

Calling Bullsh*t on Data Visualisations!

CM4125 – Topic 3

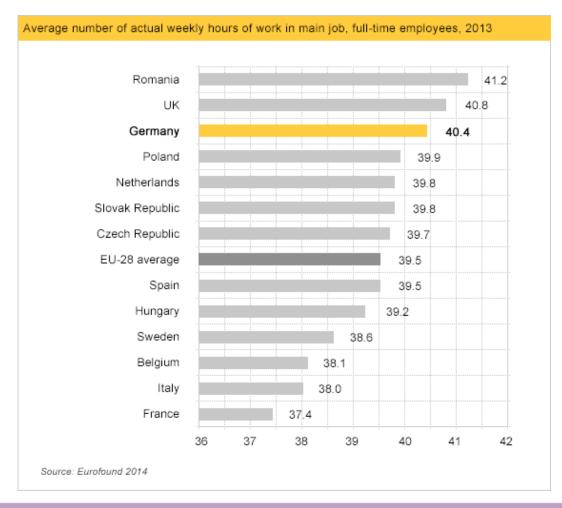


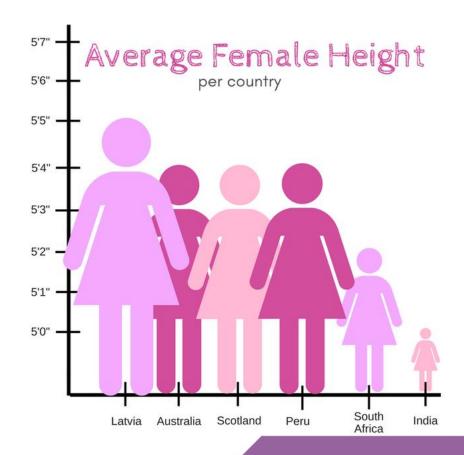
What this lecture is about?

- Based on Week 6 of the Calling Bullsh*t course of the University of Washington by Carl T. Bergstrom and Jevin West
 - https://www.callingbullshit.org/syllabus.html#Visual
- Supplementary reading:
 - Alberto Cairo (2019) *How Charts Lie: Getting Smarter about Visual Information*. W.W. Norton and Company.
 - Edward Tufte (1983) *The Visual Display of Quantitative Information*. Chapters 2 (Graphical integrity) and 5 (Chartjunk: vibrations, grids, and ducks).



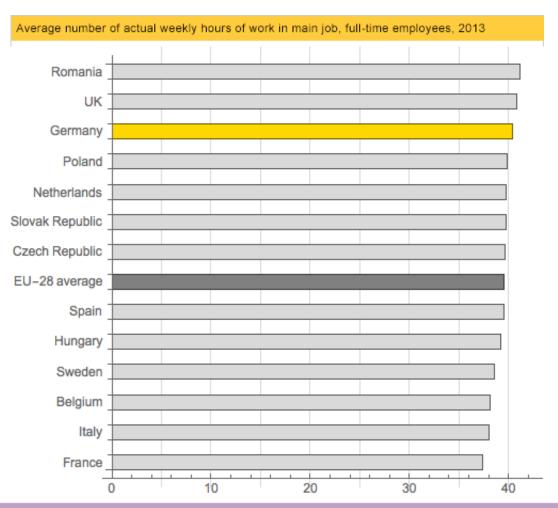
1. Bar chart y-axes <u>need to</u> include zero



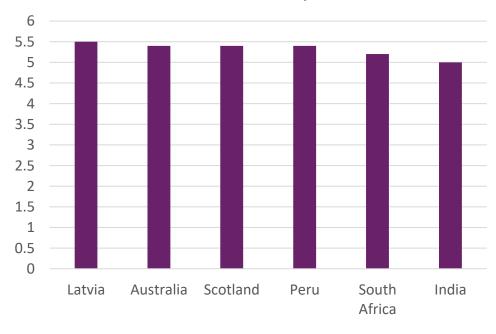




1. Bar chart y-axes **need to** include zero



Average Female Height Per country

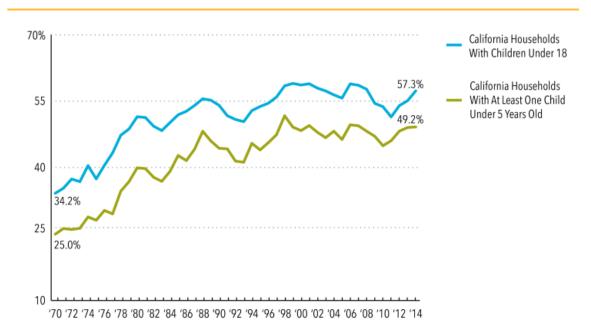




2. Line graph y-axes don't need to include zero

More California Households Have All Parents Working, Making Access to Child Care an Important Priority

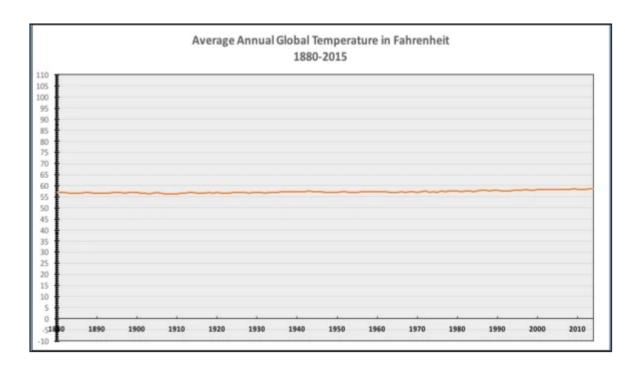
Percentage of California Households Where All Parents Work, 1970 to 2014



Note: A "household where all parents work" includes single-parent households and dual-earner households. Parents include stepparents and adoptive parents.

