

SCIENTIFIC DATA VISUALIZATION

AGENDA

Introduction

Basics of data visualization

Basic guidelines and common pitfalls

 **Break**

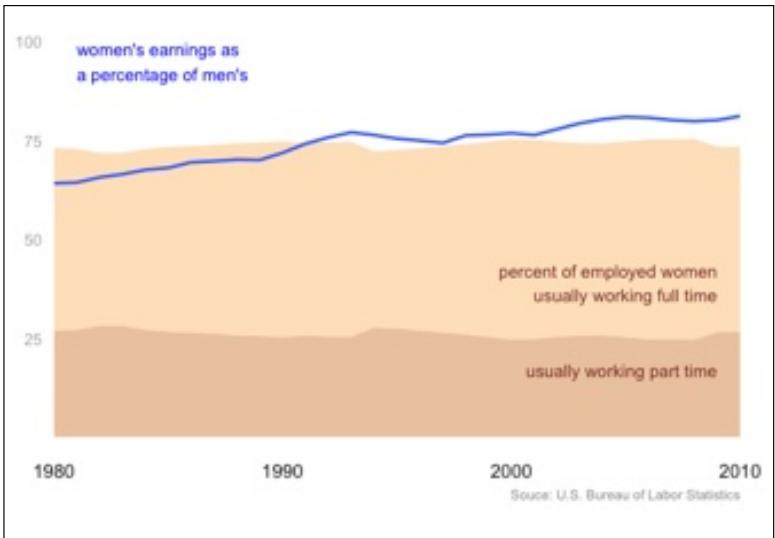
Examples and exercises in ggplot2

GEORGIOS KARAMANIS

Psychiatrist, Gender identity clinic, Uppsala University Hospital

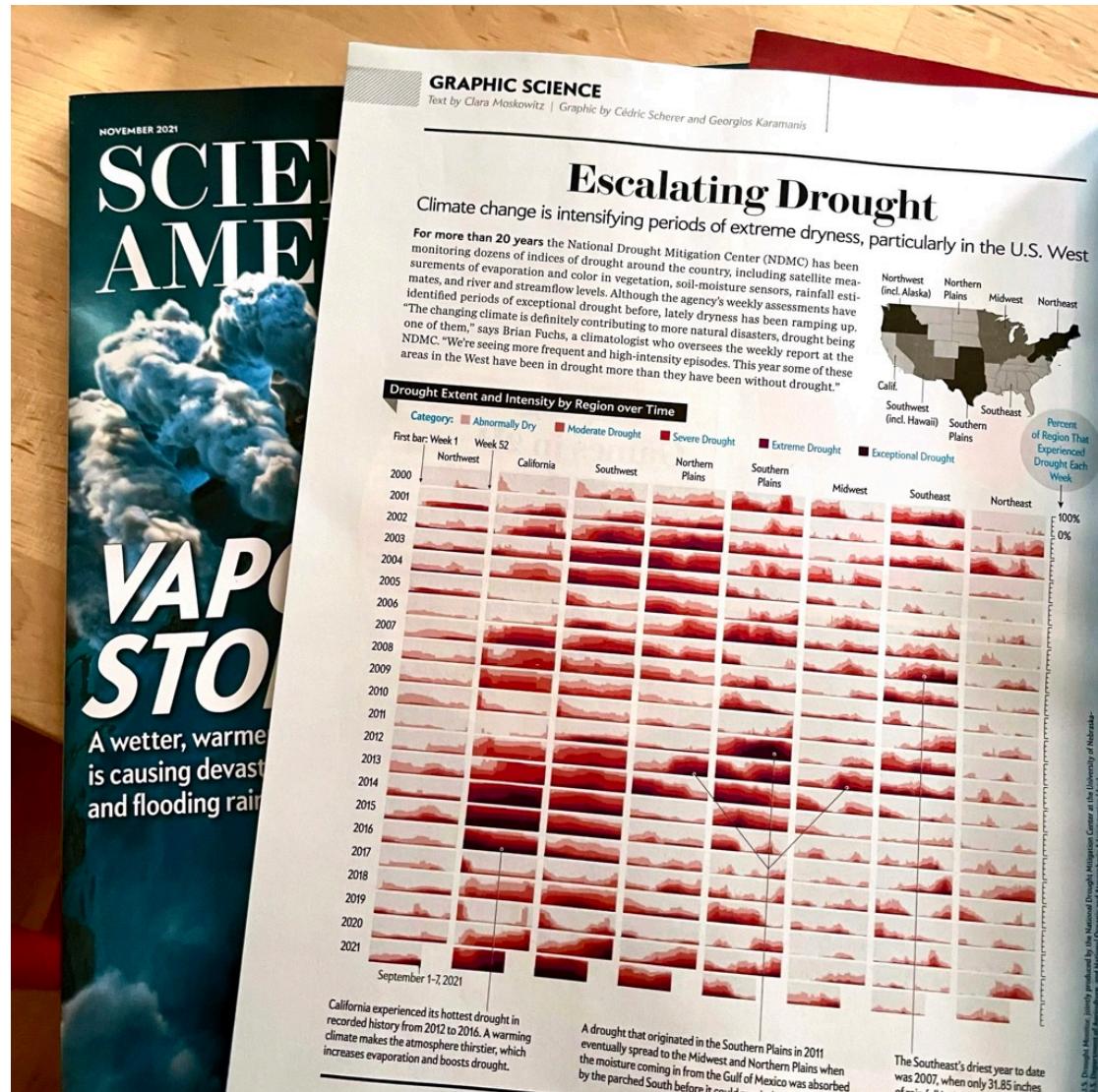
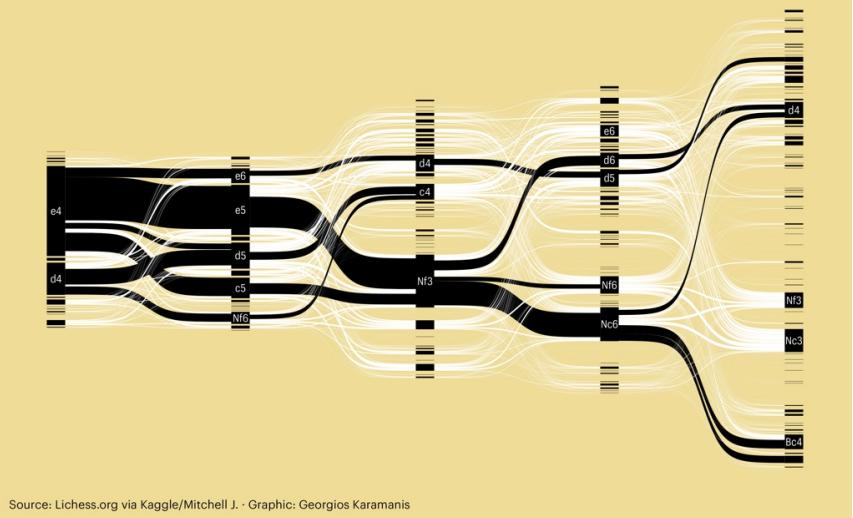
PhD student, Institution for medical sciences, Uppsala University

Data visualization designer, Explained

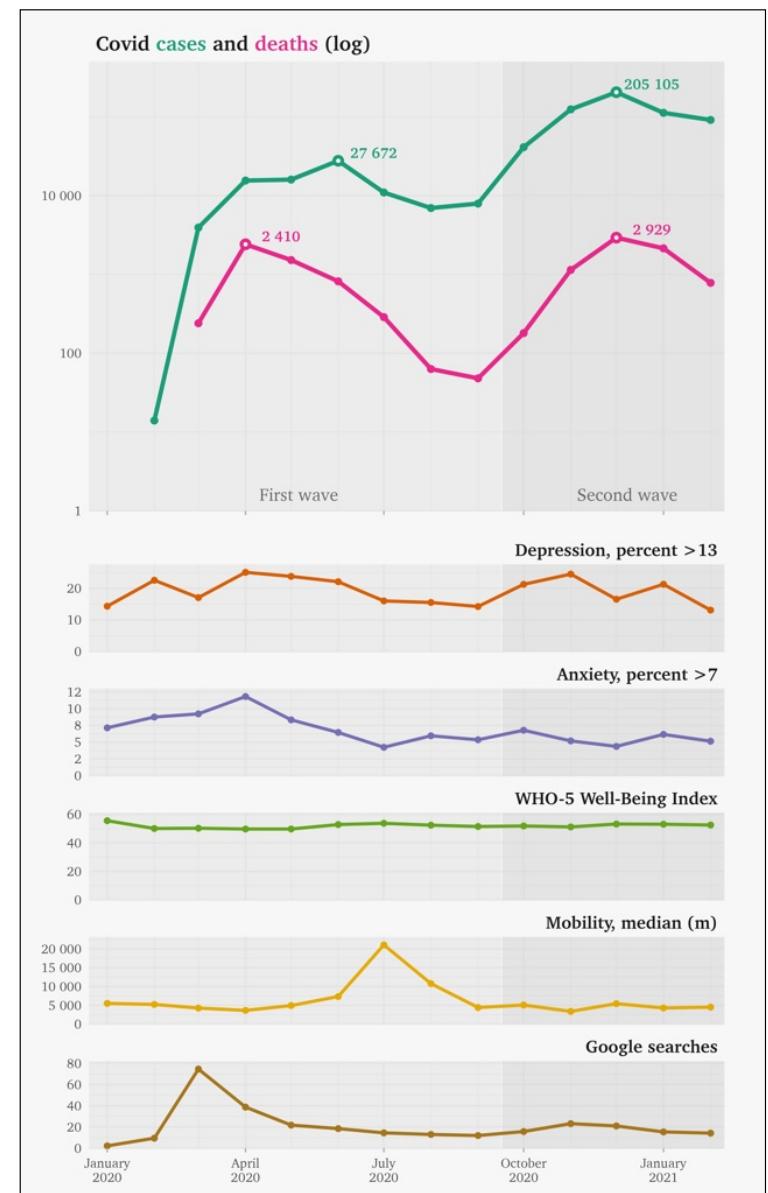
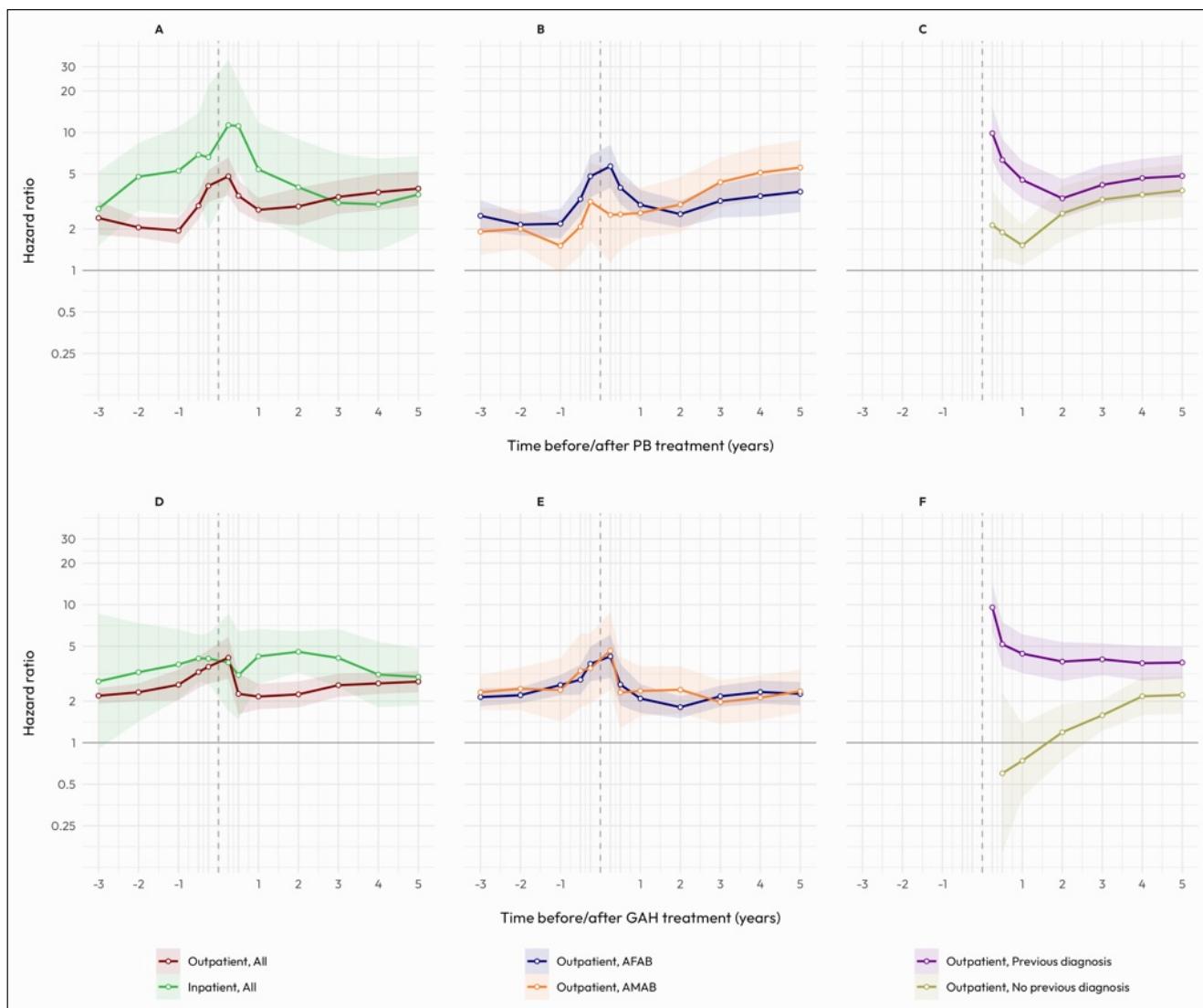


First Five Moves

Flow of opening moves in over 20 000 chess games from Lichess.org. The most popular moves are highlighted in black, revealing common opening strategies.



