# Calling Bullsh*t on Data Viz

CM4125 – Week 4

# Announcements

- Datacamp Academic Group

- (More) Coursework Clarification
  - Word count limits
  - Video from last week is embedded in Moodle
  - I will explain the clustering plot at the end of this lecture
  - I will attend more questions after the lab
    - Marking grid (i.e. what are you expected to do)

# 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).
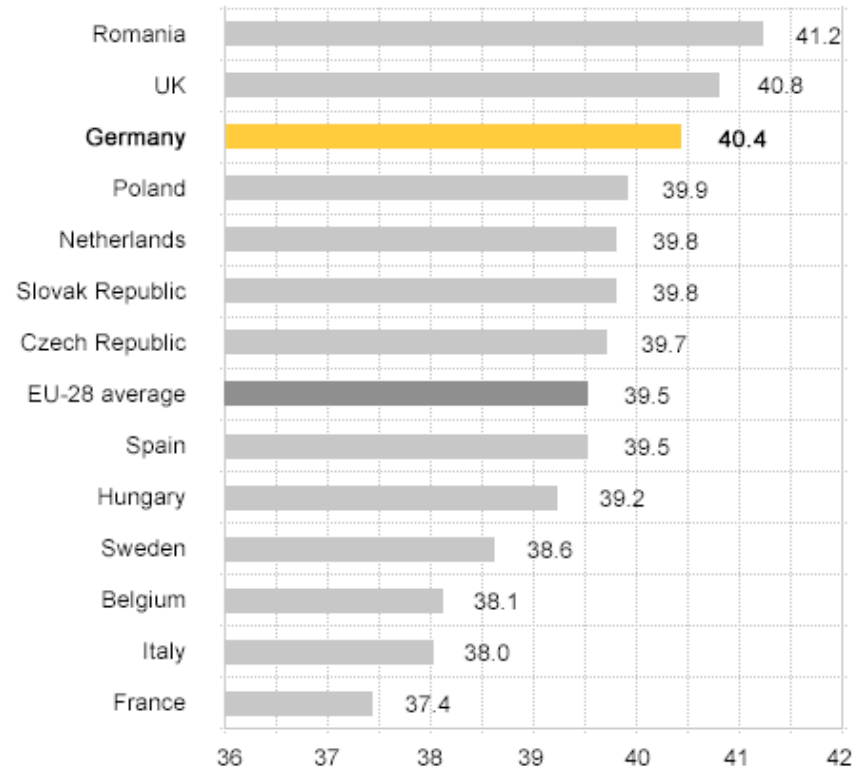
# Misleading axis
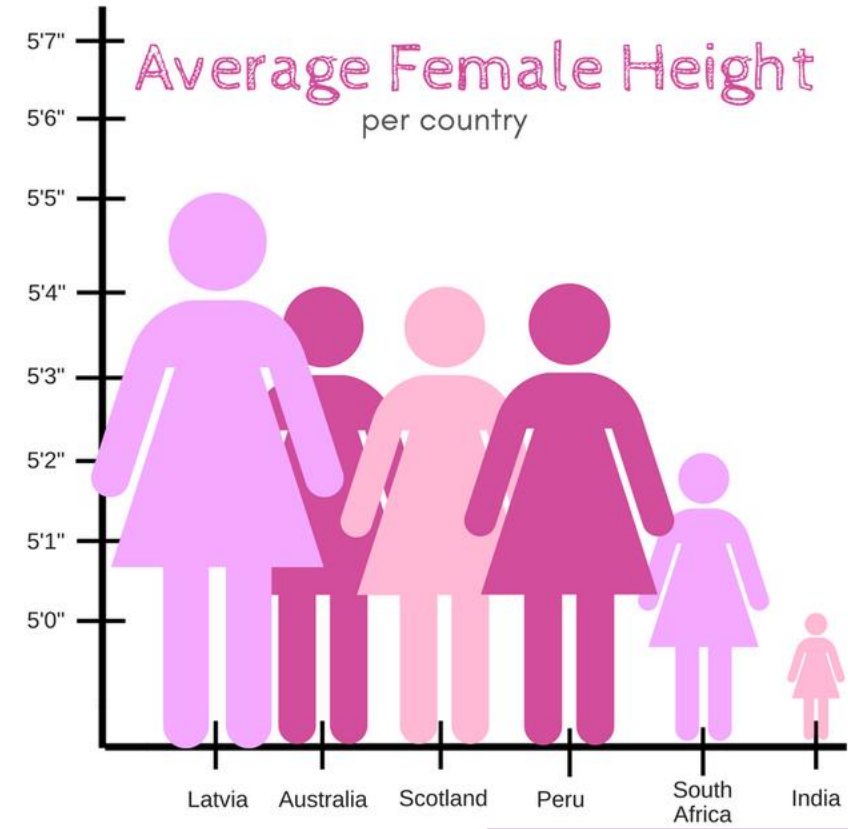
Charts can conceal or mislead if axes are set incorrectly

How to choose the range and scale of axes in a graph?

# 1. Bar chart axes **<u>should</u>** include zero



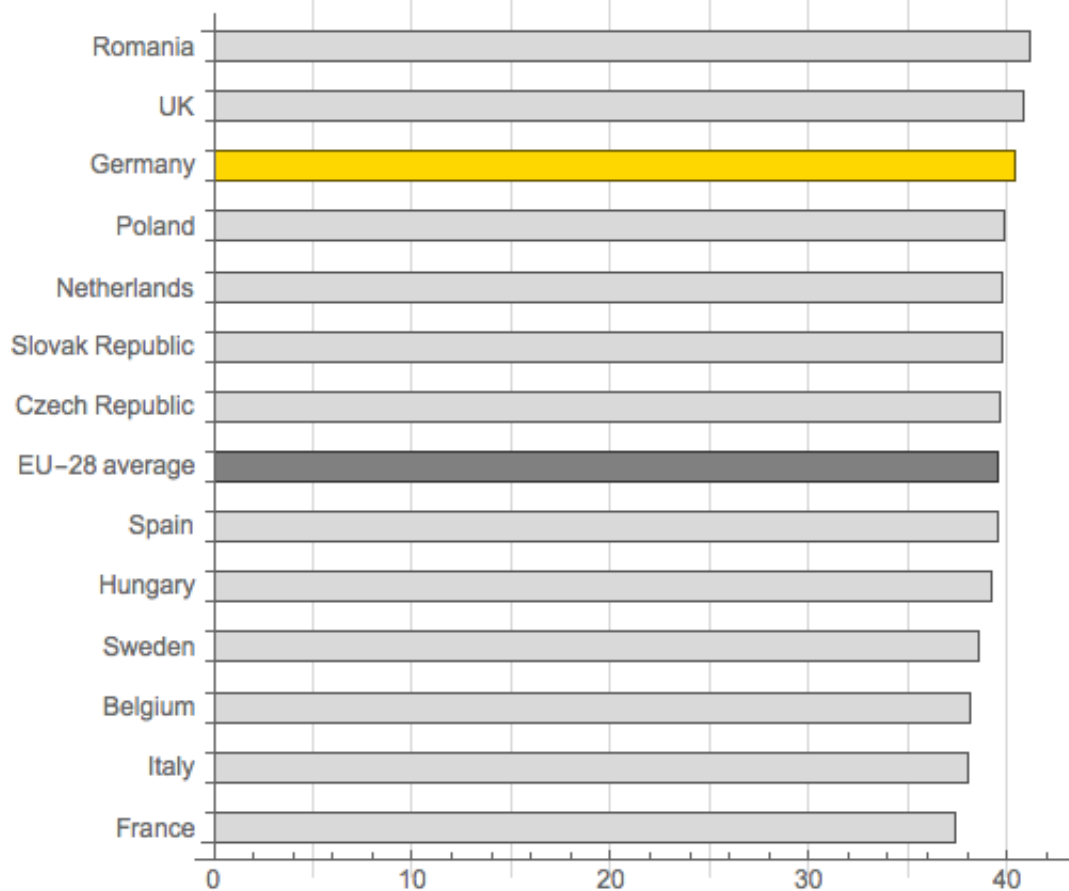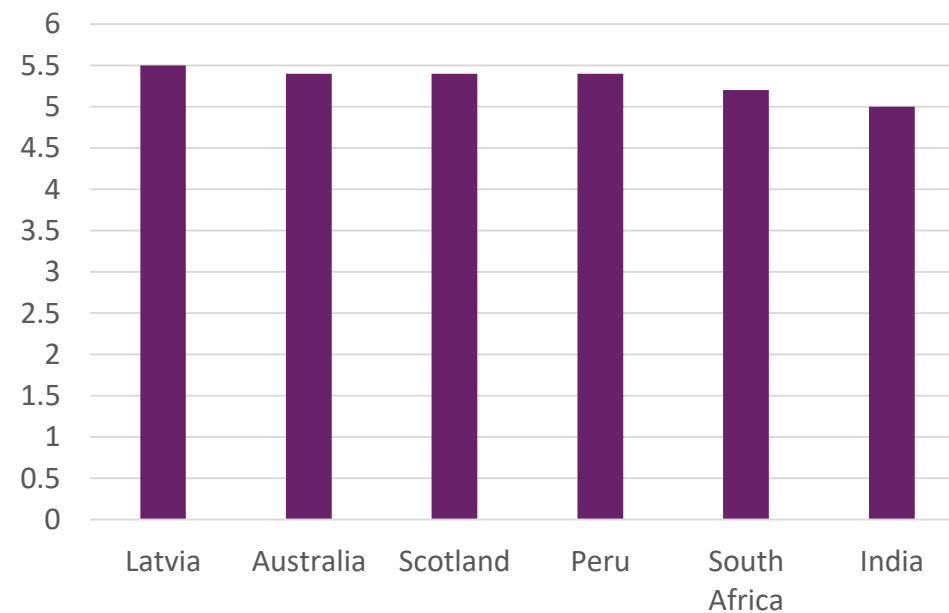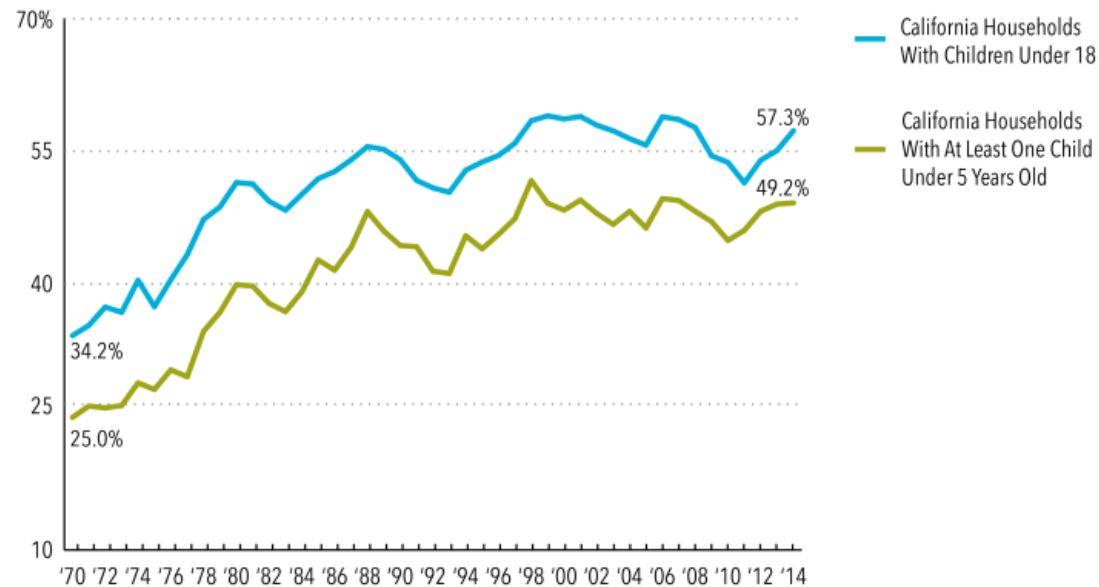Average number of actual weekly hours of work in main job, full-time employees, 2013

