

# Data Analysis and Machine Learning 4 (DAML)

**Week 2: Summarising and visualising data**

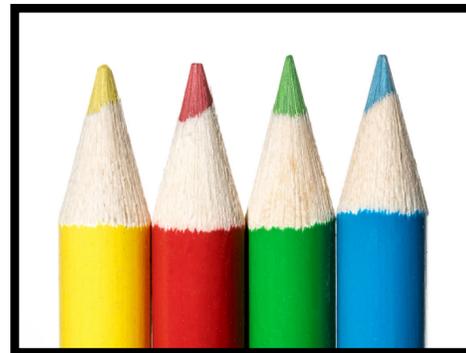
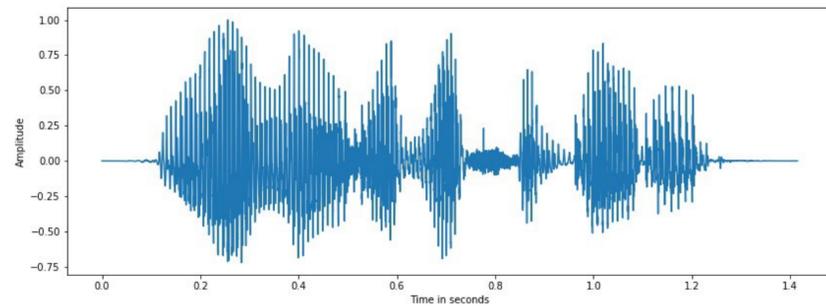
**Elliot J. Crowley, 22nd January 2024**



THE UNIVERSITY  
*of* EDINBURGH

# Recap

- We looked at different modalities of data



	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
...	...	...	...	...	...
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

- We considered variable types

iris species (nominal)



level of education (ordinal)



# Tabular data

- We will focus on this modality in this course
- It crops up a lot in real life and it is straightforward to work with

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	species
<b>0</b>	5.1	3.5	1.4	0.2	setosa
<b>1</b>	4.9	3.0	1.4	0.2	setosa
<b>2</b>	4.7	3.2	1.3	0.2	setosa
<b>3</b>	4.6	3.1	1.5	0.2	setosa
<b>4</b>	5.0	3.6	1.4	0.2	setosa
...	...	...	...	...	...
<b>145</b>	6.7	3.0	5.2	2.3	virginica
<b>146</b>	6.3	2.5	5.0	1.9	virginica
<b>147</b>	6.5	3.0	5.2	2.0	virginica
<b>148</b>	6.2	3.4	5.4	2.3	virginica
<b>149</b>	5.9	3.0	5.1	1.8	virginica

# Summarising Data

# World Happiness Report

- Produced by a non-profit of the United Nations
- What do you want to know when you see this?

Country or region	Score	GDP per capita	Social support	Healthy life expectancy	Freedom to make life choices	Generosity	Perceptions of corruption
Guatemala	6.436	0.800	1.269	0.746	0.535	0.175	0.078
Yemen	3.380	0.287	1.163	0.463	0.143	0.108	0.077
Netherlands	7.488	1.396	1.522	0.999	0.557	0.322	0.298
...	...	...	...	...	...	...	...
Libya	5.525	1.044	1.303	0.673	0.416	0.133	0.152
Jamaica	5.890	0.831	1.478	0.831	0.490	0.107	0.028
United States	6.892	1.433	1.457	0.874	0.454	0.280	0.128

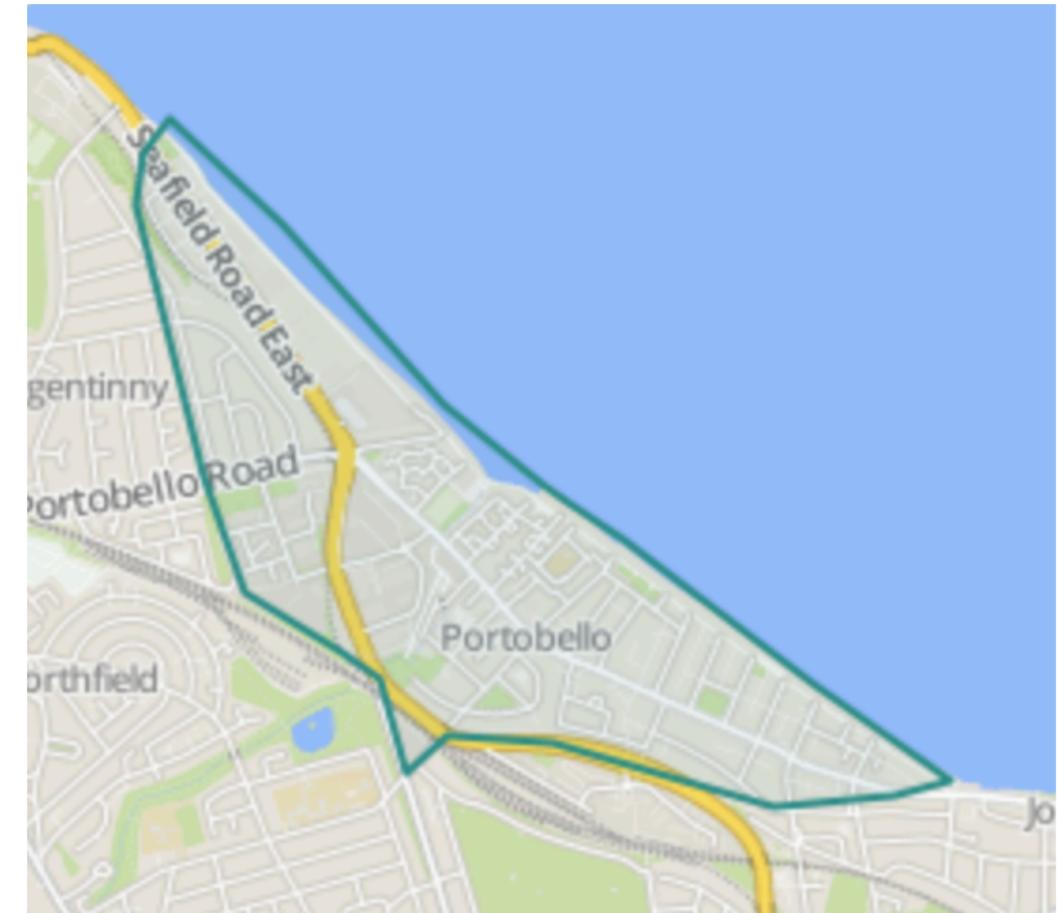
# Extreme values

- Take **maximum** of score: Finland
- Take **minimum** of perceived corruption: Moldova



# House buying

- Let's say I'm considering buying a property in Portobello
- What do I need to know?



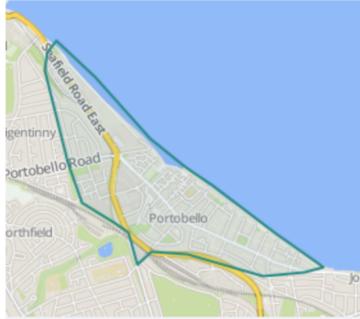
# Central values

- Good to know the mean house price
- Or median?

### House Prices in Portobello

Properties in Portobello had an overall average price of £308,327 over the last year.

Overall, sold prices in Portobello over the last year were 10% up on the previous year and 11% up on the 2008 peak of £276,604.



[Explore the map view](#) →

#### Properties sold

Filter:

2,414 sold properties Date sold ▾

2, Ormelie, Brunstane Road North, Edinburgh, Mid EH15 2DJ		
Unknown		
<b>£1,050,000</b>	31 Jan 2022	
£846,648	27 Jul 2020	
No other historical records		

#### Who provides this information?

Scottish house price data is publicly available information produced by the Registers of Scotland. Please note the dates shown below relate to the property's registered date not sold date. This material was last updated on 7 March 2022.

Average price in this area:  
**£308,327**   
↑ 10% since 2019

#### How much can I borrow?

Check your affordability and learn how much you can borrow, based on your monthly income and outgoings.

[Try Nationwide's affordability calculator](#)

Advertisement

Nationwide pays Rightmove a fee for each completed mortgage. It's up to you if you choose Nationwide, or a different lender, to suit your mortgage needs and circumstances.

What is your property worth?

# Summary statistics

- Most people will not scroll through a table!
- Summary statistics let us convey information as simply as possible
- We will now look at some (sample) statistics of (random) variables

WORLD >

## 99% of the world is breathing poor-quality air, WHO says

APRIL 4, 2022 / 3:01 PM / AP

f t

**Salaries in London Area**

Location:  or Find a Specific Employer:   Sort:

Company	Average Base Salary in (GBP)	Range
 <b>Accenture</b> London 4.1 ★ 21 salaries <a href="#">See 21 salaries from all locations</a>	£54,608 /yr	£32K - £102K
 <b>Deloitte</b> London 4.0 ★ 19 salaries <a href="#">See 20 salaries from all locations</a>	£58,219 /yr	£31K - £105K
 <b>Barclays</b> London 4.0 ★ 16 salaries <a href="#">See 16 salaries from all locations</a>	£52,872 /yr	£18K - £109K
 <b>University College London</b> London 4.3 ★ 10 salaries <a href="#">See 10 salaries from all locations</a>	£39,800 /yr	£18K - £57K

Source: Glassdoor

# Mode

- Suitable for summarising ordinal, nominal, and discrete variables
- Let's denote our (random) variable as  $X$
- We have measurements of that variable
- The mode is the measurement that occurs the most

Favourite Colour	
0	red
1	blue
2	red
3	red
4	blue
5	yellow

3 red, 2 blue, 1 yellow

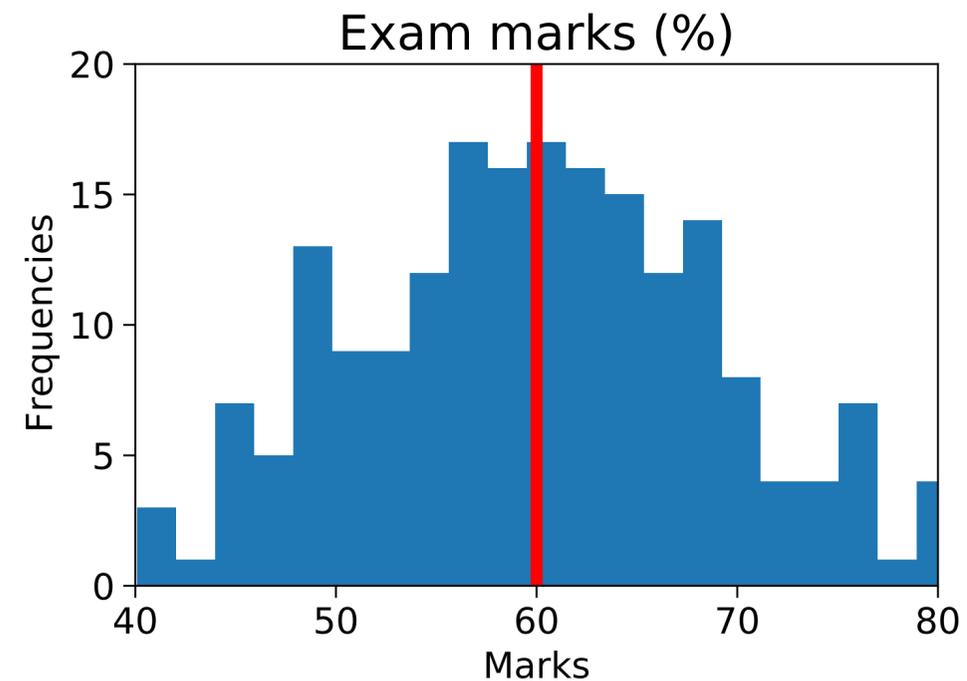
The mode is red

# Mean

- Denote as  $\mu$ . Suitable for summarising numerical variables
- For variable  $X$  we have  $N$  measurements  $\{x^{(n)}\}_{n=1}^N$
- The measurements are just  $x^{(1)}, x^{(2)}, \dots, x^{(N)}$

$$\mu_x = \frac{1}{N} \sum_{n=1}^N x^{(n)}$$

	Mark (%)
0	60
1	40
2	45
...	...



# Variance and Standard Deviation

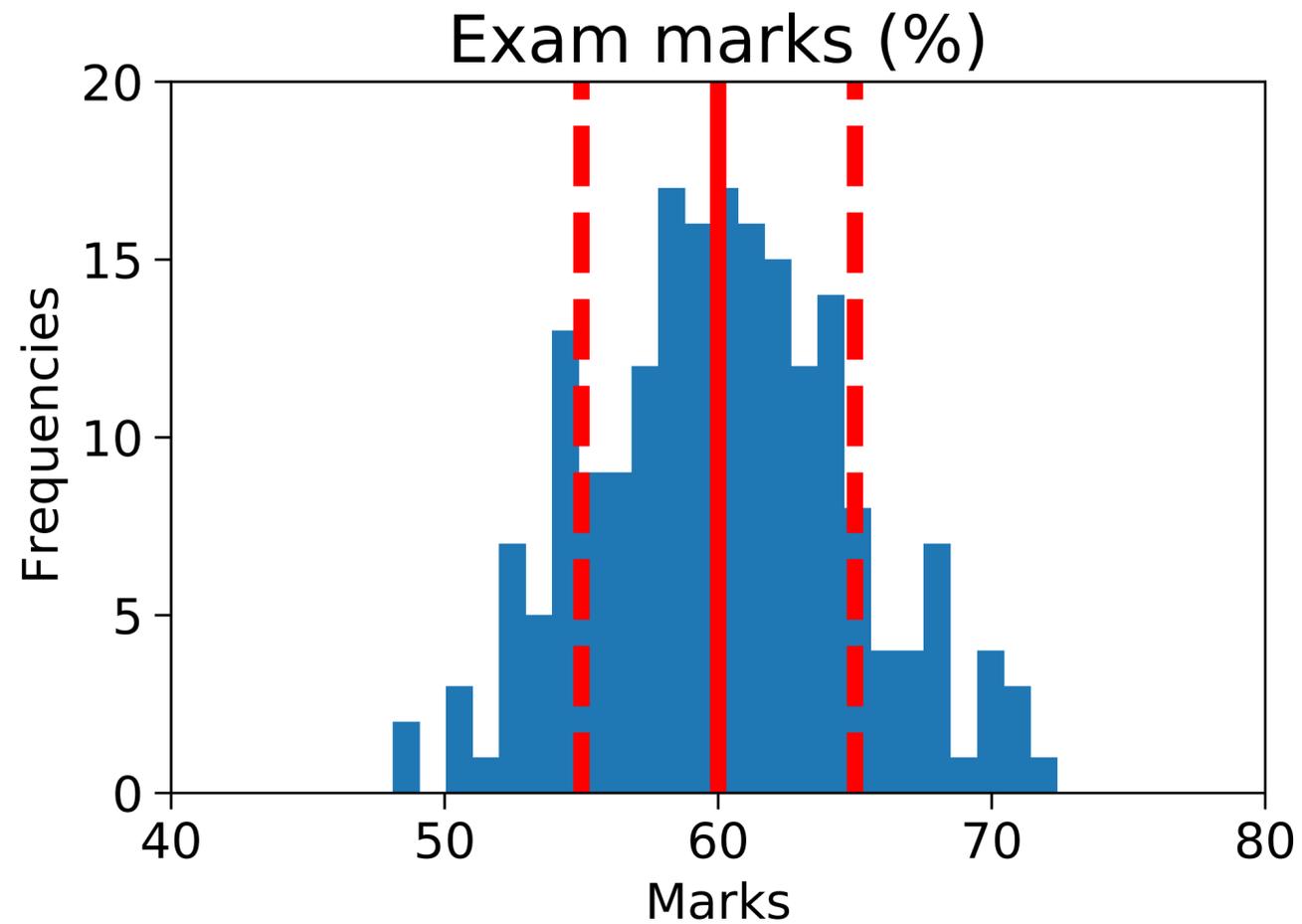
- Let's denote variance as  $\sigma^2$  and Standard deviation (SD) as  $\sigma$
- For variable  $X$  we have  $N$  measurements  $\{x^{(n)}\}_{n=1}^N$