<https://www.scientificamerican.com/article/climate-change-drives-escalating-drought/#>

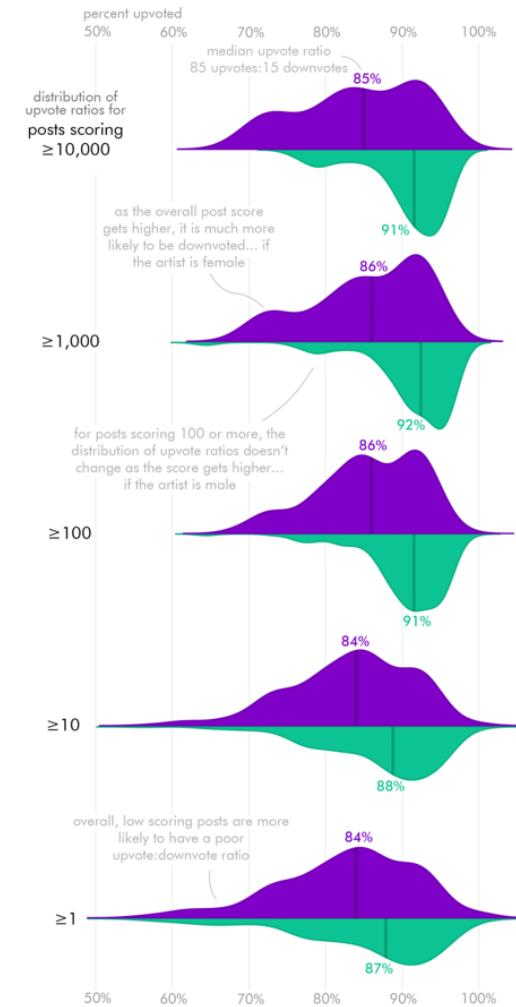


INTRODUCTION

DATA VISUALIZATION

Data visualization blends art and science to convey stories from data via graphical representations.

Controversy of art posts on r/pics
posts featuring **women** posing with their art are more likely to be downvoted than posts with **men** posing



more at erdavis.com | @erindataviz

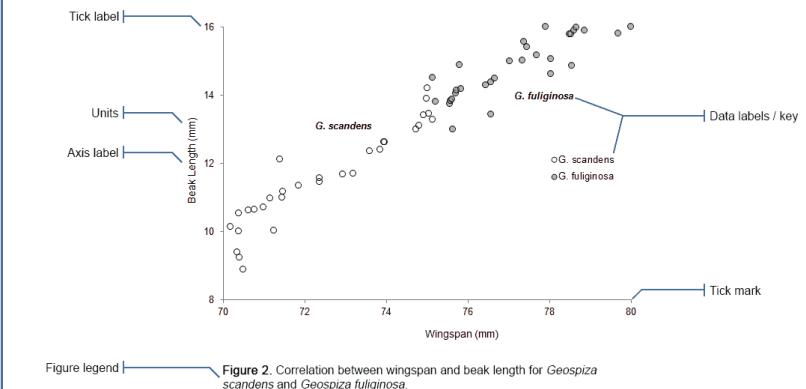
<https://erdavis.com/2021/06/14/do-women-who-pose-with-their-art-on-reddit-get-more-upvotes/>

SCIENTIFIC CHARTS

Figures in scientific publications are critically important because they often show the data supporting key findings

Weissgerber TL, Milic NM, Winham SJ, Garovic VD (2015) Beyond Bar and Line Graphs: Time for a New Data Presentation Paradigm. PLoS Biol 13(4): e1002128. <https://doi.org/10.1371/journal.pbio.1002128>

Common features of a graph

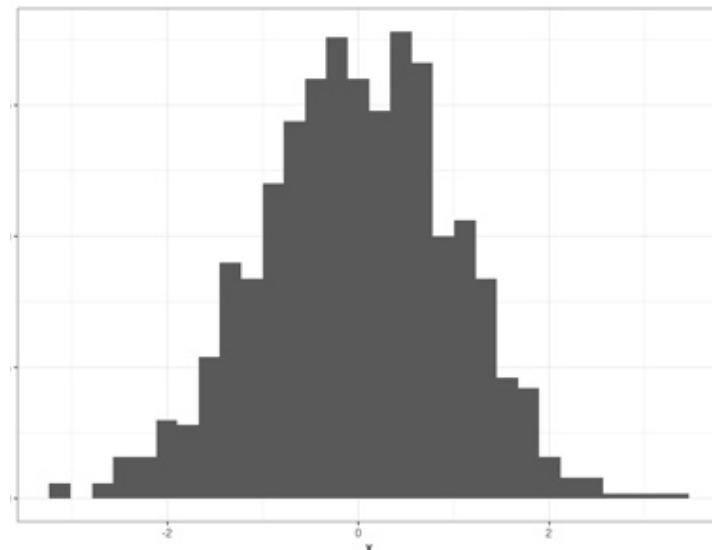


<https://www.clips.edu.au/displaying-data/>

WHY CHARTS?

Easier overview of data
(especially complex)

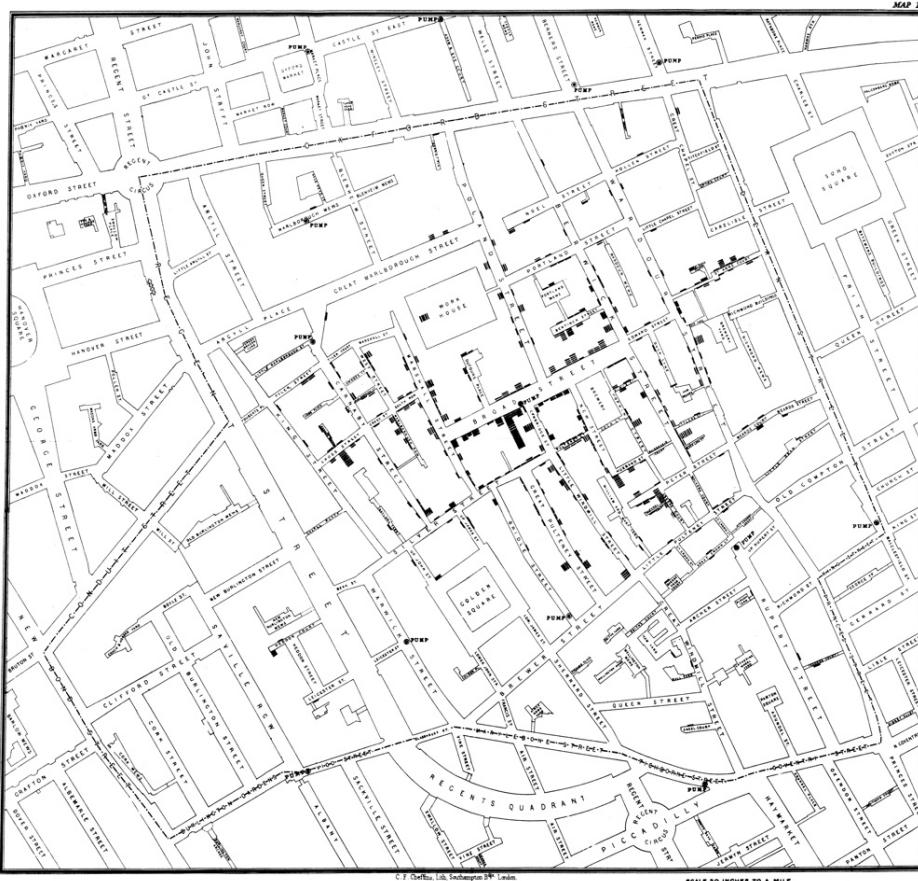
x
1 0.376935655
2 -0.603182135
3 0.659734512
4 -1.597222895
5 -0.645683665
6 -0.929817315
7 1.401684999
8 0.520906260
9 0.369003441
10 0.573936602
11 1.304568093
12 -2.537849064
13 0.920903182
14 1.419055702
15 -0.288114527
16 1.842139177
17 0.647224489
18 1.793445537
19 -2.002214213
20 0.885130196
21 -0.369039857
22 -0.136589866



WHY CHARTS?

Find new patterns

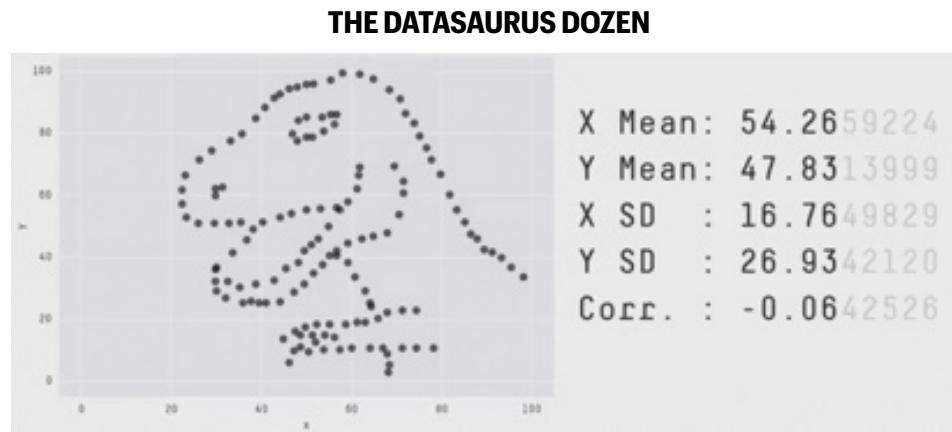
BROAD STREET PUMP OUTBREAK, JOHN SNOW, 1854



https://en.wikipedia.org/wiki/1854_Broad_Street_cholera_outbreak

WHY CHARTS?

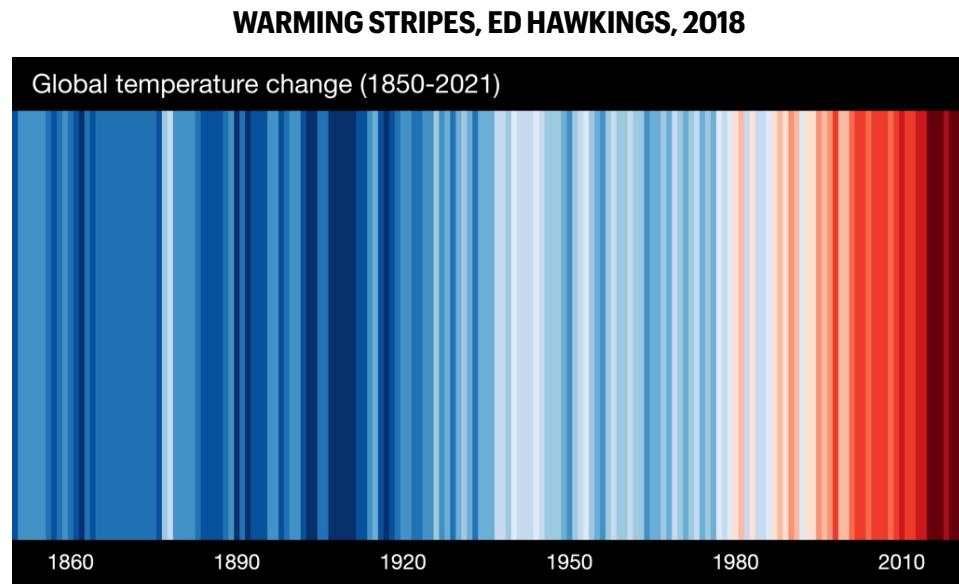
Easier for comparisons of data and their statistics



<https://www.autodesk.com/research/publications/same-stats-different-graphs>

WHY CHARTS?

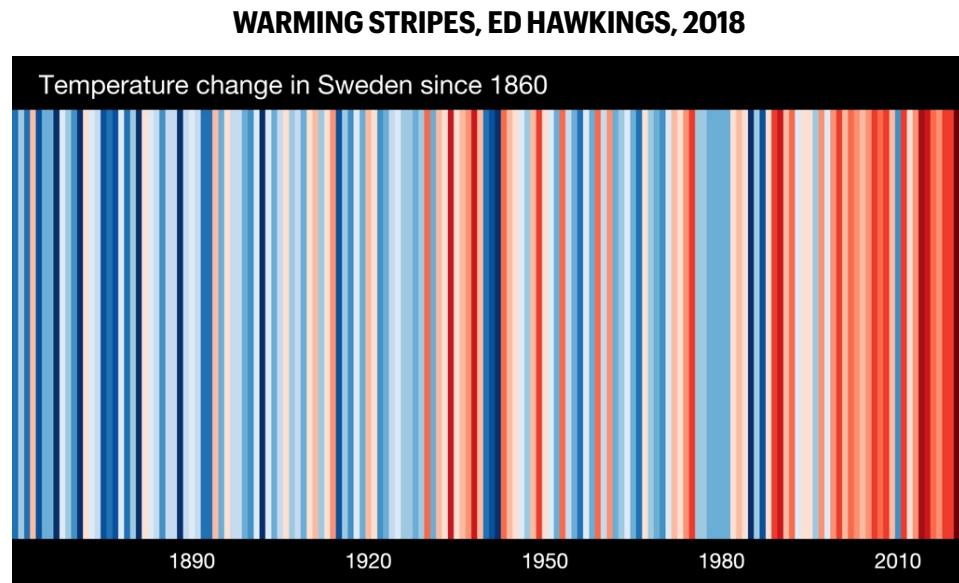
To better communicate your point



<https://www.climate-lab-book.ac.uk/2018/warming-stripes/>

WHY CHARTS?