| Country | Hours |
|---|---|
| Romania | 41.2 |
| UK | 40.8 |
| Germany | 40.4 |
| Poland | 39.9 |
| Netherlands | 39.8 |
| Slovak Republic | 39.8 |
| Czech Republic | 39.7 |
| EU-28 average | 39.5 |
| Spain | 39.5 |
| Hungary | 39.2 |
| Sweden | 38.6 |
| Belgium | 38.1 |
| Italy | 38.0 |
| France | 37.4 |

Source: Eurofound 2014



Average Female Height per country

# 1. Bar chart axes **should** include zero



Average number of actual weekly hours of work in main job, full-time employees, 2013
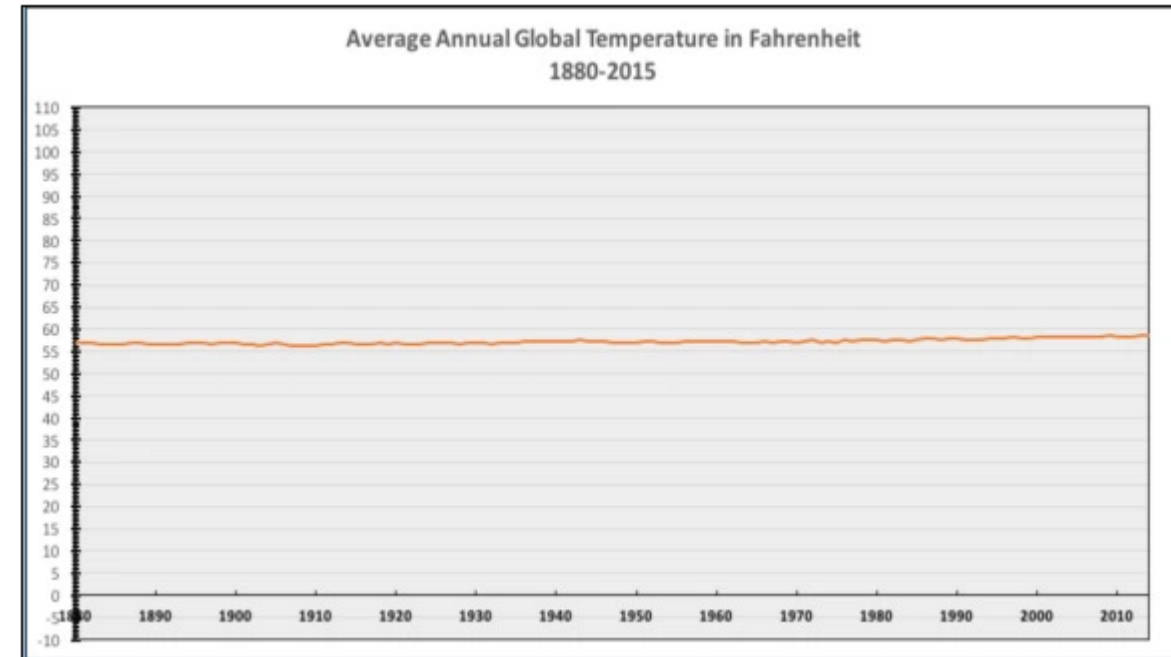