$$\sigma_x^2 = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu_x)^2$$

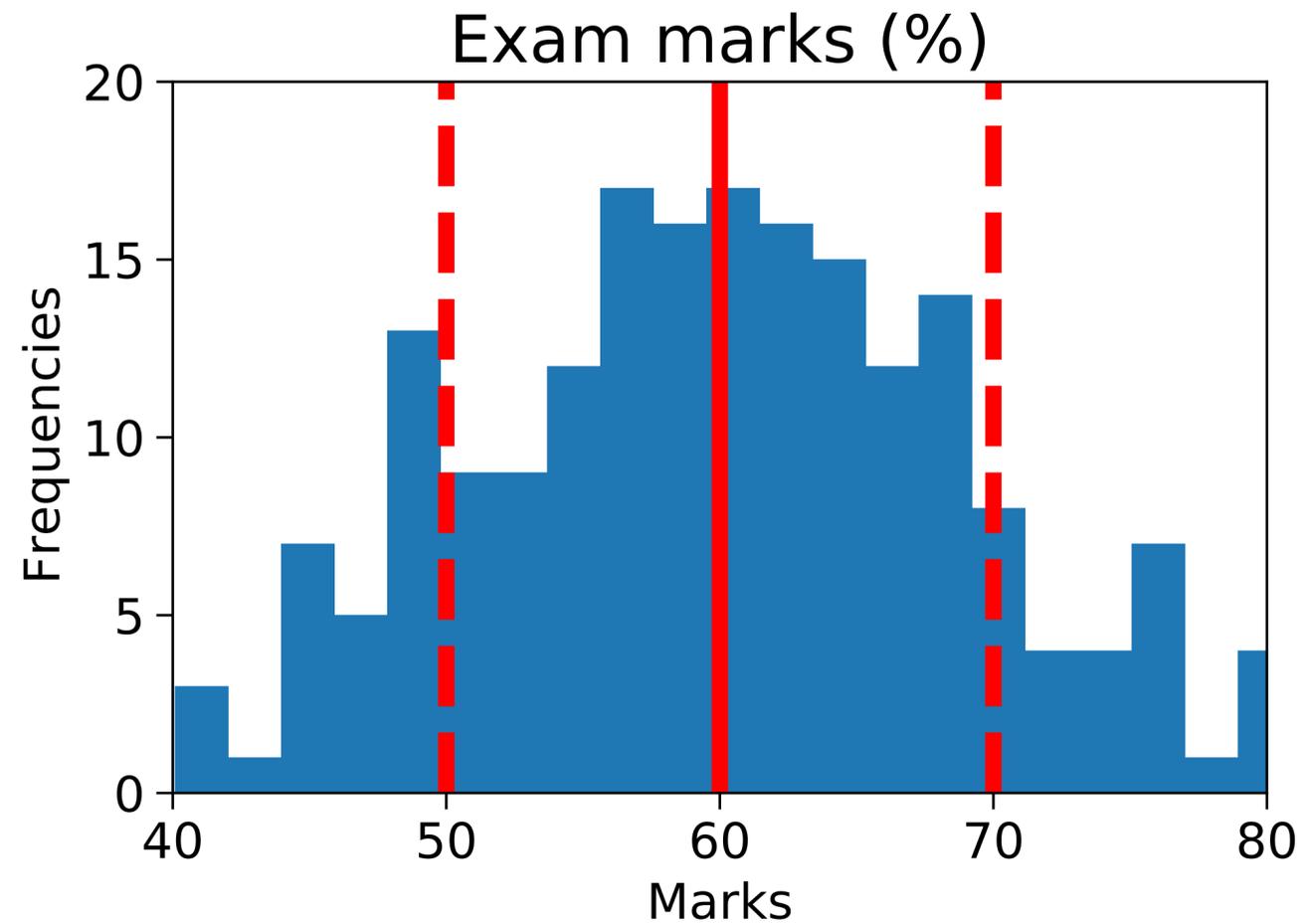
- Be aware that some definitions divide by  $N - 1$
- $N \approx N + 1$  for large  $N$  so this isn't that important!

# Standard Deviation

SD measures the extent to which measurements deviate from the mean



$$\sigma = 5$$

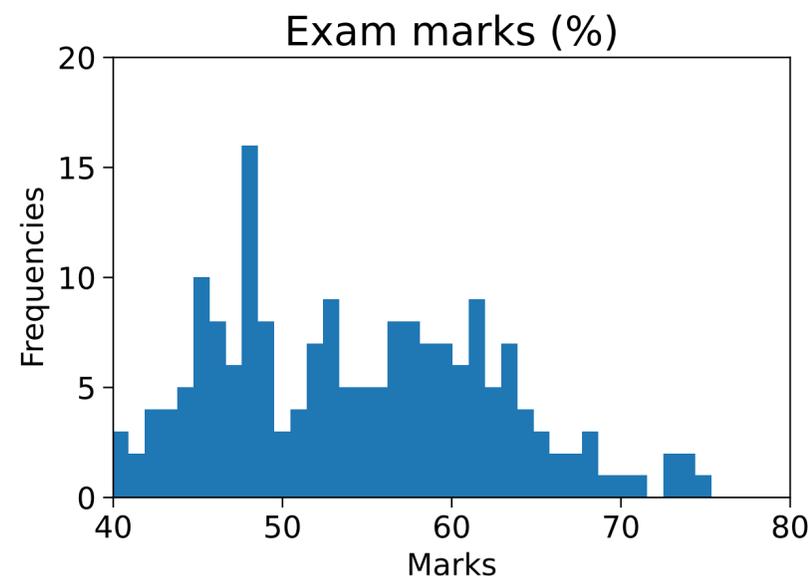


$$\sigma = 10$$

# Skewness

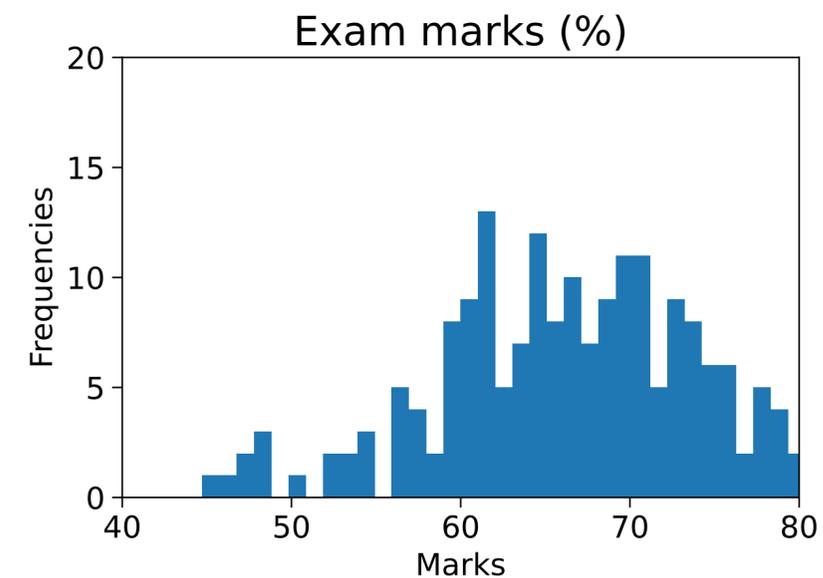
- Denote using  $s$ . For variable  $X$  we have  $N$  measurements  $\{x^{(n)}\}_{n=1}^N$

$$s_x = \frac{\frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu_x)^3}{\sigma_x^3}$$



## Positive skew

Bulk of measurements on the left  
Tail on the right



## Negative skew

Bulk of measurements on the right  
Tail on the left

# Median

- Order measurements of a numerical variable from lowest to highest
- The median is the measurement in the middle

1 2 3 **5** 8 12 17

- The median is a **robust statistic**

1 2 3 **5** 8 12 1700000000

# Medians are robust to outliers

Median salary is more meaningful than mean salary

## Bet365 boss Denise Coates sees pay jump to £221m

8 January



BET365/PA MEDIA

By Lora Jones

Business reporter, BBC News

The boss of Bet365 was paid around £221m during its last financial year, despite the gambling giant reporting a significant loss.

<https://www.bbc.co.uk/news/business-67912483>

BUSINESS

## FTSE 100 bosses earn average UK yearly pay after only three days

A typical boss of a company in London's top-flight stock market index makes £3.8 million a year



TAKEOVER DEALS DROP TO LOWEST SINCE FINANCIAL CRISIS (IAN WEST/PA)

PA ARCHIVE

DANIEL O'BOYLE @DAN\_O\_BOYLE  
4 JANUARY 2024

<https://www.standard.co.uk/business/ftse-100-bosses-ceo-salary-high-pay-centre-earnings-inequality-b1130313.html>

# Relating variables to each other

- We may be interested in the relationship between two variables
- Does GDP per capita relate to Healthy life expectancy?

Country or region	Score	GDP per capita	Social support	Healthy life expectancy	Freedom to make life choices	Generosity	Perceptions of corruption
Guatemala	6.436	0.800	1.269	0.746	0.535	0.175	0.078
Yemen	3.380	0.287	1.163	0.463	0.143	0.108	0.077
Netherlands	7.488	1.396	1.522	0.999	0.557	0.322	0.298
...	...	...	...	...	...	...	...
Libya	5.525	1.044	1.303	0.673	0.416	0.133	0.152
Jamaica	5.890	0.831	1.478	0.831	0.490	0.107	0.028
United States	6.892	1.433	1.457	0.874	0.454	0.280	0.128

# Covariance and correlation

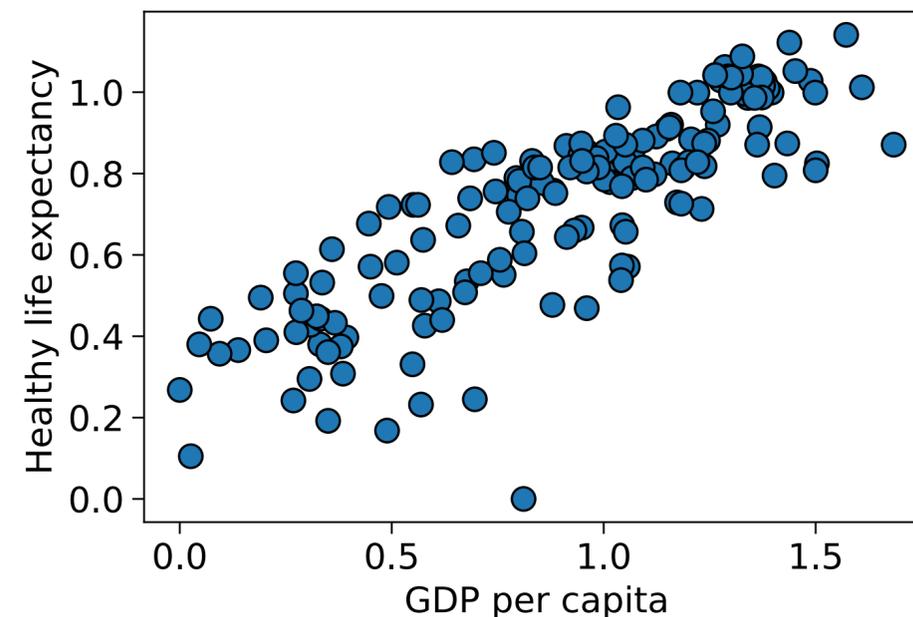
- We have two numerical variables  $X$  and  $Y$  each with  $N$  measurements
- Let's compute the means and SDs of each:  $\mu_x, \mu_y, \sigma_x, \sigma_y$
- The covariance  $\sigma_{x,y}$  and **Pearson correlation coefficient**  $\rho_{x,y}$  are given by:

$$\sigma_{x,y} = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu_x)(y^{(n)} - \mu_y)$$

