



VISUALIZATION BEST PRACTICES IN R

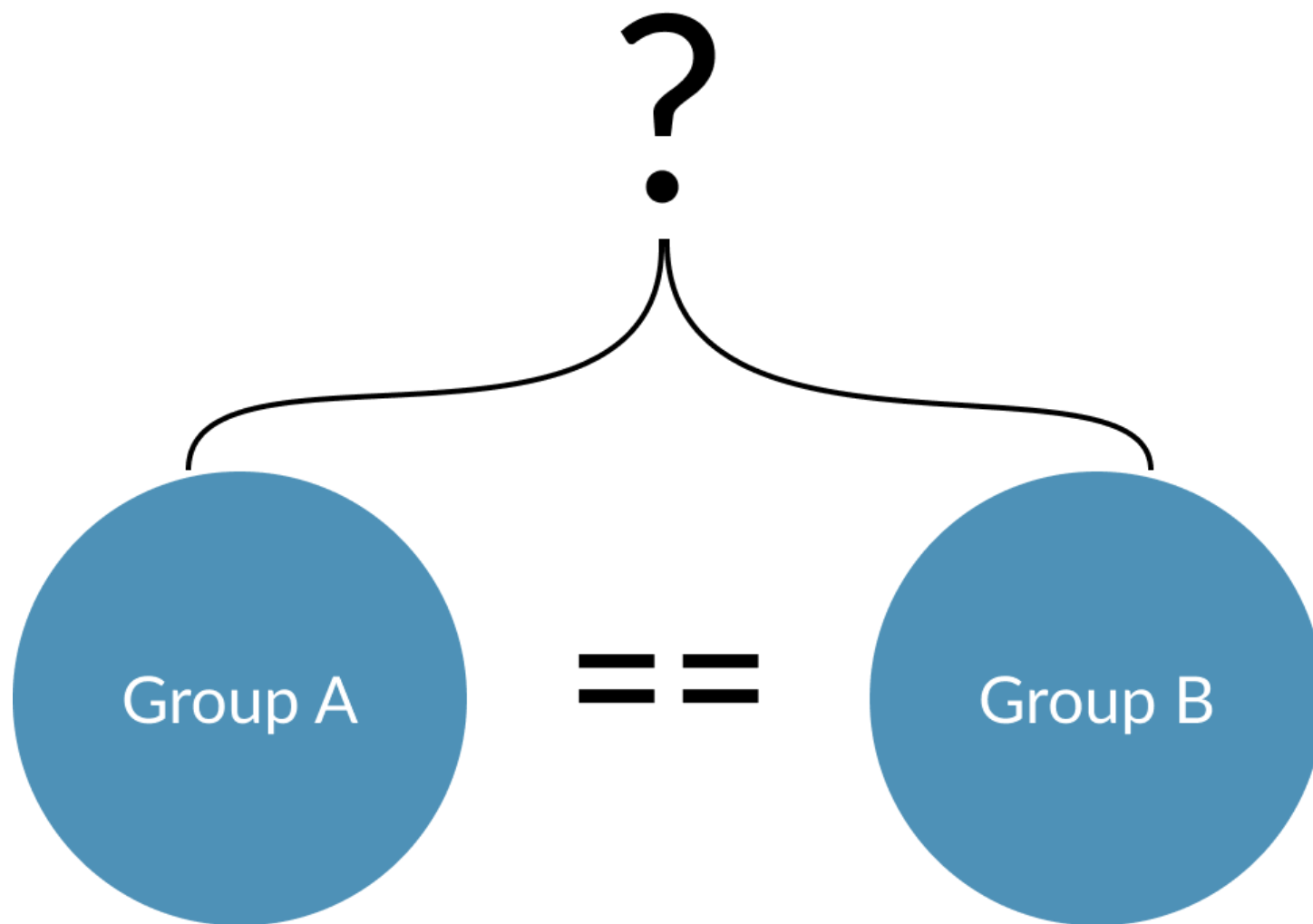
Comparing Distributions

Nick Strayer
Instructor



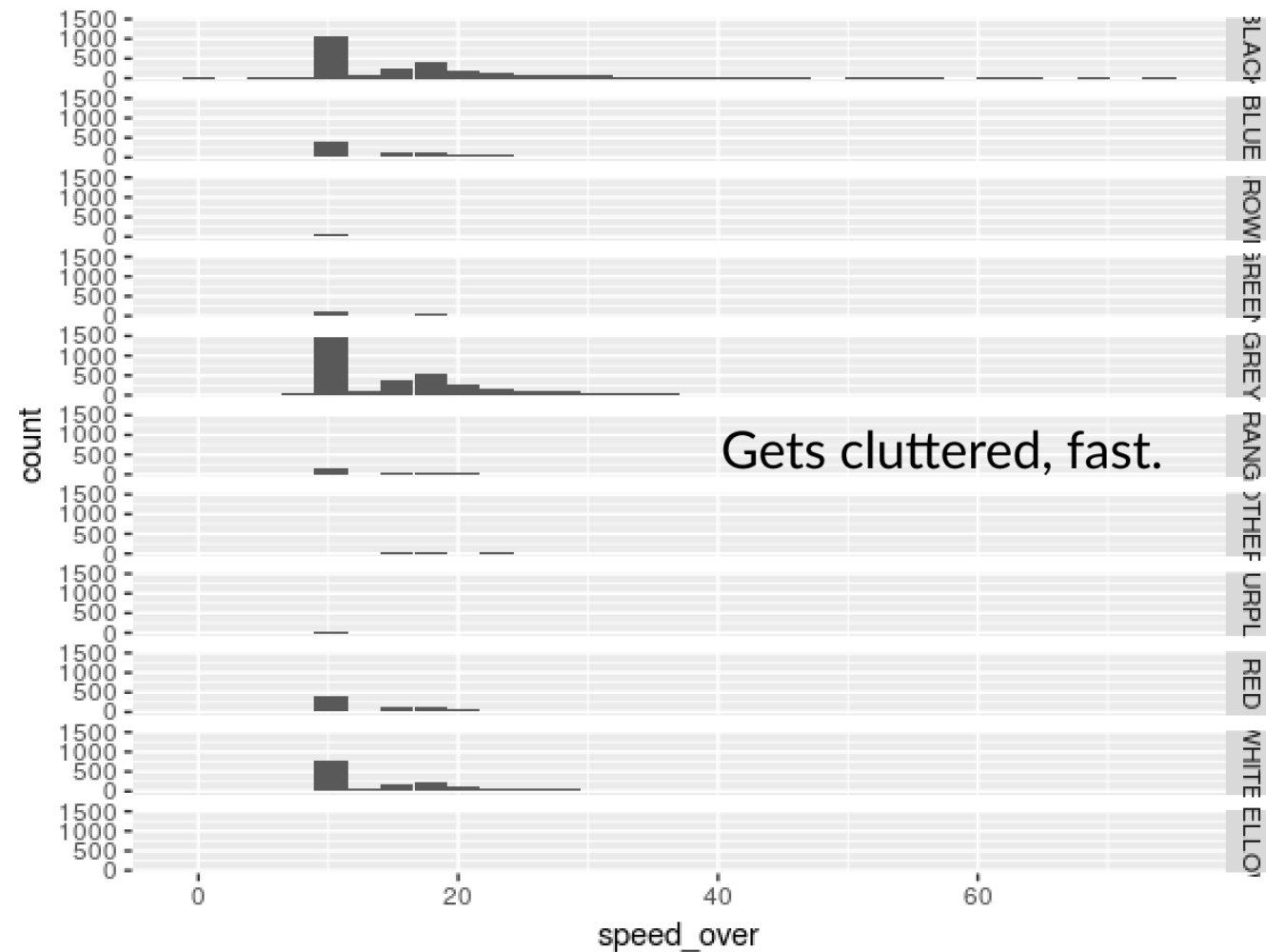
Why compare distributions?

