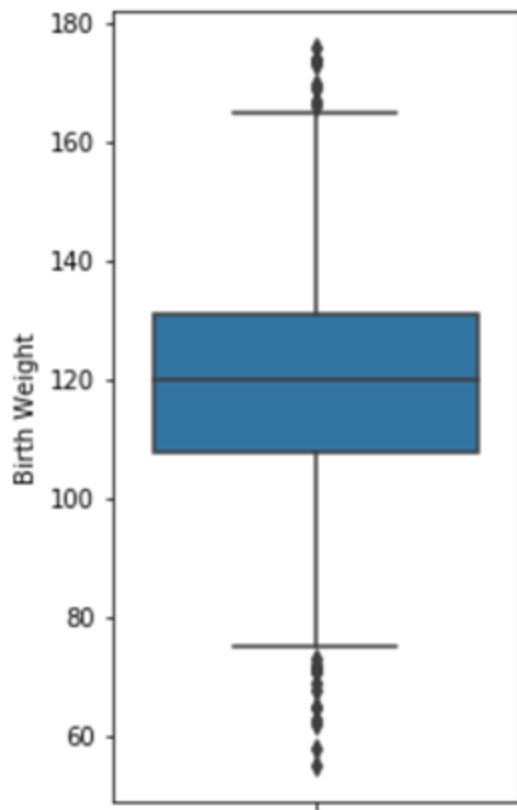


## Review

- Histograms with and without bins
- How to calculate the values versus distribution on histograms
- Density curve – as estimation of histogram
- Box plots – whiskers (IQR is the box length)
- Most of the important data in side the whiskers
- Hex Plots

# Box plots

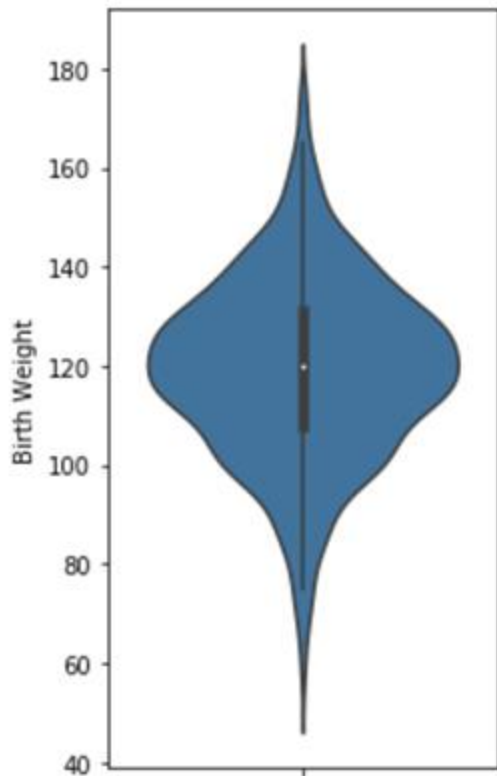


```
1 q1 = np.percentile(bweights, 25)
2 q2 = np.percentile(bweights, 50)
3 q3 = np.percentile(bweights, 75)
4 iqr = q3 - q1
5 whisk1 = q1 - 1.5*iqr
6 whisk2 = q3 + 1.5*iqr
7
8 whisk1, q1, q2, q3, whisk2
```

(73.5, 108.0, 120.0, 131.0, 165.5)

The five numbers above match what we see on the left.

# Violin plots



Violin plots are similar to box plots, but also show smoothed density curves.

- The “width” of our “box” now has meaning!
- The three quartiles and “whiskers” are still present – look closely.
- Both box plots and violin plots are useful for comparing multiple distributions, which we are about to do.
- The curve is symmetrical along the vertical axis

LECTURE 10

# Visualization, Part 2

Principles of sound visualizations; smoothing and transformations.

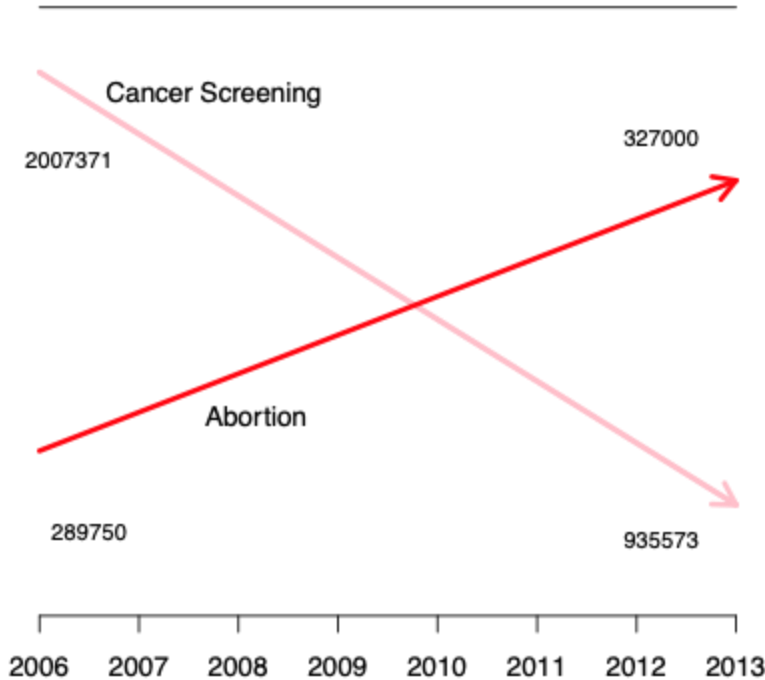
**Sean Kang**

# Overview

- In the first visualization lecture we talked about how to actually make visualizations.
- In this lecture, we will examine visualizations through the following four principles:
  - Scale.
  - Conditioning.
  - Perception.
  - Context.
- We'll also look at how kernel density estimates (KDEs) work, as a method of smoothing histograms.
- We will finish off by looking at transformations as a means to linearize relationships.

Scale

# Case Study: Planned Parenthood Hearing

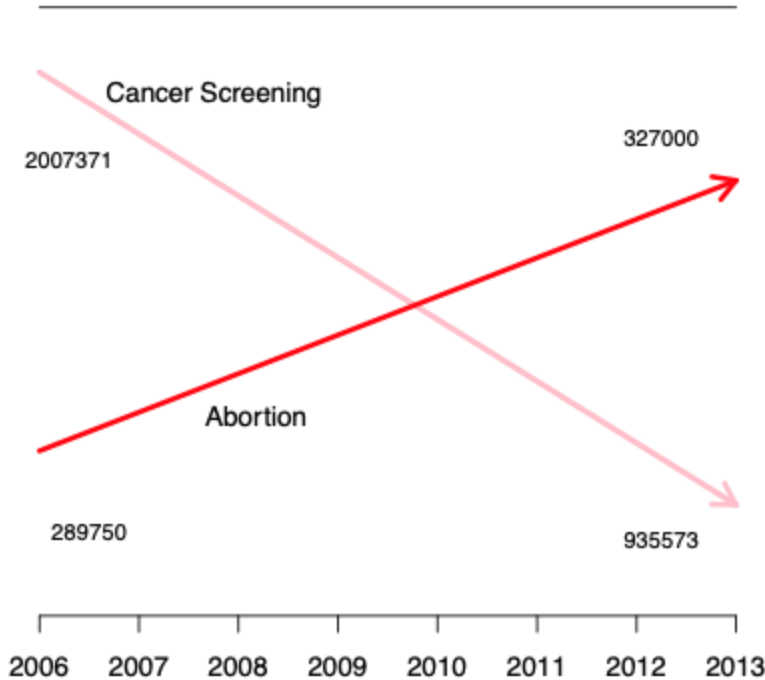


In 2015, Planned Parenthood was accused of selling aborted fetal tissue for profit.

Congressman Chaffetz (R-UT) showed this plot which originally appeared in a report by [Americans United for Life](#).

- What is this graph plotting?
- What message is this plot trying to convey?
- Is anything suspicious?

# Keep axis scales consistent



**The scales for the two lines are completely different!**

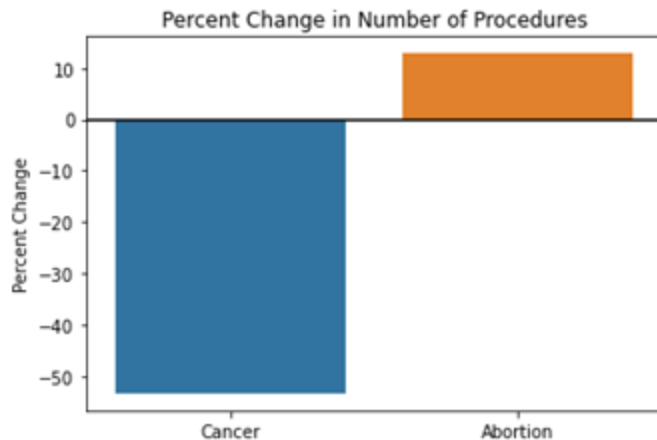
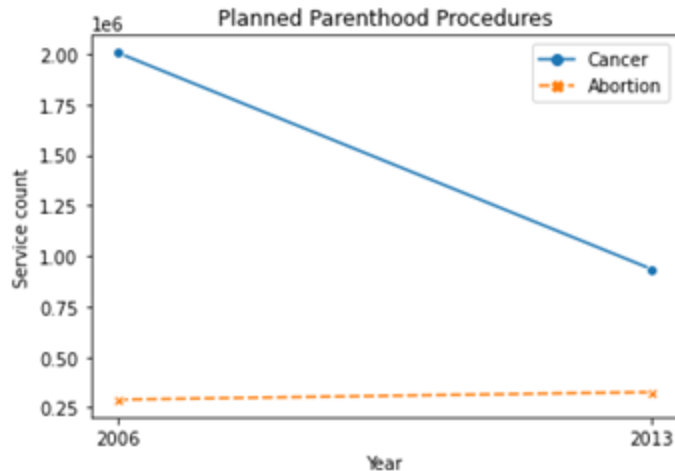
- 327000 is smaller than 935573, but appears to be way bigger.
- **Do not use two different scales for the same axis!**



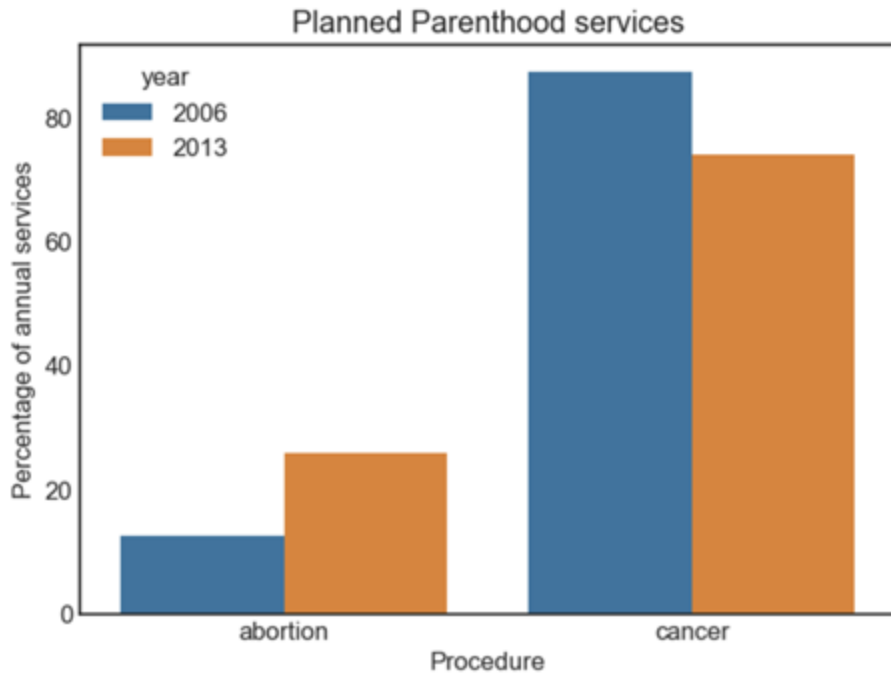
# Consider scale of the data

The top plot draws all of the data on the same scale.