To better communicate your point

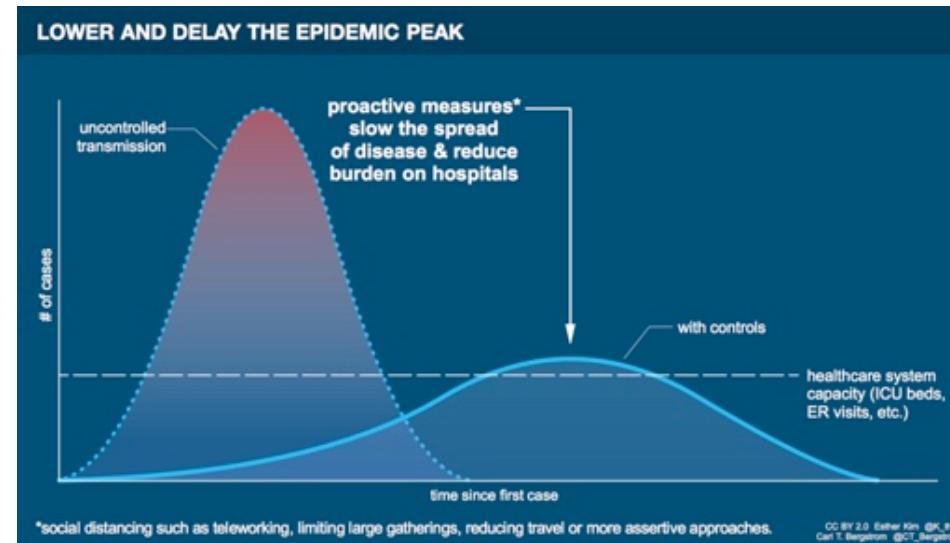


<https://showyourstripes.info/l/europe/sweden/all>

WHY CHARTS?

They look good, are memorable and shareable

ESTHER KIM & CARL T. BERGSTROM, 2019



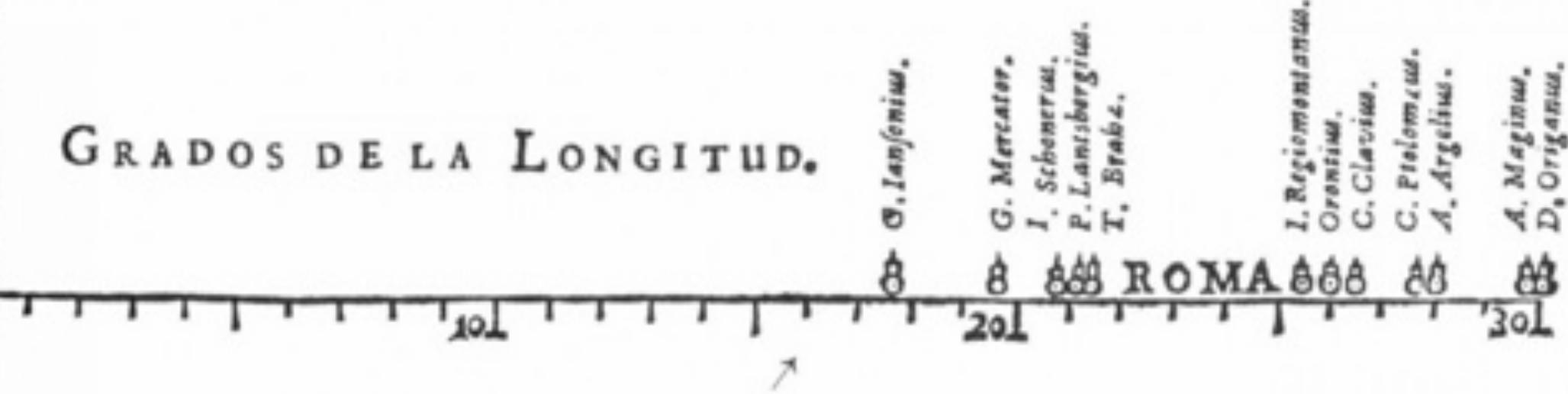
<http://ctbergstrom.com/covid19.html>



**SCIENTIFIC CHARTS
HAVE A LONG HISTORY**

HOA TOLEDO.

GRADOS DE LA LONGITUD.



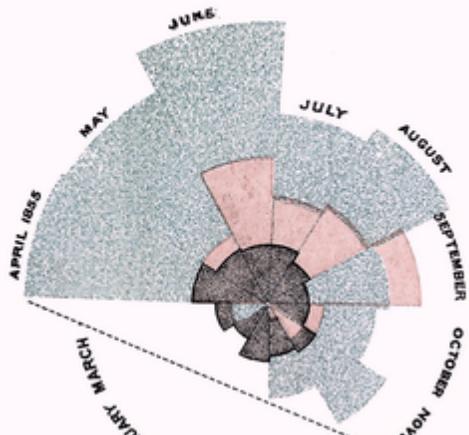
<https://historyofinformation.com/detail.php?id=3415>

DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST, FLORENCE NIGHTINGALE, 1858

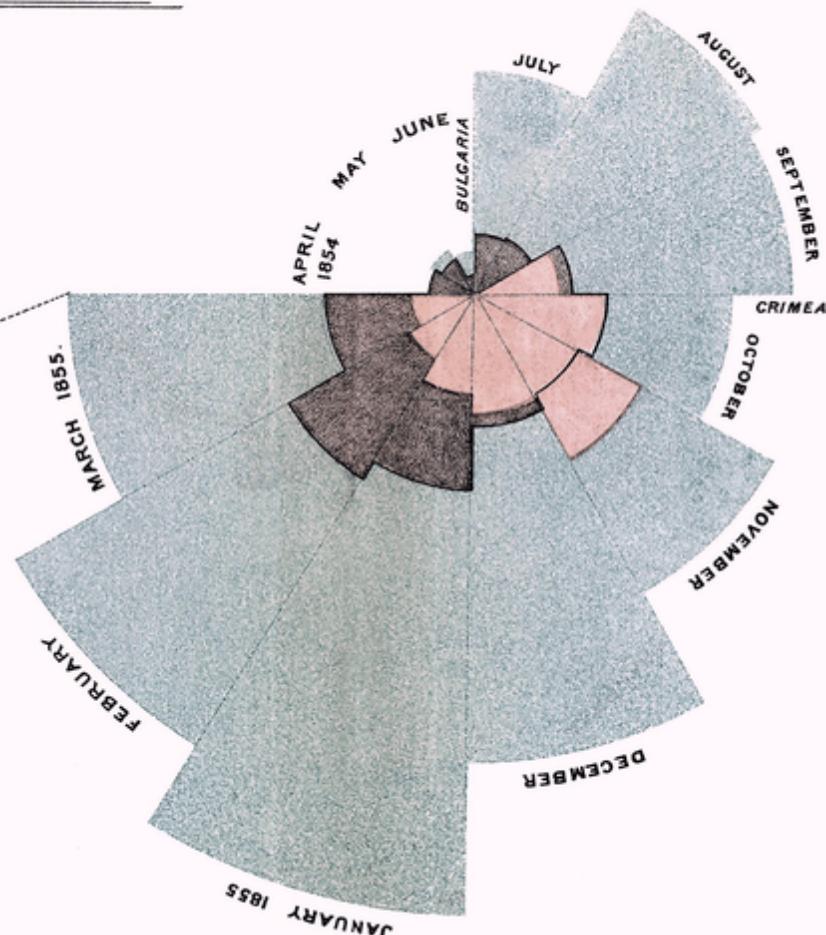
DIAGRAM OF THE CAUSES OF MORTALITY

IN THE ARMY IN THE EAST.

2.
APRIL 1855 TO MARCH 1856.



1.
APRIL 1854 TO MARCH 1855.



red, & black wedges are each measured from
common vertex.
red from the centre of the circle represent area
from Preventible or Mitigable Zymotic diseases; the
d from the centre the deaths from wounds, & the
red from the centre the deaths from all other causes.
the red triangle in Nov^r 1854 marks the boundary
all other causes during the month.
in 1855, the black area coincides with the red;
in February 1856, the blue coincides with the black.
be compared by following the blue, the red & the
them.



https://en.wikipedia.org/wiki/Florence_Nightingale

**BAD CHARTS ARE HARD TO READ,
MUDDLE OR HARM OUR MESSAGE
AND EVEN OUR CREDIBILITY**

CLINICAL → MEANINGFUL ← BEAUTIFUL

John Burn-Murdoch
Making meaningful graphics

BASICS OF DATA VISUALIZATION

VISUAL ELEMENTS

position

length

angle

direction

shape

area

volume

saturation/lightness

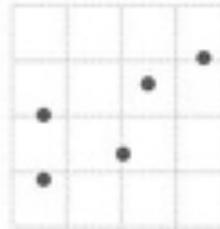
color hue

Visual cues

When you visualize data, you encode values to shapes, sizes, and colors.

Position

Where in space the data is



Length

How long the shapes are



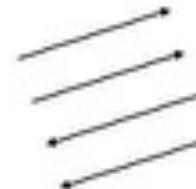
Angle

Rotation between vectors



Direction

Slope of a vector in space



Shapes

Symbols as categories



Area

How much 2-D space



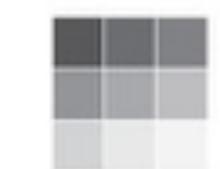
Volume

How much 3-D space



Color saturation

Intensity of a color hue



Color hue

Usually referred to as color

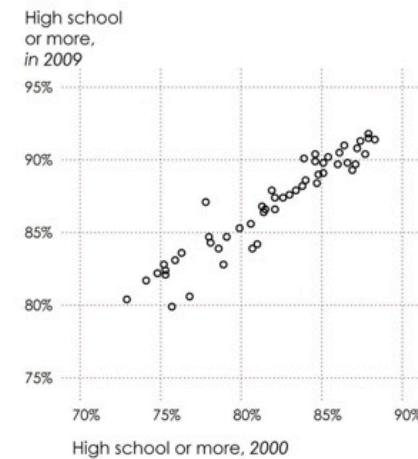


<https://flowingdata.com/data-points/DataPoints-Ch3.pdf>

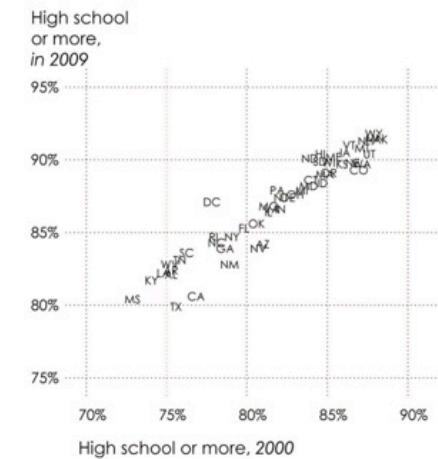
VISUAL ELEMENTS

combinations

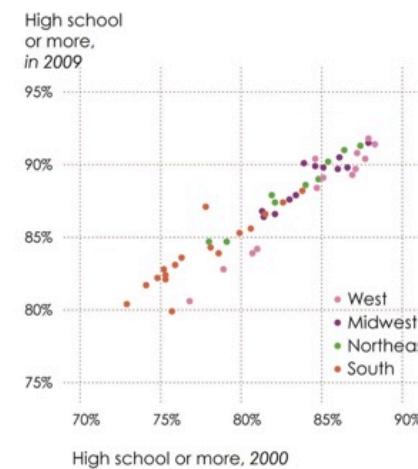
Position



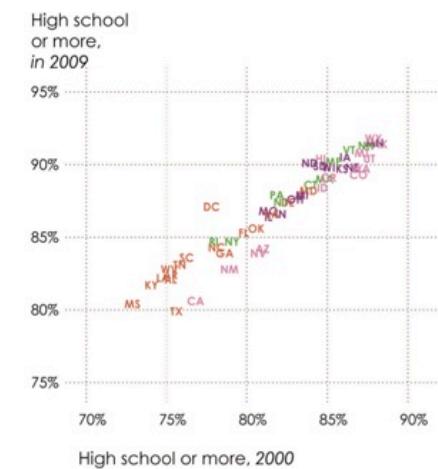
Position + Symbols



Position + Color

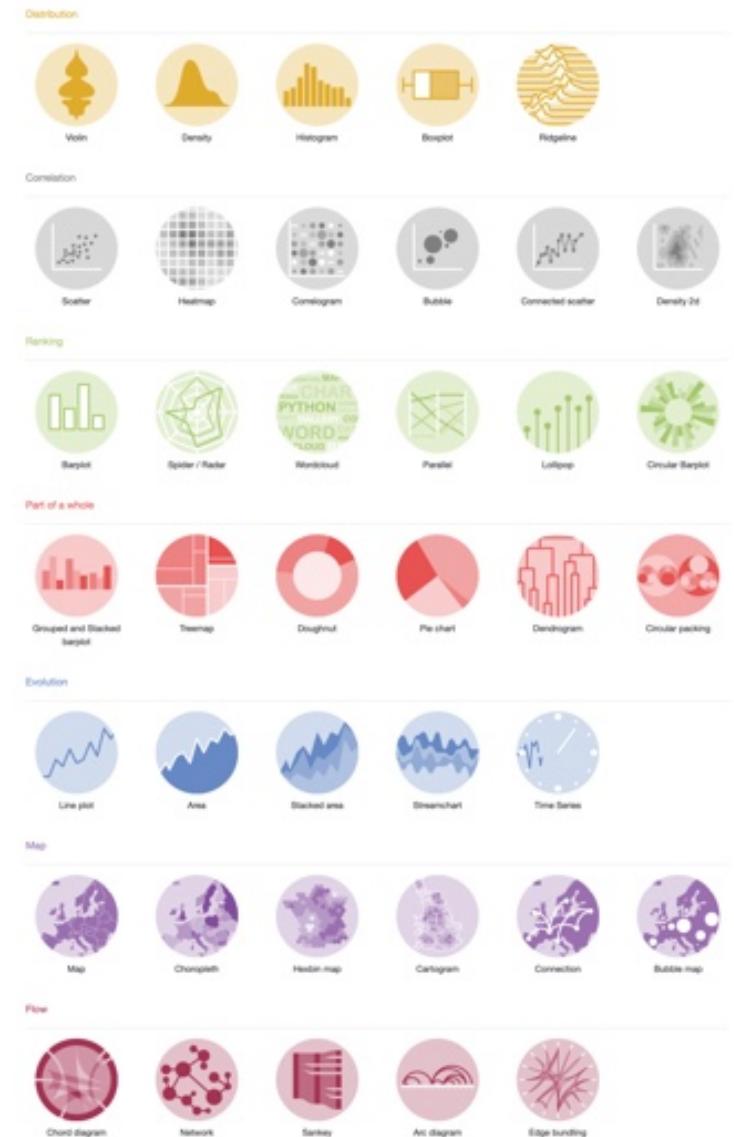


Position + Symbols + Color



<https://flowingdata.com/data-points/DataPoints-Ch3.pdf>

BASIC TYPES OF CHARTS



<https://r-graph-gallery.com>

SCATTERPLOT

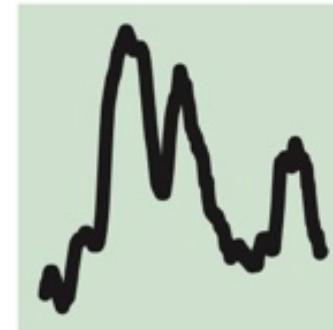
Scatterplot



The standard way to show the relationship between two continuous variables, each of which has its own axis.