- Verify balanced groups
- For comparisons sake

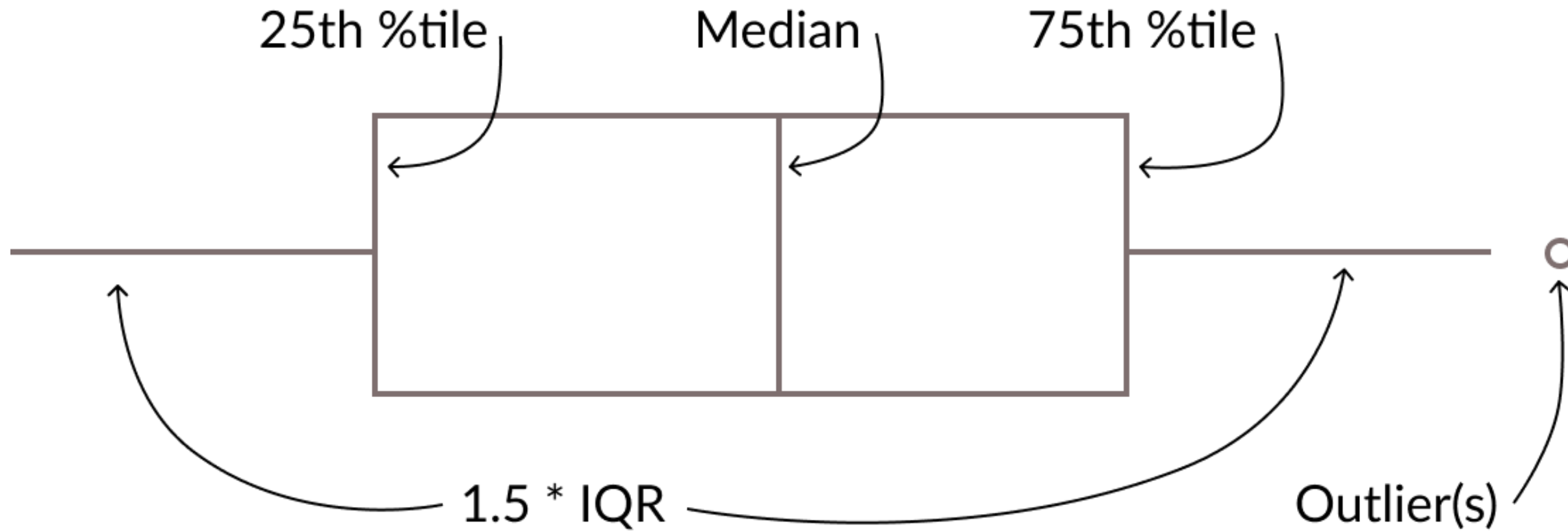


Why not facet histograms?

```
ggplot(md_speeding, aes(x = speed_over)) +  
  geom_histogram() +  
  facet_grid(vehicle_color~.)
```



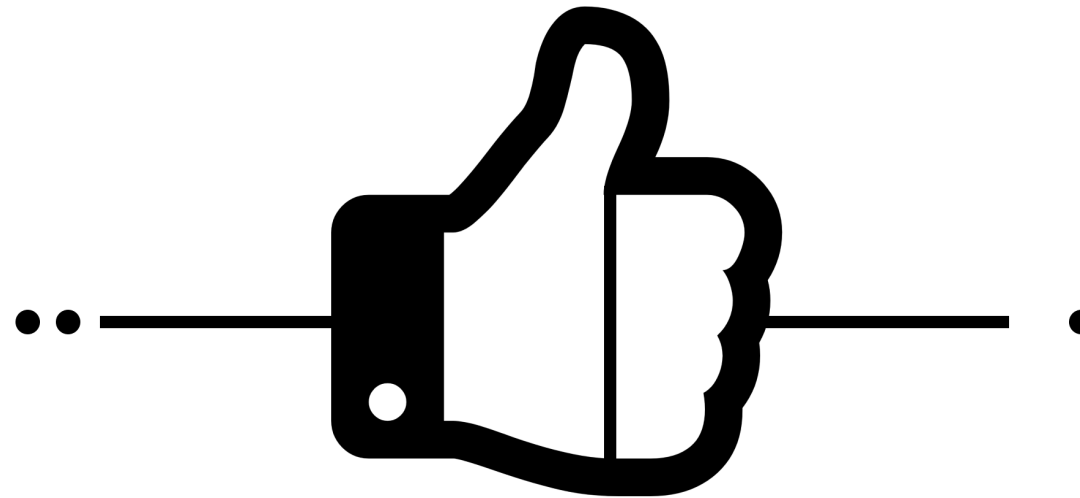
The box plot





Box plot pros

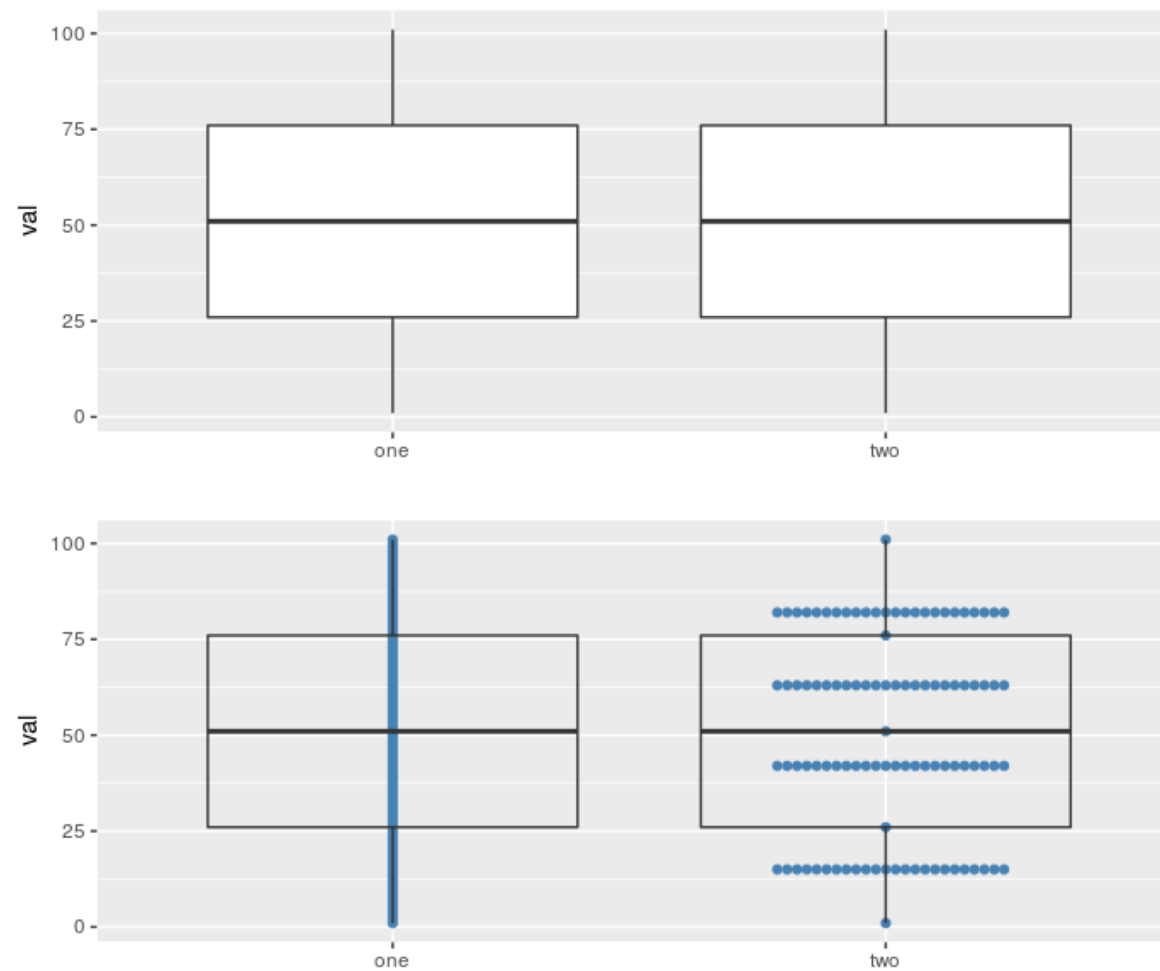
- Familiar
- Lots of good summary statistics





Boxplot cons

- Show me the data!