- It clearly shows there was a dramatic drop in cancer screenings by PP.
- But there are still far more cancer screenings than abortions.
- Can plot percentage change instead of raw counts (bottom). This shows that cancer screenings have decreased and abortions have increased, without being misleading.



# Consider scale of the data



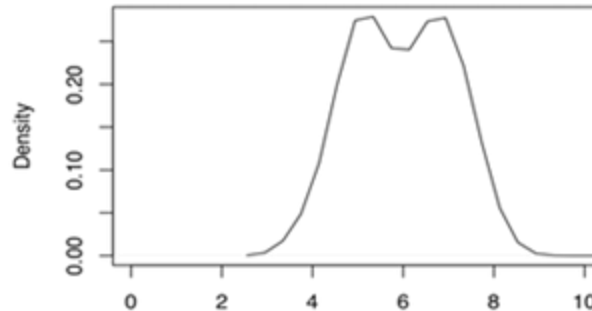
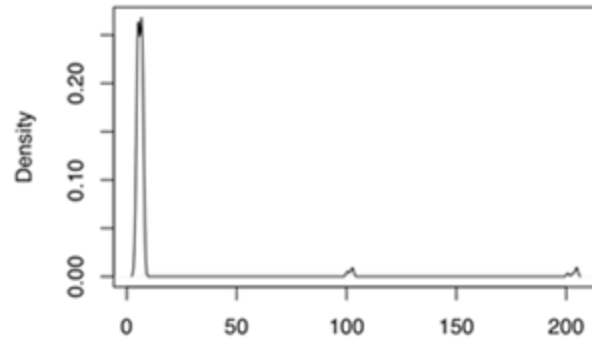
We could also visualize abortions and cancer screenings as a percentage of total procedures.

- Abortions increased from 13% to 26% of total procedures.

# Reveal the data

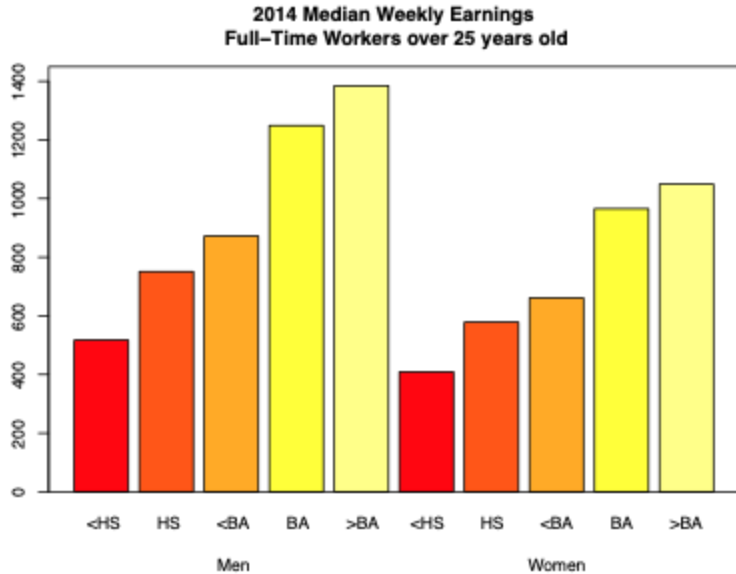
- Choose axis limits to fill the visualization.
- If necessary:
  - Zoom in on the bulk of the data.
  - Create multiple plots to show different regions of interest.

On the left, the bulk of the data is in the  $[0, 10]$  range on the x-axis, whereas the other data 15+ is un-interesting.



# Conditioning

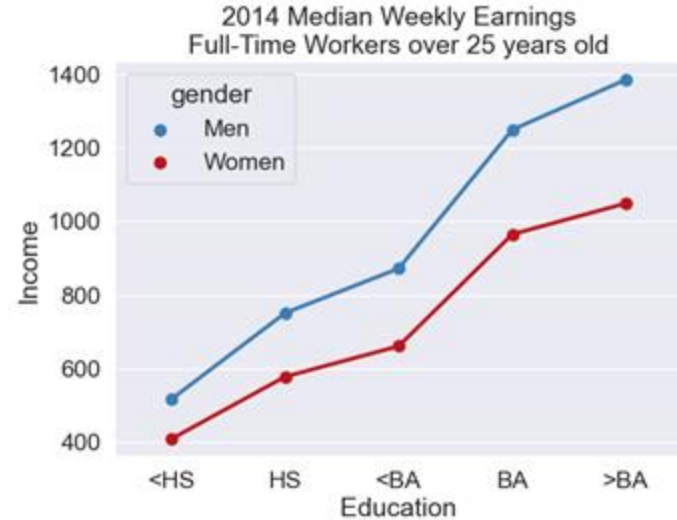
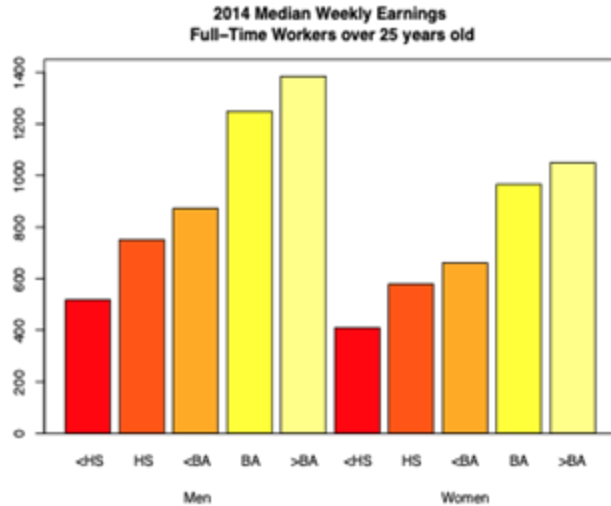
# Case Study: Median Weekly Earnings



This data comes from the [Bureau of Labor Statistics](#), who oversees surveys regarding the economic health of the US. They have plotted median weekly earnings for men and women by education level.

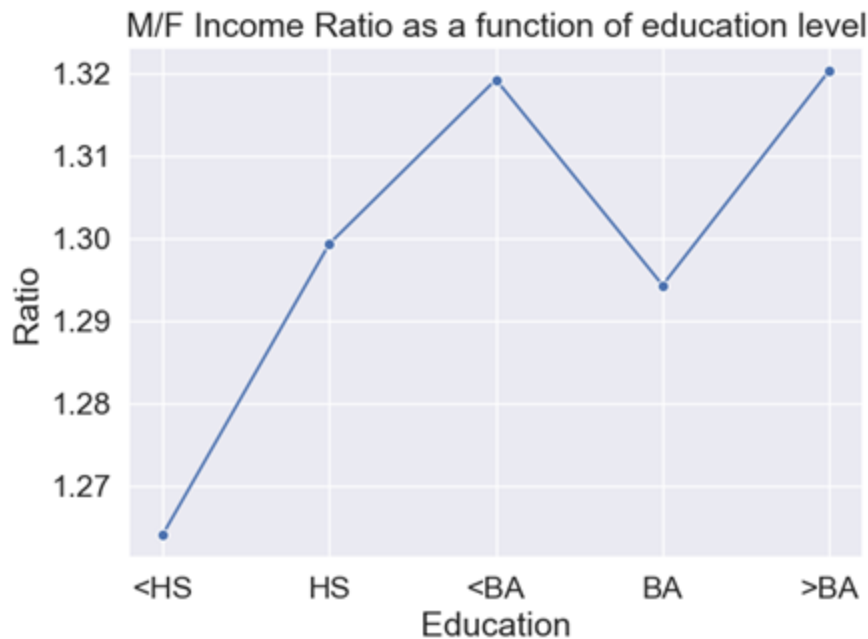
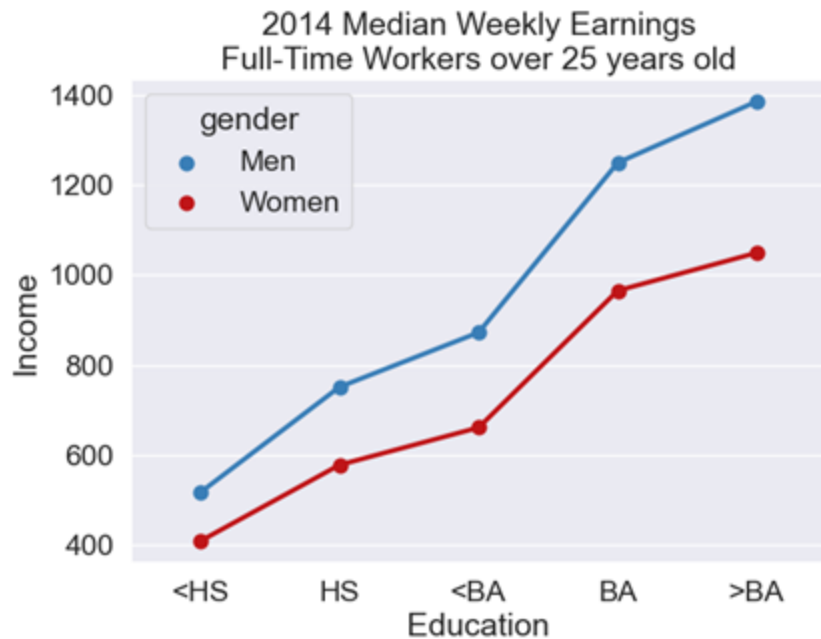
- What comparisons are made easily with this plot?
- What comparisons are most interesting and important?

# Use conditioning to aid comparison



- Lines make it easy to see the large effect of having a BA on weekly earnings.
- Having two separate lines makes clear the wage difference between men and women.
  - It also highlights the fact that the wage difference increases, as education level does.

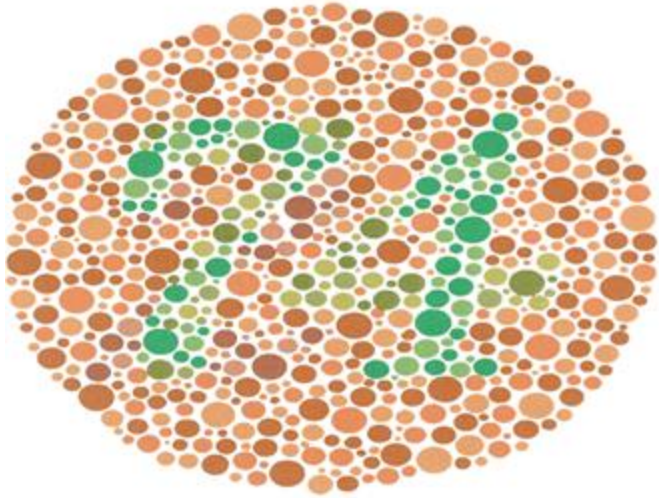
# How does the income gap increase with education?