Source: Budget Center analysis of US Census Bureau data

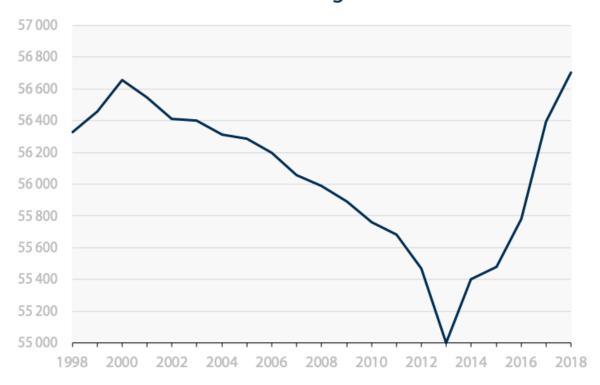




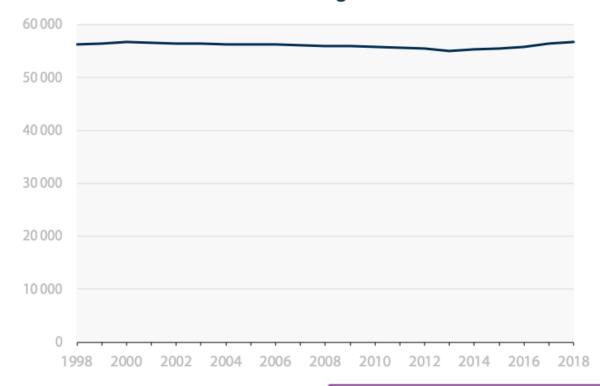


Well...

Sales of widgets

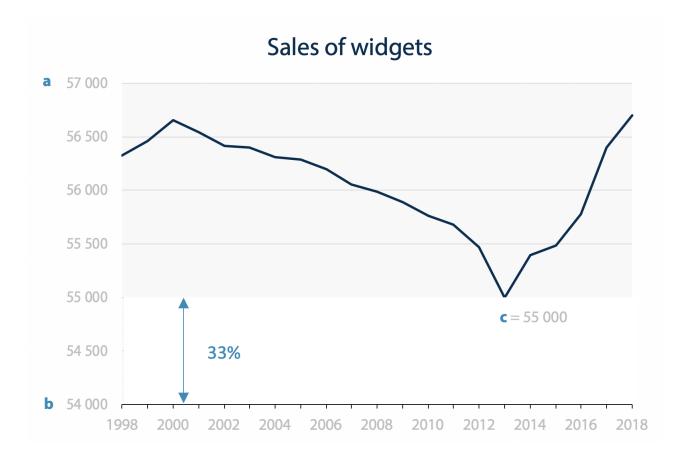


Sales of widgets





Solution



Empty space = Minimum data value - Minimum scale value

Maximum scale value - Minimum scale value

$$x = \frac{(c-b)}{(a-b)}$$

$$b = \frac{3c - a}{2}$$

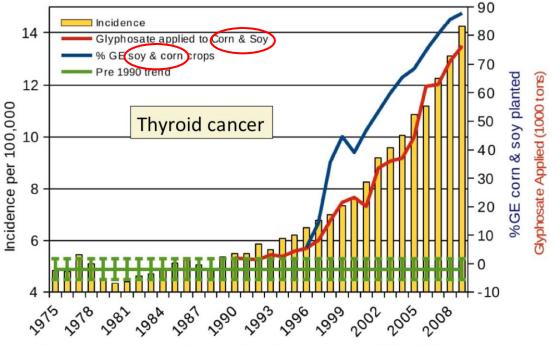
$$b = \frac{3 * 55,000 - 57000}{2} = 54\,000$$



3. No multiple axes on a single graph

Thyroid Cancer Incidence Rate (age adjusted)

plotted against glyphosate applied to U.S. corn & soy (R = 0.988, p <= 7.612e-09) along with %GE corn & soy crops R = 0.9377, p <= 2.152e-05 sources: USDA:NASS; SEER

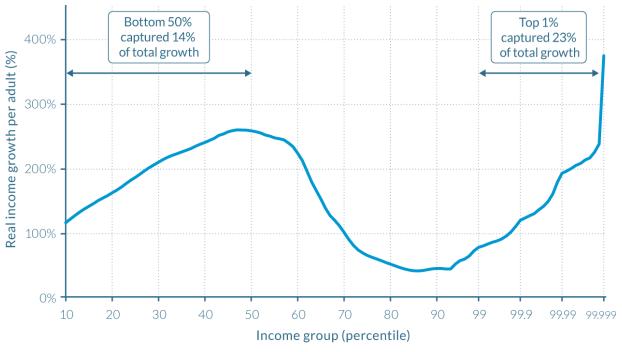


*Figure 10, Swanson et al. Journal of Organic Systems 2014; 9(2):6-37.

