

Representing data

Experimental Research

Professor Velez (11/1/22)

Recap

- Learned fundamentals regarding experimental design
- Designed and analyzed a non-human subjects experiment
- Created an intervention that we will examine in the context of a class-wide experiment
- We are currently waiting for the IRB review



This week

- Revisit concepts we covered in the measurement week
- Create summaries and data visualizations
- **Crucial** early step in the research pipeline
 - Identity errors in coding
 - Highlight expected patterns
 - Uncover unexpected patterns

Summarizing data

- We've already been working with data summaries
- Means and standard deviations
 - Capture the central tendency of a variable
 - Capture the typical observed variation
- Statistics help simplify the fairly complex information we collect in the social sciences

Example

Find every '5'

9176867991960386255930486551443
9353652502752141394912668766013
5095444663473545993604636448078
6425620171155906710121069218701
0584448205561385902503573068845
1172102391879794238892669764599
6987227951015948926017759166604
8436204036527029291707742717051
3280506428553158863136868380421
4055135906770329783671748392874

Asking questions about our data

- How is the variable distributed?
 - What is the most likely value for this variable? **Mean**
 - What is the value at the midpoint? **Median**
 - What is the most common value? **Mode**
 - Are observations dispersed across the range of the variable? **Variance/Standard Deviation**
- How are variables related to each other?
 - When X is high (low), is Y also high (low)? **Covariance/correlation**

Different averages

- Mean
 - Sum of values divided by number of values
 - Sensitive to outliers
- Median
 - Midpoint of a variable (50th percentile)
 - Less sensitive to outliers
- Mode
 - Most common value
 - Can be helpful when working with nominal data

Polling Data

| Polling Data | | | | | | |
|-------------------------------|---------------|---------|-----------------|---------------|------------------|--|
| Poll | Date | Sample | Republicans (R) | Democrats (D) | Spread | |
| RCP Average | 10/12 - 10/28 | -- | 48.0 | 45.1 | Republicans +2.9 | |
| CBS News Battleground Tracker | 10/26 - 10/28 | 1899 LV | 47 | 45 | Republicans +2 | |
| Data for Progress (D)** | 10/26 - 10/27 | 1217 LV | 49 | 45 | Republicans +4 | |
| Trafalgar Group (R) | 10/25 - 10/27 | 1089 LV | 48 | 42 | Republicans +6 | |
| Rasmussen Reports | 10/23 - 10/27 | 2500 LV | 49 | 42 | Republicans +7 | |
| InsiderAdvantage | 10/25 - 10/25 | 750 LV | 48 | 44 | Republicans +4 | |
| Economist/YouGov | 10/22 - 10/25 | 1114 LV | 45 | 49 | Democrats +4 | |
| Politico/Morning Consult | 10/21 - 10/23 | 2005 RV | 42 | 47 | Democrats +5 | |
| USA Today/Suffolk | 10/19 - 10/24 | 1000 LV | 49 | 45 | Republicans +4 | |
| Democracy Corps (D) | 10/19 - 10/23 | LV | 50 | 48 | Republicans +2 | |
| Emerson | 10/18 - 10/19 | 1000 LV | 46 | 41 | Republicans +5 | |
| NBC News | 10/14 - 10/18 | LV | 48 | 47 | Republicans +1 | |
| Monmouth | 10/13 - 10/17 | 756 RV | 50 | 44 | Republicans +6 | |
| CNBC | 10/13 - 10/16 | 800 RV | 48 | 46 | Republicans +2 | |
| Harvard-Harris | 10/12 - 10/13 | LV | 53 | 47 | Republicans +6 | |

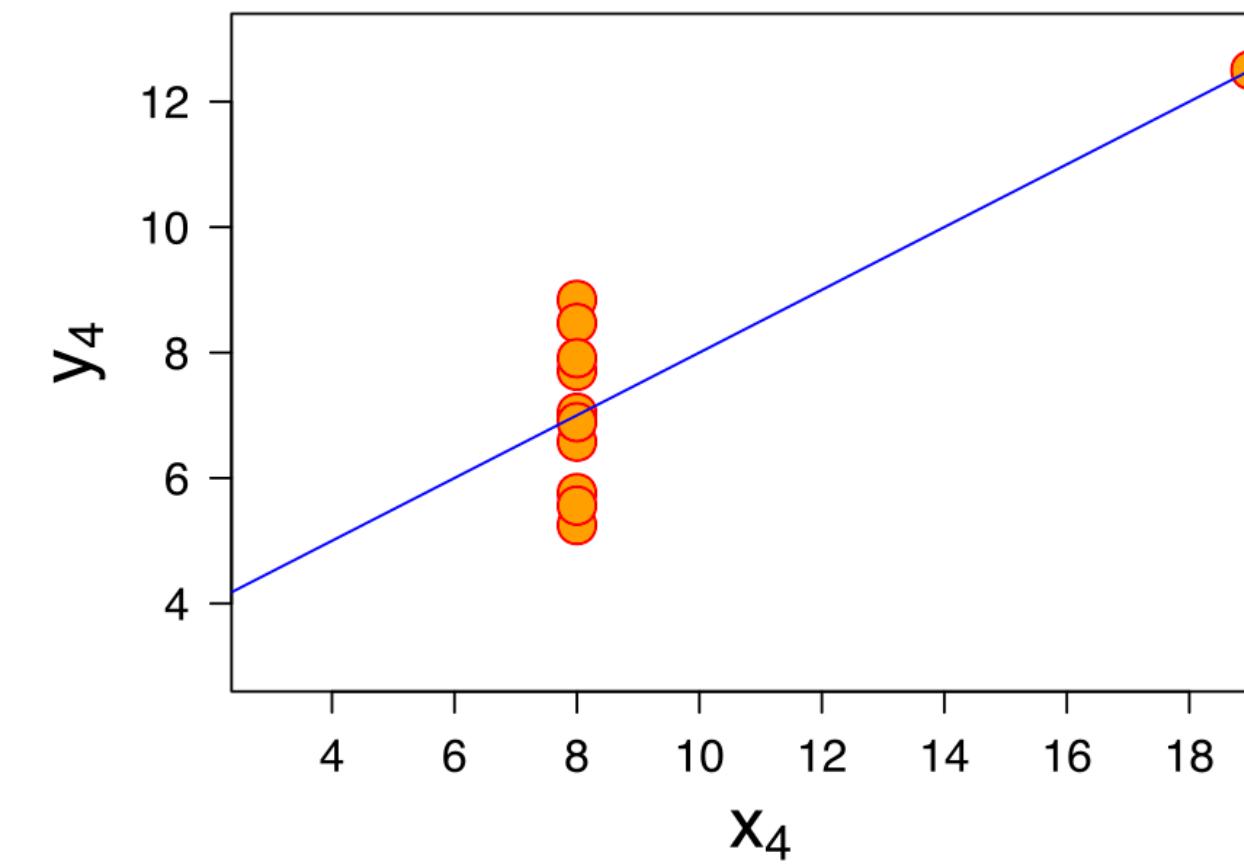
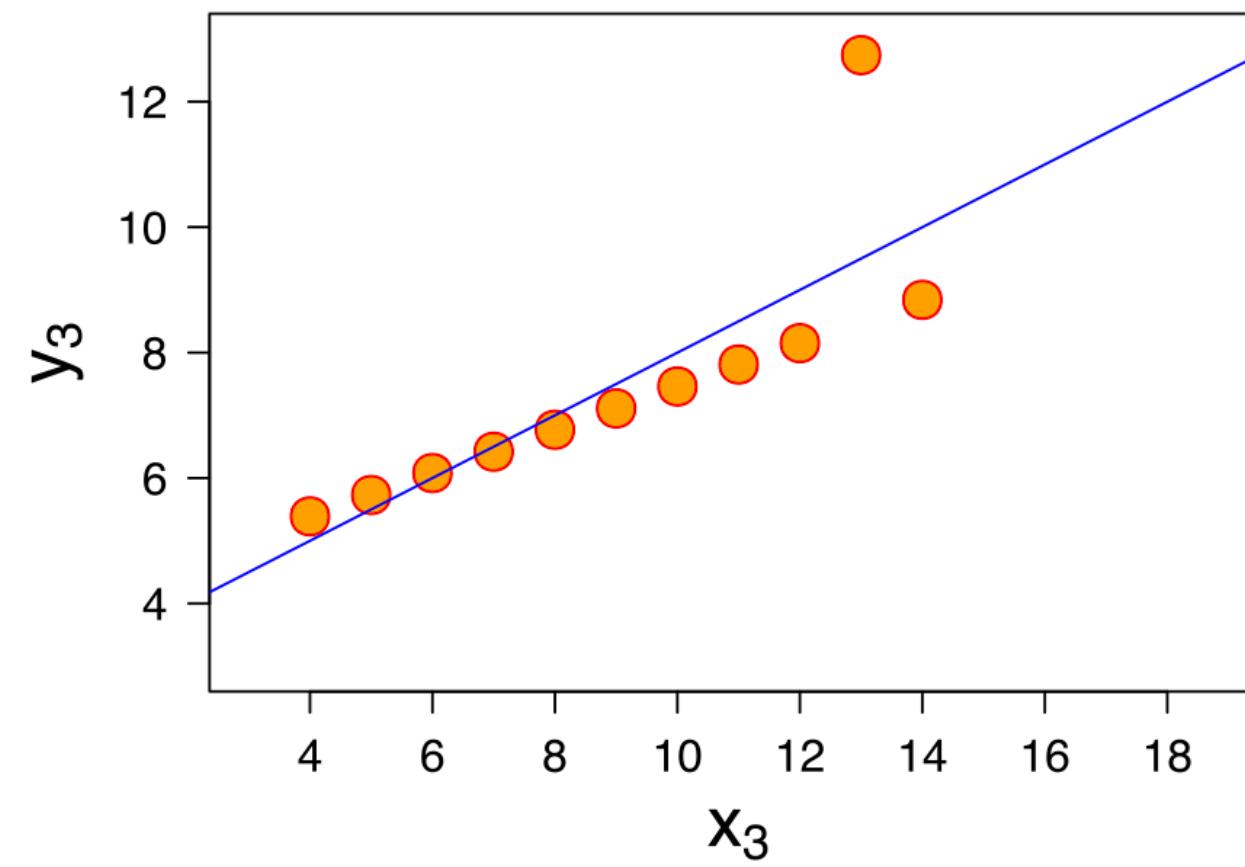
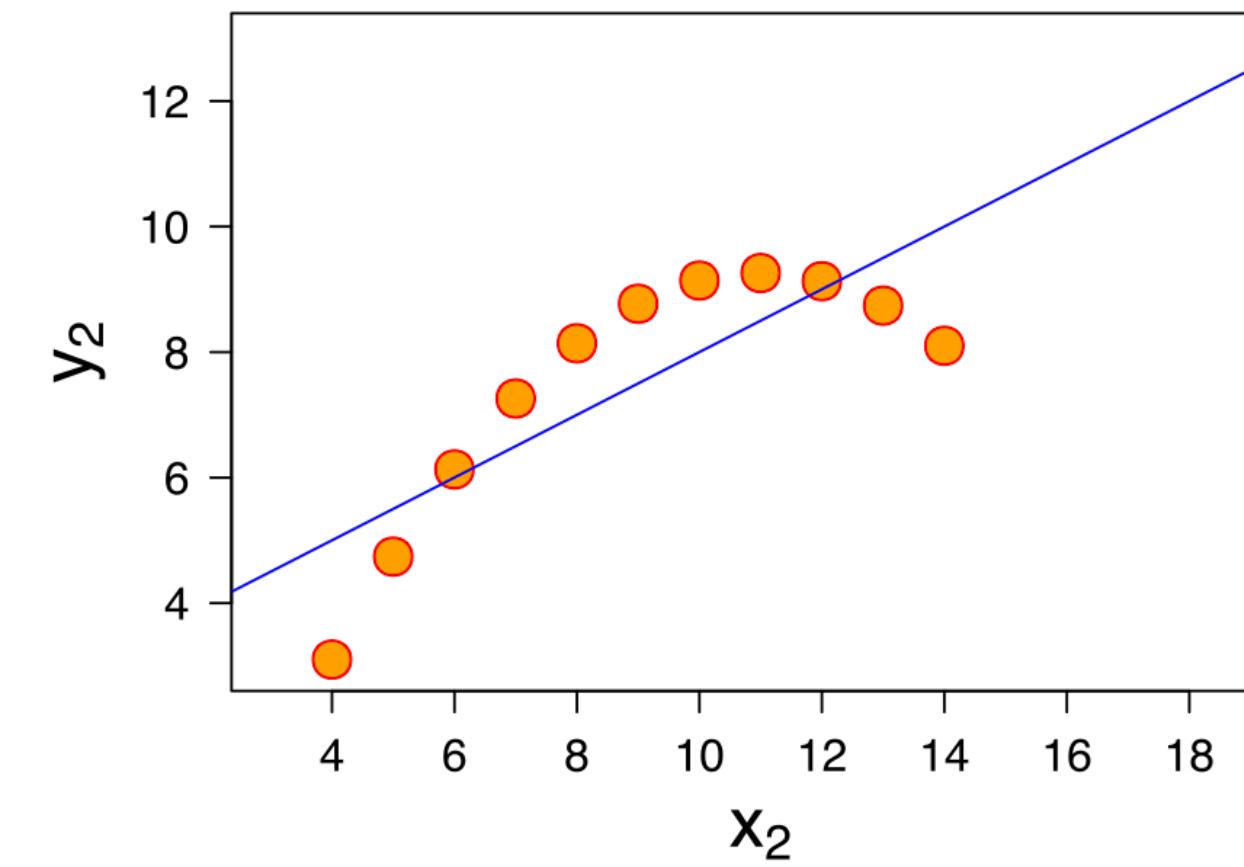
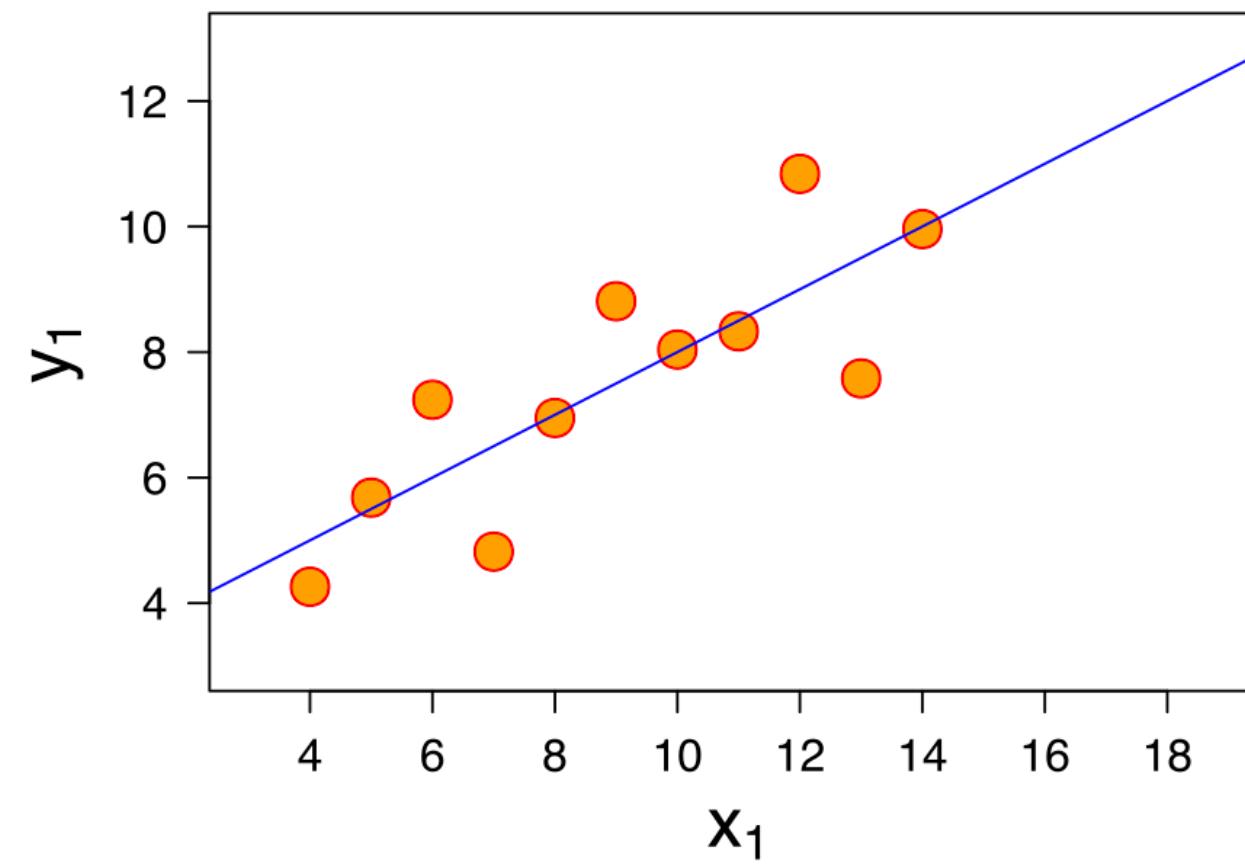
**Sometimes, summary statistics
obscure patterns that are
hidden in the data**