See notebook for how to get this figure with groupby!

# Perception

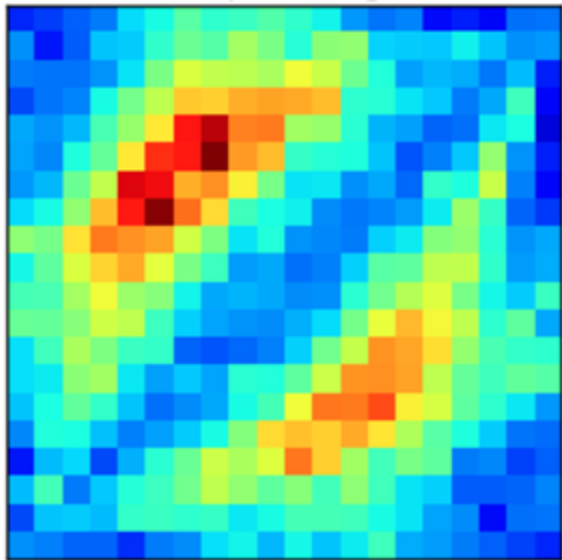




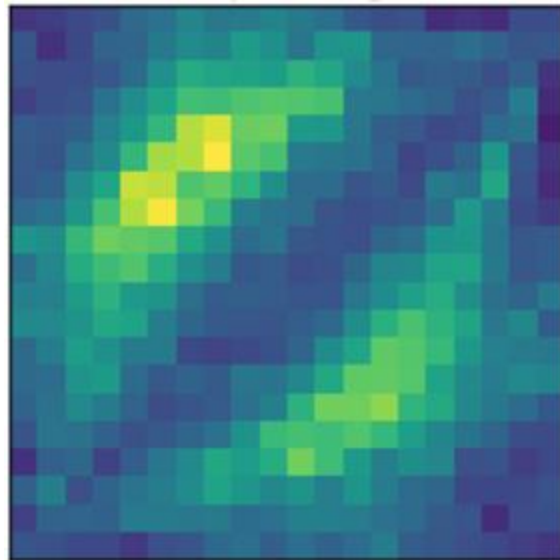
Choosing a set of colors which work together is a challenging task!

## Perception of Color

# Colormaps

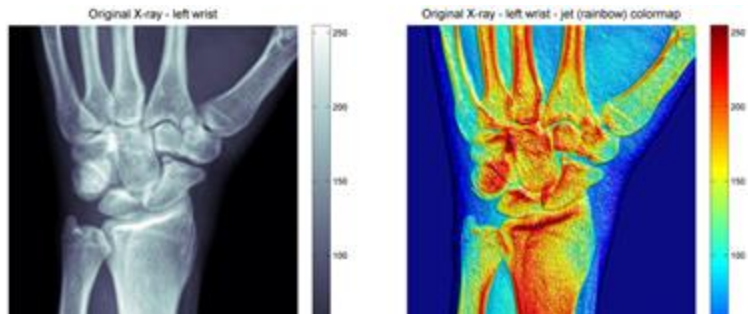


**Jet (older) aka Rainbow**

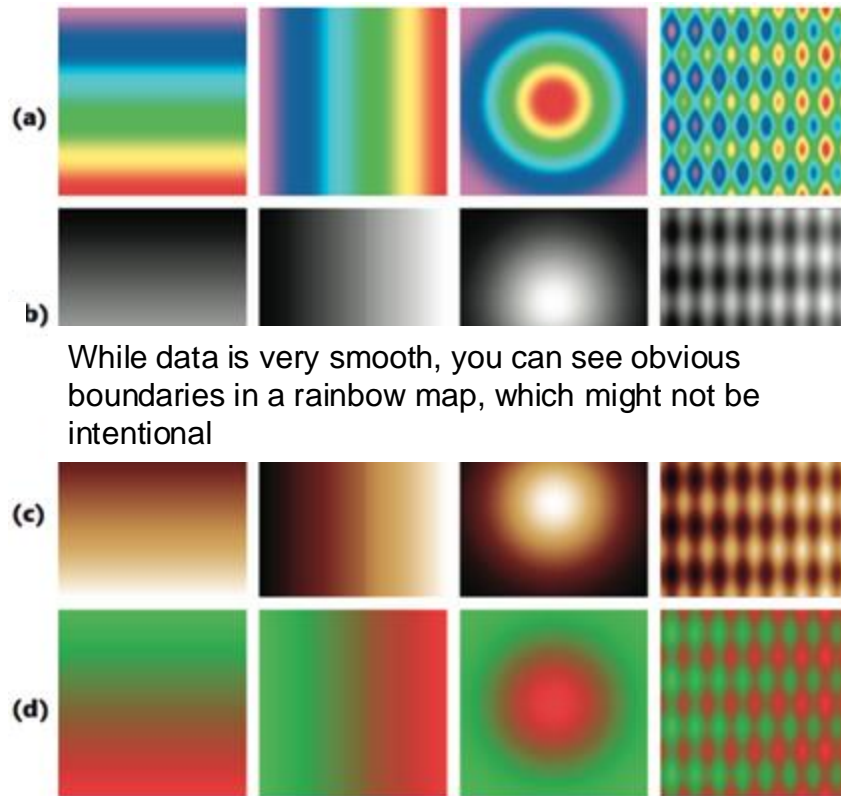


**Viridis (new)**

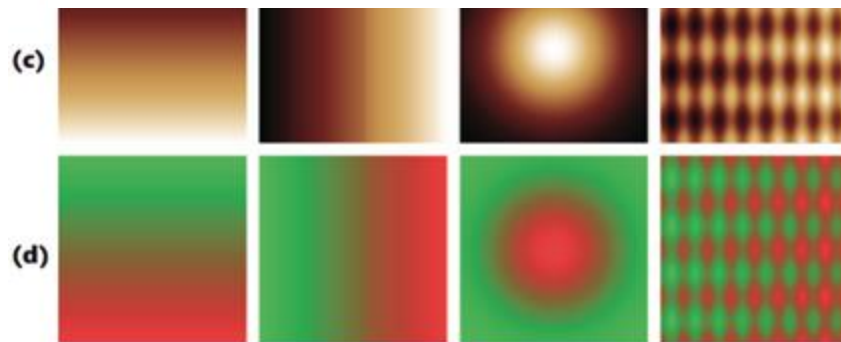
# The jet/rainbow colormap actively misleads



Eyes are drawn very strongly to certain areas of sharp visual contrast which are not necessarily your features



While data is very smooth, you can see obvious boundaries in a rainbow map, which might not be intentional



# Use a perceptually uniform colormap! -> Smoother

- **Perceptually uniform colormaps** have the property that if the data goes from 0.1 to 0.2, the **perceptual change** is the same as when the data goes from 0.8 to 0.9.
- Jet, the old matplotlib default, was far from uniform.
- Viridis, the new default colormap, is.
  - It was created by folks at the Berkeley Institute of Data Science!
  - <https://bids.github.io/colormap/>
- Avoid combinations of red and green, due to red-green color blindness. Viridis is easier on the eye.

Except when not :) The Google Turbo Colormap -> Smoother



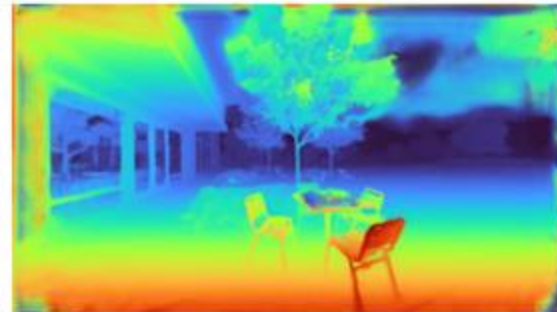
Turbo



Jet



Inferno

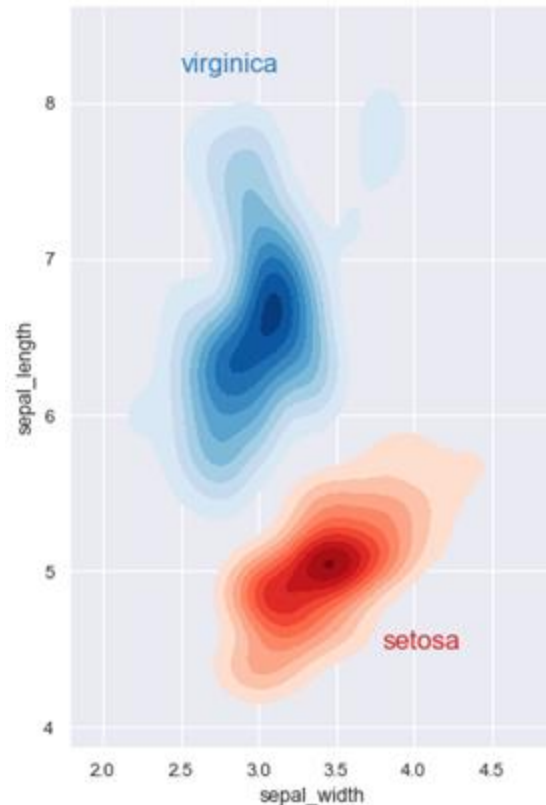


Turbo

# Use color to highlight data type

- **Qualitative:** Choose a qualitative scheme that makes it easy to distinguish between categories.
  - One category isn't "higher" or "lower" than another.
- **Quantitative:** Choose a color scheme that implies magnitude – in this case, the altitude for topographical elevation contour maps
- The plot on the right has both!

