# Week Six

## Open Science

### Philip Leftwich

### 25.10.2021

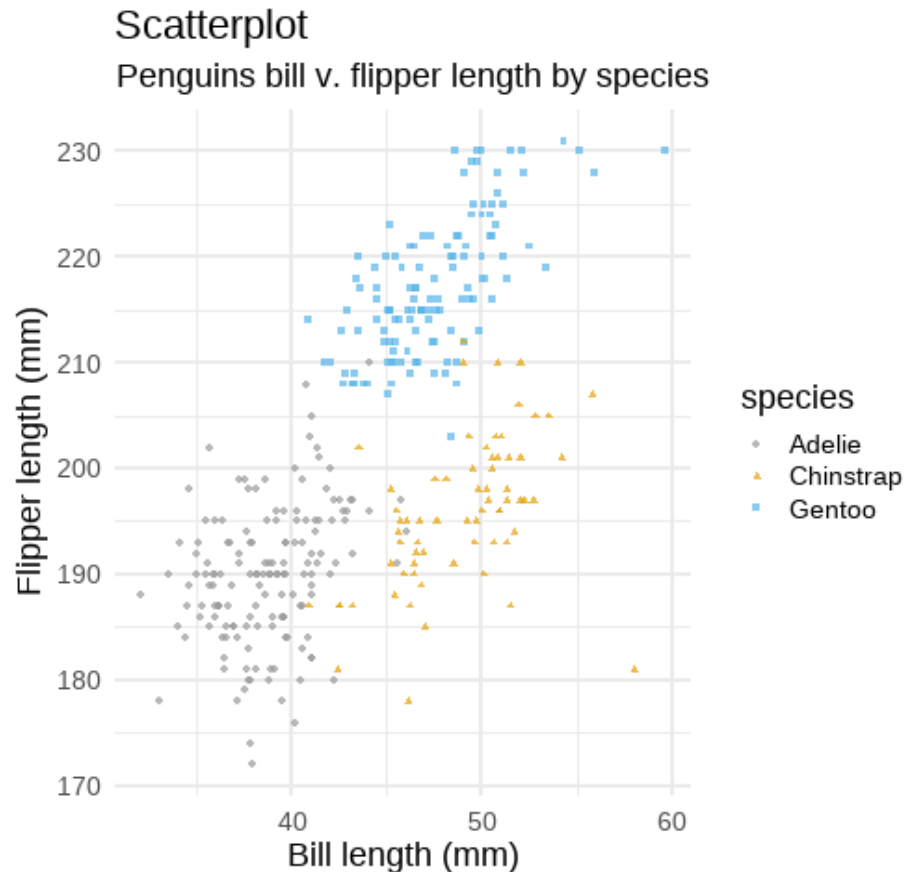# Hello! 👋

- Let me know how you are today in the Q&A on today's Slido!

- Go to Slido.com #944385

# Good data visualisations

```r
penguins %>%
  remove_missing() %>%
  ggplot(aes(x = bill_length_mm, y = f
        color = species, shape = specie
  geom_point(alpha = 0.7) +
  labs(x = "Bill length (mm)",
     y = "Flipper length (mm)",
    title = "Scatterplot",
    subtitle = "Penguins bill v. flipper le
    caption = "Source: https://github.
```