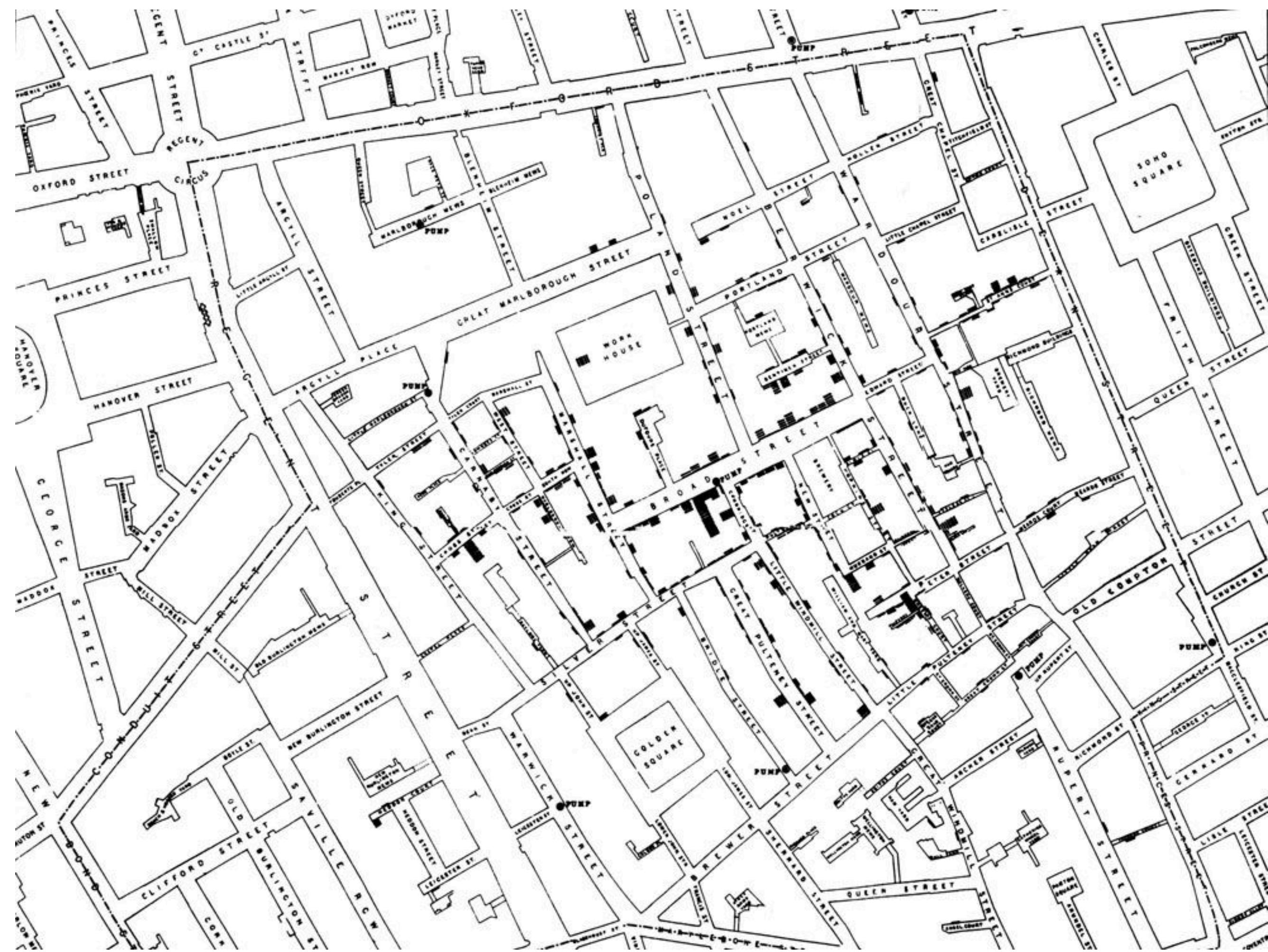
Example

- Anscombe's Quartet
 - Four data sets
 - Near identical on several metrics

| Property | Value | Accuracy |
|---|---------------------|---|
| Mean of x | 9 | exact |
| Sample variance of x : s_x^2 | 11 | exact |
| Mean of y | 7.50 | to 2 decimal places |
| Sample variance of y : s_y^2 | 4.125 | ± 0.003 |
| Correlation between x and y | 0.816 | to 3 decimal places |
| Linear regression line | $y = 3.00 + 0.500x$ | to 2 and 3 decimal places, respectively |
| Coefficient of determination of the linear regression : R^2 | 0.67 | to 2 decimal places |