How to solve this?



4. An axis should **not** change scales midstream



Use log scales!

Why not use the golden ratio?!

Source: WID.world (2017). See wir 2018.wid.world/methodology.html for data series and notes.

On the horizontal axis, the world population is divided into a hundred groups of equal population size and sorted in ascending order from left to right, according to each group's income level. The Top 1% group is divided into ten groups, the richest of these groups is also divided into ten groups, and the very top group is again divided into ten groups of equal population size. The vertical axis shows the total income growth of an average individual in each group between 1980 and 2016. For percentile group p9p99.1 (the poorest 10% among the world's richest 1%), growth was 77% between 1980 and 2016. The Top 1% captured 23% of total growth over this period. Income estimates account for differences in the cost of living between countries. Values are net of inflation.



Example (of what NOT to do)



03 October 2023 10

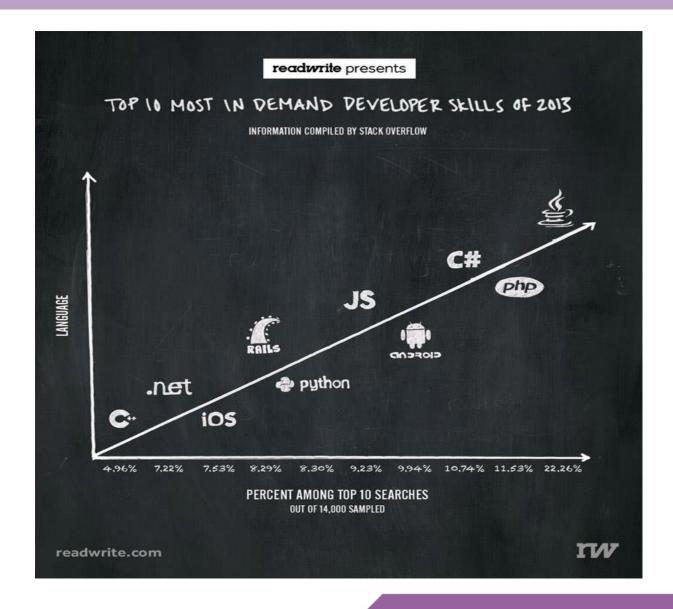




03 October 2023 11



5. An axis should have something on it

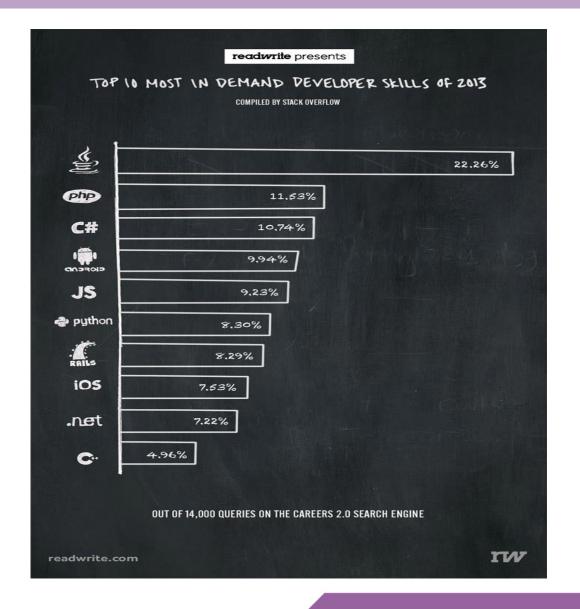




Like this

Tools such as

WebPlotDigitizer can
help you get the
missing info



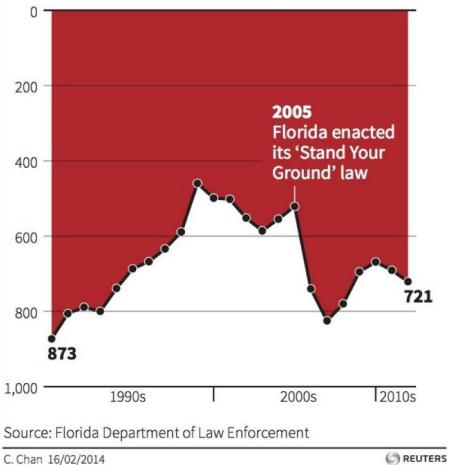


6. **Don't** invert the axis(?)

The author defended it! What's your stance?

