

# Week 2: Graphics and Visualization

## MATH-517 Statistical Computation and Visualization

Linda Mhalla

September 29th 2023

# Graphics

*“The simple graph has brought more information to the data analyst’s mind than any other device.” – John W. Tukey*

*“The greatest value of a picture is when it forces us to notice what we never expected to see.” – John W. Tukey*

**One can think of graphics (and also models, for that matter) as a low-dimensional representation for data**

# Principle of Visualization

The most common purposes of a visualization is

- presentation
  - result communication
  - decision making
- data insight
  - large data
  - detect patterns
  - find strange observations

⇒ a good choice of axes, axis limits, labels and symbols can facilitate substantially the extraction of information from the data

Datasets used for illustration are described in details in the [Lecture Notes 2](#)

# Anscombe's Quartet

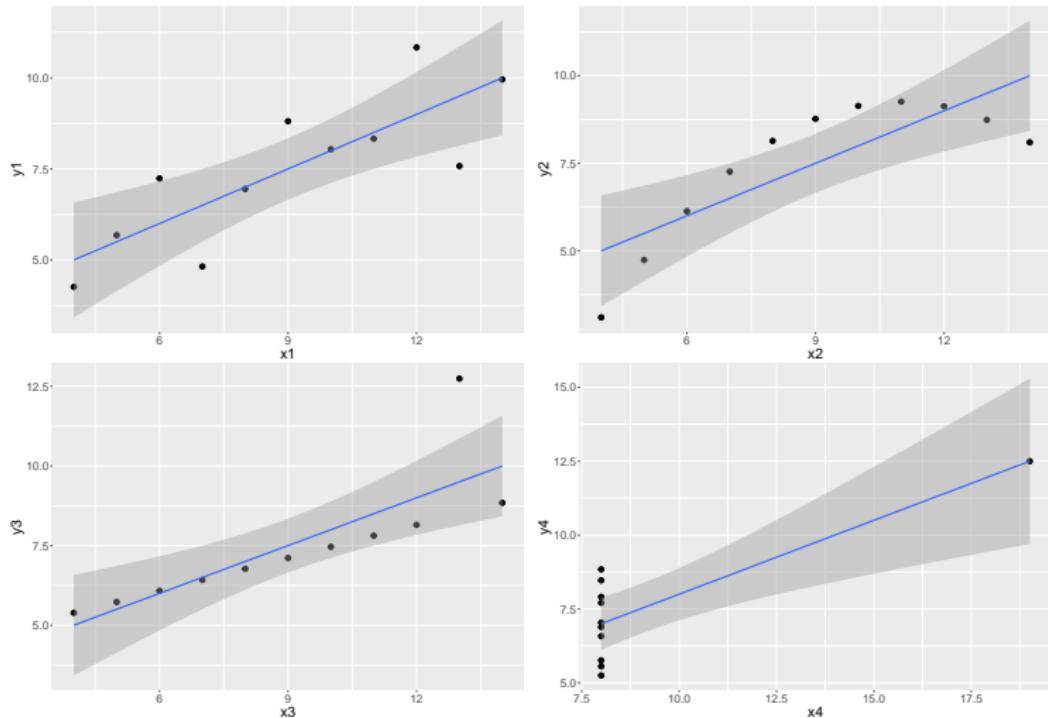
Four data sets with very similar descriptive statistics, each with

- one response variable  $y$
- one regressor  $x$

```
##      (Intercept)          x  R-squared
## lm1     3.000091 0.5000909 0.6665425
## lm2     3.000909 0.5000000 0.6662420
## lm3     3.002455 0.4997273 0.6663240
## lm4     3.001727 0.4999091 0.6667073
```

⇒ how to fool the linear regression model ...

# Anscombe's Quartet



⇒ qualitatively different structures

# Human Height



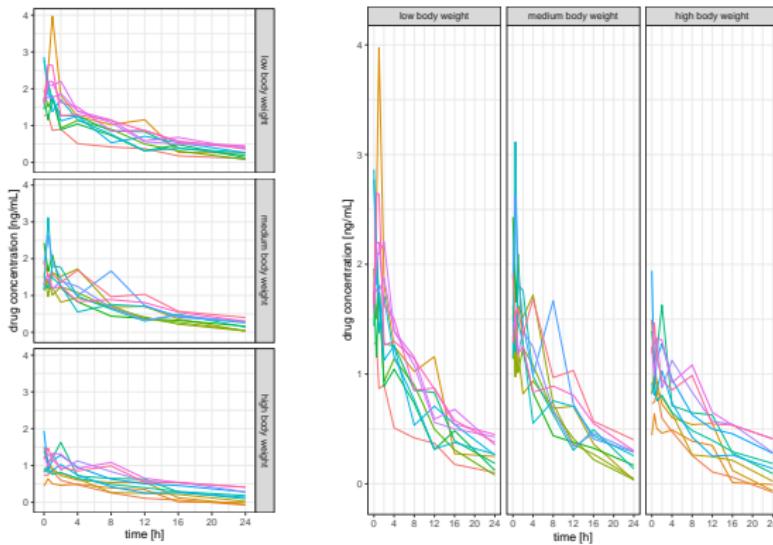
## Section 1

### Elements of Visualization

# Elements of Visualization: Layout

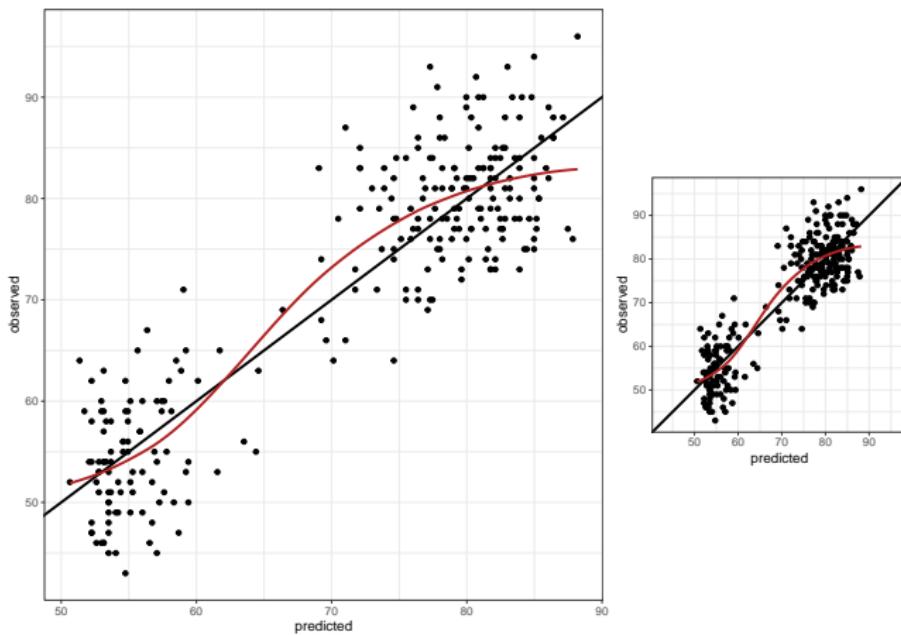
The arrangement of multiple panels is used for efficient comparison. To compare data on

- the y-axis, use a single y-axis with all panels aligned horizontally
- the x-axis, stack panels sharing a single x-axis



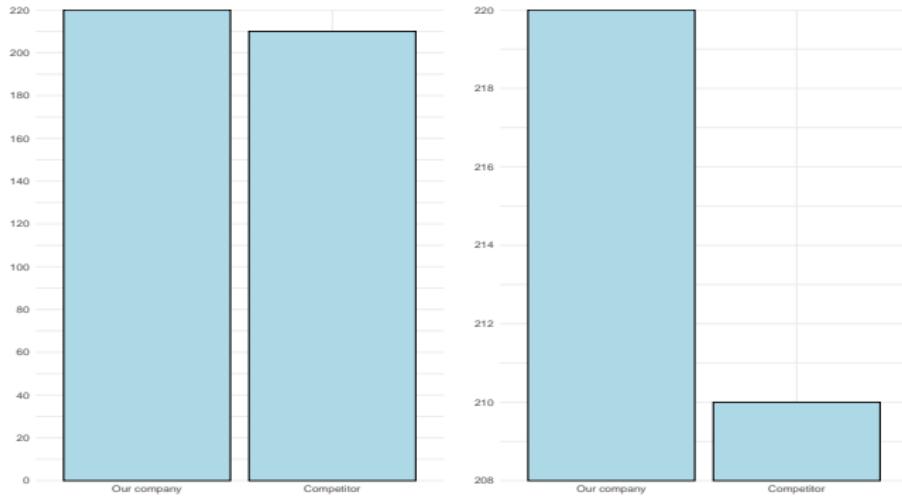
# Elements of Visualization: Aspect Ratio

- If measurements on both axes reflect the same quantity (e.g. before vs after treatment, observed vs modelled), a square figure (1:1 aspect ratio) could avoid visual bias



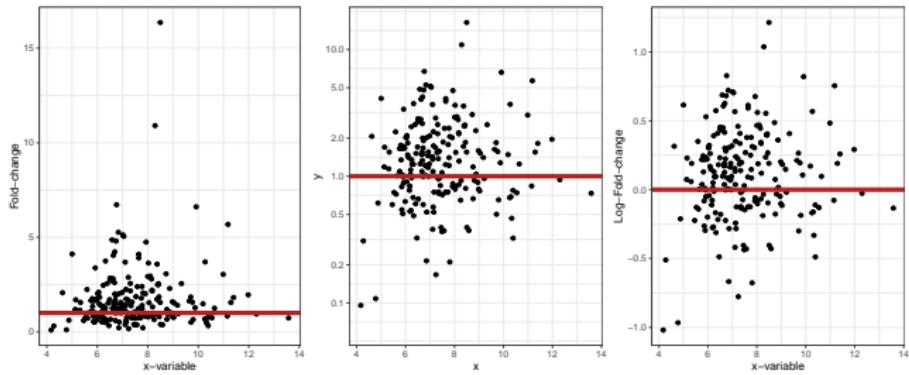
# Elements of Visualization: Axes

- If values are positive, axes should start at zero (unless there is a good reason for different choice) and not contain negative values
- If values on both axes are comparable, axes' limits should be the same (to make distances consistent and comparable)



# Elements of Visualization: Axes

- If values are ratios or relative changes, axes should be logarithmic and symmetric around the point of no change (e.g., 0)



# Elements of Visualization: Colours and Shapes

Chart features (ggplot's arguments) for points (and similarly for lines)

- colour: help identify different groups (do not use for pure decoration)

Tufte pointed out that because “they do have a natural visual hierarchy, varying shades of gray show varying quantities better than color”, and “the shades of gray provide an easily comprehended order to the data measures. This is the key”.