Example

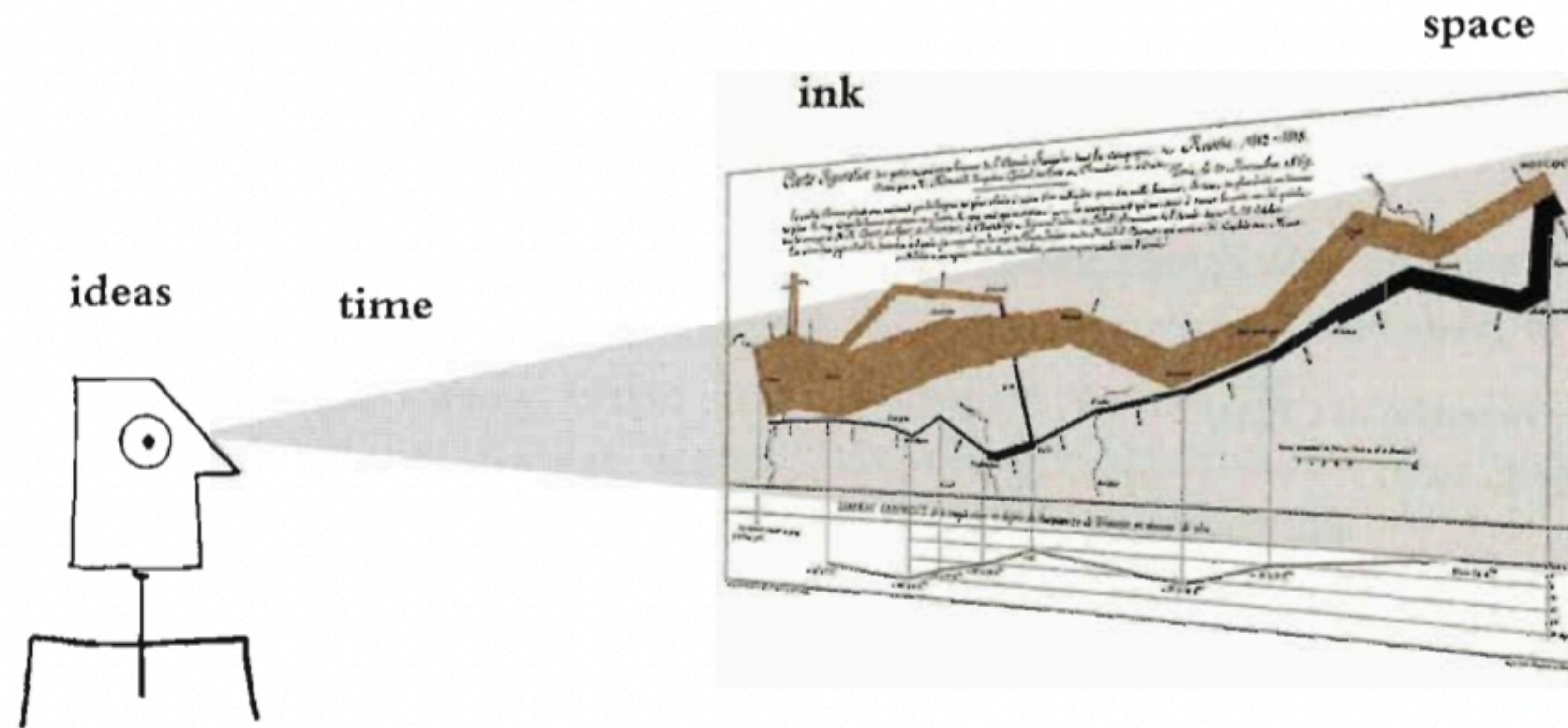




Data visualization

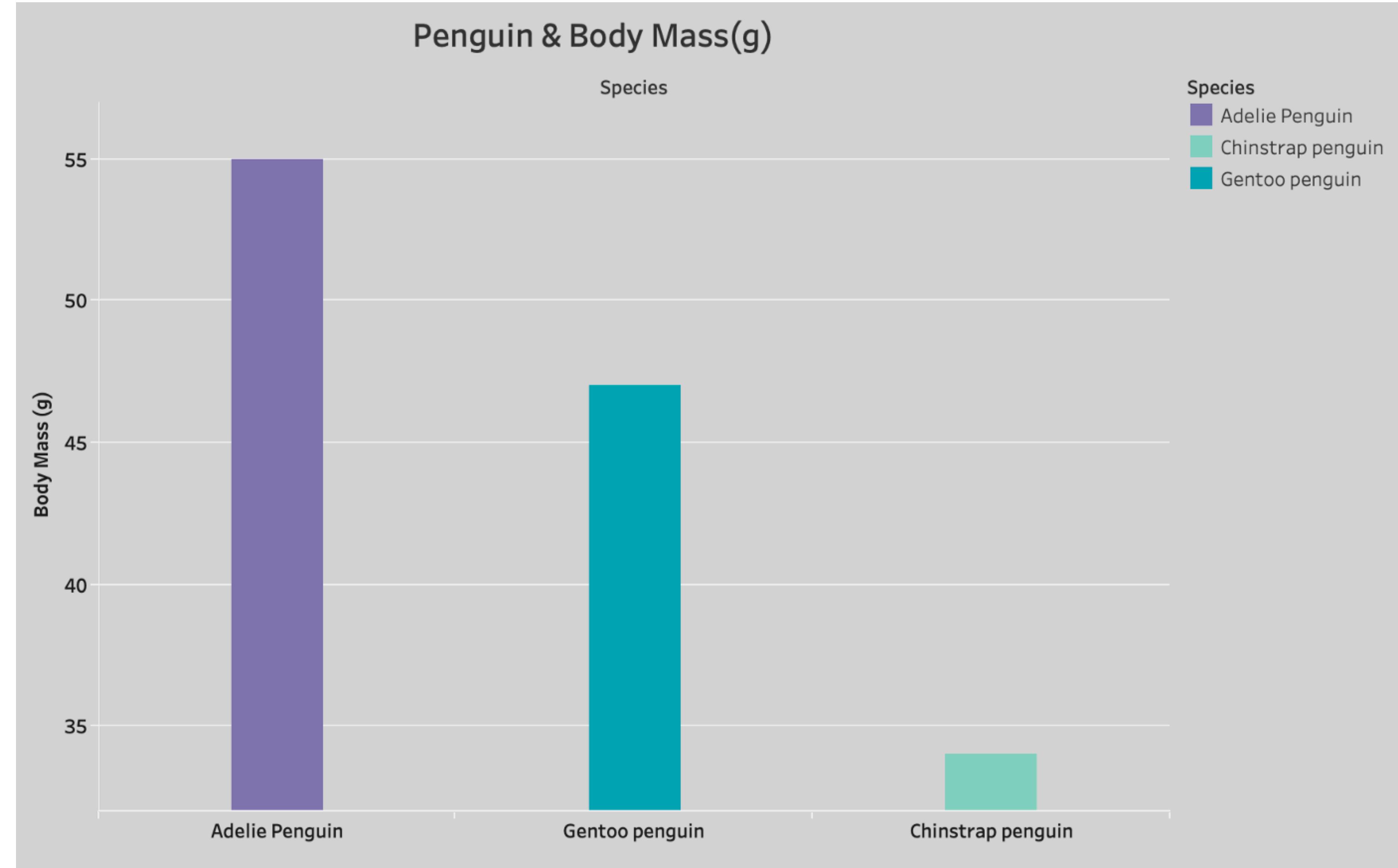
- What is data vis?
 - A visual representation of information
 - Graphics *reveal* data - Edward Tufte
- Principles of graphical excellence
 - Communicate complex ideas with clarity, precision, and efficiency
 - Viewer should get the greatest number of ideas in the shorter amount of time with the least ink in the least amount of space

Data visualization

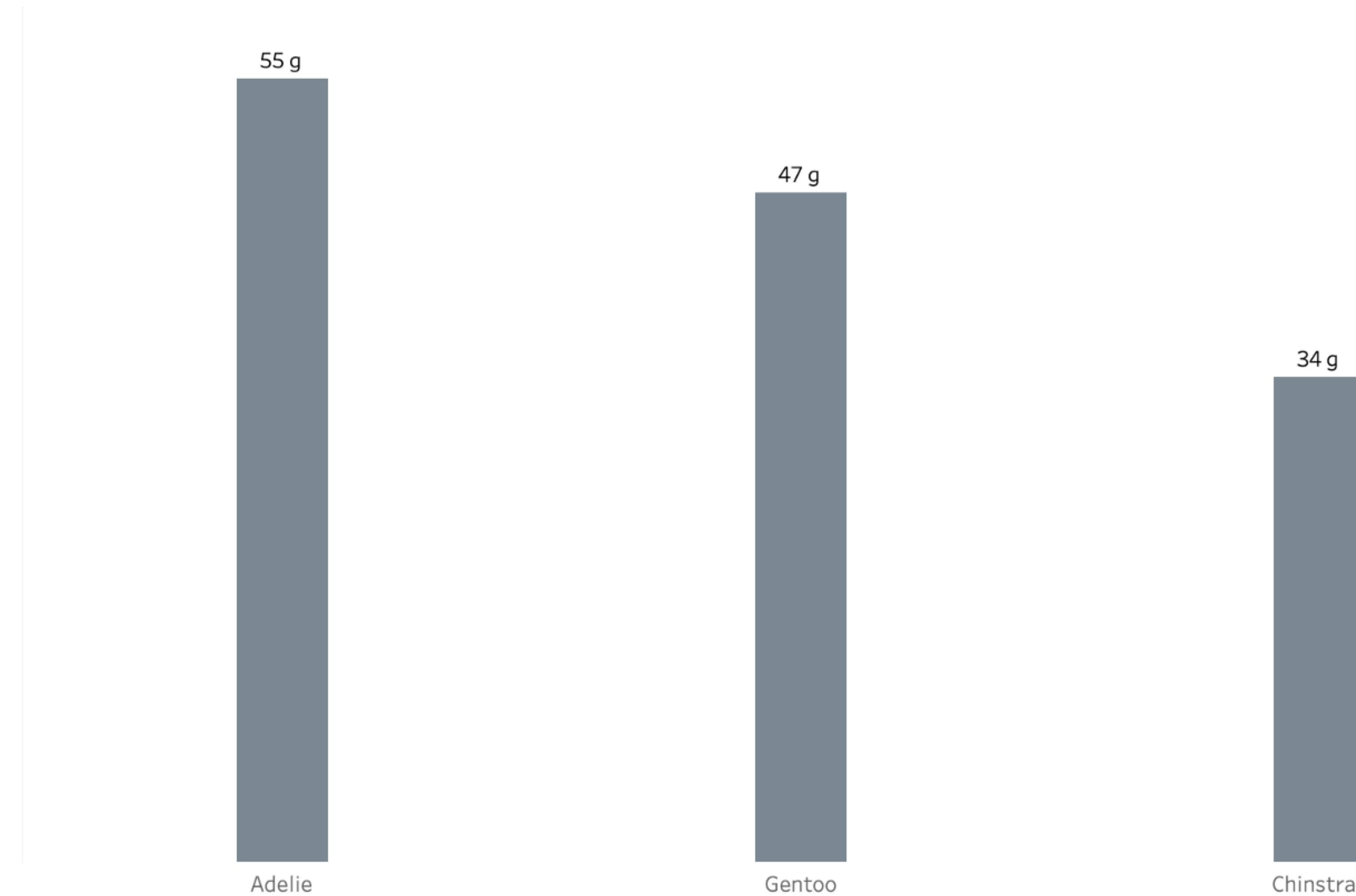


Data-ink Ratio

- Tufte argues that good data vis should maximize of the data-ink ratio
 - Data-ink: portion of ink that we can attribute to the data itself
 - Total-ink: comprised of data-ink and non-data ink.
 - Given today's technology, we can replace ink with pixels
- This is a principle of parsimony. Only visualize what is essential.



