

Visual Analytics— Visualization 101

Dr. Ab Mosca (they/them)

Slides based off slides courtesy of Jordan Crouser (<https://jcrouser.github.io/>)

Plan for Today

- Visualization overview
- Graphical primitives
- Visual dimensions
- Common visualization techniques
- *Time permitting*: Tableau demo

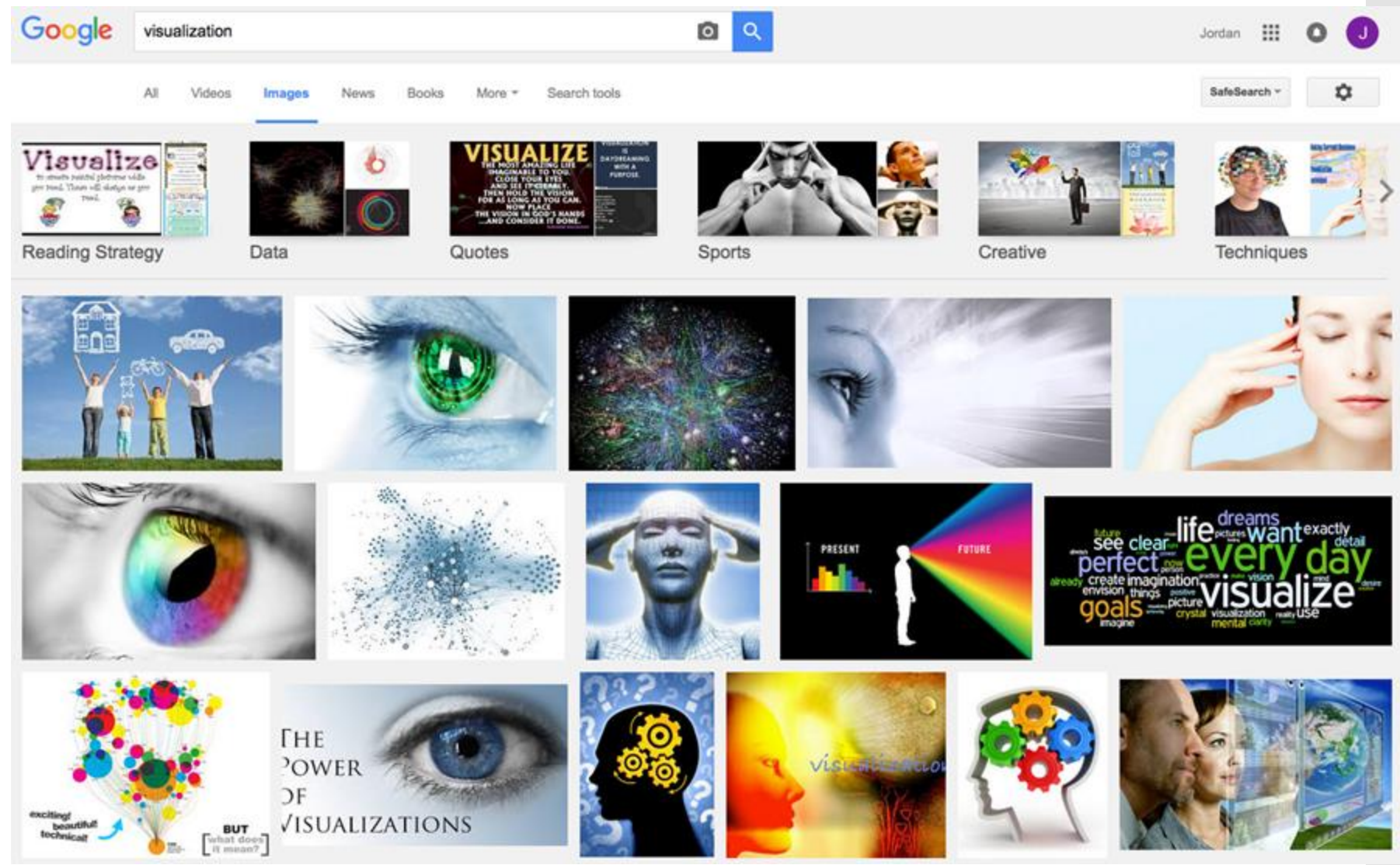
Reminder

- hwo1 is released today!
(<https://amoscao1.github.io/SDS-CS235/> > Homework > hwo1)
 - Due next Thursday (09/19) at midnight
 - You have extensions if you need them, but you MUST tell me if you're taking one
 - Revise and resubmit also exists!
 - Work with a small group (3-4) – I recommend finding people with *complementary skillsets* to work with
- **I won't always remind you there is a homework released/due, make sure you stay up to date with the course schedule! ****

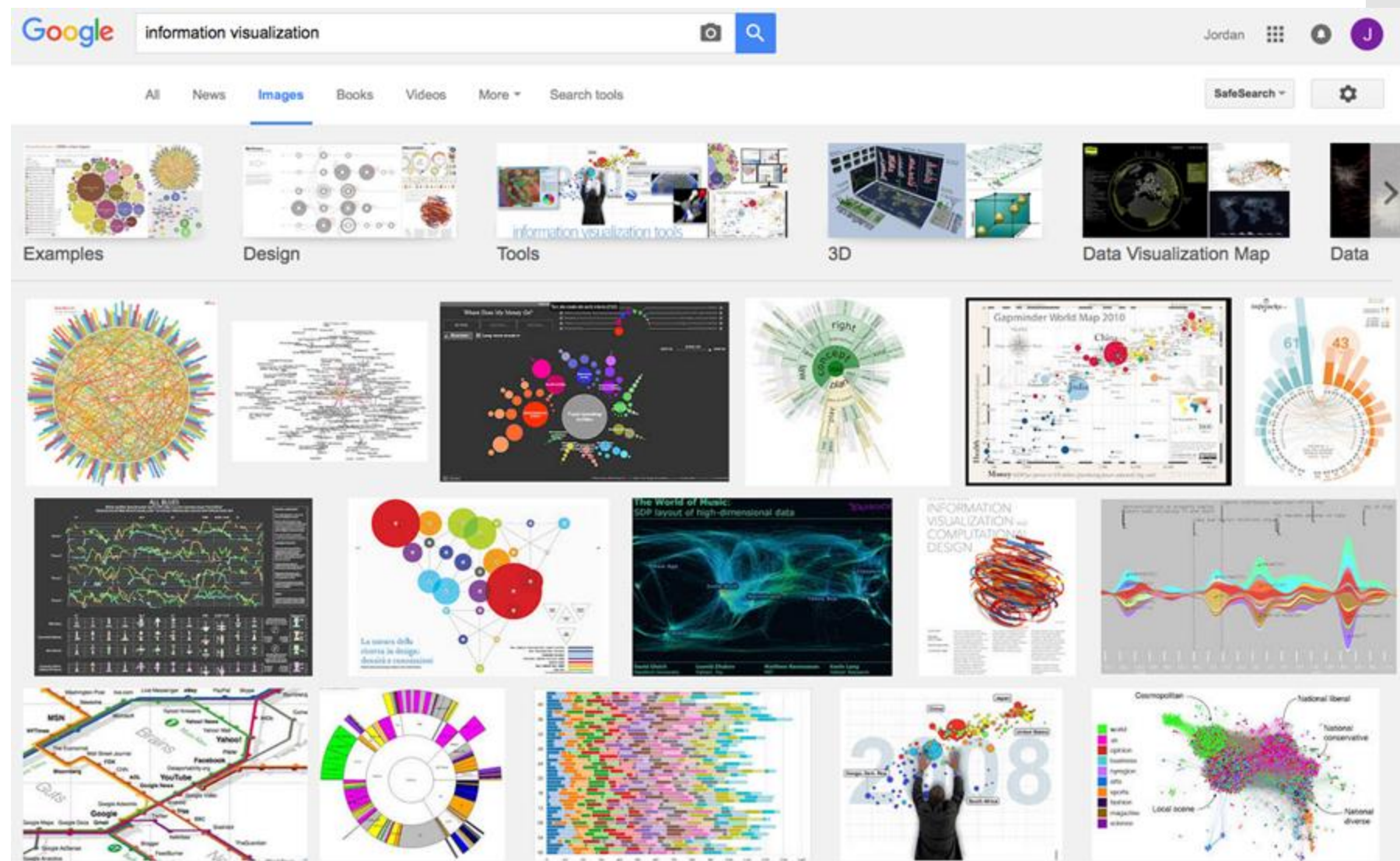
Visual Analytics

Visualization \cap Analysis

What is visualization?



What is visualization?



Perhaps a
more helpful
question:

What are some ways
a “visualization” can be **useful**?

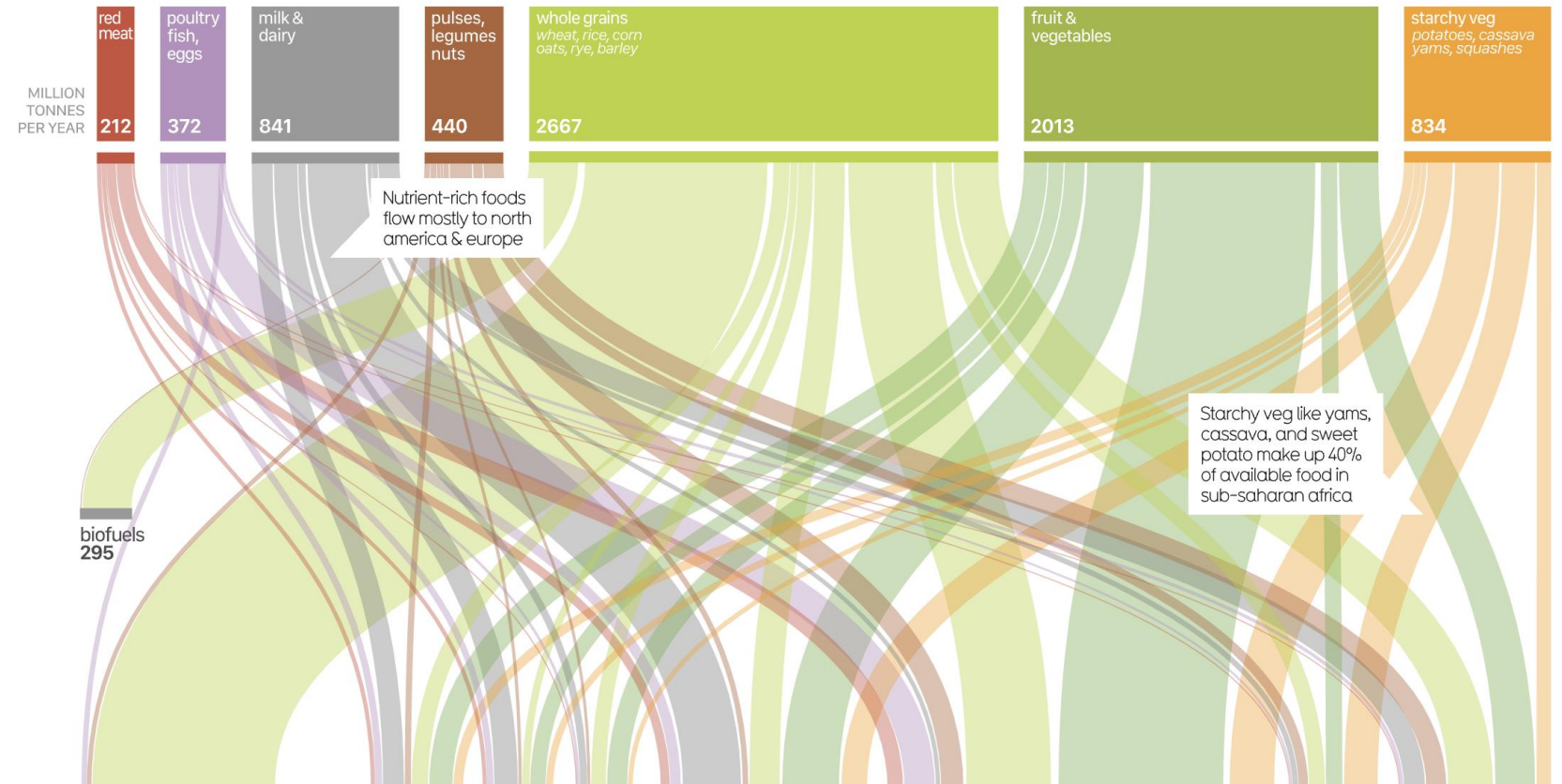
Does it help
you spot
trends?



