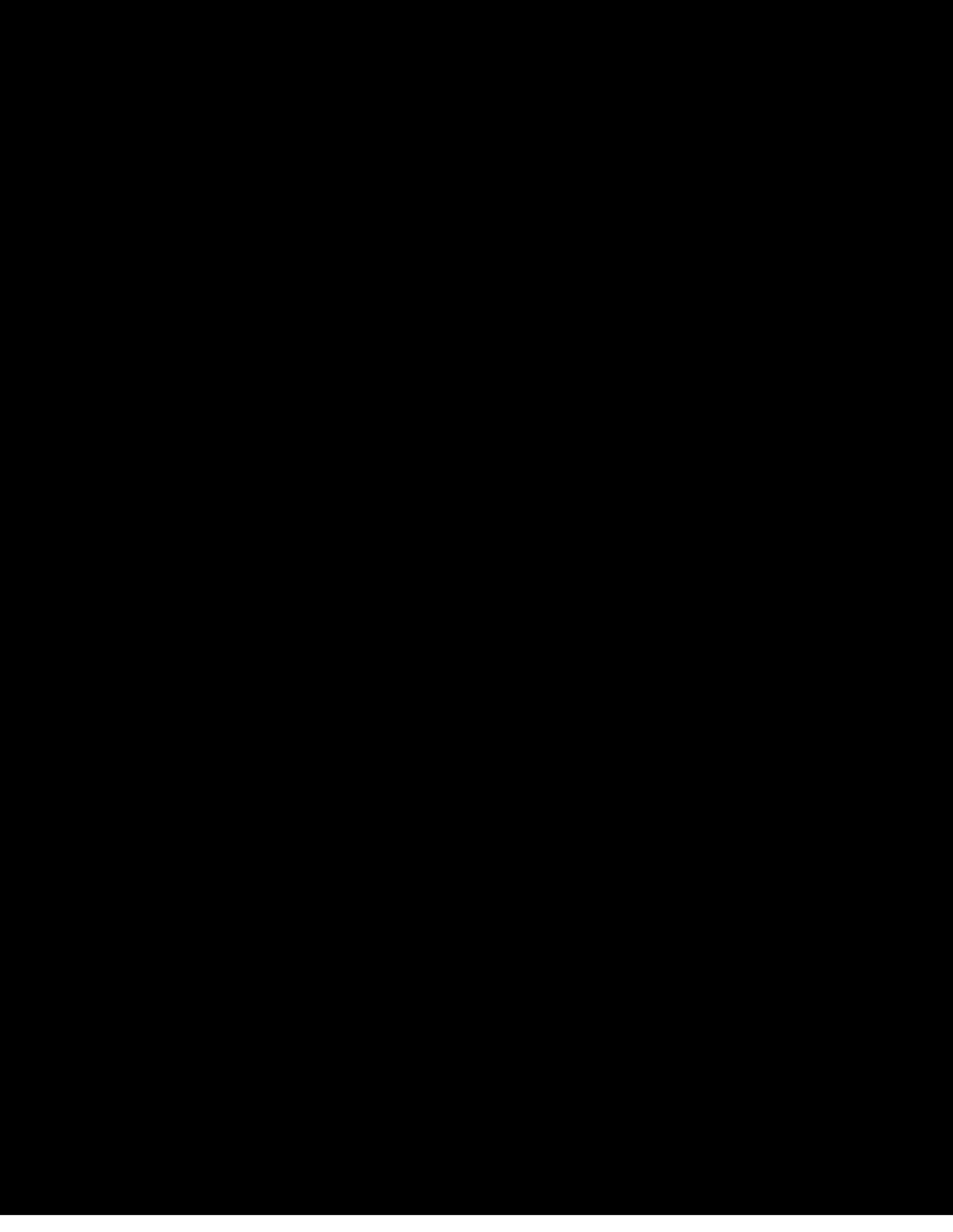
<http://had.co.nz/ggplot2>

<http://had.co.nz/ggplot2/book>

<http://groups.google.com/group/ggplot2>

<http://learnr.wordpress.com>

<http://ggplot2.wik.is>



work is licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License. To view a copy of this license visit <http://creativecommons.org/licenses/by-nc/3.0/> or send a letter to Creative Commons, 530 Second Street, Suite 300, San Francisco, California, 94105, USA.