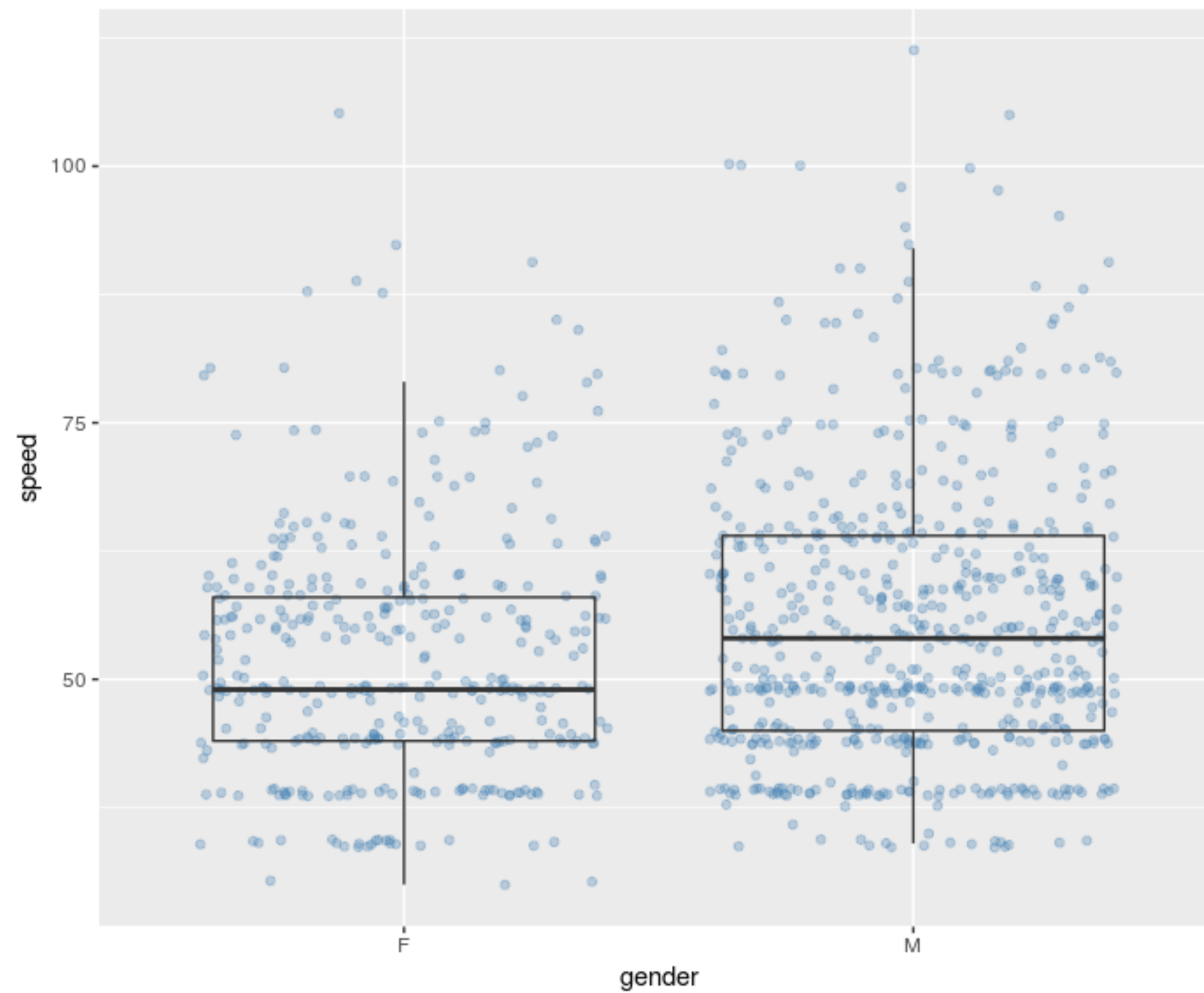
A simple addition

- `geom_jitter()` shows raw points jostled to avoid overlap.
- **Layer under your** `geom_boxplot`.

```
md_speeding %>%  
  filter(vehicle_color == 'BLUE') %>%  
  ggplot(aes(x = gender, y = speed)) +  
  # Draw points behind  
  geom_jitter(alpha = 0.3, color = 'steelblue') +  
  geom_boxplot(alpha = 0) + # make transparent  
  labs(title = 'Distribution of speed for blue cars by gender')
```



Distribution of speed for blue cars by gender





VISUALIZATION BEST PRACTICES IN R

**Let's compare some
distributions**



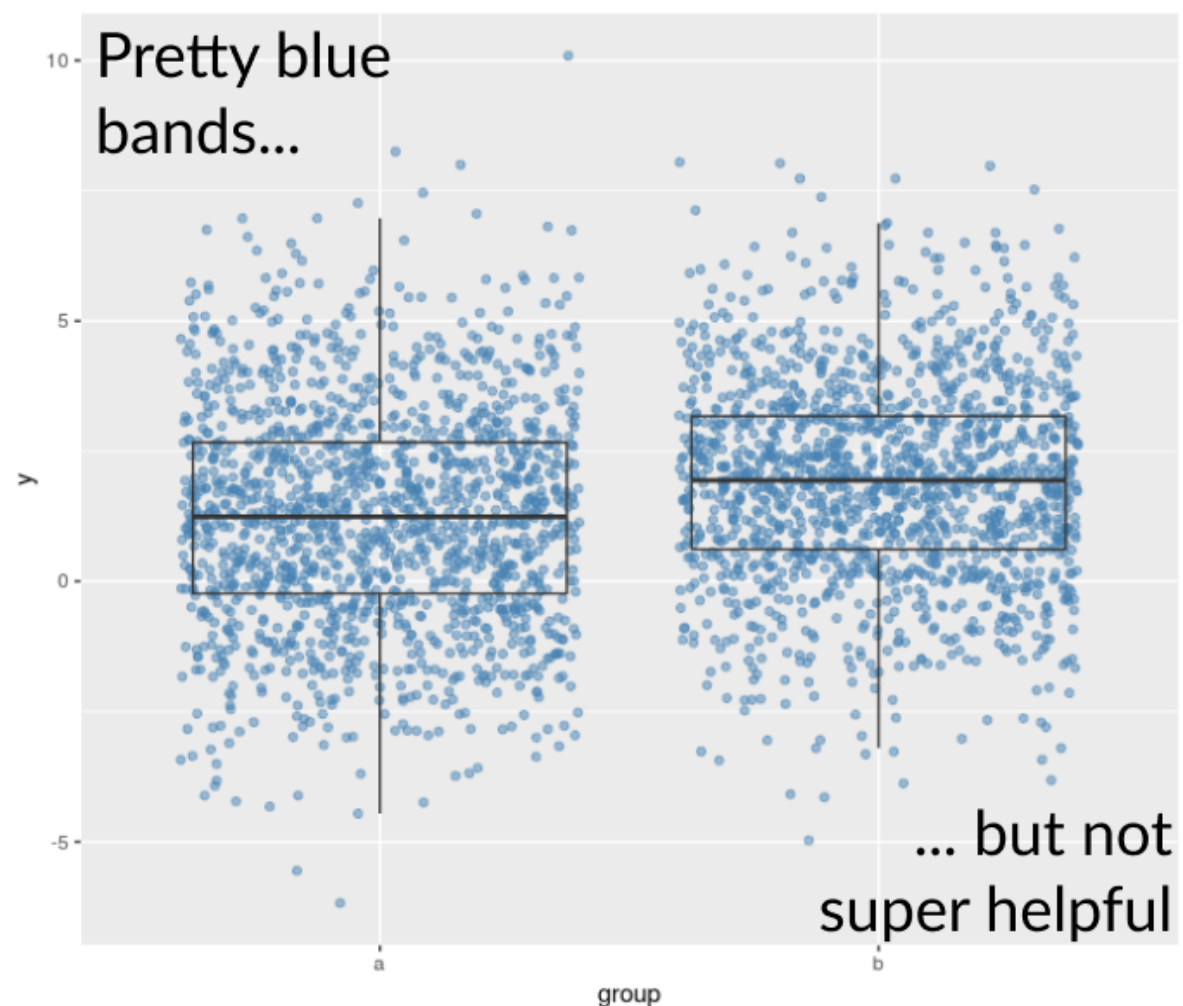
VISUALIZATION BEST PRACTICES IN R

Boxplot alternatives

Nick Strayer
Instructor

Limitations of the boxplot w/ jitter

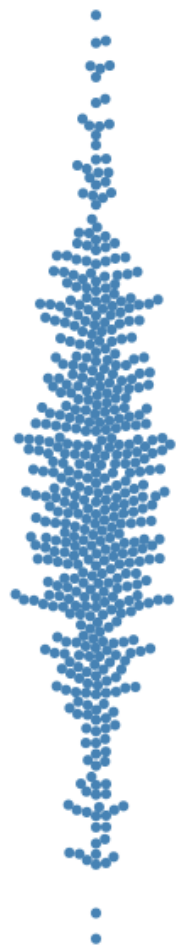
- Josteling points can only deal with so much overlap
- Hard to get an idea of data density





What are some other options?