Slido.com #944385

## Scatterplot
### Penguins bill v. flipper length by species



Source: https://github.com/allisonhorst/palmerpenguins
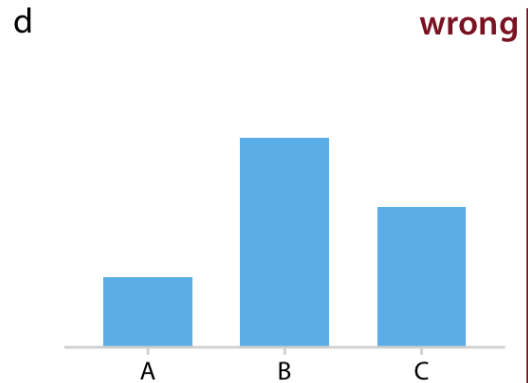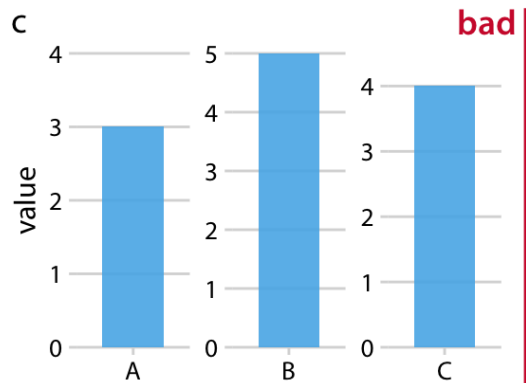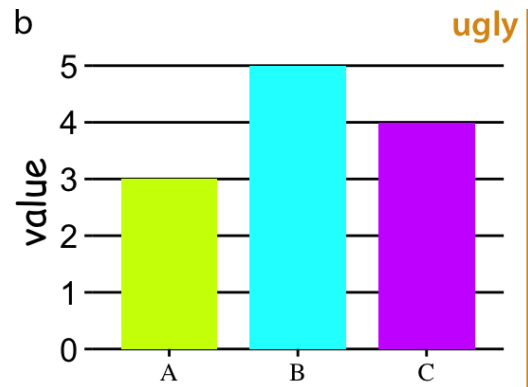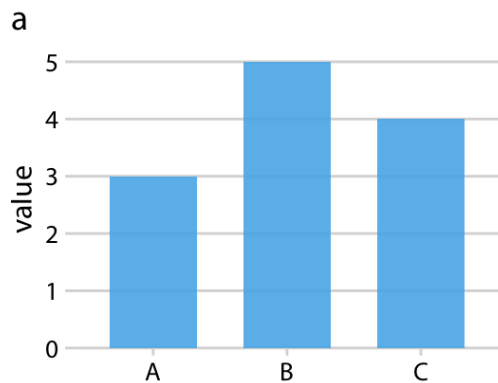
# Requirements of a good graph

- Visualisations must accurately reflect the data

- Tell a story

- Look professional

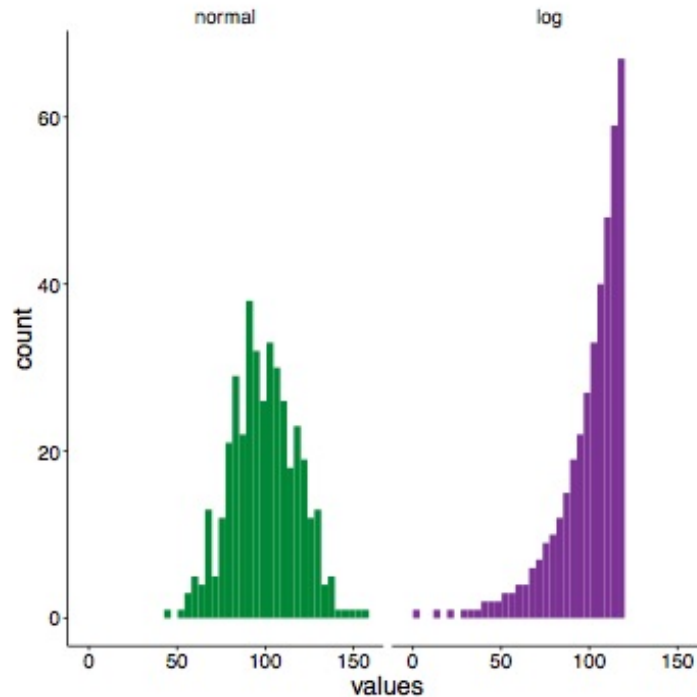# The Good, the Bad & the Ugly

# Choosing the right visual for your data

# Choosing the right visual for your data

# Choosing the right visual for your data

# Choosing a data visual

- Choosing the right data visual *requires* understanding your data

- You **must** clearly explain any non-obvious features

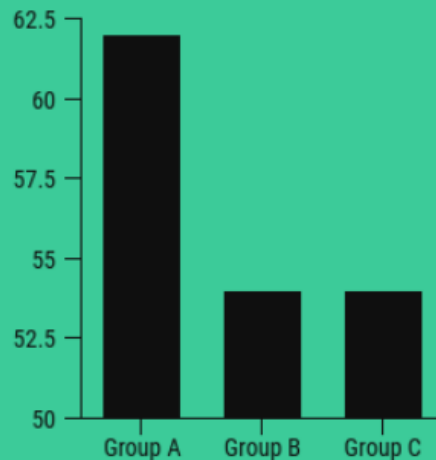- We will cover the different types of figures over the next few weeks

📊

# Five common ways graphs can mislead you
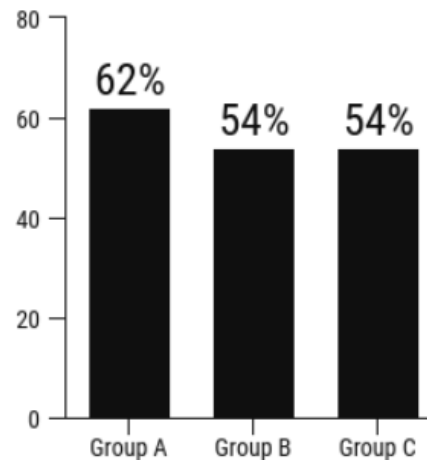
📊

# OMITTING THE BASELINE

**1** ▶

In most cases, the baseline for a graph is 0. But writers can skew how data is perceived by making the baseline a different number. This is known as a "truncated graph".



