# Quantitative Tutorials – Session 2

# Data visualization
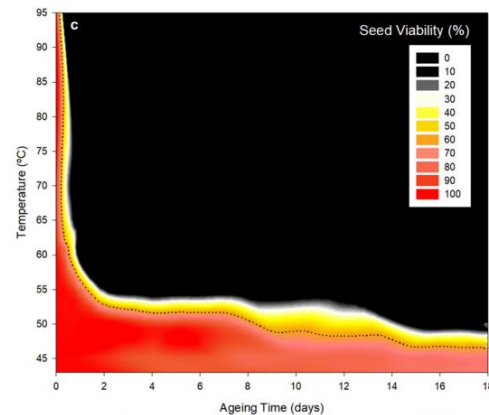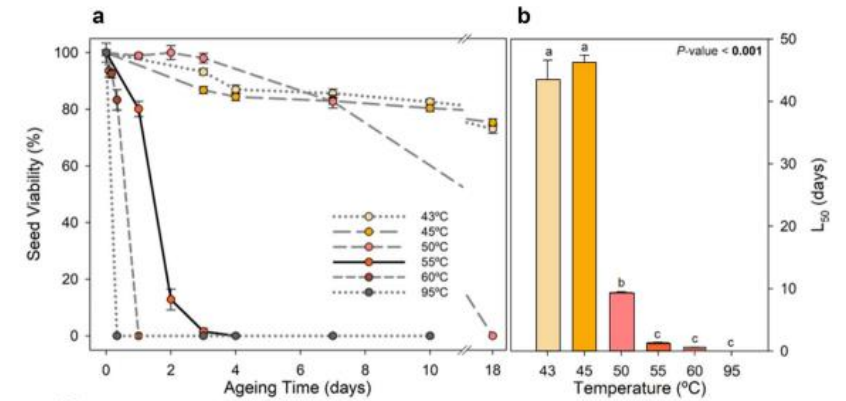## 23/11/22



Dr. Erola Fenollosa

# Data visualization objectives

Non-exclusive data visualization aims
- Communicate
- Explore
- Data in its context
- Find patterns and outliers

Visual analysis purposes

– Detect trends, patterns and outliers
– Compare
– Establish relationships

Visual codes:
position, shape, colour,
movement

But, the human eye has limitations!

# Visualizations for data analysis

## Alluvial Plot

## Spatial analysis

Today we know that cholera is spread through water, but in the early 1800s people weren't sure. John Snow's cholera map helped to show that contaminated wells were at the center of outbreaks. His research helped save countless lives and set the foundation for the field

# Good visualizations
– Simple and effective
– High information with low ink
– Intuitive (colours, axis)
– Honest



distribution of height of 2 durum wheat varieties

# CASE 1: What do you think about this visualization?



https://at.tumblr.com/badvisualisations/this-is-from-britannias-investor-presentation/o7hi1vwlyu96

# CASE 2: What do you think about this visualization?



**TECH STOCK THROWBACK?**
NETFLIX (JUN 2011-NOV 2013)   TESLA (JAN 2019-TODAY)

X axis missing

Do not manipulate the data to fit an idea. Use from-0 and proportional axis

# CASE 3: What do you think about this visualization?



**TOTAL SIXES**
ICC CRICKET WORLD CUP HISTORY

| | | | |
|---|---|---|---|
| ICC CWC 2019 | 186 SIXES | 21 MATCHES | 8.85 PER MATCH |
| ICC CWC 2015 | 463 SIXES | 48 MATCHES | 9.64 PER MATCH |
| ICC CWC 2011 | 258 SIXES | 49 MATCHES | 5.26 PER MATCH |
| ICC CWC 2007 | 373 SIXES | 51 MATCHES | 7.31 PER MATCH |
| ICC CWC 2003 | 266 SIXES | 52 MATCHES | 5.11 PER MATCH |
| ICC CWC 1999 | 153 SIXES | 42 MATCHES | 3.64 PER MATCH |
| ICC CWC 1995 | 148 SIXES | 36 MATCHES | 4.11 PER MATCH |
| ICC CWC 1992 | 93 SIXES | 39 MATCHES | 2.38 PER MATCH |
| ICC CWC 1987 | 126 SIXES | 27 MATCHES | 4.66 PER MATCH |
| ICC CWC 1983 | 77 SIXES | 27 MATCHES | 2.85 PER MATCH |
| ICC CWC 1979 | 28 SIXES | 14 MATCHES | 2.00 PER MATCH |