# Elements of Visualization: Colours and Shapes

- colour
- shape: if data are ordered, use ordered symbols (number of vertices)

`pch = _`

1 ○ 6 ▽ 11 ■ 16 ● 21 ○  
2 △ 7 □ 12 ■ 17 ▲ 22 □  
3 + 8 \* 13 ■ 18 ◆ 23 ◇  
4 × 9 ♦ 14 □ 19 ● 24 △  
5 ◇ 10 ♦ 15 ■ 20 ● 25 ▽

- size
- alpha (opacity/transparency): can be used to clarify plots with many points

These features can be used

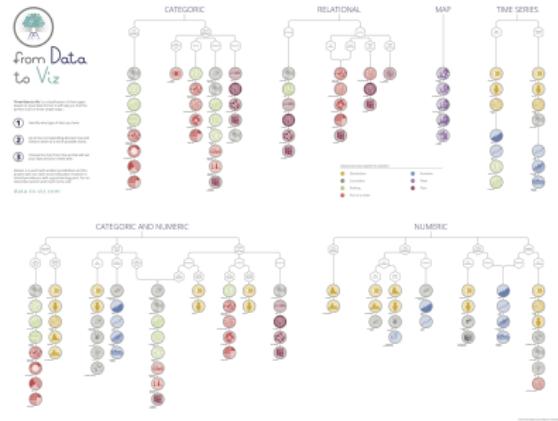
- to include additional information (or dimensions, i.e., to include additional variables) in a scatterplot
- to combat overplotting

## Section 2

### Visualization types

# What type of chart?

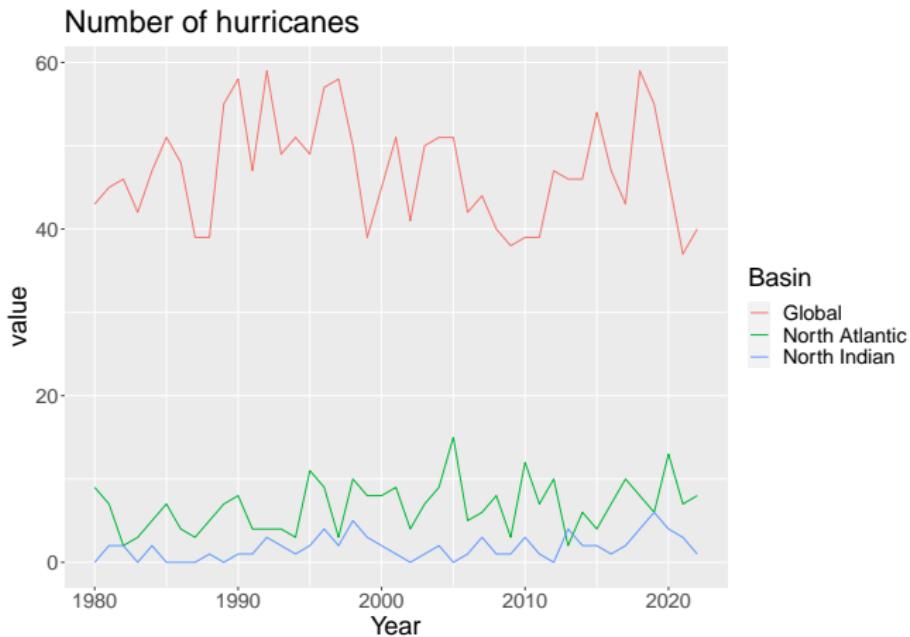
- Be clear about the intended purpose of a data visualisation (compare groups, highlight correlation, show evolution over time, ...)
- Choose the right chart according to the type of data and purpose of visualization: [From Data to Viz](#)



- Further general guidance can be found in [Robbins' \*Creating Better Graphs\*](#)

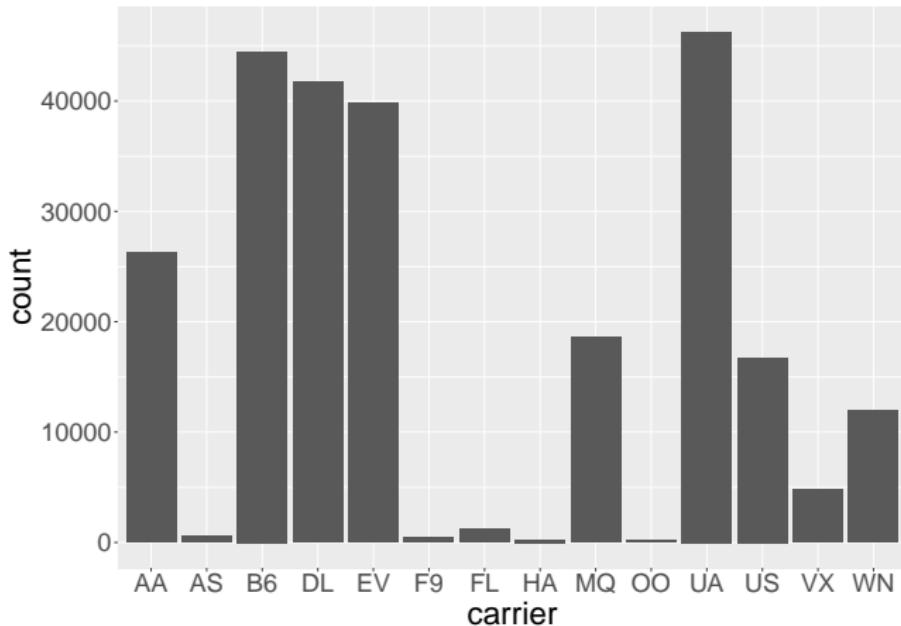
# Line Plot

- there needs to be a linearly ordered variable, typically time
- if multiple groups with an inherent order, choose line types with an order (e.g., of thickness or dash density)



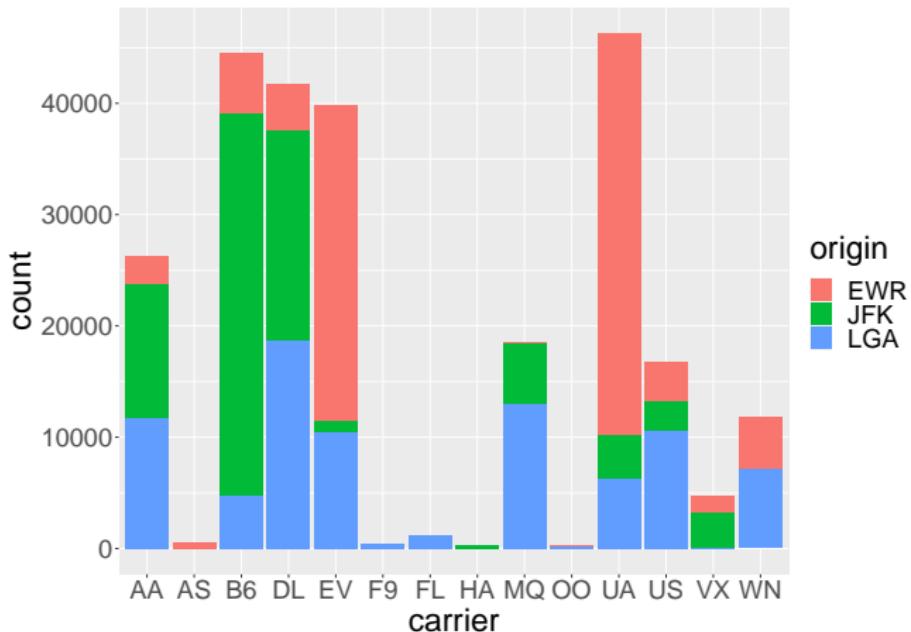
# Bar Plot

- shows the relationship between a numeric and a categorical variable



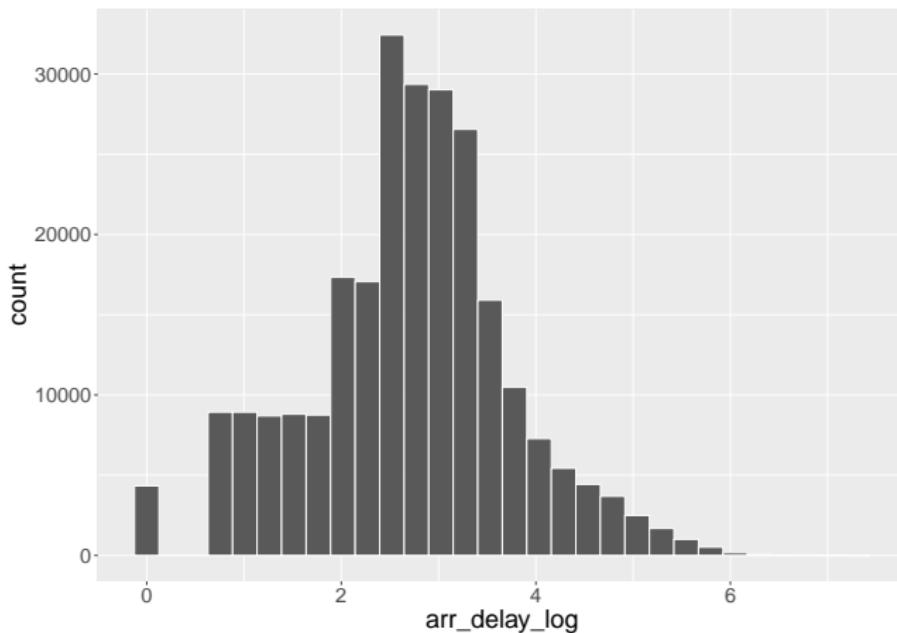
# Bar Plot

- shows the relationship between a numeric and a categorical variable
- or displays values for several grouping levels