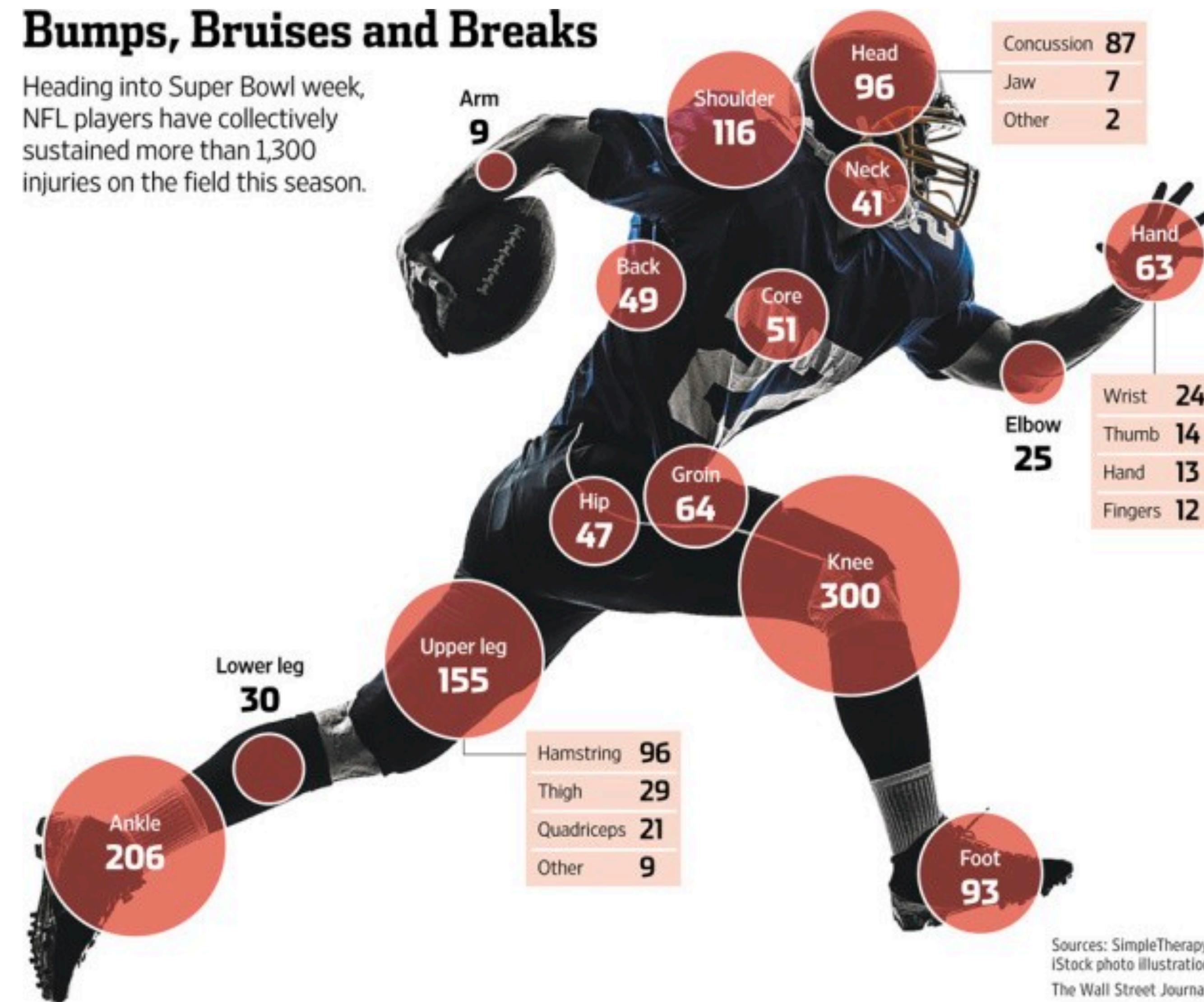
Penguin Species & Body Mass



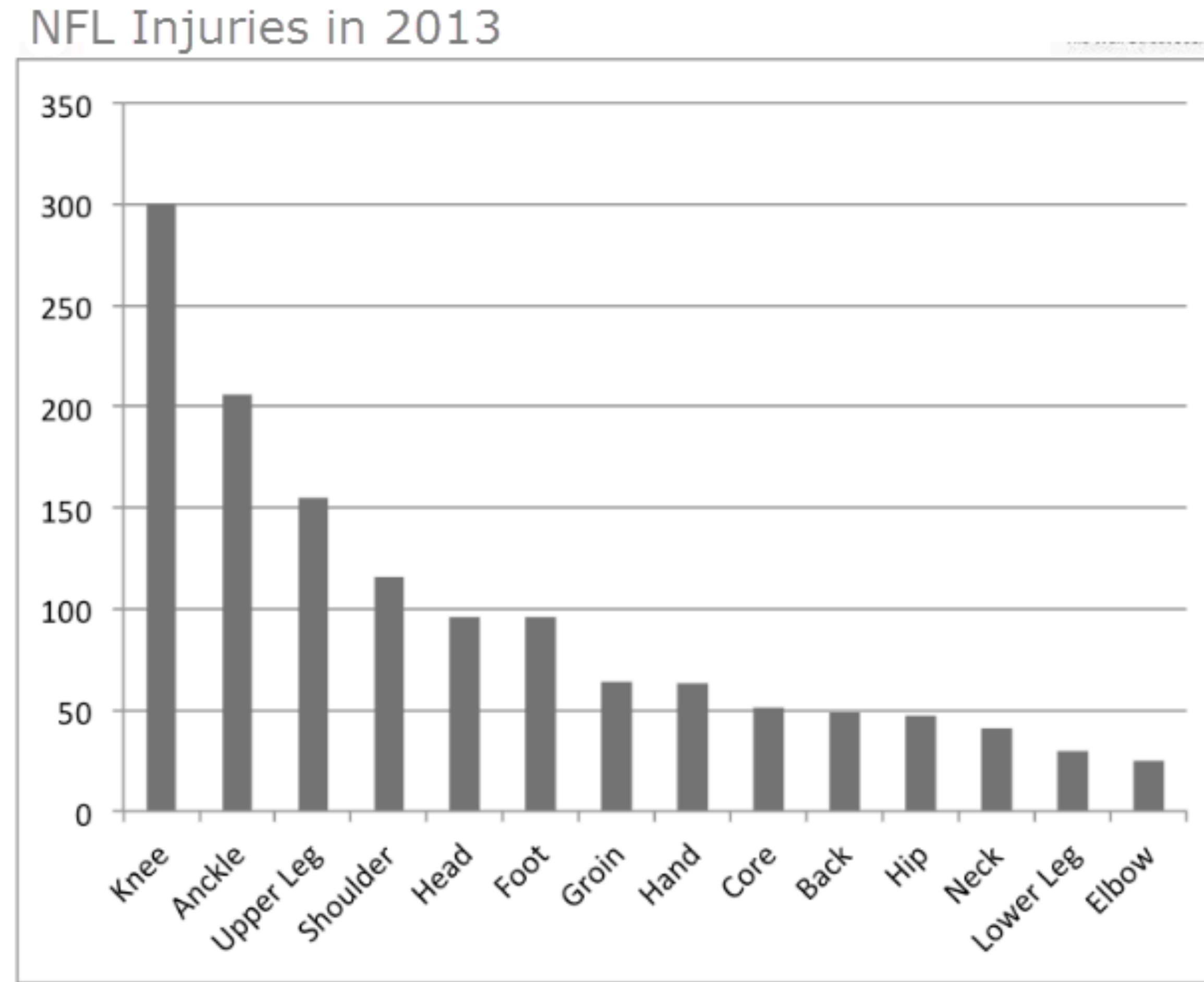
Minimize “chart junk”

Bumps, Bruises and Breaks

Heading into Super Bowl week, NFL players have collectively sustained more than 1,300 injuries on the field this season.



Minimize “chart junk”



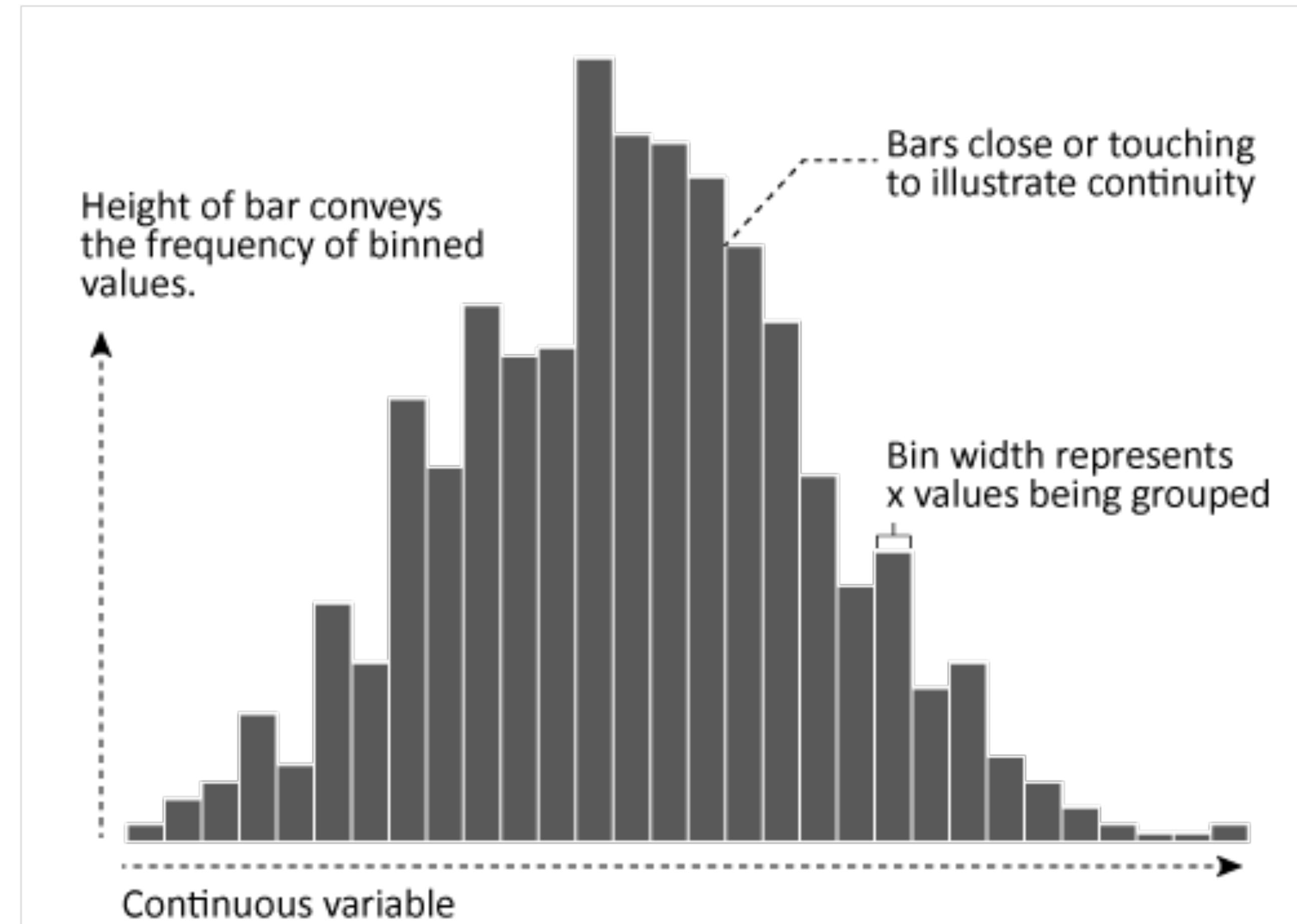
Different Visualizations

- Now that we have some criteria to evaluate visualizations, what do different visualizations look like? What do they accomplish?