Average Female Height Per country

# 2. Line graph axes **need not** include zero
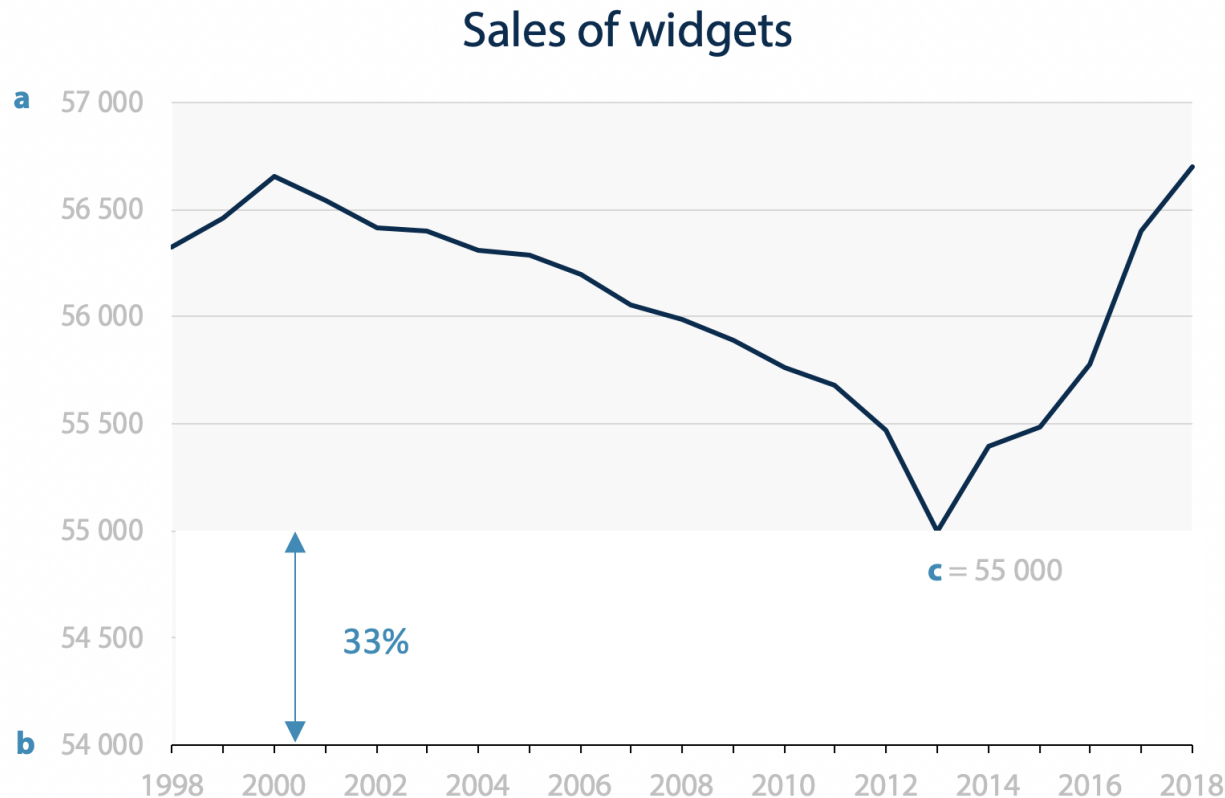
# Well…

# Solution

Empty space = $\dfrac{\text{Minimum data value - Minimum scale value}}{\text{Maximum scale value - Minimum scale value}}$
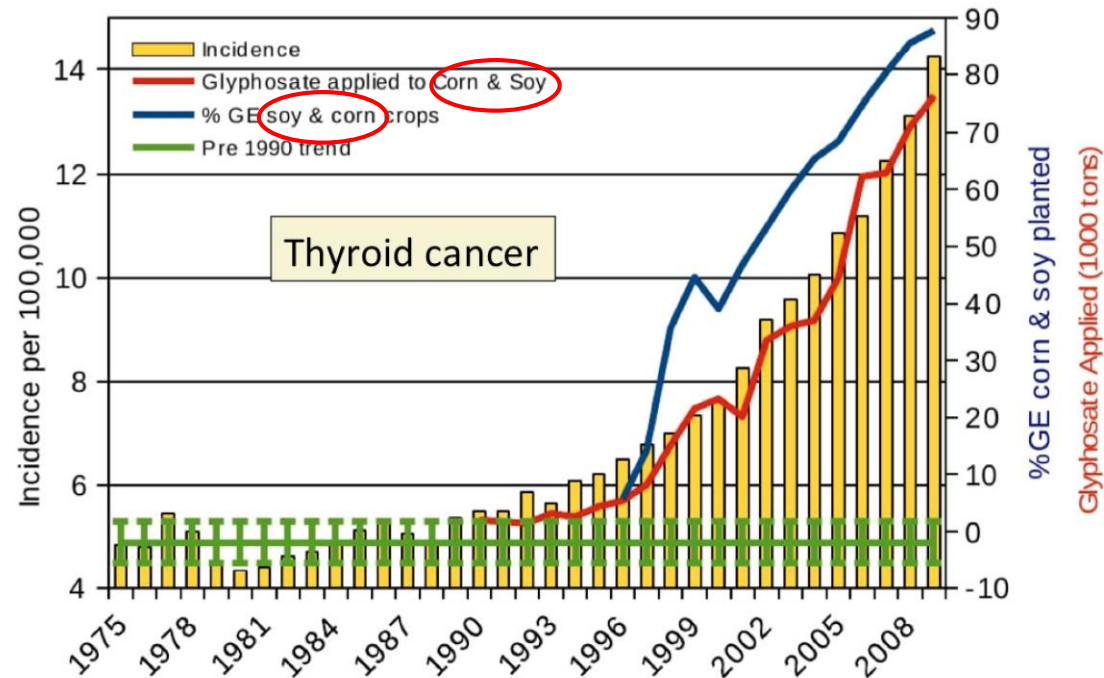


Sales of widgets

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

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

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

**Gagnon, F. (2018). A golden ratio for line charts with truncated y-axis**

# 3. **No** multiple axes on a single graph
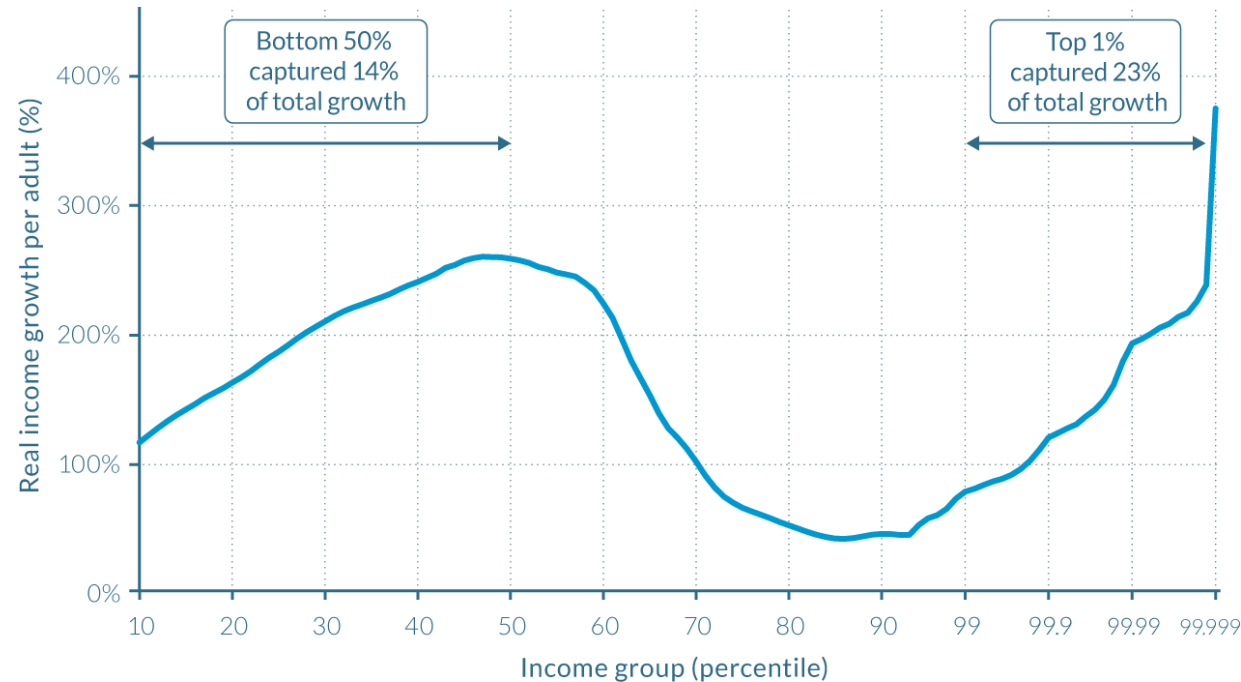


Thyroid Cancer Incidence Rate (age adjusted)

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

Incidence
Glyphosate applied to Corn & Soy
% GE soy & corn crops
Pre 1990 trend

Thyroid cancer

Incidence per 100,000

%GE corn & soy planted    Glyphosate Applied (1000 tons)

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

**How to solve this?**

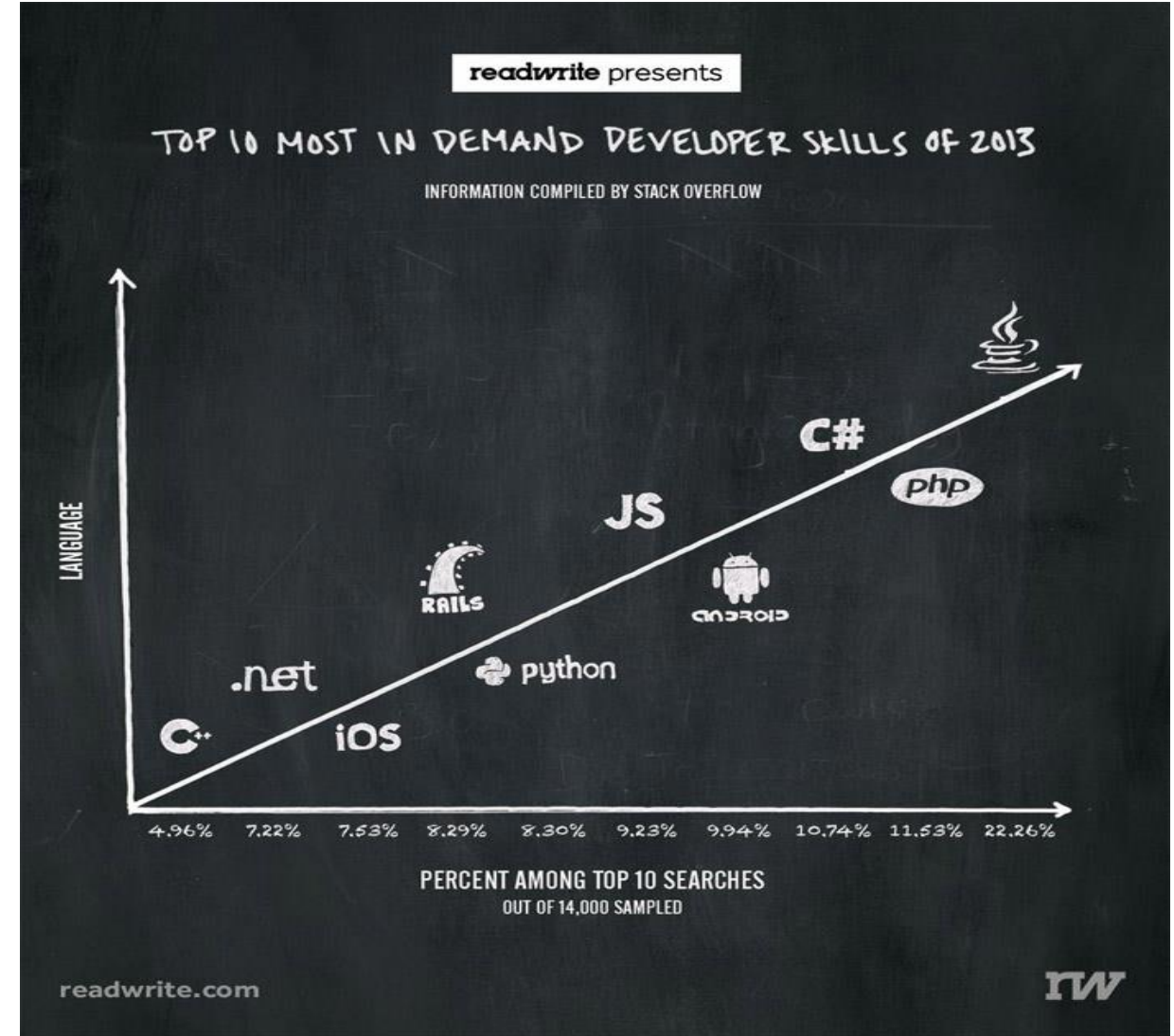# 4. An axis should **<u>not</u>** change scales midstream