More info here: http://en.wikipedia.org/wiki/1854_Broad_Street_cholera_outbreak

Does it help you explore?

How much do we make?



Who gets the food?

Farmed animals eat

<https://informationisbeautiful.net/visualizations/global-food-supply-where-does-all-the-worlds-food-go/>

Does it tell a
story?



Visualization
(def.)

Visual
representations
of data that
reinforce human
cognition



Wait... what is
“data”?



Data: a definition

Data is a set of *variables* that capture various aspects of the world:



*Tuition rates, enrollment numbers,
public vs. private, etc.*

Data: a definition

A dataset also contains a set of *observations* (also called *records*) over these variables. For example:



tuition = \$46,288, *enrollment* = 2,563,
private, etc.

Data: a definition

A dataset also contains a set of *observations* (also called *records*) over these variables. For example:



tuition = \$16,115, *enrollment* = 28,635,
public, *etc.*

One way to
think about this:

		VARIABLES			
		Tuition	Enrollment	Public vs. Private	...
OBSERVATIONS	Smith College	\$46,288	2,563	private	
	UMass Amherst	\$16,115	28,635	public	
	Hampshire College	\$48,065	1,400	private	
	Mount Holyoke College	\$43,886	2,189	private	
	Amherst College	\$50,562	1,792	private	
	⋮				

Data

- Remember...

country	year	cases	population
Afghanistan	1999	1725	19987071
Afghanistan	2000	1666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	1272915272
China	2000	216766	128012583

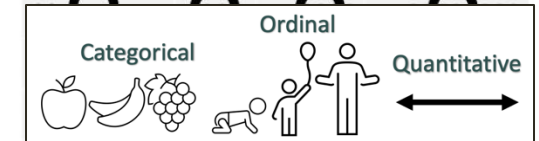
variables

country	year	cases	population
Afghanistan	1999	1725	19987071
Afghanistan	2000	1666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	1272915272
China	2000	216766	128012583

observations

country	year	cases	population
Afghanistan	1999	1725	19987071
Afghanistan	2000	1666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	1272915272
China	2000	216766	128012583

values



Data → Visuals

- Remember...

country	year	cases	population
Afghanistan	1999	21666	19987071
Afghanistan	2000	2666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127291272
China	2000	216766	128012583

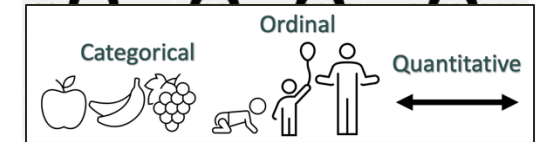
variables

country	year	cases	population
Afghanistan	1999	21666	19987071
Afghanistan	2000	2666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127291272
China	2000	216766	128012583

observations

country	year	cases	population
Afghanistan	1999	21666	19987071
Afghanistan	2000	2666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127291272
China	2000	216766	128012583

values



- **Big idea behind visualization**

- Data have dimensions
- Visualizations have dimensions, too
- To build good visualizations, we need to **map data dimensions to visual dimensions** in a principled way

Data → Visuals

- Remember...

country	year	cases	population
Afghanistan	1999	216745	19987071
Afghanistan	2000	21666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127215272
China	2000	216766	128023583

variables

country	year	cases	population
Afghanistan	1999	216745	19987071
Afghanistan	2000	21666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127215272
China	2000	216766	128023583

observations

country	year	cases	population
Afghanistan	1999	216745	19987071
Afghanistan	2000	21666	20095360
Brazil	1999	31737	17206362
Brazil	2000	80488	17404898
China	1999	212258	127215272
China	2000	216766	128023583

values



- Big idea behind visualization

- Data have dimensions
- Visualizations have dimensions, too
- To build good visualizations, we need to

map data dimensions to visual dimensions in a principled way

Data → Visuals

Data

country	year	cases	population
Afghanistan	1999	1866	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

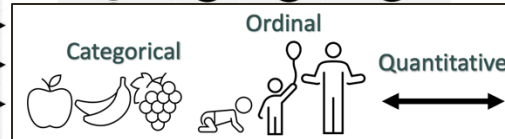
variables

country	year	cases	population
Afghanistan	1999	1866	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

observations

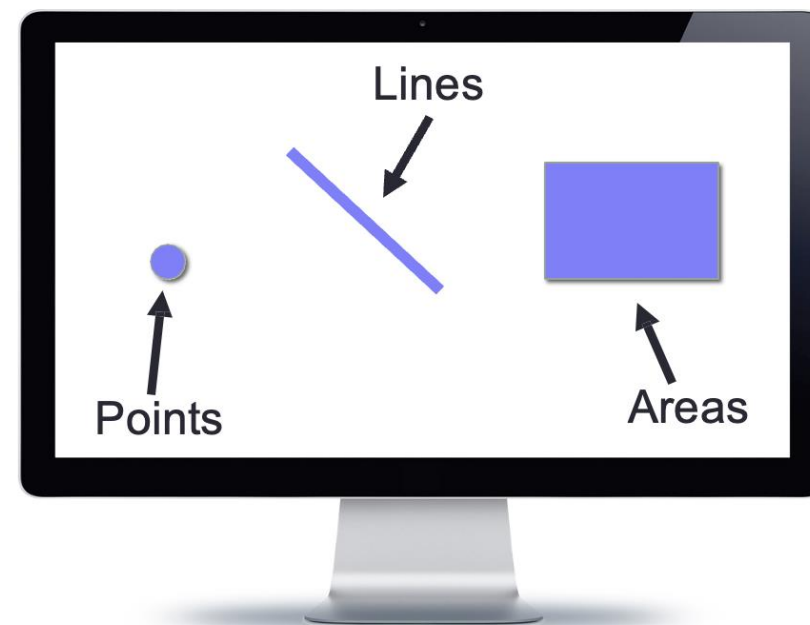
country	year	cases	population
Afghanistan	1999	1866	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

values



Visuals

- **Marks**
- The “ink”



Data → Visuals

Data

country	year	cases	population
Afghanistan	1999	1845	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

variables

country	year	cases	population
Afghanistan	1999	1845	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

observations

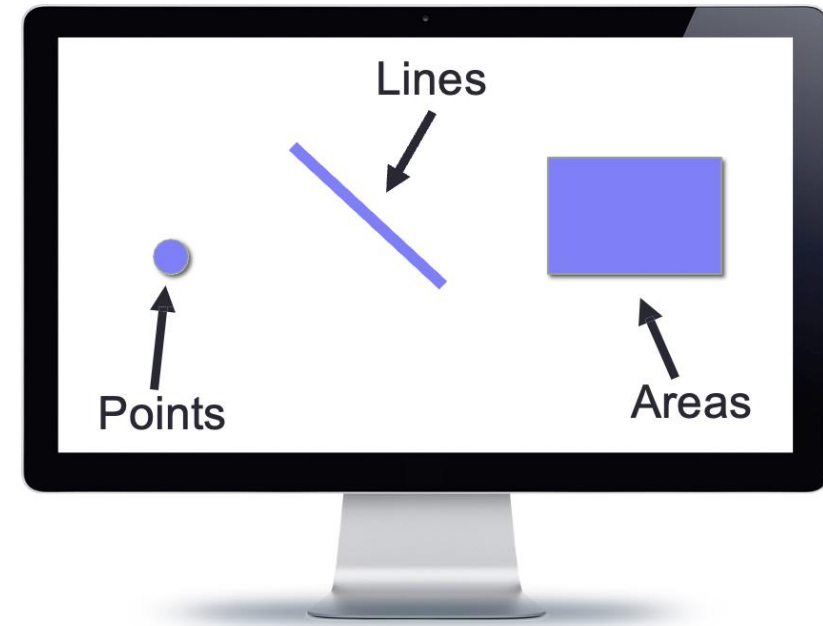
country	year	cases	population
Afghanistan	1999	1845	19987071
Afghanistan	2000	1866	20095360
Brazil	1999	30737	17206362
Brazil	2000	80488	17404898
China	1999	210258	1272015272
China	2000	210766	128602583

values



Visuals

- **Marks**
 - The “ink”

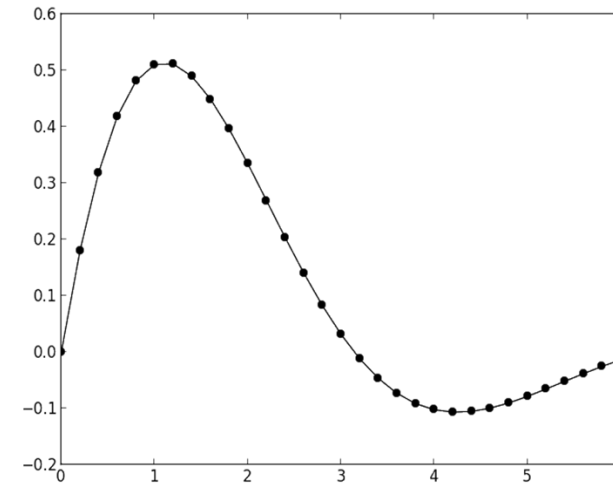
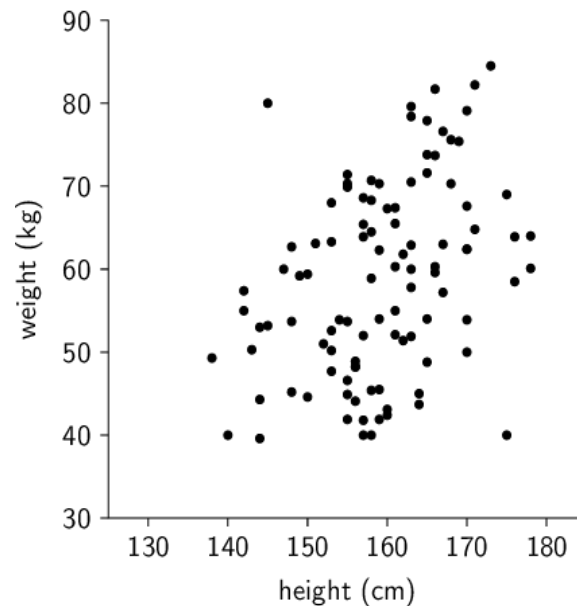


- **Channels or dimensions**
 - *How the marks show up on the page*

Visual Channels / Dimensions

Position

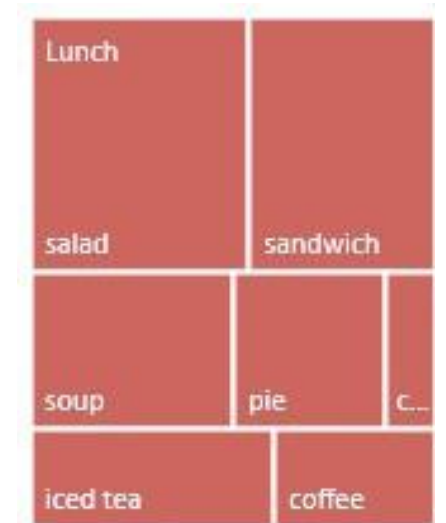
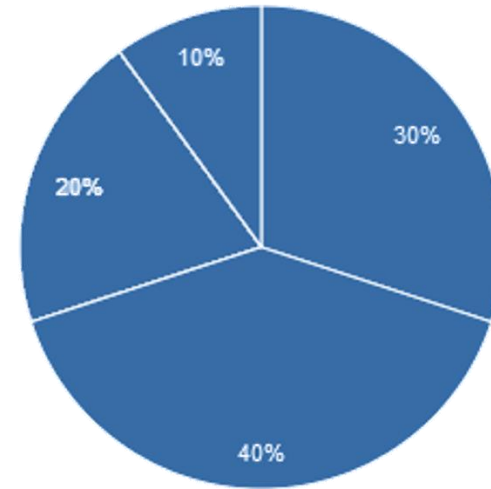
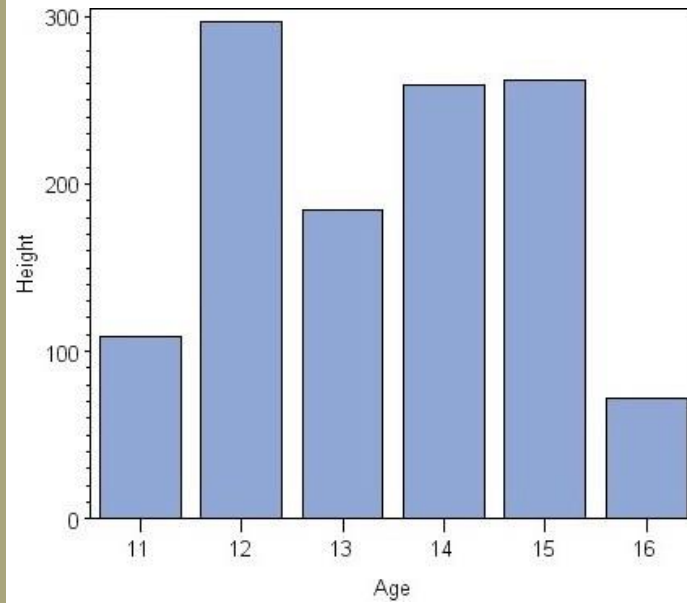
- Encode information using *where* mark is drawn
- Ex.