Gun deaths in Florida

Number of murders committed using firearms

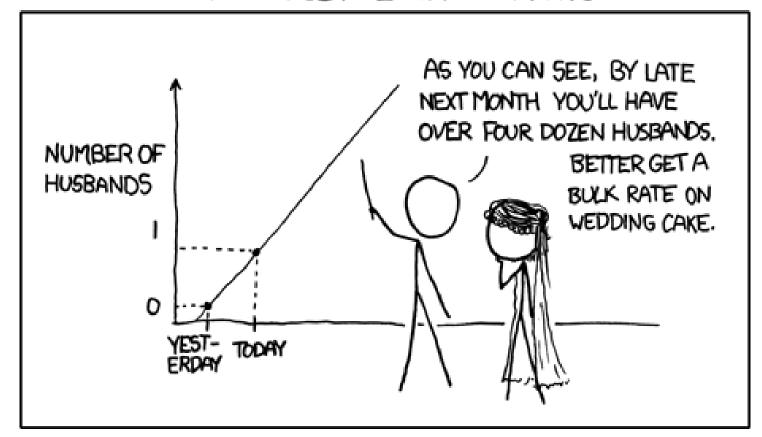






7. Avoid extrapolation!

MY HOBBY: EXTRAPOLATING



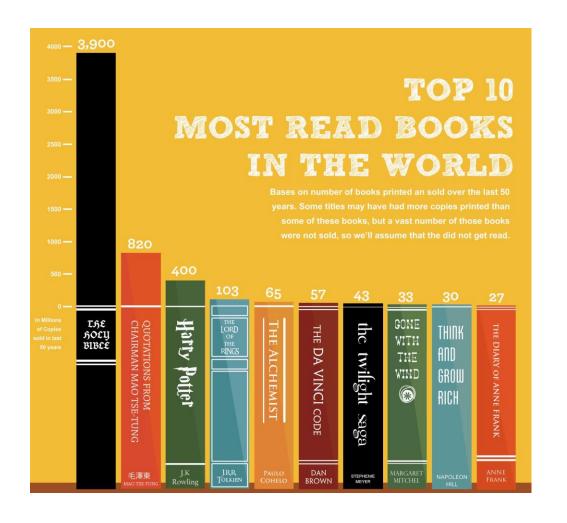


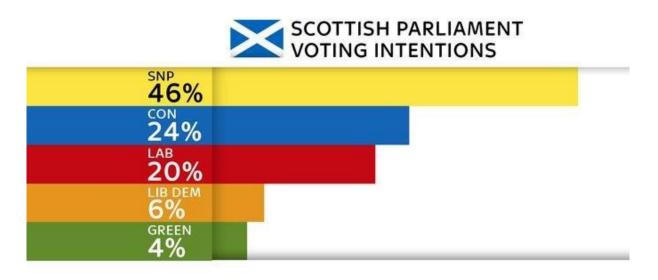
8. Proportional Ink

"When a shaded region is used to represent a numerical value, the area of that shaded region should be directly proportional to the corresponding value"

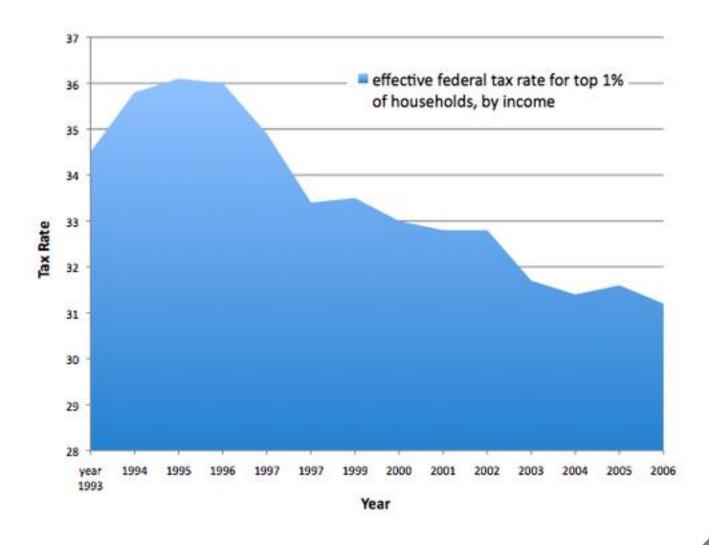
Extends the argument for misleading axes