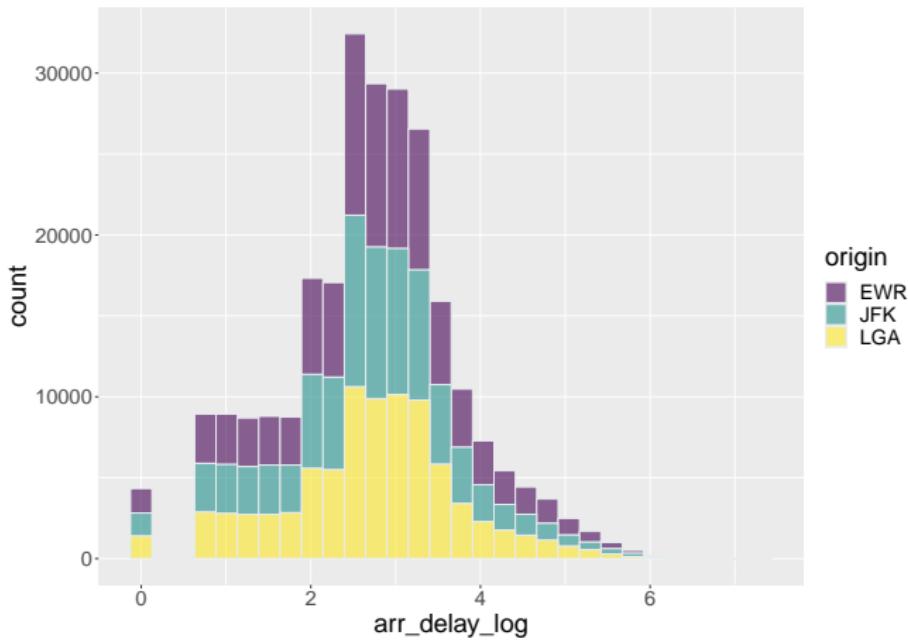
# Histogram

- the variable is cut into several bins, and the number of observations per bin is represented by the height of the bar
- graphical representation of the distribution of a numeric variable
- allows to check the distribution for mistakes, outliers, or extremes



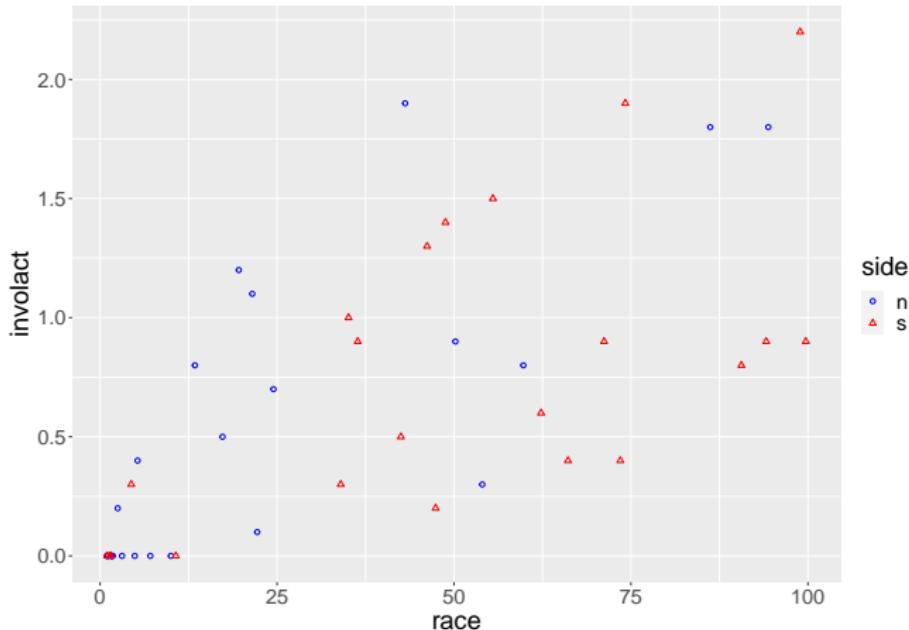
# Histogram

- histogram allows to compare the distribution of several variables (not too many...)



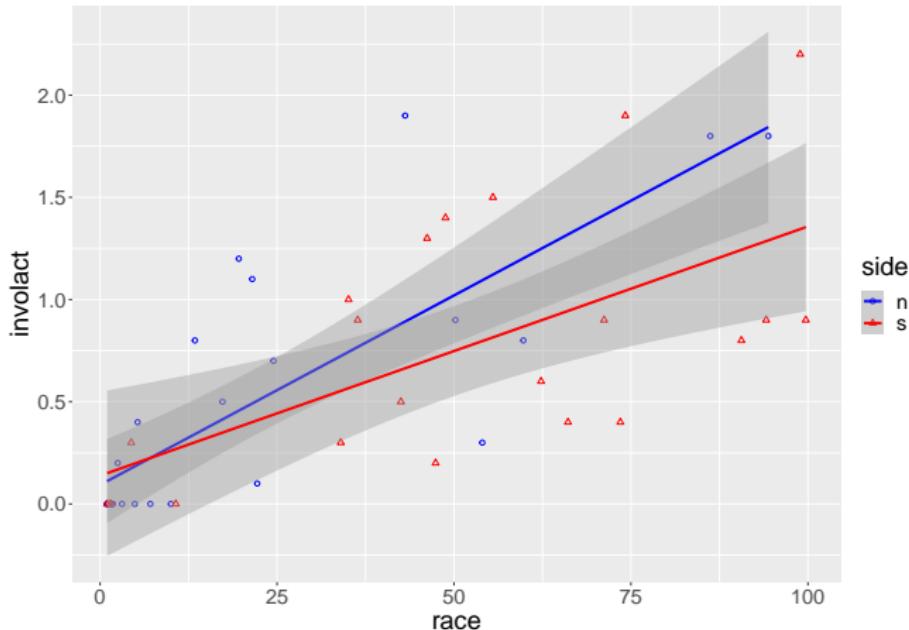
# Scatterplot

- allows to study the relationship between two numeric variables
- subgroups can be added for more insights into hidden patterns

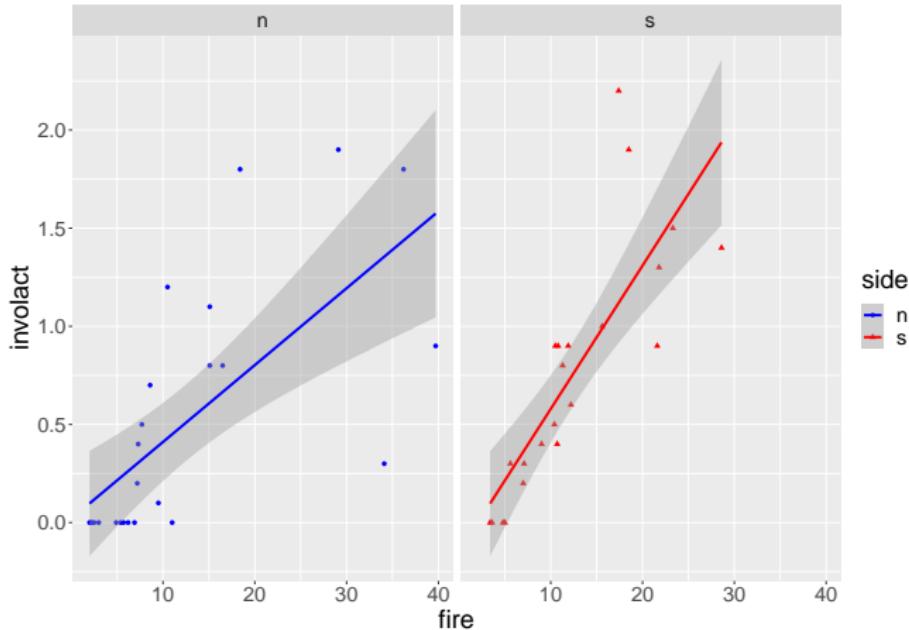


# Scatterplot with Regression Lines

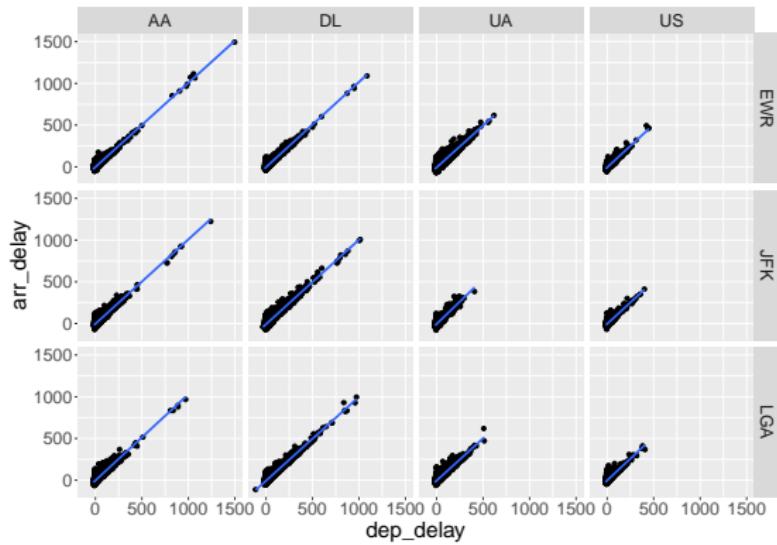
- different types of relationships can be detected by adding fitted regression curves