Visual Channels / Dimensions

Size

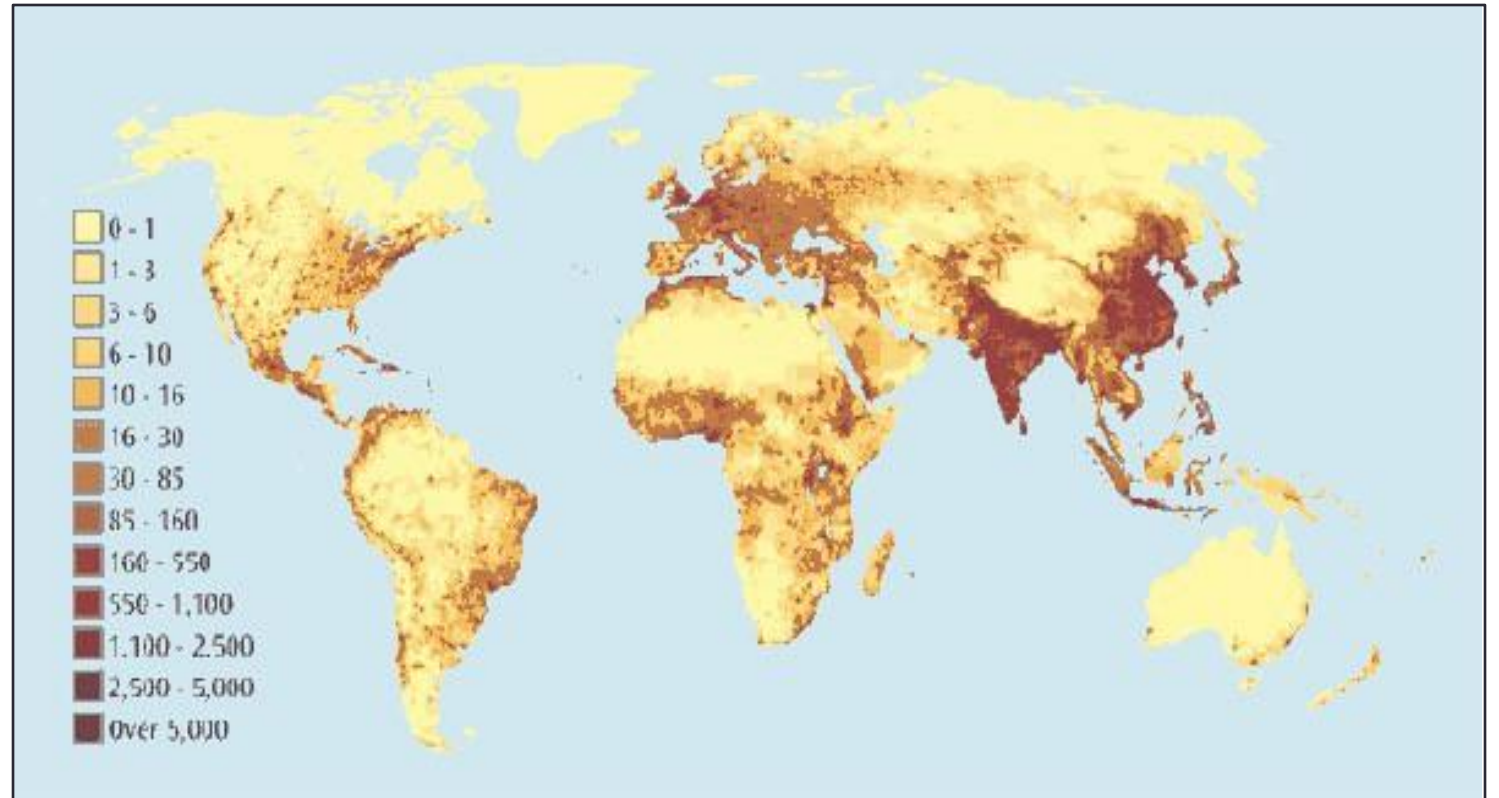
- Encode information using *how big* mark is drawn
- Ex.



Visual Channels / Dimensions

Value

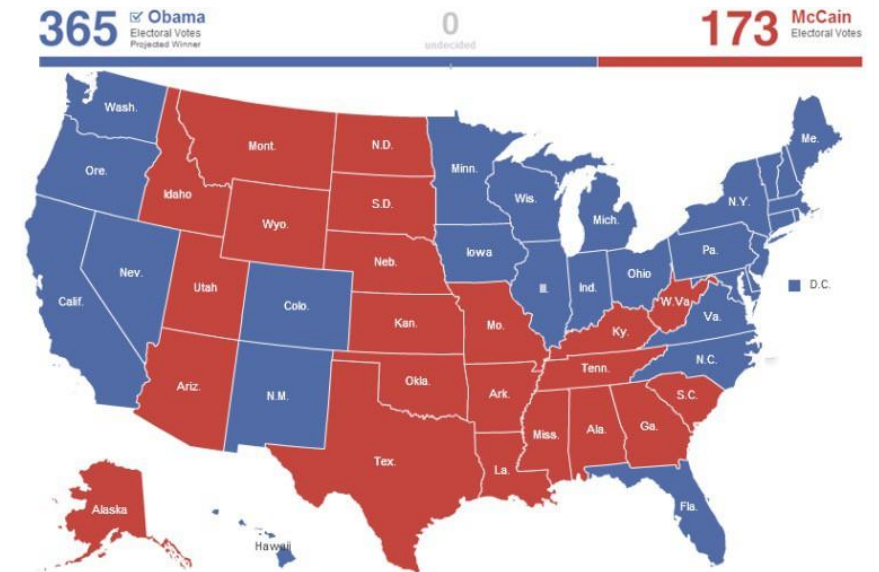
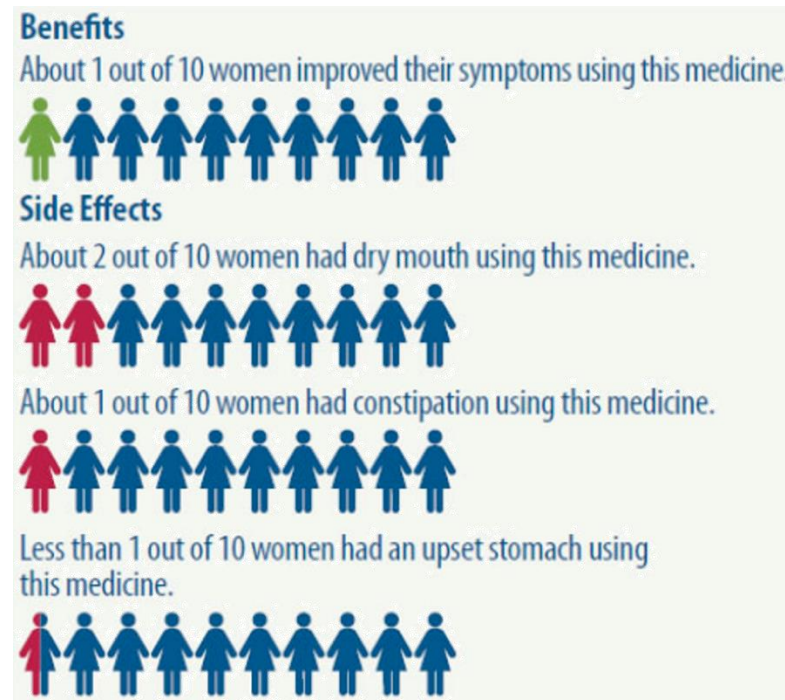
- Encode information using *how dark* mark is drawn
- Ex.



Visual Channels / Dimensions

Color

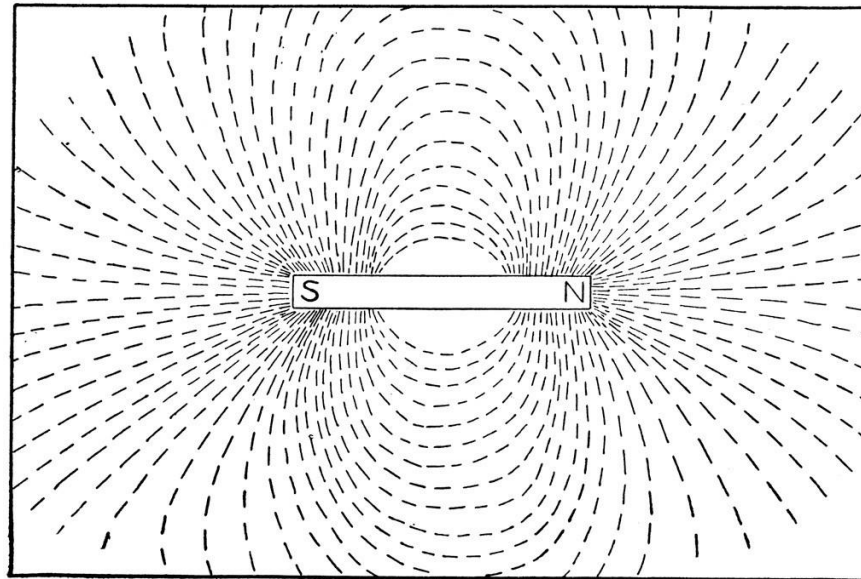
- Encode information using *hue* of mark
- Ex.



Visual Channels / Dimensions

Orientation

- Encode information using how mark is *rotated*
- Ex.



Visual Channels / Dimensions

Shape

- Encode information using how mark is *shaped*
- Ex.