## ☹ MISLEADING

VS

## ACCURATE ☺

- Starting the vertical axis at 50 makes a small difference between groups seem massive

- Group A looks much larger than Groups B and C
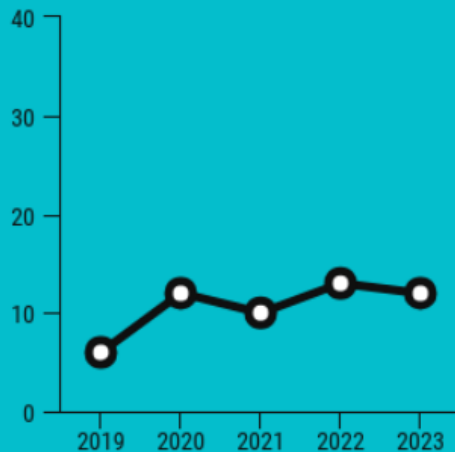
- Starting the vertical axis at 0 offers a more accurate depiction of the data

- The difference between the groups does not seem as dramatic
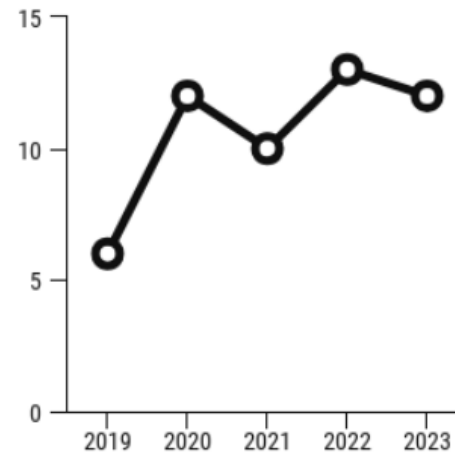
## 2 ▸ MANIPULATING THE Y-AXIS

Expanding or compressing the scale on a graph can make changes in data seem more or less significant than they actually are.

☹ **MISLEADING**

**VS**

**ACCURATE** ☺

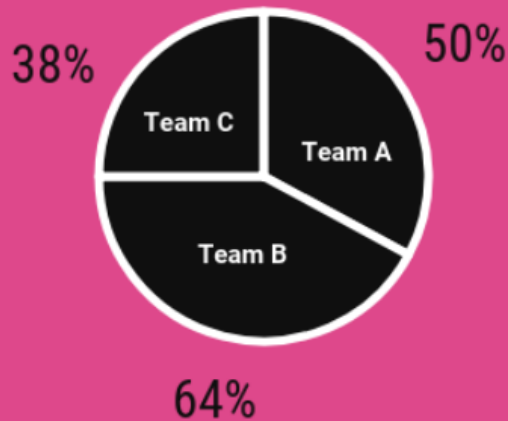- The scale is disproportionate to the data, making the change over time seem small

- The scale is proportionate to the data, showing a greater change over time
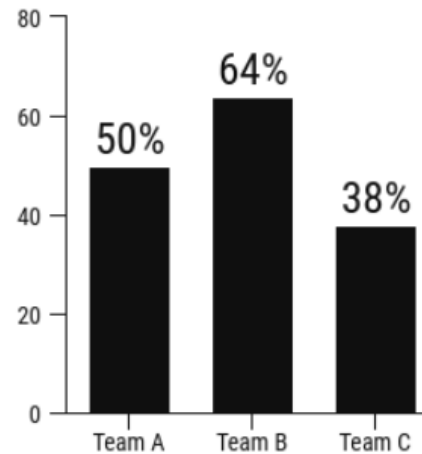
# 4 ▶ USING THE WRONG GRAPH

The type of graph you use should depend on the type of data you want to visualize.
Using the wrong type of graph can skew the data. Writers will sometimes use
the wrong type of graph on purpose.

