# BMS PhD Bitesize Training
# Introduction to Statistics

30/05/23

Andrew Mason

andrew.mason@york.ac.uk
asmasonomics.github.io
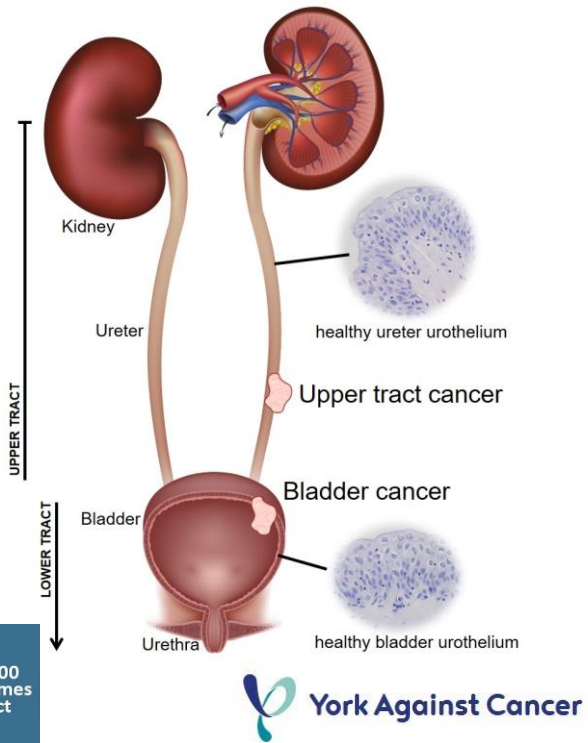@asmasonomics

**Andrew Mason**

Lecturer in Cancer Informatics
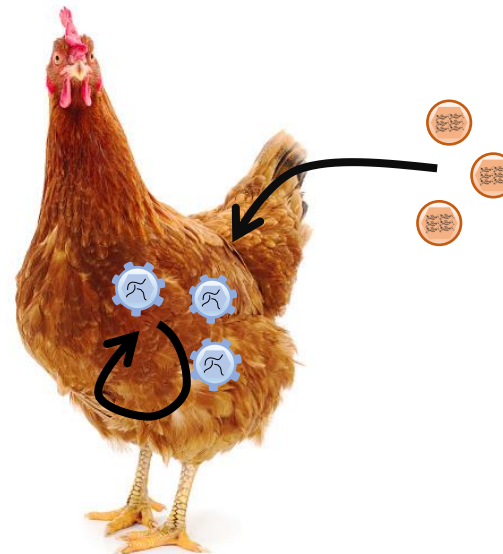
Elixir Data Stewardship Fellow

Bladder cancer bioinformatic lead for the 100,000 Genomes Project

Endogenous retrovirus lead for chicken reference genome, pangenome and diversity consortia