**Use log scales wisely!**

# 5. An axis should have *something* on it

# Like this

Tools such as [WebPlotDigitizer](WebPlotDigitizer) can help you get the missing info
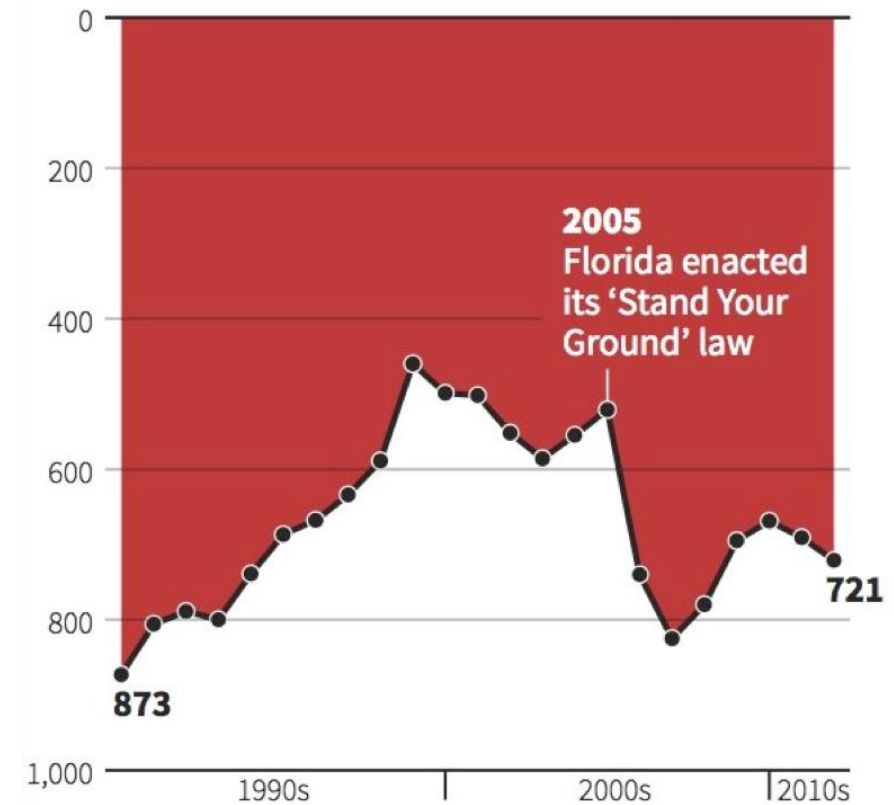
# 6. **Don't** invert the axis(?)

[The author defended it](#)!
What's your stance?



**Gun deaths in Florida**

Number of murders committed using firearms

2005
Florida enacted
its 'Stand Your
Ground' law

873

721

1990s    2000s    2010s

Source: Florida Department of Law Enforcement

C. Chan 16/02/2014                    REUTERS
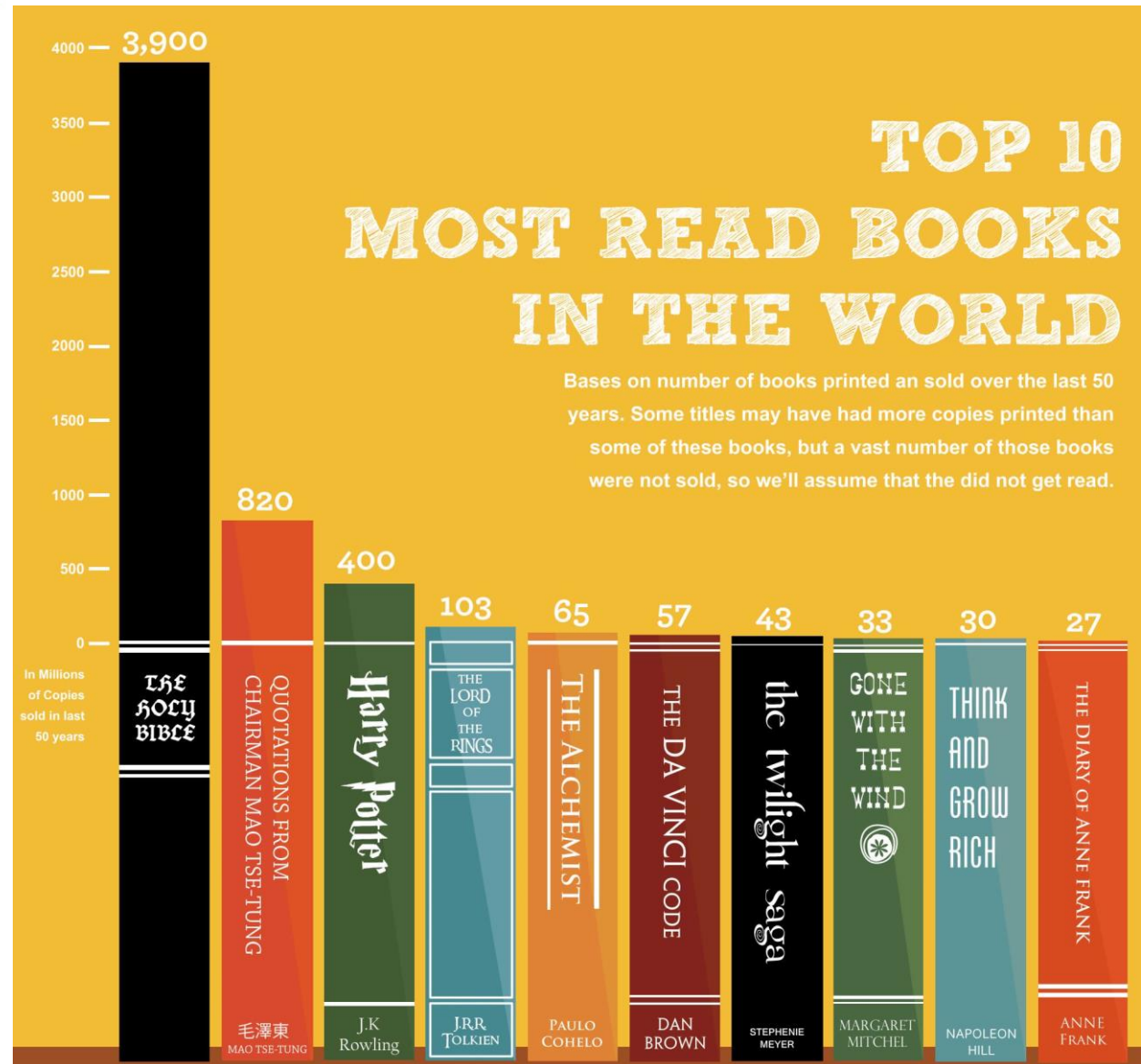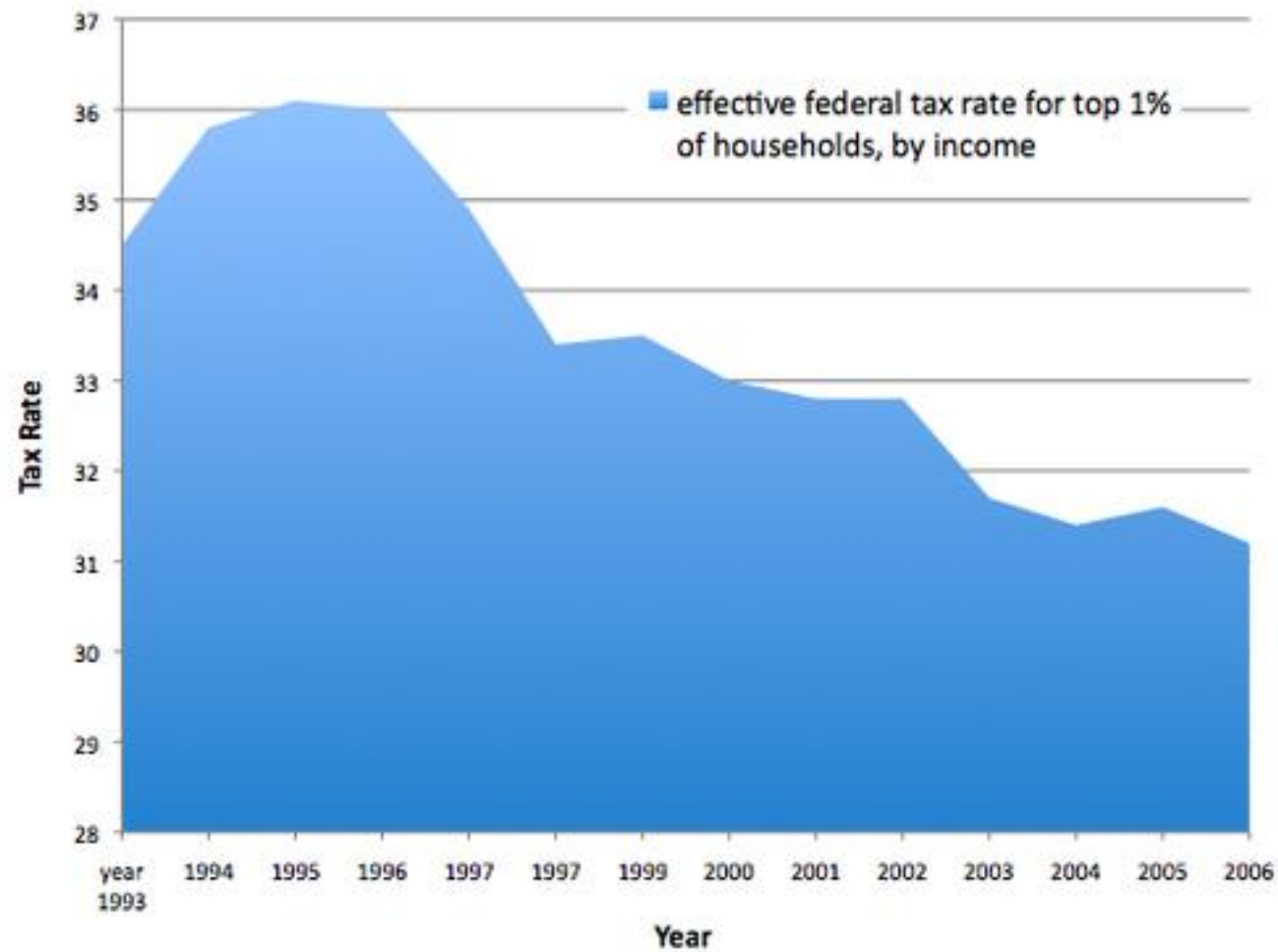
# Proportional Ink