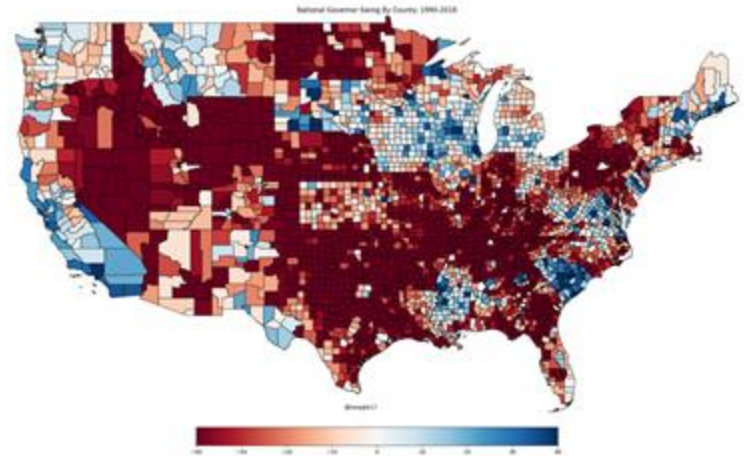
Flower Data



# Sequential vs. diverging colormaps for quantitative data



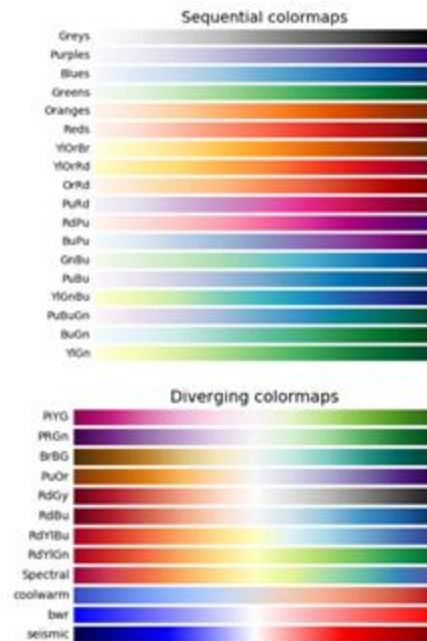
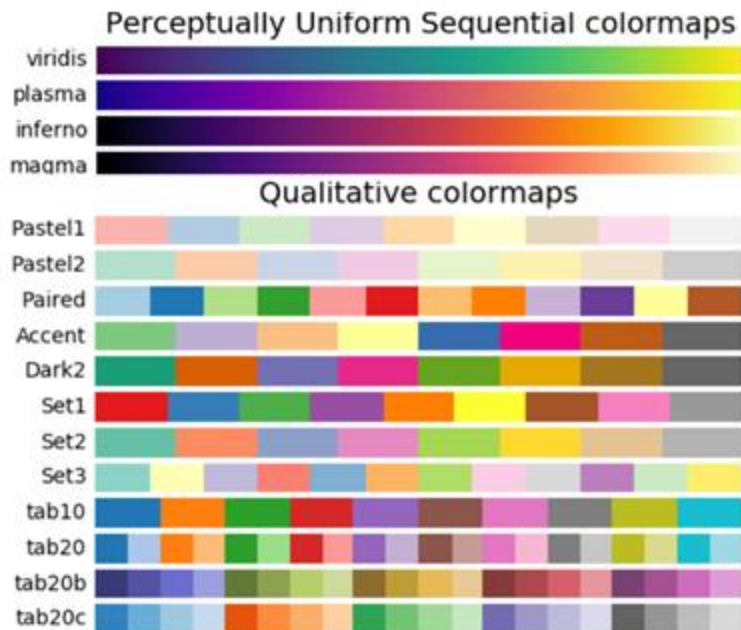
If the data progresses from low to high, use a **sequential** scheme where lighter colors are for more extreme values.



If low and high values deserve equal emphasis, use a **diverging** scheme where lighter colors represent middle values. (ie Political Group)

# Default matplotlib colormaps

Default is Viridis



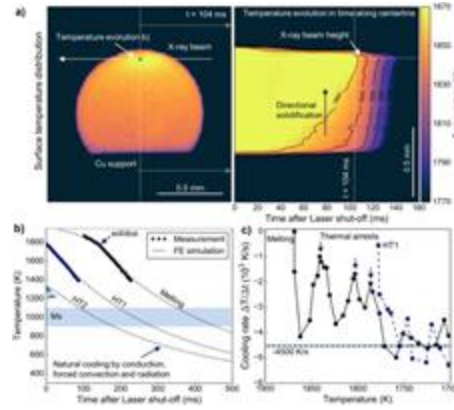
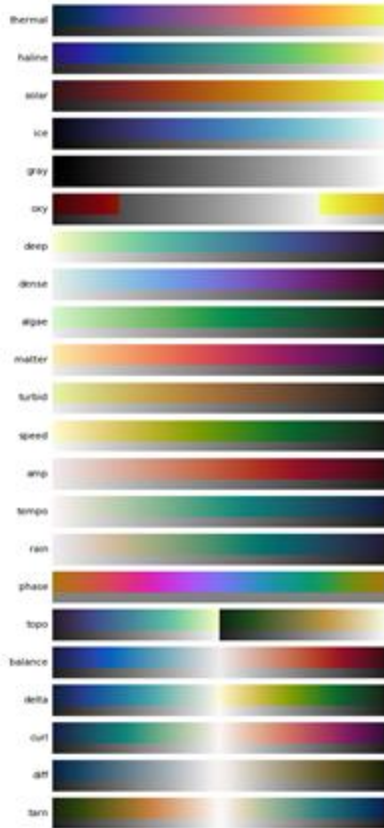
Taken from [matplotlib documentation](https://matplotlib.org/1.3.2/colormaps/colormaps.html).

Many Options Exists

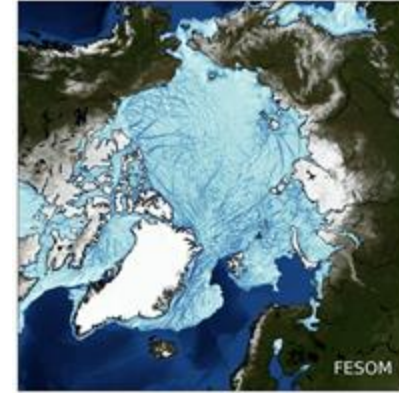


# Domain specific colormaps: [cmocean](#)

(beautiful colormaps for oceanography, by [Kristen Thyng](#))

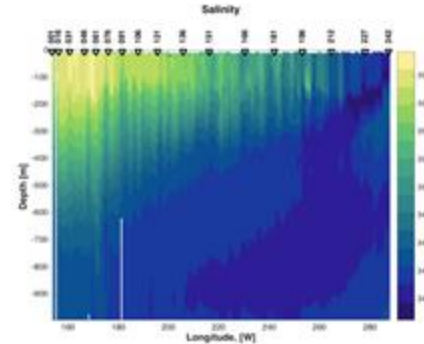


Thermal



Ice

The dark blue can mean the thickness

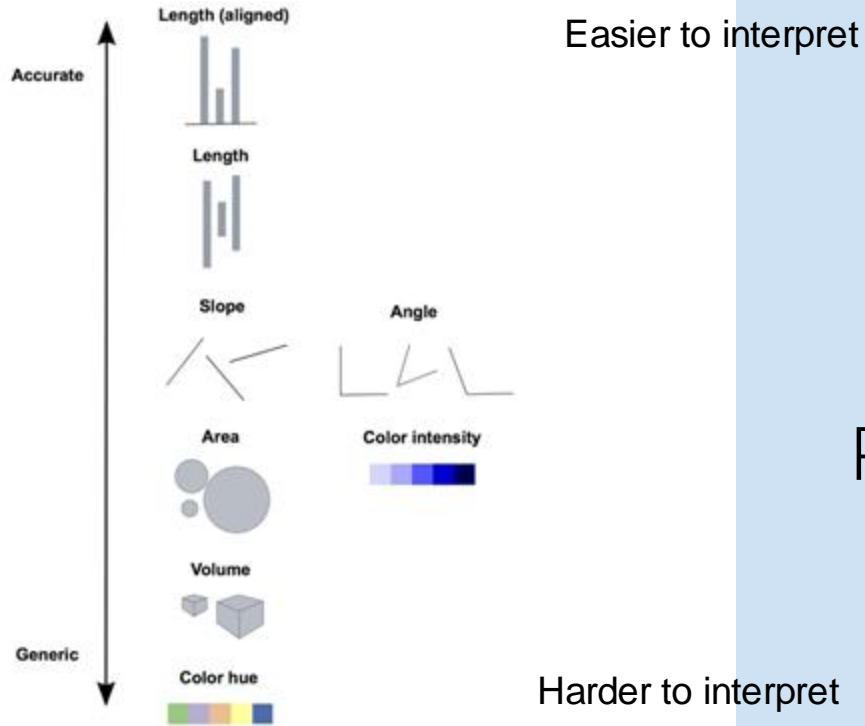


Haline

# Extra reading

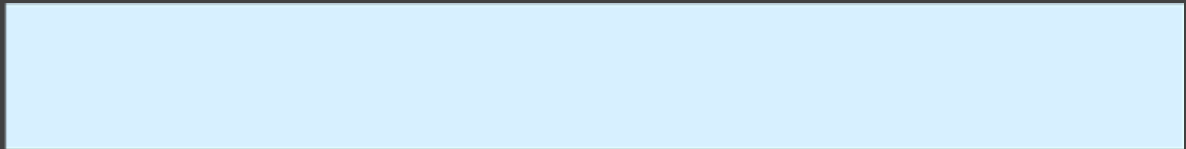
You may want to refer to these articles, which also discuss colormaps.

- Rainbow Colormap (Still) Considered Harmful – look in notes on this slide.
- <https://eagereyes.org/basics/rainbow-color-map - maybe it isn't so bad>
- <https://blog.datawrapper.de/diverging-vs-sequential-color-scales/>
- [https://web.natur.cuni.cz/~langhamr/lectures/vtfg1/mapinfo\\_2/barvy/colors.html](https://web.natur.cuni.cz/~langhamr/lectures/vtfg1/mapinfo_2/barvy/colors.html)

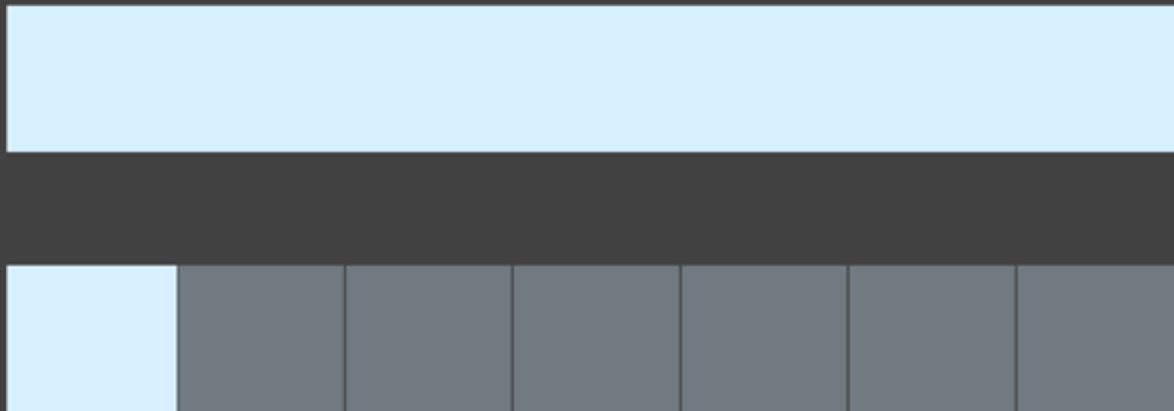


## Perception of Markings

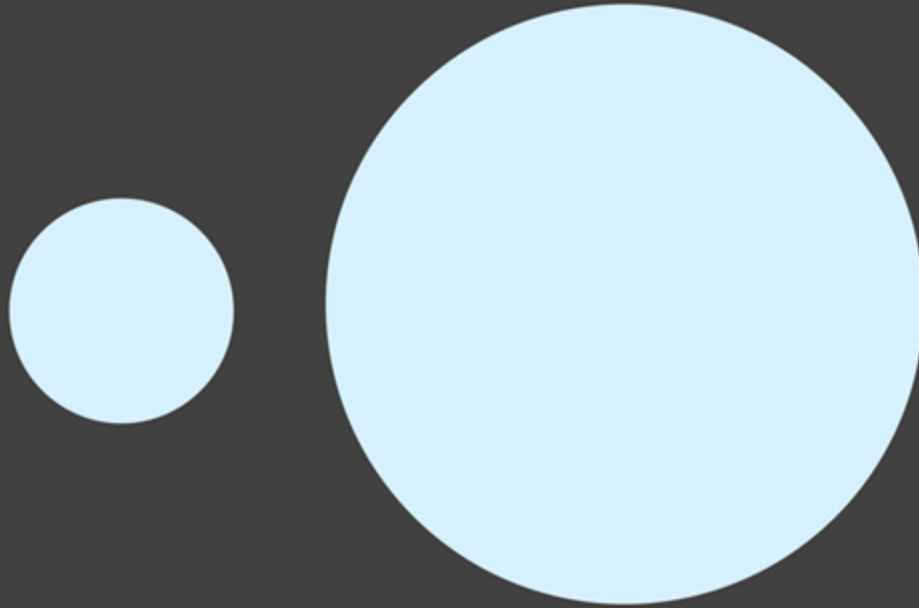
The accuracy of our judgements depend on the type of marking.