38% 50%

Team C

Team A

Team B

64%

80

64%

50%

60

38%

40

20

0

Team A    Team B    Team C

**VS**

## 🙁 MISLEADING

- Pie charts are used to compare parts of a whole, not the difference between groups

- A different type of graph should be used to compare the three teams
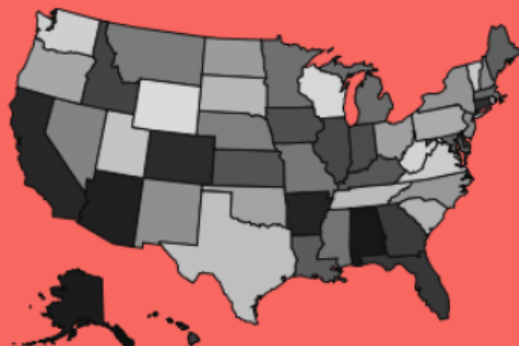
## ACCURATE 🙂

- Bar graphs are better for showing the differences between groups

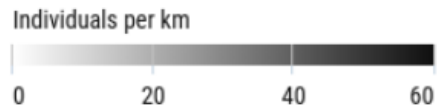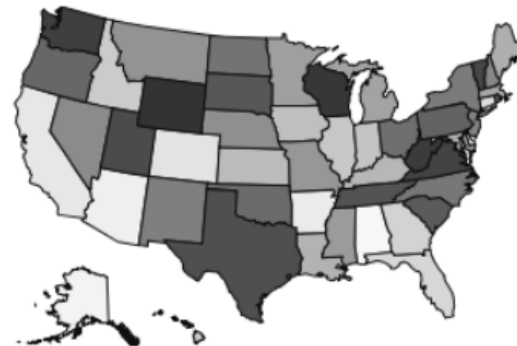- This chart is a better visualization of the data