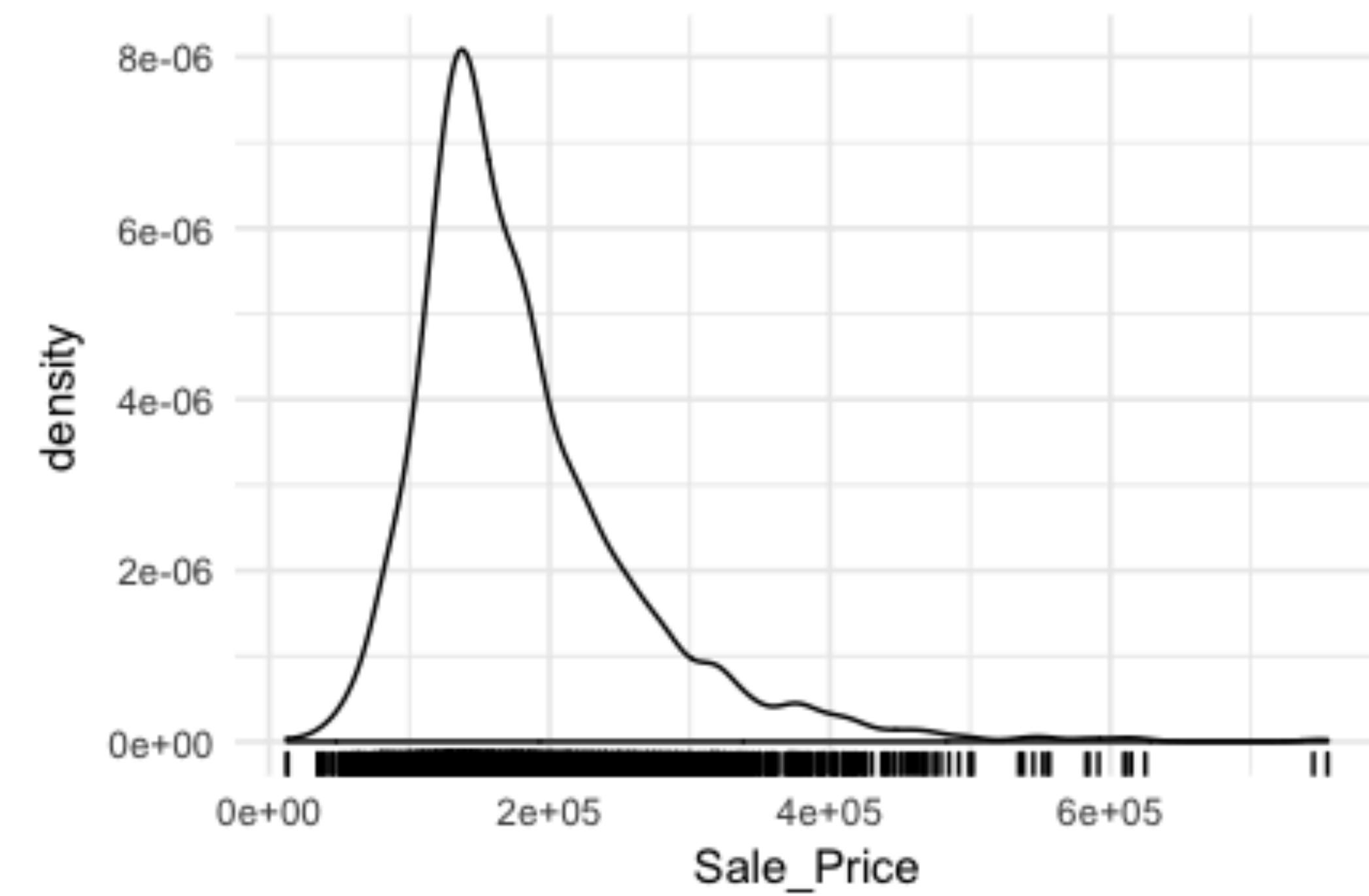
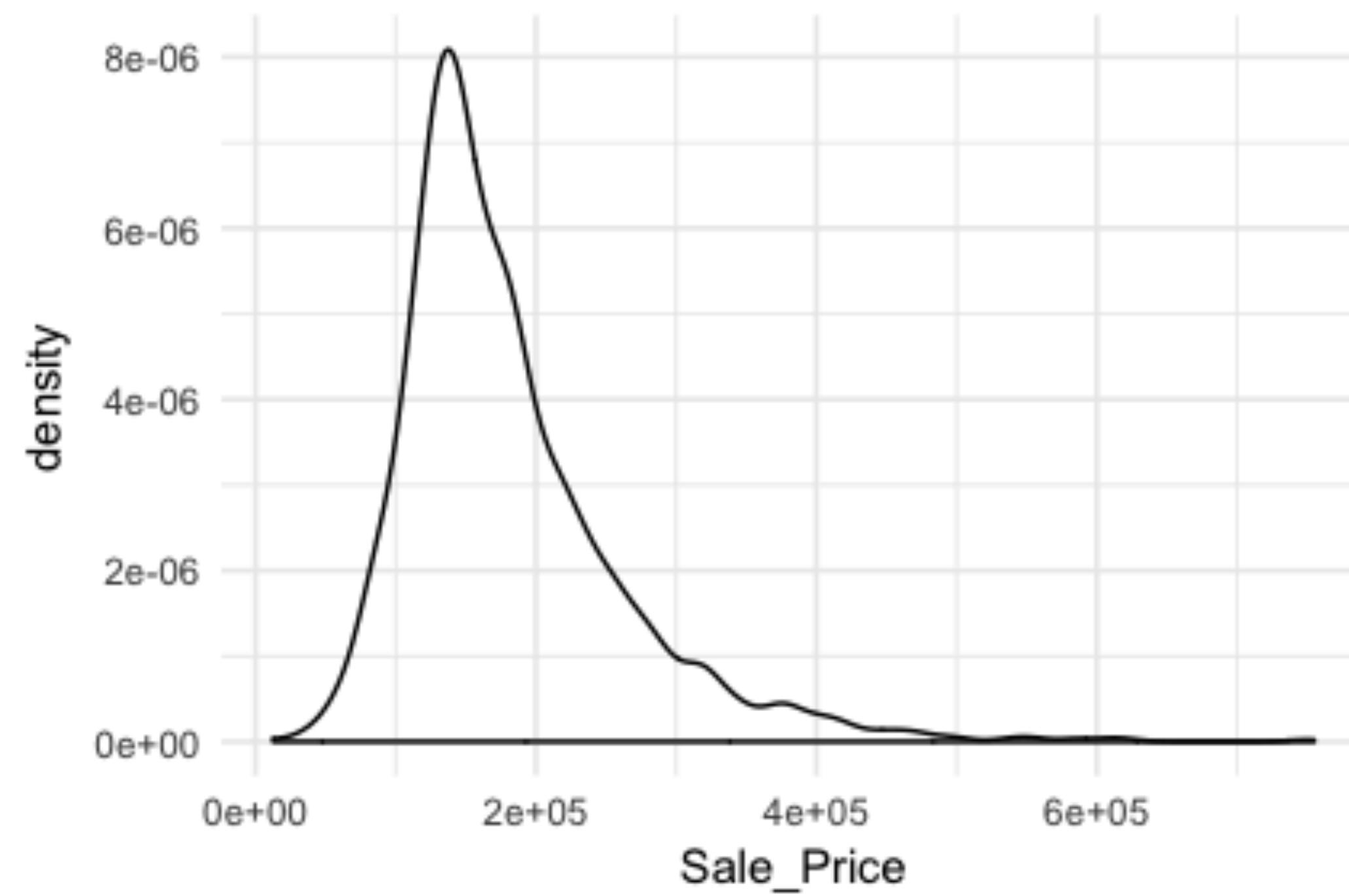
Univariate Plots

- The following plots help us summarize a single variable
- How is it distributed?
- Which values are most common?

