

# Data Visualization with ggplot2

Session 3: Visualizing Diverse Data

Nathan Barron

✉ [nathanbarron@ou.edu](mailto:nathanbarron@ou.edu)

# Different data need different visualizations

data	mpg	cyl	am	cyl x am							
	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3	1
Duster 360	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3	4
Merc 240D	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4	2
Merc 230	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4	2
Merc 280	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4	4

# Types of data

## Levels of Data Measurement

- Categorical
  - Binary
  - Non-binary
- Continuous

[.](#) [Cat.: Binary](#) [Cat.: Non-binary](#)

[Continuous](#)

---

*Click through each header to read more about each type*

## Storing Data in R

- Boolean
- Character
- Numerical
- Factor

[.](#) [Boolean](#) [Character](#) [Numerical](#)

[Factor](#) [Dates](#)

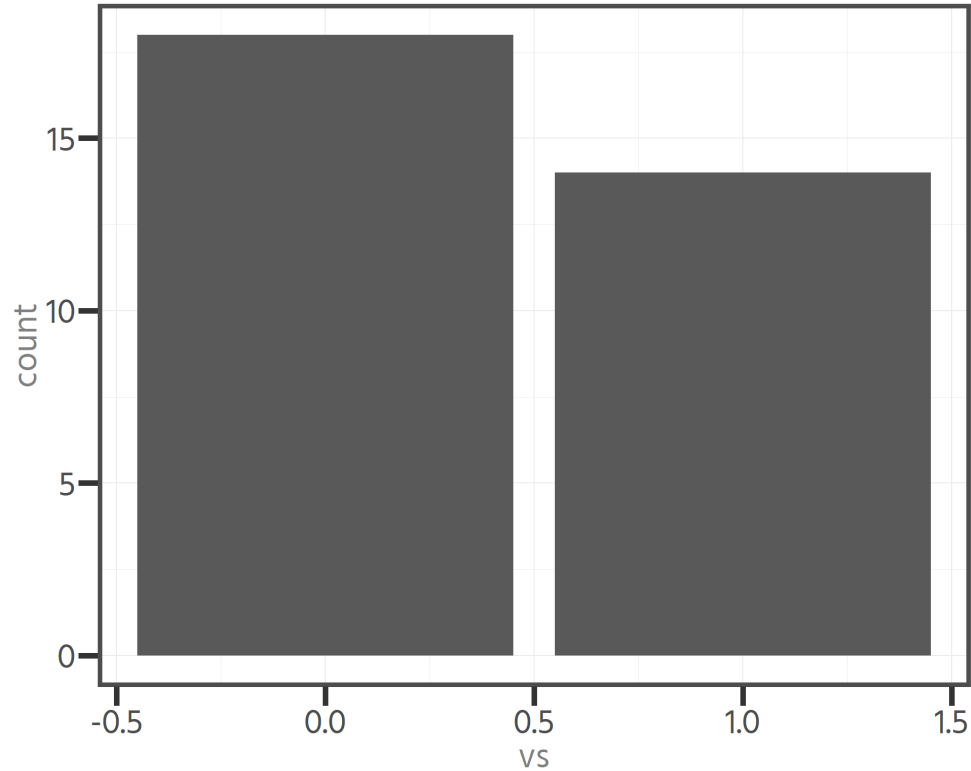
---

*Click through each header to read more about each type*

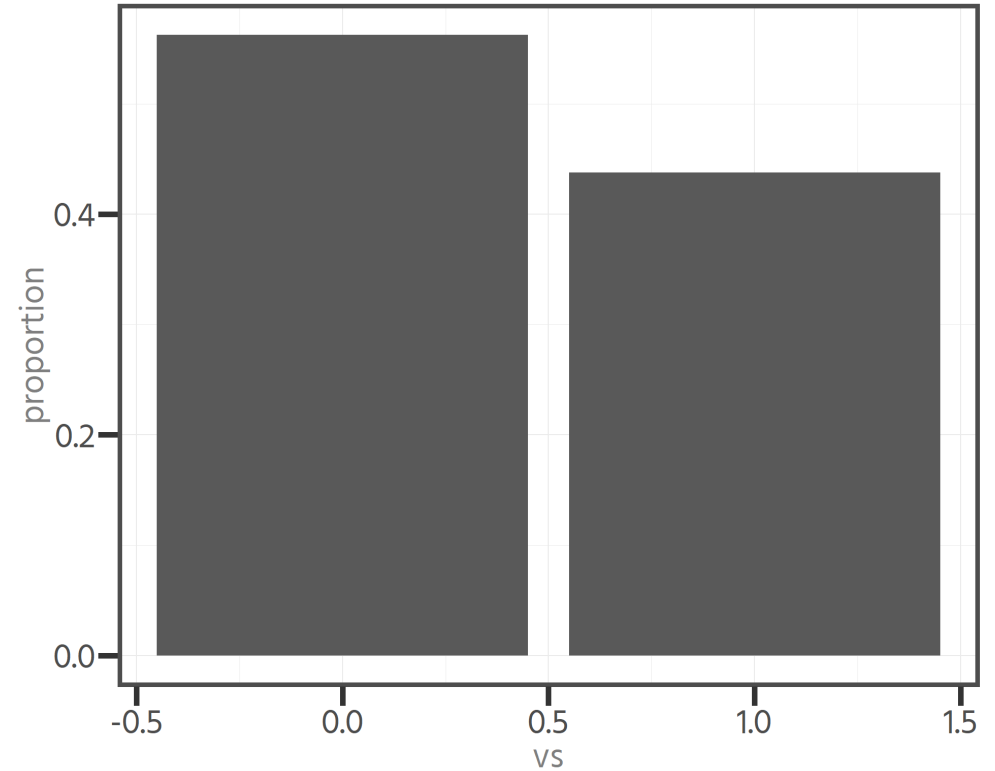
# Visualizing a single variable

# Binary

Barplot (count)



Barplot (proportion)

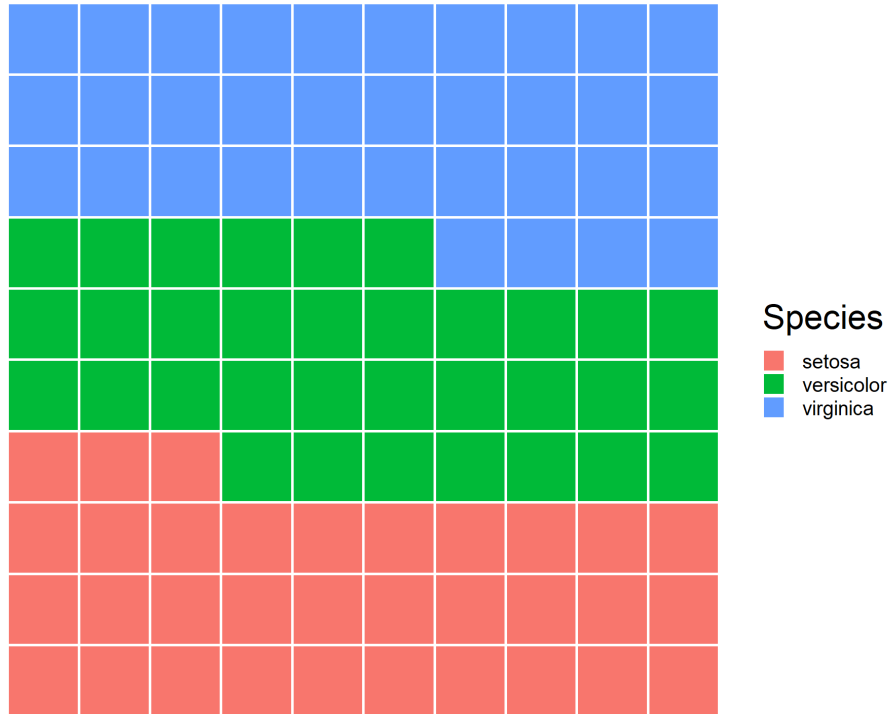


```
1 # Barplot with frequency on y-axis
2 ggplot(mtcars) +
3   geom_bar(aes(x=vs))
```

```
1 # Barplot with proportion on y-axis
2 ggplot(mtcars) +
3   geom_bar(aes(x=vs, y=..count../sum(..count..)))
```