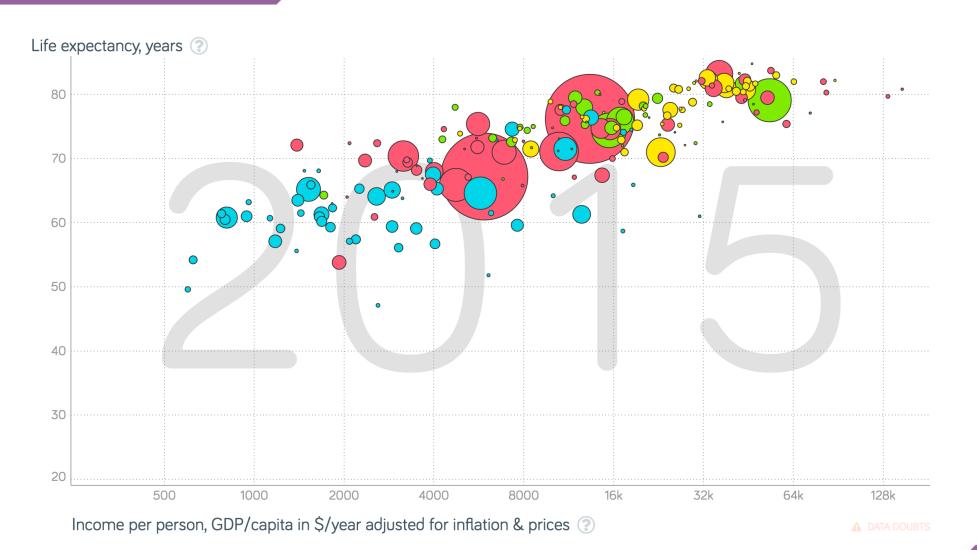




SOURCE: OPINIUM / SKY NEWS

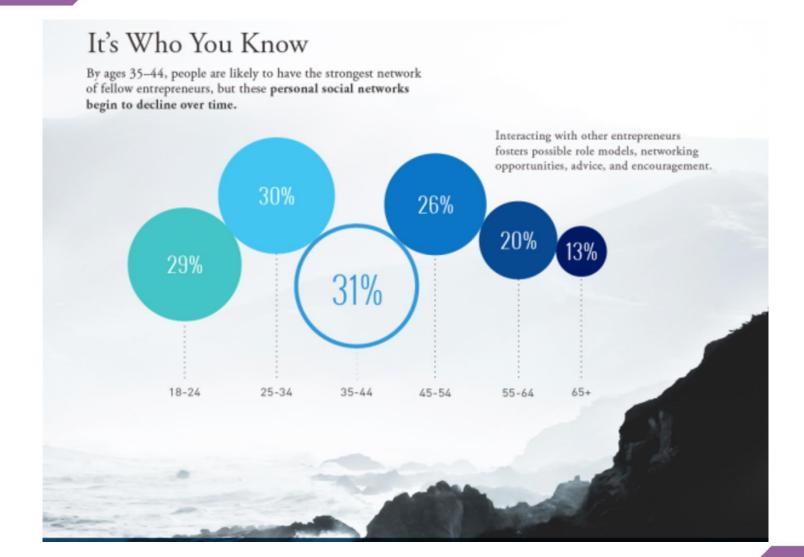






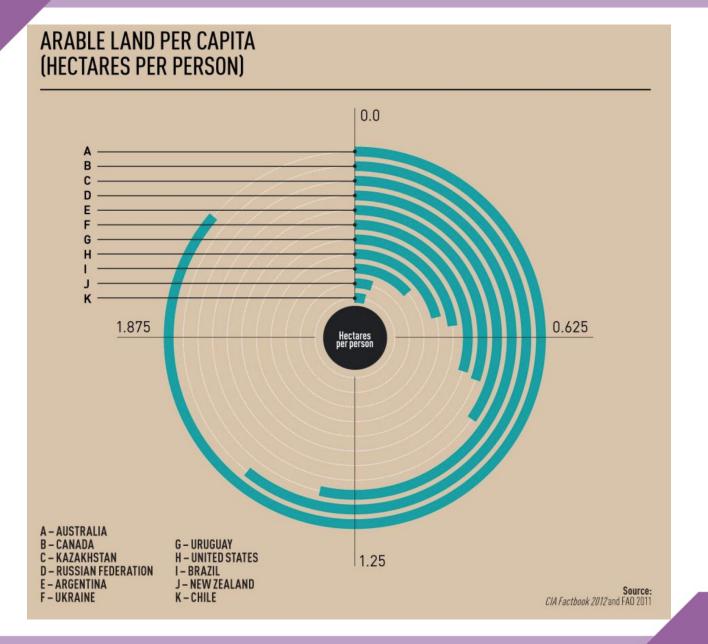
Should we violate the principle?