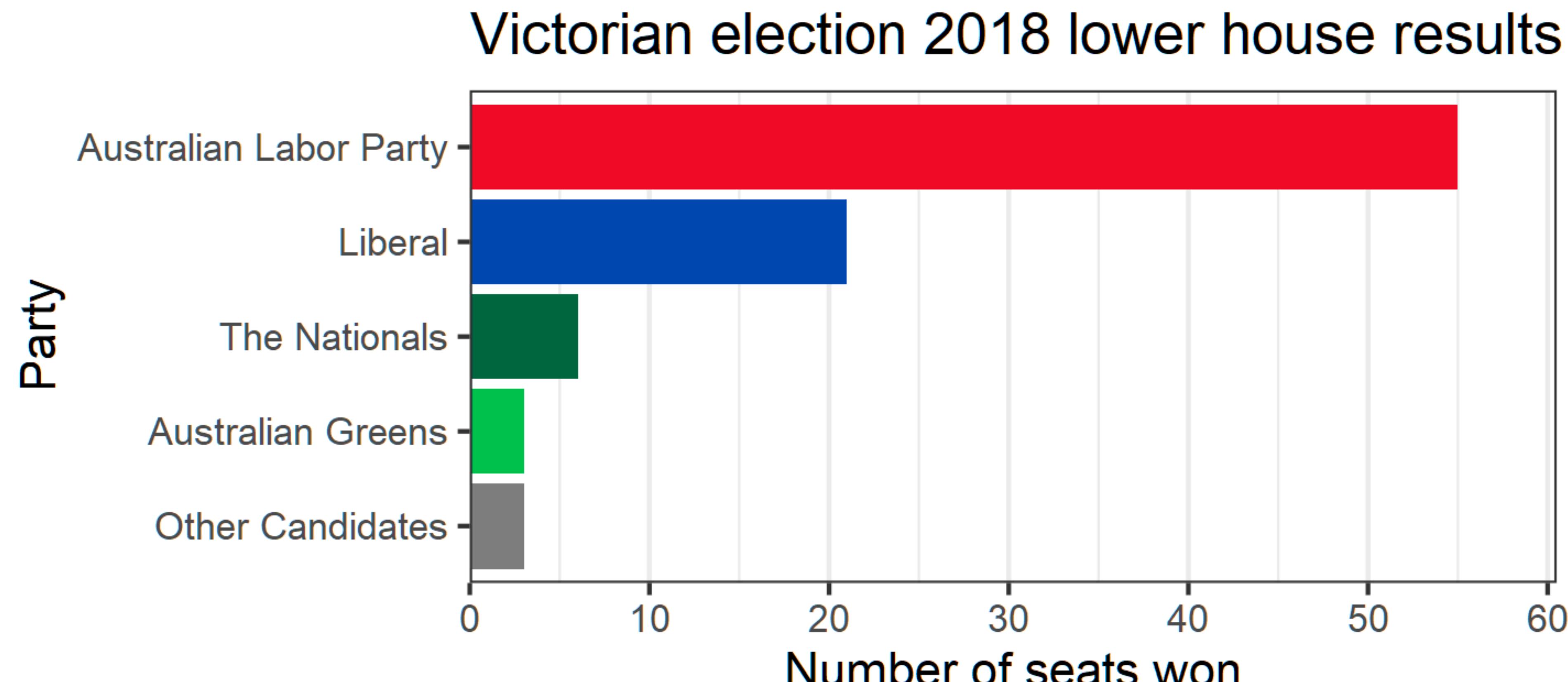
Histogram



Density

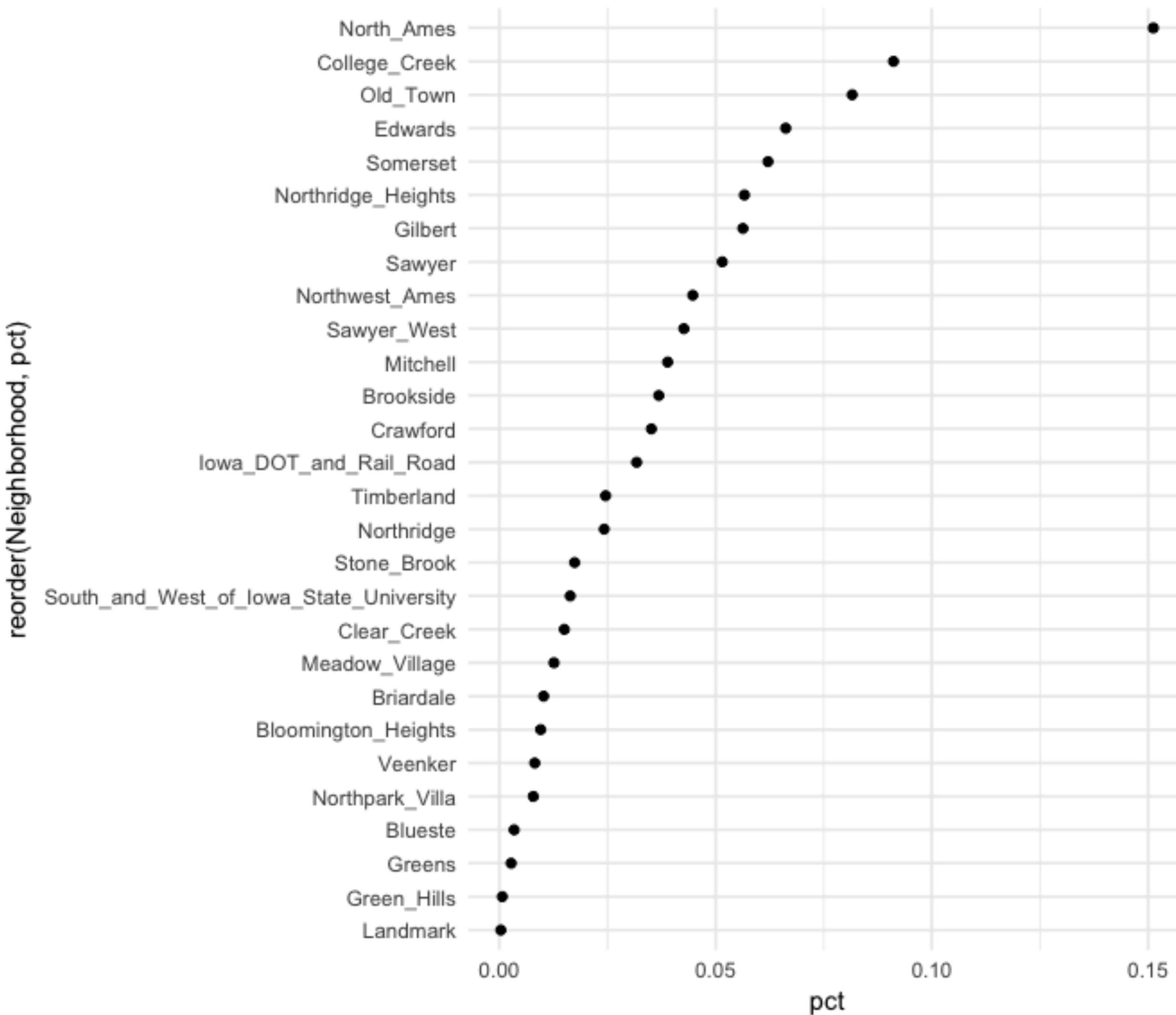


Bar Plots

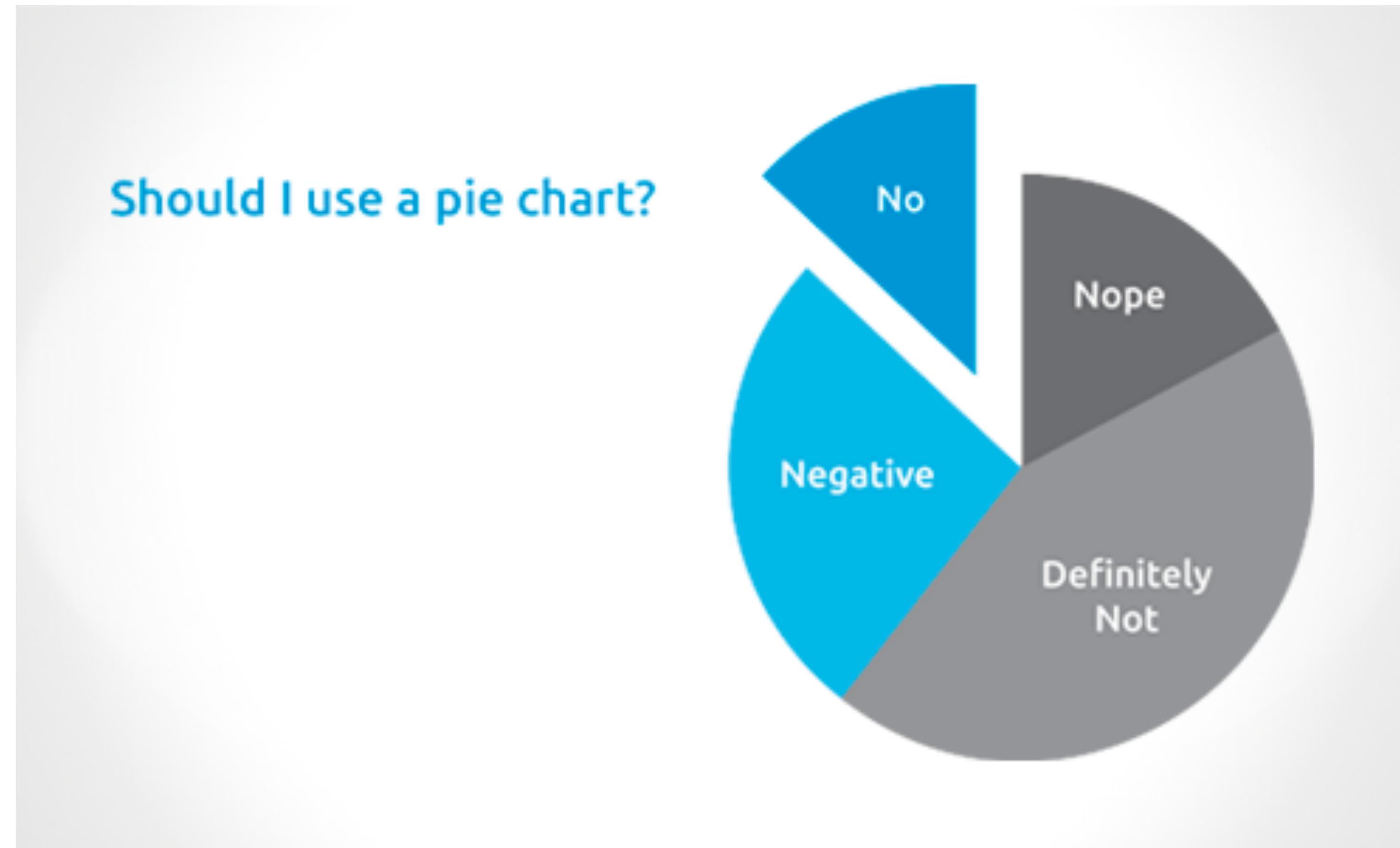


Data source: Victorian Electoral Commission

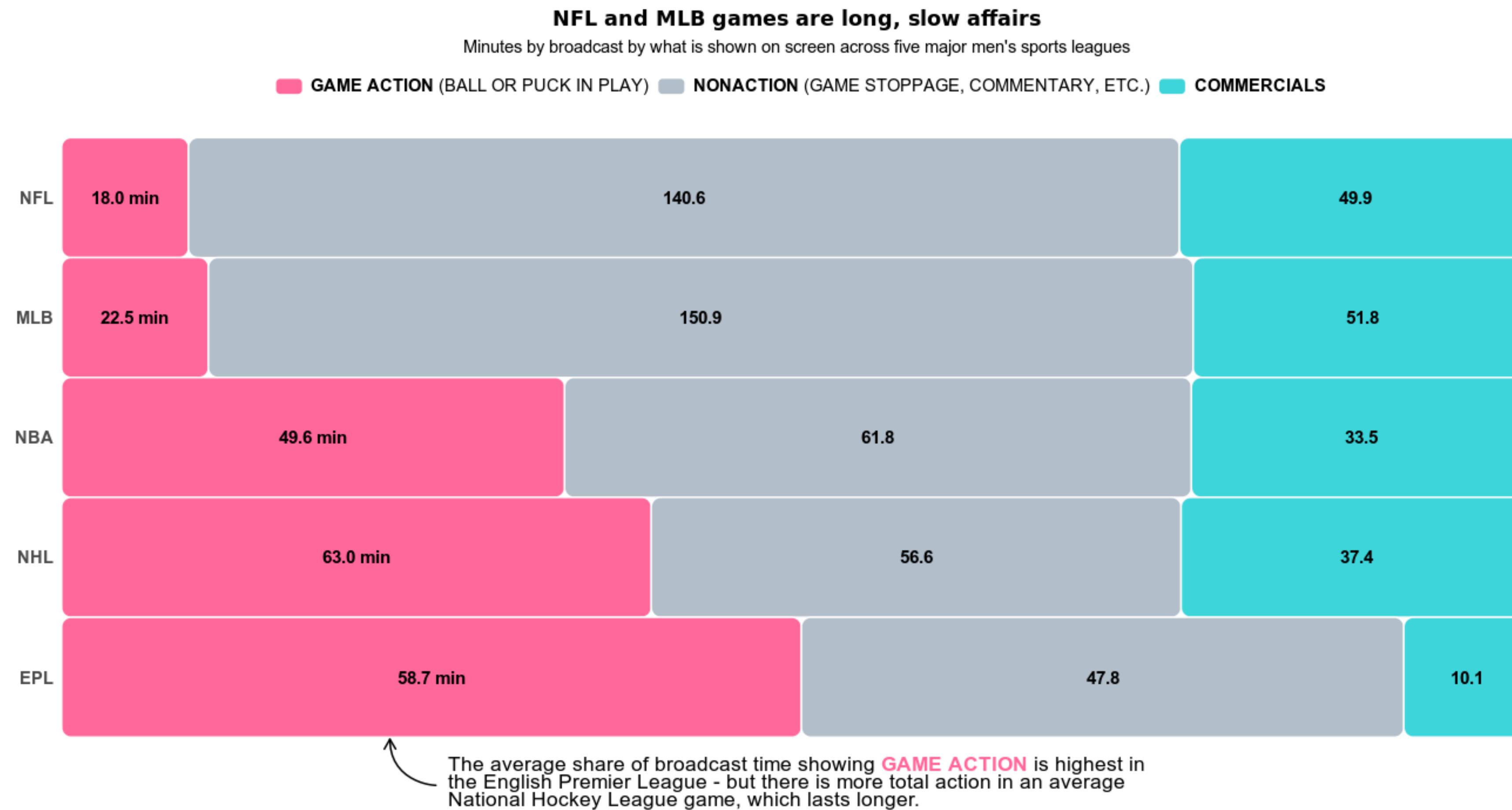
Dot Plot



Pie Chart



Stacked Bar Chart



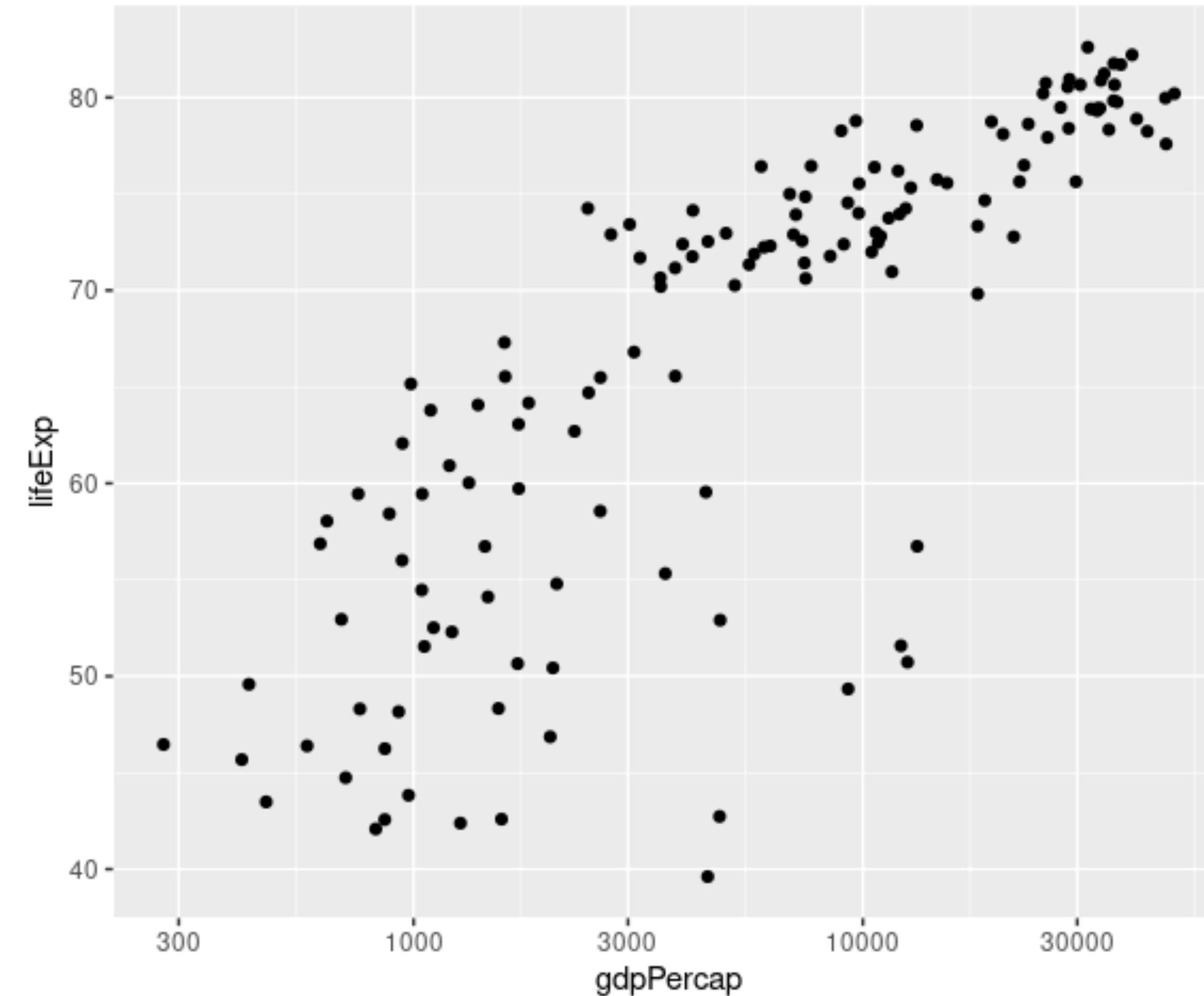
Games that we included: 10 NFL regular-season games between Nov. 7 and Nov. 18, 2019; 17 MLB postseason games, including all the games in the 2019 ALCS, NLCS, and World Series; 10 NBA regular-season games between Nov. 6 and Nov. 15, 2019; 10 NHL regular-season games between Nov. 5 and Nov. 19, 2019, including three overtime games; and seven English Premier League games between Nov. 9 and Nov. 23, 2019. NBA game action includes free throws, so the action time exceeds the game time.