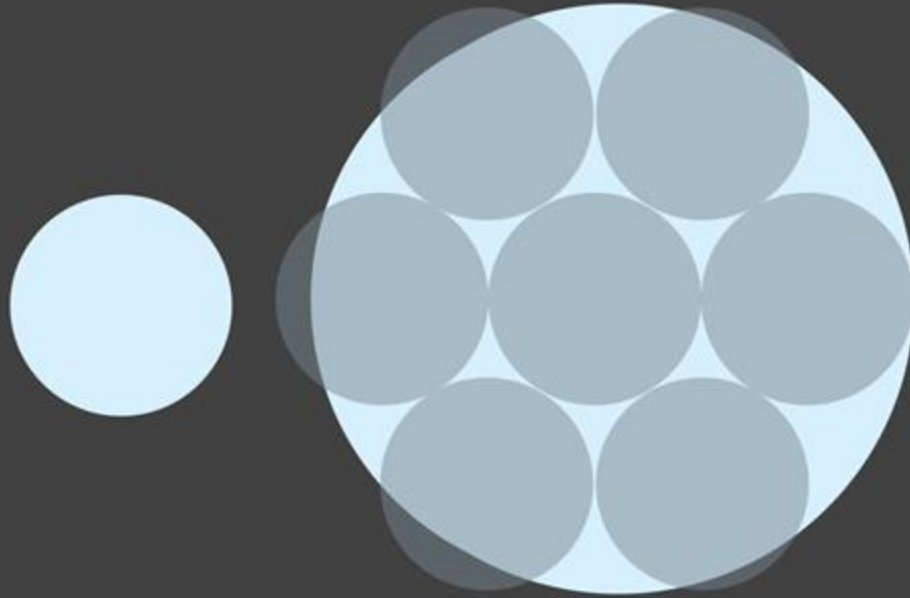
How much longer is the top bar?



The top bar is 7 times longer than the bottom bar.

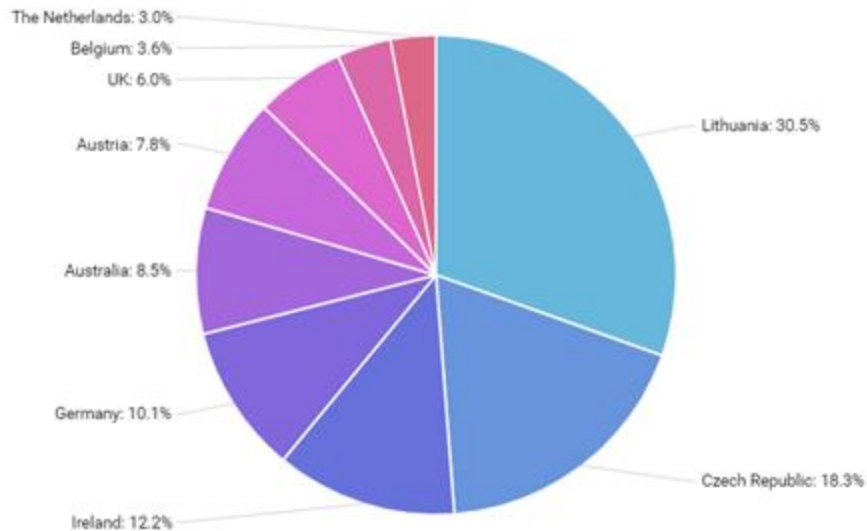
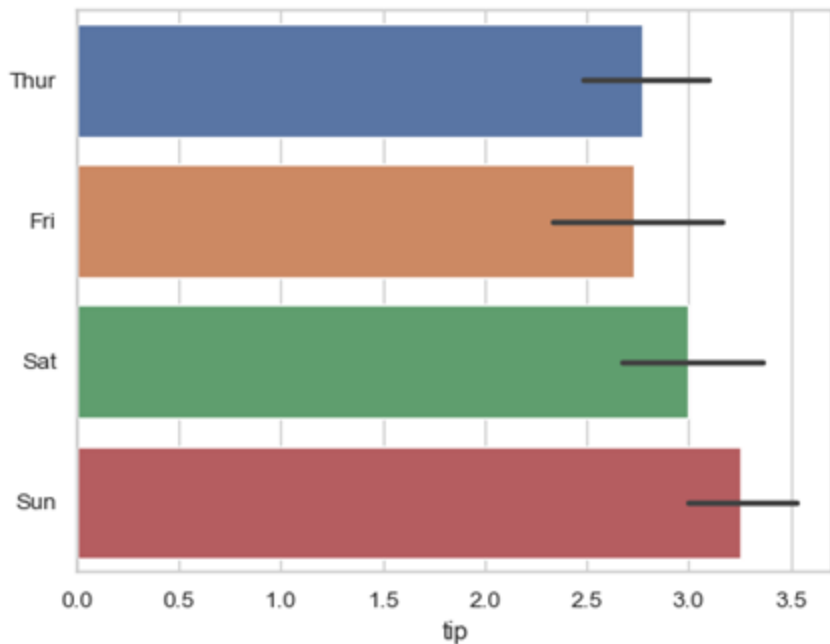


**How much bigger is the big circle?**



The area of the big circle is 7 times larger than the area of the small circle.

Lengths are easy to distinguish; angles (portions) are hard



**Don't use pie charts!** Angle judgements are inaccurate, and misleading



# Areas are hard to distinguish

## African Countries by GDP

### TOP COUNTRIES BY GDP IN U.S. \$ BILLIONS

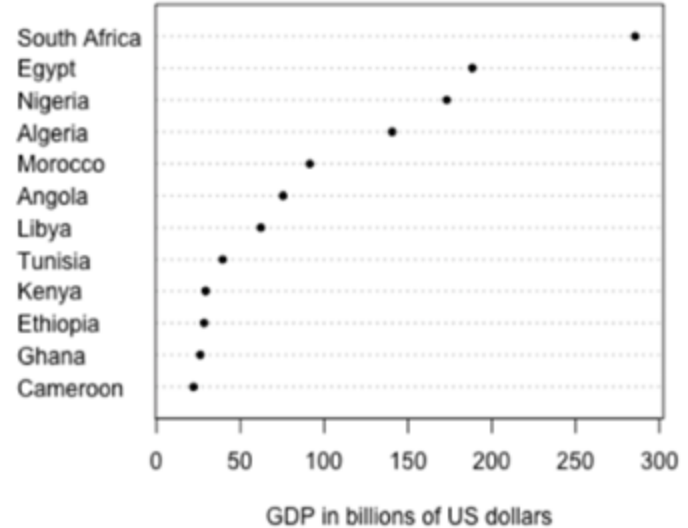
Shows domestic product (GDP) refers to the market value of all final goods and services produced within a country in a given period (2007 - 2008).

### GDP CALCULATION

private consumption + gross investment + government spending + exports - imports

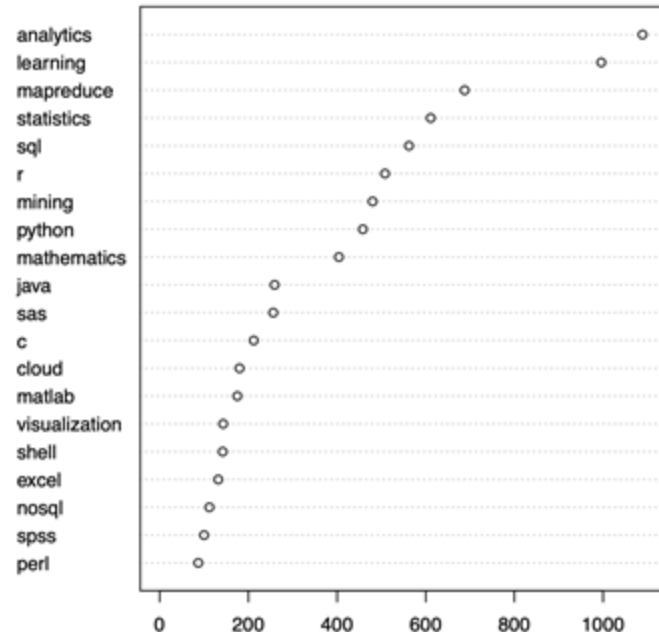
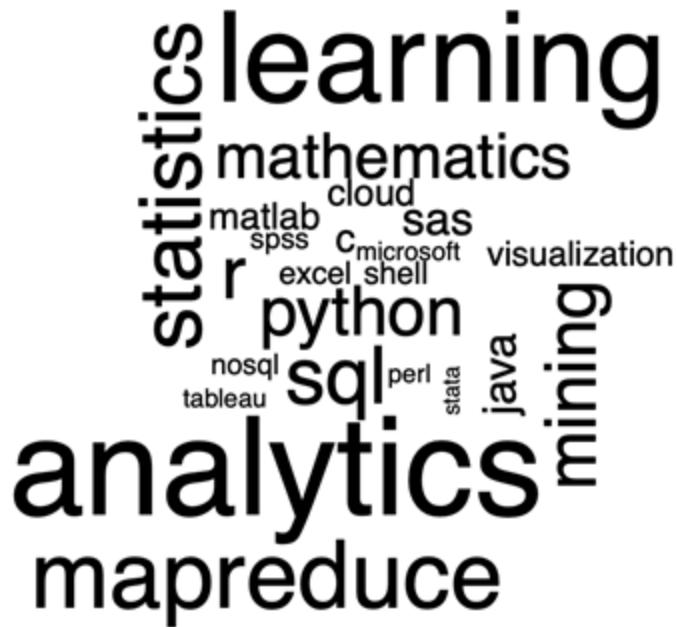


## African Countries by GDP



**Avoid area charts!** Area judgements are inaccurate. (For instance, South Africa has twice the GDP of Algeria, but that isn't clear from the areas. But easy to see which is bigger )

Areas are hard to distinguish



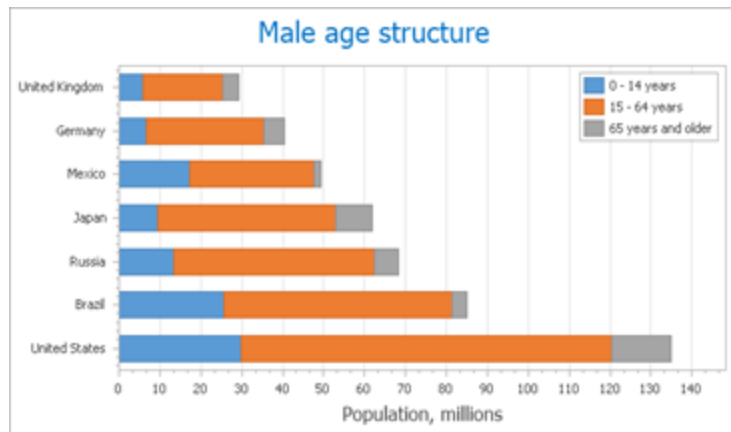
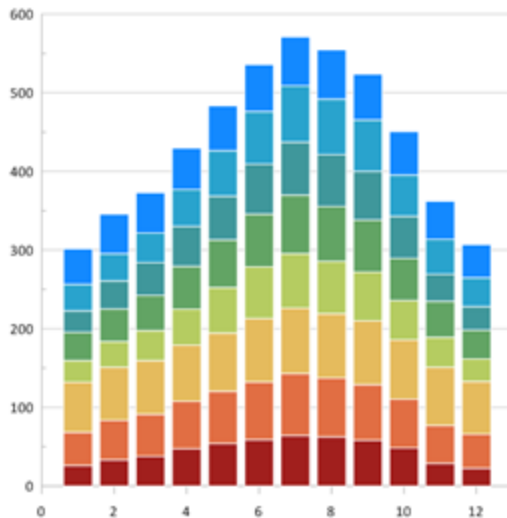
**Avoid word clouds too!** It's hard to tell the area taken up by a word.