# 5 ▶ GOING AGAINST CONVENTIONS

Over time, we have developed standards for how data is visualized. Flipping those conventions can make a graph confusing or misleading to readers.

Individuals per km

0    20    40    60

Individuals per km

0    20    40    60

## 😦 MISLEADING

**VS**

## ACCURATE 🙂

- Normally, darker shades are associated with density on a map but here, dark has been used to depict lower population density

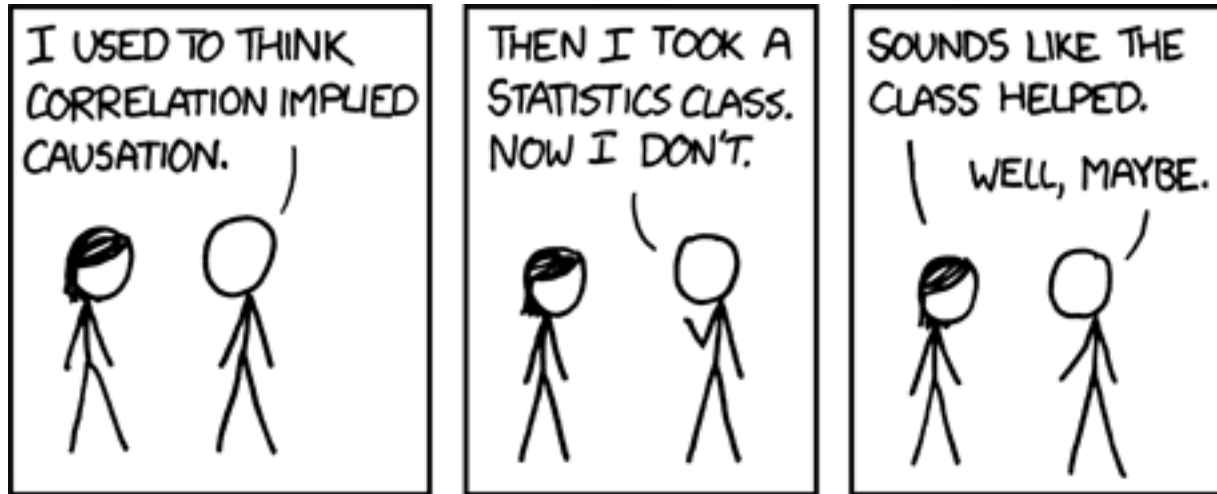- This graph can confuse and mislead readers, who expect dark to represent
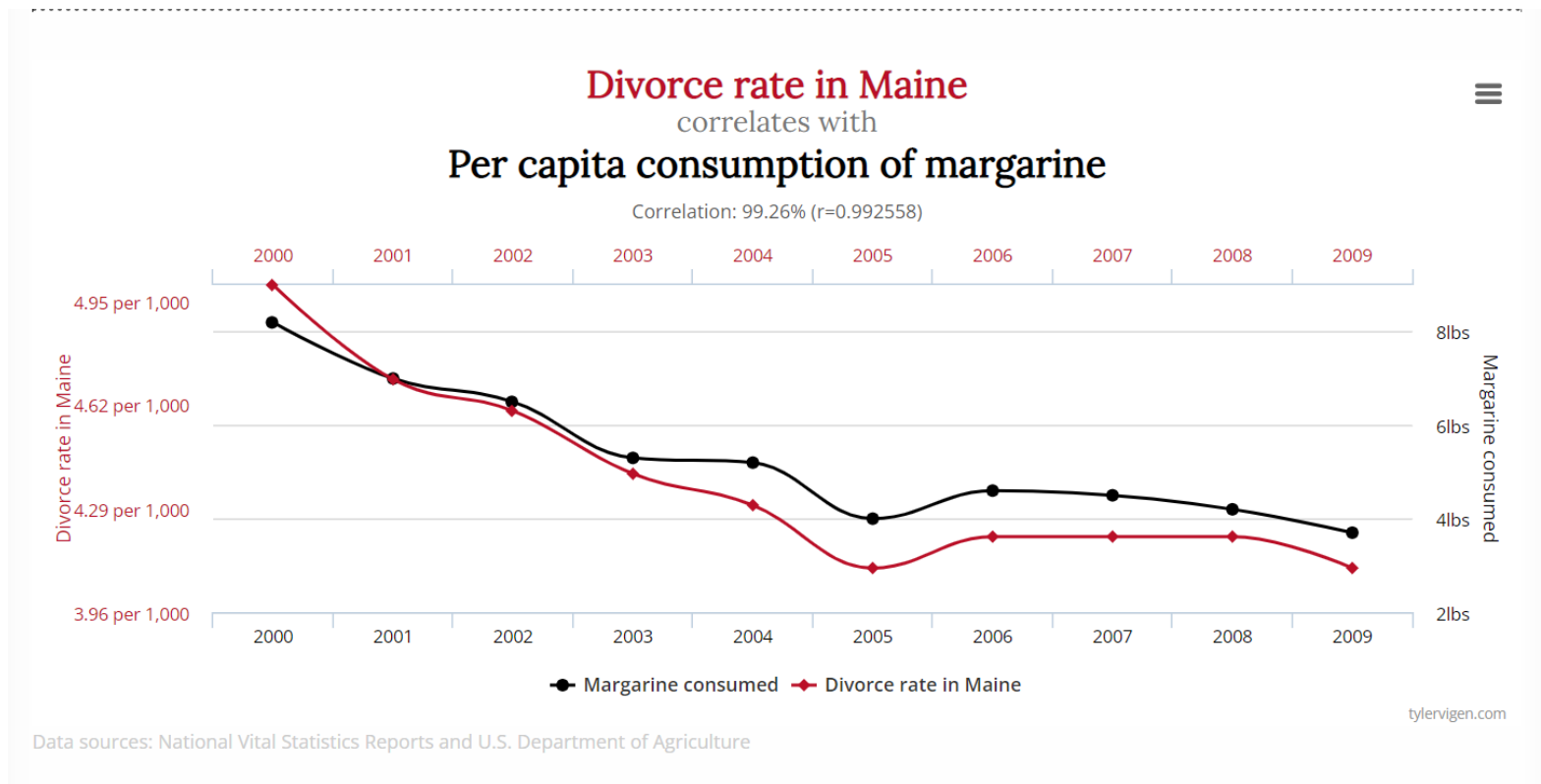
- This map follows the convention of using lighter shades for lighter density and darker shades for higher density