LINE CHART

Line



The standard way to show a changing time series. If data are irregular, consider markers to represent data points.

HISTOGRAM

Histogram



The standard way to show a statistical distribution - keep the gaps between columns small to highlight the 'shape' of the data.

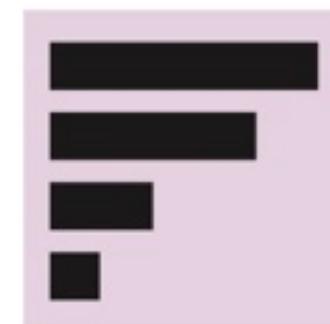
BAR CHART

Column

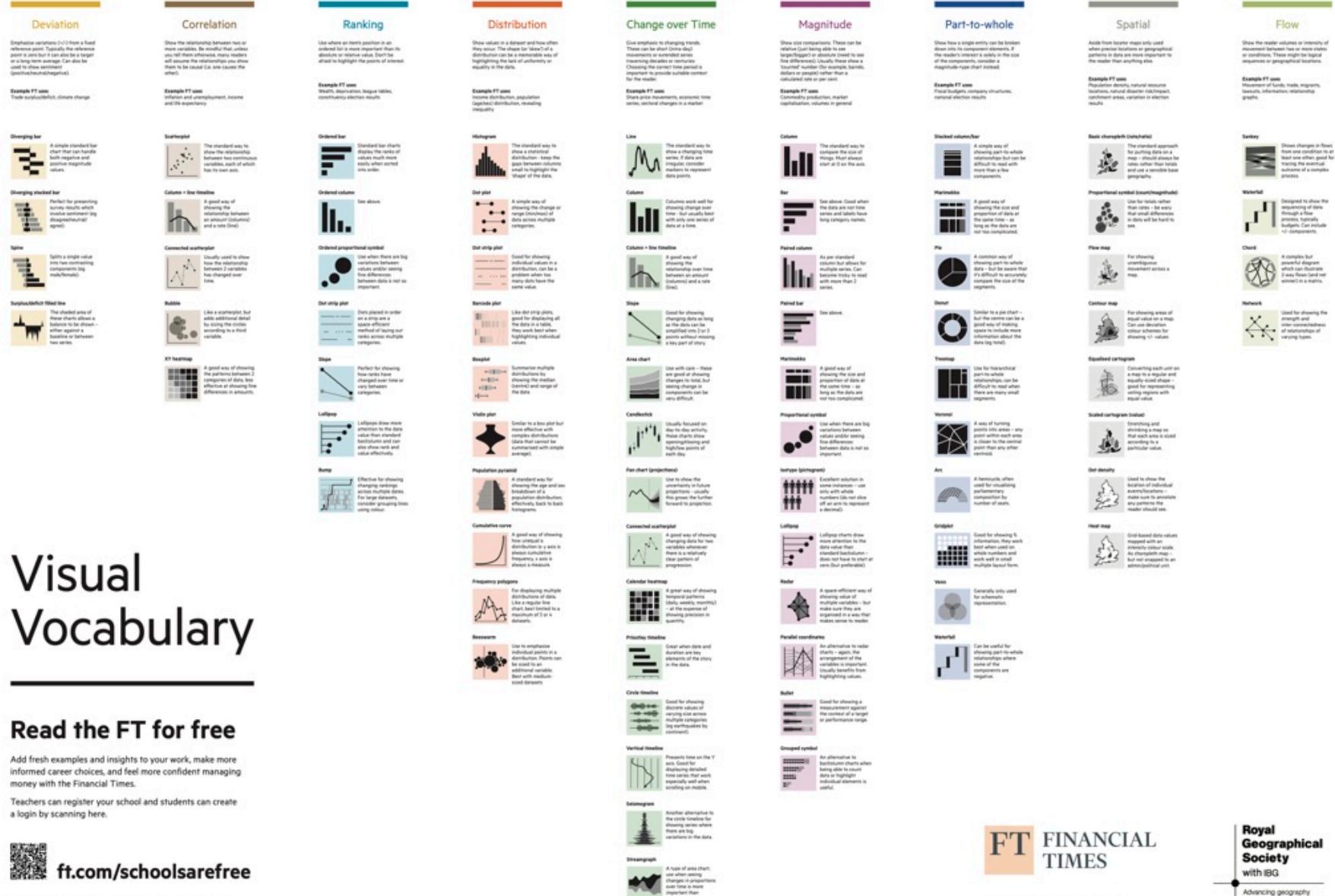


The standard way to compare the size of things. Must always start at 0 on the axis.

Bar



See above. Good when the data are not time series and labels have long category names.



BASIC GUIDELINES AND COMMON PITFALLS



Short communication

Ten guidelines for effective data visualization in scientific publications

Christa Kelleher Thorsten Wagener

Show more

Share Cite

<https://doi.org/10.1016/j.envsoft.2010.12.006>

Get rights and content

Ten guidelines for effective data visualization in scientific publications

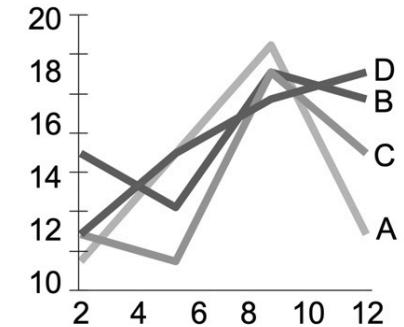
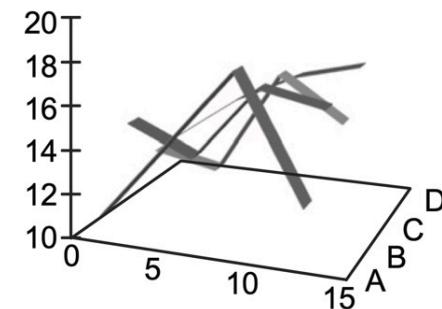
Abstract

Our ability to visualize scientific data has evolved significantly over the last 40 years. However, this advancement does not necessarily alleviate many common pitfalls in visualization for scientific journals, which can inhibit the ability of readers to effectively understand the information presented. To address this issue within the context of visualizing environmental data, we list ten guidelines for effective data visualization in scientific publications. These guidelines support the primary objective of data visualization, i.e. to effectively convey information. We believe that this small set of guidelines based on a review of key visualization literature can help researchers improve the communication of their results using effective visualization. Enhancement of environmental data visualization will further improve research presentation and communication within and across disciplines.

Christa Kelleher, Thorsten Wagener, Ten guidelines for effective data visualization in scientific publications, Environmental Modelling & Software, Volume 26, Issue 6, 2011, Pages 822-827, ISSN 1364-8152, <https://doi.org/10.1016/j.envsoft.2010.12.006>.

GUIDELINE 1

Create the simplest graph that conveys the information you want to convey



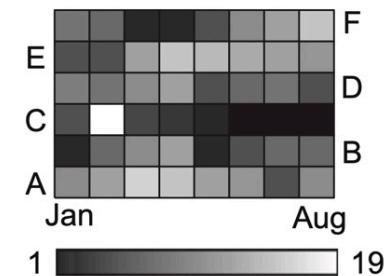
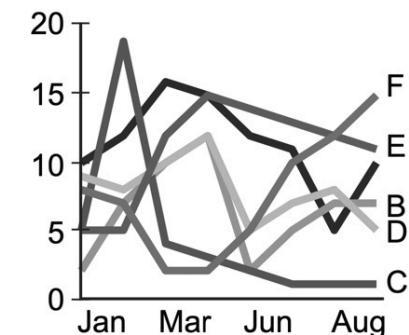
GUIDELINE 2

Consider the type of encoding object and attribute used to create a plot

Value encoding attribute			
	Length	Width	Orientation
Form			/
	Size	Shape	Curvature
	● ● ● ●		○ ○ ○ ○
	● ● ● ●	■	○ ○ ○ ○
Color	● ● ● ●		○ ○ ○ ○
	● ● ● ●	● ● ● ●	● ● ● ●
	Hue	Intensity	Transparency
	● ● ● ●	● ● ● ●	● ● ● ●
Spatial Position	2-D Position	Spatial Grouping	Density
	● ● ● ●	● ● ● ●	○ ○ ○ ○
Motion	Direction	Pathway	Pathway
	▲ ▲ ▲	● ● ●	○ ○ ○ ○

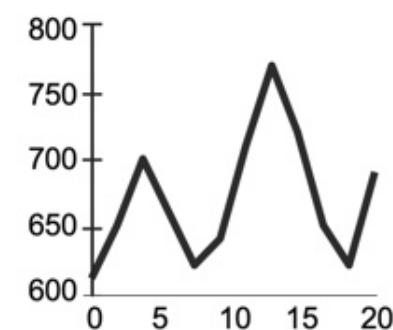
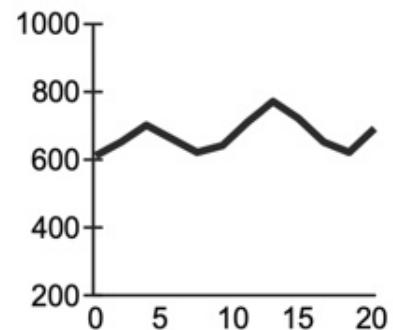
GUIDELINE 3

Focus on visualizing patterns or on visualizing details, depending on the purpose of the plot



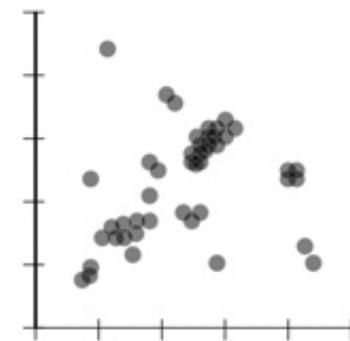
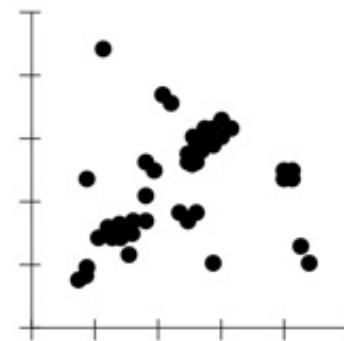
GUIDELINE 4

Select meaningful axis ranges



GUIDELINE 6

Plot overlapping points in a way that density differences become apparent in scatter plots



COMMON PITFALLS

and their alternatives

More at

<https://www.data-to-viz.com/caveats.html>

CAVEATS

A collection of dataviz caveats by data-to-viz.com

Show all Top 10 Improvement Misleading Map Bar



Order your data

When displaying the value of several entities, ordering them makes the graph much more insightful.



To cut or not to cut?

Cutting the Y-axis is one of the most controversial practice in data viz. See why.



The spaghetti chart

A line graph with too many lines becomes unreadable: it is called a spaghetti graph.



Pie chart

The human eye is bad at reading angles. See how to replace the most criticized chart ever.



Bar



Box plot



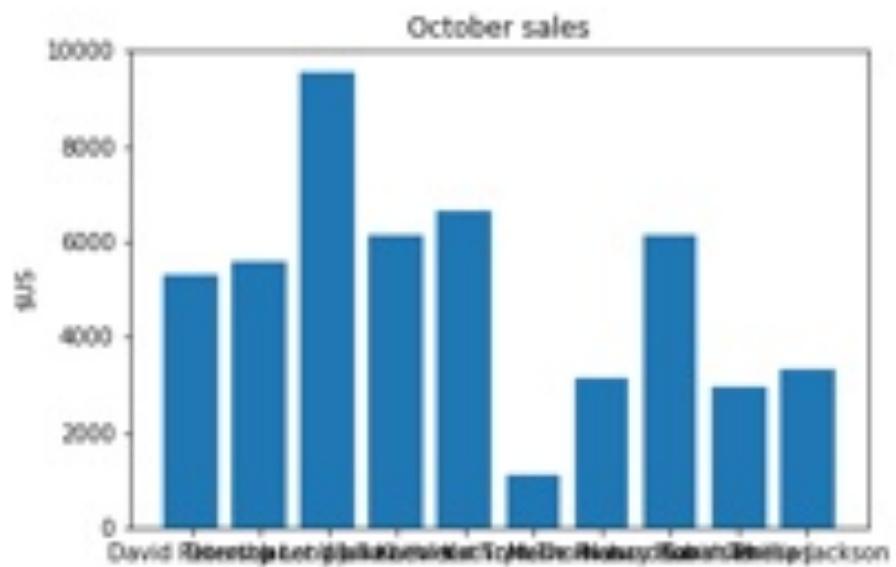
Histogram



Line chart

VERTICAL AXIS TEXT

Common problem

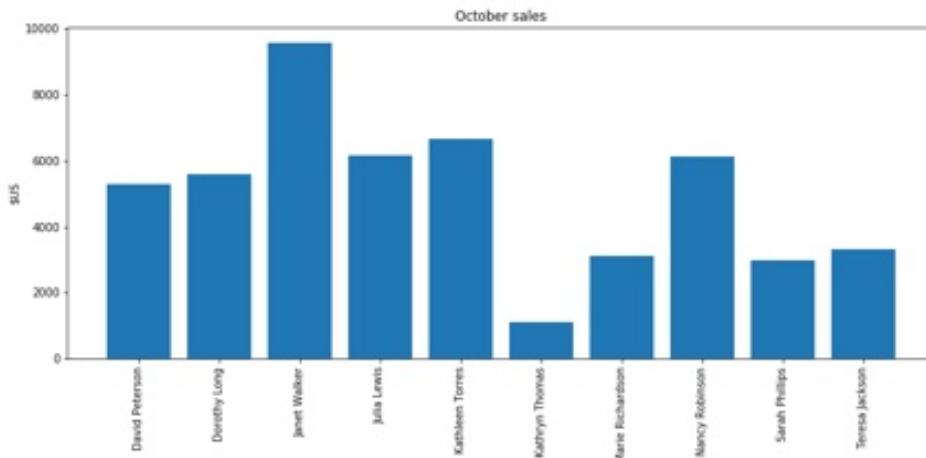


<https://gorelik.net/2017/11/23/how-to-make-a-graph-less-readable-rotate-the-text-labels/>

VERTICAL AXIS TEXT

Solved?