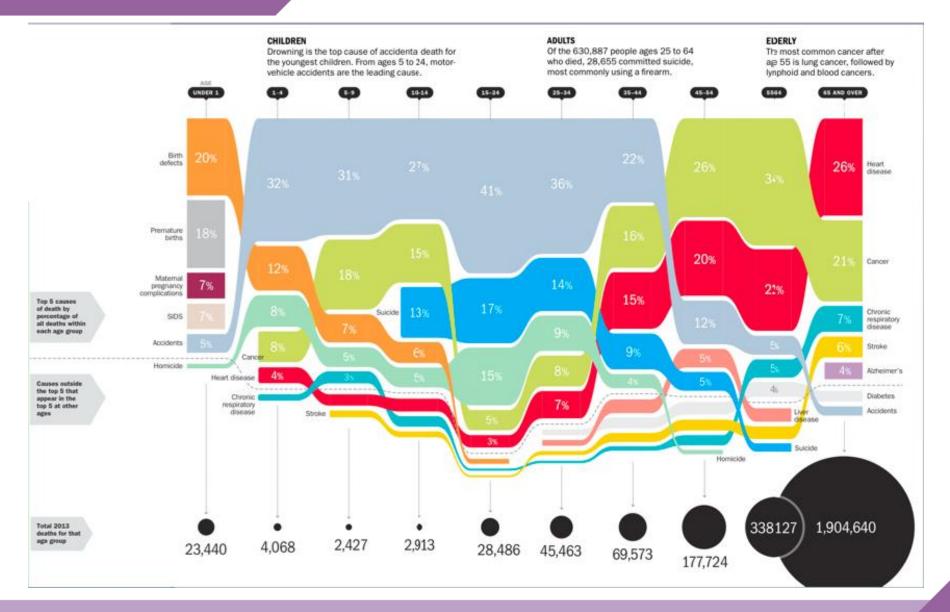


Radius ≠ Area





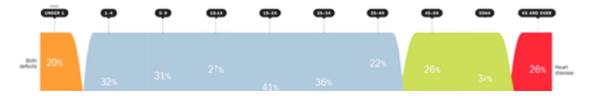
ROBERT GORDON UNIVERSITY ABERDEEN

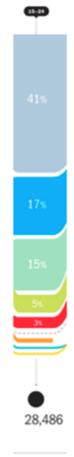




Along vertical slices, ink is proportional to value because shaded areas represent the fraction of a fixed number of deaths (here 28,486) from each cause.

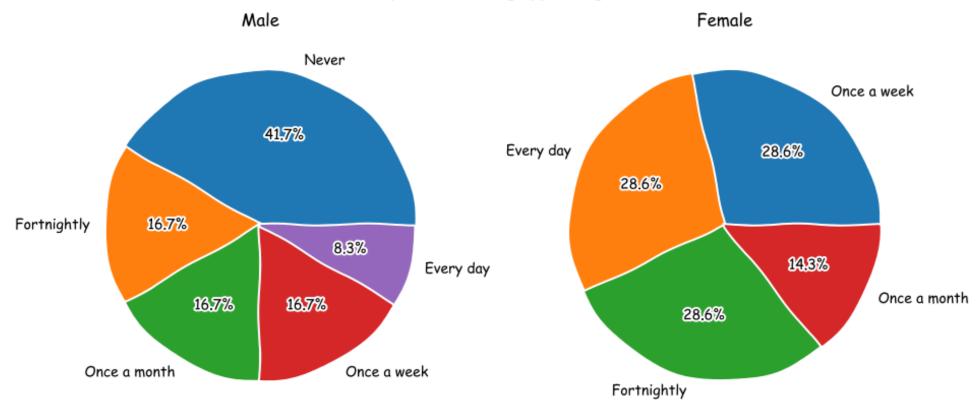
Along horizontal slices, ink is not proportional because total deaths differ widely by age group. Far more people 65 and older died of heart disease (red) than children age 1-4 die of accidents (blue-gray), but the latter takes more ink because it represents a larger percentage of the (relatively few) total deaths at that age.







How often do you use a dating app during a month?



What is wrong in this one?



9. Perspective



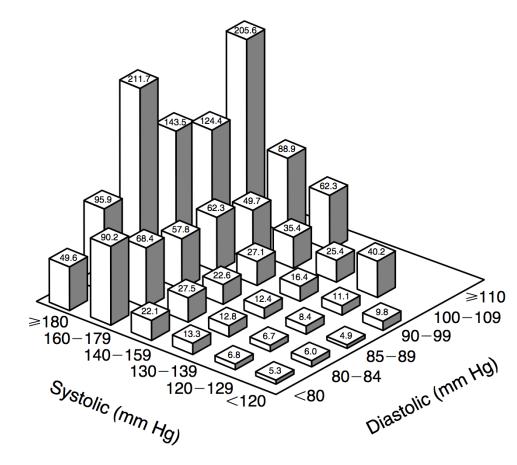
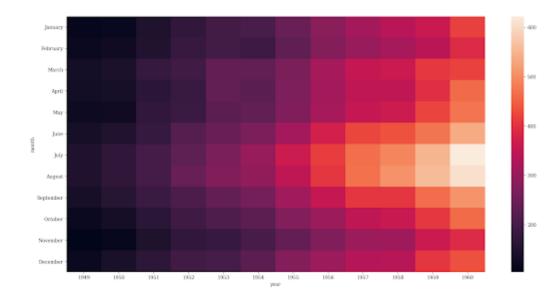
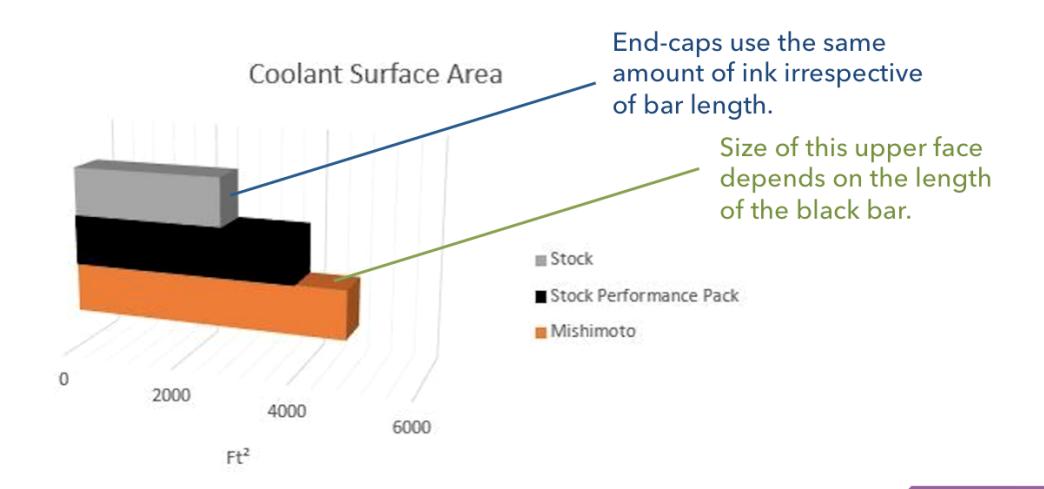


Figure 2. Age-Adjusted Rate of End-Stage Renal Disease Due to Any Cause per 100,000 Person-Years, According to Systolic and Diastolic Blood Pressure in 332,544 Men Screened for MRFIT.