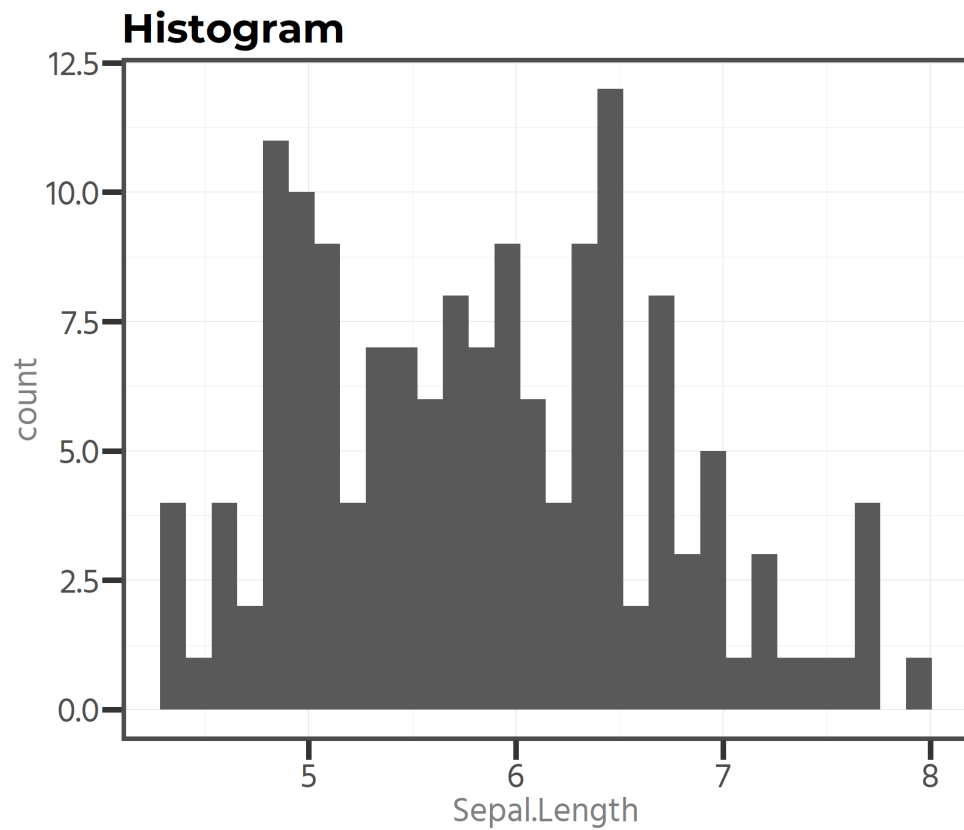
# Non-binary Categorical

## Waffle plot

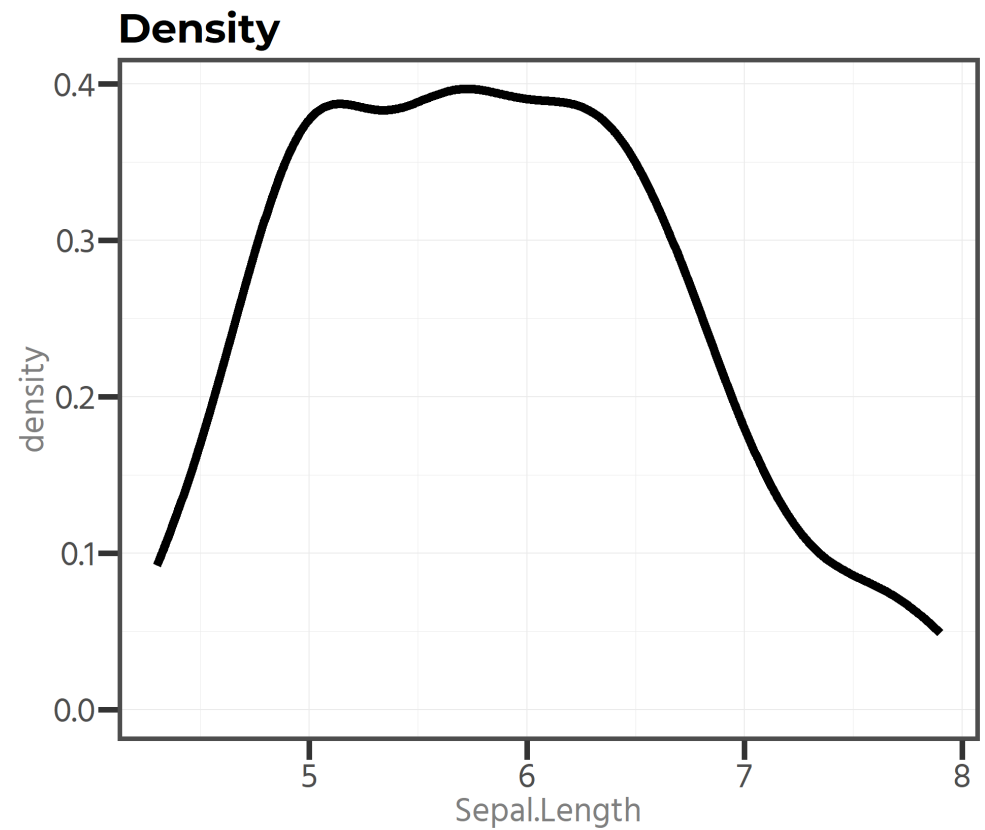


```
1 # install.packages("waffle",
2 #                   repos = "https://cinc.rud.is")
3
4 library(waffle)
5 library(dplyr)
6
7 iris %>%
8   count(Species) %>%
9   ggplot(aes(fill = Species, values = n)) +
10  geom_waffle(size = 1,
11              colour = "white",
12              na.rm=TRUE,
13              flip = TRUE,
14              make_proportional = TRUE) +
15  theme_void() +
16  coord_equal()
```