Beeswarm plots



Violin plots

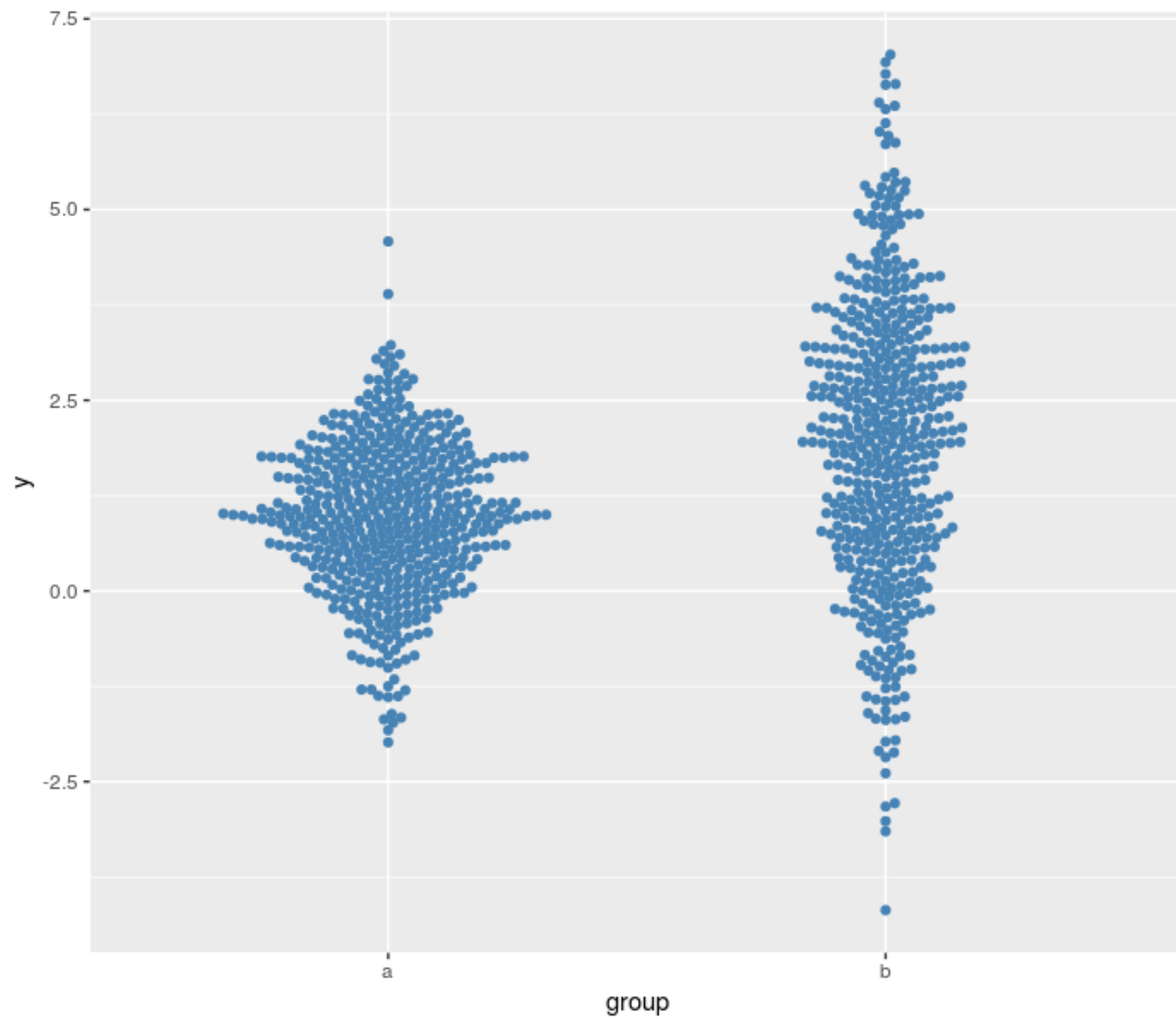




Beeswarm plots

- 'Smart' jittering
- Individual points are clumped together as close to the axis as possible
- Handily included as `geom_beeswarm` in the `ggbeeswarm` package.

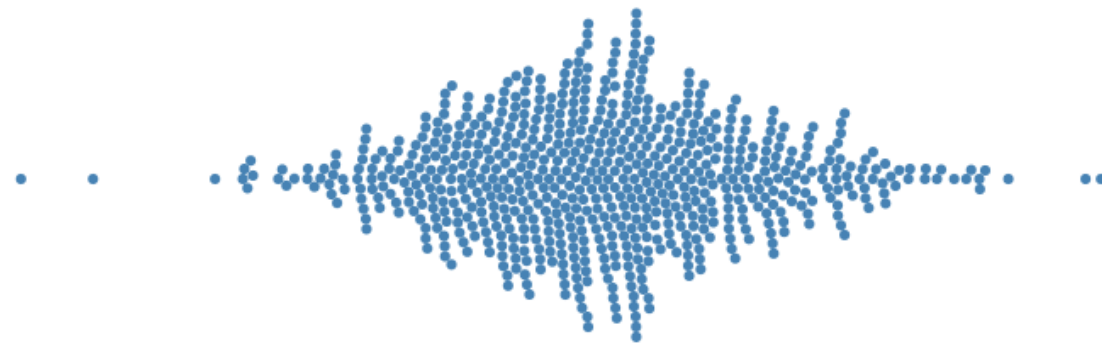
```
library(ggbeeswarm)
ggplot(data, aes(y = y, x = group)) +
  geom_beeswarm(color = 'steelblue')
```





Beeswarm pros

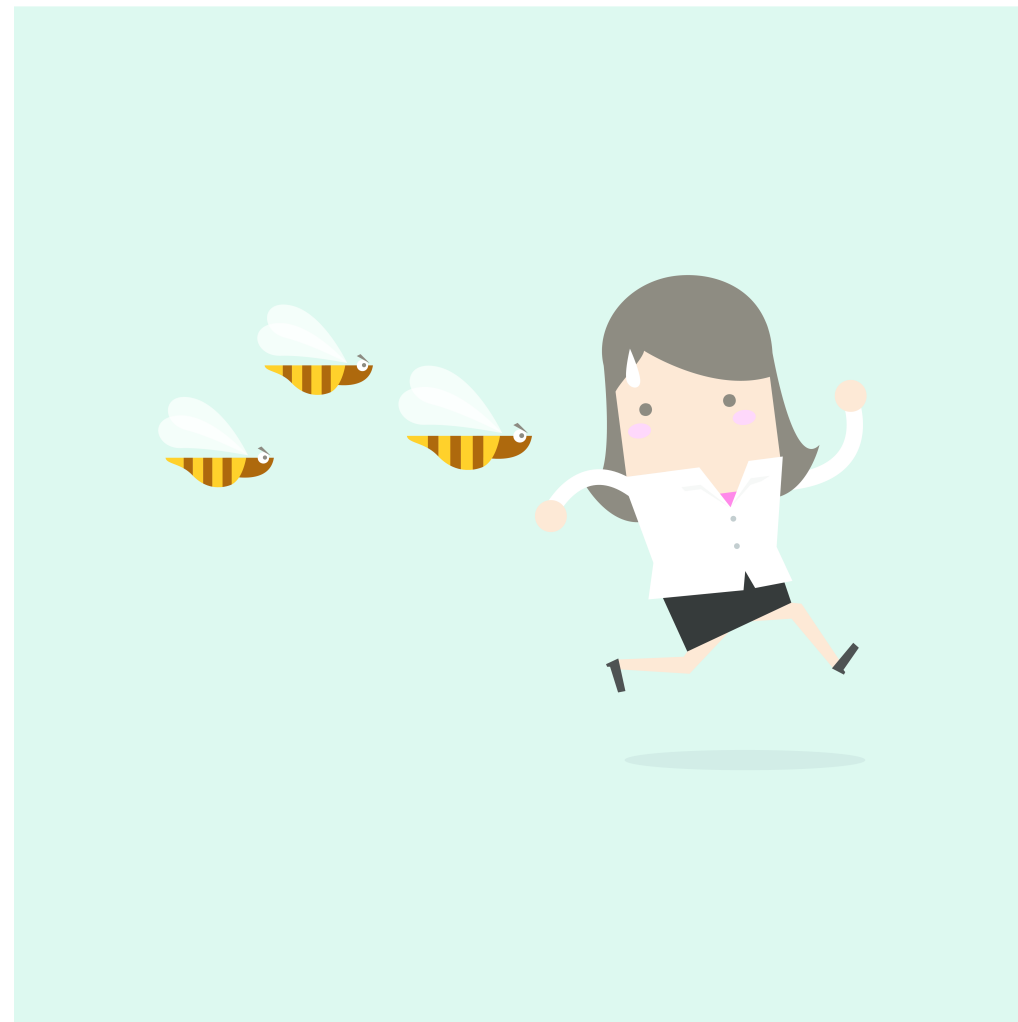
- Individual datapoints
- Distributional shape