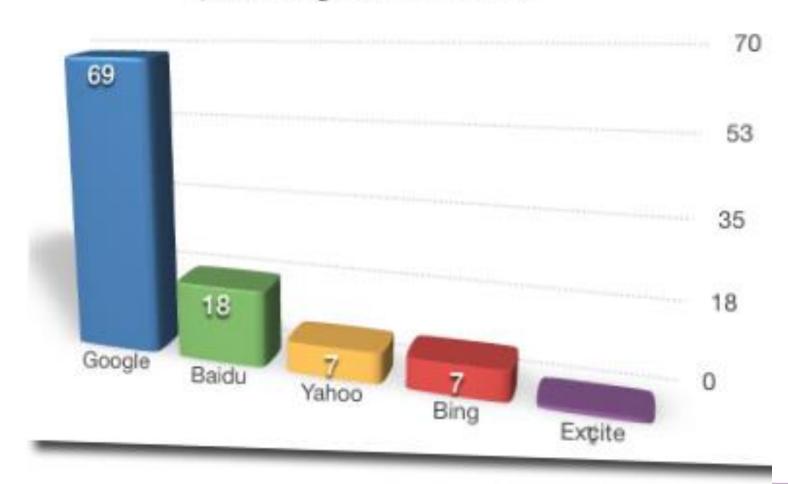


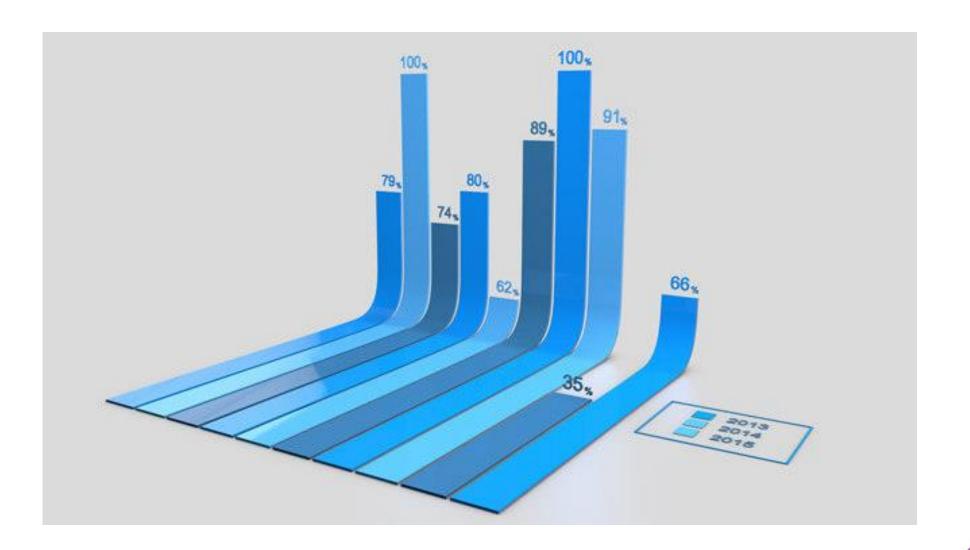




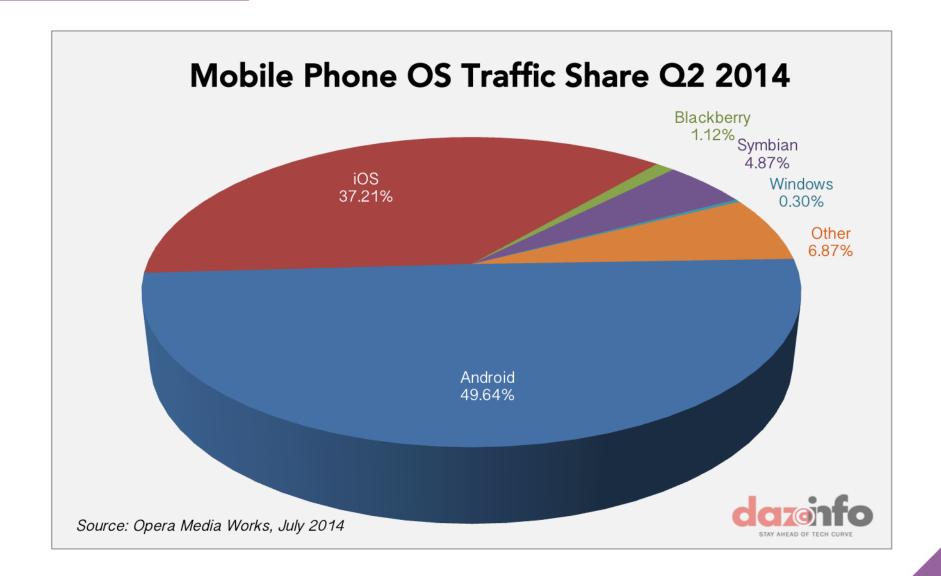


Search Engine Market Share











Is there a way to choose the "right" plot?