*"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
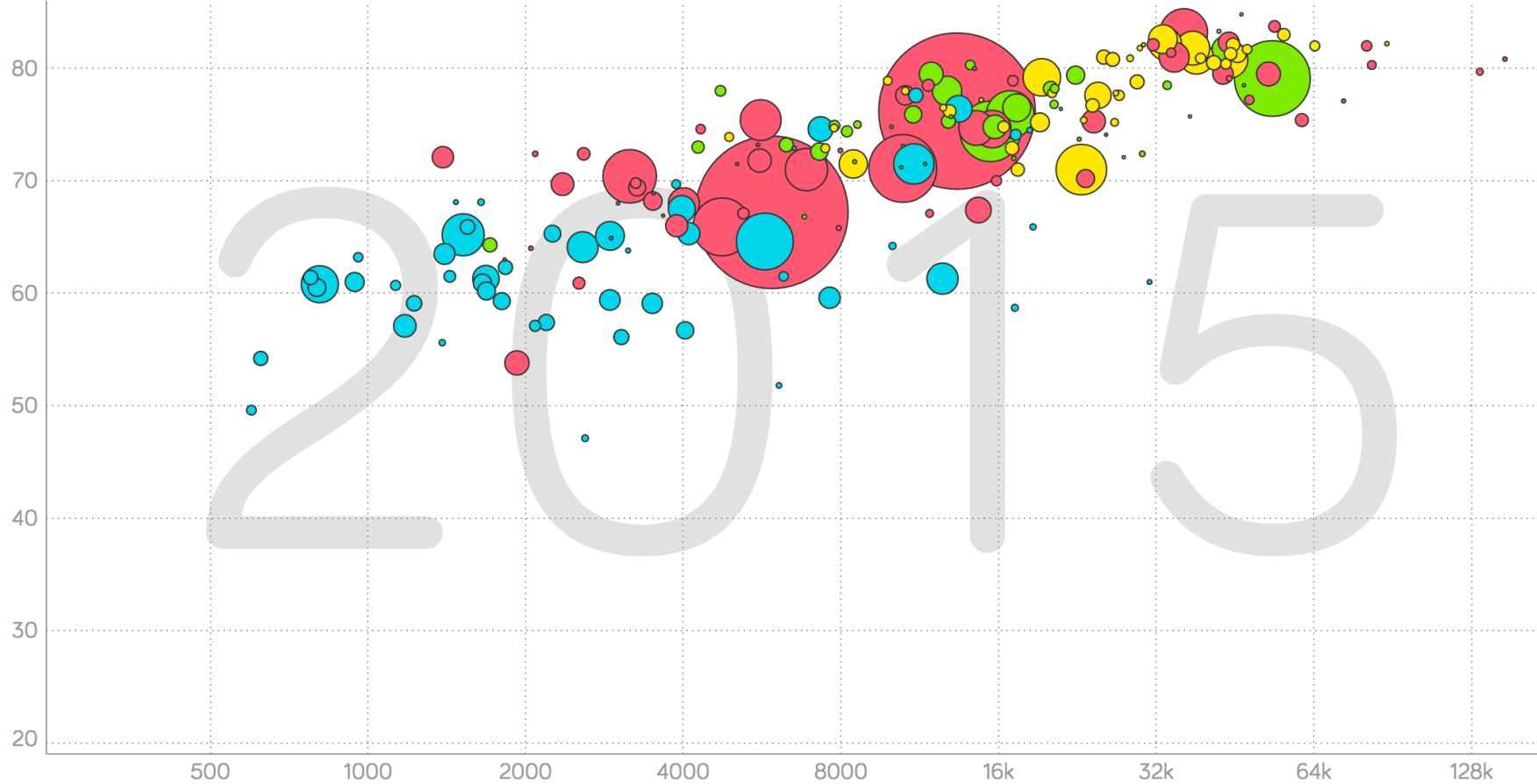
Life expectancy, years

80

70

60

50

40

30

20

500    1000    2000    4000    8000    16k    32k    64k    128k

Income per person, GDP/capita in $/year adjusted for inflation & prices
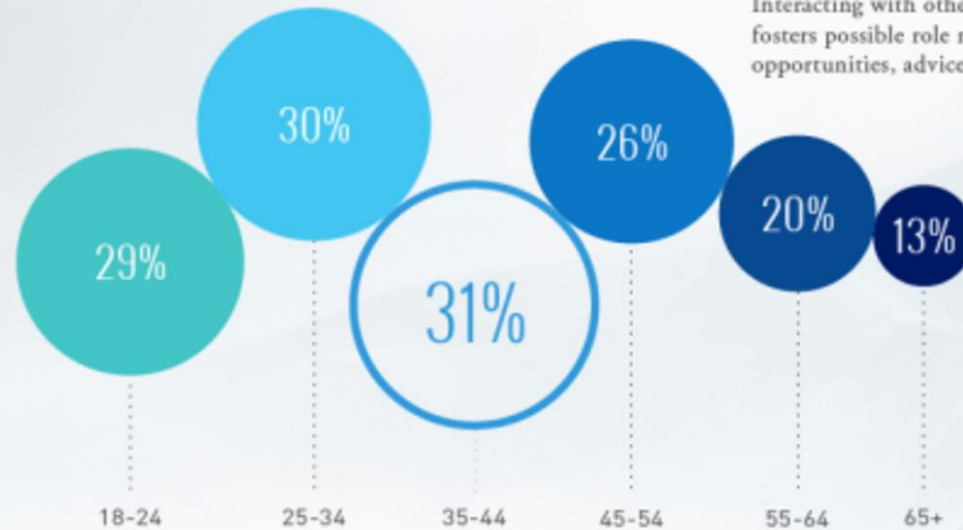
**Should we violate the principle?**