Mind the ink/data ratio

https://at.tumblr.com/badvisualisations/this-is-a-screengrab-from-star-sports/1p97g4l7r6hj

# CASE 4: What do you think about these visualizations?



THE WORLD CUP'S BIG GUNS
% OF TEAM'S RUNS SCORED BY TOP SCORER

- WILLIAMSON 30.23
- ROHIT 29.05
- SHAKIB 28.25
- WARNER 25.02
- BABAR 24.51
- DU PLESSIS 21.06
- POORAN 20.01
- ROOT 19.07
- KUSAL PERERA 18.16
- RAHMAT SHAH 14.8

espncricinfo

BOULT AND SOUTHEE'S SHARE OF TEST WICKETS AMONG NEW ZEALAND BOWLERS

TRENT BOULT — 23%, 24%, 26%, 27%

TIM SOUTHEE — 24%, 24%, 26%, 26%

- % OF WKTS
- % OF WKTS IN WINS
- % OF WKTS AT HOME
- % OF WKTS IN HOME WINS

**Pie charts are a bad choice**

**Mind the color**

https://www.espncricinfo.com/story/_/id/27143430/kane-williamson-hand-steadies-new-zealand-ship

http://www.thecricketmonthly.com/story/1181193

# CASE 5: What do you think about this visualization?



**The Hindi belt scores low, while the south does better**

Female labour force participation rate (%)

- 4.1 - 15.4%
- 15.5 - 19.9%
- 20 - 27.0%
- 27.1 - 33.5%
- 33.6 - 51.2%

Source: NSSO · Get the data · Created with Datawrapper

Do not be manipulative

# CASE 6: What do you think about this visualization?



The present economic situation of the country is.....

Good   So-so   Bad   No response

**Do not treat numerical as categorical**

# CASE 7: What do you think about this "pizza" visualization?



Add value with your visualization

Add legend if necessary

https://observablehq.com/@mbostock/2019-h-1b-employers

# CASE 8: What do you think about this visualization?



## 2017 This Is What Happens In An Internet Minute

- facebook — 900,000 Logins
- Google — 3.5 Million Search Queries
- NETFLIX — 70,017 Hours Watched
- $751,522 Spent Online
- 1.8 Million Snaps Created
- 15,000 GIFs Sent via Messenger
- 120 New Accounts Created (LinkedIn)
- 50 Voice-First Devices Shipped (amazon echo)
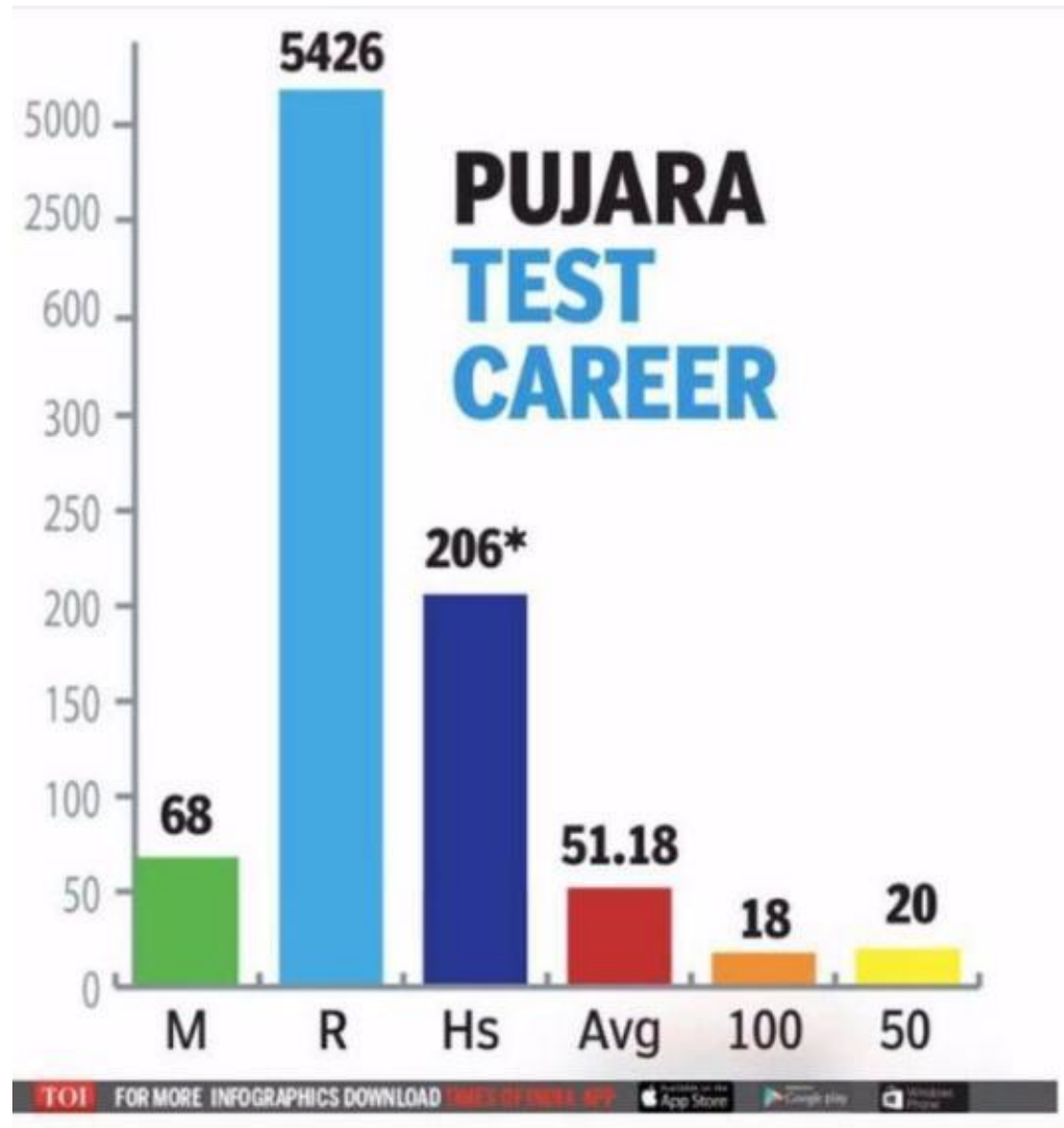- 40,000 Hours Listened (Spotify)
- 16 Million Text Messages
- YouTube — 4.1 Million Videos Viewed
- 342,000 Apps Downloaded
- 46,200 Posts Uploaded (Instagram)
- 452,000 Tweets Sent
- 990,000 Swipes (tinder)
- 156 Million Emails Sent
- Created By: @LoriLewis @OfficiallyChadd