# Plot by Factor

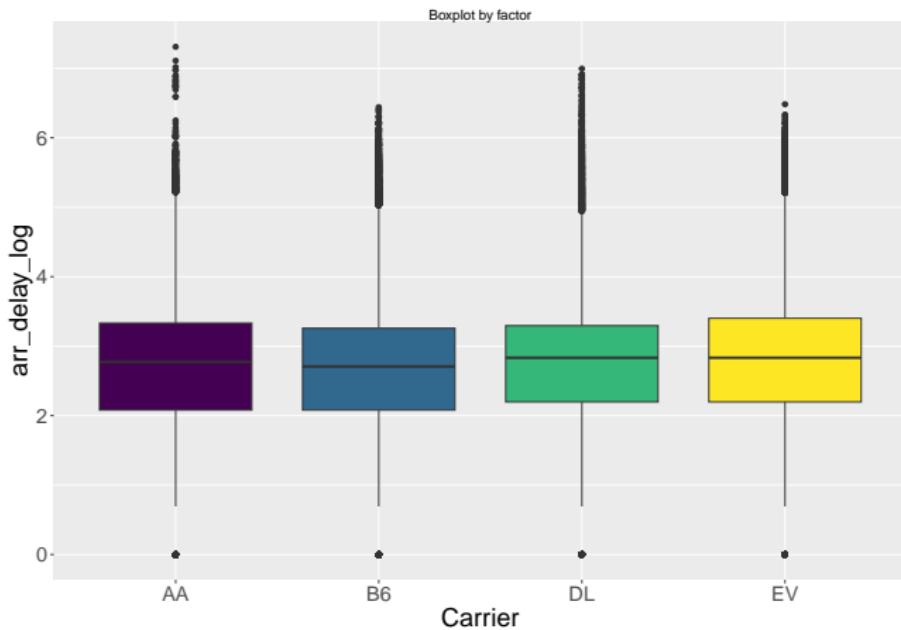


# Plot by Two Factors

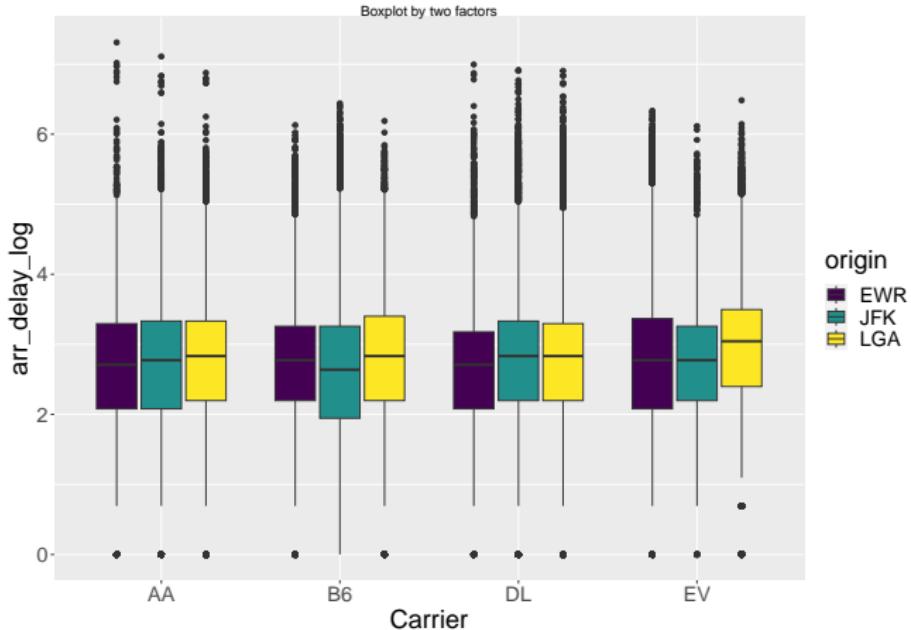


# Boxplot

- summarizes the distribution of a numeric variable for several groups (presence of outliers, symmetry, dispersion, . . . )

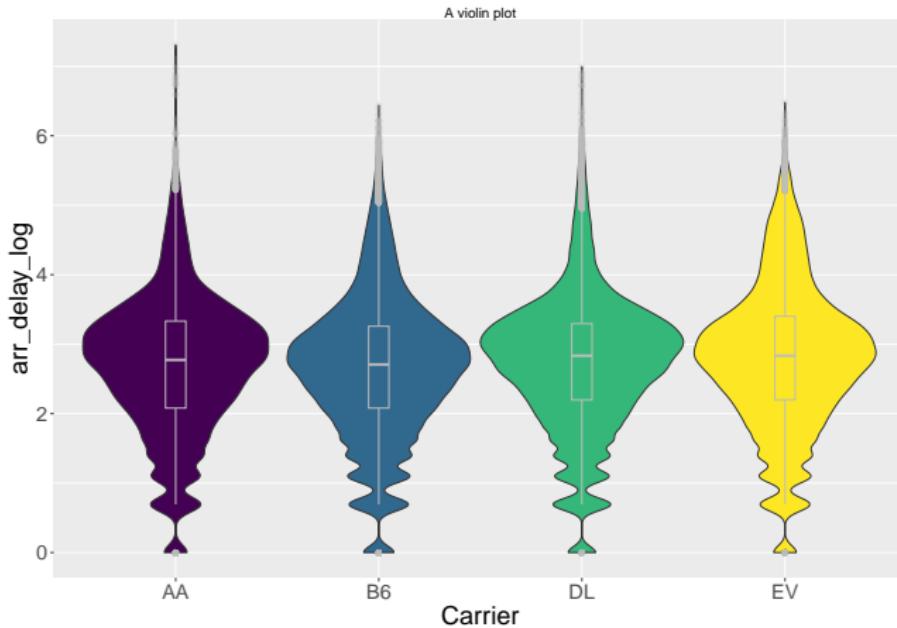


# Boxplot by Two Factors



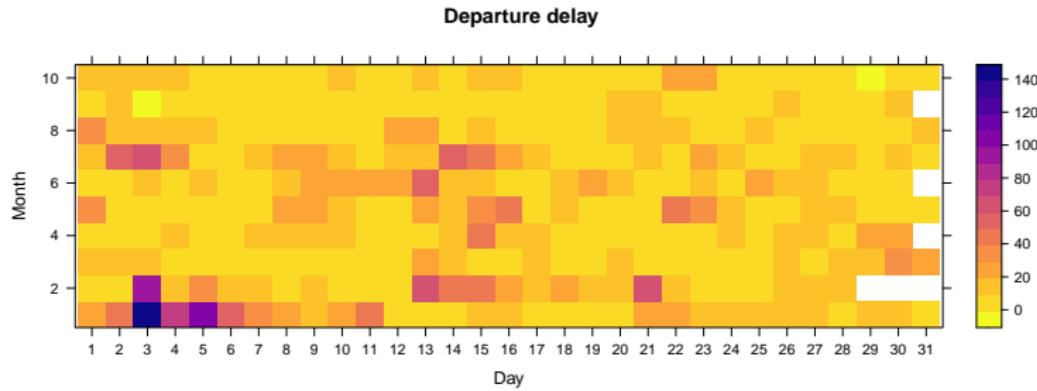
# Improved version of a boxplot: The violin plot

- adds a representation of the distribution (a smooth version of the histogram)



# Heatmaps

- represents a large matrix of data with colours reflecting their values (widely used for gene expression data)
- displays the output of hierarchical clustering (adding dendrogram)

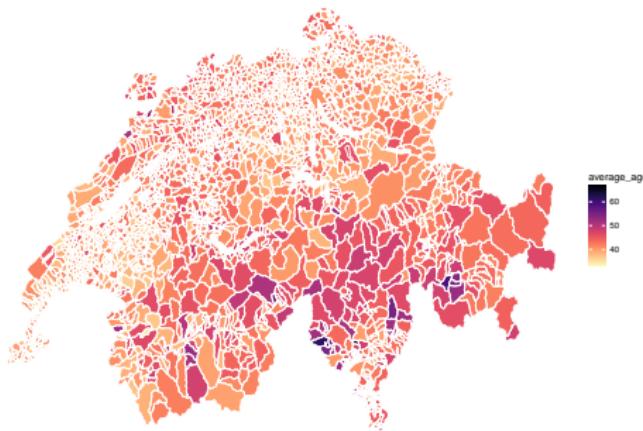


# Visualization of Spatial Data

Spatial data are complicated due to different

- data structures (vector or raster data)
- data sources
- data processing packages
- visualization packages

(Not so) short course about visualizing spatial data [here](#) (only if interested)

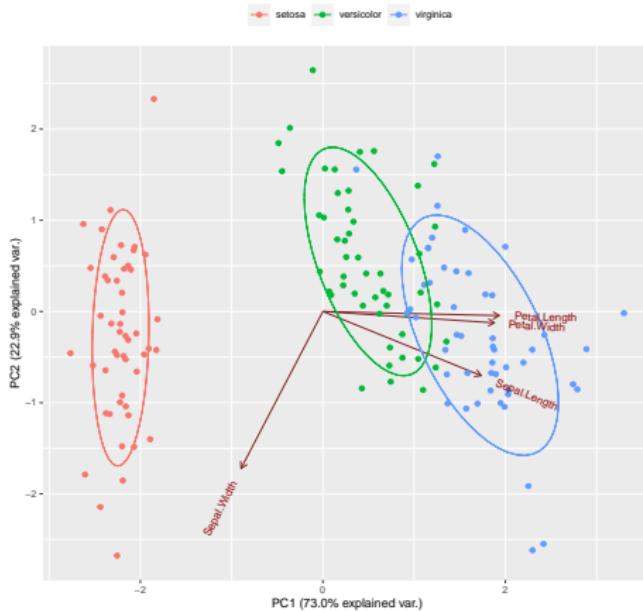


inspired by [this blogpost](#)

# Visualization of High-dimensional Data

Exploration of high-dimensional data can be done using

- clustering methods: clusters are formed based on similarities between features (e.g.,  $K$ -means, hierarchical clustering, ...)
- PCA: dimension reduction technique that preserves most of the data variability