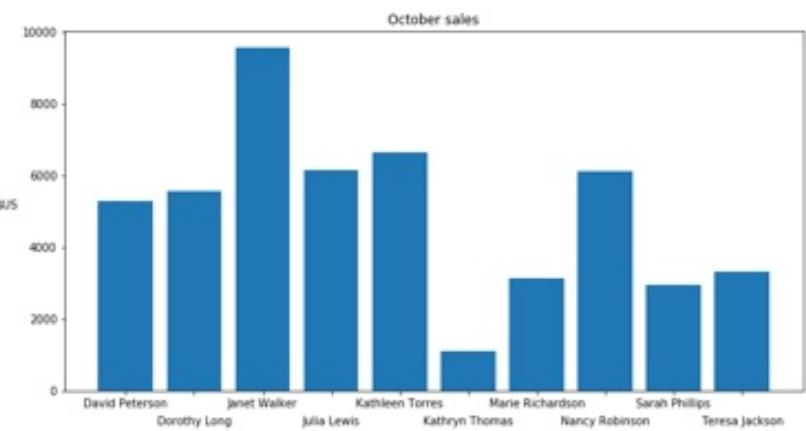
Difficult to read



<https://gorelik.net/2017/11/23/how-to-make-a-graph-less-readable-rotate-the-text-labels/>

VERTICAL AXIS TEXT

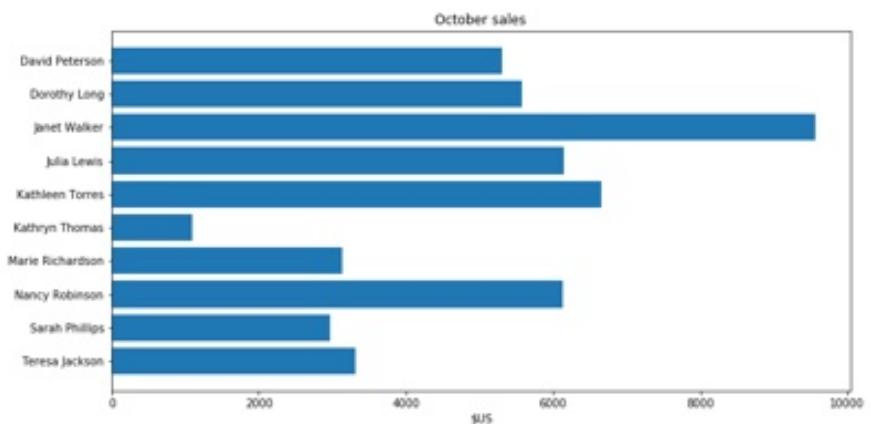
Consider moving the labels



<https://gorelik.net/2017/11/23/how-to-make-a-graph-less-readable-rotate-the-text-labels/>

VERTICAL AXIS TEXT

Rotate the bars

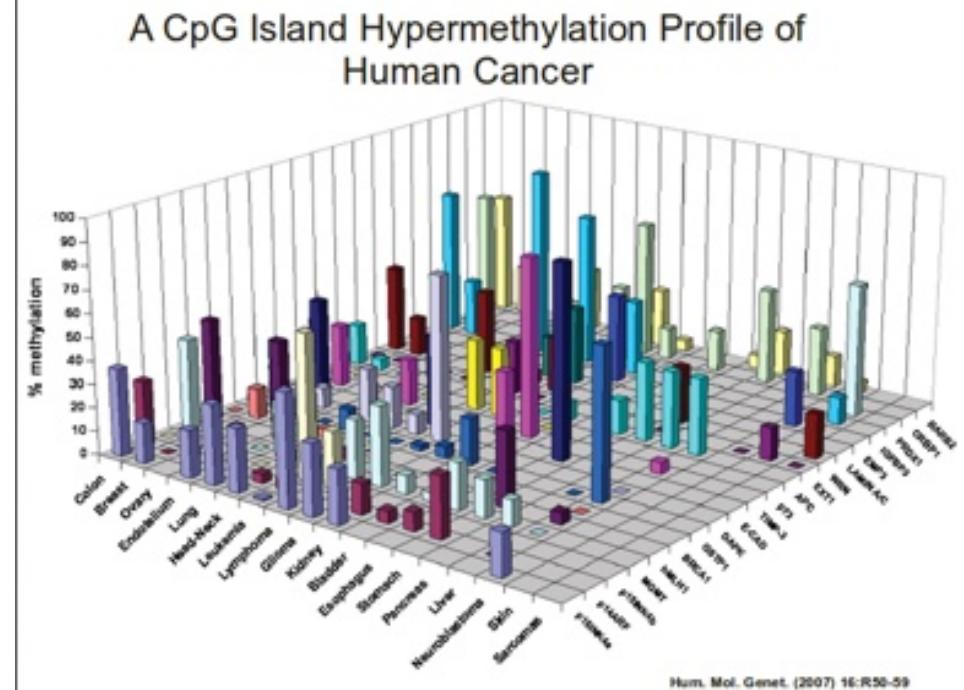


<https://gorelik.net/2017/11/23/how-to-make-a-graph-less-readable-rotate-the-text-labels/>

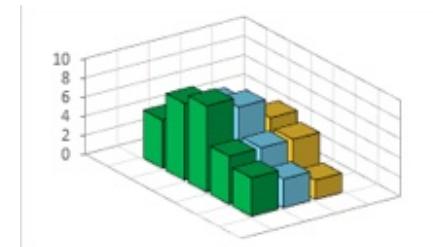
3D

Avoid

Very difficult to read and compare, hidden data

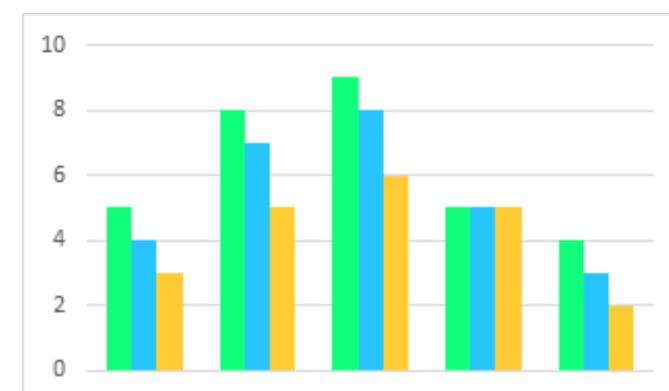


3D



<https://peltiertech.com/3d-bar-chart-alternatives/>

Consider other chart types



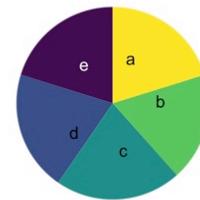
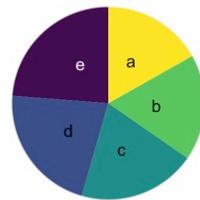
PIE CHARTS

Avoid

Perception problems (but not always)

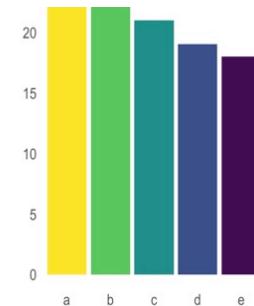
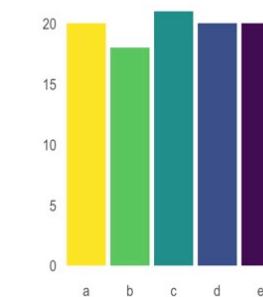
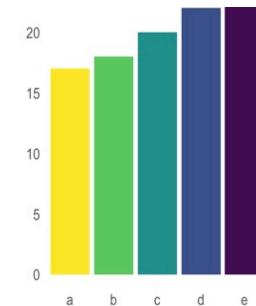
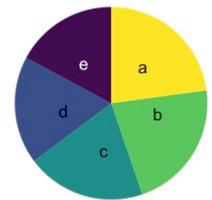
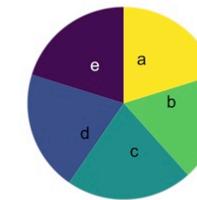
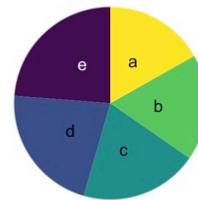
OK for few categories

Can be difficult to read



PIE CHARTS

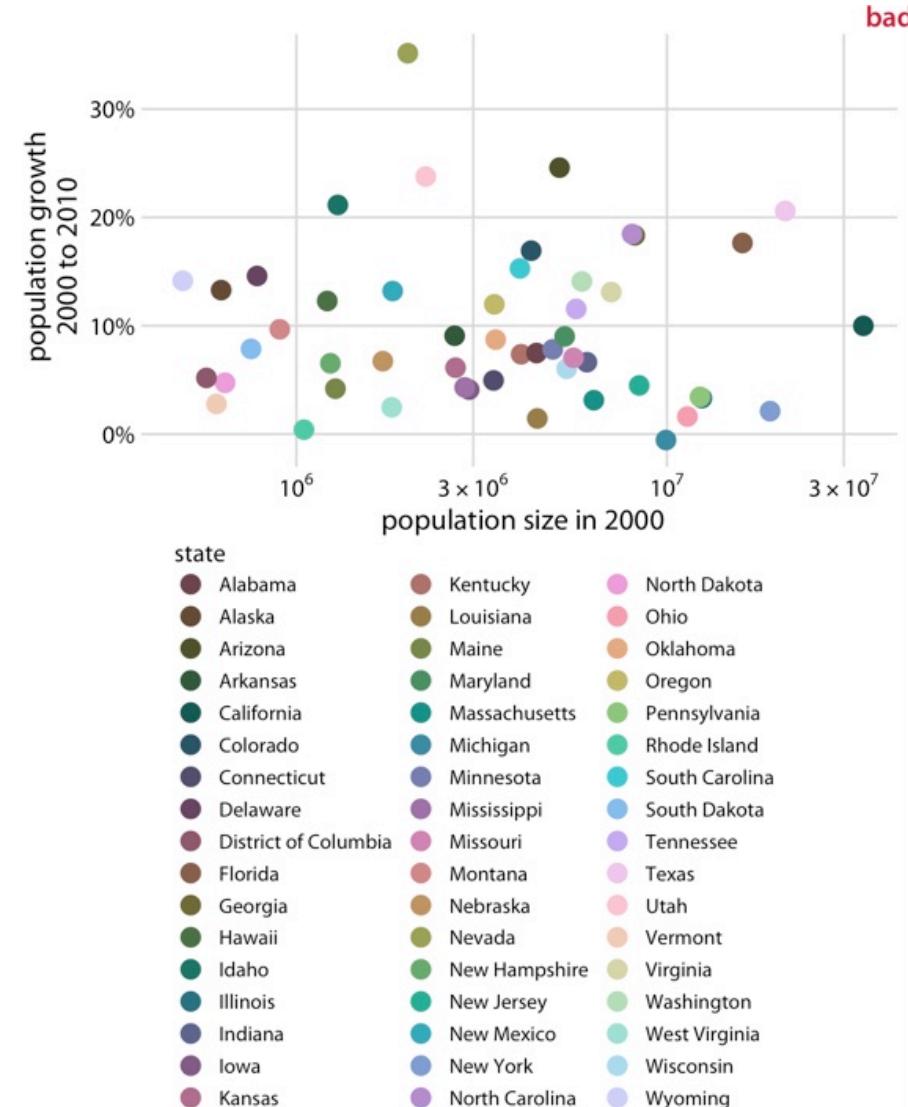
Alternative: bar chart



<https://www.data-to-viz.com/caveat/pie.html>

TOO MANY COLORS

OK to have less than 6-8 categories



<https://clauswilke.com/dataviz/color-pitfalls.html>

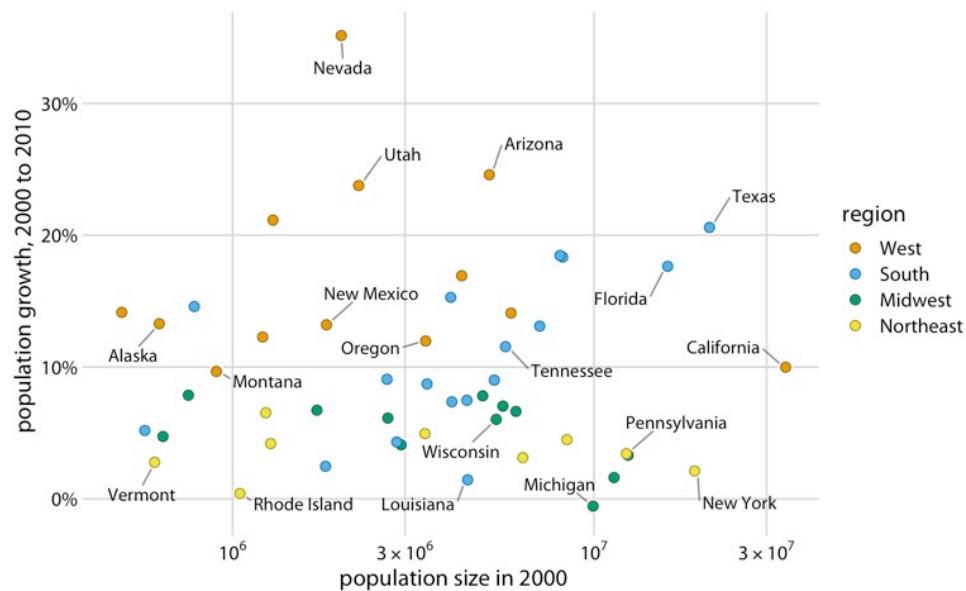
TOO MANY COLORS

Alternatives:

Direct labelling

Bar chart or other chart type

10 ways to use fewer colors in your data visualizations

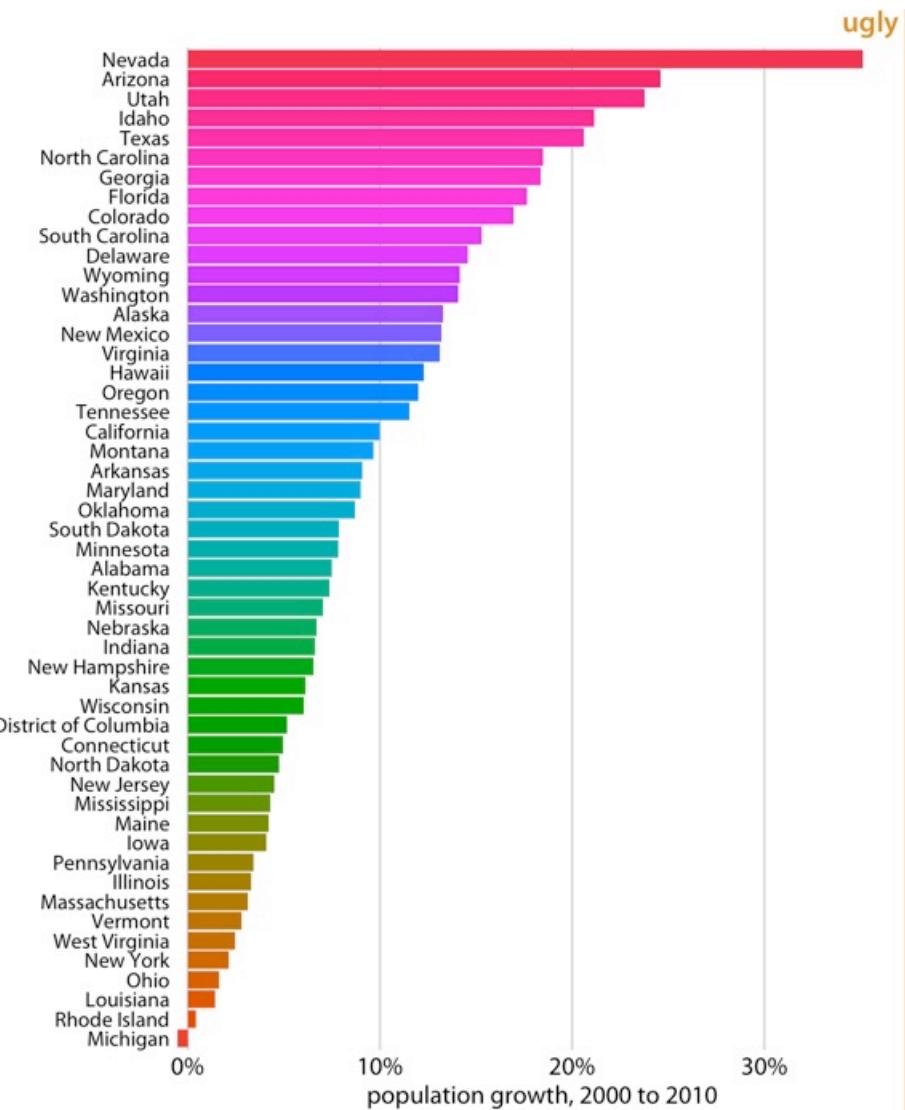


<https://clauswilke.com/dataviz/color-pitfalls.html>

UNNECESSARY COLOR

When color has no information

Use one color instead

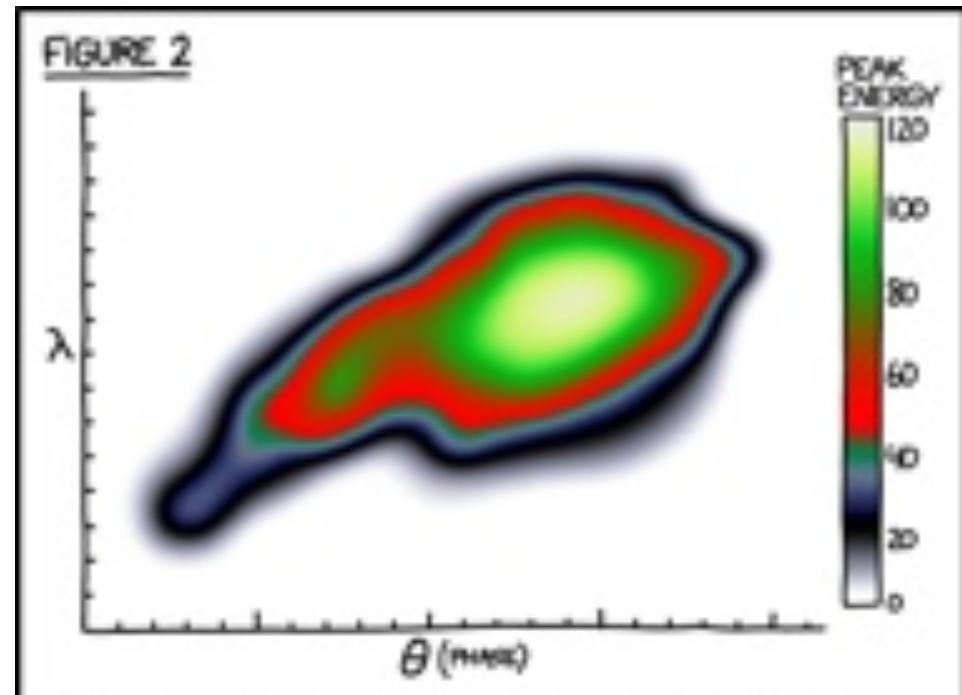


<https://clauswilke.com/dataviz/color-pitfalls.html>

RAINBOW COLOUR SCALES

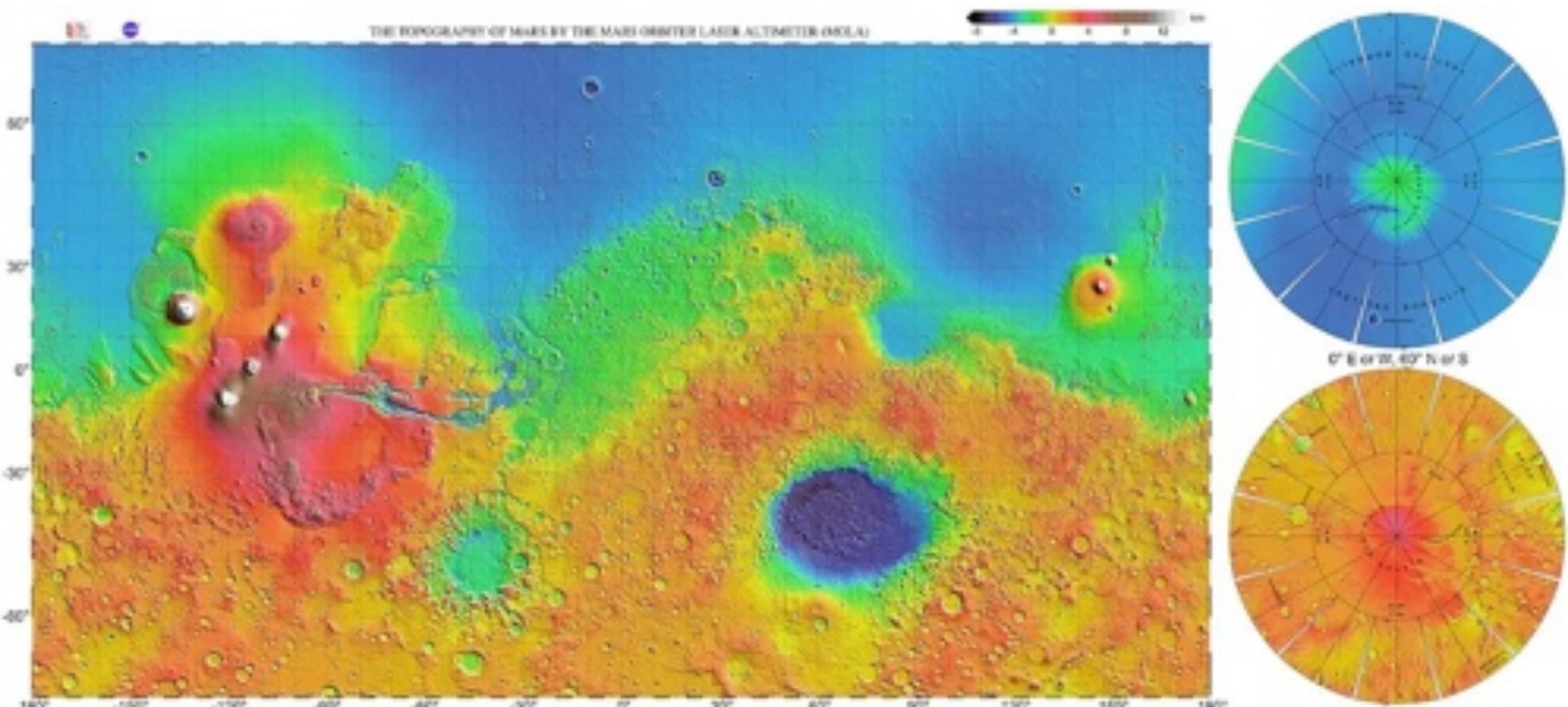
Caution when using the default rainbow colour scales

How rainbow colour maps can distort data and be misleading



EVERY YEAR, DISGRUNTLED SCIENTISTS COMPETE FOR THE RAINBOW AWARD FOR WORST COLOR SCALE.