# Avoid “jiggling” the baseline

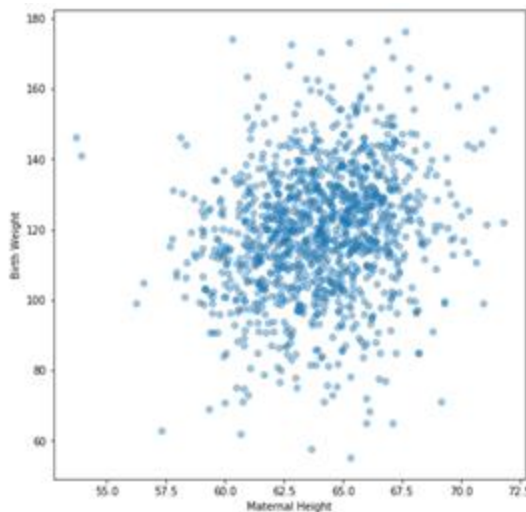
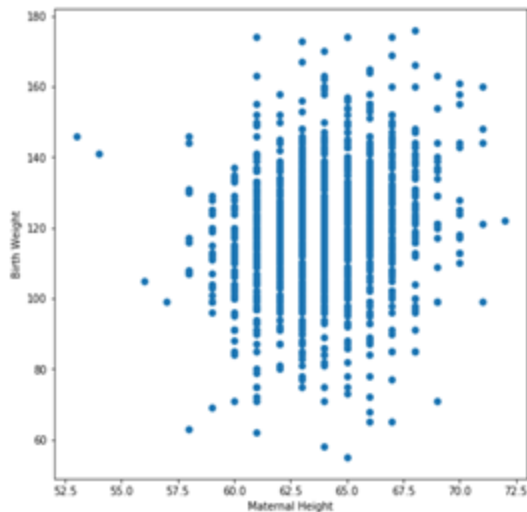
Stacked bar charts, histograms, and area charts are hard to read because the baseline moves.

- In the first plot, the top blue bars are all roughly of the same length. But that’s not immediately obvious!
- In the second plot, comparing the number of 15-64 year old males in Germany and Mexico is difficult.

The blue bar has a similar baseline, but others have jiggling baselines



# Related – overplotting



In the plot on the left, it's hard to tell exactly how many points are being visualized.

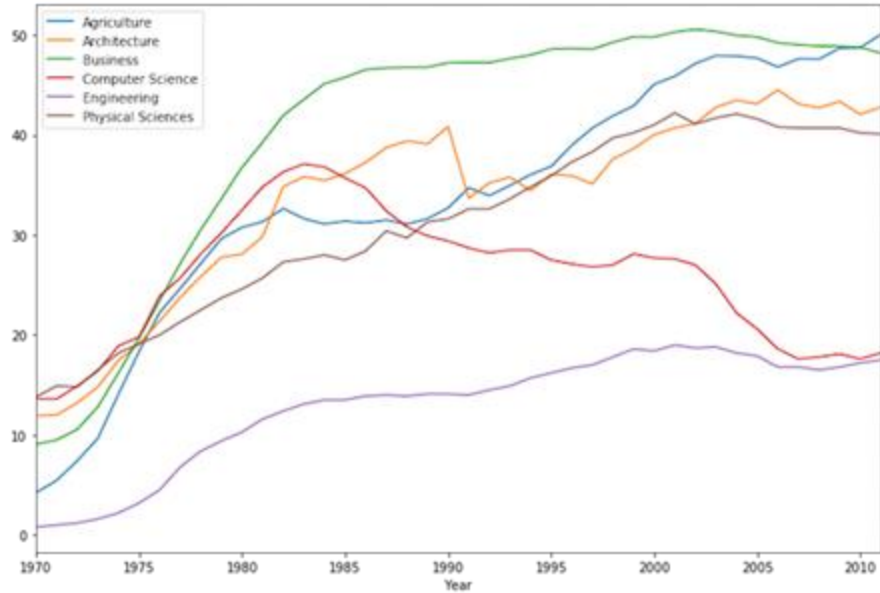
- Many on top of one another.
- Observations only on lattice points (integers)

## Some solutions:

- Add small random noise to both x and y (“jittering”).
- Make points smaller (wouldn't help here though).

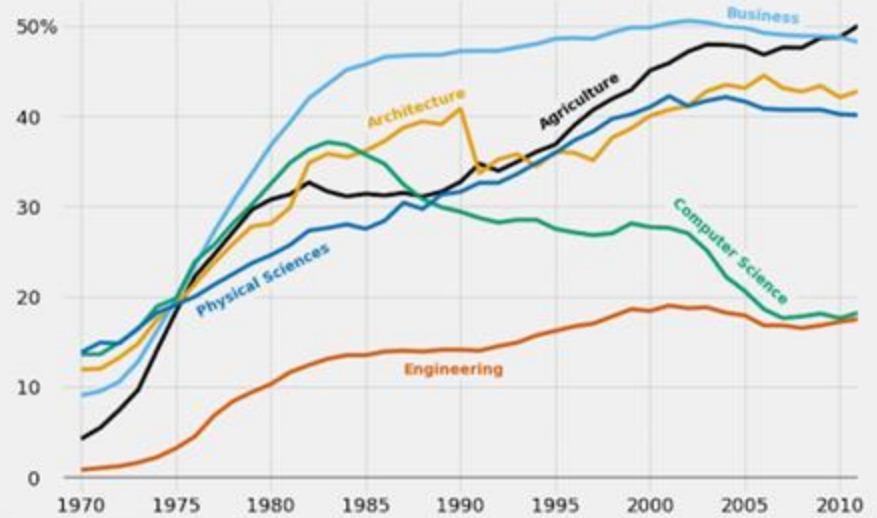
Context

Which graph is easier to read?



### The gender gap is transitory - even for extreme cases

Percentage of Bachelors conferred to women from 1970 to 2011 in the US for extreme cases where the percentage was less than 20% in 1970



©DATAQUEST

Source: National Center for Education Statistics

Python cannot directly put legend text on the graph. Putting description text helps.

# Add context directly to plot

A publication-ready plot needs:

- Informative title (takeaway, not description).
  - “Older passengers spend more on plane tickets” instead of “Scatter plot of price vs. age”.
- Axis labels.
- Reference lines, markers, and labels for important values.
- Legends, if appropriate.
- Captions that describe the data.

The plots you create in this class always need titles and axes labels, too.

# Captions

A picture is worth a thousand words, but not all thousand words you want to tell may be in the picture. In many cases, we need captions to help tell the story.

Captions should be:

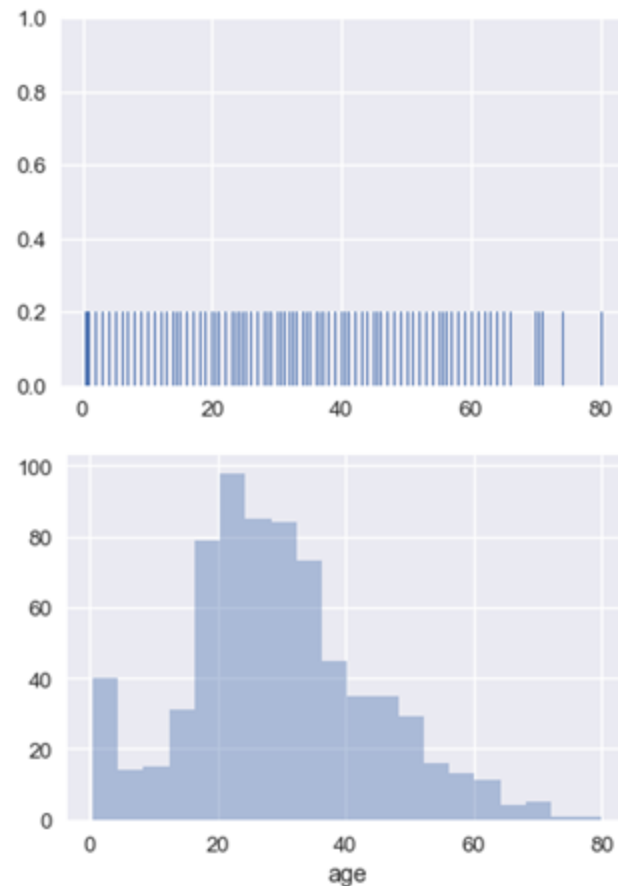
- Comprehensive and concise.
- Describe what has been graphed.
- Draw attention to important features.
- Describe conclusions drawn from graph.



Smoothing

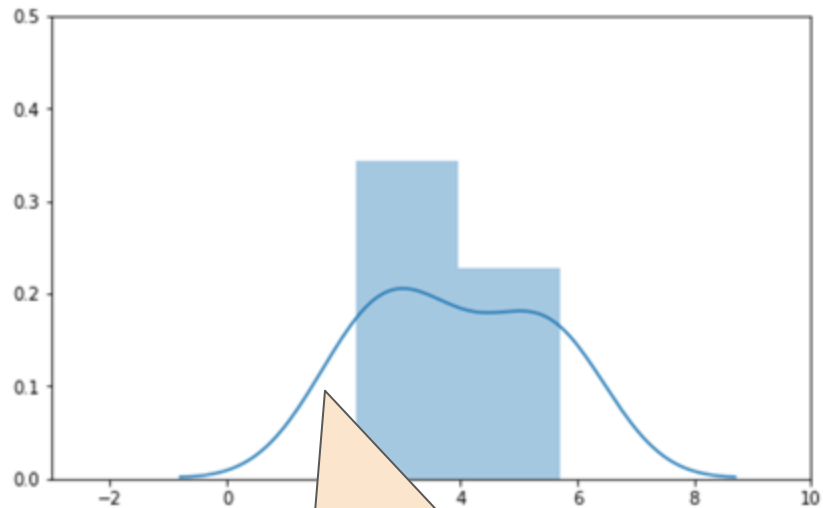
# Smoothing

- Histograms are a smoothed version of rug plots.
- We smooth if we want to focus on general structure rather than individual observations.