# Good Visualization Practices

- context: exploratory vs. explanatory analysis/graphics
  - exploratory ... helps you understand the patterns
  - explanatory ... designed to communicate your understanding
- provide context (in text **and** in caption)
- seek simplicity, clarity, etc.
- gray scale often preferable
  - color-blindness (friendly palettes, e.g. [Cools](#))
- axes (scale, gaps, etc.)
  - text of appropriate size
- publication-specific conditions
- be artistic!
  - sometimes bend the rules (responsibly and justifiably)

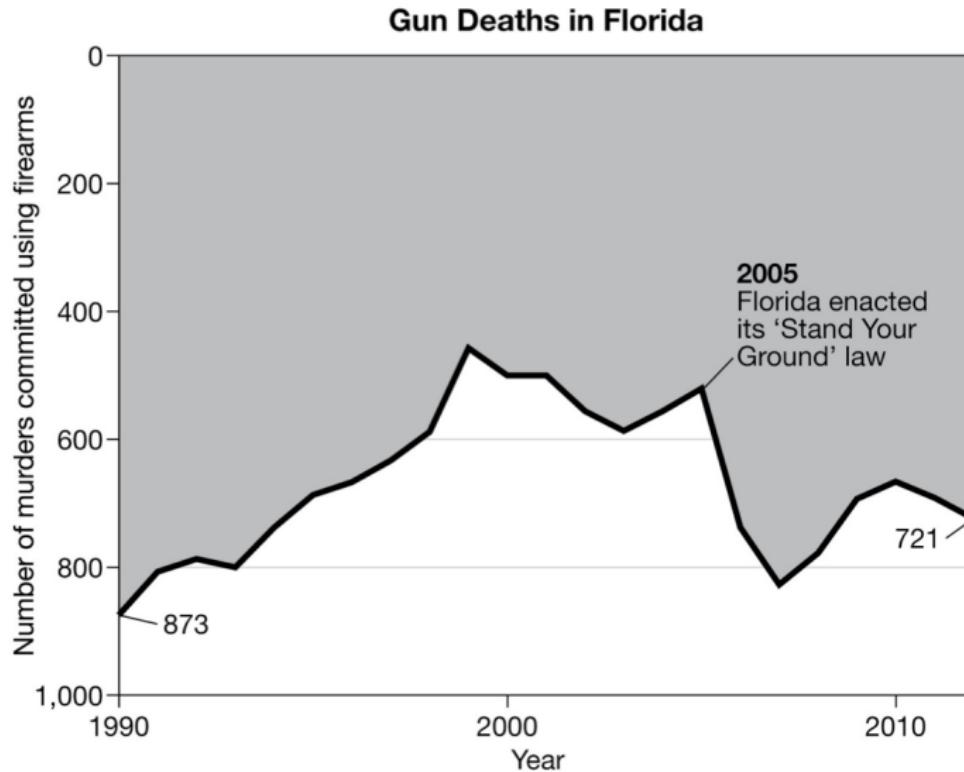
Find inspiration in [The R Graph Gallery](#).

Beware when exporting graphics.