FivethirtyEight SOURCE: UNIVERSITY OF TEXAS AT AUSTIN SPORTS ANALYTICS COURSE

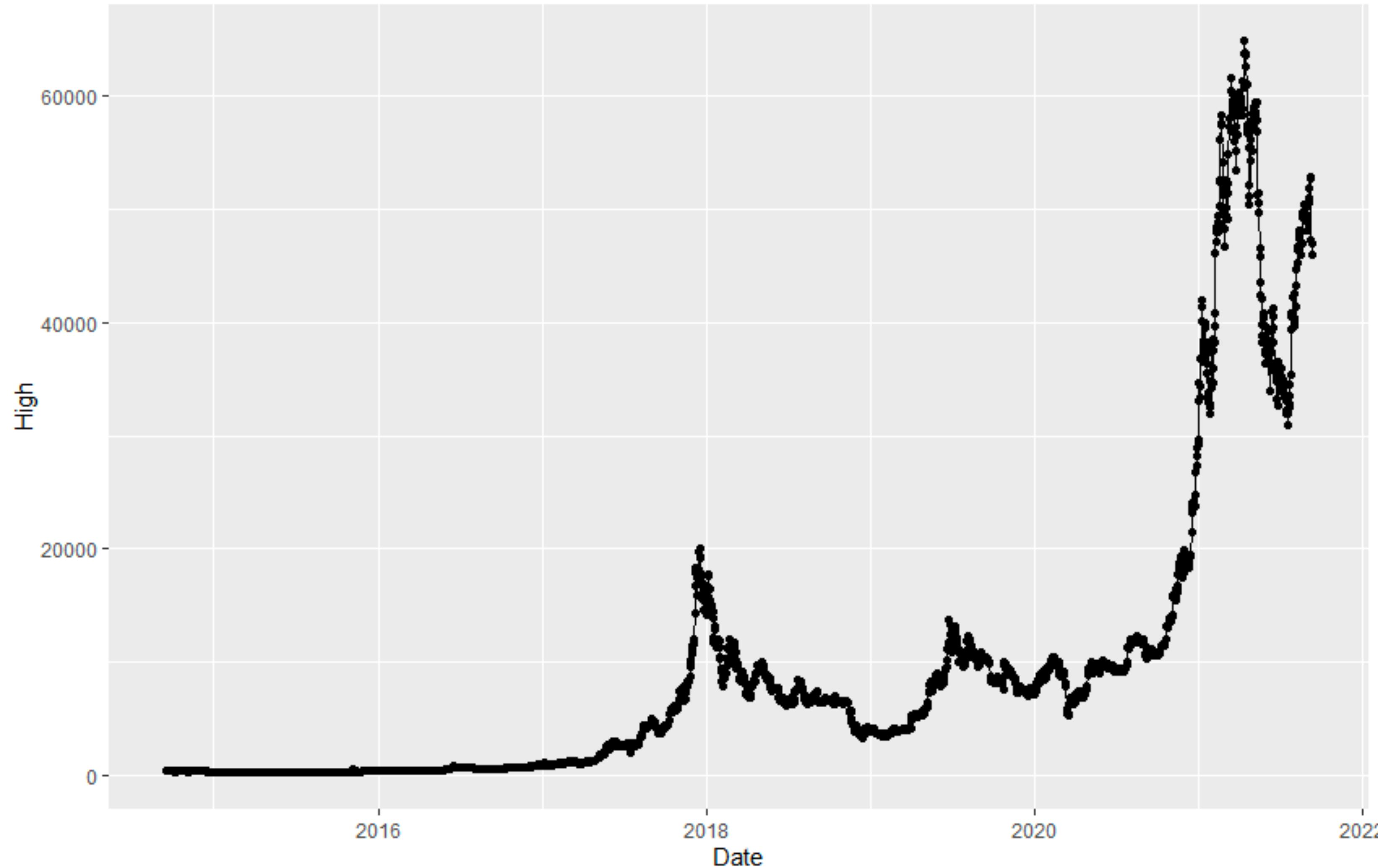
Multivariate plots

- Plots that summarize the relationship between two or more variables

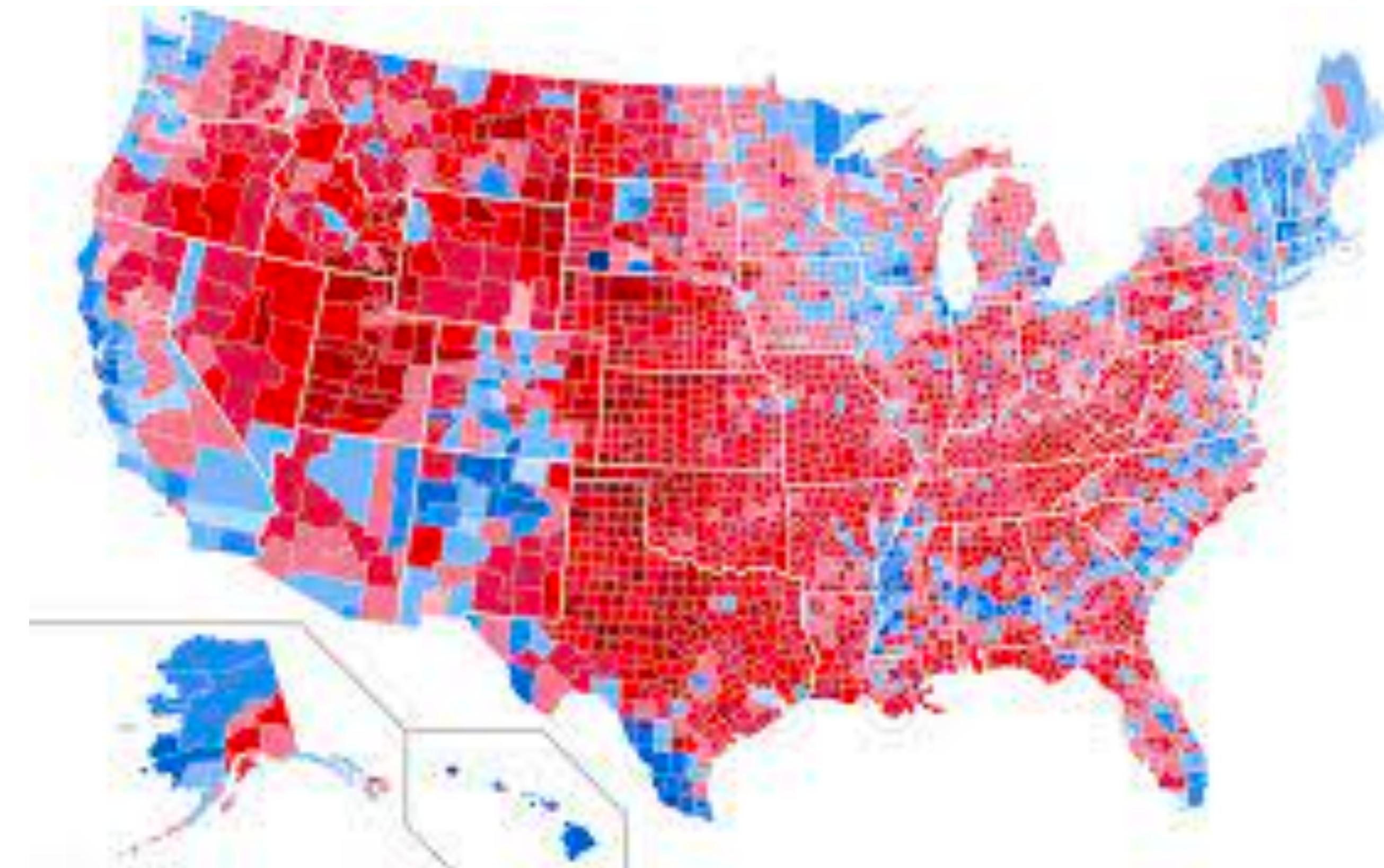
Scatterplot