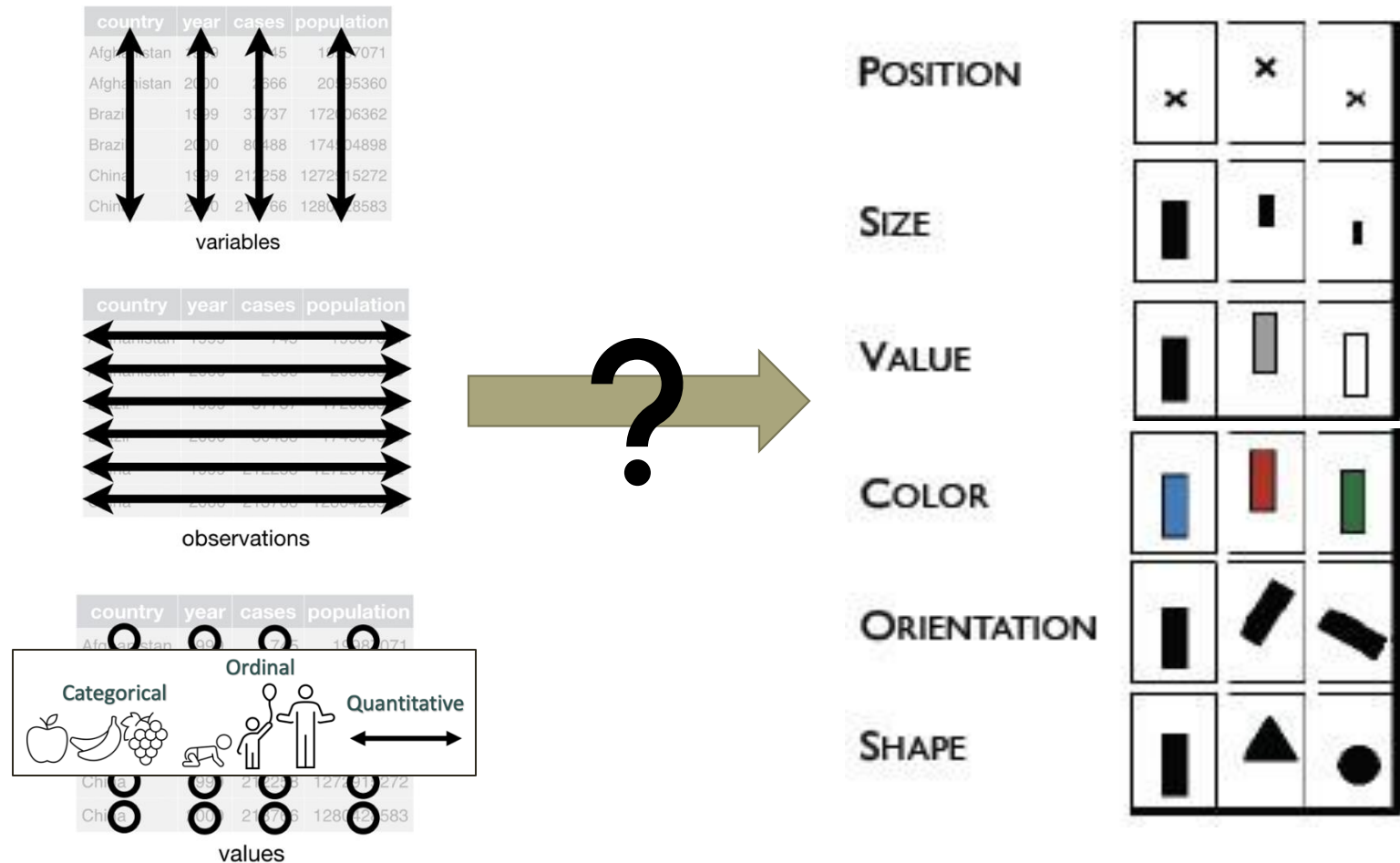


Data → Visuals

- Remember... **Big idea behind visualization**
 - Map data dimensions to visual dimensions in a principled way

Data → Visuals

- Remember... **Big idea behind visualization**
 - Map data dimensions to visual dimensions in a principled way



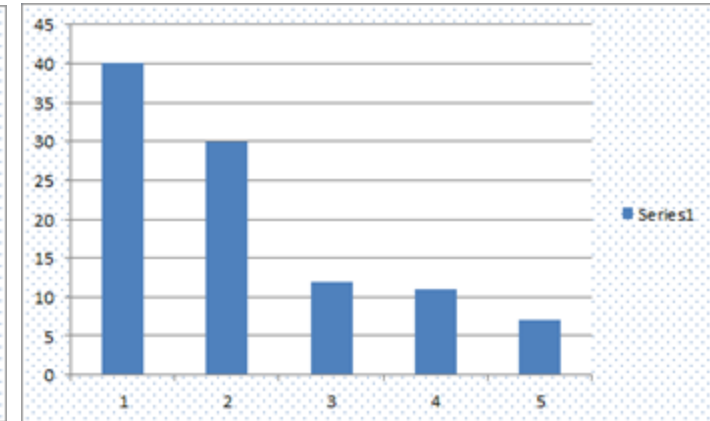
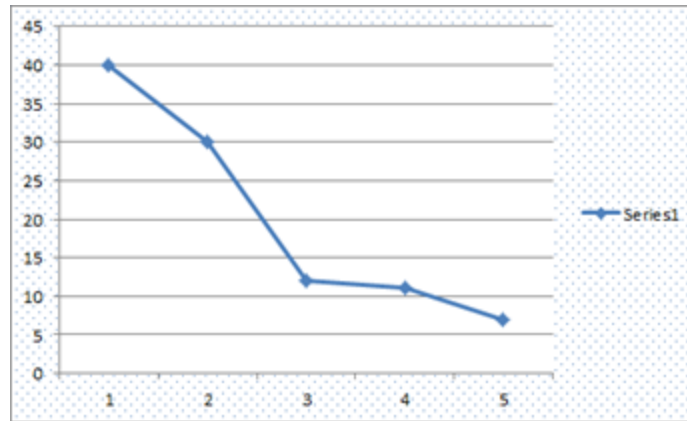
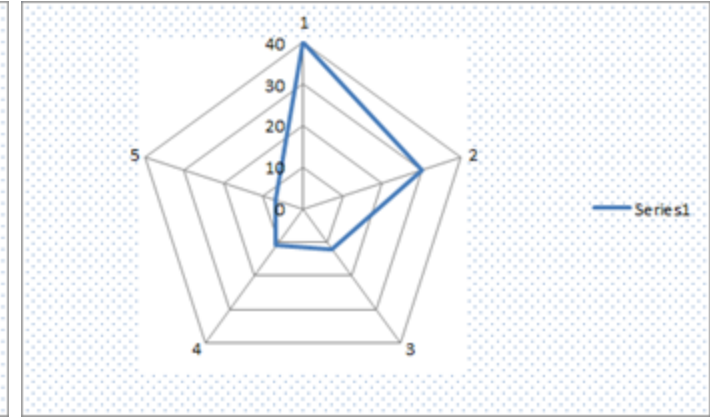
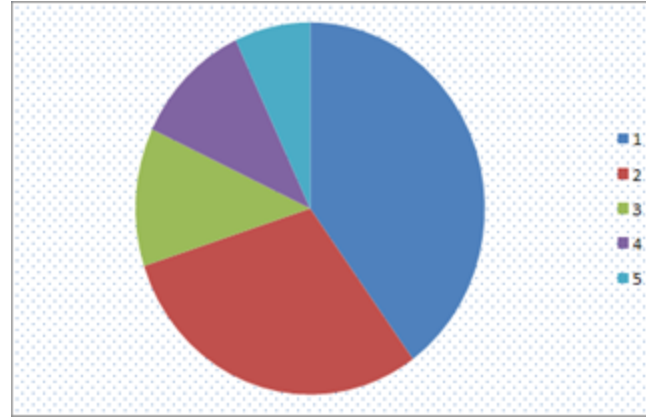
Data → Visuals

- Remember... **Big idea behind visualization**
 - Map data dimensions to visual dimensions in a principled way
 - Not all visual dimensions can represent all data types

	<div>Categorical<div>Apple, Banana, Grapes</div></div> <div>Ordinal<div>Child, Adult, Elderly</div></div> <div>Quantitative<div>Double-headed arrow</div></div>		
POSITION			
SIZE			
VALUE			
COLOR			
ORIENTATION			
SHAPE			

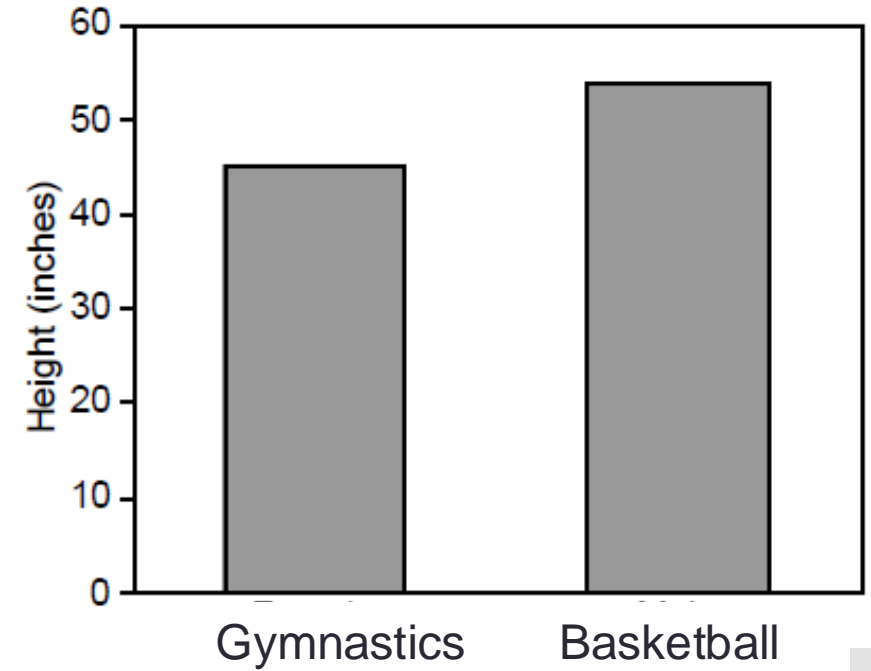
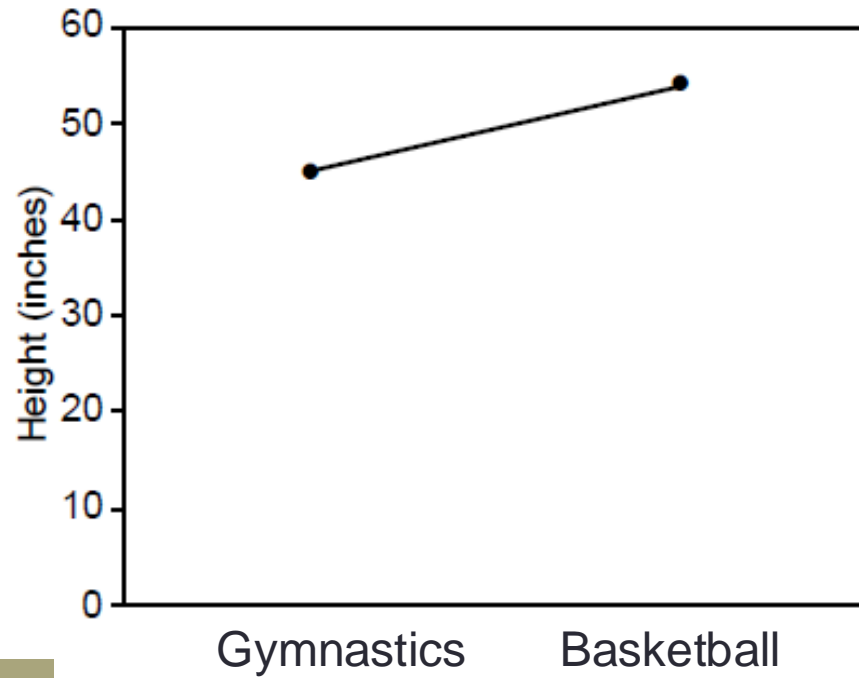
Key question
for this course

Which **data dimension** should be mapped
to which **visual dimension**?