# Kernel density estimation (KDE)

Kernel Density Estimation is used to estimate a **probability density function (PDF)** (or density curve) from a set of data.



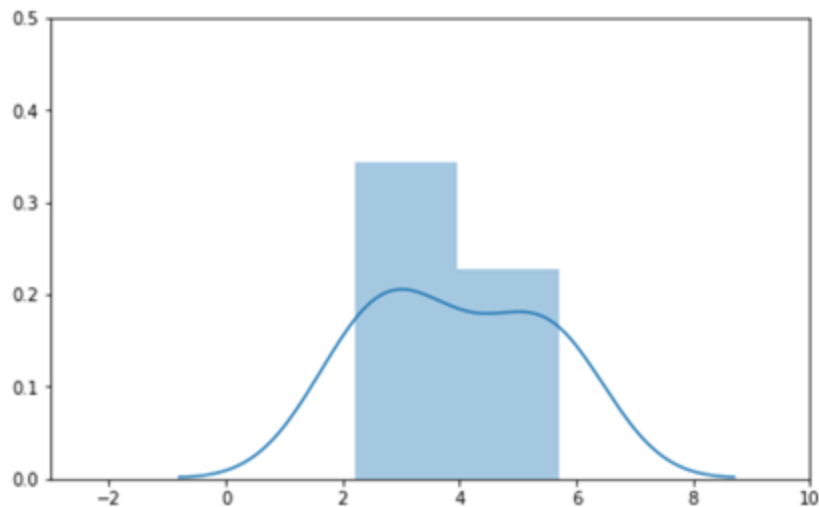
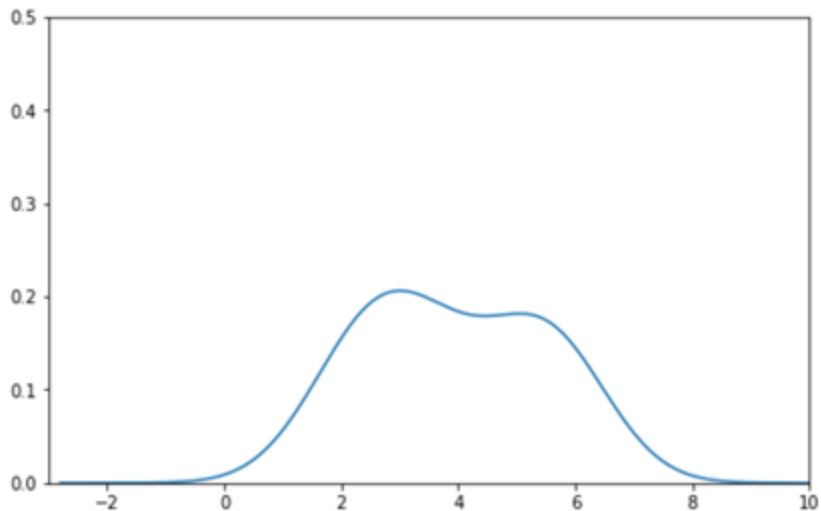
Our goal is to recreate this smooth curve ourselves.

# Kernels

- A kernel (for our purposes) is a valid density function. That means it:
  - Must be non-negative for all inputs.
  - Must integrate to 1.
- The most common kernel is the **Gaussian** kernel.

# Kernel density estimates

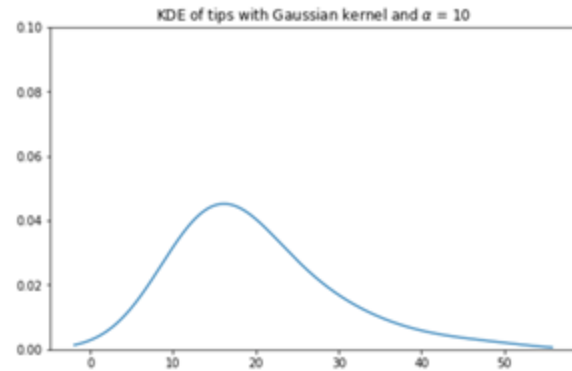
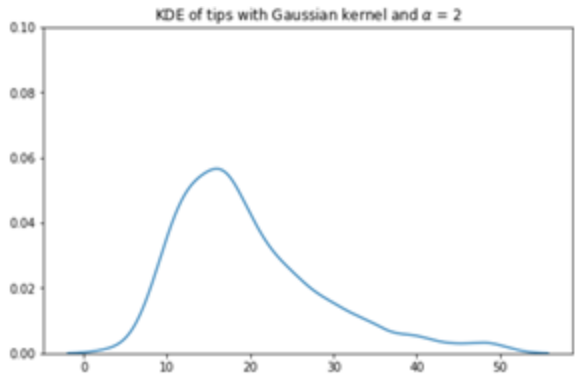
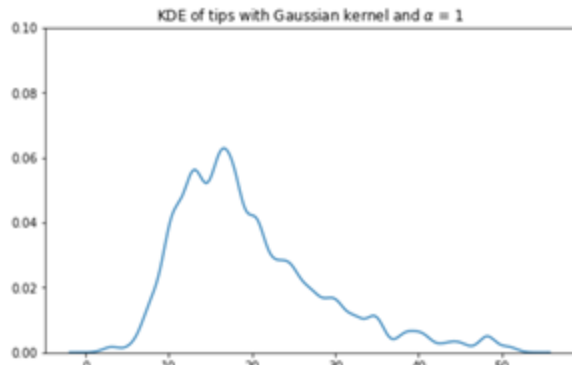
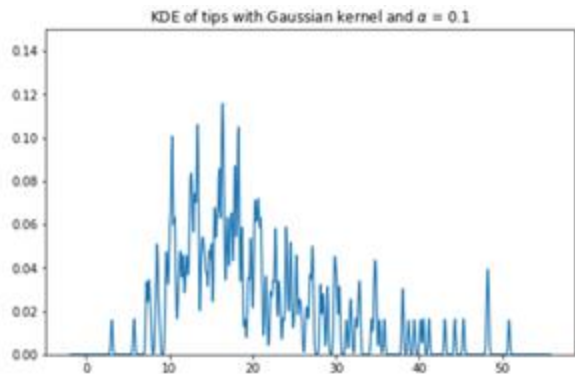
The curve we manually created (left) exactly matches the one that `sns.distplot` creates for us (right)!



# Effect of bandwidth on KDEs

Bandwidth is analogous to the width of each bin in a histogram.

- As  $\alpha$  (alpha) increases, the KDE becomes more smooth.
- Simpler to understand, but gets rid of potentially important distributional information.
- We call  $\alpha$  a **hyperparameter**. Be familiar with this term! This means it is manually set, not determined from the data.



Bimodal

Unimodal

Over-smoothing

## Summary of KDE

$$f_{\alpha}(x) = \frac{1}{n} \sum_{i=1}^n K_{\alpha}(x, x_i)$$

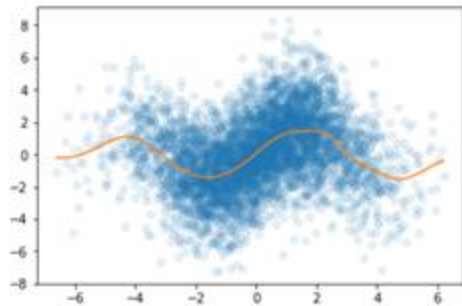
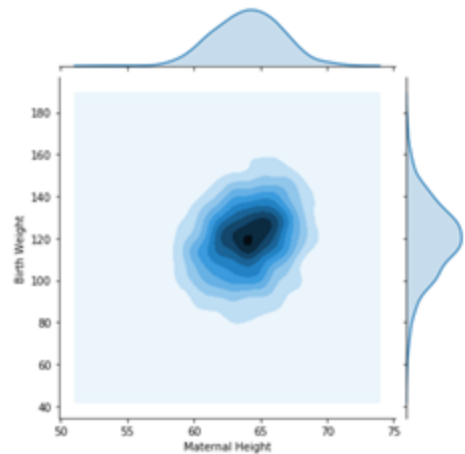
The “KDE formula” is above  
(Not required to memorize)

<https://towardsdatascience.com/kernel-density-estimation-explained-step-by-step-7cc5b5bc4517>

To make visually appealing graphs of the PDF for any dataset.

# Extensions

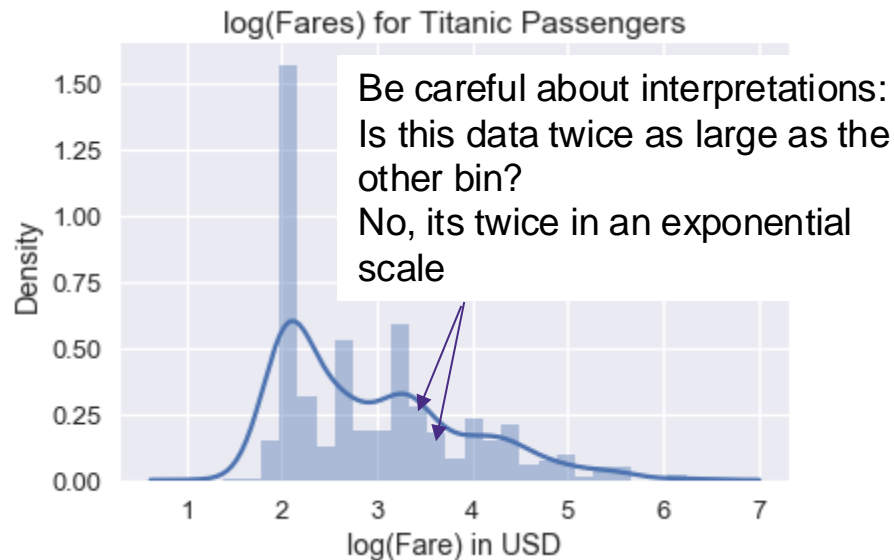
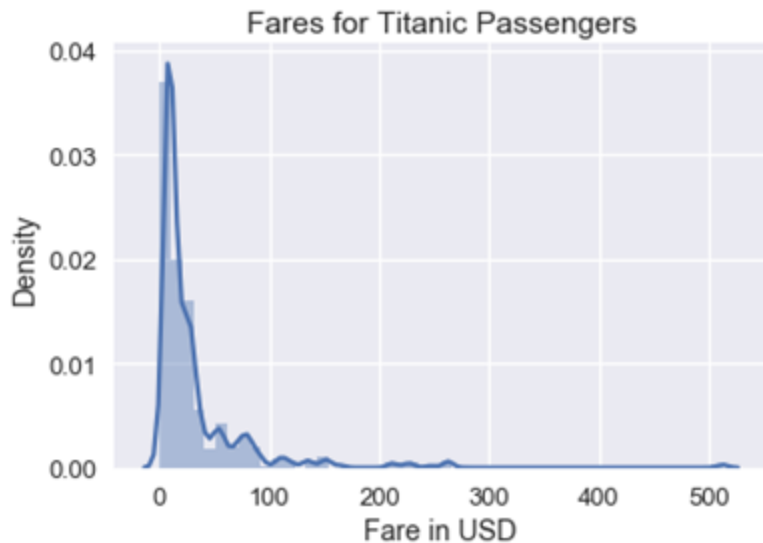
- One can extend the idea of kernel density estimation to two dimensions.
  - A contour plot is a two dimensional KDE (top).
- One can also use kernels to create smoothed versions of scatterplots (bottom).
  - Won't do that in this course





# Transformations

# Transforming data can reveal patterns

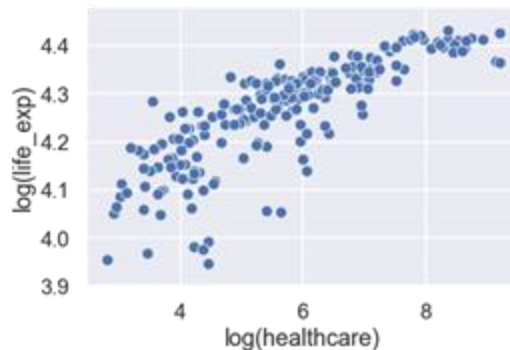
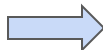
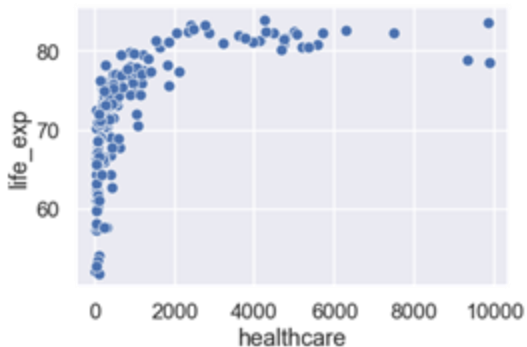


When a distribution has a large dynamic range, it can be useful to take the log.