03 October 2023 31

ROBERT GORDON UNIVERSITY ABERDEEN

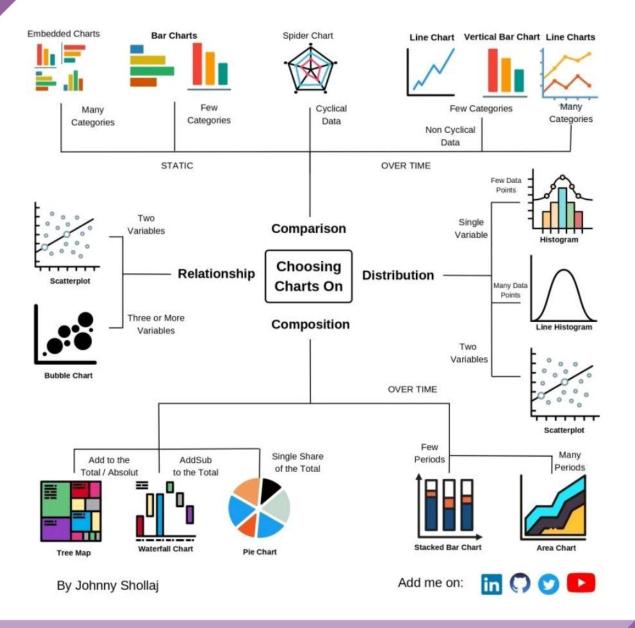
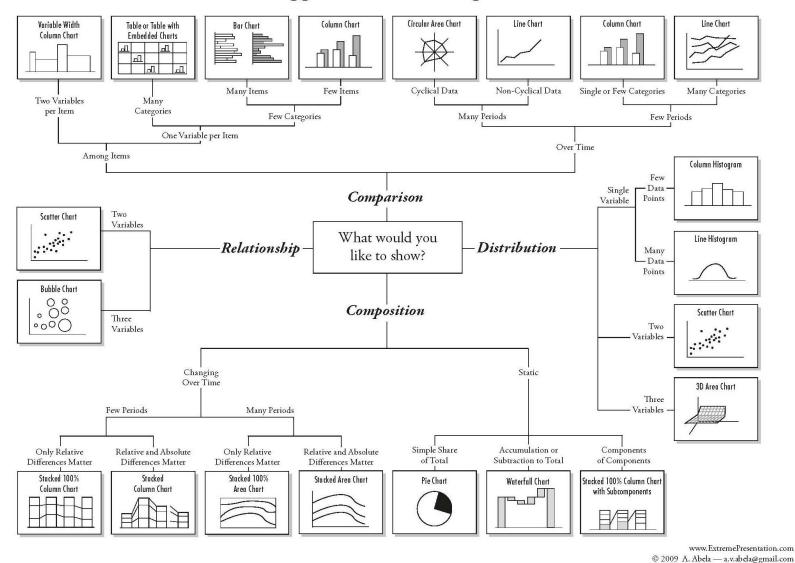




Chart Suggestions—A Thought-Starter





Deviation

Emphasise variations (+/-) from a fixed reference point. Typically the reference point is zero but it can also be a target or a long-term average. Can also be used to show sentiment (positive/neutral/negative).







Correlation

Show the relationship between two or more variables. Be mindful that, unless you tell them otherwise, many readers will assume the relationships you show them to be causal (i.e. one causes the







regories of data, less ective at showing fine ferences in amounts.

Ranking

Lollipops draw more attention to the dat value than standard bar/column and can

changing rankings across multiple dates. For large datasets, consider grouping lines using colour.



Barcode plot

Dot strip plot





















Distribution Change over Time

Show values in a dataset and how often they occur. The shape (or 'skew') of a distribution can be a memorable way of highlighting the lack of uniformity or equality in the data. Give emphasis to changing trends. These can be short (intra-day) movements or extended series traversing decades or centuries. Choosing the correct time period is important to provide suitable context for the reader.

showing the change or range (min/max) of data across multiple





































Magnitude

Show size comparisons. These can be relative (just being able to see larger/bigger) or absolute (need to see fine differences). Usually these show a 'Counted' number (for example, barrels, dollars or people) rather than a calculated rate or per cent.



























Part-to-whole





















Spatial















Visual vocabulary

Designing with data

There are so many ways to visualise data - how do we know which one to pick? Use the categories across the top to decide which data relationship is most important in your story, then look at the different types of chart within the category to form some initial ideas about what might work best. This list is not meant to be exhaustive, nor a wizard, but is a useful starting point for making informative and meaningful data visualisations.



Flow















If there's time...

Check some examples from last lab