# Continuous



```
1 ggplot(iris) +  
2   geom_histogram(aes(x=Sepal.Length))
```

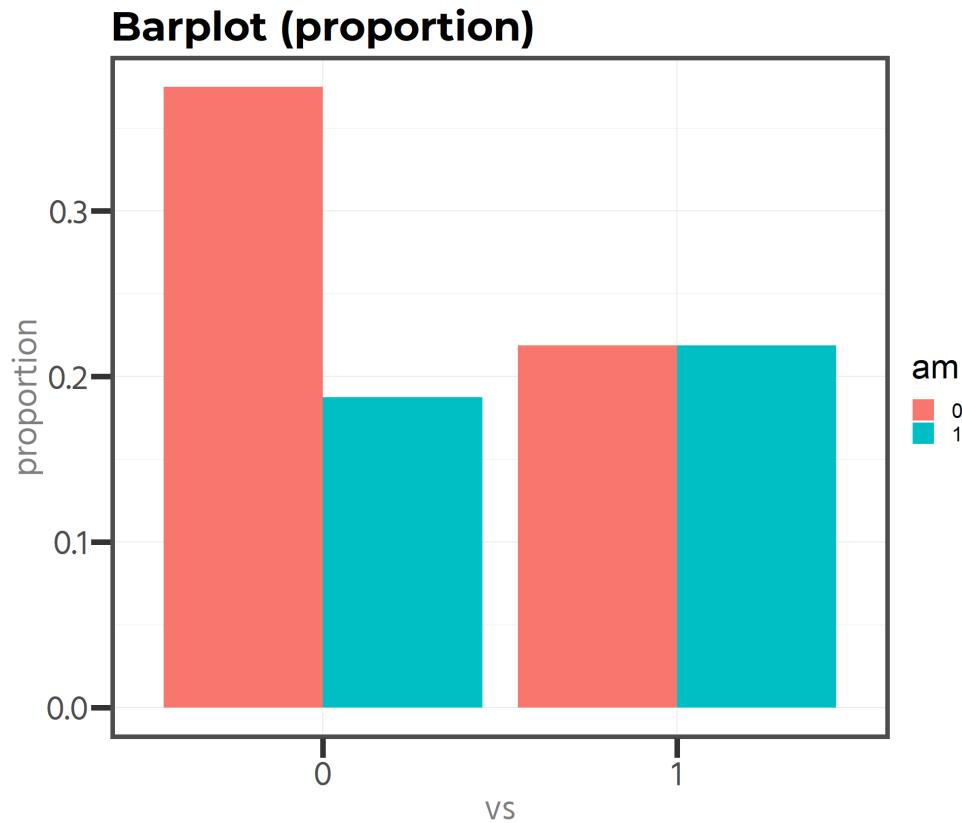


```
1 ggplot(iris) +  
2   geom_density(aes(x=Sepal.Length))
```