# Why straighten relationships?

Now, we will look at how to **linearize** the scatter plot of two variables. Why?

- If we know what transformation made our plot of  $y$  vs.  $x$  linear, we can “backtrack” to figure out the exact relationship between  $x$  and  $y$ .
- Linear relationships are particularly simple to interpret.
  - We know what slopes and intercepts mean.



# Log of y-values

If we take the log of our y-values and notice a linear relationship, we can say (roughly) that

$$\log y = ax + b$$

Working backwards:

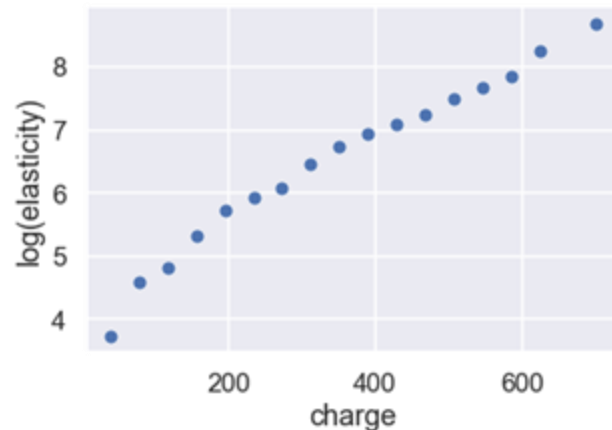
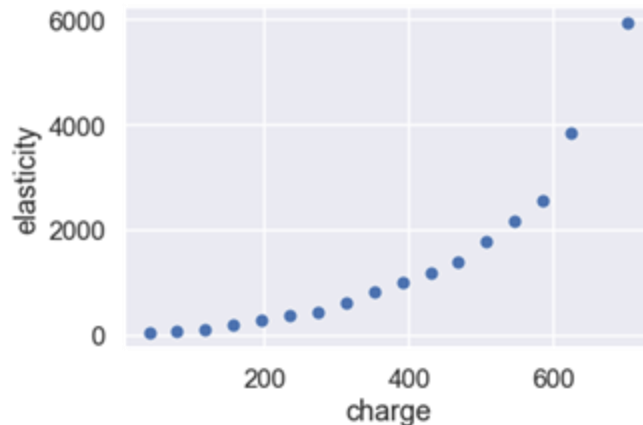
$$\log y = ax + b$$

$$y = e^{ax+b}$$

$$y = e^{ax} e^b$$

$$y = Ce^{ax}$$

This implies an **exponential** relationship in the original plot.



# Log of both x and y-values

If we take the log of both axes and notice a linear relationship, we can say (roughly) that

$$\log y = a \cdot \log x + b$$

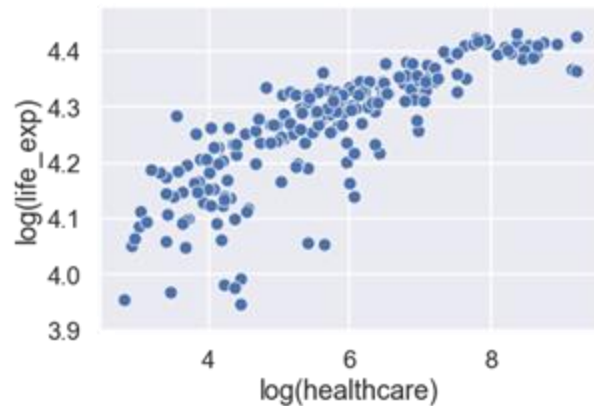
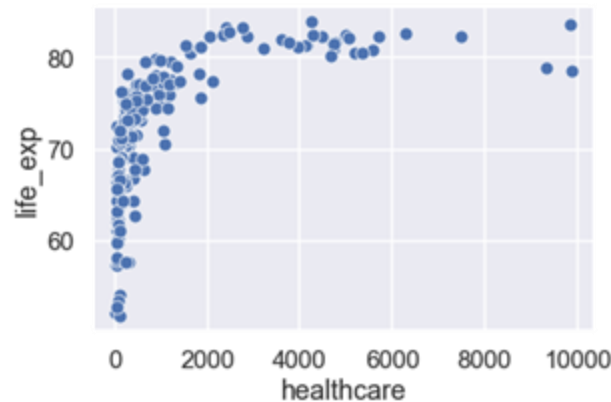
Working backwards:

$$y = e^{a \cdot \log x + b}$$

$$y = C e^{a \cdot \log x}$$

$$y = C x^a$$

This implies a **power** relationship in the original plot (a one-term **polynomial**)



Log transform as a “Swiss army knife”

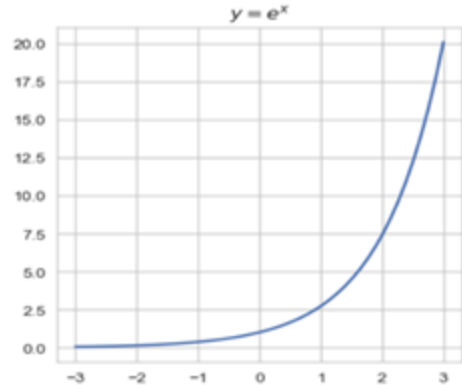
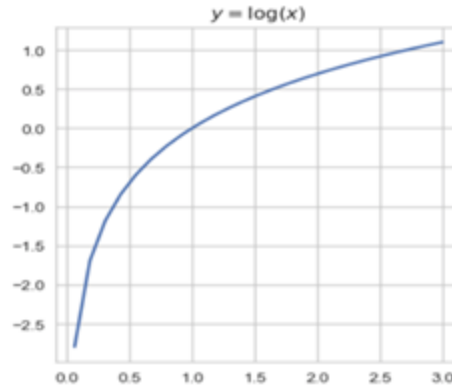
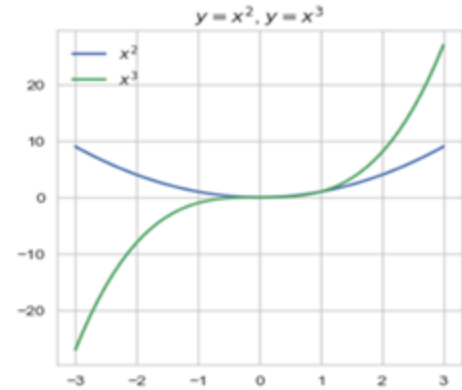
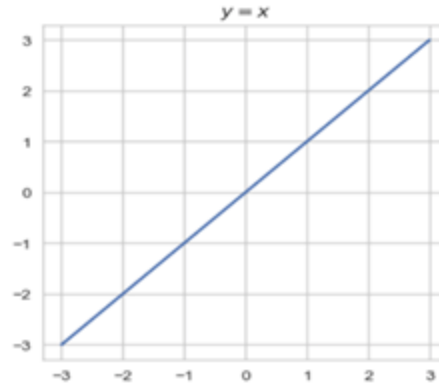
$$y = a^x \rightarrow \log(y) = x \log(a)$$

$$y = ax^k \rightarrow \log(y) = \log(a) + k \log(x)$$

Properties of logarithms make them very powerful!

# Basic functional relations

Knowing the general shapes of polynomial (quadratic and cubic), exponential, and logarithmic curves (regardless of base) will go a long way.

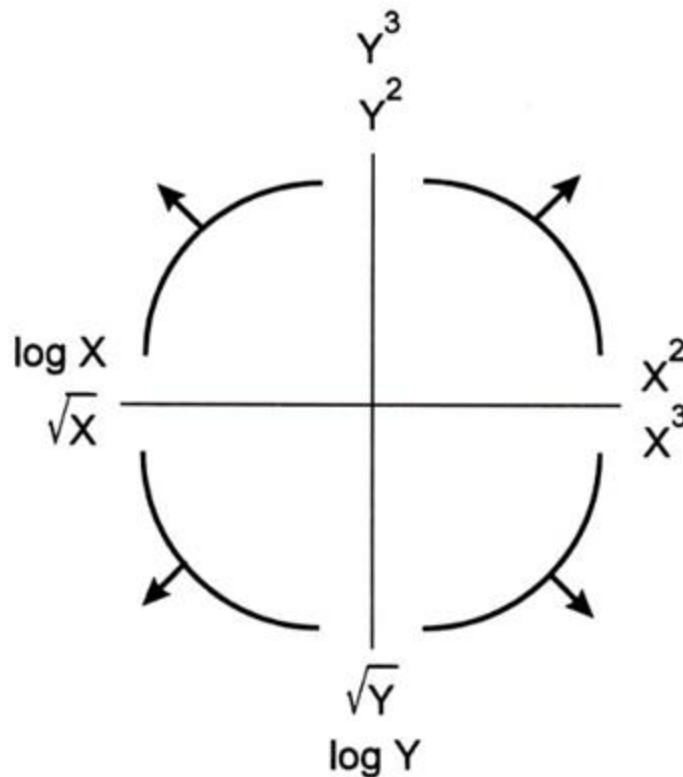


# Tukey-Mosteller Bulge Diagram

This diagram can help us choose which transformation(s) to apply to our data in order to linearize it.

- There are multiple solutions. Some will fit better than others.
- sqrt and log make a value “smaller”. Raising to a value to a power makes it “bigger”.
- Each of these transformations equates to increasing or decreasing the scale of an axis.

**(demo in jupyter)**





# Summary

- Choose appropriate scales.
- Condition in order to make comparisons more natural.
- Choose colors and markings that are easy to interpret correctly.
- Add context and captions that help tell the story.
- Smoothed estimates of distributions help with big-picture interpretation.
  - Kernel Density Estimates are a method of smoothing data.
- Transforming our data can linearize relationships.
  - Helpful when we start linear modeling next lecture.
- **More generally – reveal the data!**
  - Eliminate anything unrelated to the data itself – “chart junk.”
  - It’s fine to plot the same thing multiple ways, if it helps fit the narrative better.