- Readers will intuitively know how to interpret the data

# BONUS: Spurious correlation

# BONUS: Spurious correlation



Divorce rate in Maine correlates with Per capita consumption of margarine. Correlation: 99.26% (r=0.992558). Data sources: National Vital Statistics Reports and U.S. Department of Agriculture. tylervigen.com

http://tylervigen.com/spurious-correlations
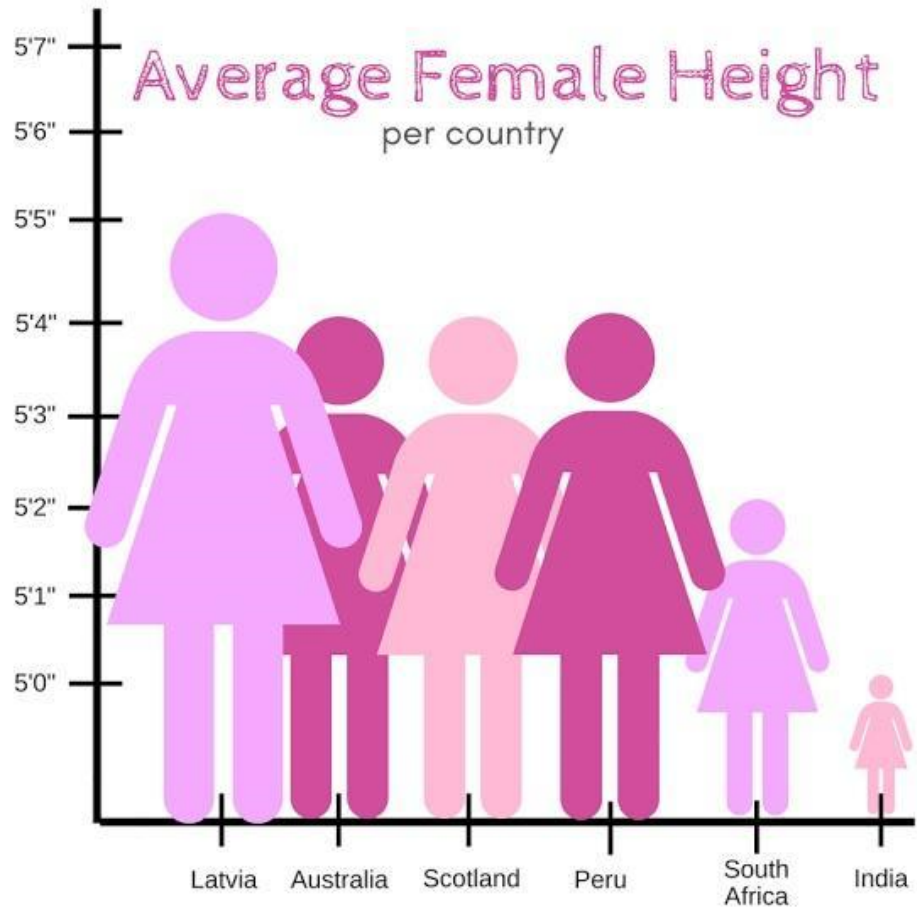
# BONUS: Should have been a log scale?