## 2018 This Is What Happens In An Internet Minute

- facebook — 973,000 Logins
- Google — 3.7 Million Search Queries
- NETFLIX — 266,000 Hours Watched
- $862,823 Spent Online
- 2.4 Million Snaps Created
- 25,000 GIFs Sent via Messenger
- 38 Million Messages
- 67 Voice-First Devices Shipped (amazon echo)
- 936,073 Views (twitch)
- 18 Million Text Messages
- YouTube — 4.3 Million Videos Viewed
- 375,000 Apps Downloaded
- 174,000 Scrolling Instagram
- 481,000 Tweets Sent
- 1.1 Million Swipes (tinder)
- 187 Million Emails Sent

**Add value with your visualization**

# CASE 9: What do you think about this visualization?



Do not mix data types

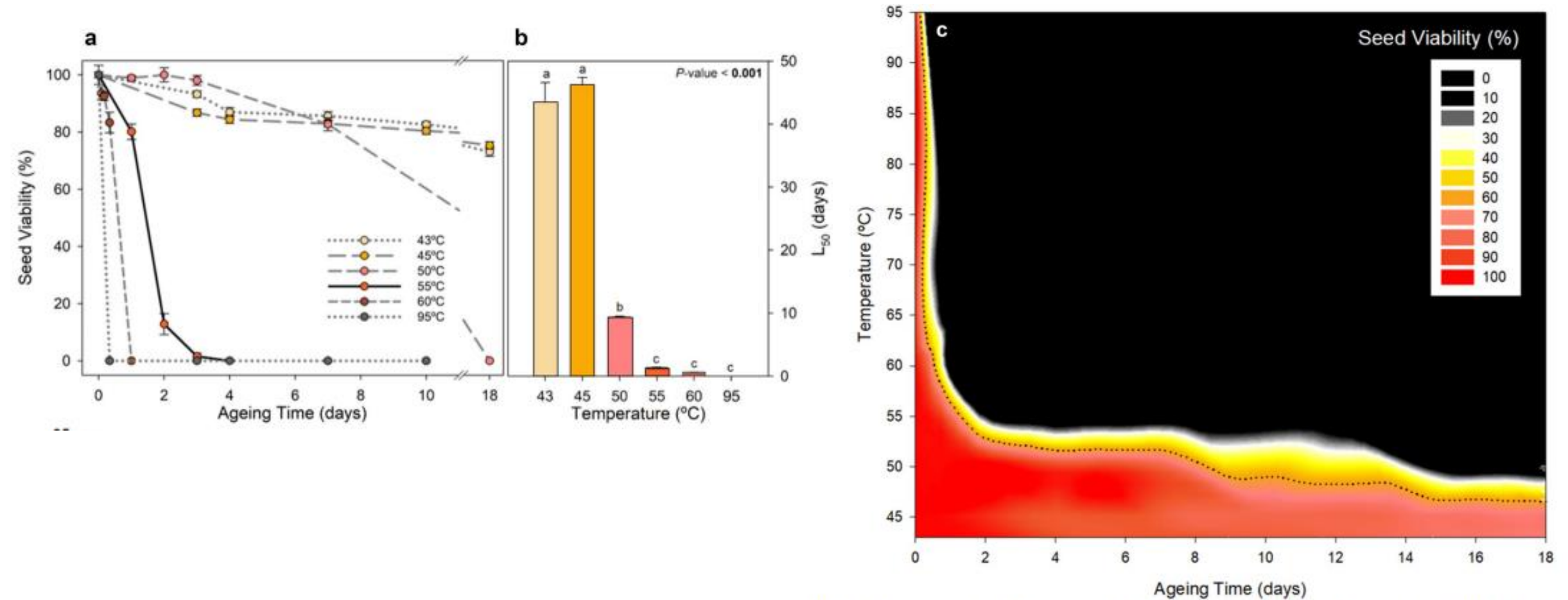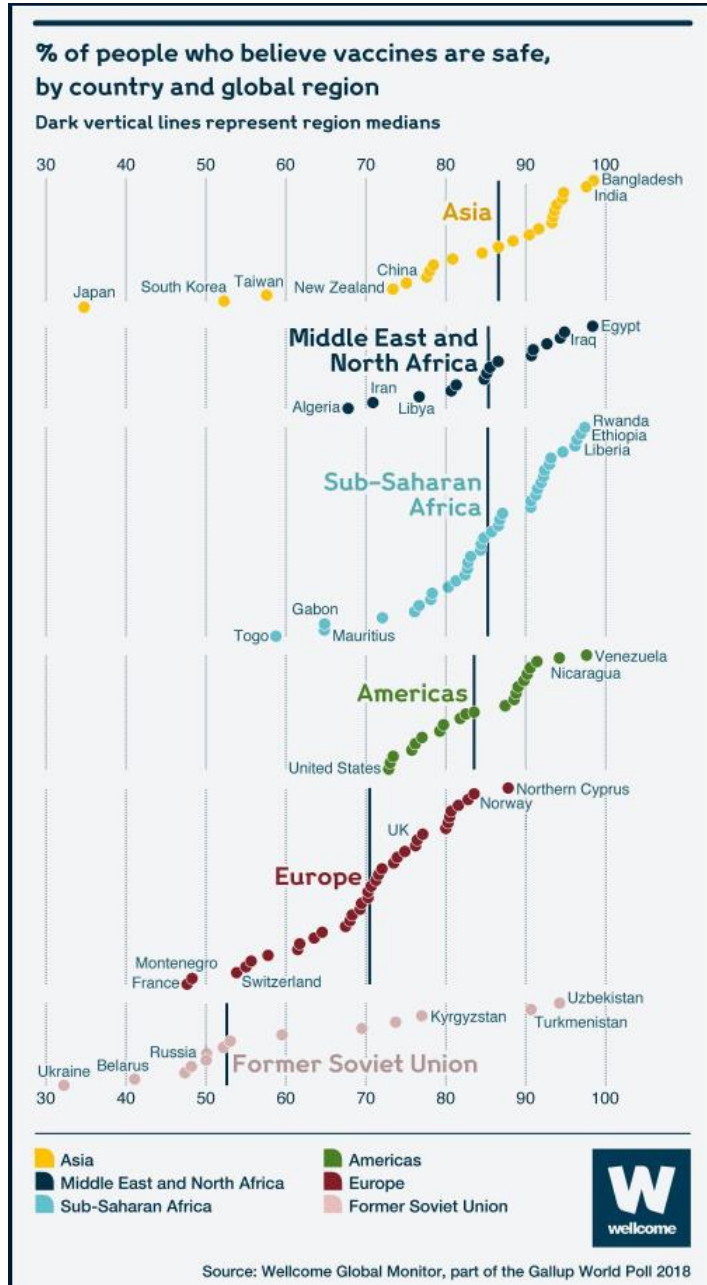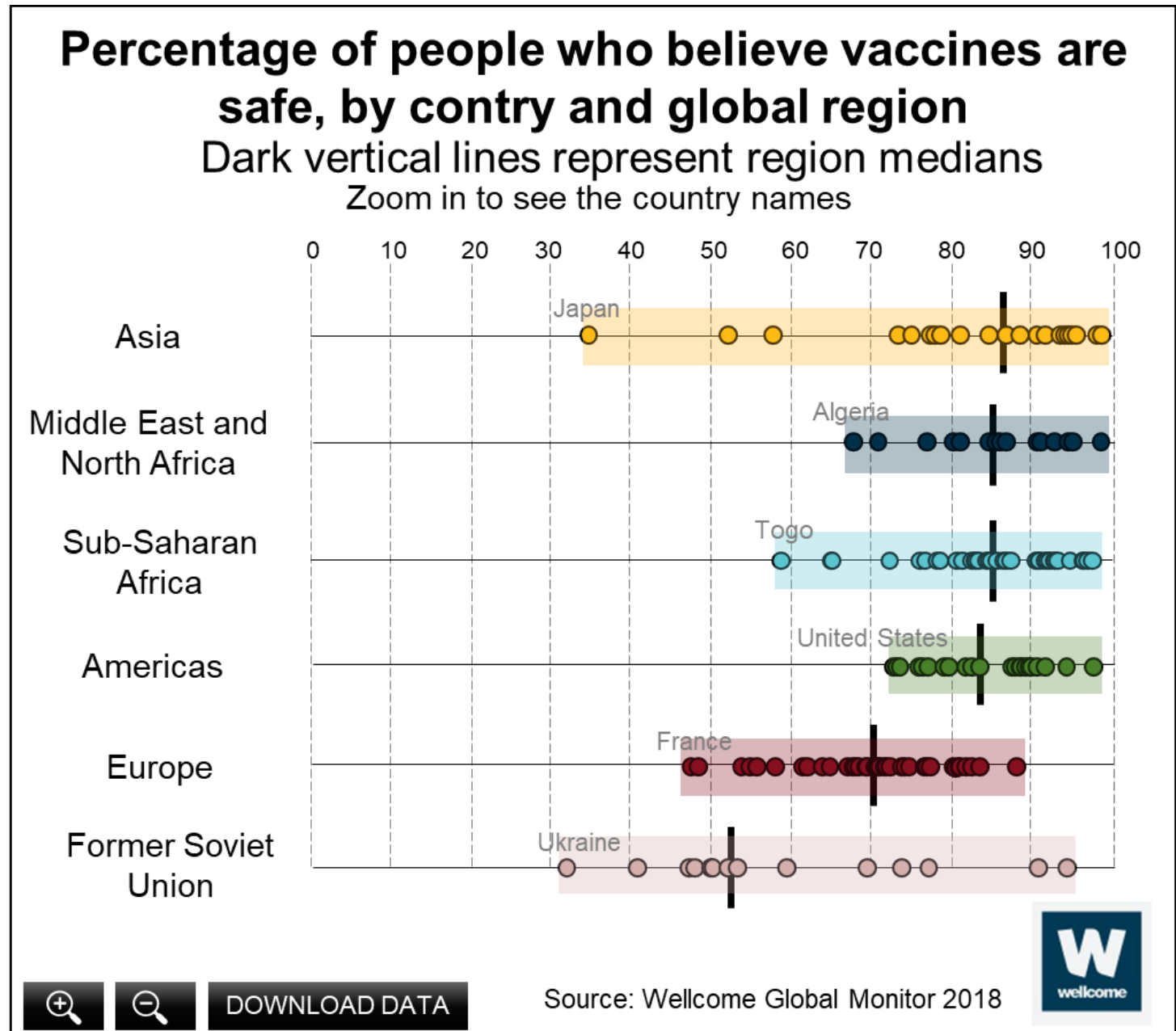# CASE 10: What do you think about these visualizations?