## Section 3

### Bad Visualization Practices

# Reverted Axis

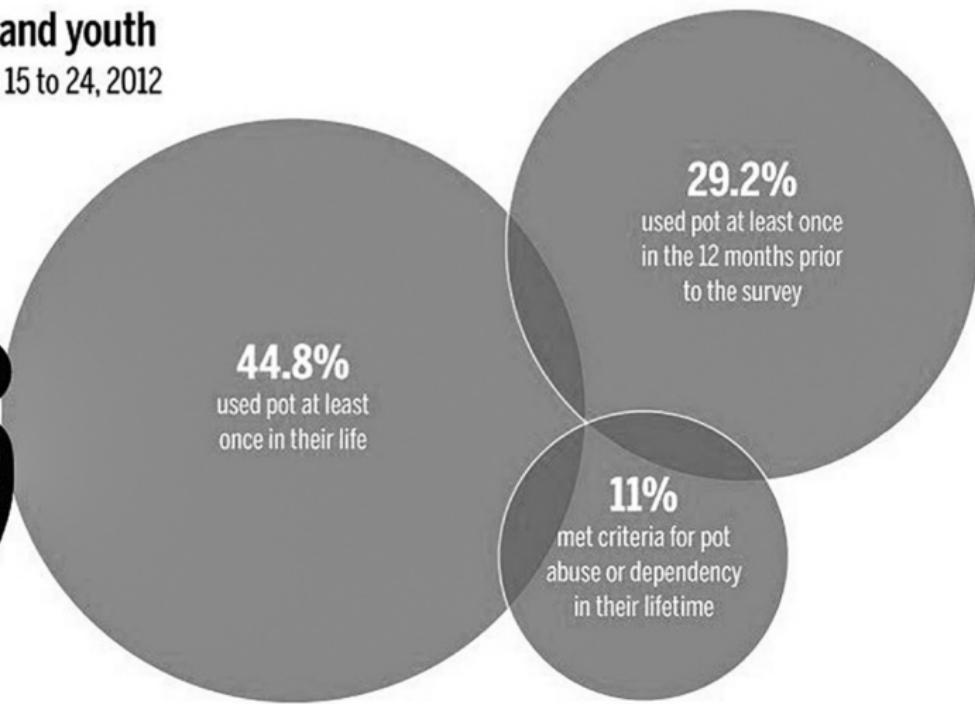


Source: Florida Department of Law Enforcement

# False Venn's Diagrams

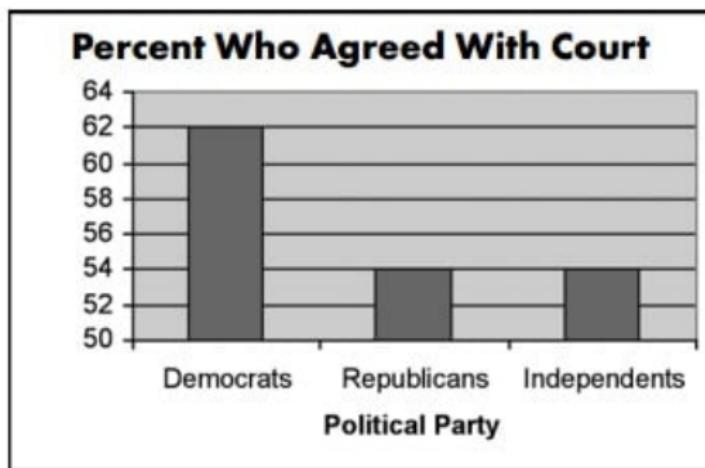
## Marijuana and youth

Canadians age 15 to 24, 2012



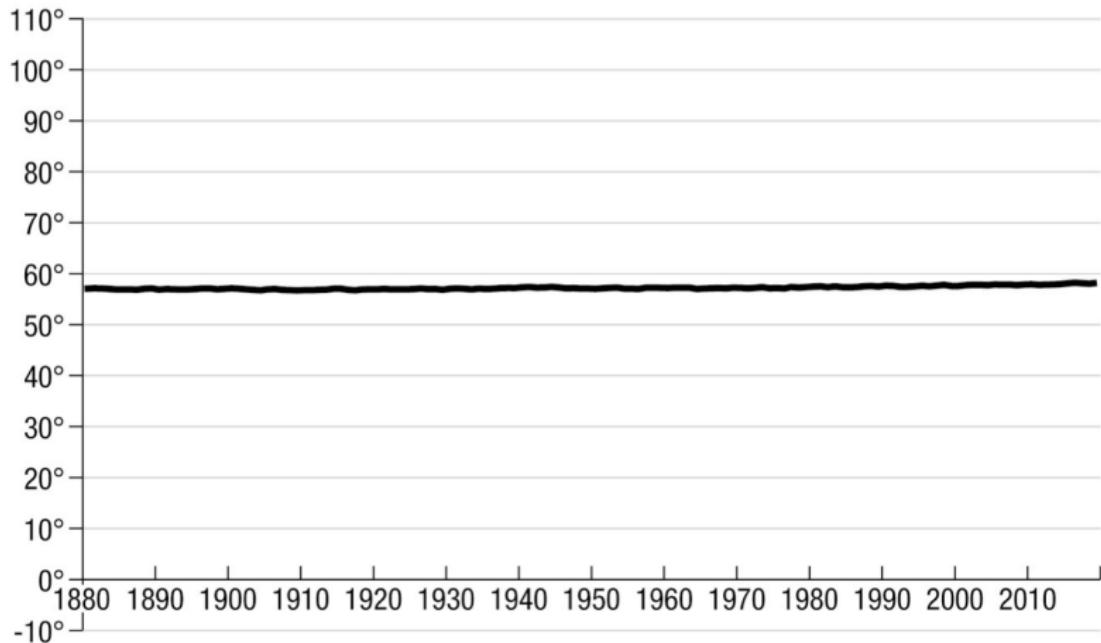
SOURCE: STATISTICS CANADA

# Missing Baseline



# Axis Starting at Zero

Average Annual Global Temperature in Fahrenheit, 1880–2019

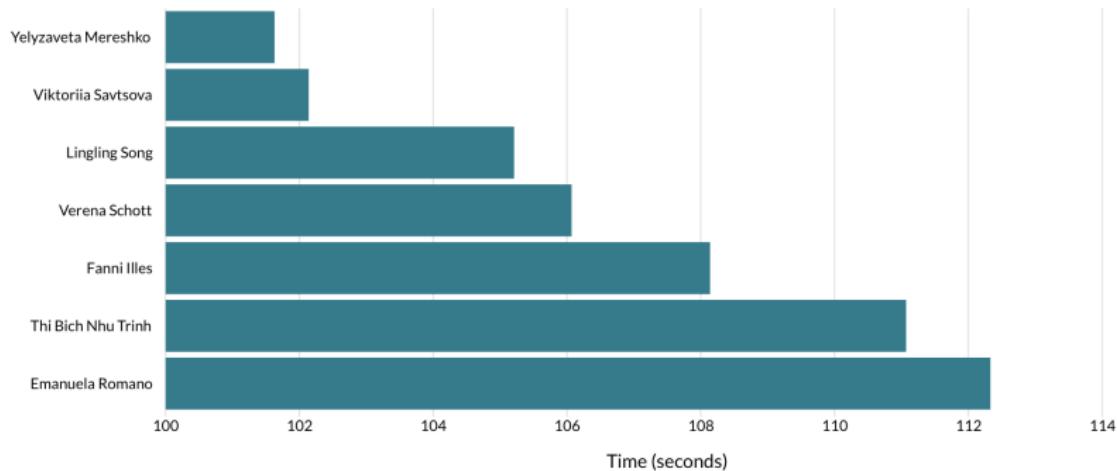


## Average Global Temperature by Year



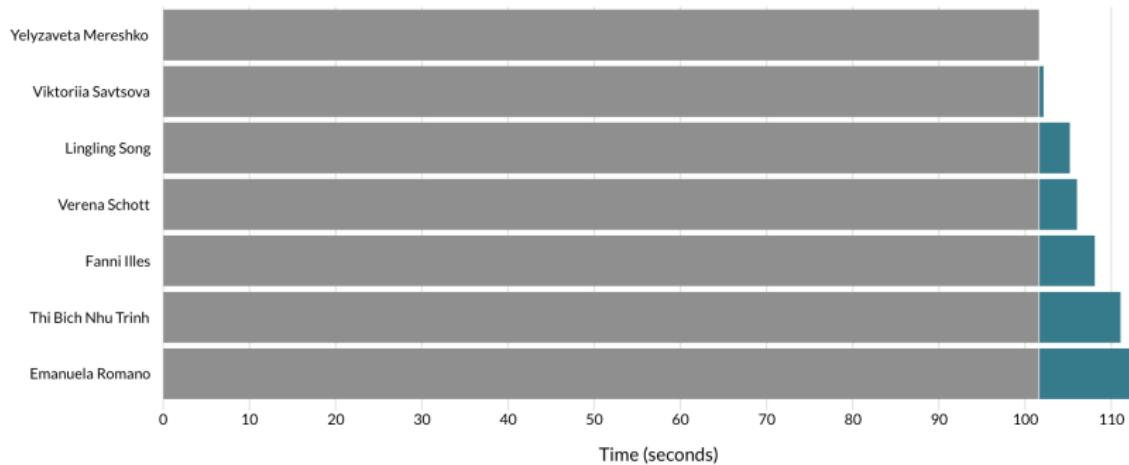
# Over-exaggeration with axes

Women's 100m breaststroke SB5, Rio 2016

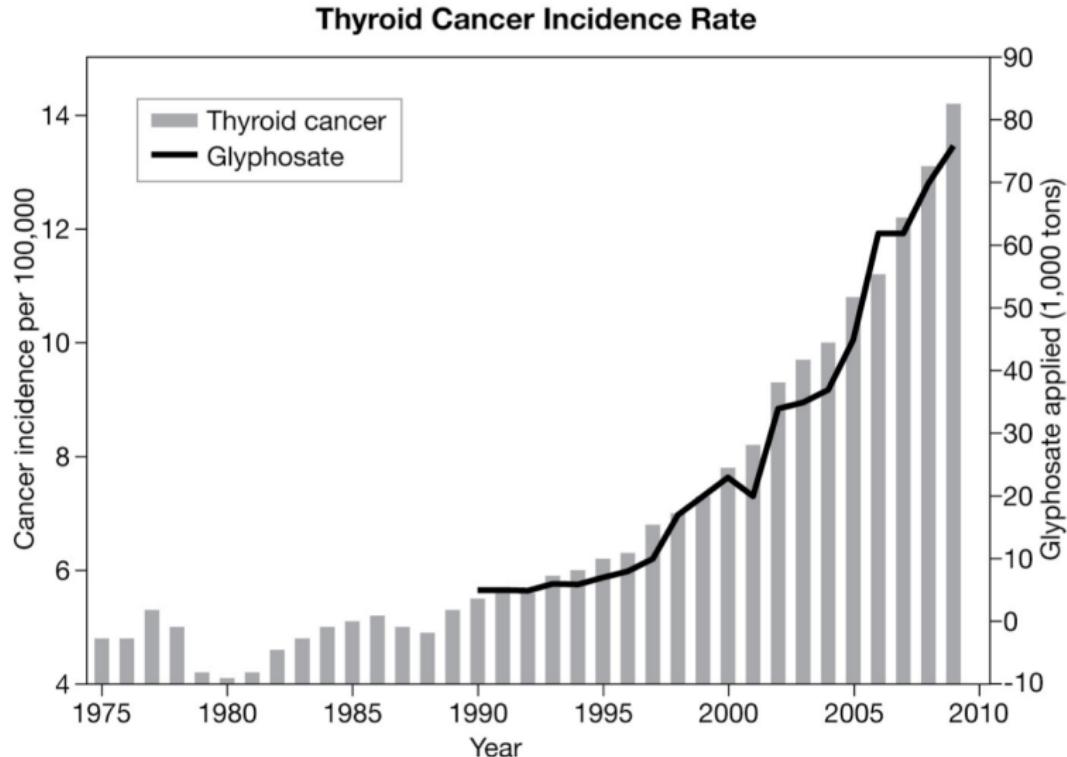


# Better

Women's 100m breaststroke SB5, Rio 2016



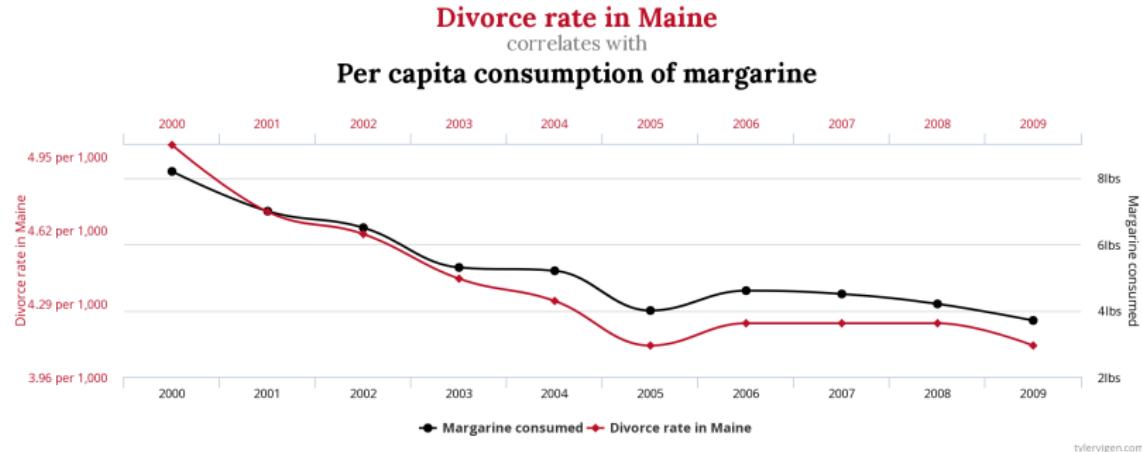
# Double Axes



- This is actually quite good, but double axes are usually problematic.