# BONUS: Should have been a log scale?

Caveat - does everyone undestand logs? More on this in future sessions

# Graph Crimes

# Graph Crimes



Average Female Height
per country

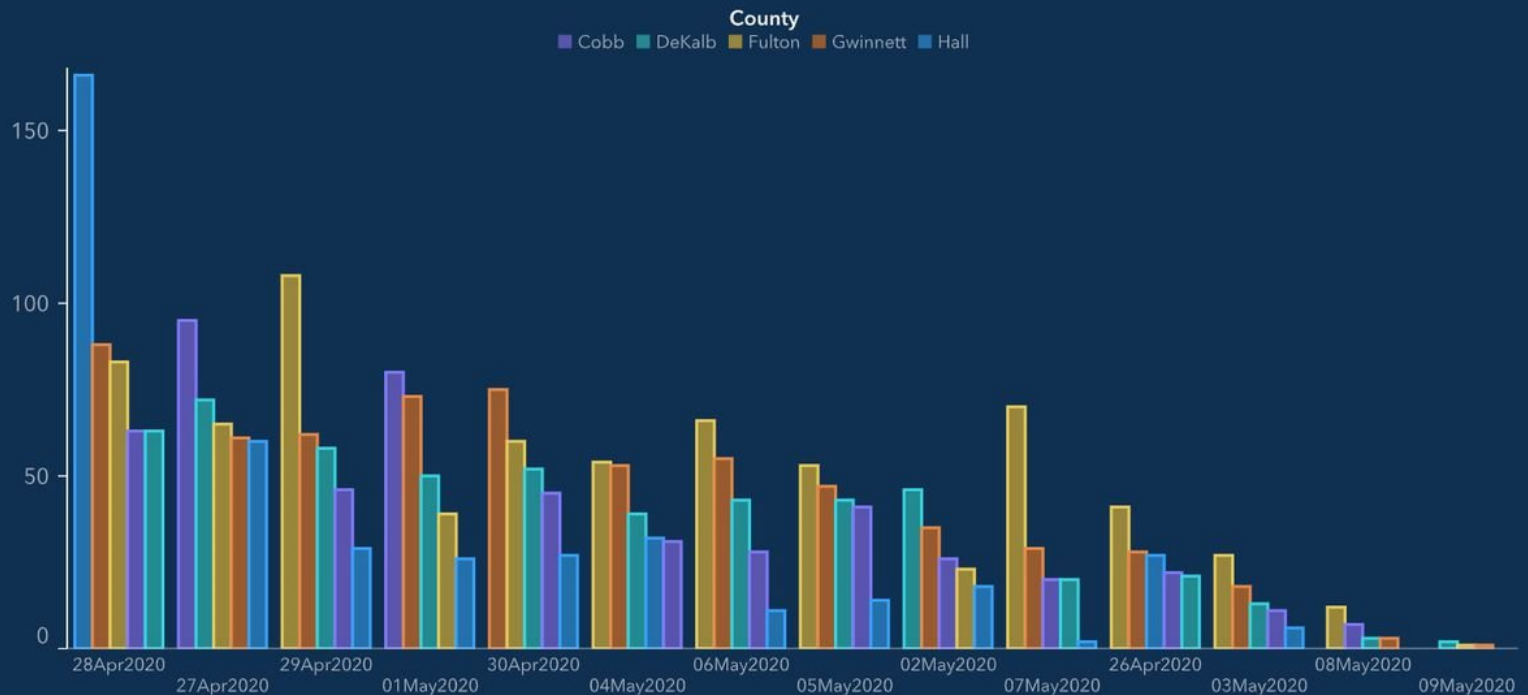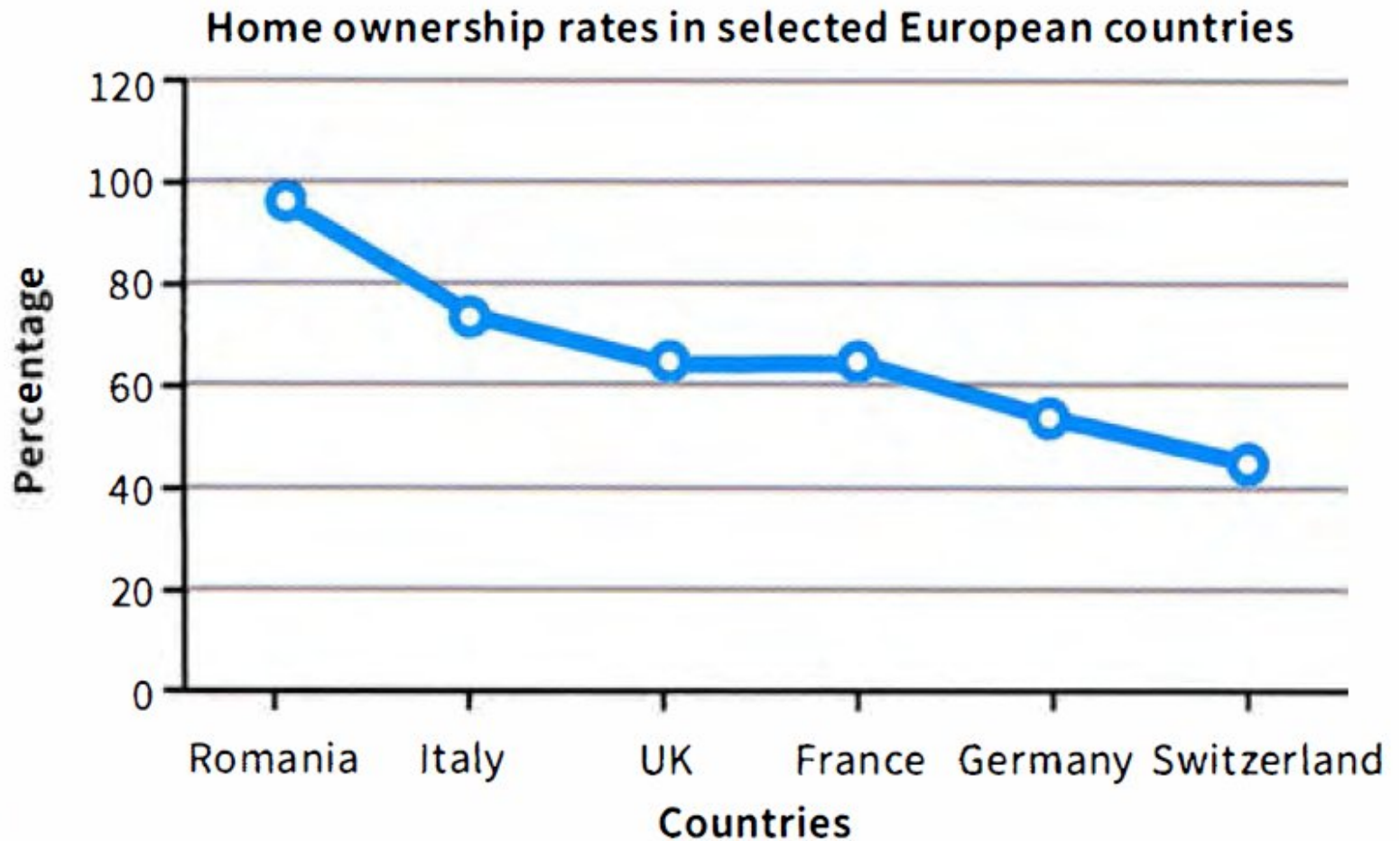Latvia  Australia  Scotland  Peru  South Africa  India