Answer: it depends

Average Height for Youth Sports Participants



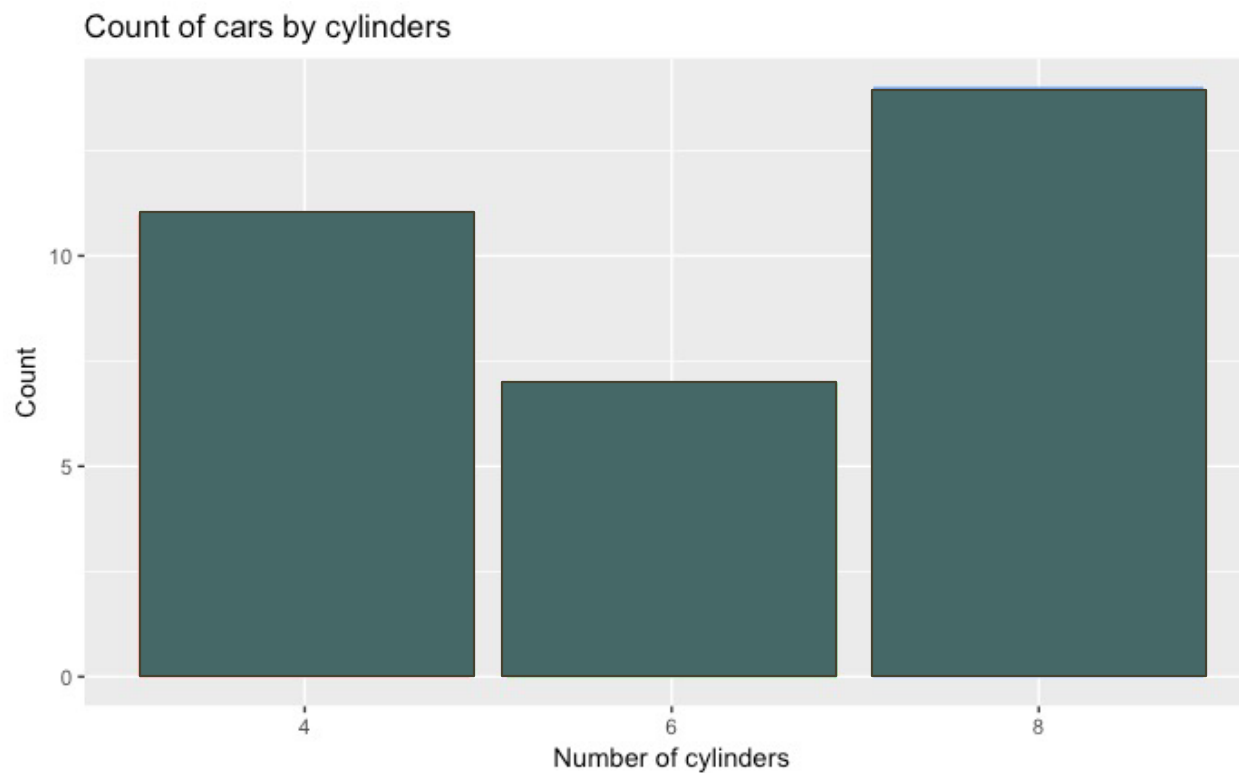
A sampling of visualization techniques



Bar chart

- Used to
 - **Compare**
 - Highlight
 - Grouped groups
- variable
ies
to sub-

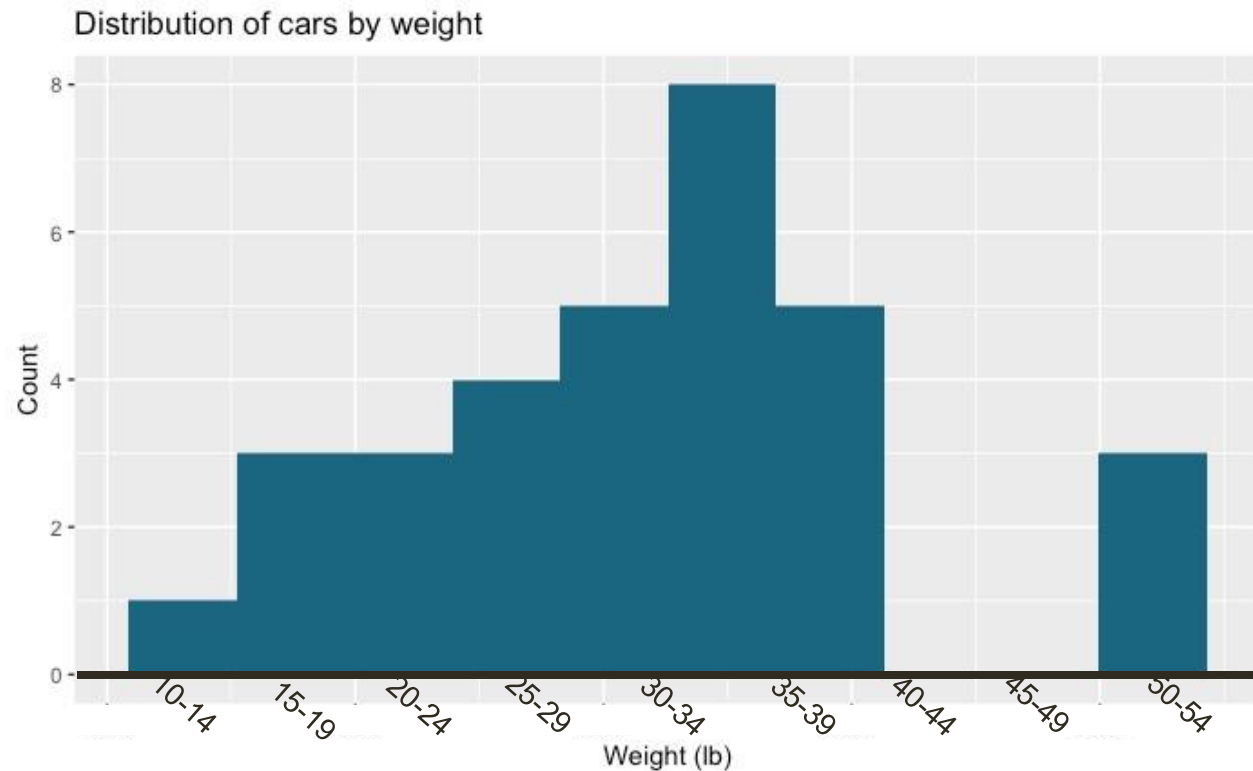
What is the data visual mapping in this chart?



Histogram

- Used for
- Looks like a bar chart, but the bars represent size ranges
- Y-axis represents count or frequency
- Highlights distribution
- Note: bin size makes a big difference!

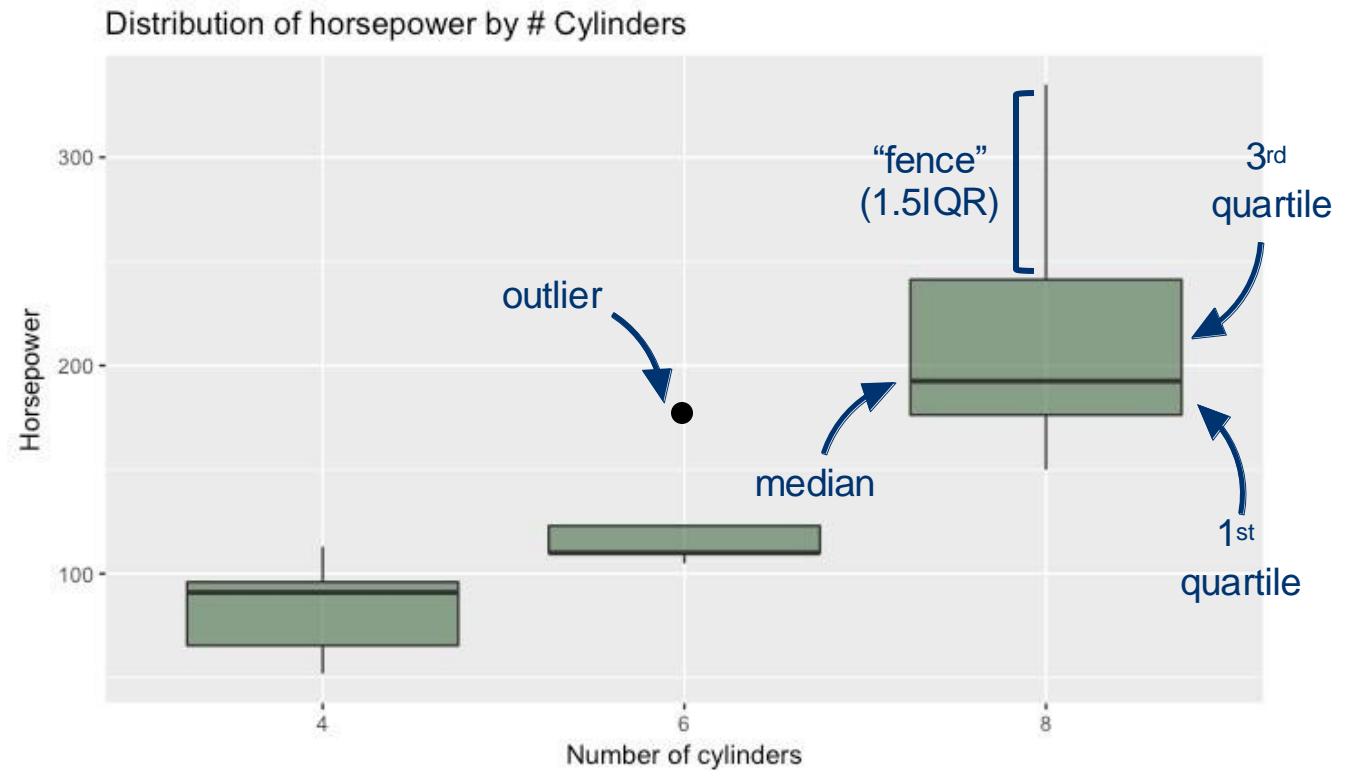
What is the data visual mapping in this chart?



Boxplot

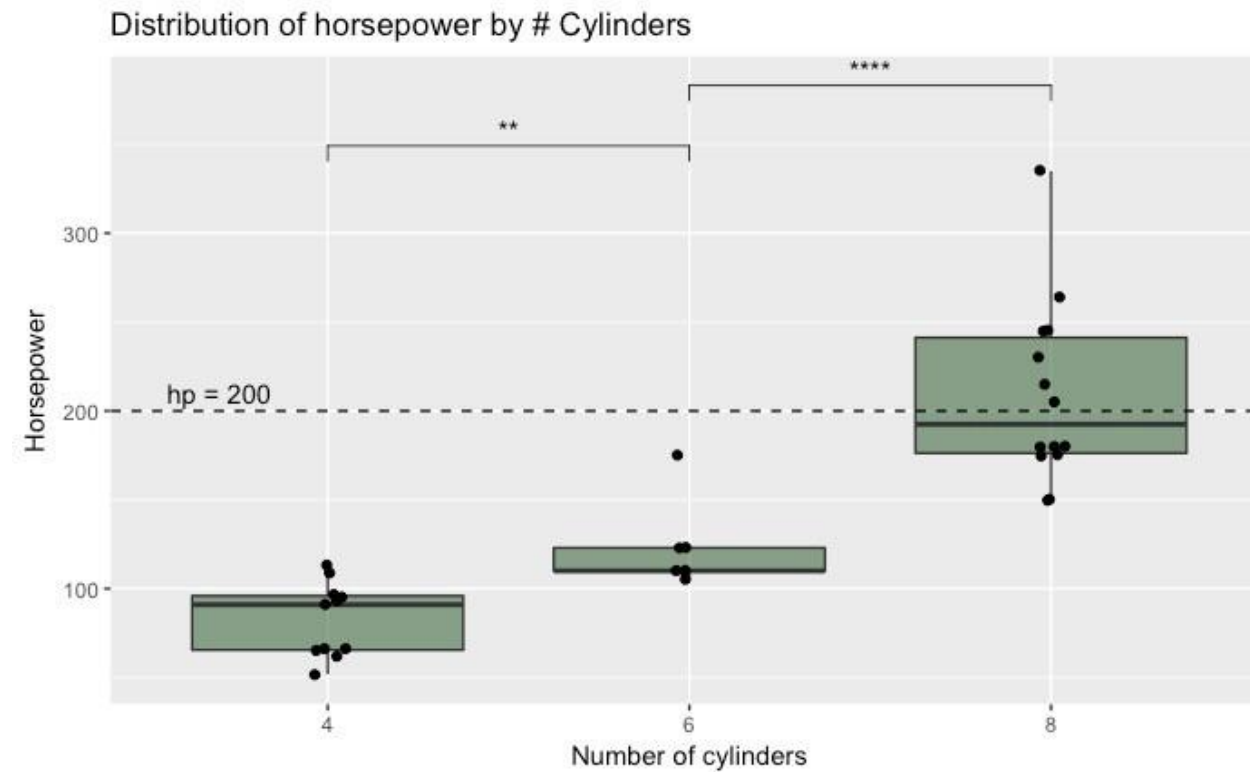
- Used for
- Also used
- Calls out
 - median
 - 1st & 3rd quartiles
 - "fences"
 - outliers

What is the data visual mapping in this chart?