**Fig. 2** Influence of temperature and aging time in artificially aged seeds of *C. edulis*. **a** Viability loss of *C. edulis* seeds at the different tested temperatures and 87% RH. Data is represented as Mean ± SE (n = 6). **b** L$_{50}$ (loss of 50% viability) at the different tested temperatures. Different letters indicate statistically significant differences (*P*-value < 0.05). **c** Contour plot of seed viability considering temperature and aging days. Dotted line represents the L$_{50}$, 50% viability loss
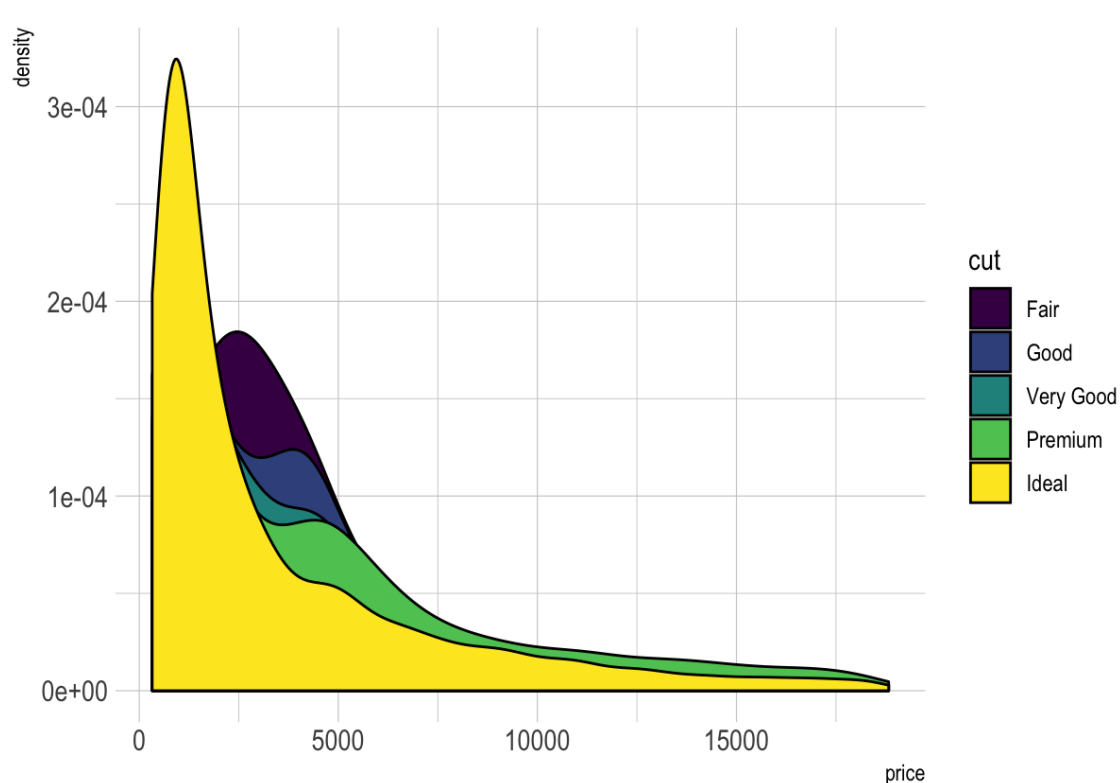
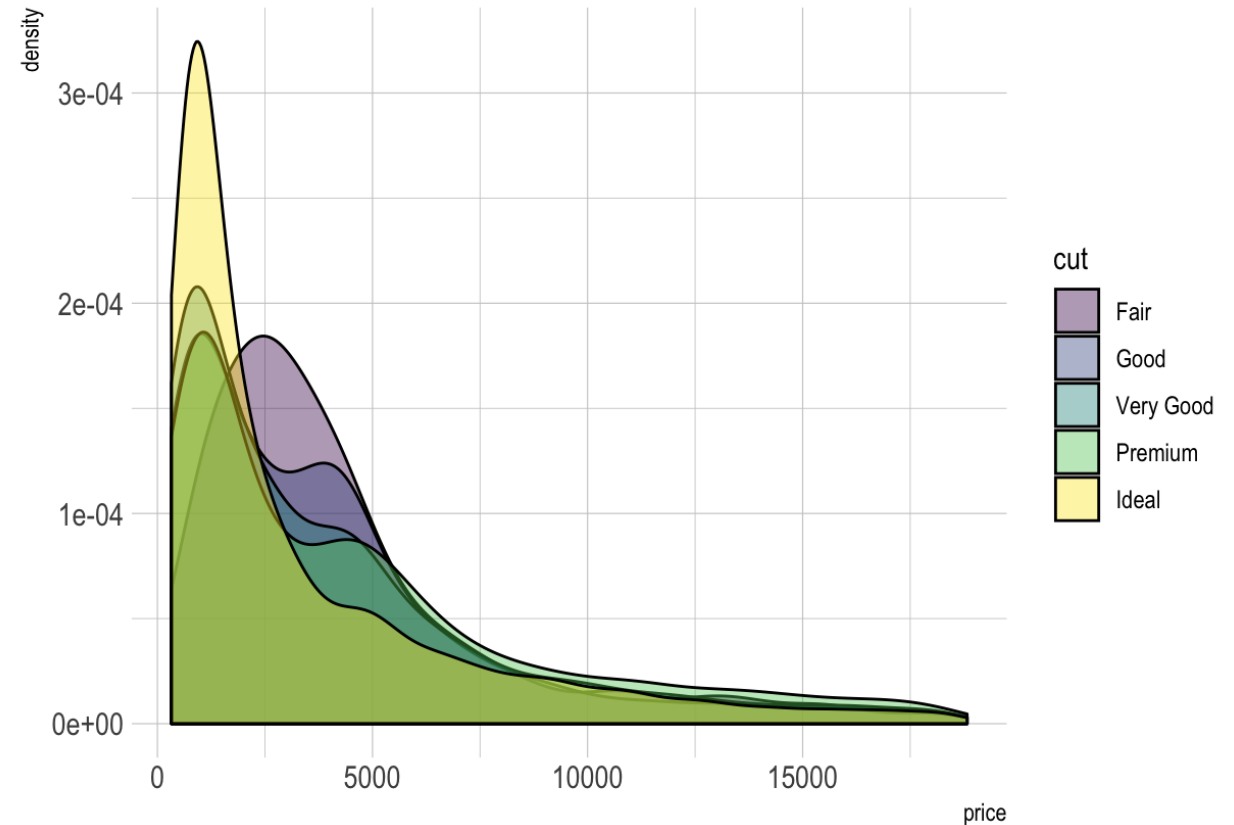Fenollosa et al., 2020 Plant Methods