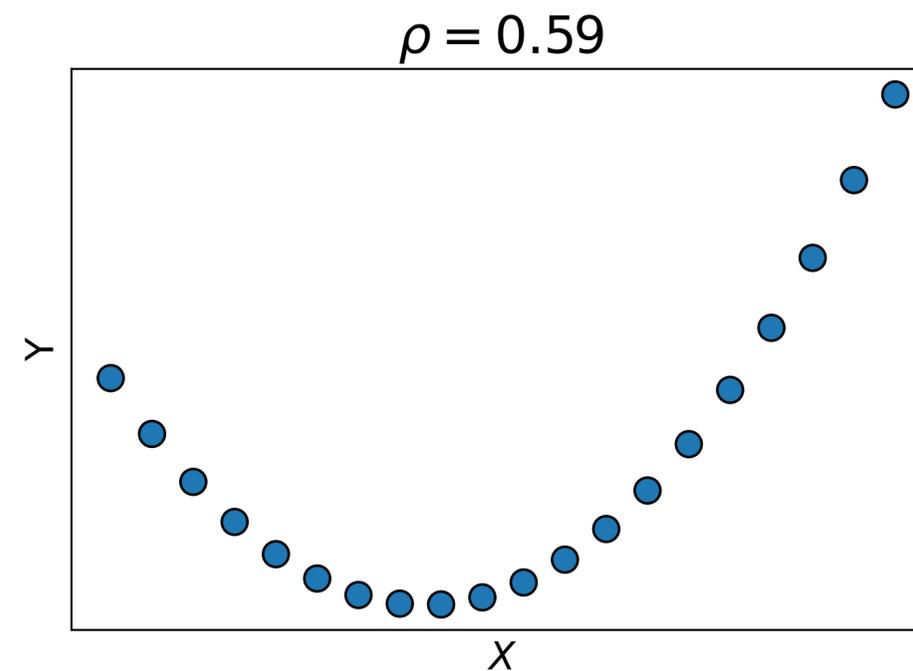
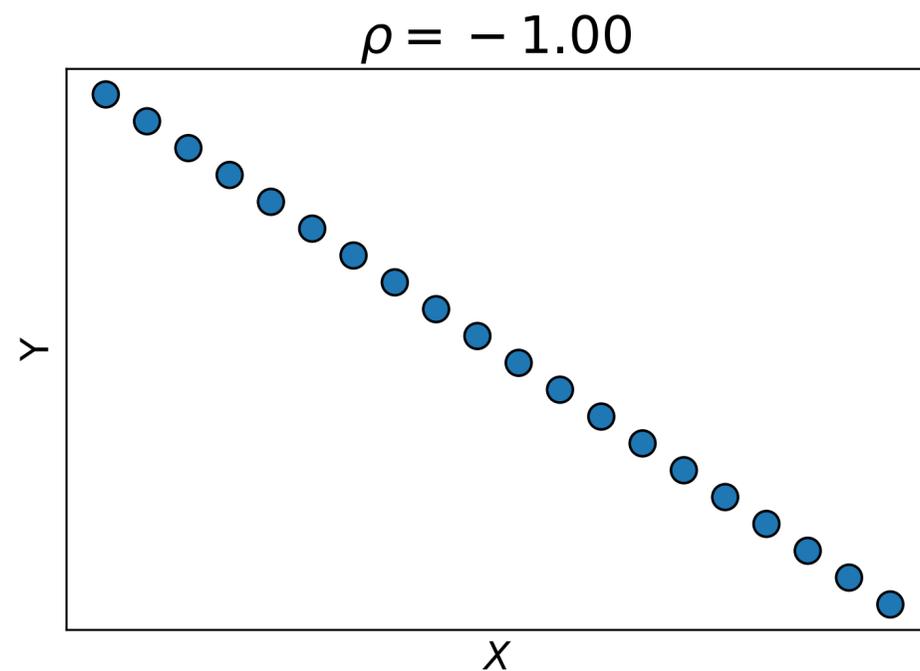
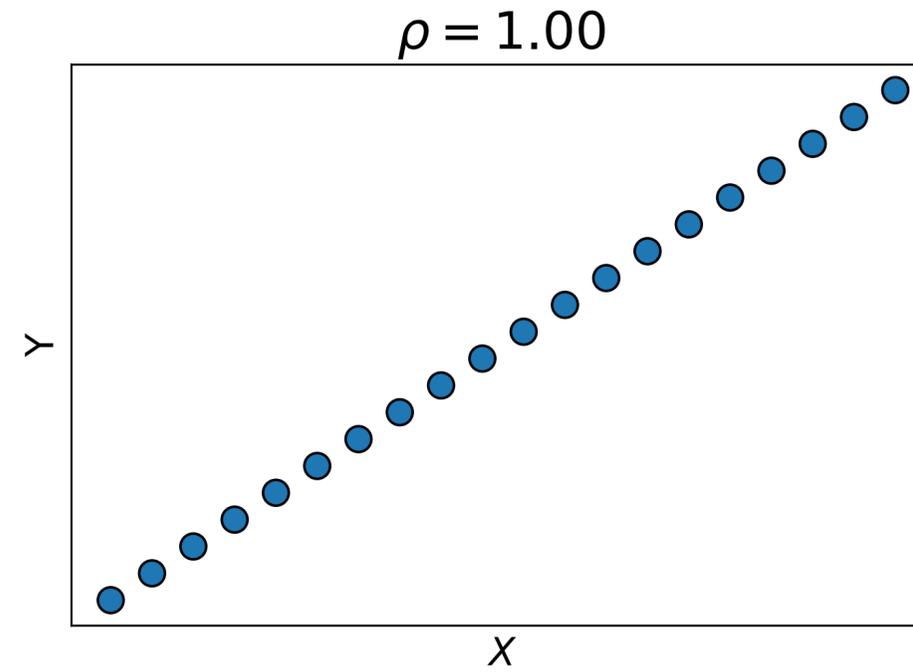
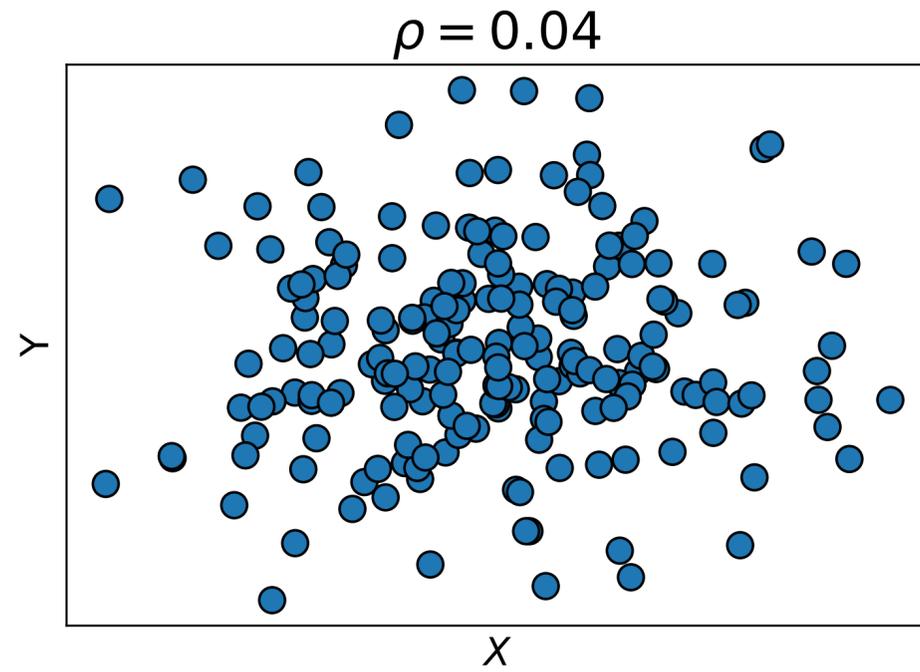
$$\rho_{x,y} = \frac{\sigma_{x,y}}{\sigma_x \sigma_y}$$

# Pearson correlation coefficient

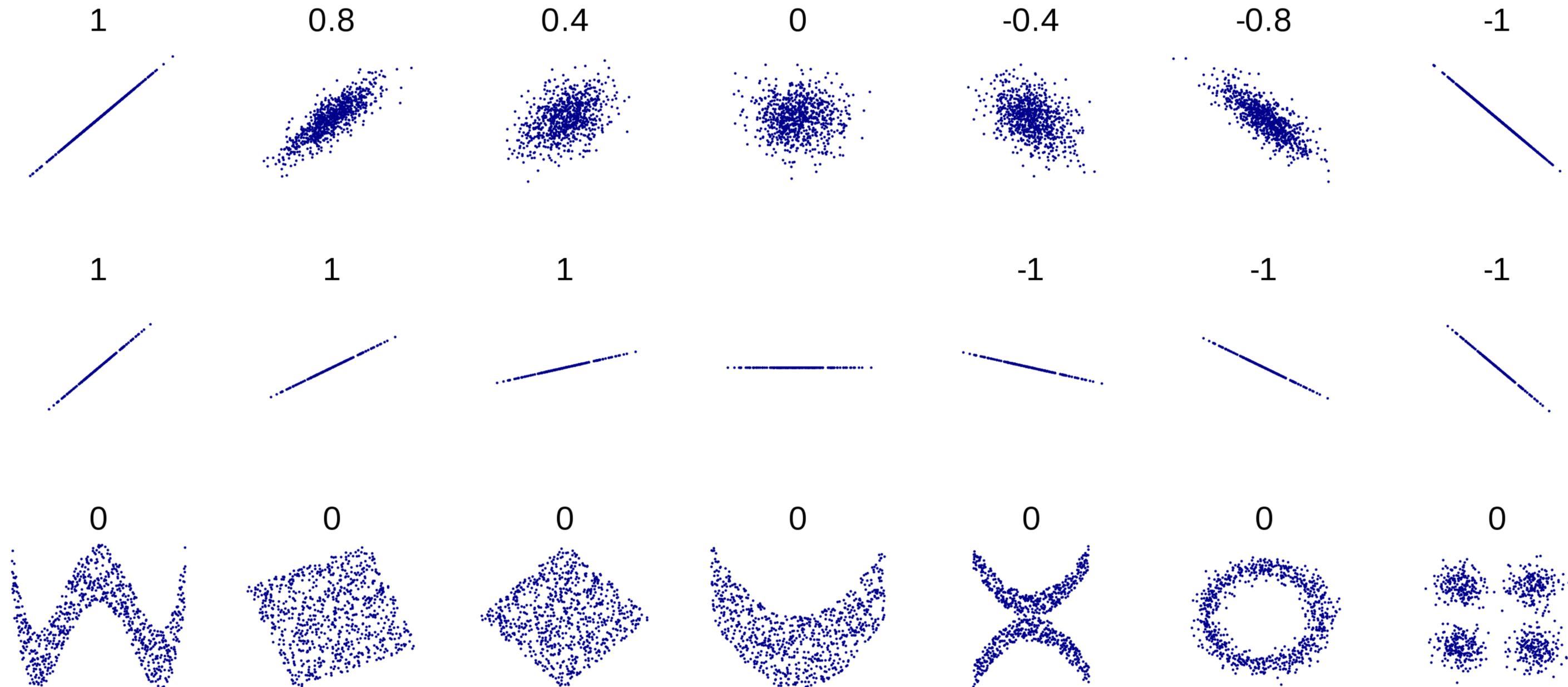
- $\rho_{x,y}$  has a value between  $-1$  and  $+1$
- It gives a measure of how linear the relationship between  $X$  and  $Y$  is
- That is, the extent to which we can use a line to map one to the other
- 0.84 for GDP per capita and Healthy life expectancy



# Pearson correlation coefficient

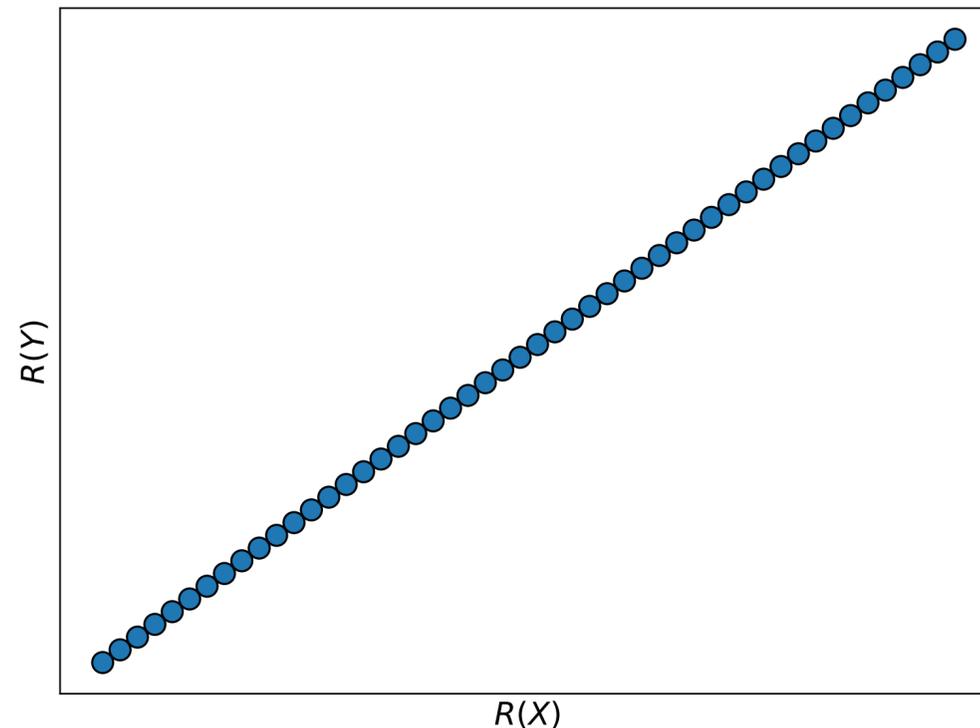
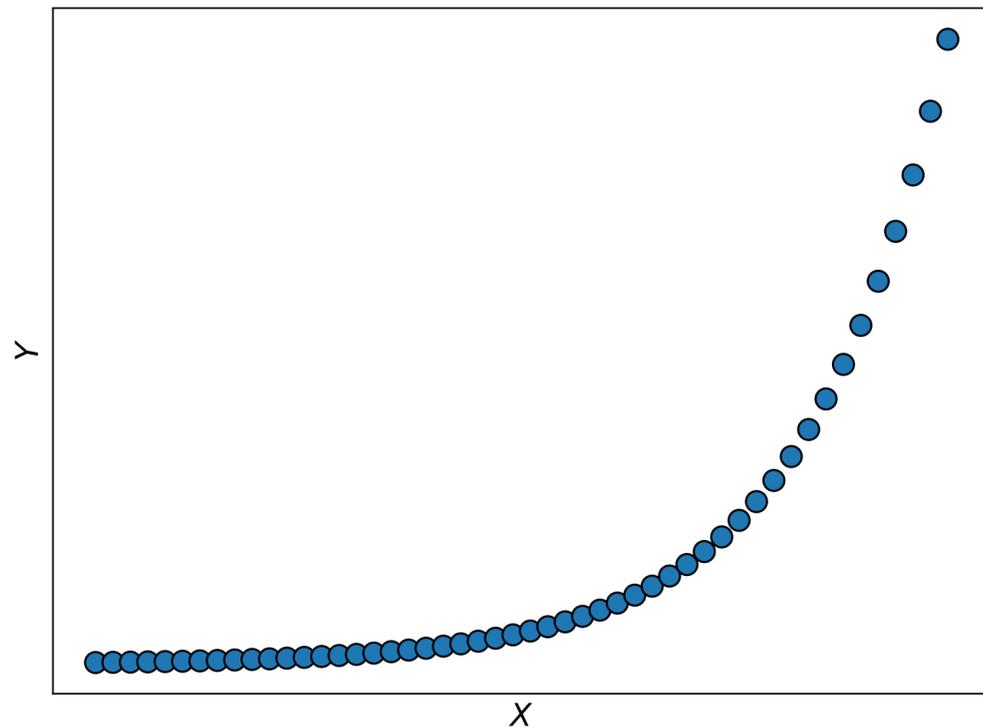