Beeswarm cons

- Get hard with lots of data
- Arbitrary stacking

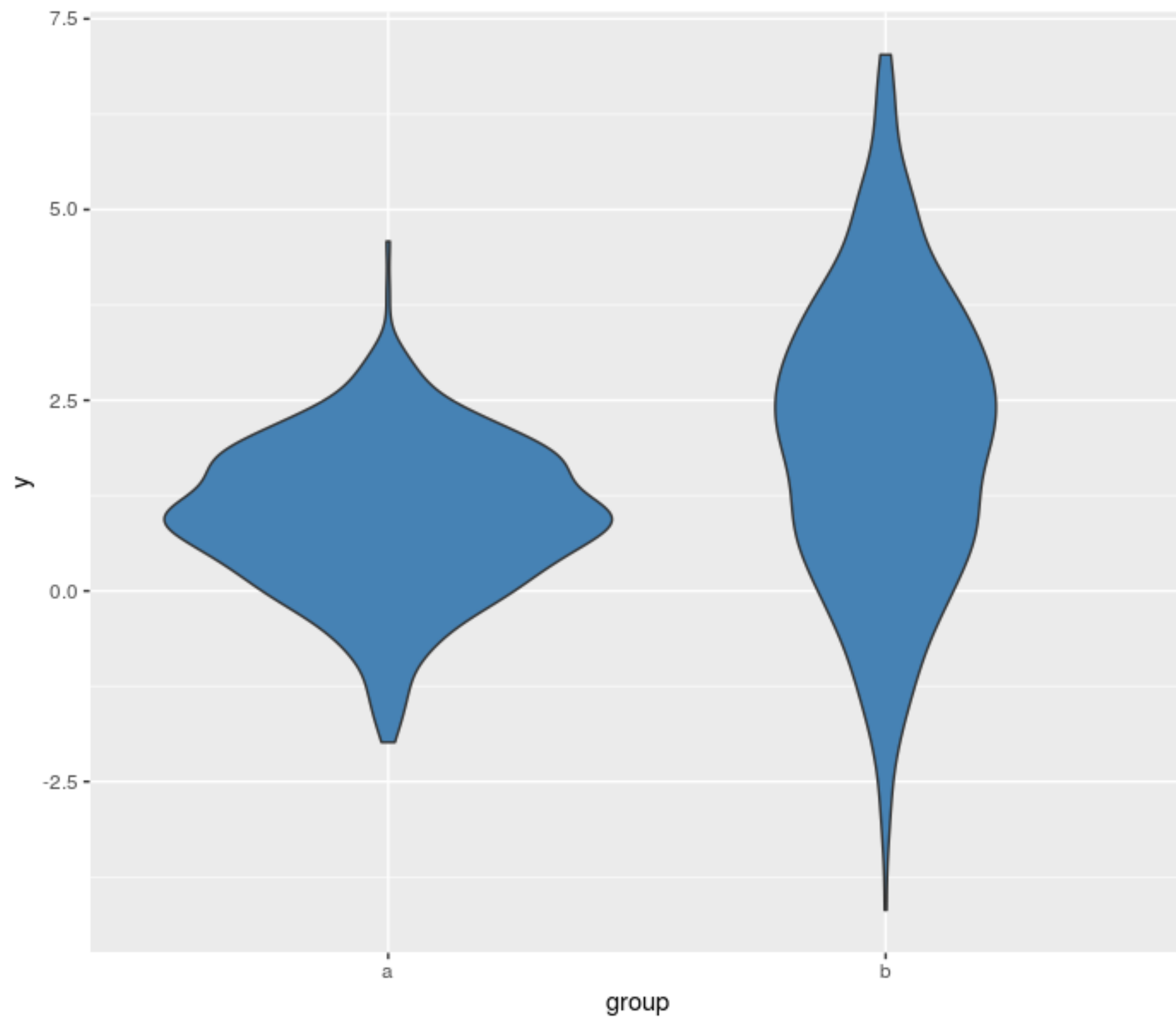




Violin plots

- KDE reflected to be symmetric
- Just replace `geom_boxplot` with `geom_violin`.

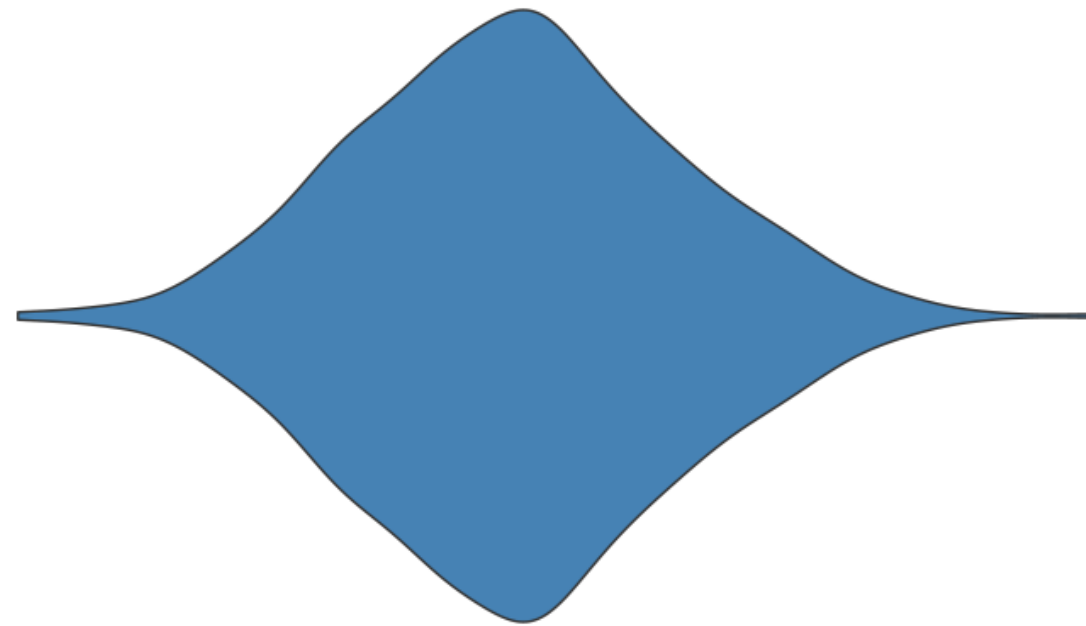
```
ggplot(data, aes(y = y, x = group)) +  
  geom_violin(fill = 'steelblue')
```





Violin pros

- Every datapoint is heard
- Not every datapoint is seen, so good for lots of data.





Violin cons

- Kernel width choice
- Not every datapoint is seen





VISUALIZATION BEST PRACTICES IN R

**Let's try some more
advanced comparisons!**



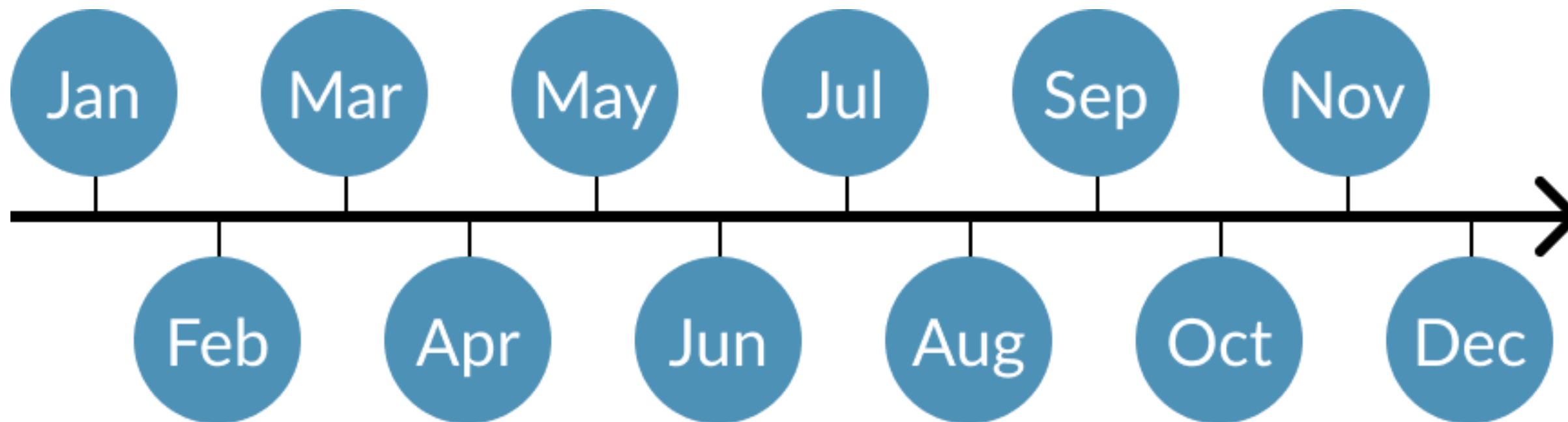
VISUALIZATION BEST PRACTICES IN R

Comparing spatially related distribution

Nick Strayer
Instructor

What are 'spatially connected axes'?