# Transforming a visualization

**My proposal**

# What's the problem with this plot?

## Which is the best visualization for the diamont dataset?

# How would you improve now your own creations?

# Statistics supporting ideas

By now we are working with just two variables of interest (categorical and numerical or both numerical):

-  Both numerical: Correlation significance and adjustment to test relationships

- Numerical and categorical: ANOVA to compare groups

# Next sesion

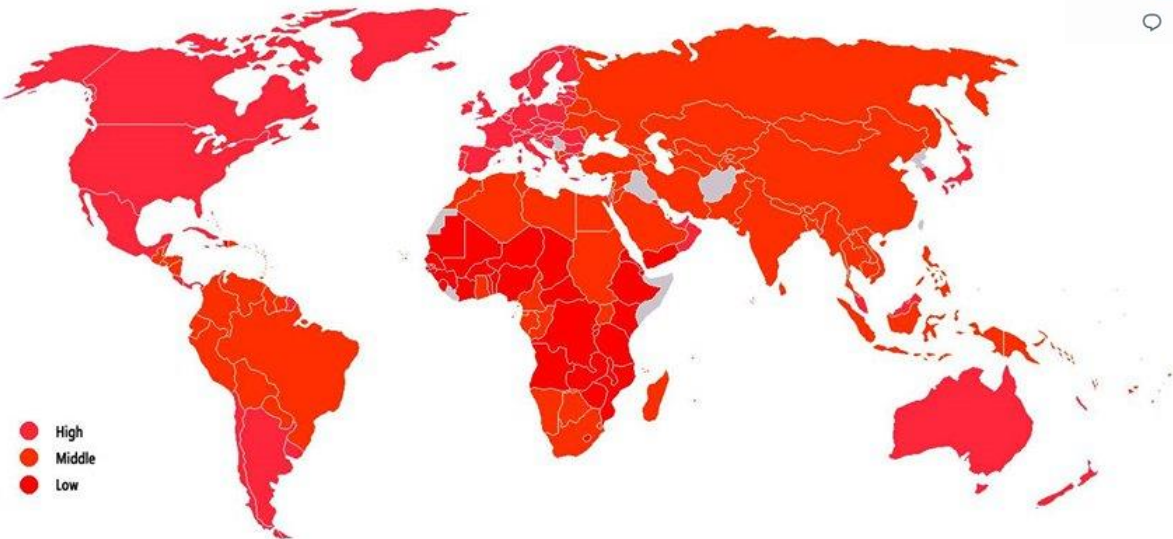We will cover one of the main methods for contrasting groups: ANOVA (linear models)

1. Choose two numerical and one categorical variable from your dataset

2. Think about a research question and writte it down

3. Test with a a linear model (lm, aov, glm, etc.)

    – Perform an ANOVA to compare categorical–numerical

    – Perform a Pearson correlation for the two numericals.