# Pearson correlation coefficient



# Spearman's rank correlation coefficient $r_s$

- This is the Pearson correlation coefficient of the **ranks** of two variables
- e.g. if  $X = \{1, 100, 1000, 10000\}$ ,  $R(X) = \{4, 3, 2, 1\}$
- It measures how monotonic the relationship between the variables is

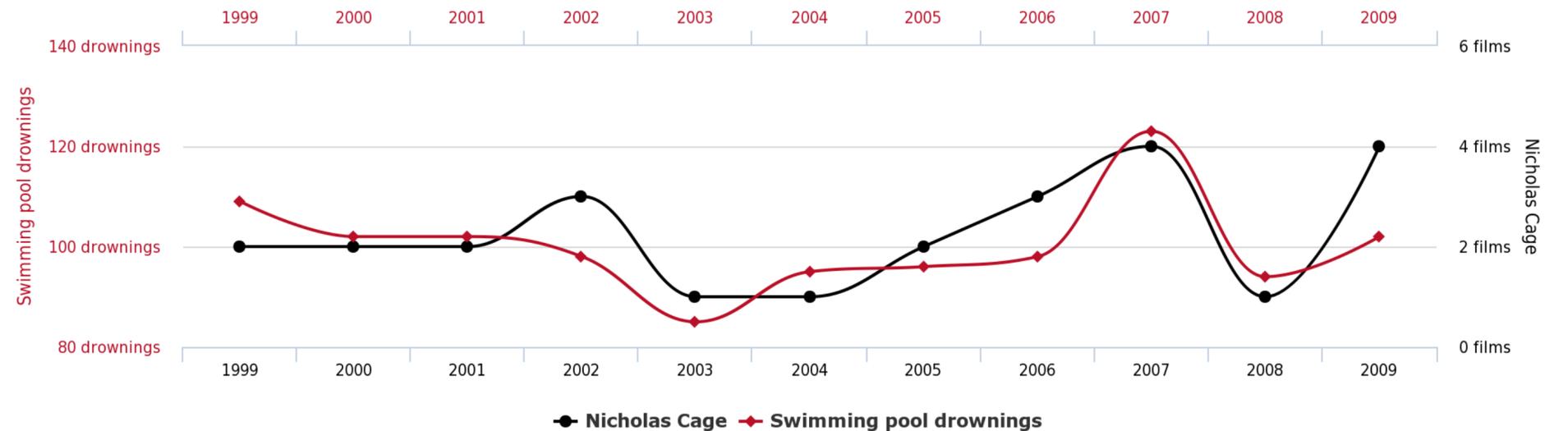


The Pearson correlation here is 0.82  
The Spearman correlation is 1

# Correlation does not imply causation



**Number of people who drowned by falling into a pool**  
correlates with  
**Films Nicolas Cage appeared in**



# Rubbish in, rubbish out

If your data is rubbish then anything you extract from it is also rubbish

- You might not have enough data points
- The process for collecting data might be flawed (e.g. biased)
- Measurements might be recorded incorrectly
- The variables chosen might not be useful



# Misleading statistics

Can be nefarious, or just stupidity



**BBC NEWS** **LIVE** **BBC NEWS CHANNEL**

Last Updated: Wednesday, 17 January 2007, 02:45 GMT

[E-mail this to a friend](#) [Printable version](#)

## Colgate warned over '80%' boast

**The maker of Colgate toothpaste has been warned not to repeat its famous advertising claim that "more than 80% of dentists recommend Colgate".**



Colgate's claim on posters was "misleading"

The Advertising Standards Authority concluded the claim on Colgate posters was "misleading" after investigating the phone survey behind the boast.

It found the dentists surveyed were allowed to name more than one brand.

But the ASA said the advertising claim implied 80% of dentists recommended Colgate to the exclusion of its rivals.

In fact, the ASA's inquiry found another competitor's brand was recommended almost as much as Colgate was by those dentists who were surveyed.

It added the survey "did not make clear the poll was on behalf of Colgate".

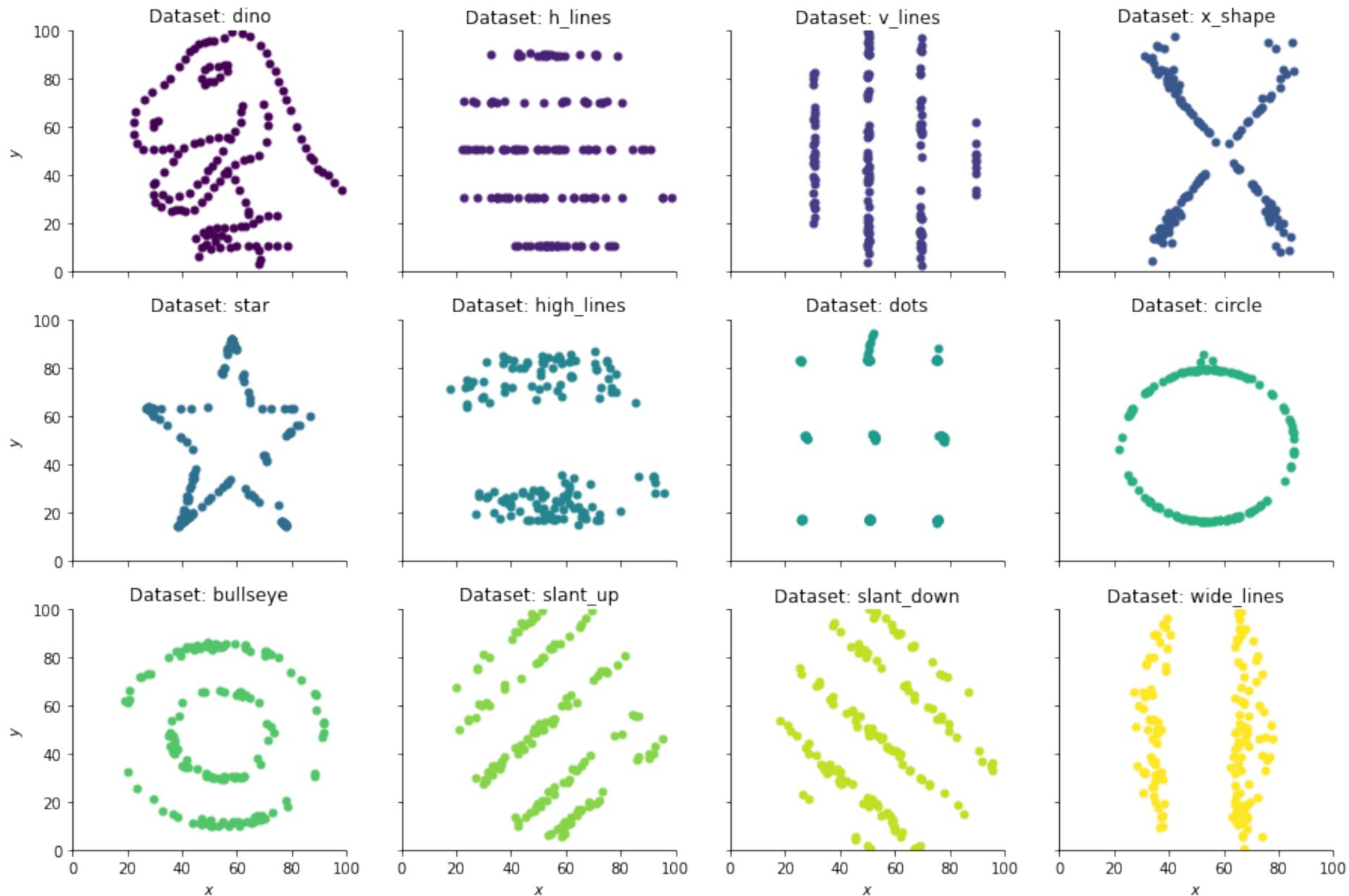
**RELATED BBC SITES**  
SPORT  
WEATHER  
CBBC NEWSROUND

## HANLON'S RAZOR

*Never attribute to malice  
that which is adequately explained  
by stupidity*



# The limitations of summary statistics



All 12 of these datasets have the same  $\mu_x, \mu_y, \sigma_x, \sigma_y, \rho_{x,y}$

# Visualising Data

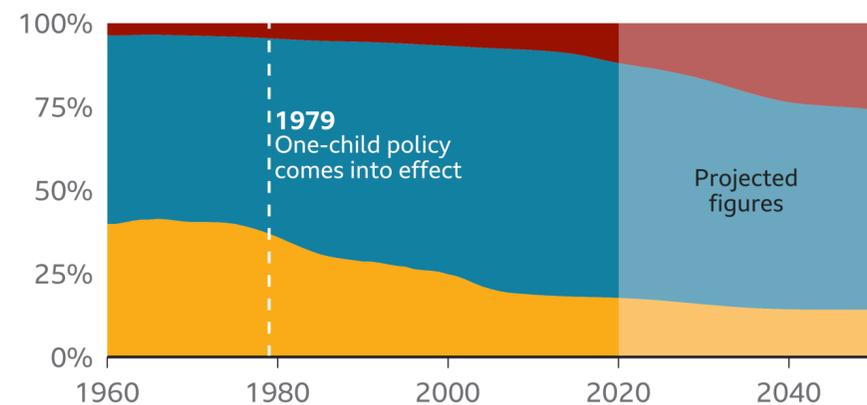
# Visualising data for presentation

- Conveying information as **simply**, and **clearly** as possible
- It is an art form, combining data analysis with graphic design

## China's population by age group

Proportion of total population (1960-2050)

0-14 years 15-64 65+



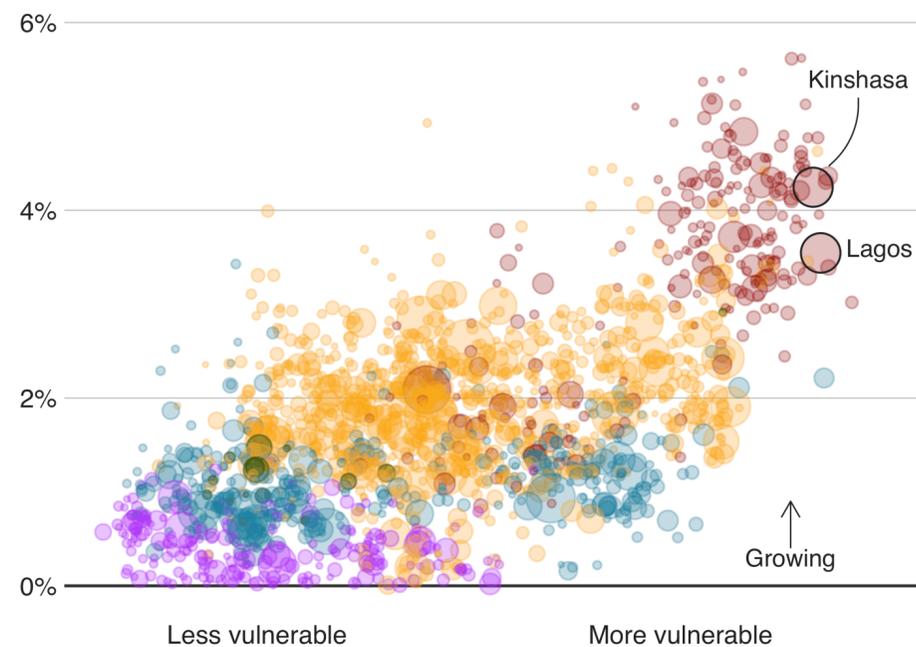
Source: The World Bank

BBC

## Fast-growing cities face worse climate risks

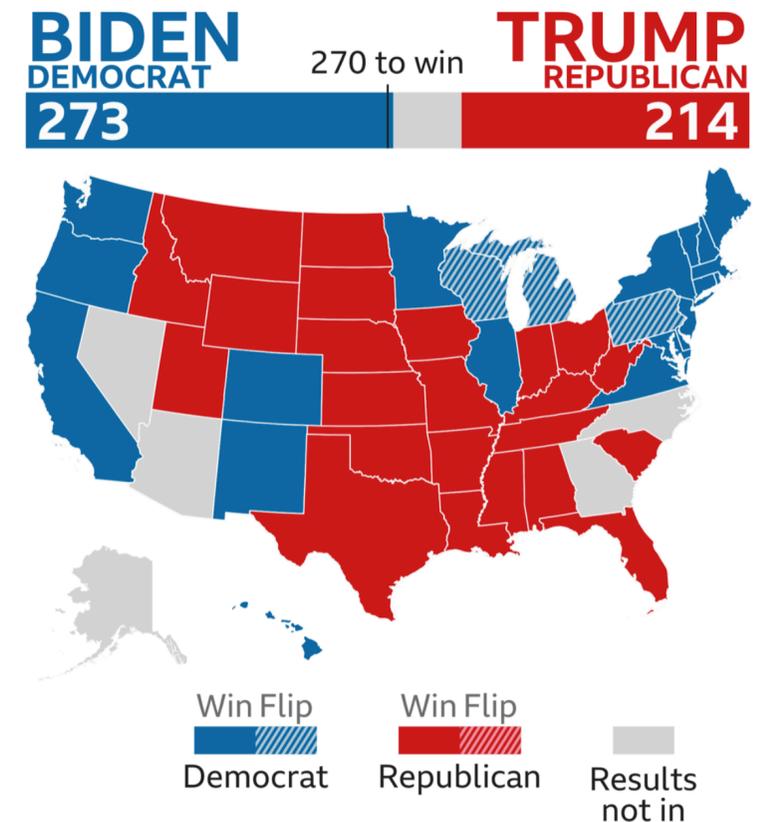
Population growth 2018-2035 over climate change vulnerability

Africa Asia Americas Europe Oceania



Source: Verisk Maplecroft. Circle size represents current population.

