# Data Visualisation

Dr Andrew J. Stewart
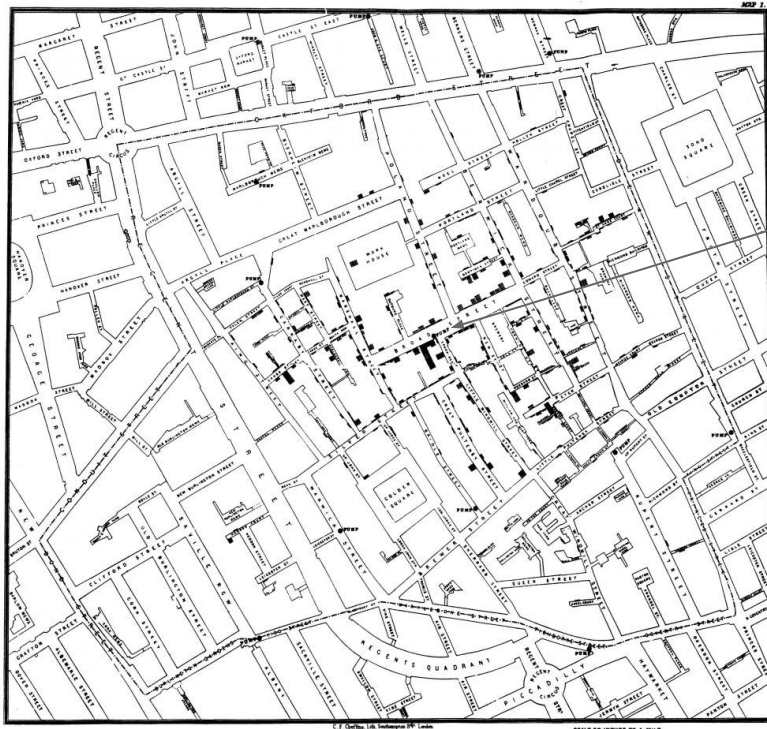
E: drandrewjstewart@gmail.com
T: @ajstewart_lang
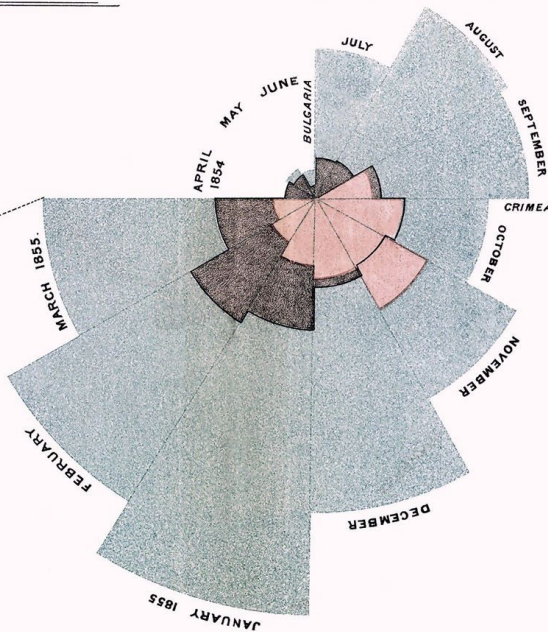G: ajstewartlang

# Some Classic Visualisations



John Snow's map of the cholera outbreak in 1854 in London showed that the outbreak was centred around a contaminated water pump in Broad Street.
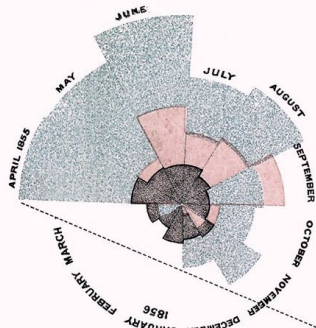
# Some Classic Visualisations



This Rose Chart or Coxcomb by Florence Nightingale (yes, that Florence Nightingale!) was used to capture the causes of death of "The Army in the East".
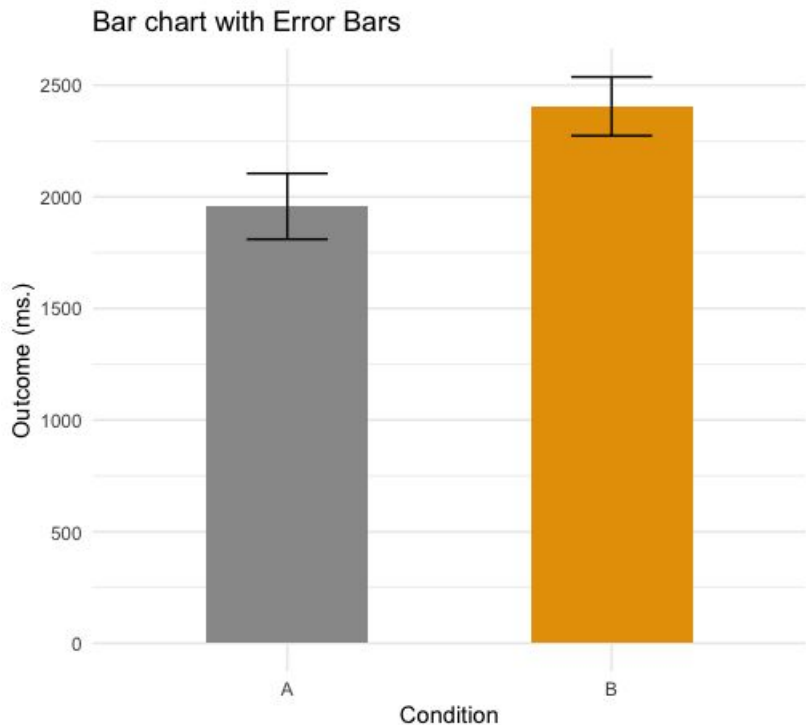
# The Classic Book

"The Visual Display of Quantitative Information" by Edward Tufte is the classic book on data visualisation - hugely influential and contains lots of examples of the best (and worst kinds) of data visualisations.

And speaking of one of the worst kinds...

SECOND EDITION

The Visual Display
of Quantitative Information

EDWARD R. TUFTE

# Visualisations that Can Mislead



Bar chart with Error Bars

Bar graphs tend to be quite limited in terms of what they communicate. Here they communicate the means for two conditions and information about variability around the means. But they don't tell us anything about the **distribution** of the data.

# Anscombe's Quartet



The data underlying each of the four plots on the left are all different, but each has the same mean and standard deviation for their x-values, the same mean and standard deviation for their y-values, the same correlation coefficient, and the same regression line. If we reported only these things, we'd think the four datasets were identical (whereas they are clearly not!)

# Plots Based on Aggregated Data Can Mislead



You might make one set of inferences based on this boxplot - maybe a median around 1,250 with the 25th and 75th percentiles being ~480 to ~1,980

# But look more closely at the actual data...



The data are clearly bimodal.

Distribution shape matters and we need to capture that in our data visualisations.

# The ggplot2 Package

The ggplot2 package is part of the Tidyverse and the function `ggplot()` forms the basis for data visualisations using Tidyverse packages and verbs.

There are three basic components when using the `ggplot()` function to build a visualisation:
- Data
    - The raw data that you want to plot
- Geometries (e.g., `geom_point()` and `geom_jitter()`)
    - The geometric shapes that will represent the data.
- Aethetics (`aes()`)
    - Aesthetics of the geometric and statistical objects, such as color, size, shape and position.

# Onto the R script below...