**Radius ≠ Area**
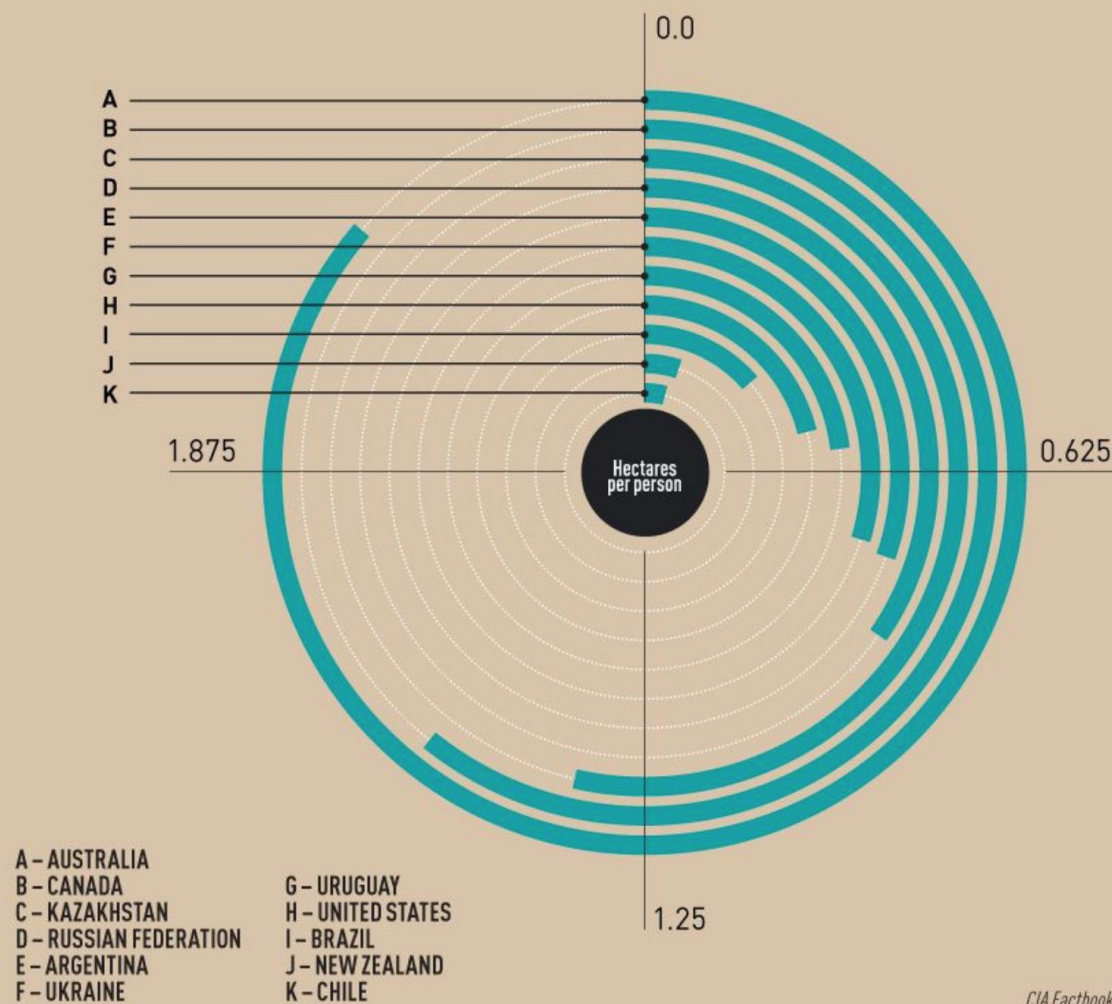


It's Who You Know

By ages 35–44, people are likely to have the strongest network of fellow entrepreneurs, but these **personal social networks begin to decline over time.**
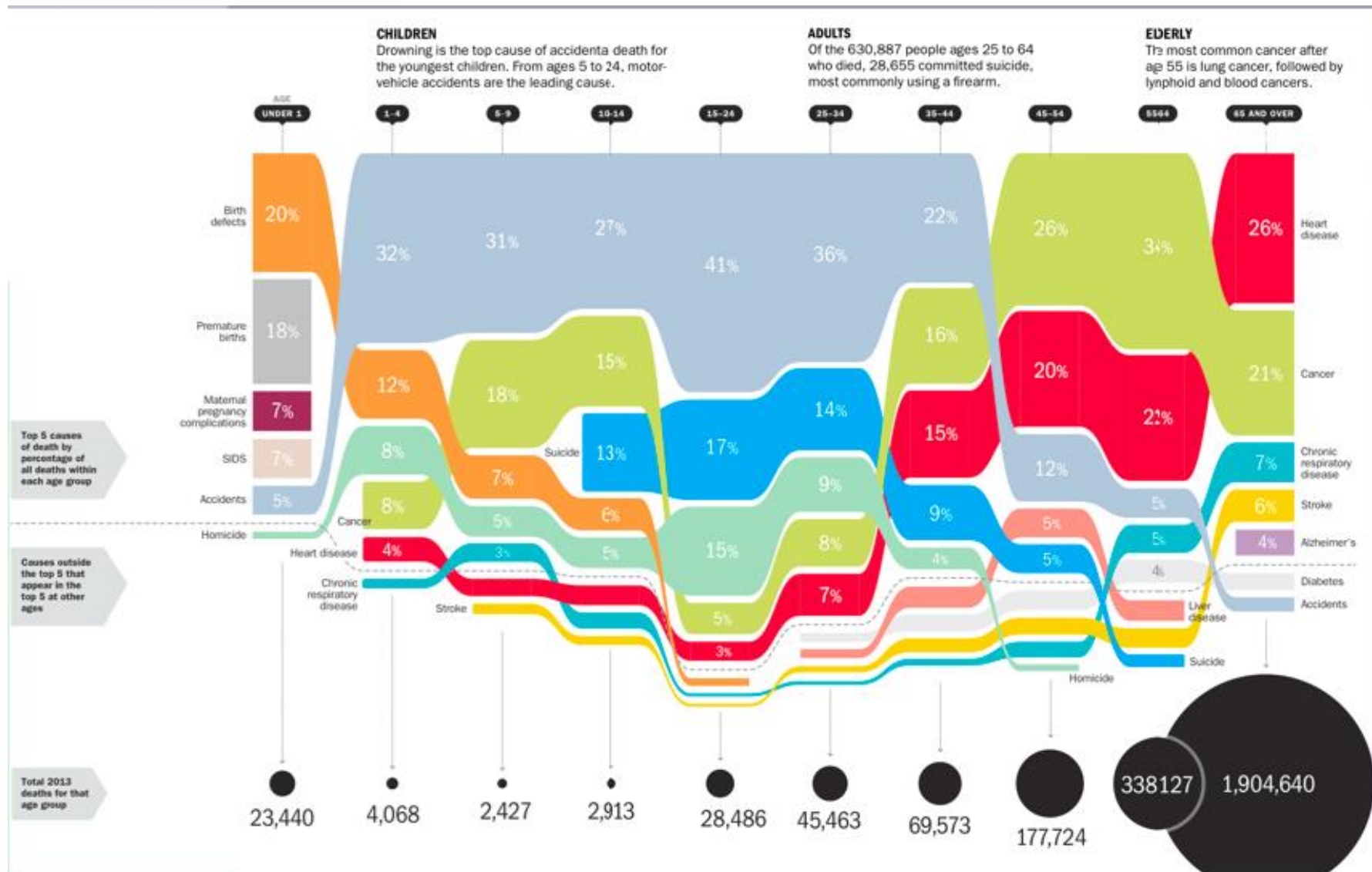
Interacting with other entrepreneurs fosters possible role models, networking opportunities, advice, and encouragement.

| 18-24 | 25-34 | 35-44 | 45-54 | 55-64 | 65+ |
| --- | --- | --- | --- | --- | --- |
| 29% | 30% | 31% | 26% | 20% | 13% |

ARABLE LAND PER CAPITA
(HECTARES PER PERSON)

A – AUSTRALIA
B – CANADA
C – KAZAKHSTAN
D – RUSSIAN FEDERATION
E – ARGENTINA
F – UKRAINE
G – URUGUAY
H – UNITED STATES
I – BRAZIL
J – NEW ZEALAND
K – CHILE

Source: CIA Factbook 2012 and FAO 2011

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.
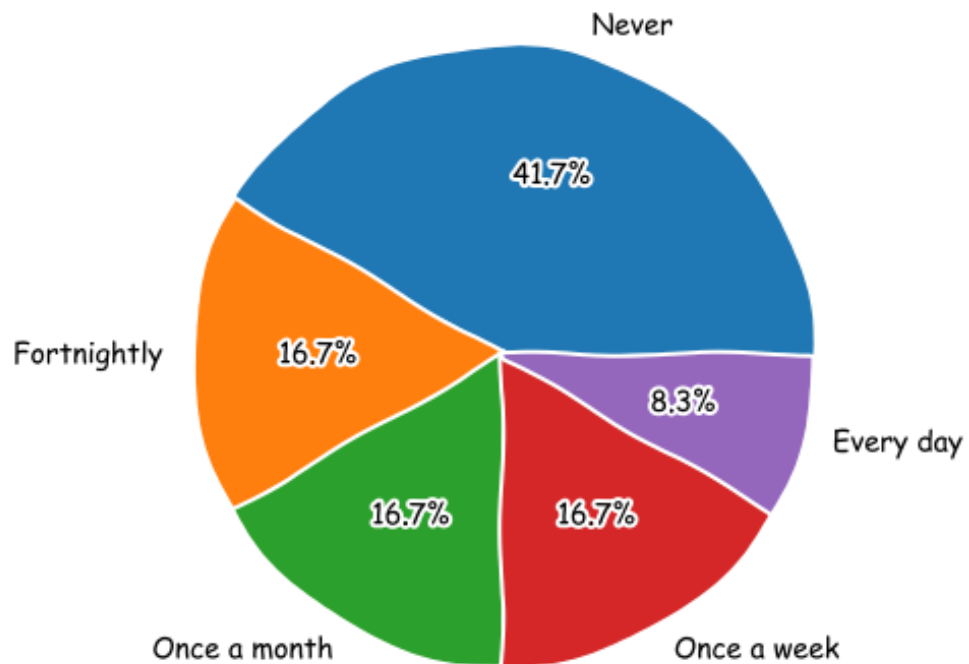
Along horizontal slices, ink is not proportional because total deaths differ widely by age group. Far more people 65 and older died of heart disease (red) than children age1-4 die of accidents (blue-gray), but the latter takes more ink because it represents a larger percentage of the (relatively few) total deaths at that age.
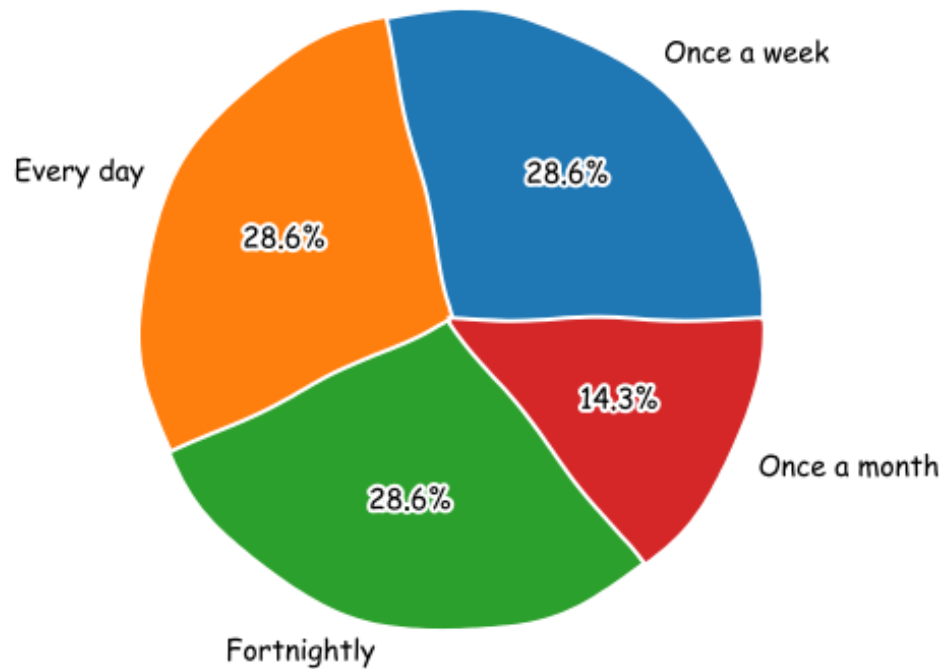
How often do you use a dating app during a month?