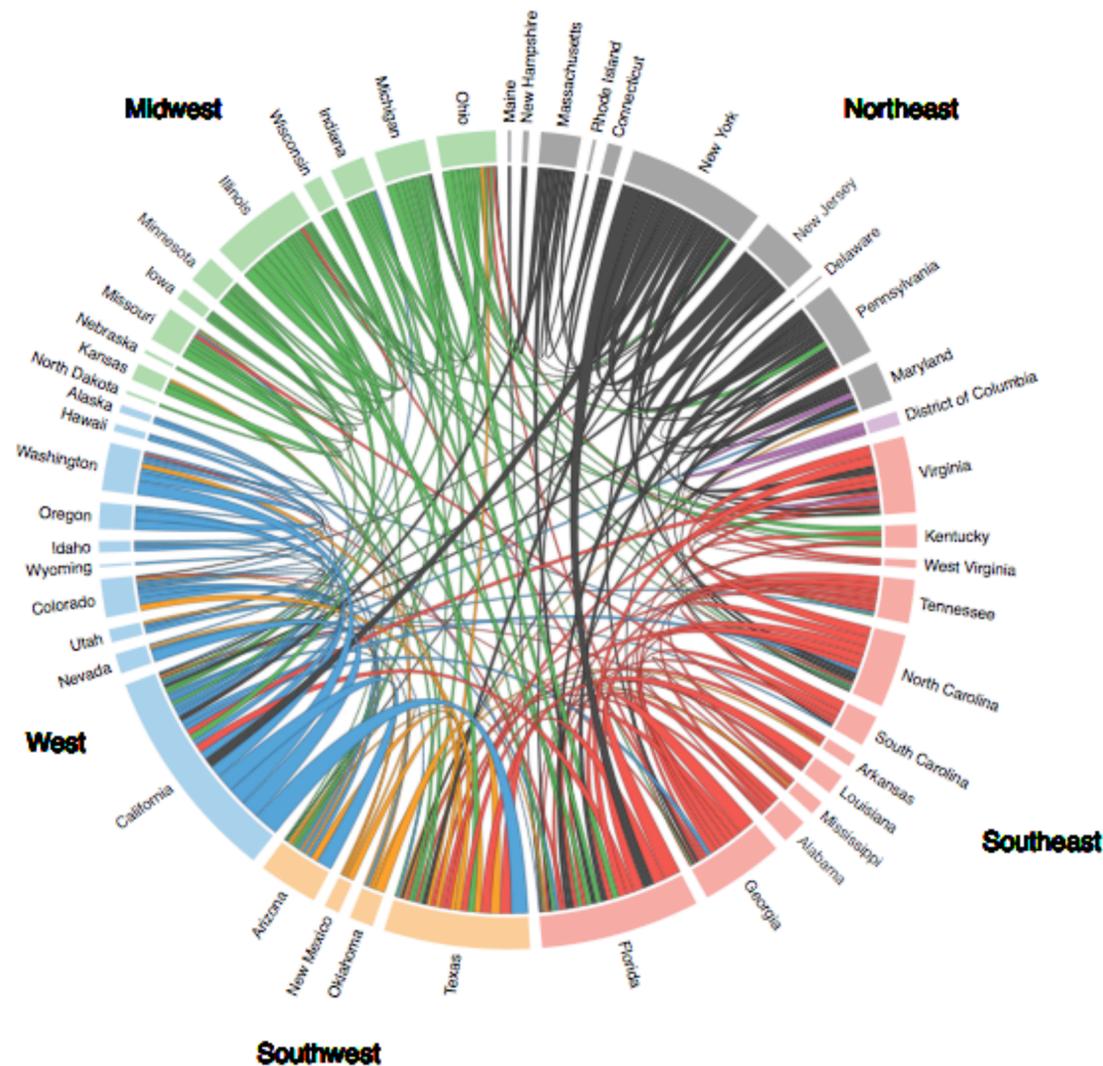
BBC



Source: BBC

# Visualising data for presentation

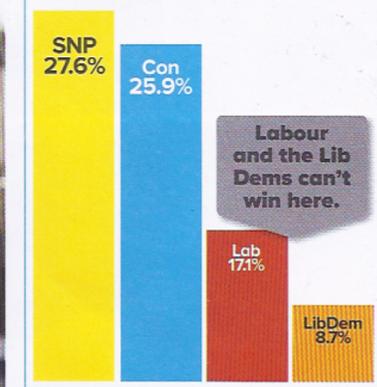
Can be done badly e.g. overcomplicated or misleading



## IAIN MCGILL: RUTH DAVIDSON'S CANDIDATE

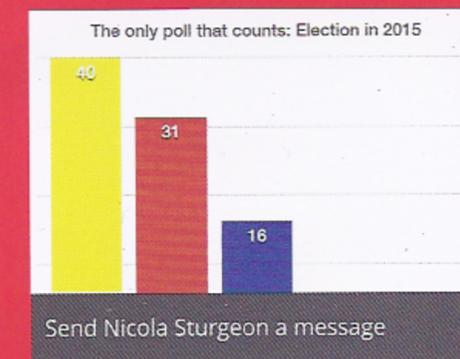


HOW WE VOTED IN  
EDINBURGH NORTH AND LEITH  
ON 4TH MAY 2017



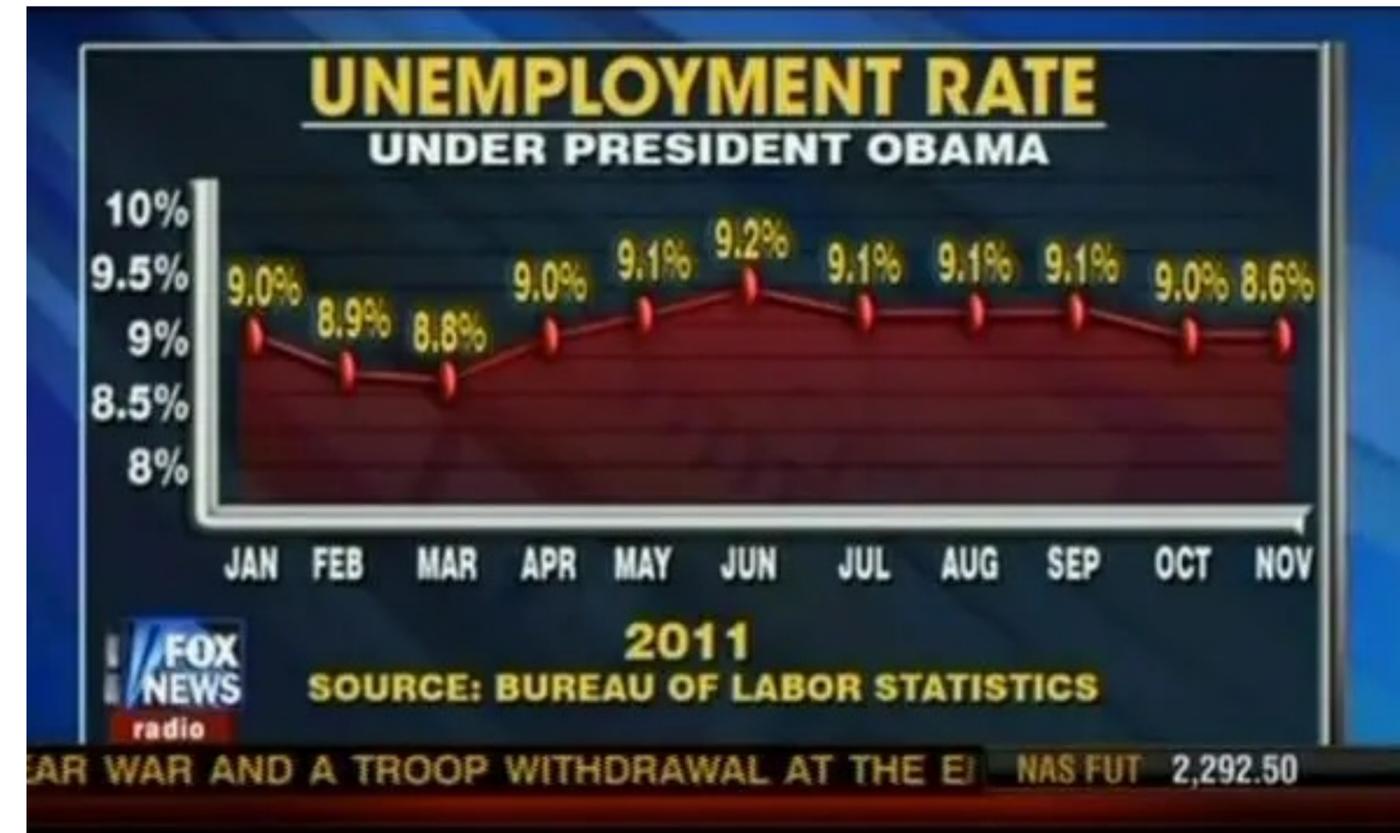
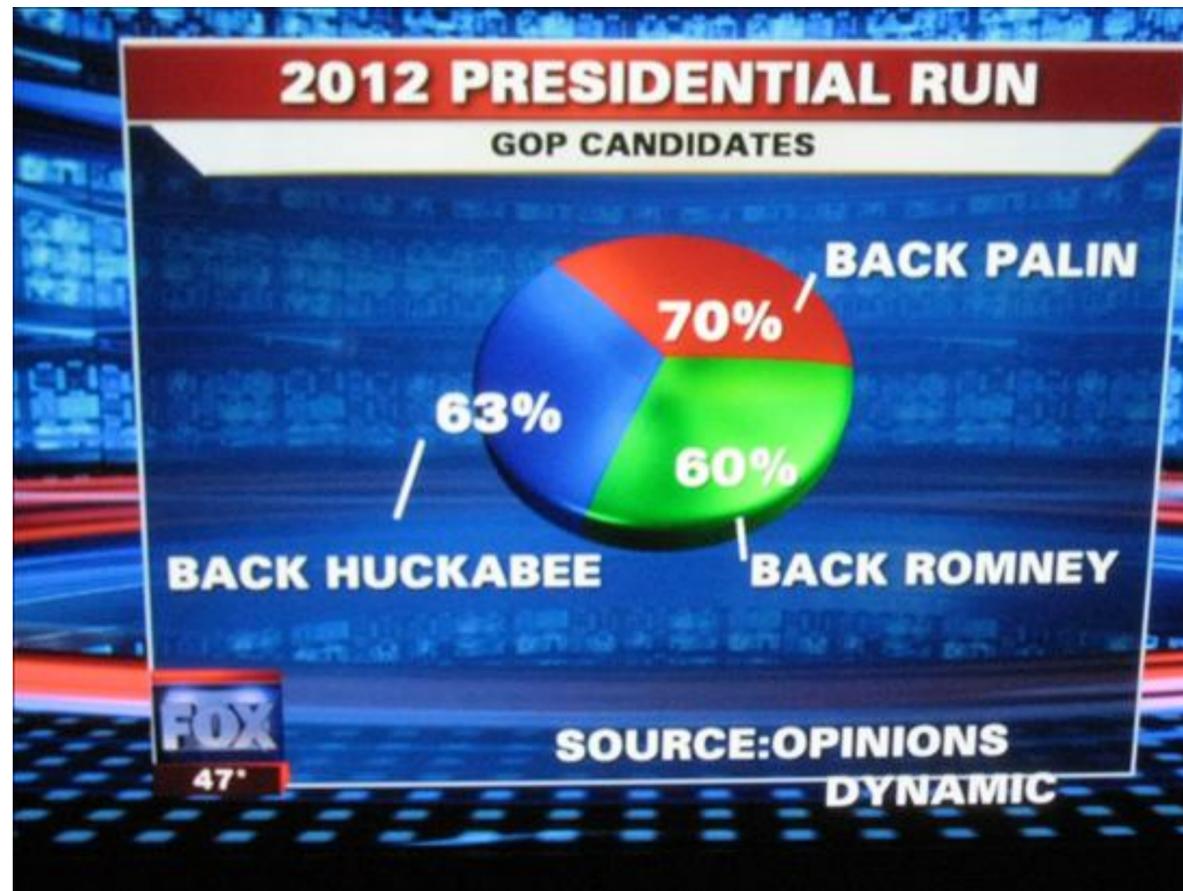
## Two Horse Race

In Edinburgh North and Leith it's a two horse race between Labour and the SNP. The only way to stop the Nationalists is to vote Labour. Only in 2015 Labour was a close second to SNP. Conservatives a poor third, with Labour double their votes. In 2015 the SNP secured half of the Scottish vote, and these official figures show that has now plummeted by 18 points.



