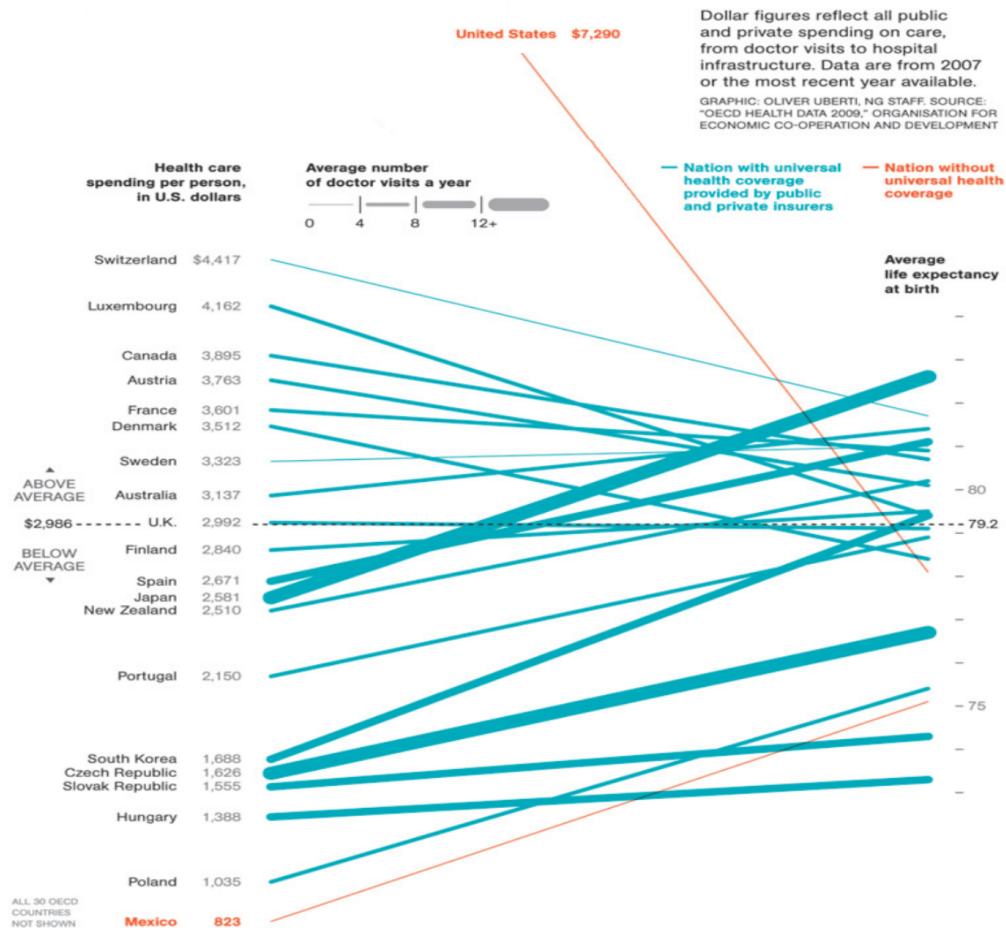


Week 3 - Discussion

Name:

Data Visualization



1. The first part of the discussion will be centered on the above visualization.

- (a) Five variables are being represented visually in this graphic. What are they and what are their types (ie qualitative, quantitative, nominal, ordinal)?

- (b) How are the variables represented in the graphic, e.g., the variable XXX is mapped to the x -axis, the variable WWW is mapped to the y -axis, the variable ZZZ is conveyed through color, etc.?

- (c) How can we figure out how to interpret the visual qualities of the plot, e.g., how do we know what a color represents?

(d) What purpose does the comment at the top right of the plot serve?

(e) Make 3 observations about the figure. Describe the feature that you are basing your observation on.

For example, South Korea's expenditure on health care is comparable to Eastern European countries (and among the lowest of all countries plotted), but the life expectancy is much higher than the Eastern European countries. In the plot we see that the left endpoint of South Korea's line segment is near the Eastern European countries, but the slope of the line segment is much steeper.

- (f) Consider the steep negative slope and narrowness of the line segment that represents the data for the United States. What systemic, social, or societal issues might explain this?

2. Name some appropriate 2D visualizations if your goal is to explore:

(a) The distribution of population for various cities.