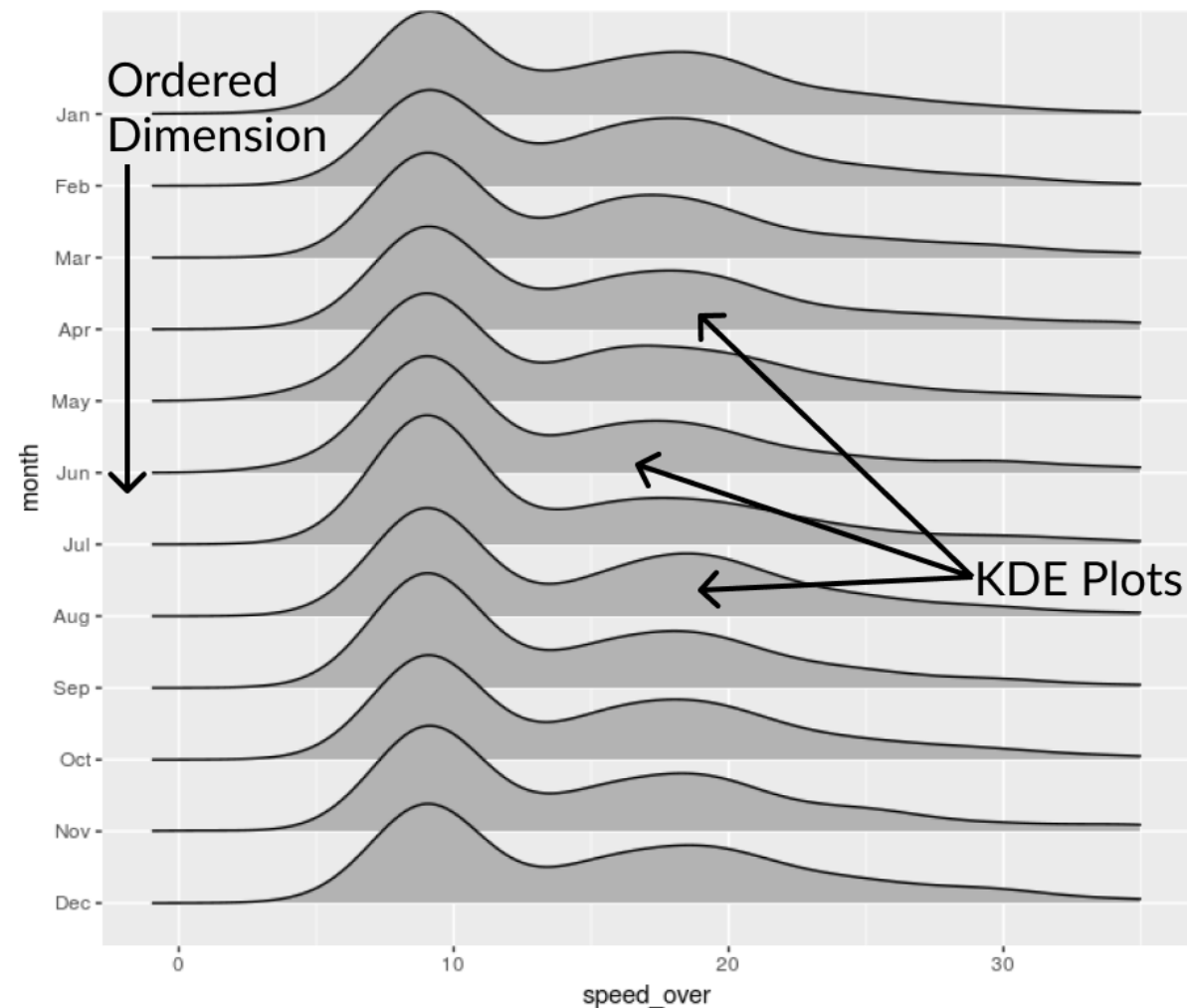
- There is an underlying ordering of the classes.
- E.g. months of the year: Jan < Feb < Mar < ...



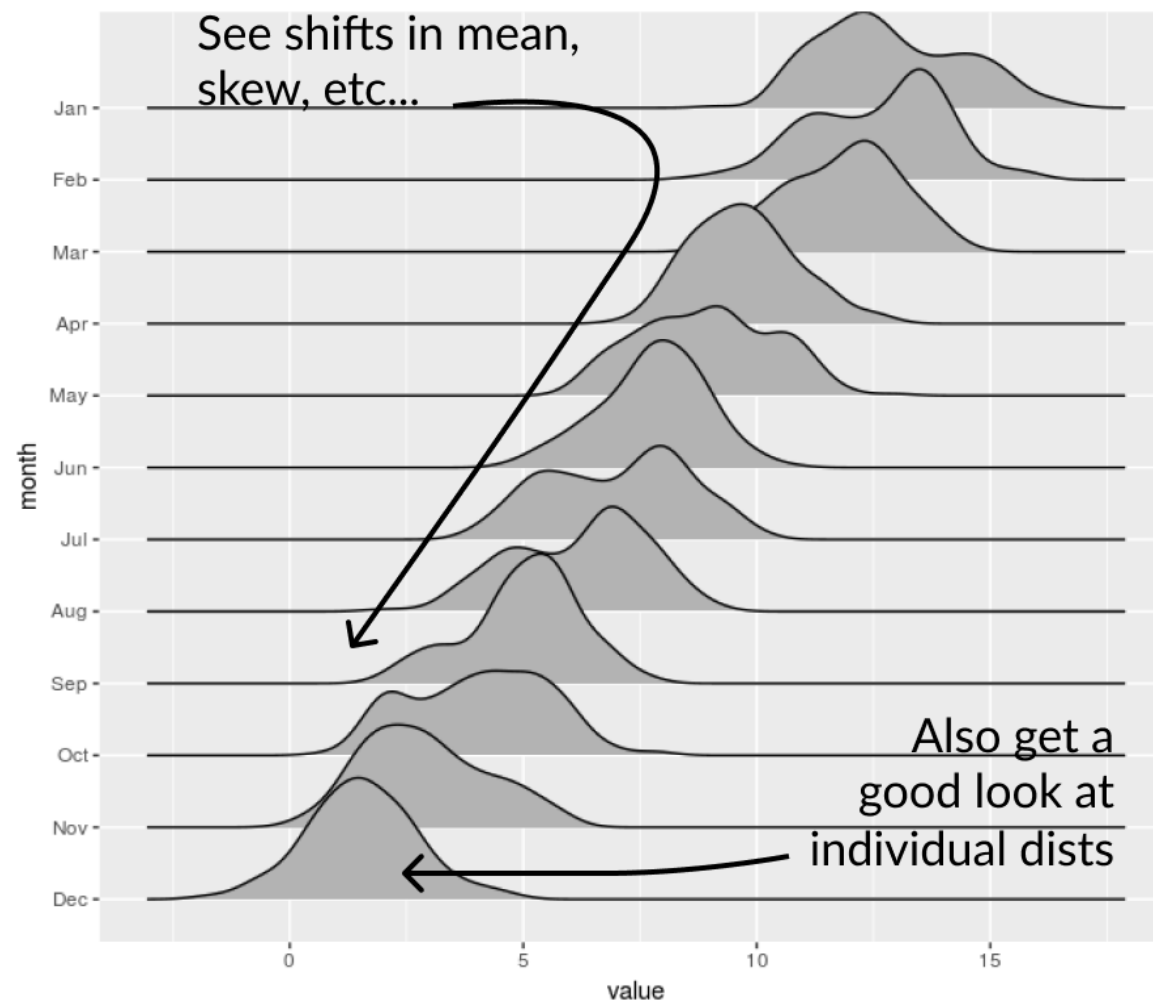


The ridgeline plot

```
library(ggribes) # gives us geom_density_ridges()
ggplot(md_speeding, aes(x = speed_over, y = month)) +
  geom_density_ridges(bandwidth = 2) +
  xlim(1, 35)
```