# Visualizing multiple variables



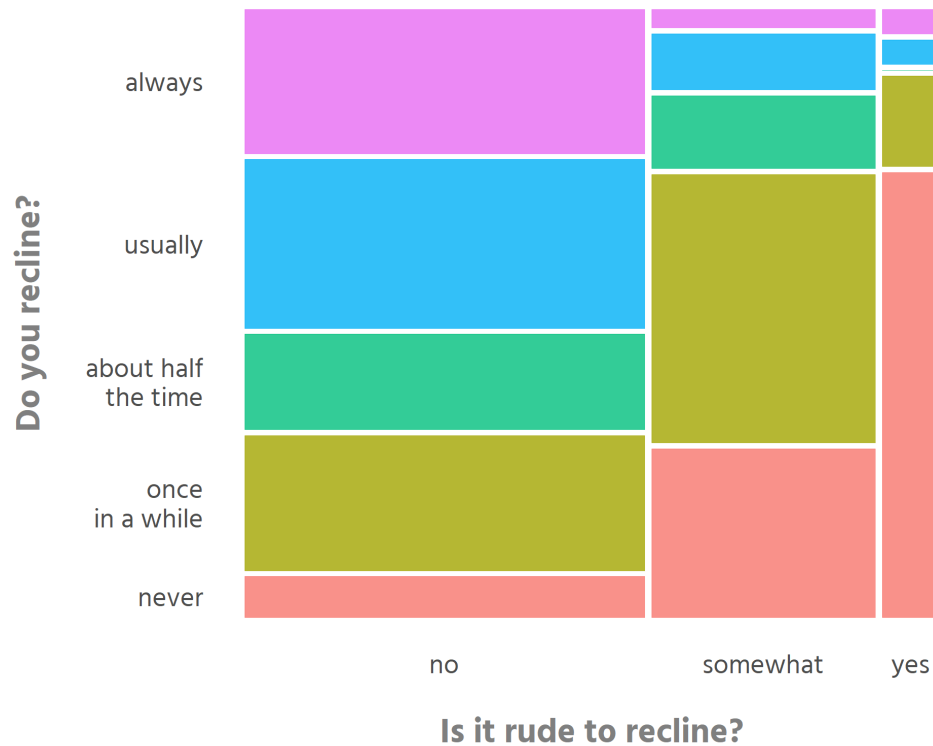
# Multiple categorical variables



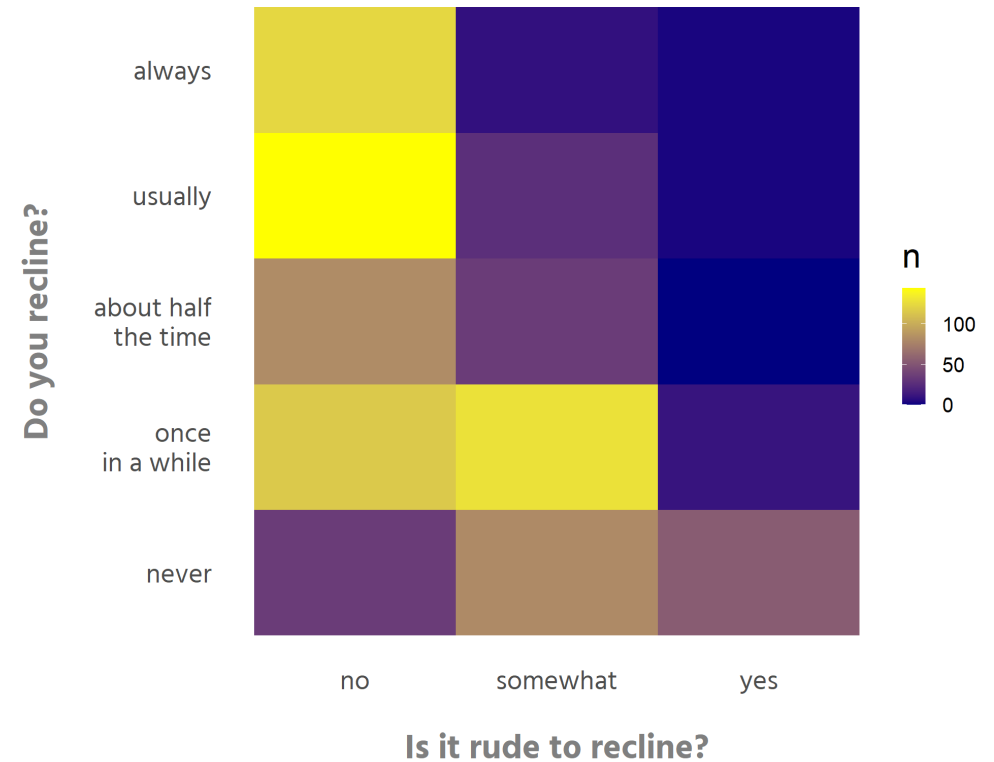
```
1 mtcars %>%  
2   mutate(vs = factor(vs), am = factor(am)) %>%  
3   ggplot(.) +  
4   geom_bar( aes(x=vs,  
5                 y=..count../sum(..count..),  
6                 group=am,  
7                 fill=am),  
8             position = 'dodge')
```