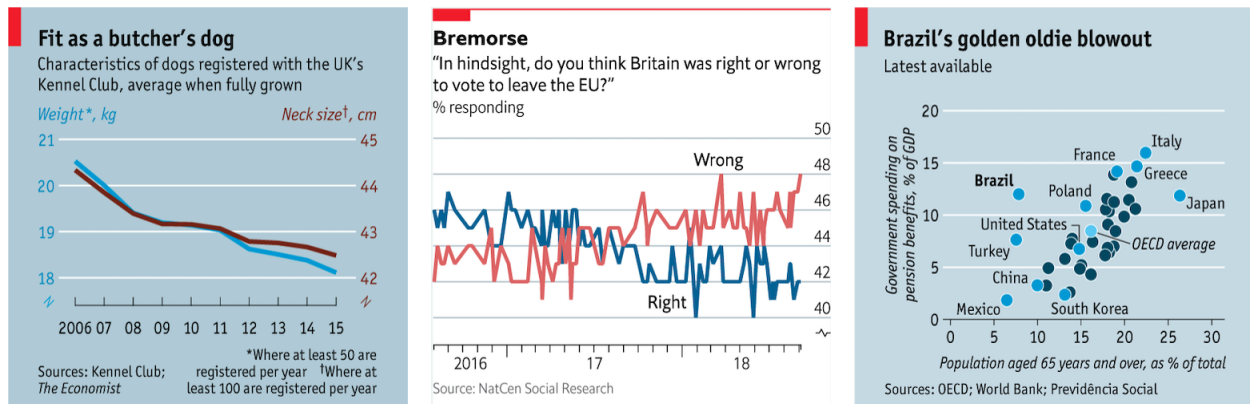
(b) The distribution of income.

(c) The relationship between income and life expectancy.

(d) The relationship between income and city.

(e) The relationship between income, life expectancy, smoking status (non-smoker, social smoker, smoker), and city.

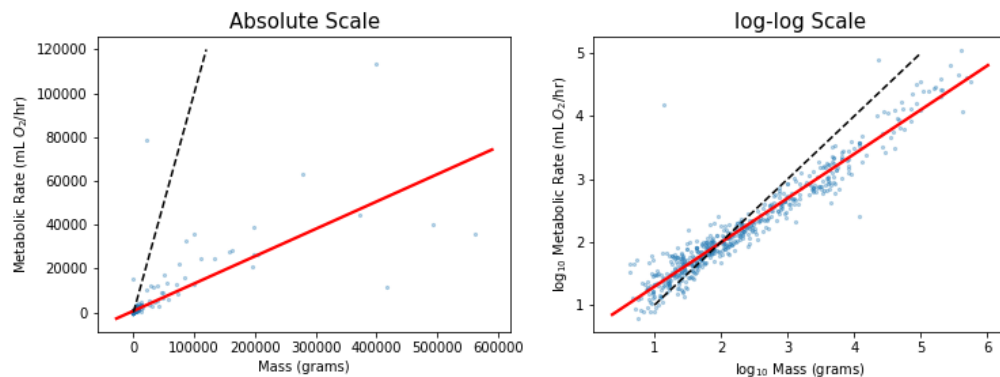
3. Creating visualizations that represent data accurately and that support the narrative we wish to create is no easy task. Even the journalists and editors at *The Economist*, a newspaper known for its compelling, data-driven articles, have been known to make blunders. Three of their ill-thought-out plots are presented below. Consider what aspects of the visualizations are misleading, and think of ways in which you can remedy them.



Hint: The datapoints in the rightmost plot are shaded based on whether or not they are labeled.

Logarithmic Transformations

4. One of your friends at a biology lab asks you to help them analyze panTHERIA, a database of mammals. They are interested in the relationship between mass, measured in grams, and metabolic rate (“energy expenditure”), measured by oxygen use per hour. Originally, they show you the data on a linear (absolute) scale, shown on the left. You notice that the values on both axes vary over a large range with many data points clustered around the smaller values, so you suggest that they instead plot the data on a log-log scale, shown on the right. The solid red line is a “line of best fit” (we’ll formalize this later in the course) while the black dashed line represents the identity line $y = x$.



- (a) Let C and k be some constants and x and y represent mass and metabolic rate, respectively. Based on the plots, which of the following best describe the pattern seen in the data? Reminder: $\log(ab) = \log(a) + \log(b)$.
- ☐ A. $y = C + kx$ ☐ B. $y = C \times 10^{kx}$ ☐ C. $y = C + k \log_{10}(x)$ ☐ D. $y = Cx^k$
- (b) What parts of the plots could you use to make initial guesses on C and k ?
- (c) Your friend points to the solid line on the log-log plot and says “since this line is going up and to the right, we can say that, in general, the bigger a mammal is, the greater its metabolic rate”. Is this a reasonable interpretation of the plot?

5. When making visualizations, what are some reasons for performing log transformations on the data?