**Male** — Never 41.7%, Fortnightly 16.7%, Once a month 16.7%, Once a week 16.7%, Every day 8.3%

**Female** — Once a week 28.6%, Every day 28.6%, Fortnightly 28.6%, Once a month 14.3%
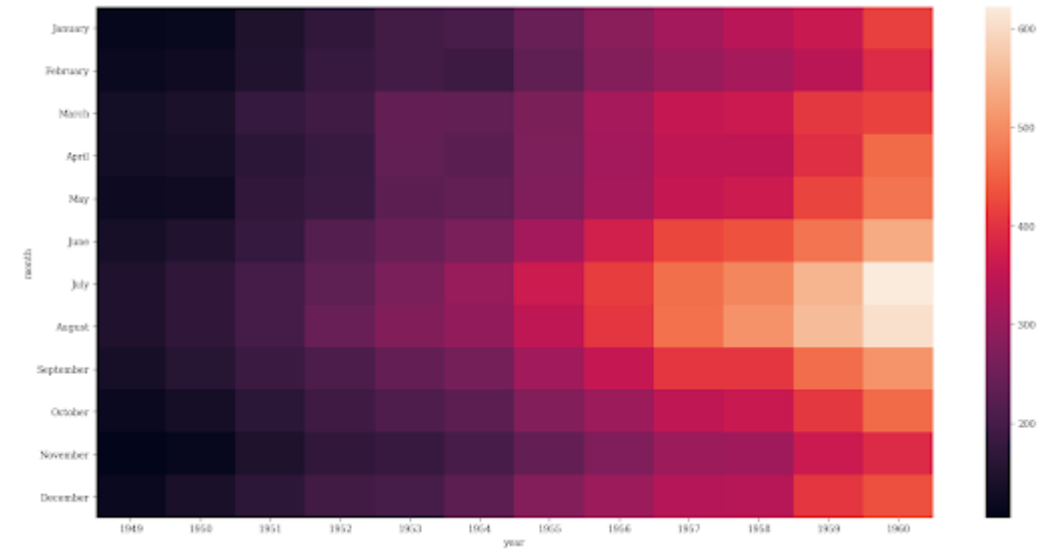
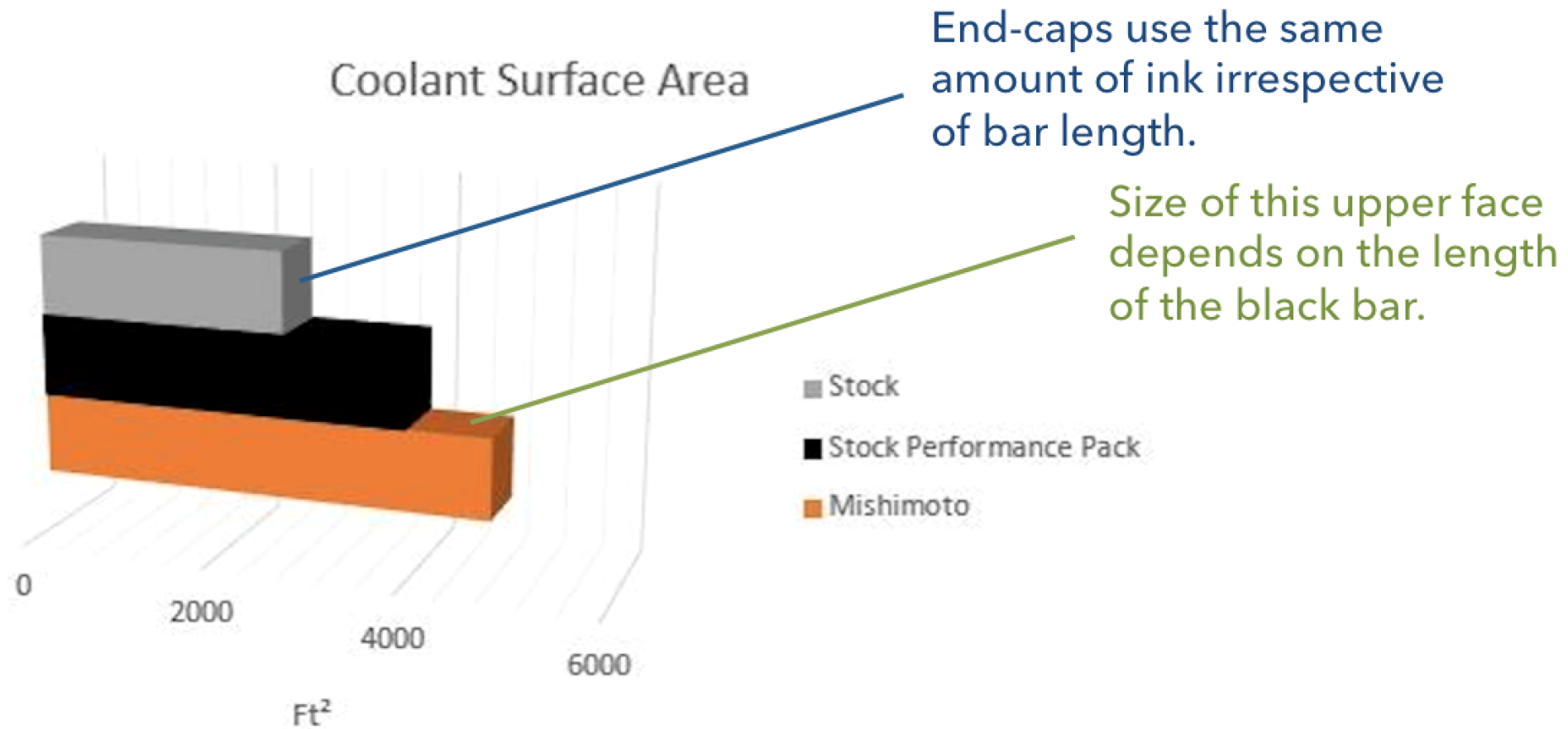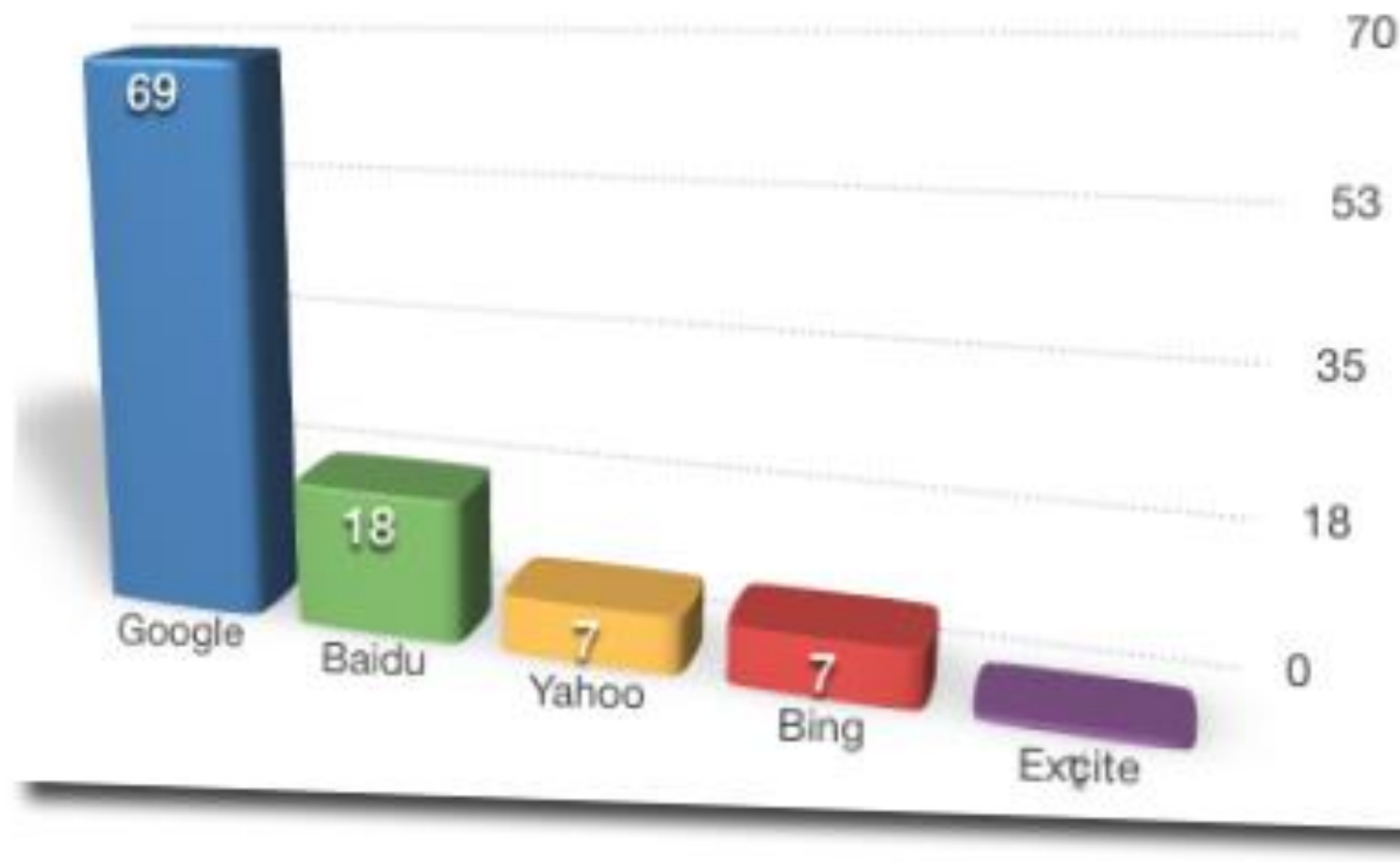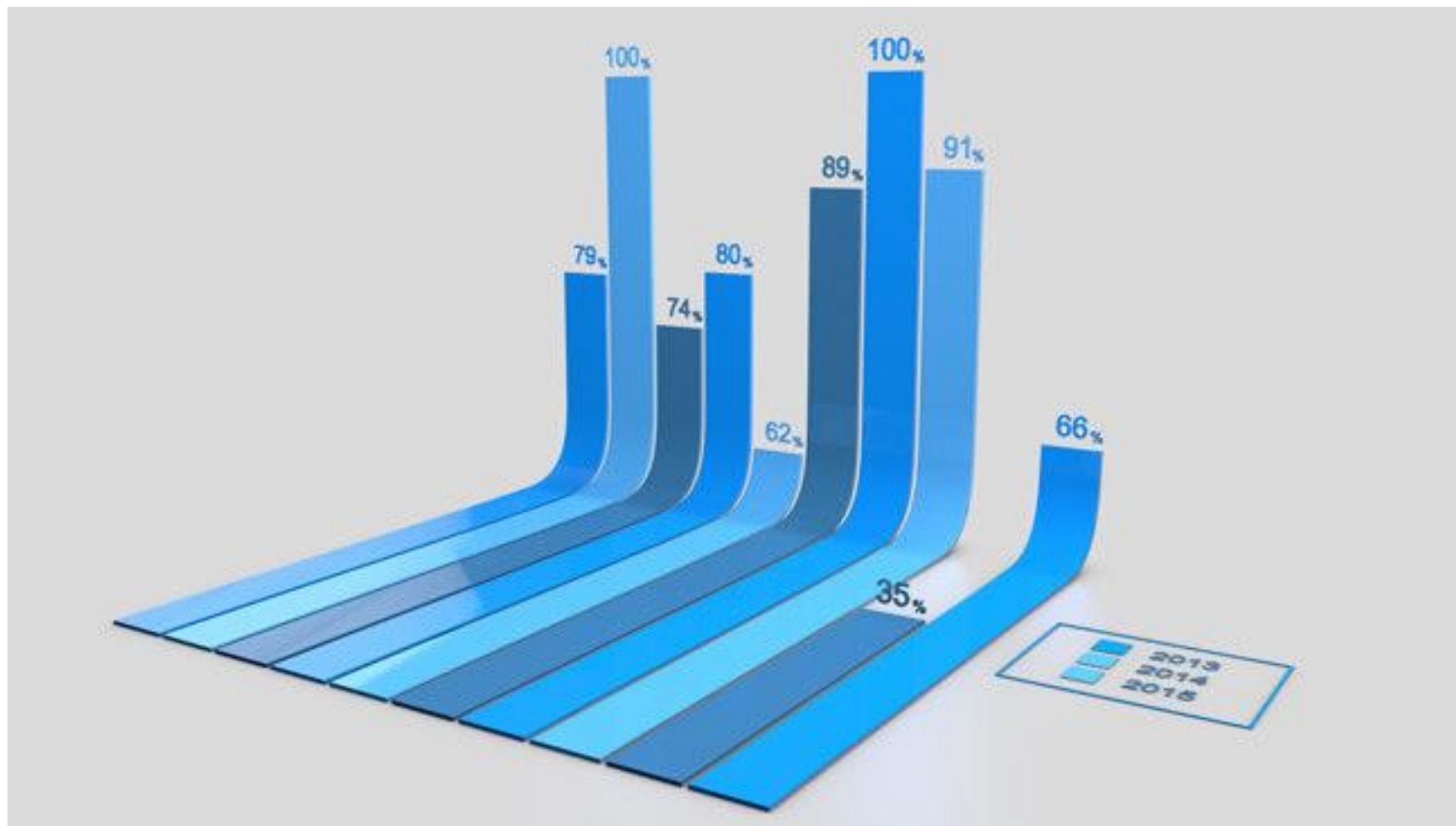**What is wrong in this one?**