# Multiple categorical variables II

**Mosaic plot**



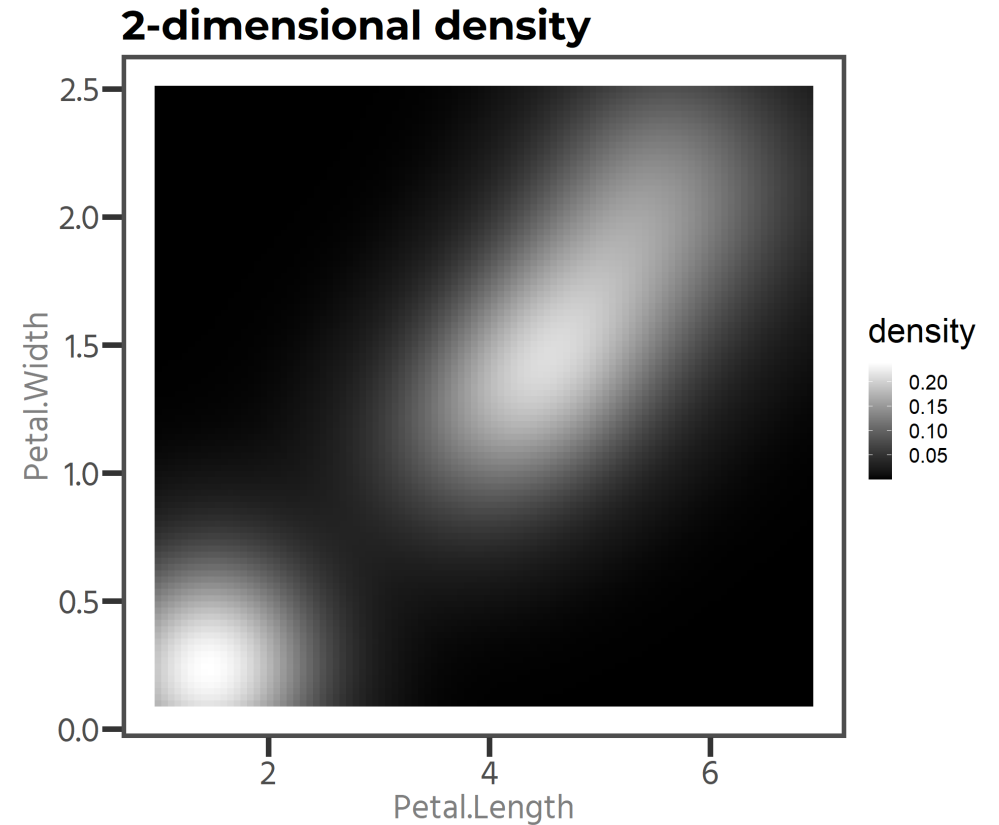
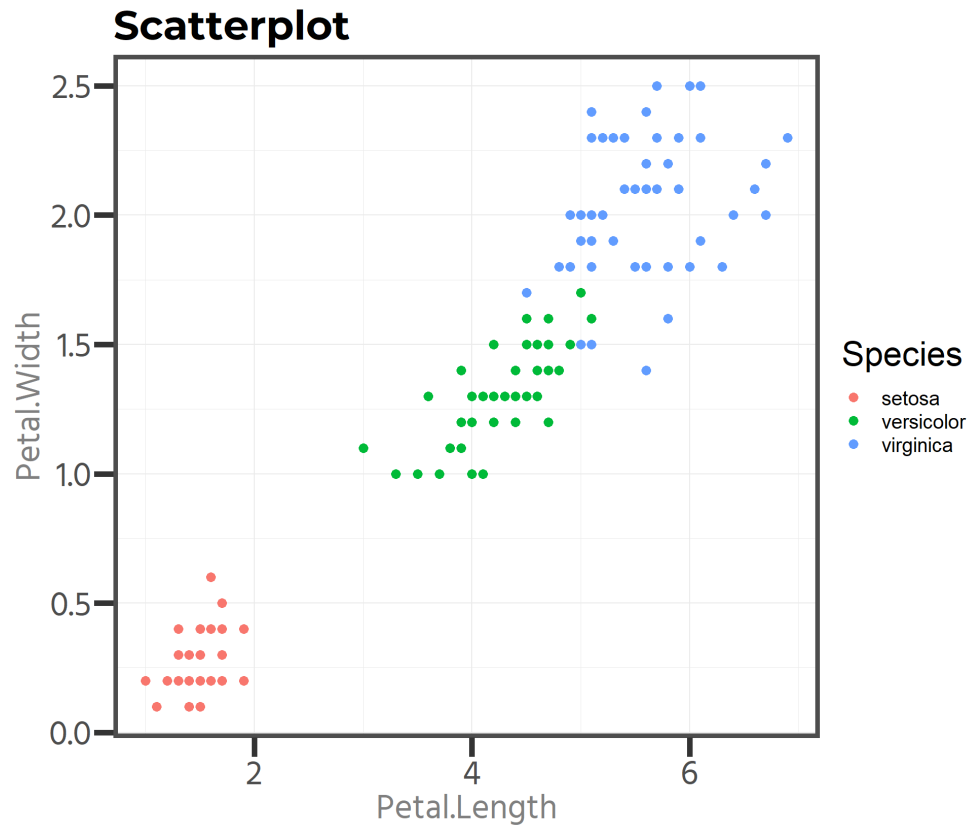
**Categorical heatmap**



```
1 # library(ggmosaic)
2 # library(tidyr)
3
4 ggplot(ggmosaic::fly) +
5   geom_mosaic(aes(x = product(do_you_recline, rude_to_recline))) +
6   theme(panel.grid = element_blank(),
7         legend.position = 'none')
```

```
1 ggmosaic::fly %>%
2   count(rude_to_recline, do_you_recline, .drop = FALSE) %>%
3   ggplot(.) +
4     geom_tile(aes(x = rude_to_recline,