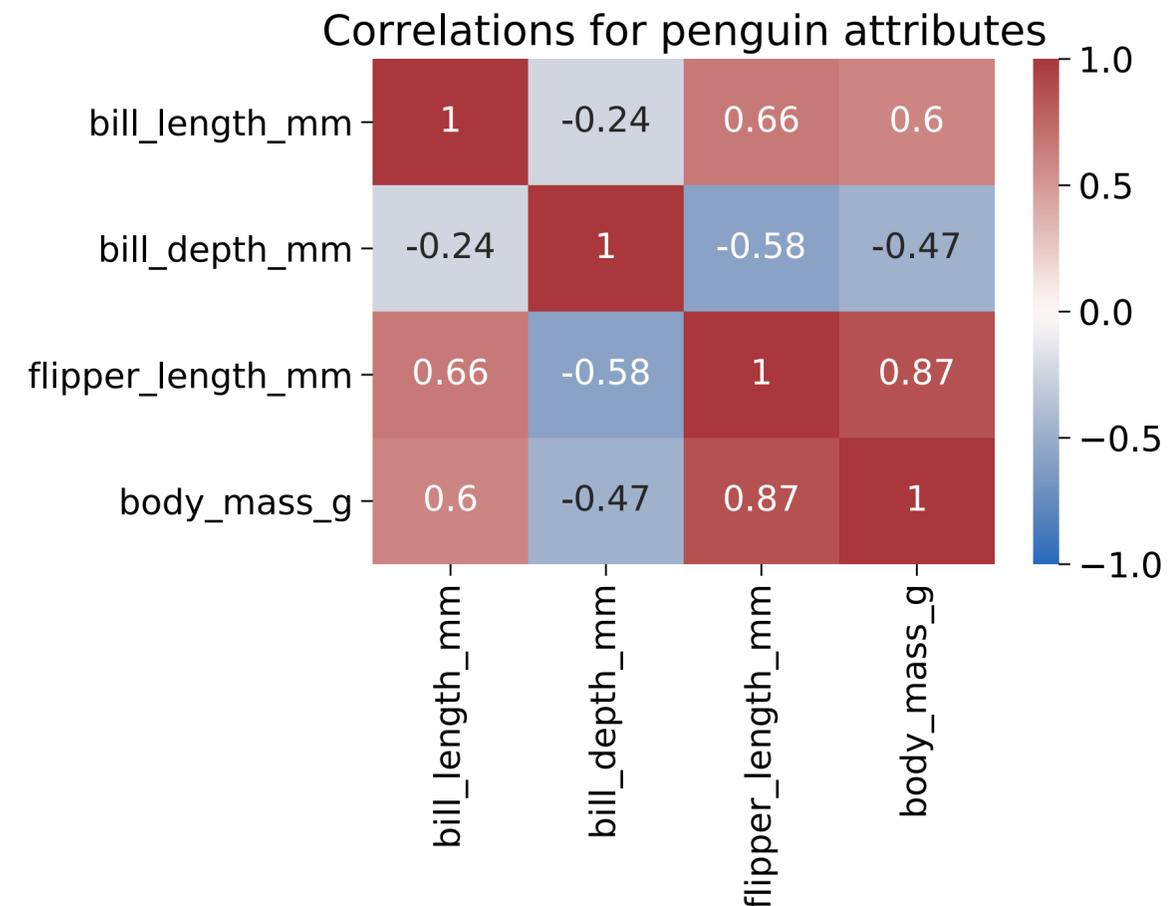
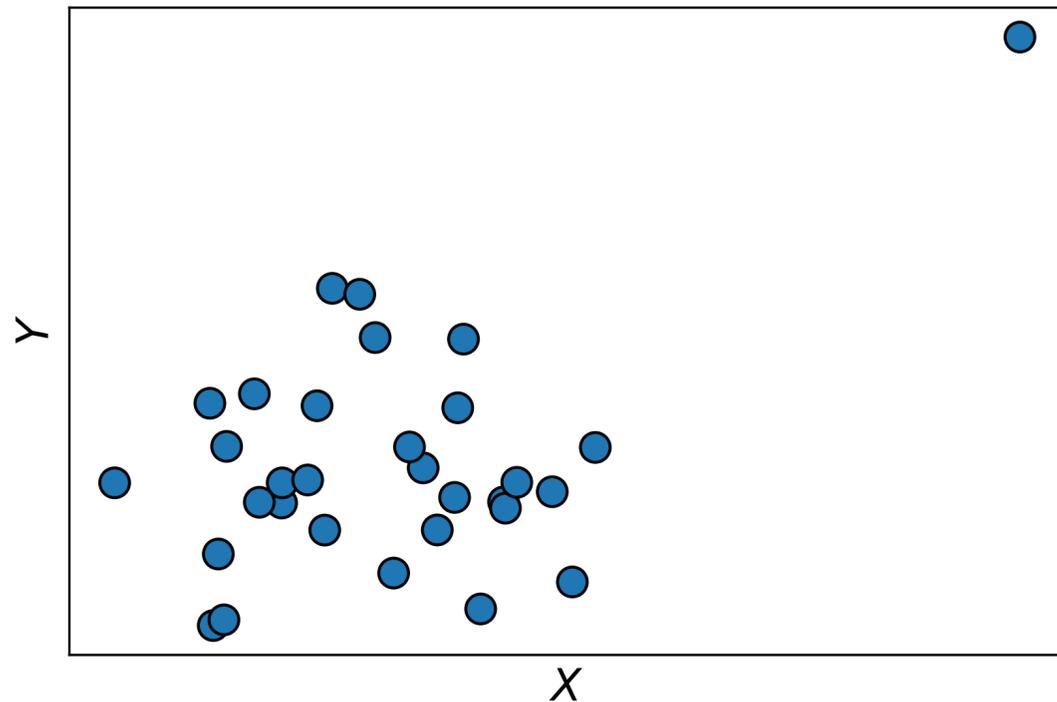
# Visualising data for presentation

Or can just be completely wrong



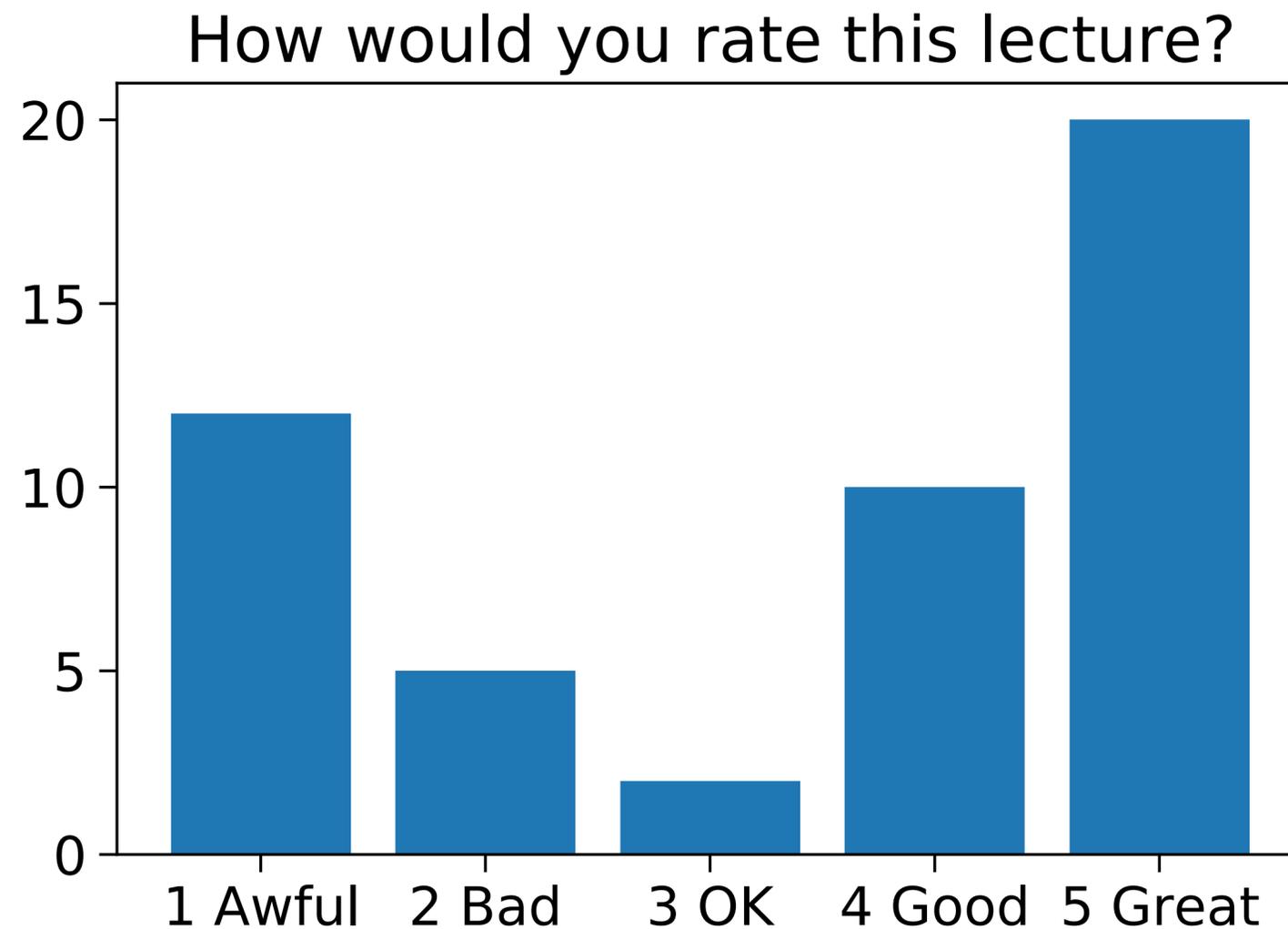
# Visualising data for exploration

- This lets us find patterns, spot outliers/errors, identify important variables...
- It helps us decide which machine learning methods to use (if any!)
- **We want to know if the data makes sense and if it is meaningful**



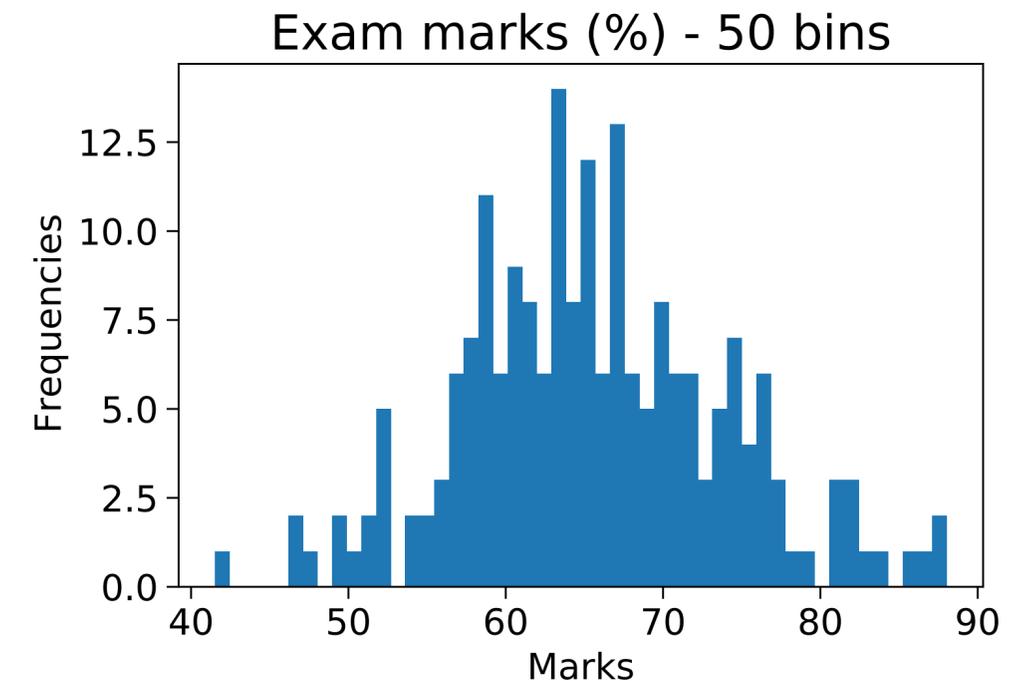
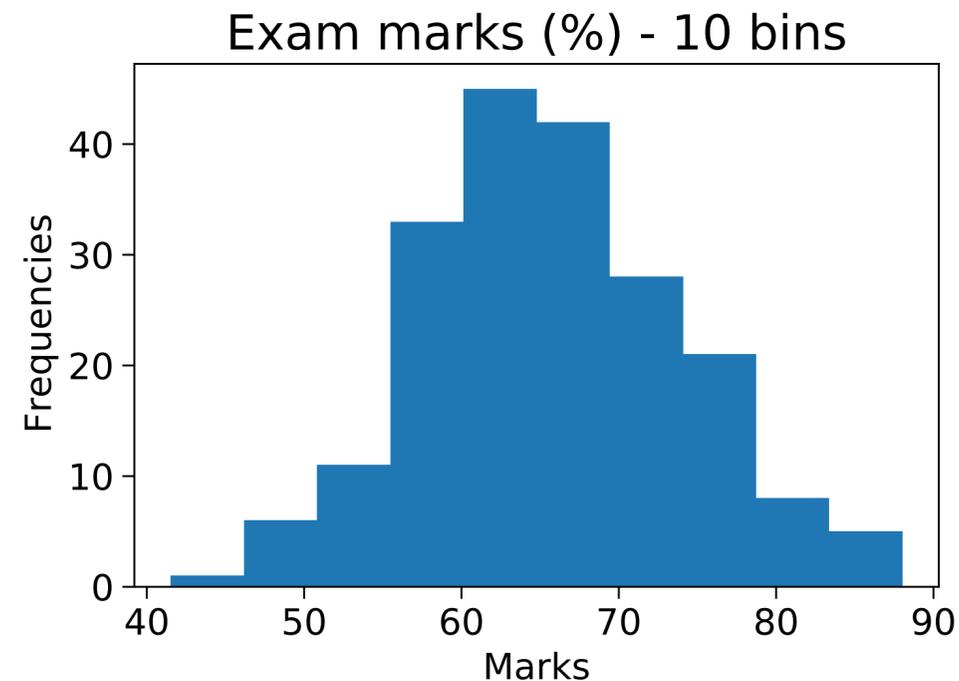
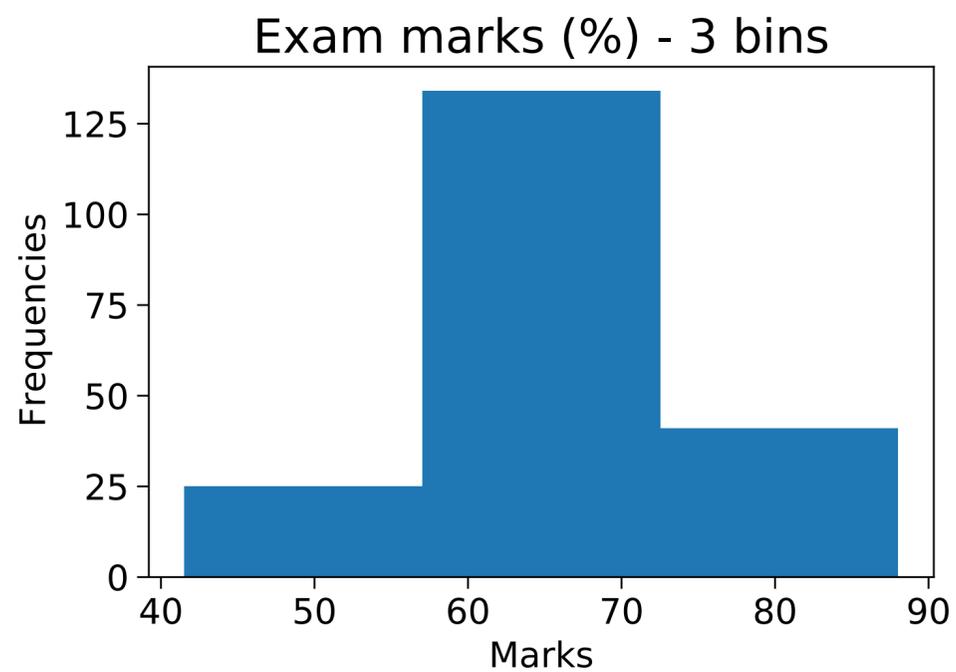
# Bar plots