# 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.

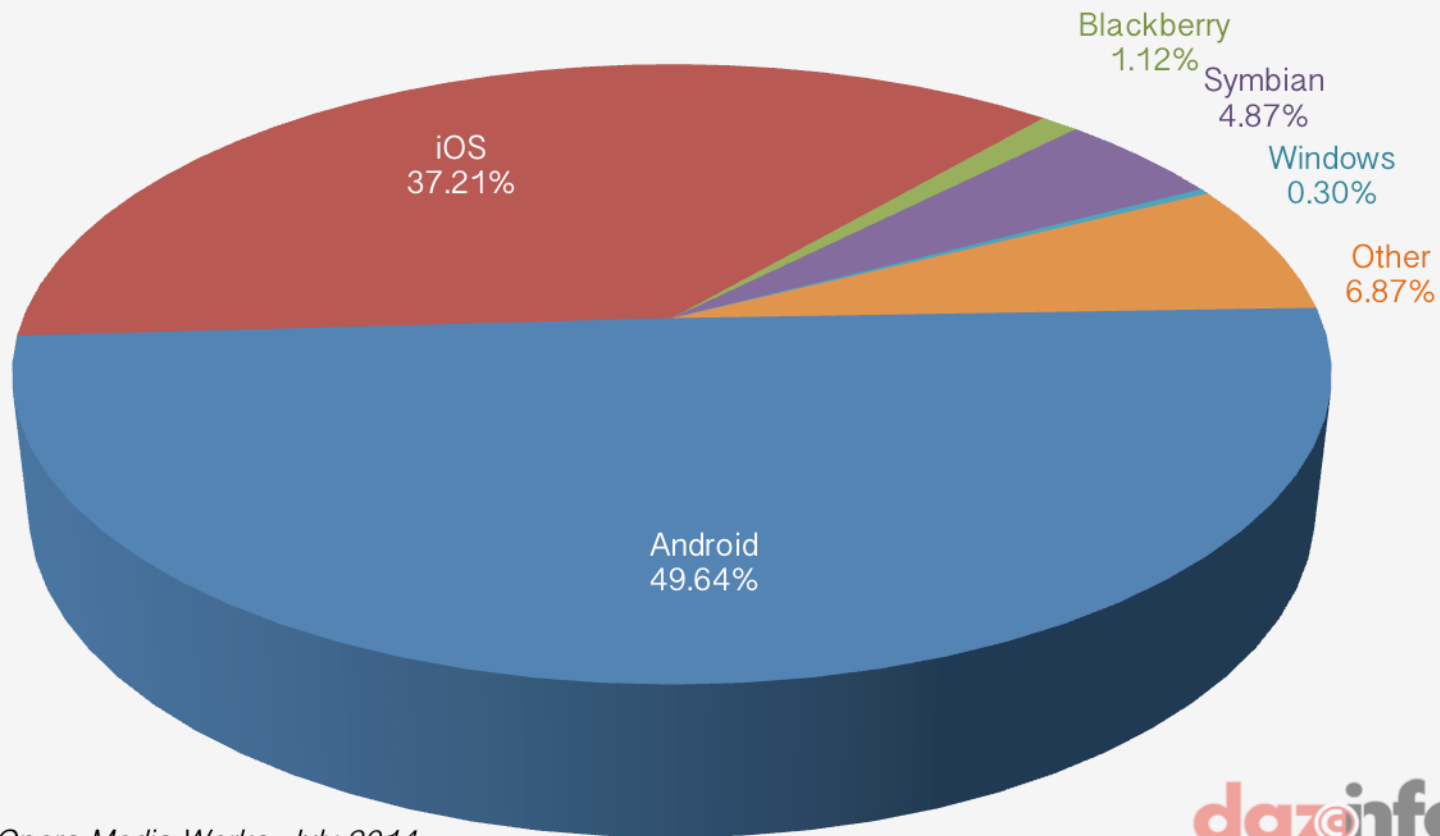Coolant Surface Area

End-caps use the same amount of ink irrespective of bar length.

Size of this upper face depends on the length of the black bar.

Stock
Stock Performance Pack
Mishimoto

0   2000   4000   6000

Ft²

Search Engine Market Share

Mobile Phone OS Traffic Share Q2 2014

Blackberry 1.12%
Symbian 4.87%
iOS 37.21%
Windows 0.30%
Other 6.87%
Android 49.64%

Source: Opera Media Works, July 2014

# The most hated visualisation of all time

# Clustering

- The **IDEA** of this plot was to show stakeholders (i.e. the people from the listed restaurants looking at this dashboard) that users/profiles can be grouped in NEW WAYS

- Therefore, we process the user/profile data with a method called *k-means*

- In short, this method analyses the **FEATURES** of all users/profiles and groups them into clusters

# Step 1: Select the Features

# Step 2: Visualise the clustering



-In this example, the algorithm has found 12 groups using the 6 features selected in the previous slide

-Each dot represents a user/<u>profile</u>

-The problem is that we selected 6 features, but since each feature corresponds to one axis, it is impossible to do 6-dimension plots

-Therefore, we use a method called PCA to reduce the dimensionality from 6 to 2

# Step 3: Understand each cluster

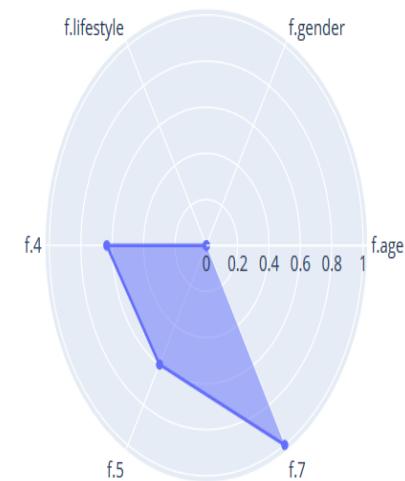-For each group found, you have a tab at the top

-This tab takes you to a radial chart, in which the **IDEA** is to explain which are the features that influence the most to the creation of a certain cluster

-In this case, we can see that for group 0 it was feature 7 (milk allergy) the one that most defines this group

-This can mean that either it is highly likely that most people in this group are either <u>all allergic</u> or <u>not allergic</u> to milk

-However, at this stage we <u>don't know which is it!</u>

# As for the marking…

What am I expecting? How will I grade?

# One very simple recommendation…