5                   y = do_you_recline,
6                   fill = n)) +
7     scale_fill_gradient(low = 'navy', high = "yellow")
```

# Multiple continuous variables



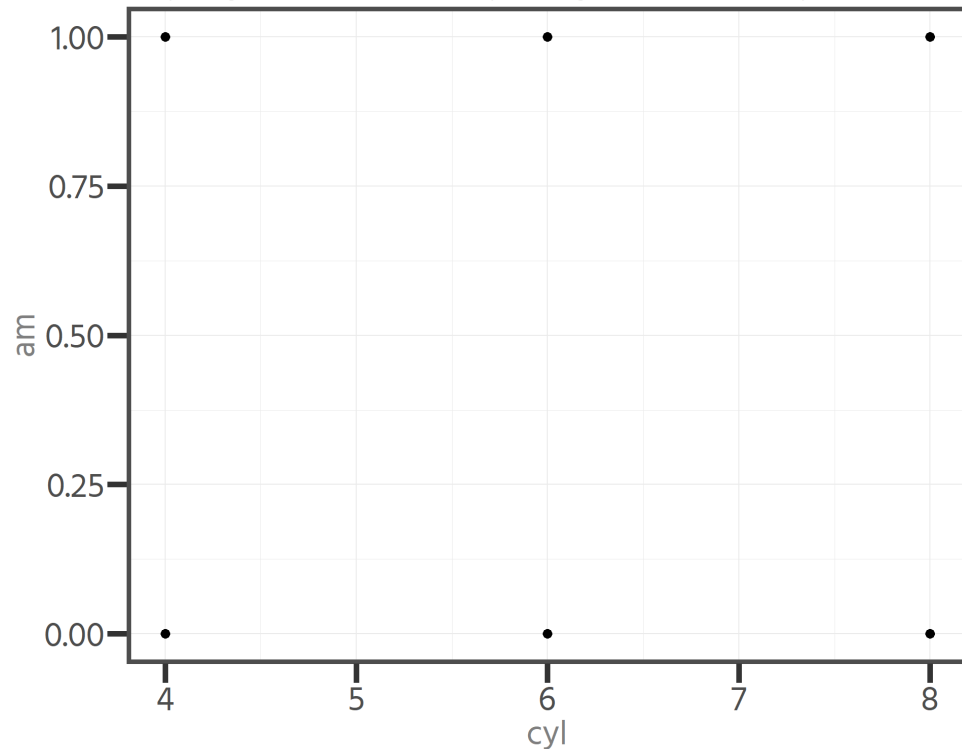
```
1 ggplot(iris) +  
2   geom_point(aes(x=Petal.Length,  
3                 y=Petal.Width,  
4                 color=Species))
```

```
1 ggplot(iris, aes(x=Petal.Length, y=Petal.Width)) +  
2   stat_density2d(aes(fill = ..density..),  
3                 geom = "raster", contour = FALSE) +  
4   scale_fill_gradient(low = "black", high = "white")
```