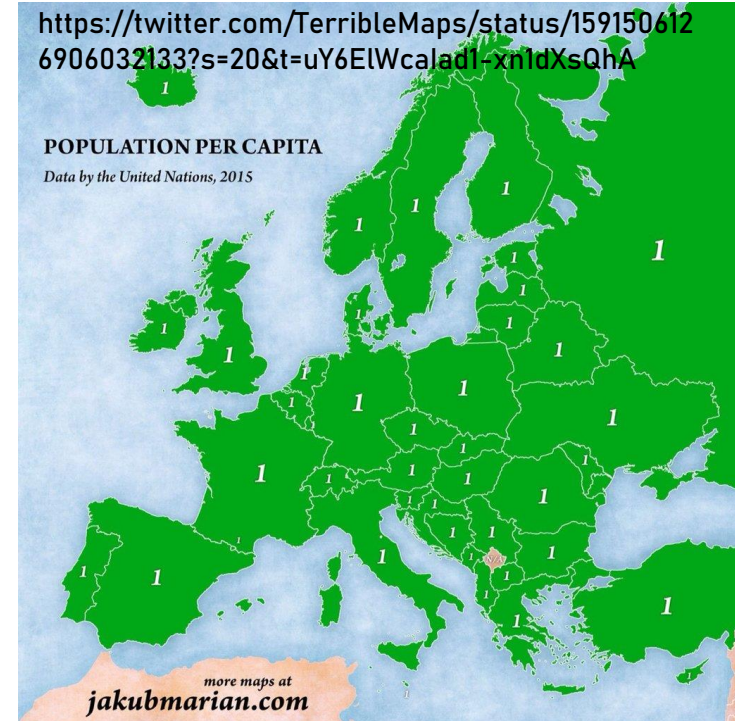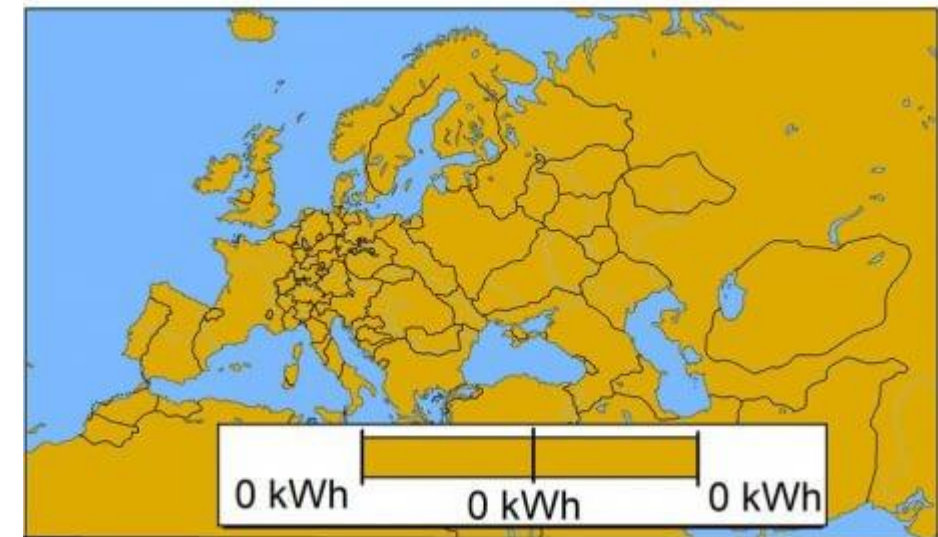Additional: Create a GitHub profile and start your portfolio

# See more:

- https://twitter.com/terriblemaps
- https://twitter.com/amazingmap
- https://www.connectedpapers.com/



https://twitter.com/TerribleMaps/status/159150612
6906032133?s=20&t=uY6ElWcalad1-xn1dXsQhA

POPULATION PER CAPITA

Data by the United Nations, 2015

more maps at
jakubmarian.com



Terrible Maps @TerribleMaps · 24 oct.
Locations Johnny Cash claims to have been in "I've been everywhere"

194        2.369        25,3 mil



# World Incidence of Color Blindness

High
Middle
Low

https://twitter.com/TerribleMaps/status/1592618803523117057?s=20&t=uY6El
Wcalad1-xn1dXsQhA



Electricity consumption in Europe in 1507

0 kWh        0 kWh        0 kWh

https://twitter.com/TerribleMaps/status/1583067072858181635?s=20&t=uY6ElWcalad1-xn1dXsQhA