# Graph Crimes

# Graph Crimes



**Top 5 Counties with the Greatest Number of Confirmed COVID-19 Cases**

The chart below represents the most impacted counties over the past 15 days and the number of cases over time. The table below also represents the number of deaths and hospitalizations in each of those impacted counties.
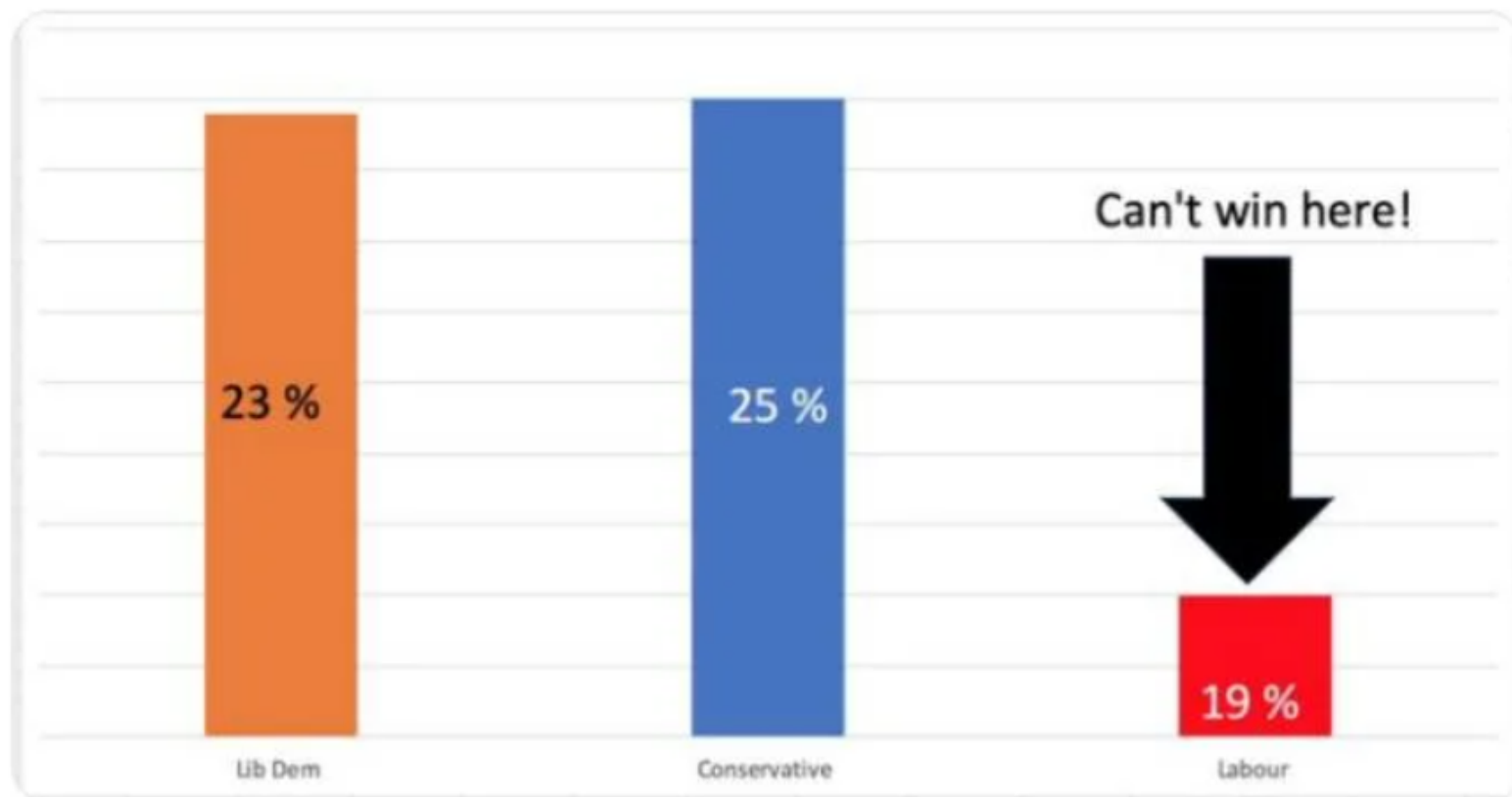
# Graph Crimes



Home ownership rates in selected European countries

# Graph Crimes

# Graph Crimes

# Graph Crimes

# Graph Crimes

# Lib Dem party *dislikes* accurate graphs

# Government Daily Briefings 2020

# Daily Briefings via the BBC



ated 568,100 people within the
nity population in England had
9 during the most recent week,
7th to 23rd October 2020. This
s to 1.04% of the population in
d or around 1 in 100 people.

ate there were around 9.52 new
-19 infections for every 10,000
per day within the community
n in England, equating to around
,900 new cases per day.

BBC NEWS
■ Further 21,915 coronavirus cases and 326 deaths across UK      18:50

# Daily Briefings via the BBC

# Essential Reading:

Fundamentals of Data Visualisation - Claus O. Wilke

R Graphics Cookbook - Winston Chang

A ggplot tutorial for beautiful plotting in R - Cédric Scherer