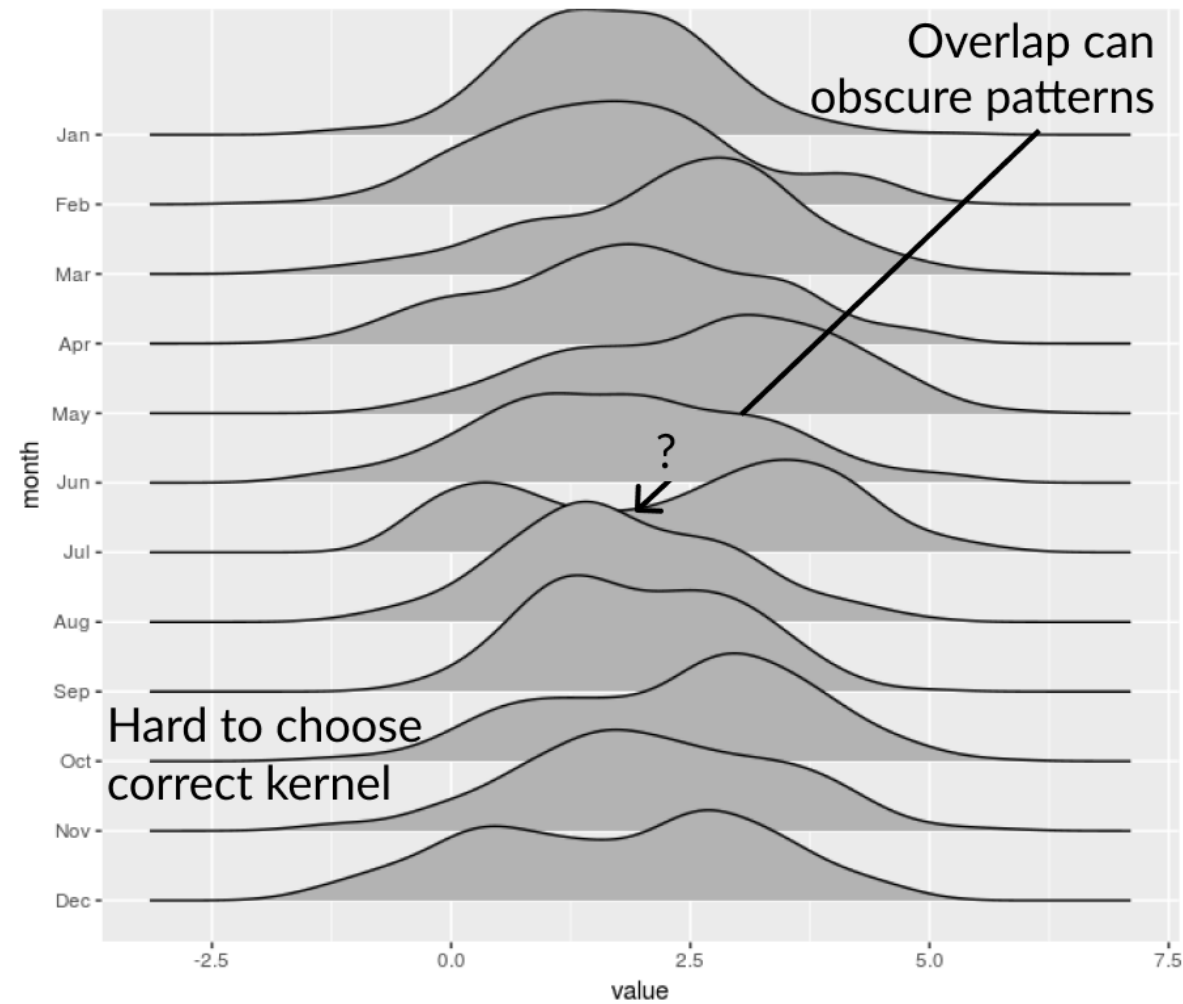


Ridgeline pros



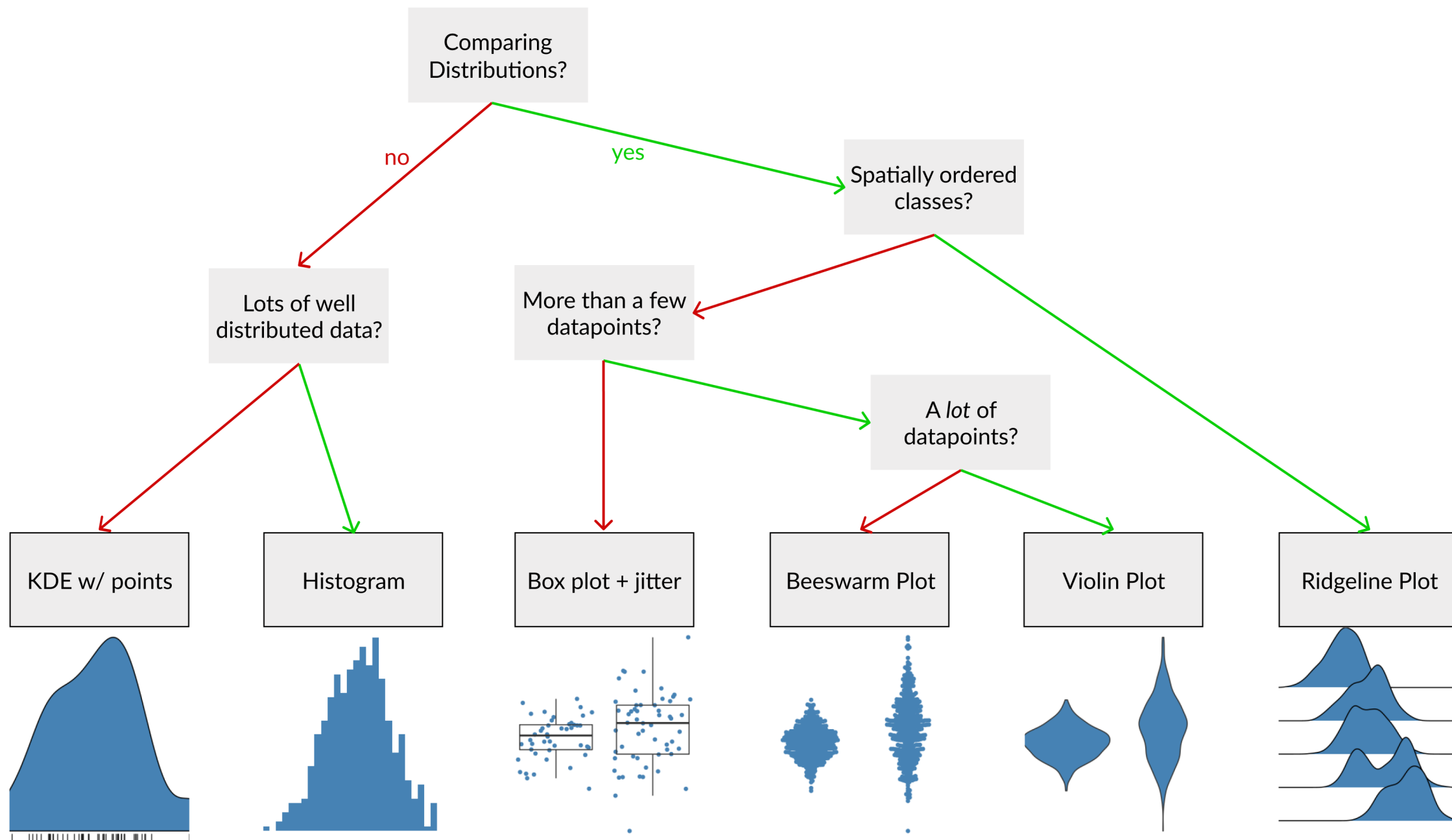


Ridgeline cons





Overview of distribution visualizations





VISUALIZATION BEST PRACTICES IN R

Let's make some ridgelines!



VISUALIZATION BEST PRACTICES IN R

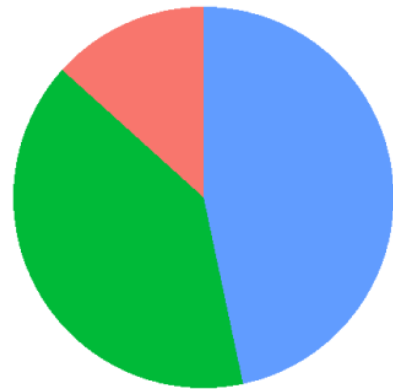
Congratulations!

Nick Strayer
Instructor

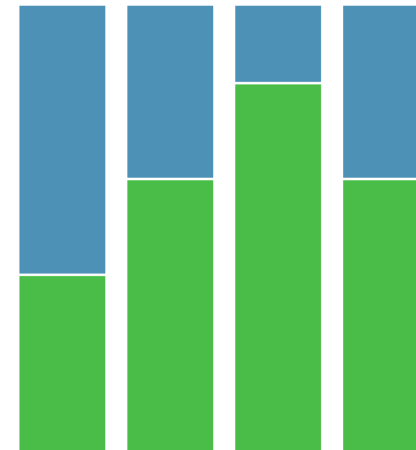


Chapter 1: Proportions

Three or less classes and precision not important?



Good for comparing values
across populations....



Need more precision and
have the space?