Boxplot

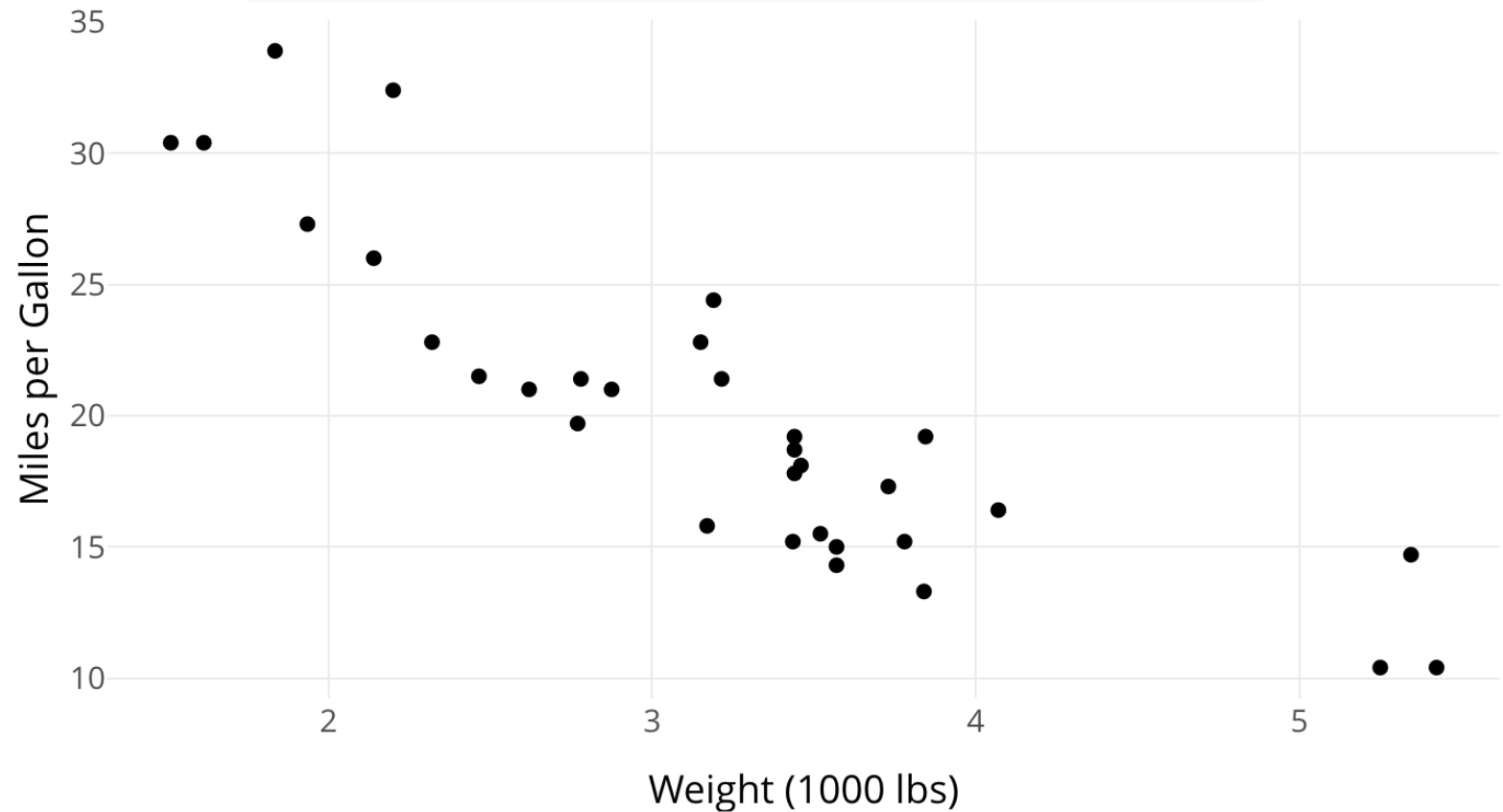
- Use “jitter” to show actual values
- Reference lines can help provide context
- Can use annotations to show statistical significance



Scatterplot

- Used to visualize the relationship between two quantitative variables
- Shows the distribution of data points
- Each point represents an observation

What is the data visual mapping in this chart?

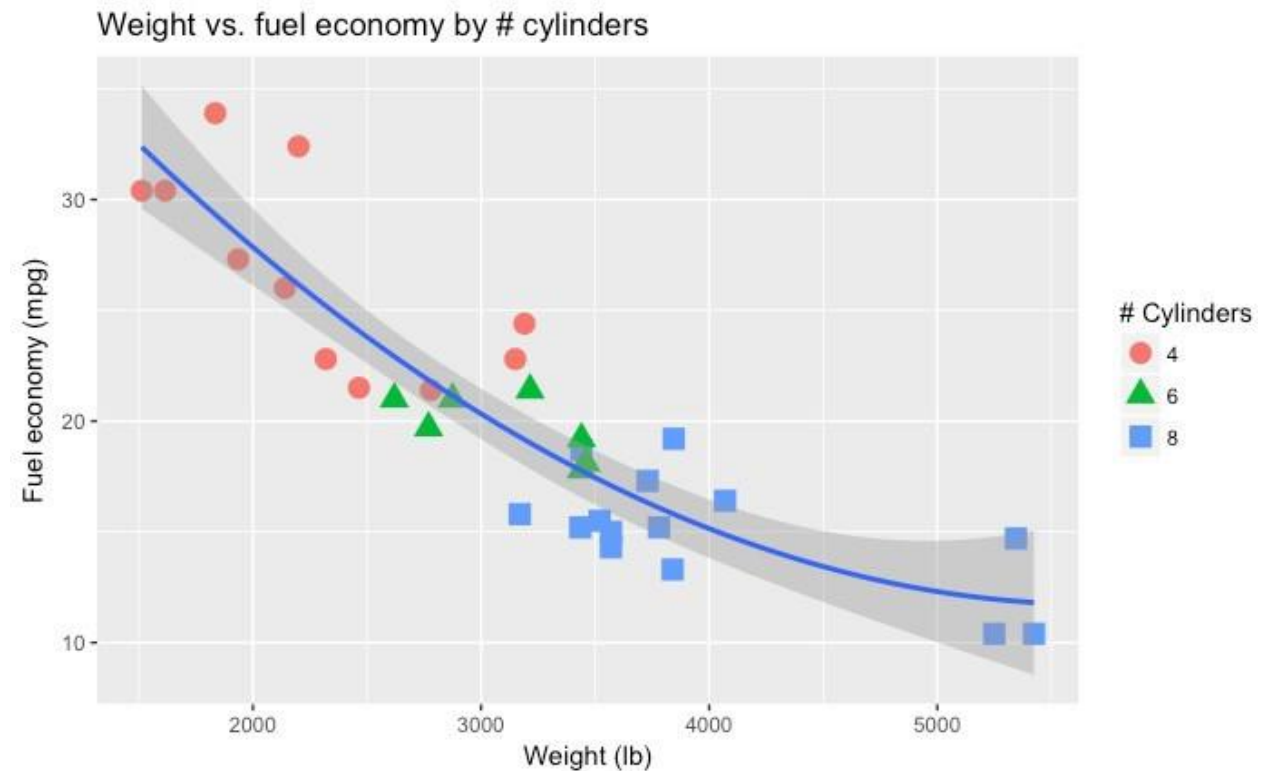


Scatterplot

- Can use
- **quantitative**
- This high
- Sometimes

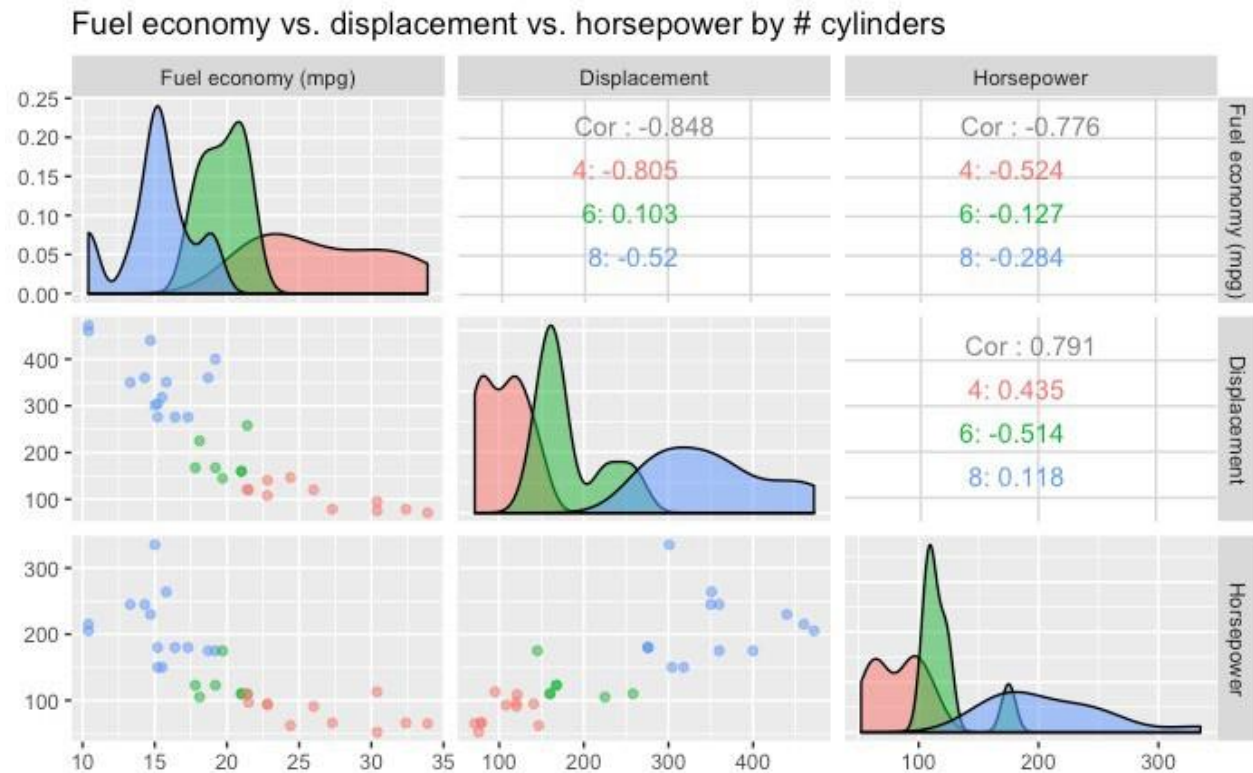
variable X one

What is the data visual mapping in this chart?