<https://xkcd.com/2537/>

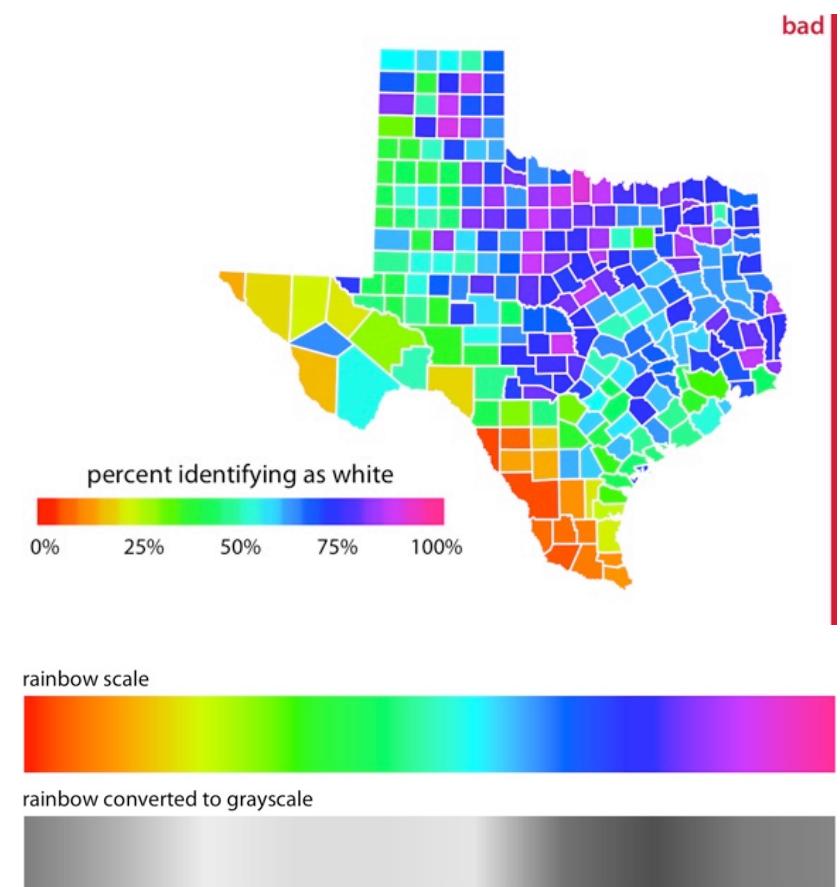


<https://attic.gsfc.nasa.gov/mola/images.html>

RAINBOW COLOUR MAPS

Difficult to compare values

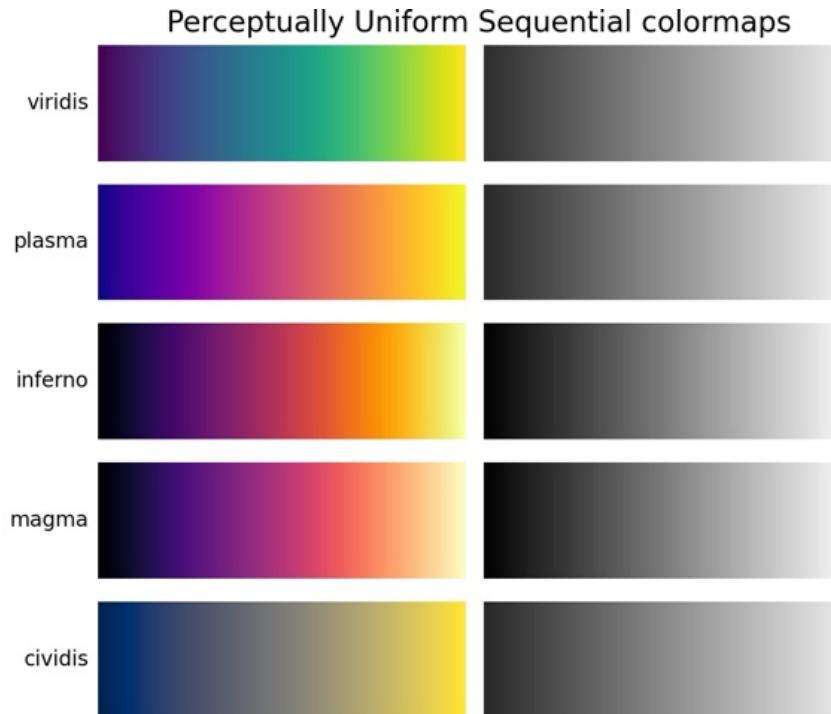
Uneven hues and lightness, misleading



<https://clauswilke.com/dataviz/color-pitfalls.html>

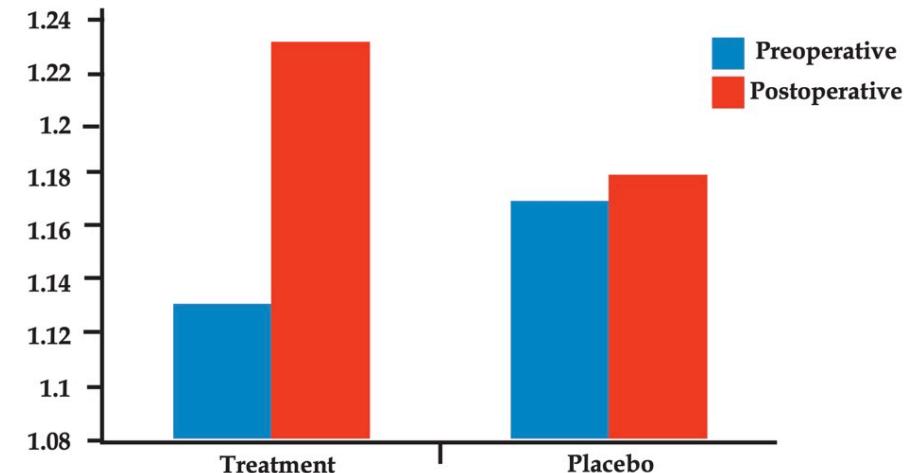
RAINBOW COLOUR SCALES

Use uniform palettes instead



<https://matplotlib.org/stable/tutorials/colors/colormaps.html>

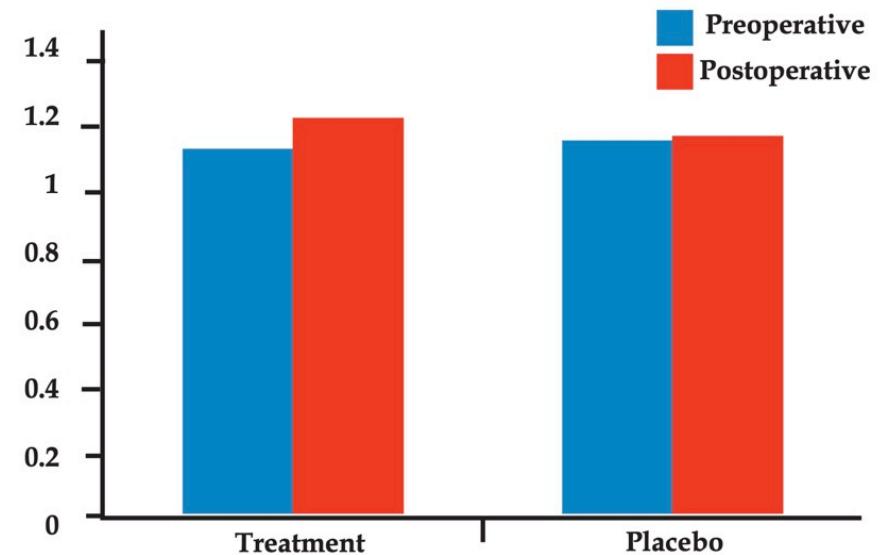
Y-AXIS DOES NOT START AT ZERO



Duquia RP, Bastos JL, Bonamigo RR, González-Chica DA, Martínez-Mesa J. Presenting data in tables and charts. An Bras Dermatol. 2014;89(2):280-5

Y-AXIS DOES NOT START AT ZERO

In bar charts y-axis should start at zero



Duquia RP, Bastos JL, Bonamigo RR, González-Chica DA, Martínez-Mesa J. Presenting data in tables and charts. An Bras Dermatol. 2014;89(2):280-5

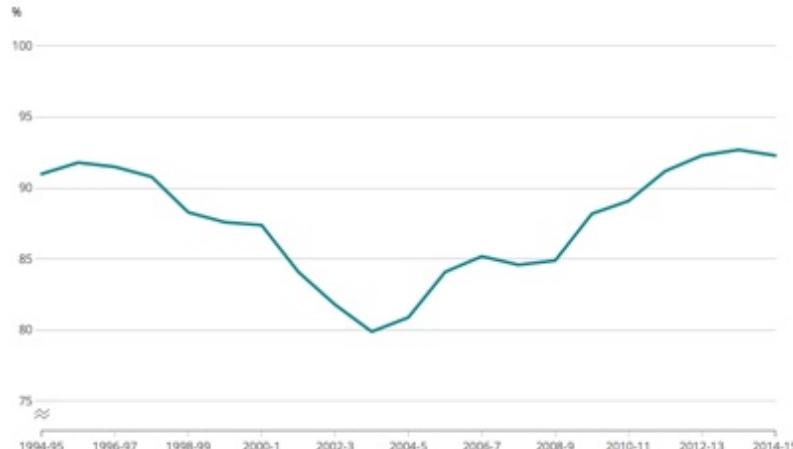
Y-AXIS DOES NOT START AT ZERO

No consensus about other charts

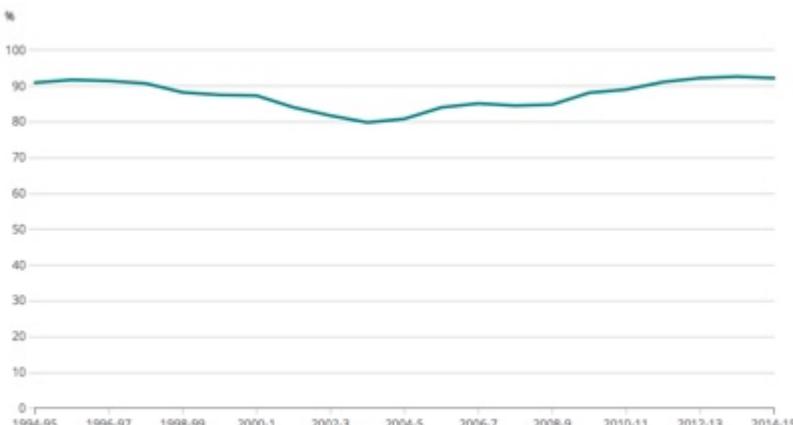
Line chart: same trends but small and important changes could be lost

Where to Start and End Your Y-Axis Scale

COMBINED MMR VACCINATION RATE, 1994-5 TO 2014-15,
ENGLAND



Source: NHS Immunisation Statistics - England, 2014-15, Table 8 and 9, HSCIC

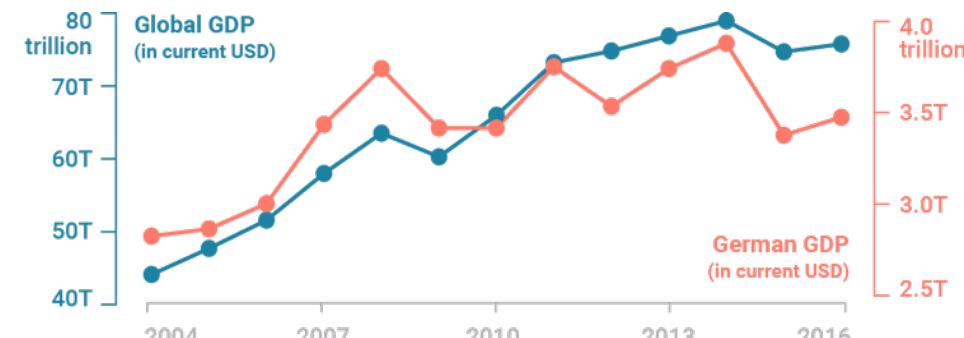


Source: NHS Immunisation Statistics - England, 2014-15, Table 8 and 9, HSCIC

<https://digitalblog.ons.gov.uk/2016/06/27/does-the-axis-have-to-start-at-zero-part-1-line-charts/>

DUAL Y-AXIS

Avoid



<https://blog.datawrapper.de/dualaxis/>

DUAL Y-AXIS

Difficult to read

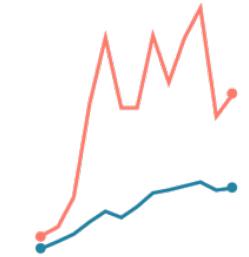
Depending on scaling and position, a correlation or other meaning is implied



Orange steady,
Blue massively increasing.



Blue steady,
Orange increasing.



Both started at the same
level, but Orange increased
far more than Blue.



Both started at the same
level, but Blue increased far
more than Orange.



Both started with the
same increase, then Blue
raced to the top.



Both steady.

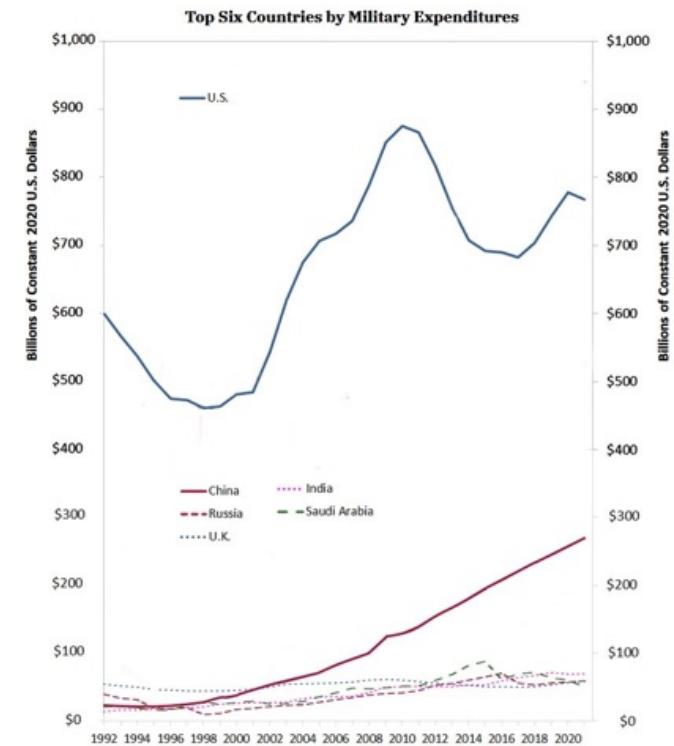
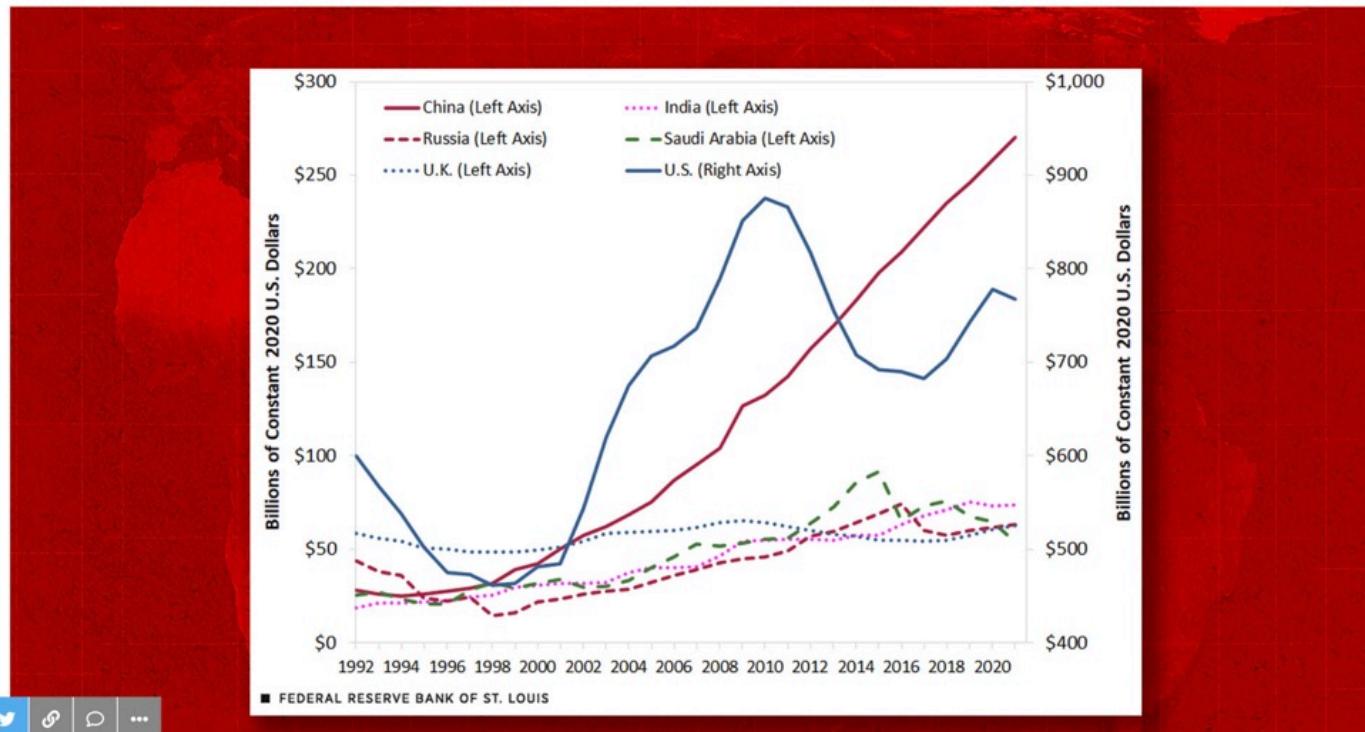
<https://blog.datawrapper.de/dualaxis/>

Exaggerating China's military spending, St. Louis Fed breaks all statistical rules with misleading graph

The Federal Reserve Bank of St. Louis published a jaw-droppingly misleading graph that portrays China as spending more on its military than the US. In reality, the Pentagon's budget is roughly three times larger.