- Good for visualising categorical variables
- If the variable is ordinal then make sure that the columns are in order



# Histograms

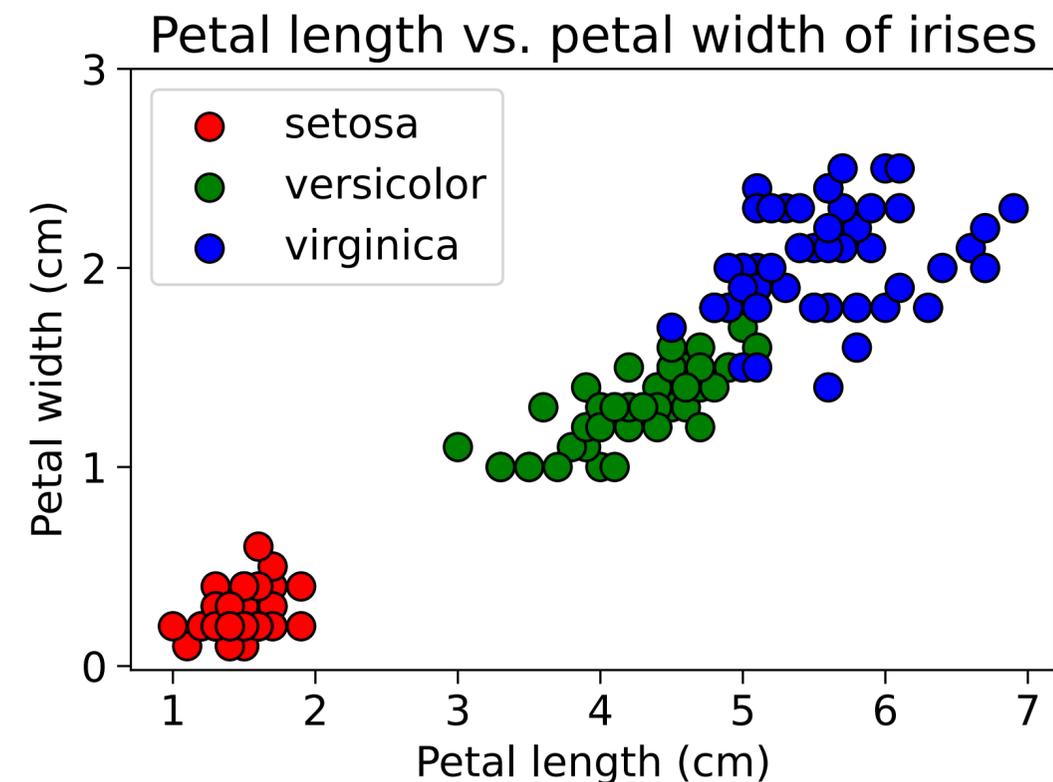
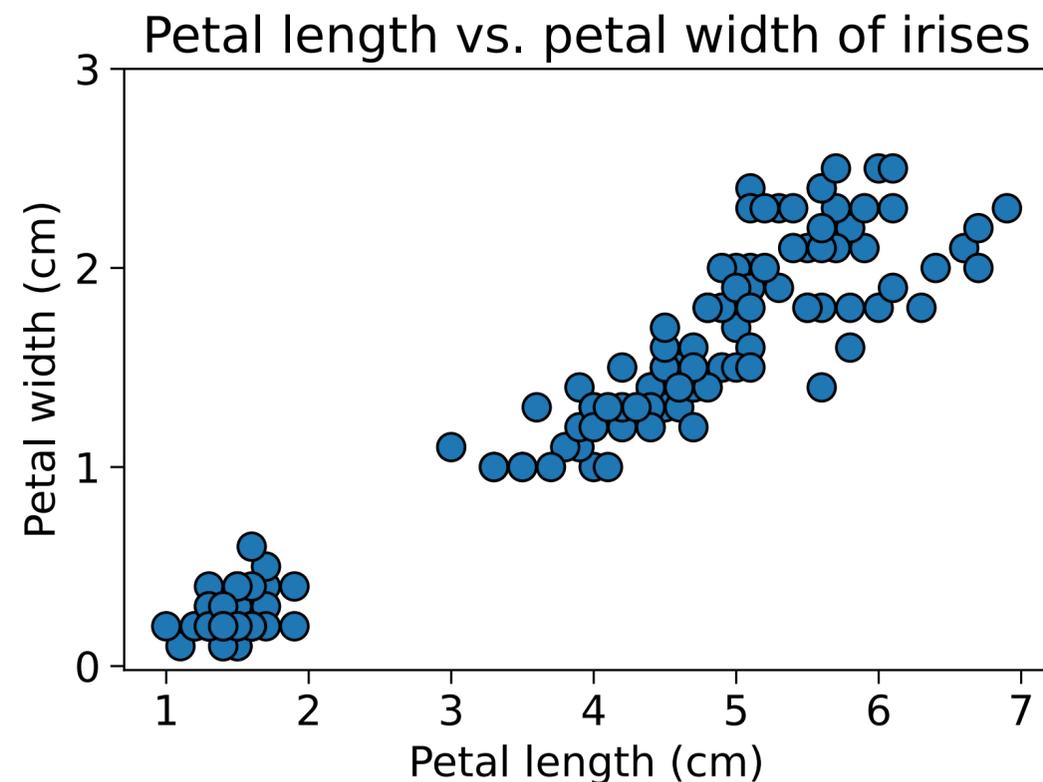
- Sorts measurements for numerical variables into equal sized bins
- The number of bins (and/or bin width) may need tweaking depending on use



There are strange y ticks on this plot.  
This can also be tweaked!

# Scatter plots in 2D

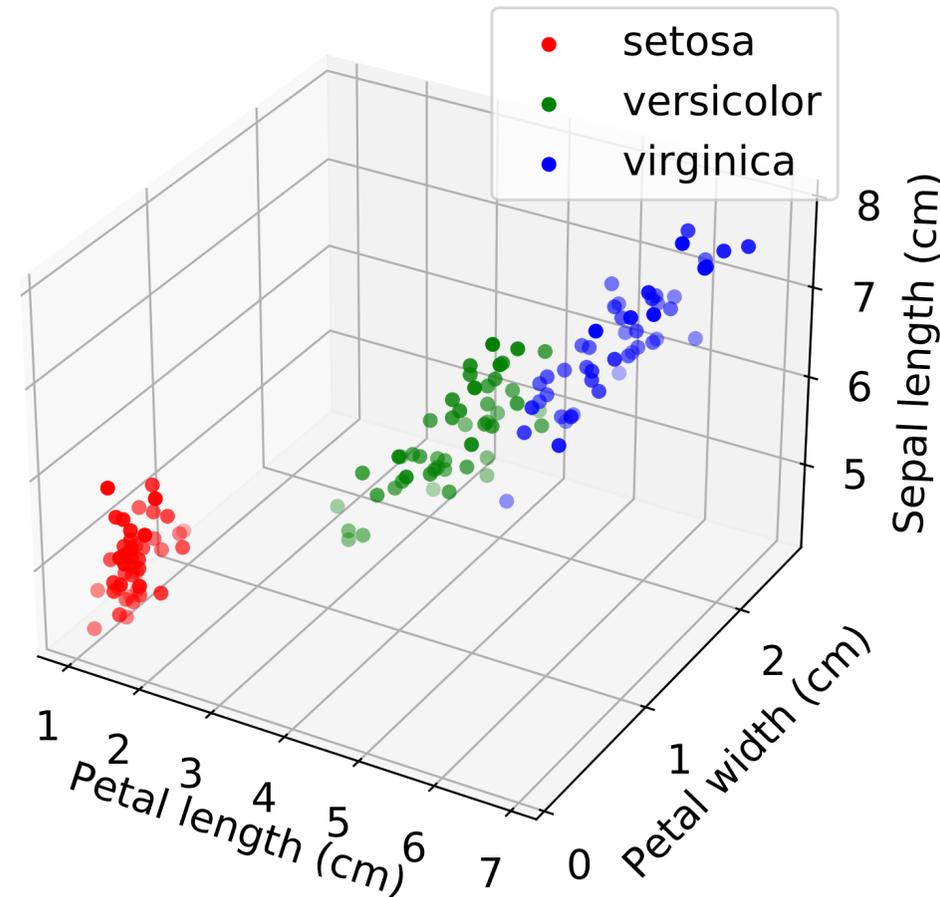
- Each point corresponds to a data item
- The  $x$ ,  $y$  values for that point are measurements of two numerical variables
- We can also distinguish points by category by using different colours/shapes



# Scatter plots in 3D

- We can have  $x, y, z$  values to show three measurements per point
- But beware: we can't see the space properly as it's only a 2D projection :(

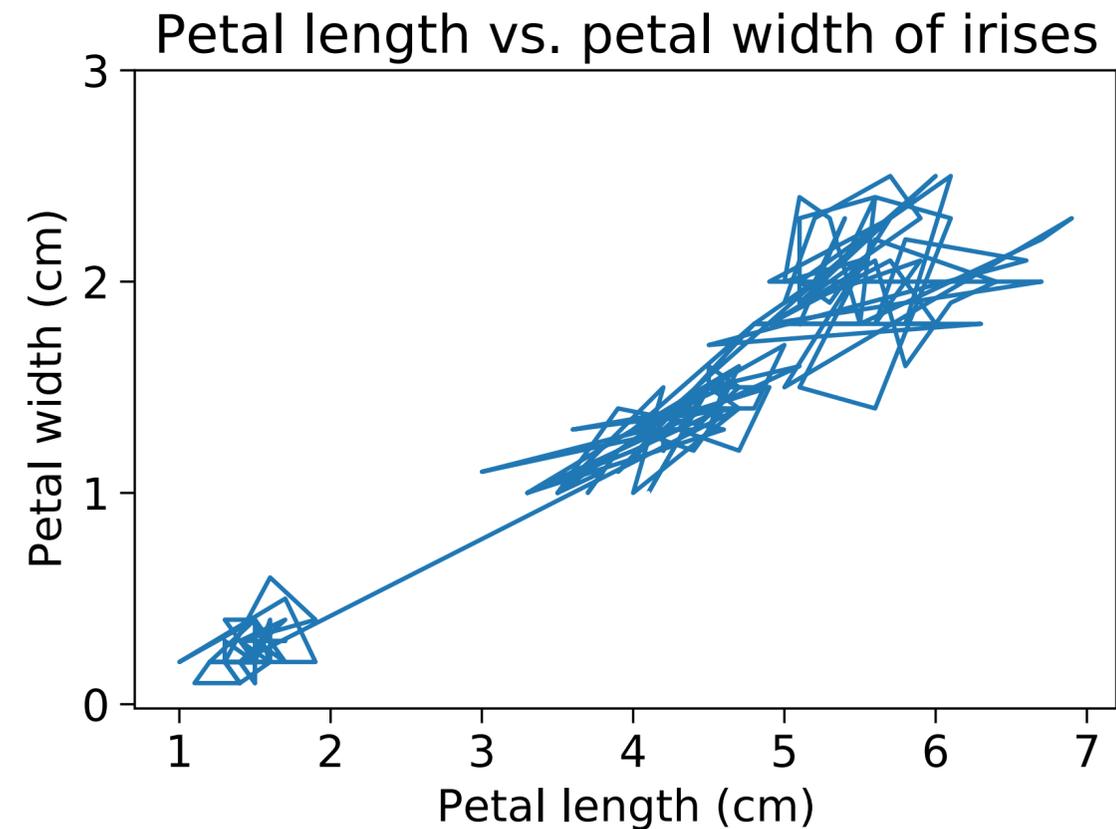
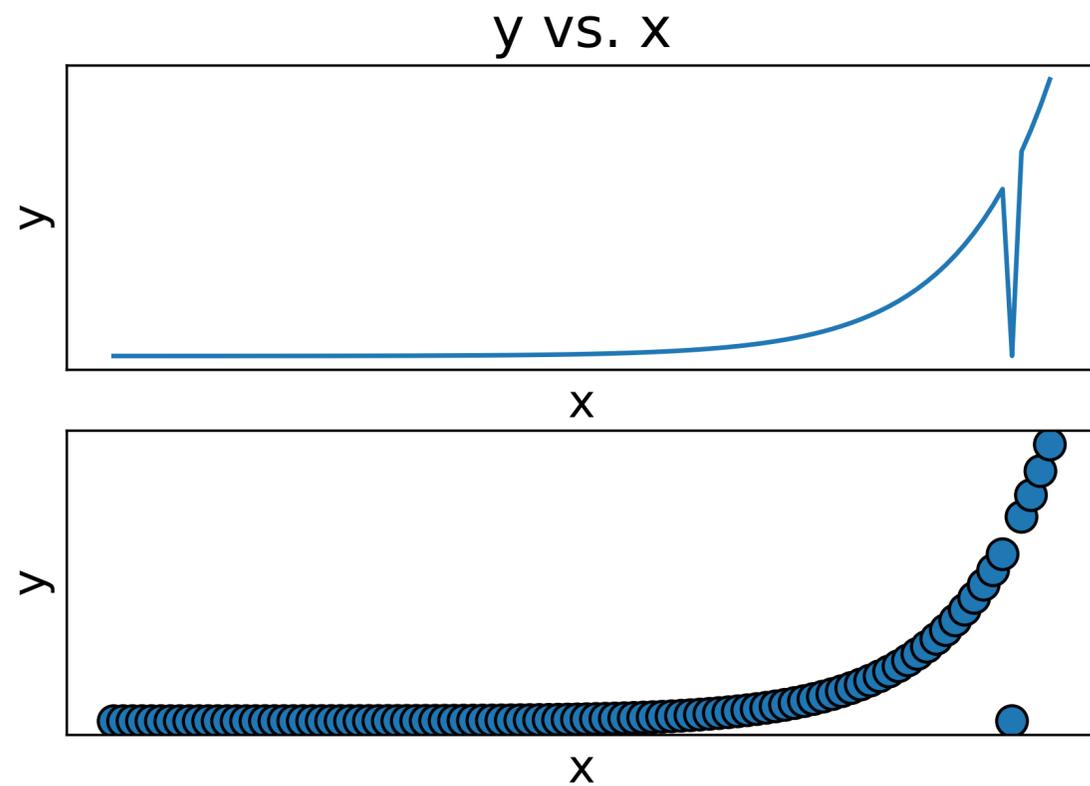
Sepal Length vs. Petal length vs. petal width of irises