# Hybrid (some categorical, some continuous)

## Overplotting

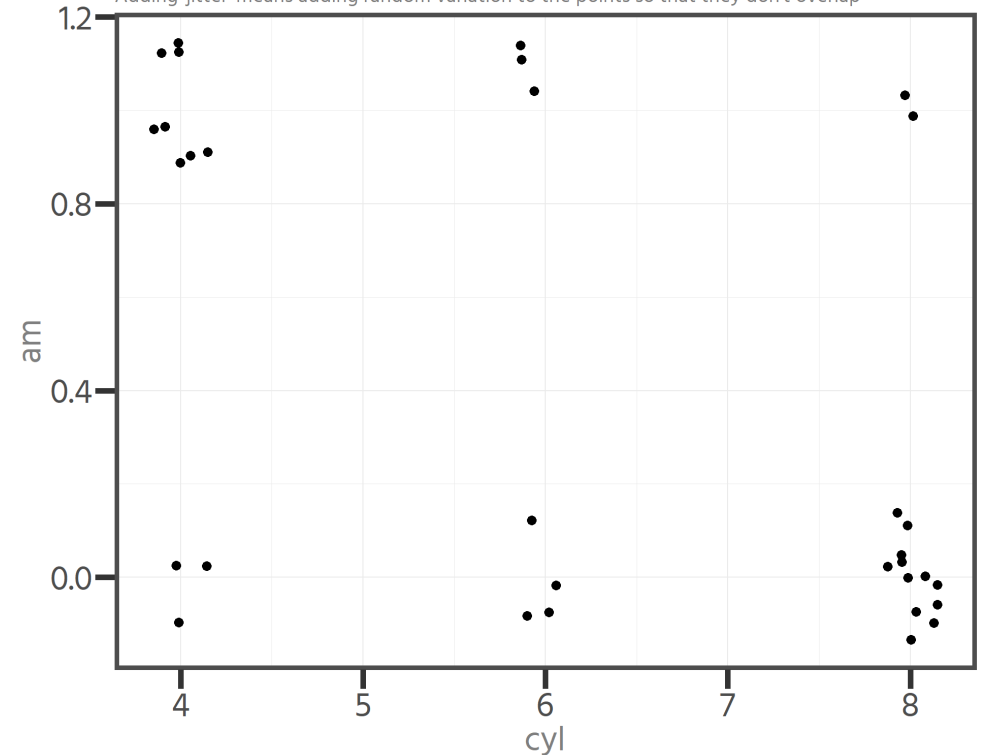
Overplotting is when data or labels overlap, meaning that the viewer loses important information



```
1 ggplot(mtcars) +  
2   geom_point(aes(x=cyl, y=am))
```

## Jitterplot

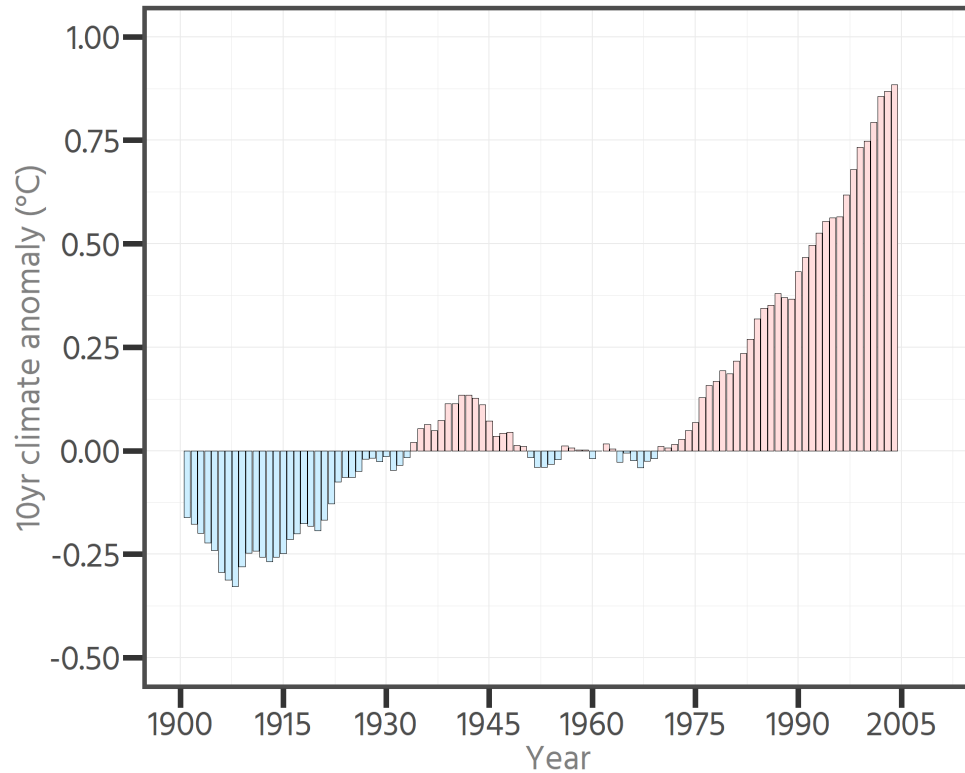
Adding 'jitter' means adding random variation to the points so that they don't overlap



```
1 ggplot(mtcars) +  
2   geom_jitter(aes(x=cyl, y=am))  
3  
4 # You can also specify 'width' and 'height'  
5 # as additional arguments to refine the  
6 # jitter. Remember, jitter is random:  
7 # jitter plots will change whenever you  
8 # re-run the code.
```