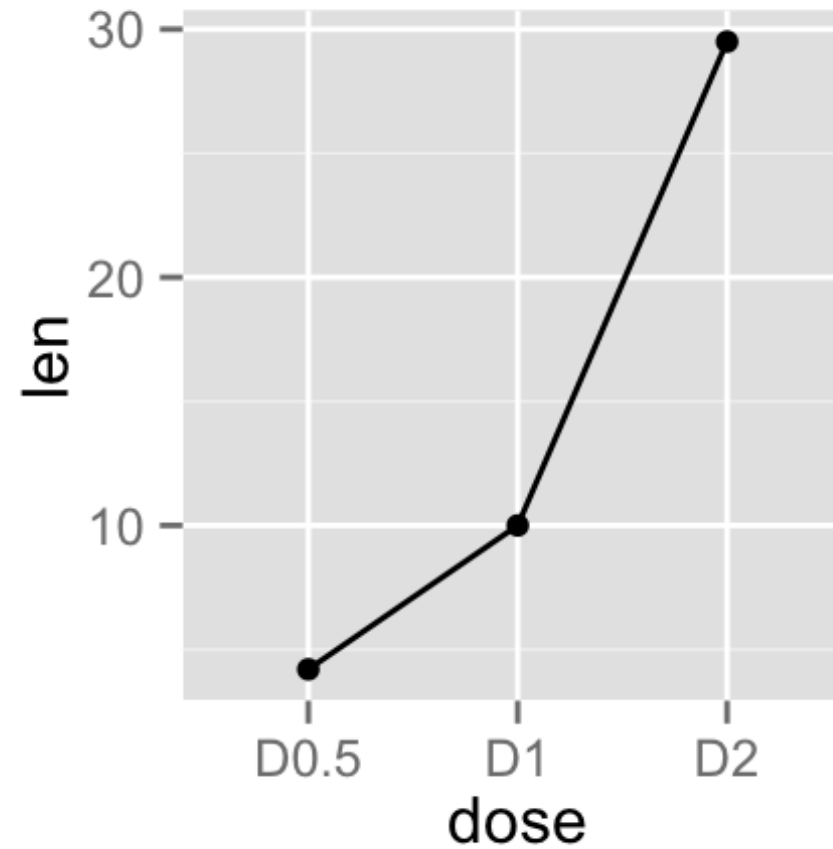
Scatterplot matrix (SPLOM)

- Can use to show **many combinations of one quantitative variable X one quantitative variable**
- Combines multiple scatterplots into a matrix to show **additional relationships**



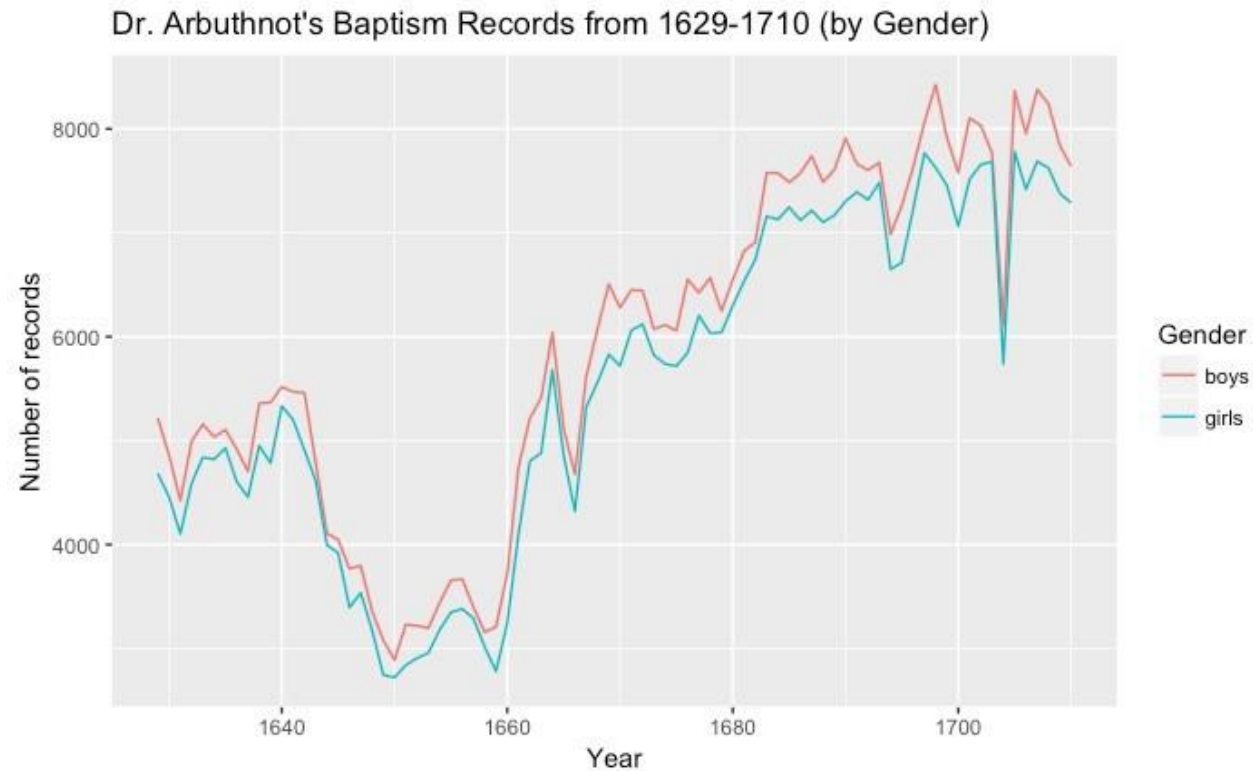
Line chart

- Shows the relationship between variables over time
- What is the data visual mapping in this chart?

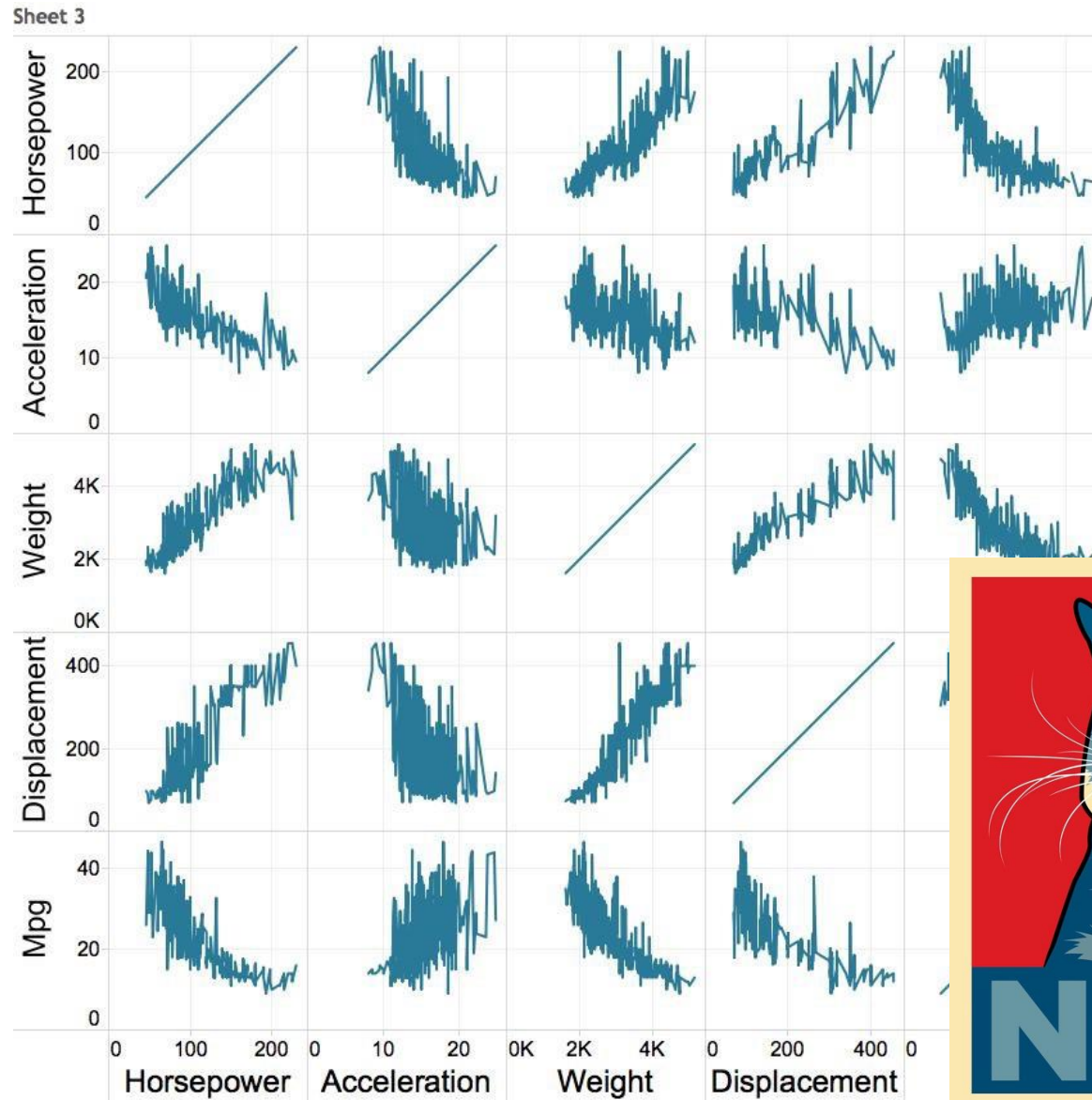


Line chart

- Multiple lines allow **comparison of trends**
- Can show one quantitative variable across groups, or multiple quantitative variables (if they have the same scale)
- Highlights “position switches”

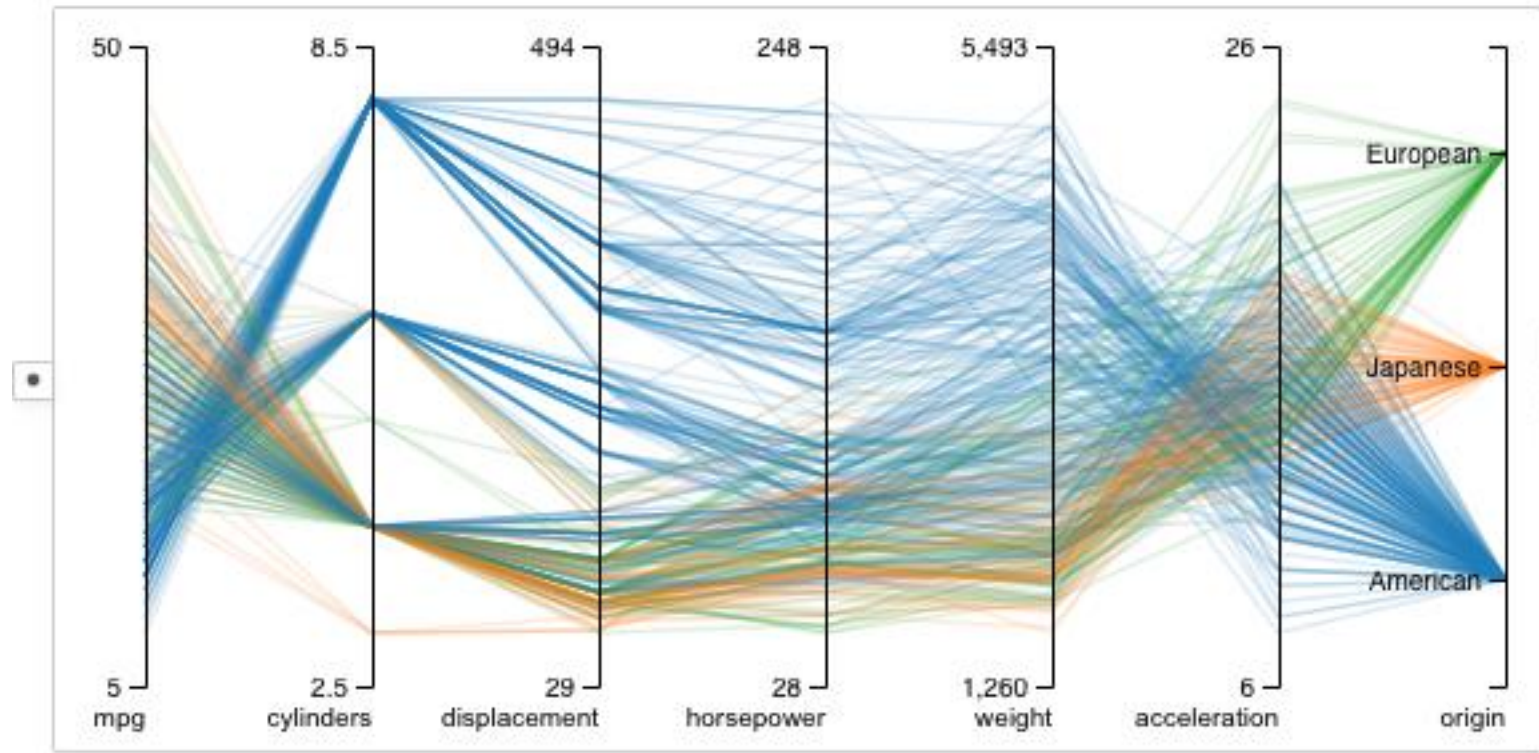


Multiple
variables: line
chart matrix?



Parallel Coordinates Plot

- Supports (pairwise) **comparison of a collection of quantitative variables**
- Each axis represents one variable
 - They may have different scales, typically you normalize them
- Each line represents one observation (connecting the associated values along each axis)
- Axis order matters!