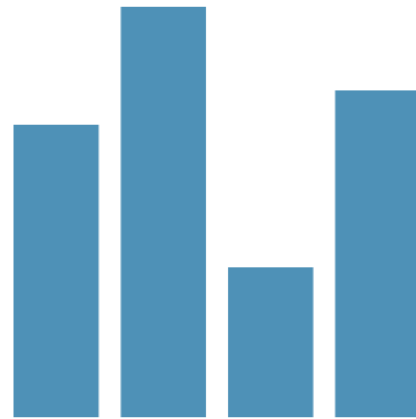


... bad for comparing
values *within* populations

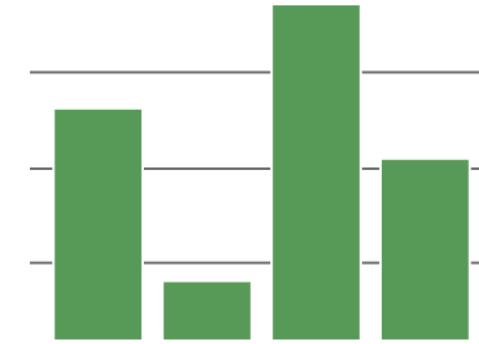


Chapter 2: Point data

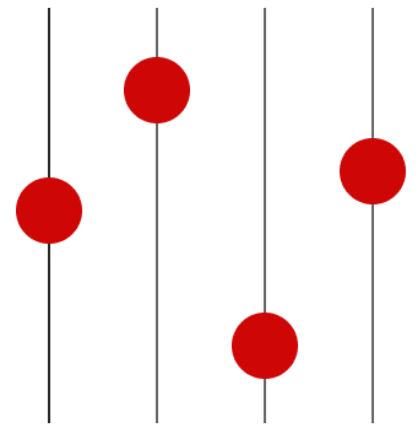
Data is a stackable quantity? E.g. dollars, counts...



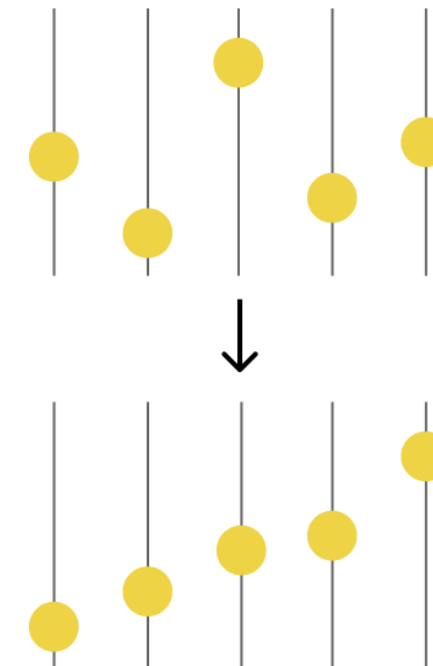
No need for vertical grid lines on bars.



Not stackable? E.g. percents, ratio...



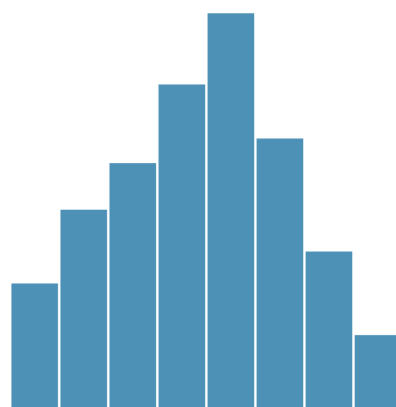
Ordering is almost always a good idea.



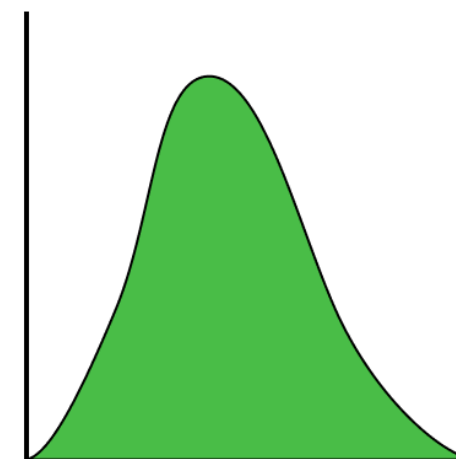


Chapter 3: Single distributions

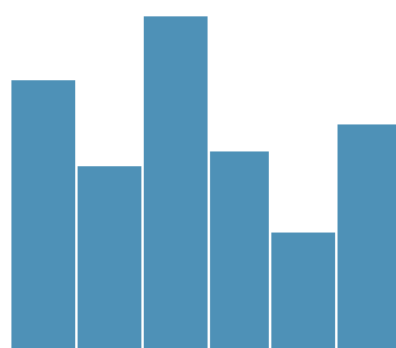
Histograms are simple and intuitive...



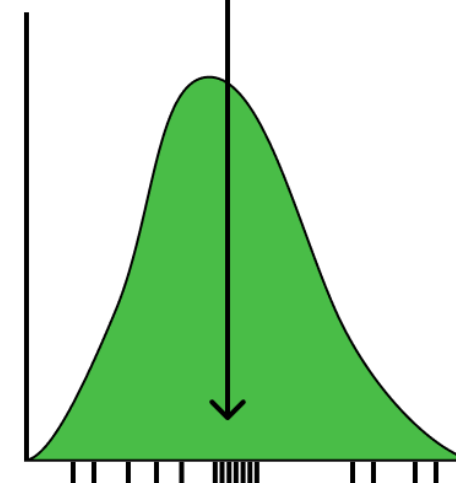
Kernel Density Estimators work with small data and produce nice smooth plots.



... but are sensitive to bin width and placement choices.



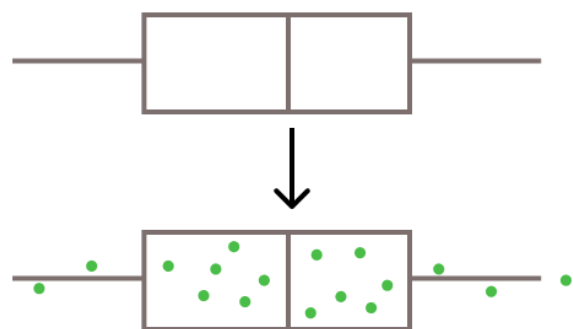
But always try and show raw data if possible to see gaps the KDE may be filling.



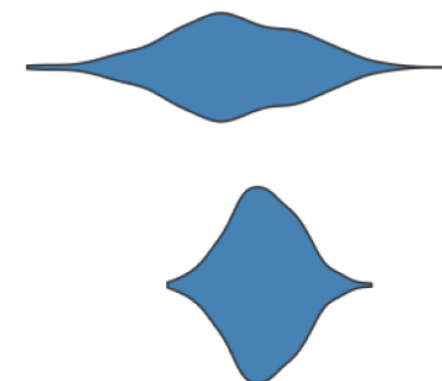


Chapter 4: Comparing distributions

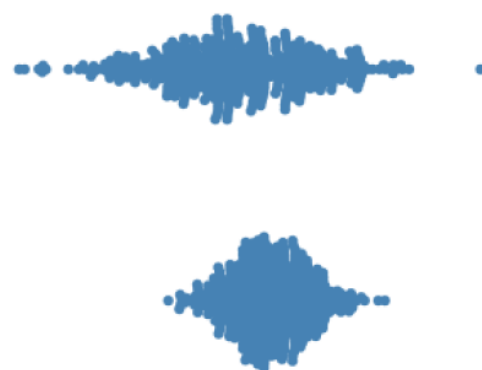
Boxplots hide a lot of data, so augment them with jittered points



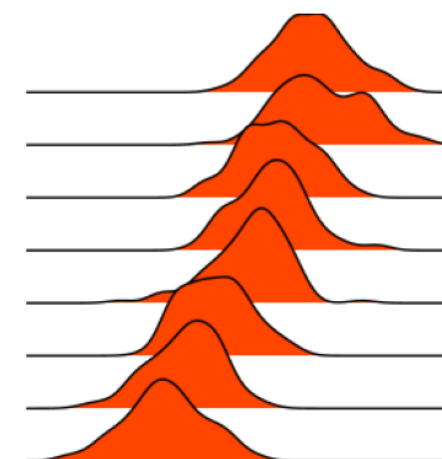
Violin plots are symmetric KDEs that work well when you have lots of points.



Beeswarm plots are an alternative 'smart' jittering that shows density.



If you have spatially ordered groups, consider the ridgeline plot.





Going further

[Flowing data](#)

Curated list of data visualizations and R-based tutorials.

[Datawrapper Blog](#)

Articles that dig deep into visualization techniques and mistakes.

[Twitter \(#datavis\)](#)

An ongoing stream of cool projects and inspiration.

Books!

- [Data Visualization](#), Andy Kirk
- [The Functional Art and The Truthful Art](#) by Alberto Cairo



VISUALIZATION BEST PRACTICES IN R

Thank You!