By Ben Norton Published 2023-01-23



DUAL Y-AXIS

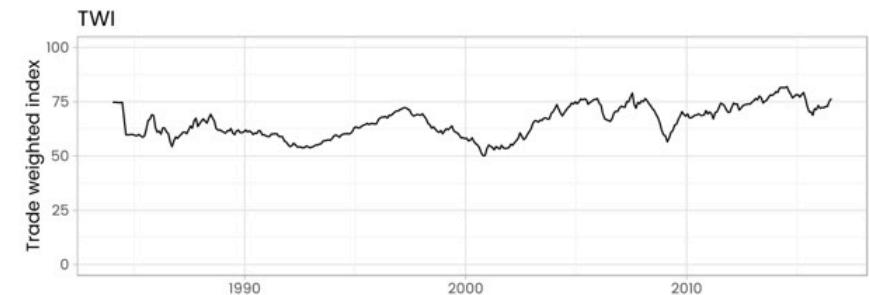
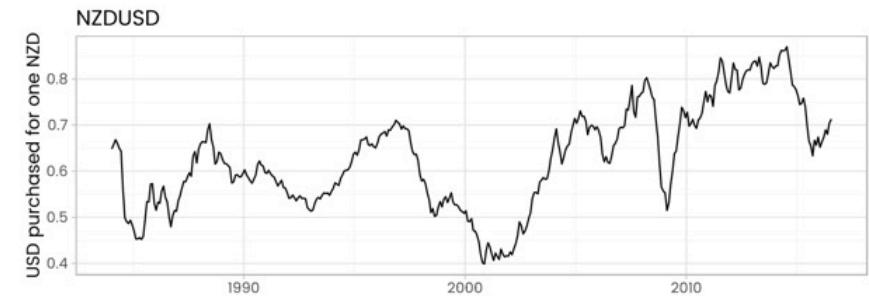
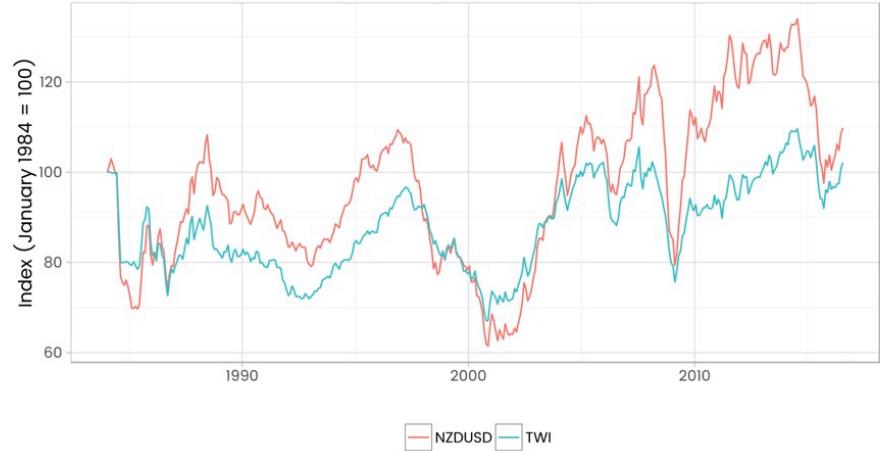
Alternatives:

Indexing

Side by side plots

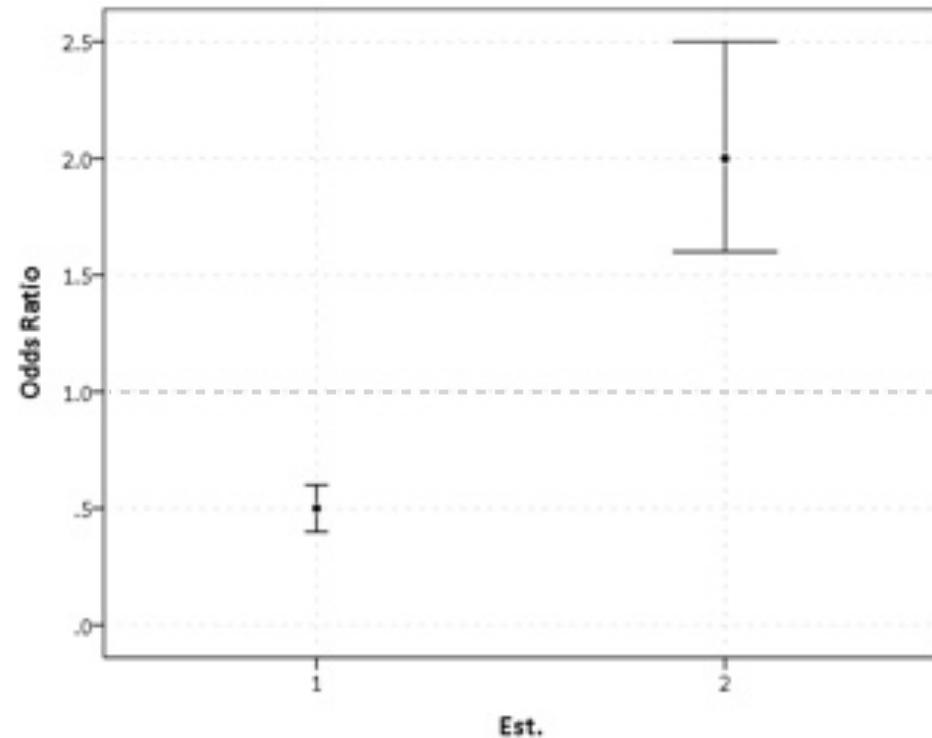
Why not to use two axes, and what to use instead

Usually accepted version of comparing two time series
Converted to an index, reference period first point in time



<http://freerangestats.info/blog/2016/08/18/dualaxes>

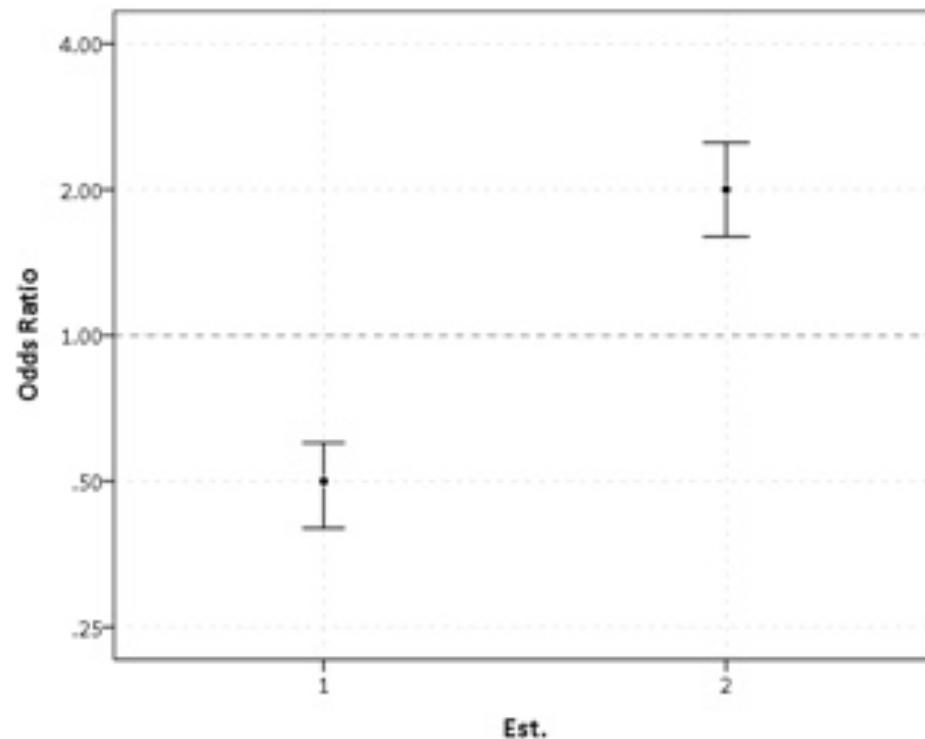
RATIOS THAT DON'T USE LOG AXIS



<https://andrewpwheeler.com/2013/10/26/odds-ratios-need-to-be-graphed-on-log-scales/>

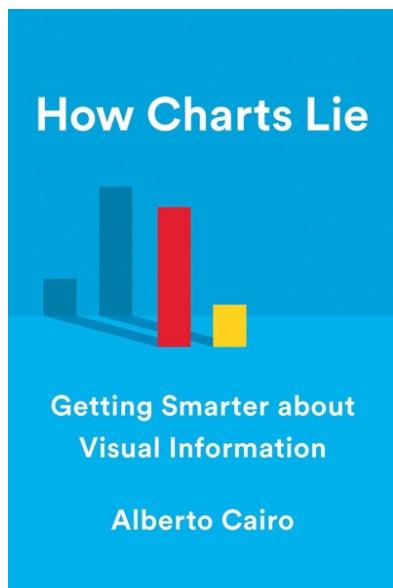
RATIOS THAT DON'T USE LOG AXIS

Use logarithmic scale



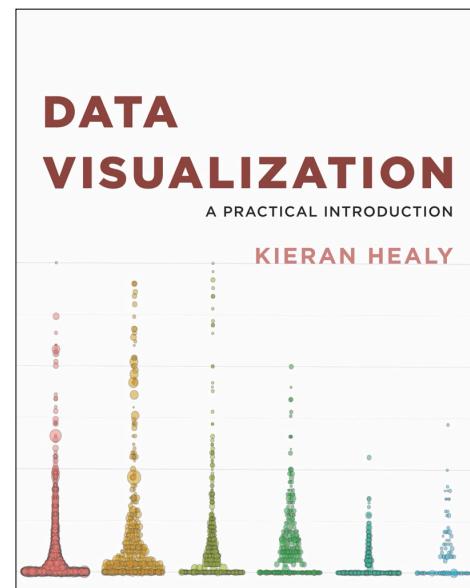
<https://andrewpwheeler.com/2013/10/26/odds-ratios-need-to-be-graphed-on-log-scales/>

RESOURCES



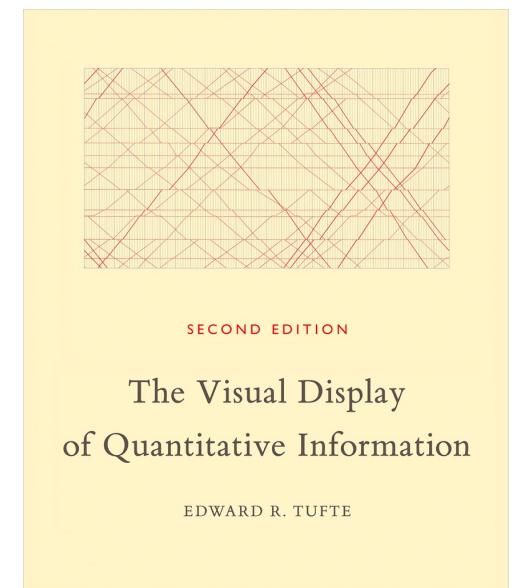
How charts lie: getting
smarter about visual
information

[UU library](#)



Data visualization: a practical
introduction

[UU library](#)
[Free online version](#)



The Visual Display
of Quantitative Information

[UU library](#)

[**https://blog.datawrapper.de**](https://blog.datawrapper.de) (theory and how-tos, showcase, stories)

[**https://nightingaledvs.com**](https://nightingaledvs.com) (tutorials, stories, industry-related news)

[**https://flowingdata.com**](https://flowingdata.com) (showcase)

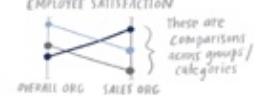
[**https://r-graph-gallery.com**](https://r-graph-gallery.com) (tutorials with code)

[**Aggplot2 tutorial for beautiful plotting in R**](#) (tutorial with code)



FIRST, LET'S REVIEW COMMON MYTHS in DATA VIZ

MYTH:
LINE GRAPHS are
for **CONTINUOUS**
DATA ONLY

The **LINES** that **CONNECT** the **POINTS** HAVE to MAKE SENSE
 Example: SURVEY DATA SLOPEGRAPH
 (a line graph w/
only the points)

 These are comparisons across groups/categories

MYTH:
BARS are
ALWAYS
BETTER

BARS are a **GOOD PLACE** to **START...** but **NOT ALWAYS THE BEST**
 ASK "what do I want my audience to see?"


Try other types of charts and decide what meets your needs

* **THIS IS TRUE** for **BAR CHARTS**



When **USING A PIE**, ASK yourself **WHY?**



Studies prove that people read pies & donuts by comparing **AREA**, not **ANGLE**

If you think pies make sense for your data & audience, test them out to see!

MYTH:
PIE CHARTS
ARE **EVIL**

MYTH:
UNBIASED
DATA EXISTS

We are **BIASING** our **DATA** at **EVERY STEP** of the **PROCESS**
RULE! DON'T LIE with DATA
 WHAT we choose to measure
 HOW we aggregate and compare
 HOW we show things

MYTH:
MORE DATA is
ALWAYS **BETTER**

BEFORE CHASING after **MORE DATA**,
 ask "WHAT will it **HELP** us **DO** or **DECIDE**?"
 Audience & context are important when it comes to the right amount of data

MYTH:
AVERAGES
ALWAYS WORK to
SUMMARIZE DATA

YOU NEED to UNDERSTAND the
DISTRIBUTION, SPREAD, and VARIABILITY



Averages can be misled by hiding a spread in a single number

MYTH:
THERE IS A SINGLE
RIGHT ANSWER when **VISUALIZING**
DATA

YOU SHOULD ALWAYS CONSIDER when
SHOWING DATA: **WHAT IS YOUR GOAL?**



**DO NOT
DECEIVE!**



 georgios.karamanis@neuro.uu.se

 <https://karaman.is>

 [geokaramanis](https://twitter.com/geokaramanis)

 <https://github.com/gkaramanis>