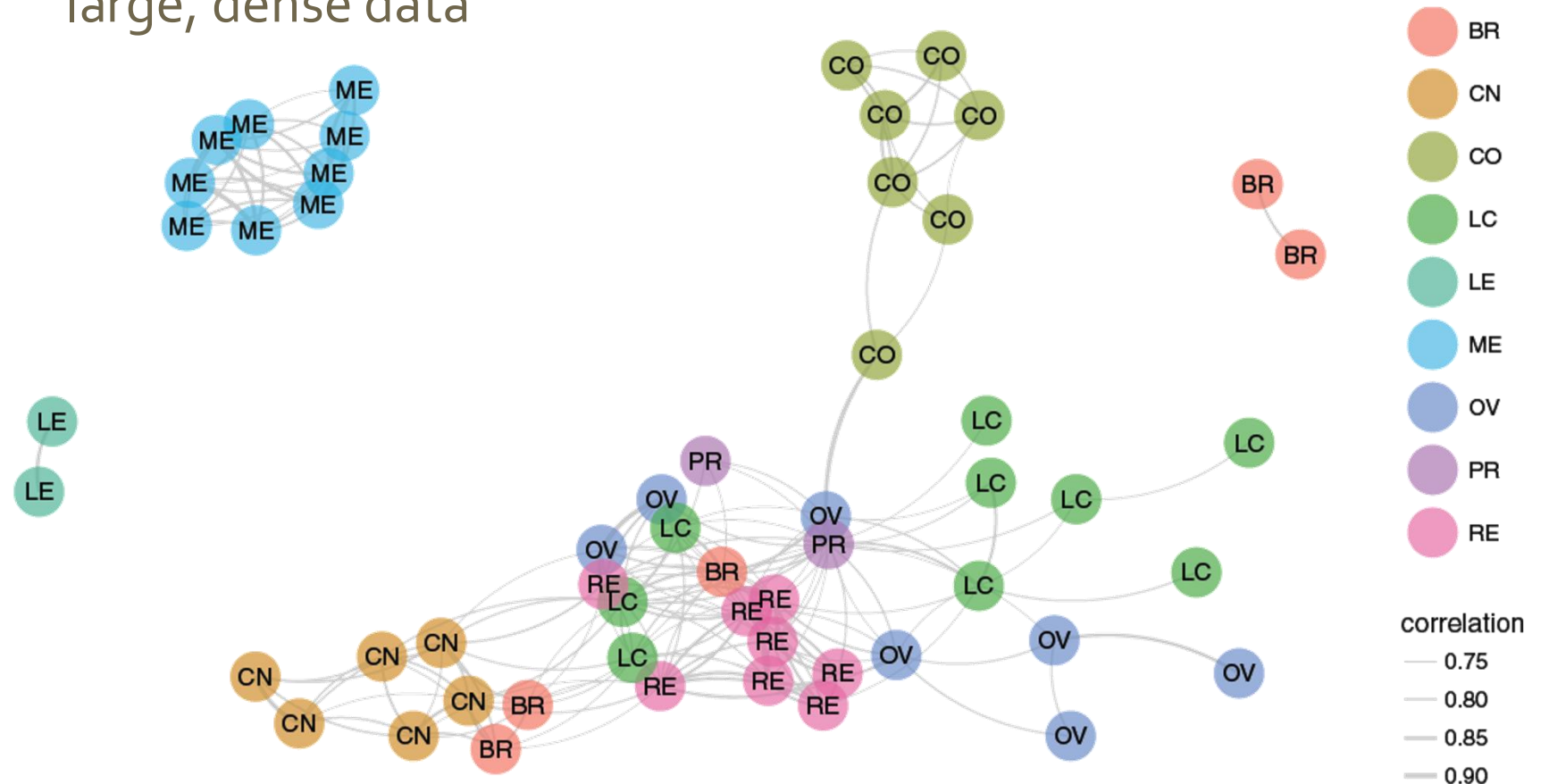
<https://visflow.org/node/visualization/parallel-coordinates.html>

[ta-to-parallel.html](#)

Network

- Shows
 - Useful
 - Can use
 - Caveat
- large, dense data

What is the data visual mapping in this chart?



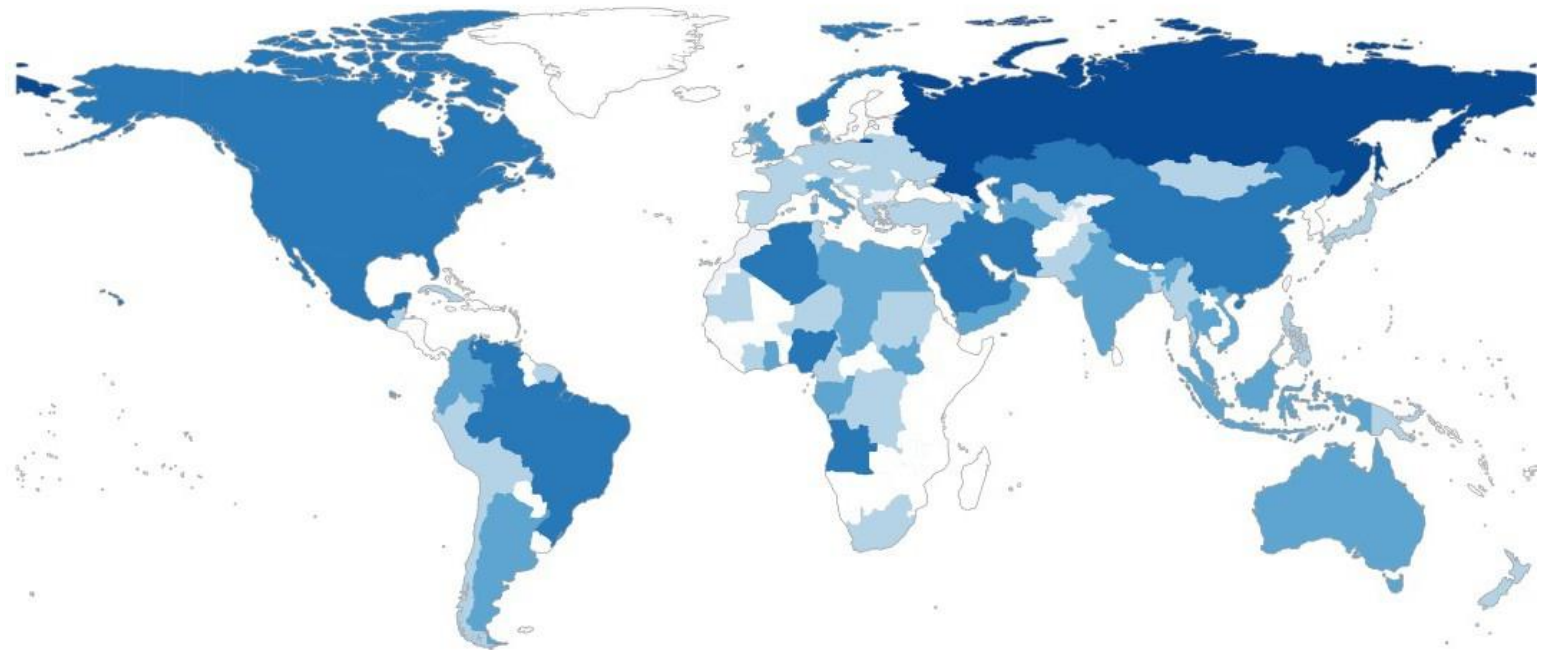
- Shows
 - Useful
 - Filled
 - Remember
- than size comparisons.

What is the data visual mapping in this chart?

ponent
humans

Oil Prod. (bbl/day)

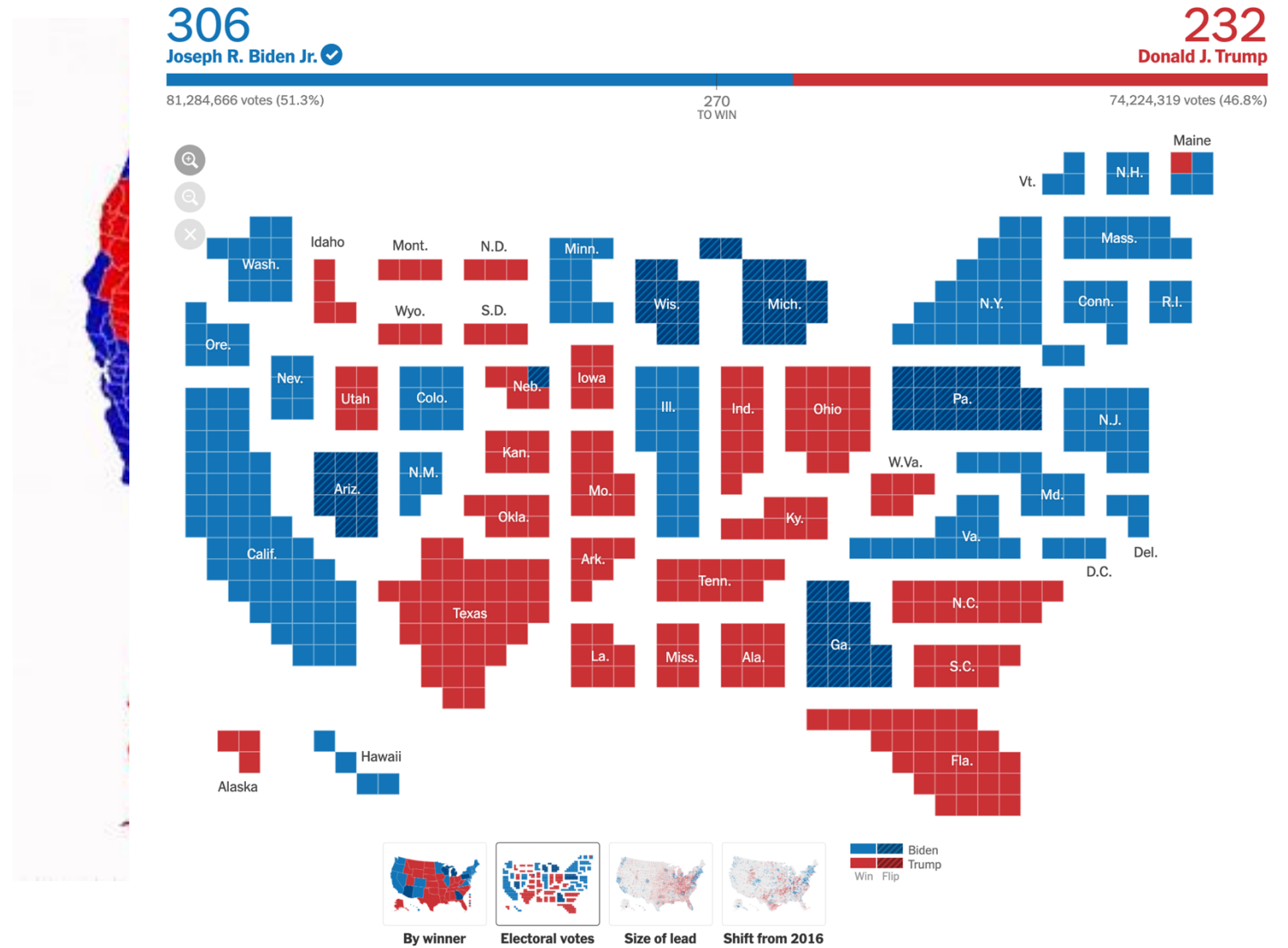
>1000	>100,000	>10 million
>10,000	>1 million	NA



Map

Map

- Remember to map the correct data to your visual channels



Your turn

- Work with 1-2 classmates
- Download `nyc_trees.csv` from the course website.
 - Data is a subset of the dataset available here:
https://data.cityofnewyork.us/Environment/2015-Street-Tree-Census-Tree-Data/uvpi-gqnh/about_data
- Using R or Python (you pick!) and `nyc_trees.csv` generate a:
 - Bar chart
 - Histogram
 - Scatterplot
 - If you have time, try a Boxplot and Linechart
- R Plotting Resources:
<https://r4ds.hadley.nz/layers>
<https://r-graphics.org/>
<https://r-graph-gallery.com/ggplot2-package.html>
- Python Plotting Resources:
<https://matplotlib.org/>
<https://plot.ly/python/>
<https://seaborn.pydata.org/tutorial.html>