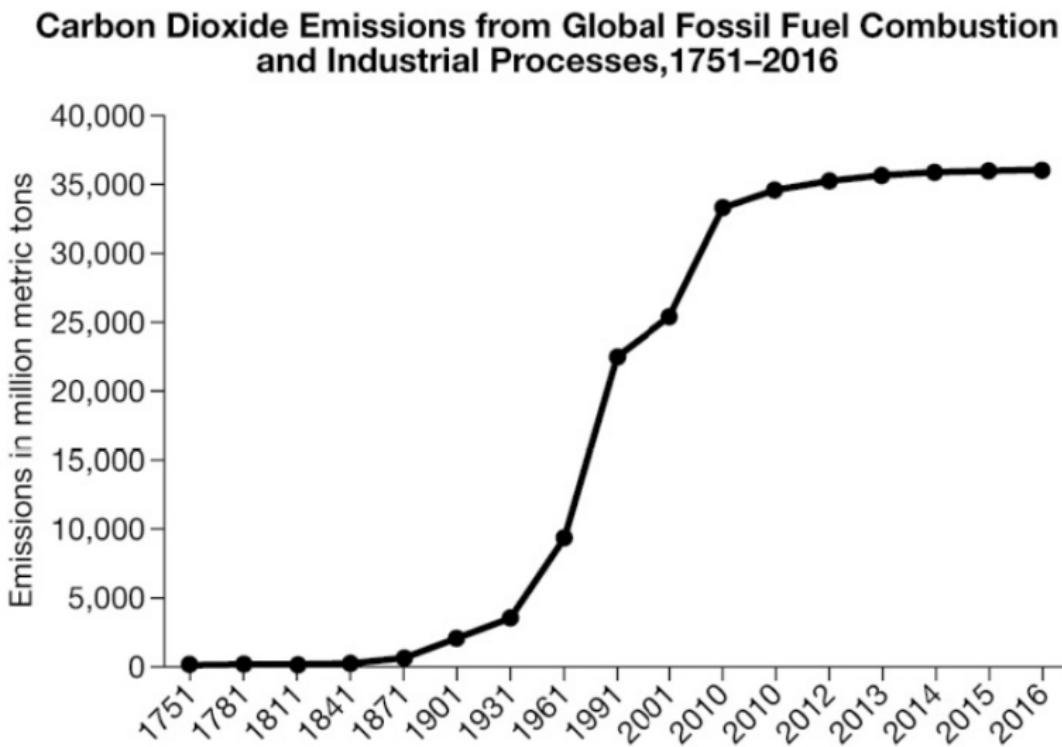
# False Causation

Correlation does not imply causation

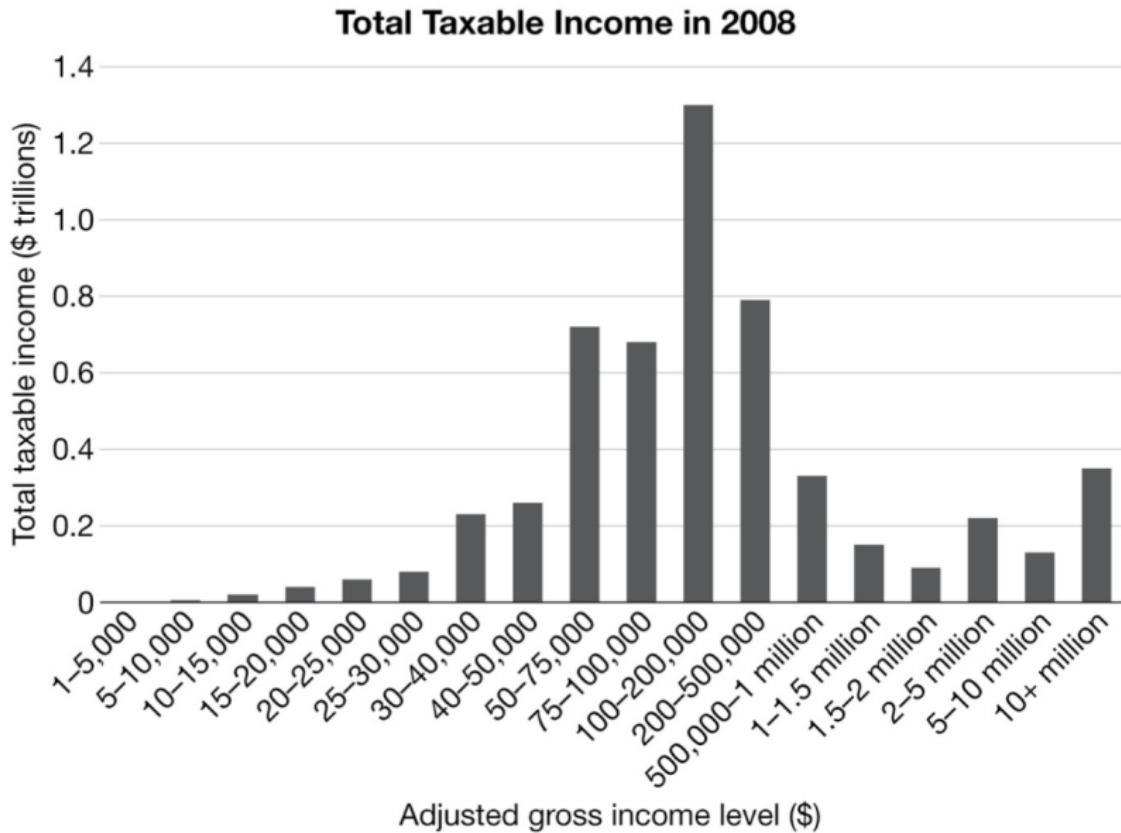


Data sources: National Vital Statistics Reports and U.S. Department of Agriculture

# Tweaking Axis

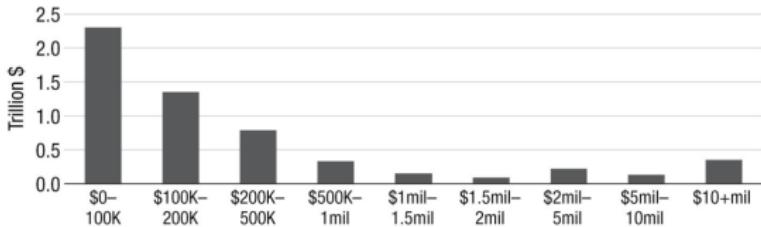


# Binning

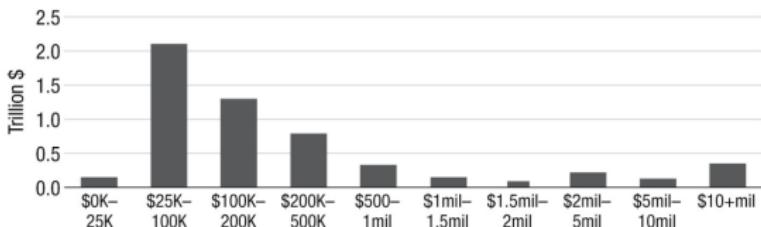


# Different Kinds of Binning

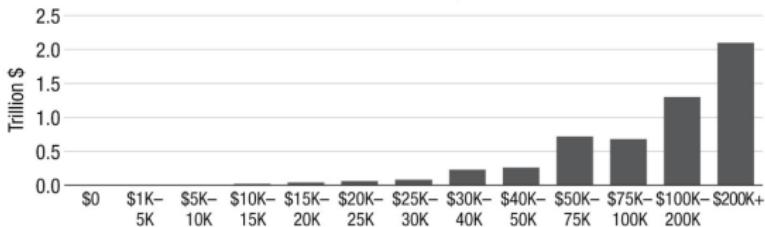
Tax the Poor!



Tax the Middle Class!

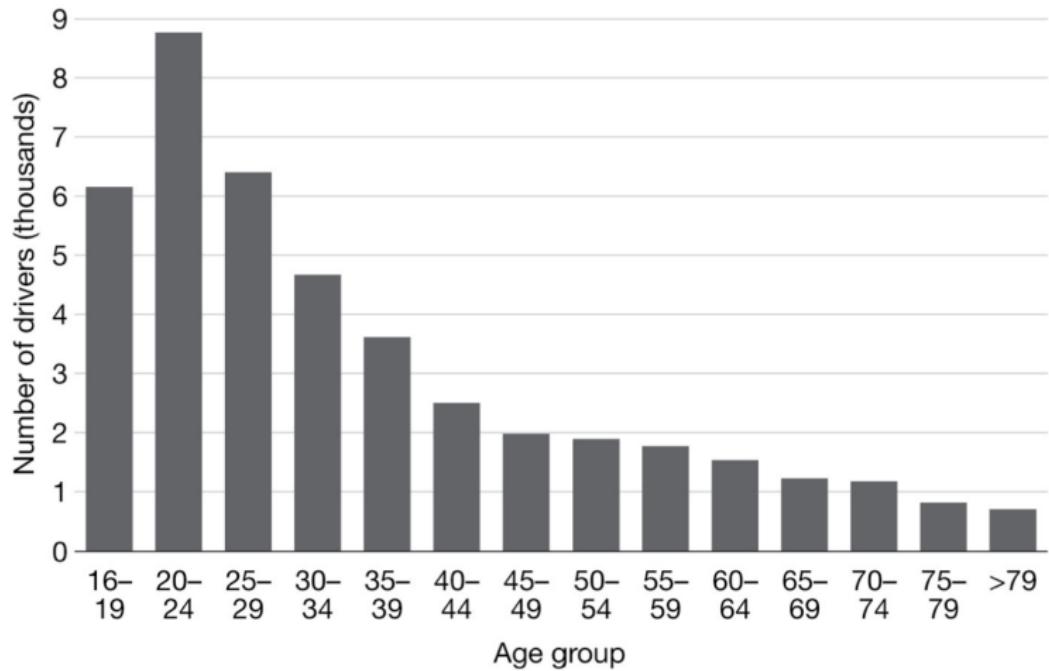


Tax the Wealthy!