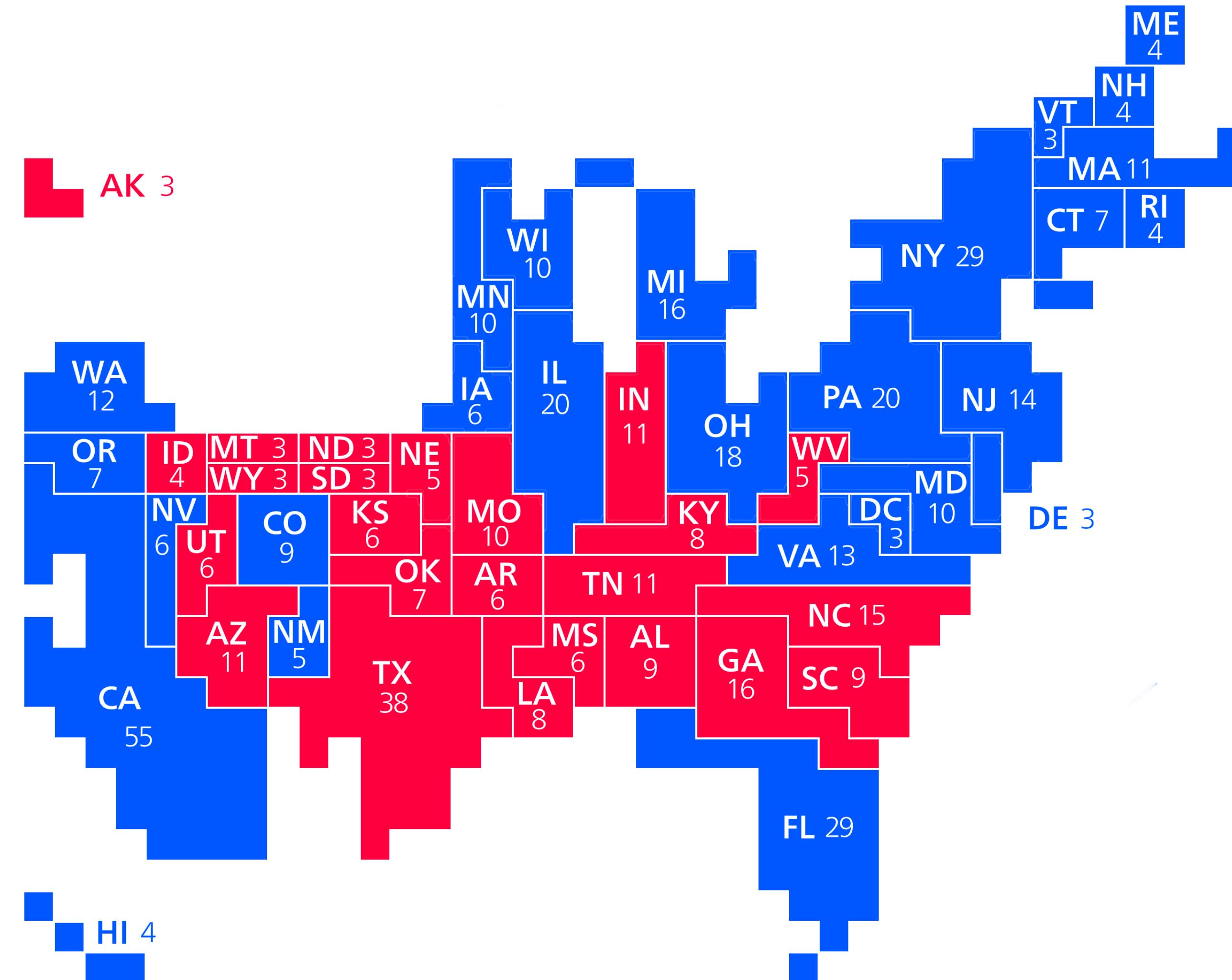
Time Series



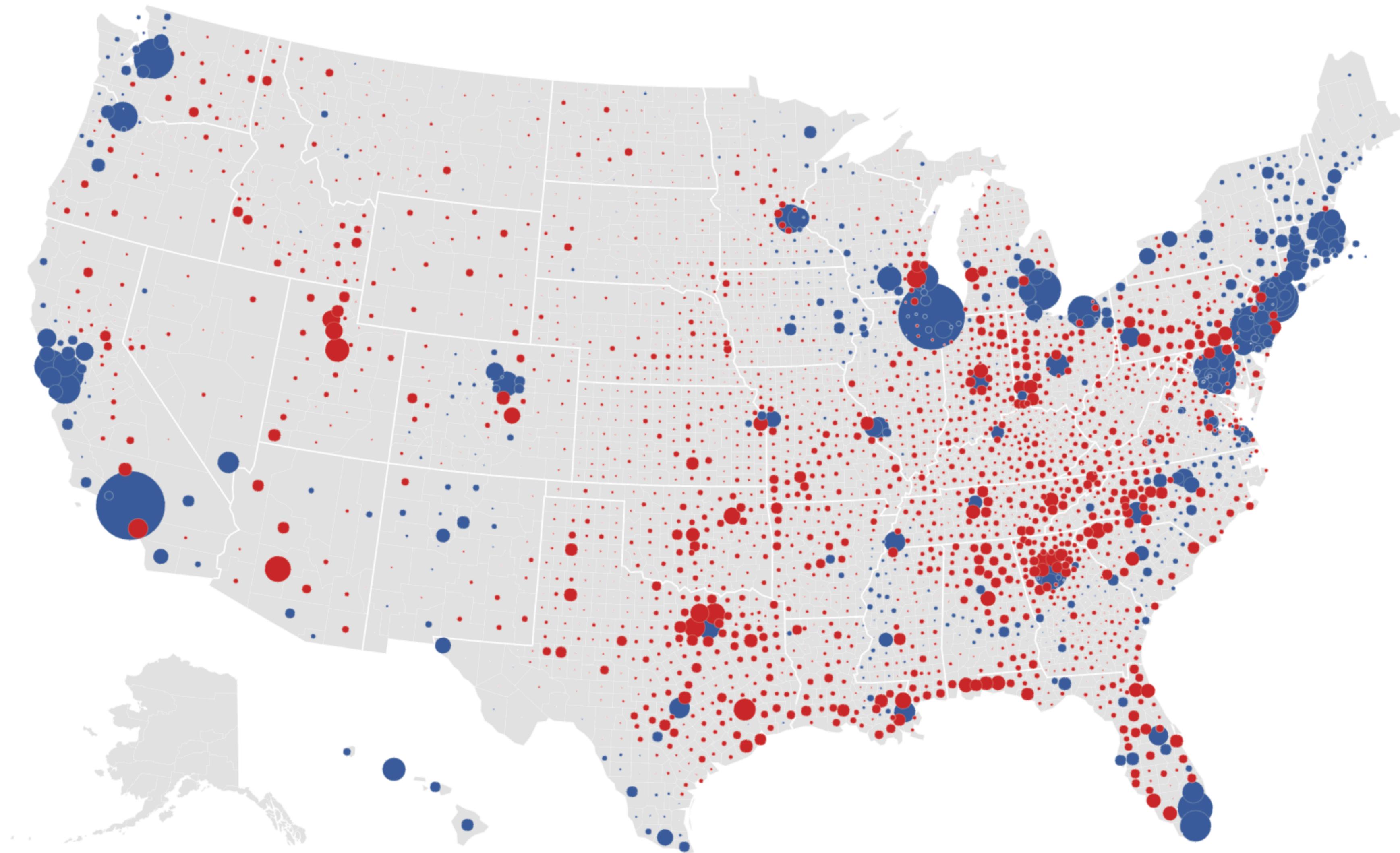
Choropleth Map



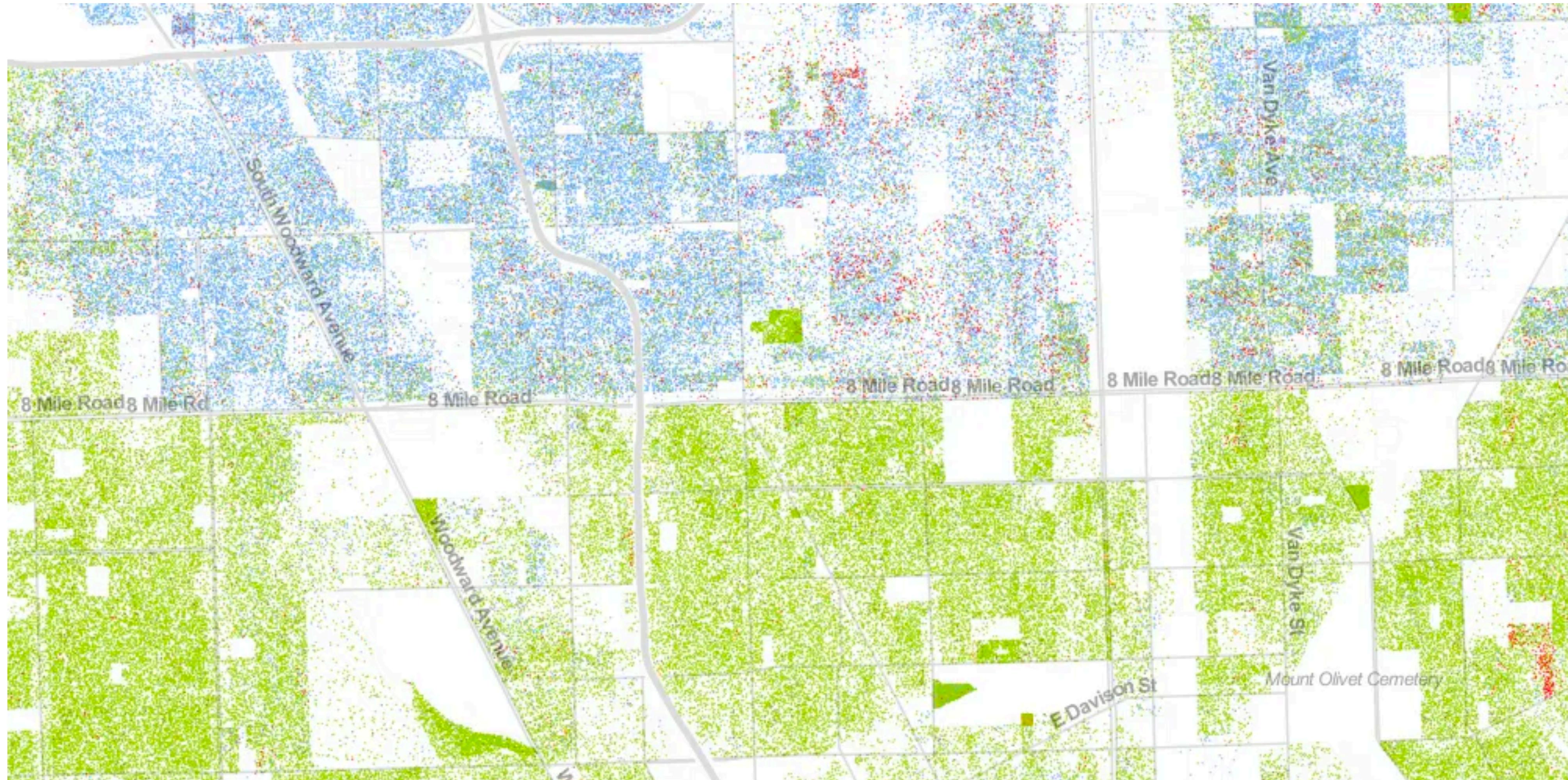
Cartogram



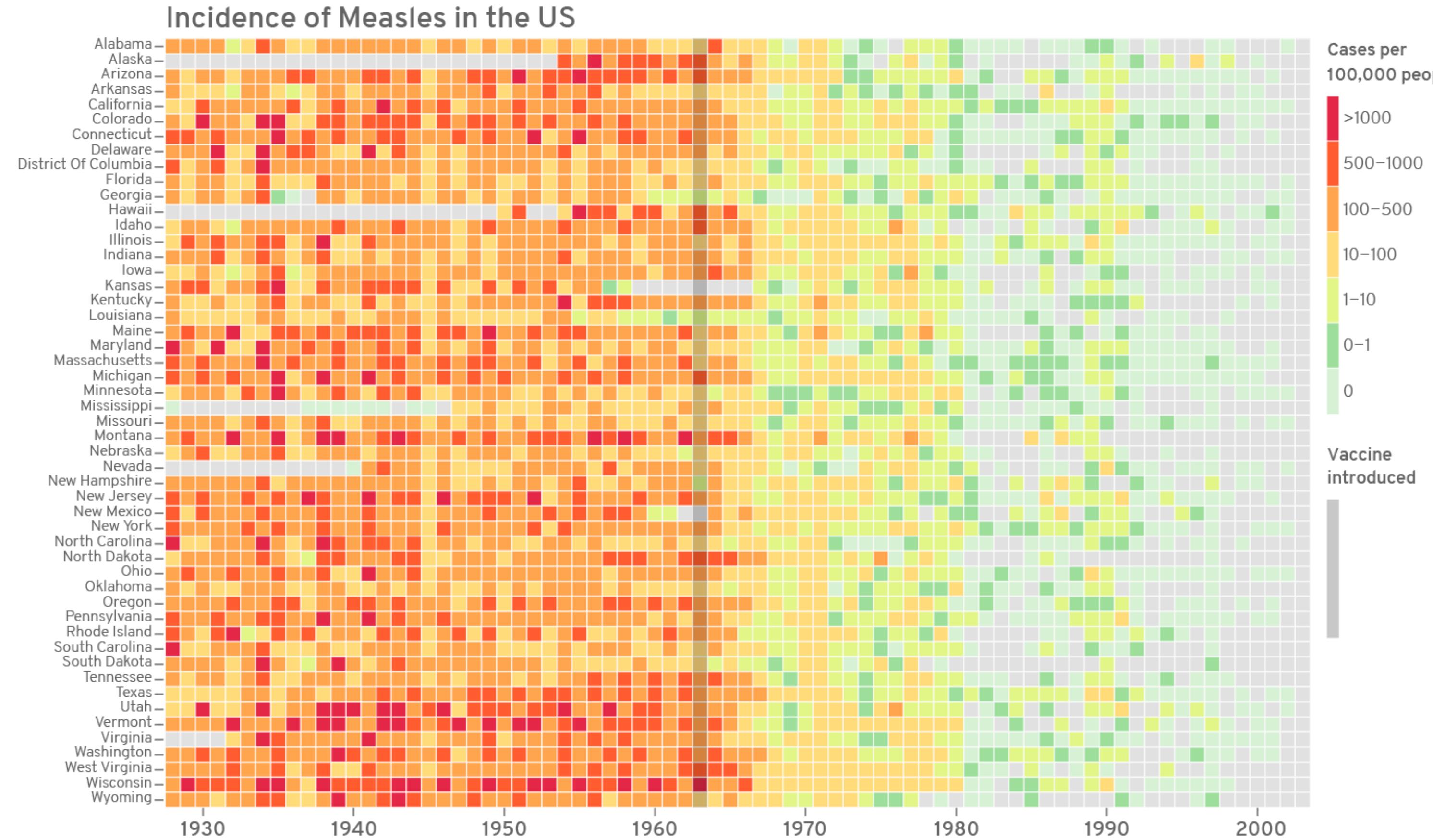
Proportional Symbol Maps



Dot Distribution Map



Heatmap



Group Assignment