I tend to avoid 3D plots where possible

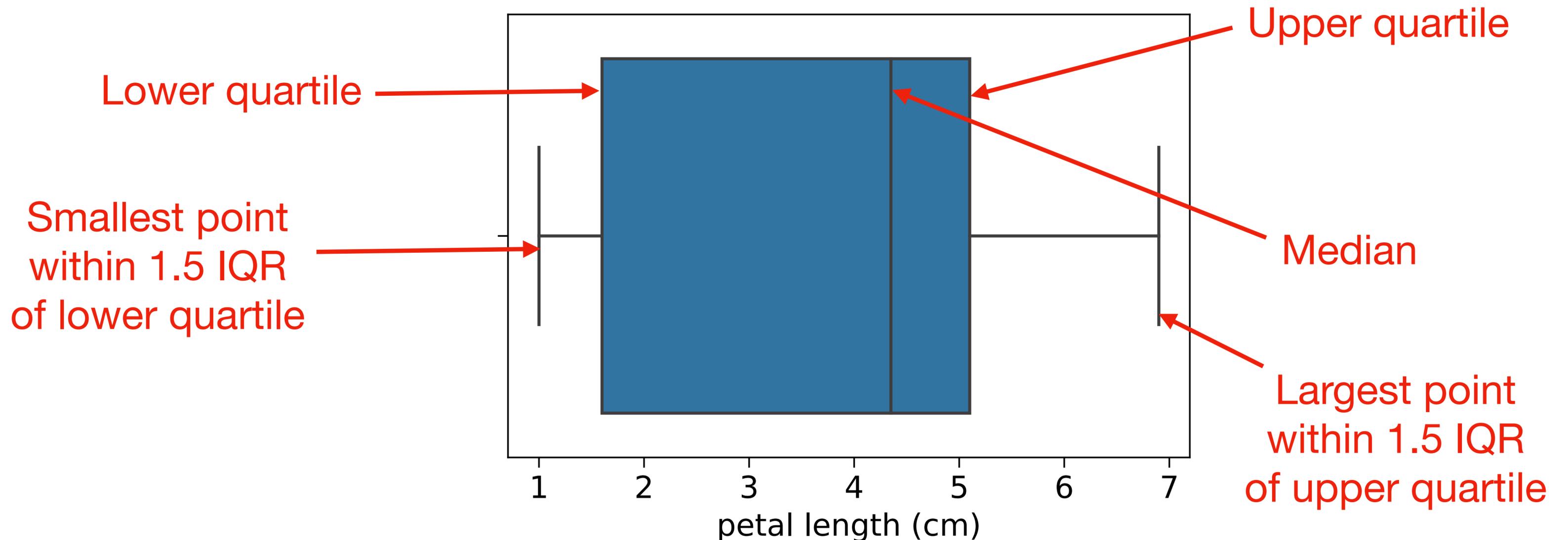
# Line plots

- Can be useful for interpolation
- But can depict a functional relationship that doesn't exist if used carelessly



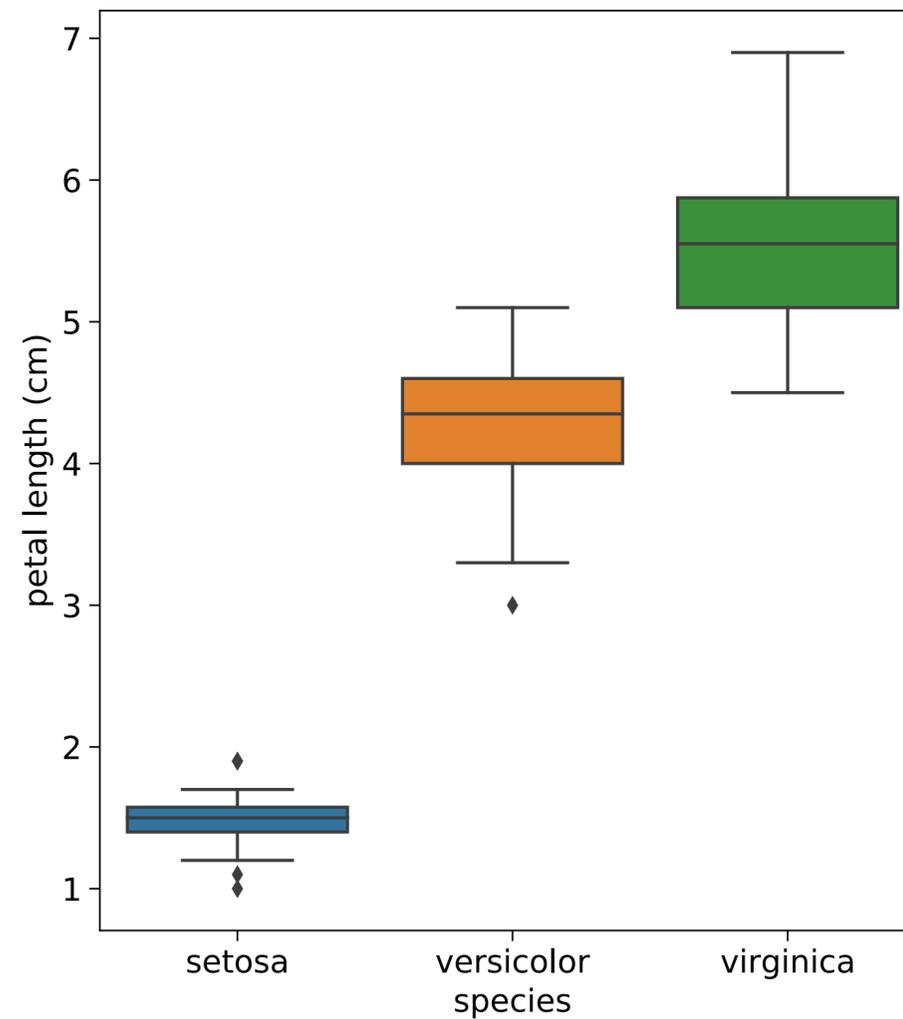
# Box plots

- Shows 5 key statistics of a variable, each being an actual measurement
- Interquartile range (IQR) = upper quartile - lower quartile



# Box plots

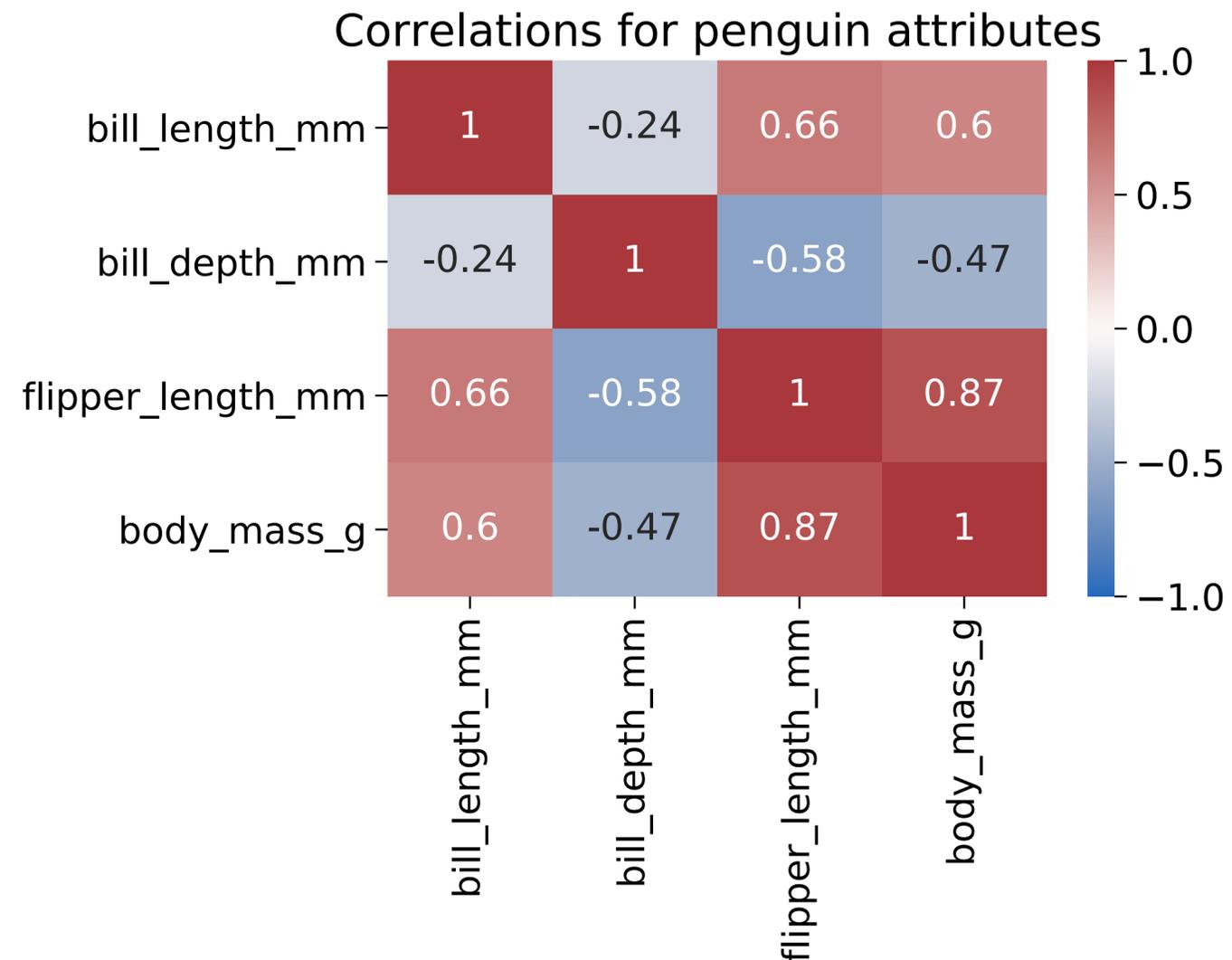
- We can view these statistics split by category
- Any points outside of the *whiskers* are plotted



Plot can be  
horizontal or  
vertical

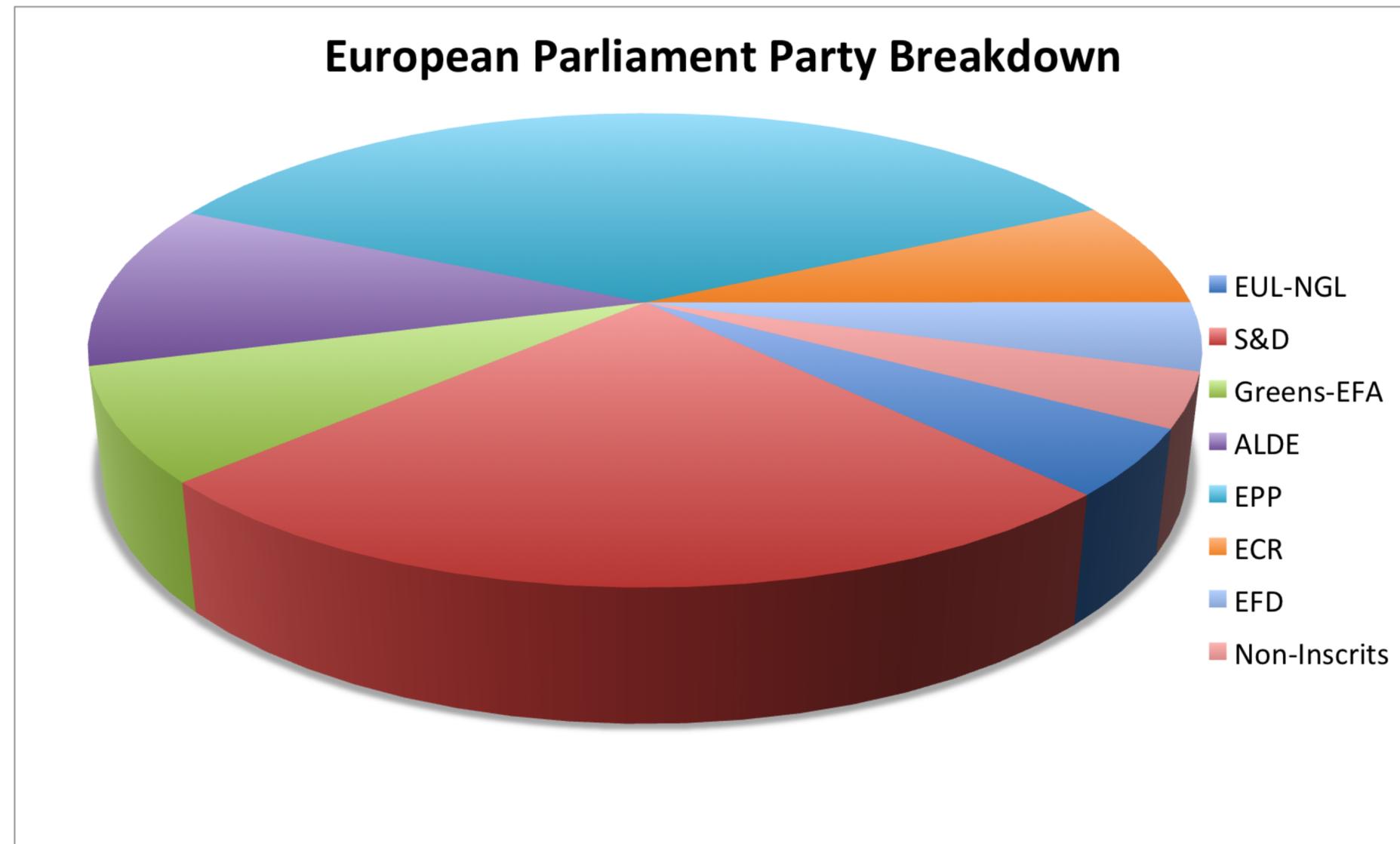
# Heat maps

- A matrix of colours representing different magnitudes of some quantity
- Here we have Pearson correlation of different attributes of penguins



# And of course ... pie charts

Avoid!



Source: <https://www.businessinsider.com/pie-charts-are-the-worst-2013-6?r=US&IR=T>

# Summary

- We have revised some statistics and seen how they can summarise data
- We have considered correlations for different pairs of variables
- We have seen examples of good and bad visualisations of data
- We have considered different ways of plotting data