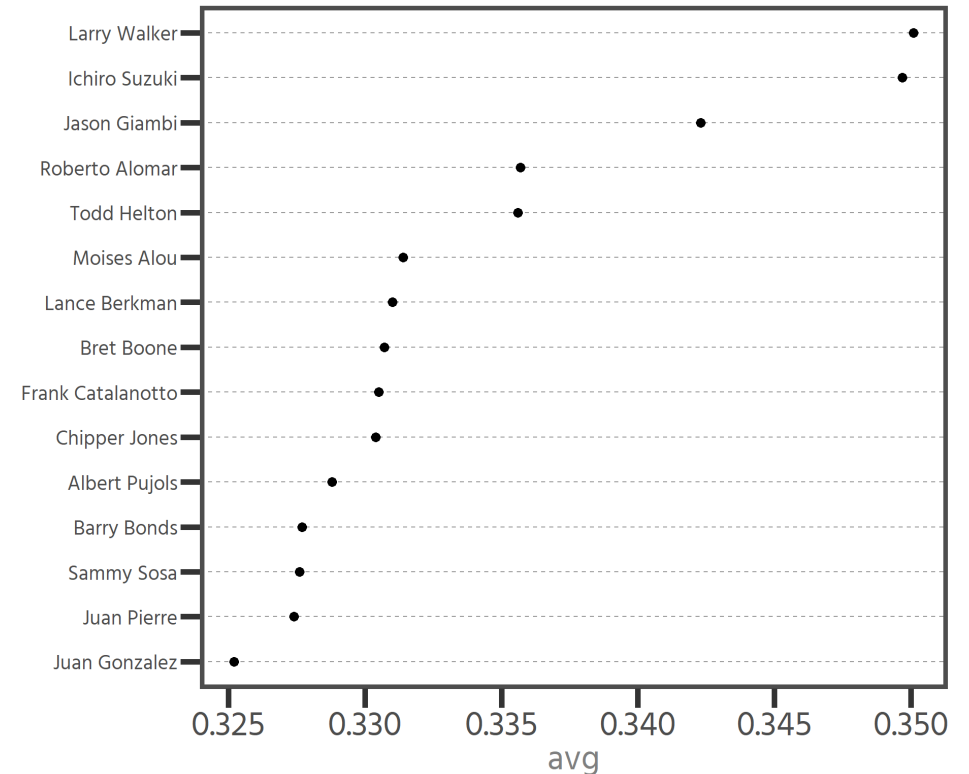
# Hybrid II

Barplot with +/- indicator



```
1 library(gcookbook)
2 climate_sub <- climate %>%
3   filter(Source == "Berkeley" & Year >= 1900) %>%
4   mutate(pos = Anomaly10y >= 0)
5
6 ggplot(climate_sub, aes(x = Year, y = Anomaly10y, fill =
7   geom_col(position = "identity", colour = "black", size
8   scale_fill_manual(values = c("#CCEEFF", "#FFDDDD"), gui
```

Cleveland dot plot



```
1 library(gcookbook)
2 tophit <- tophitters2001[1:15, ]
3
4 ggplot(tophit) +
5   geom_point(aes(x = avg, y = reorder(name, avg)))
6
7 # The `reorder` function is helpful
8 # when making plots of a non-binary categorical
9 # variable and a corresponding continuous variable.
```

There is no “one size fits all”

# Data viz is creative ...

and creativity requires inspiration.

## Check out these cool resources:

- [R graphics gallery](#)
- [Nightingale: The Magazine of the Data Visualization Society](#)
- [Daydreaming Numbers Blog](#)

# Data viz is also programming ...

and programming requires code.

## Check out these cool resources:

- [R Graphics Cookbook](#)
- [R Graph Gallery](#)