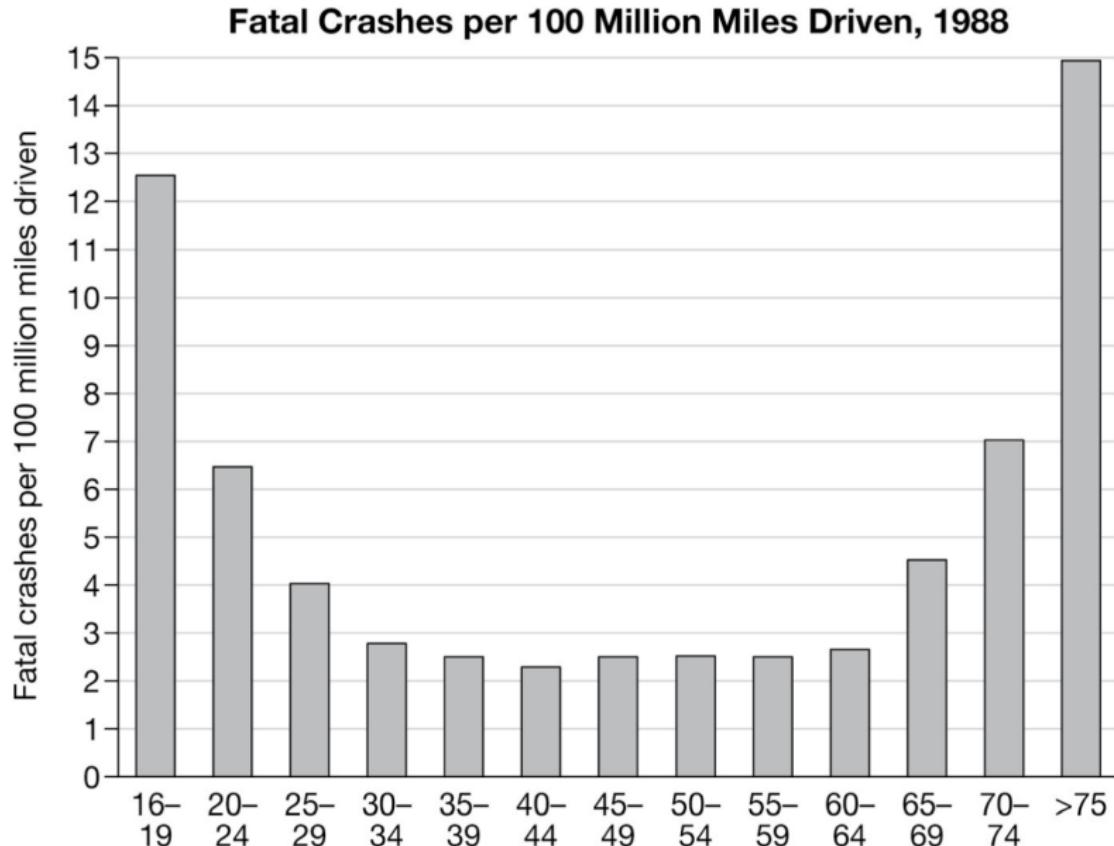


Total

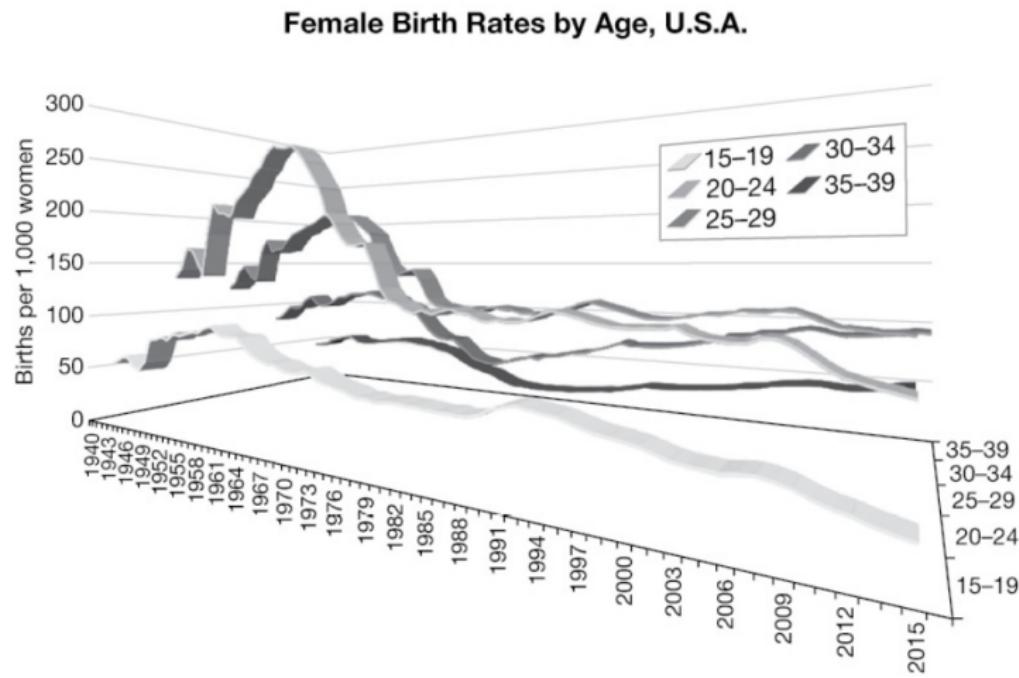
### Number of Drivers in Fatal Crashes, 1988



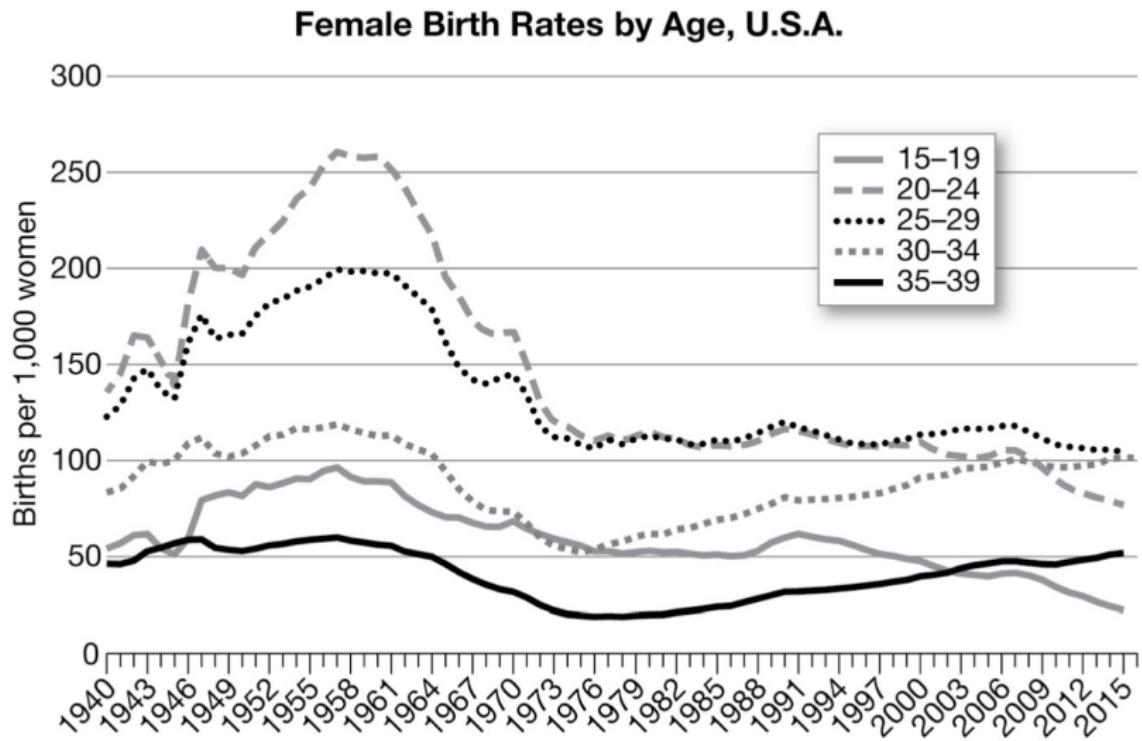
# Relative



# Useless 3D



# Better 2D

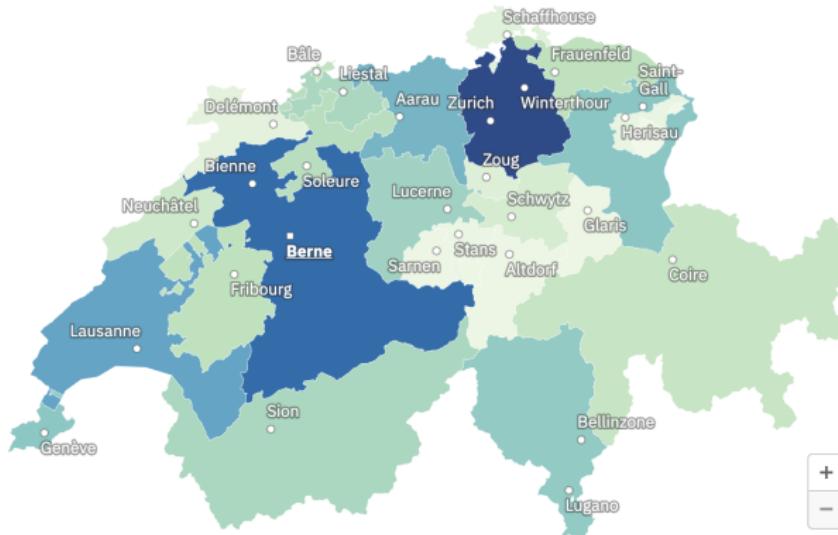


# Total vs. Relative Again

Nombre de personnes atteintes de démence dans les différents cantons suisses

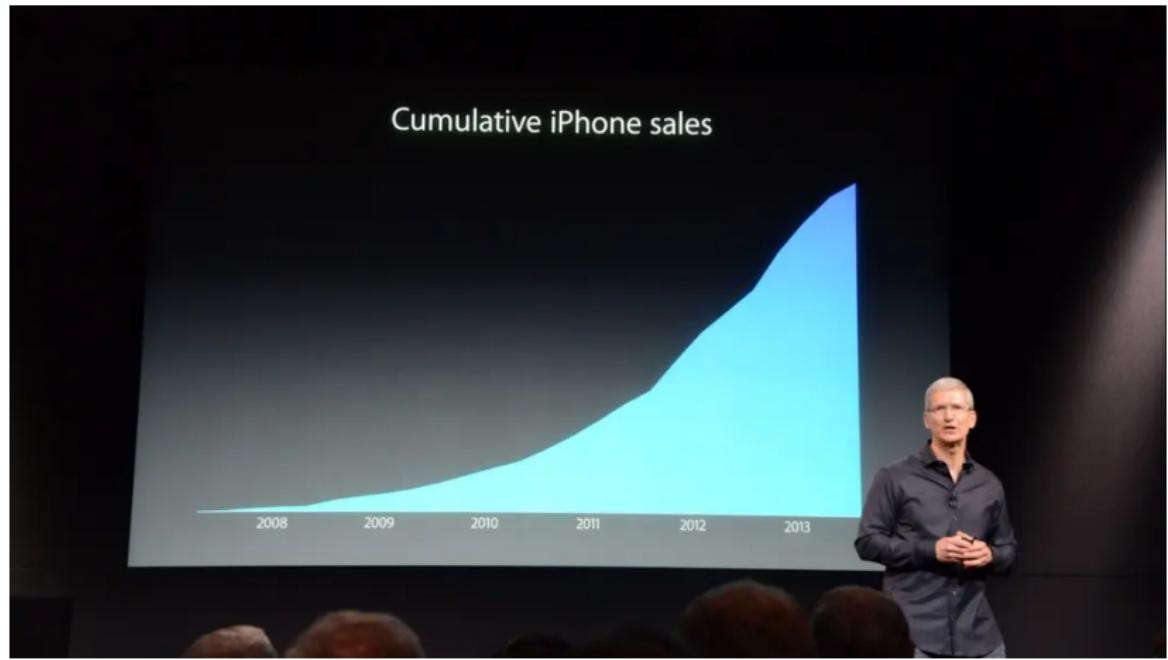
285

24380



Carte: G. Laplace. Nombre exact par canton en se positionnant dessus avec la souris; Source: Faits et chiffres (Alzheimer Suisse, 2021); [Récupérer les données](#)

# Missing Axis & Misguidance



# Assignment

**Small project [20%]**. Deadline on Week 5.

The goal of this project is *data exploration*. Details can be found on the dedicated [course page](#).

Some links to open data can be found [here](#).

# References

- Krause, Andreas, Nicola Rennie and Brian Tarran (2023) Best Practices for Data Visualisation
- Poldrack (2019) Statistical Thinking for the 21st Century
- JASA Ethical Guidelines for Statistical Practice
- Gelman (2018) Ethics in statistical practice and communication
- Wickham & Grolemund (2017) R for Data Science