Urothelial carcinoma in humans

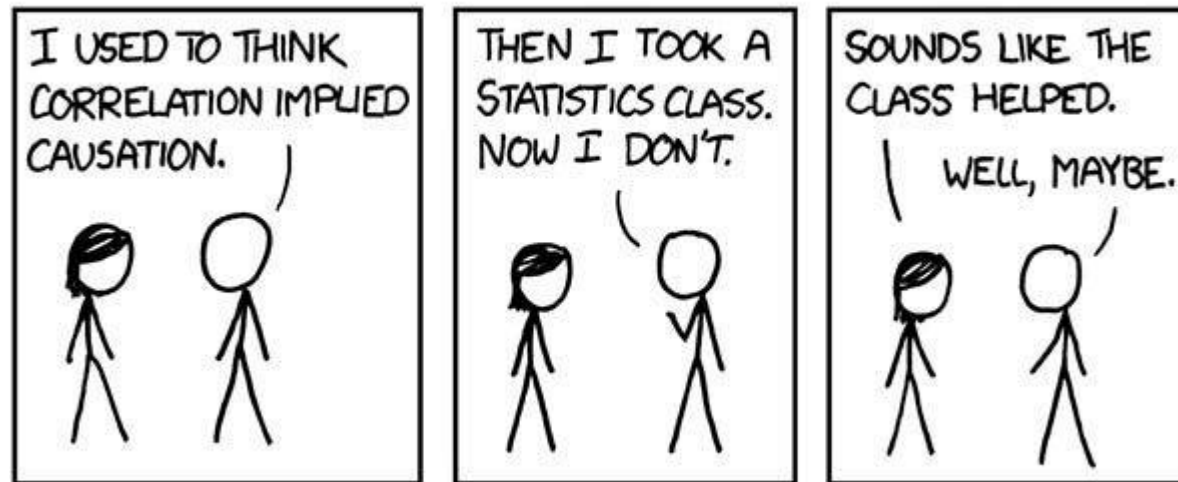Oncogenic viruses in chickens

Data science training in the life sciences: UG, PG, Academics
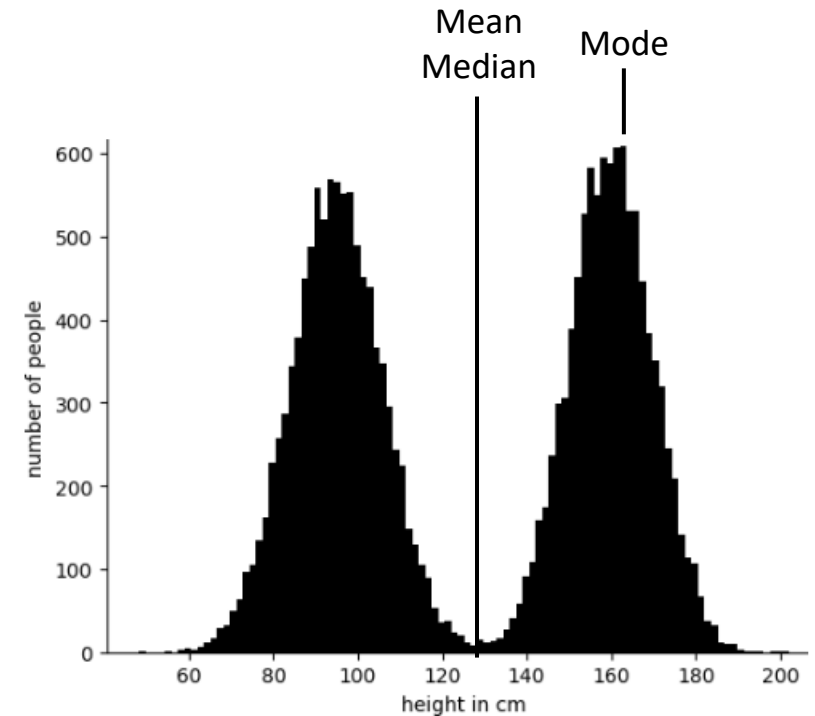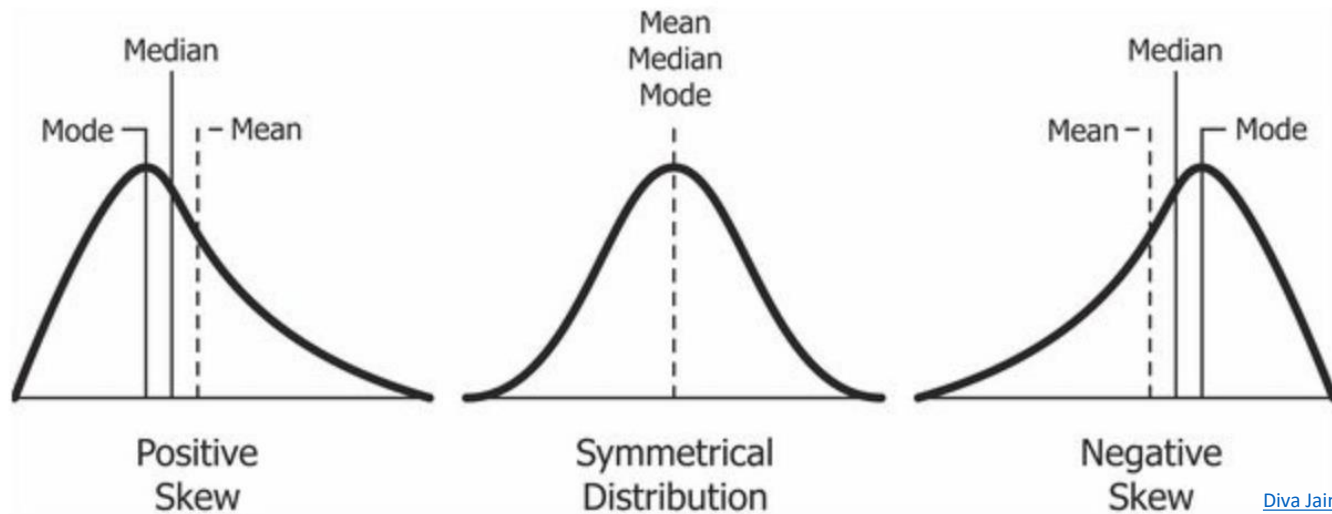
Introduction to statistics:

- Better idea of the statistical test to use

- Appreciation of some common pitfalls

- Make the most of your data!

- Nail those stats questions in your viva

- (Maybe) know more about statistics than your supervisor…
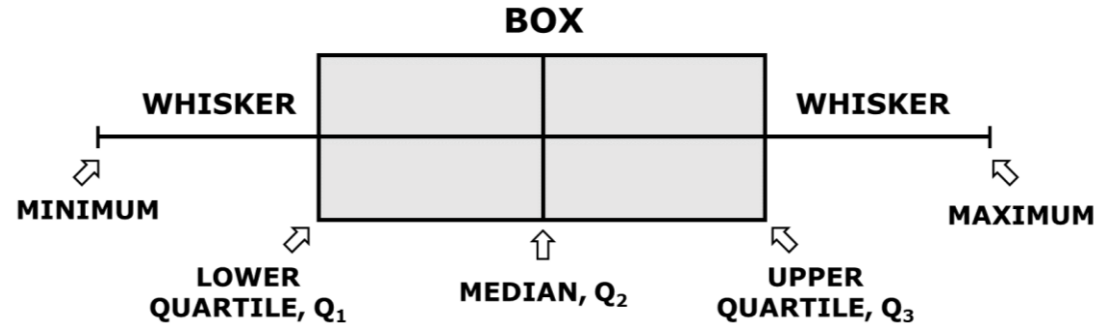
Mean – the "traditional" average (sum/n)

Median – the value in the middle of your dataset
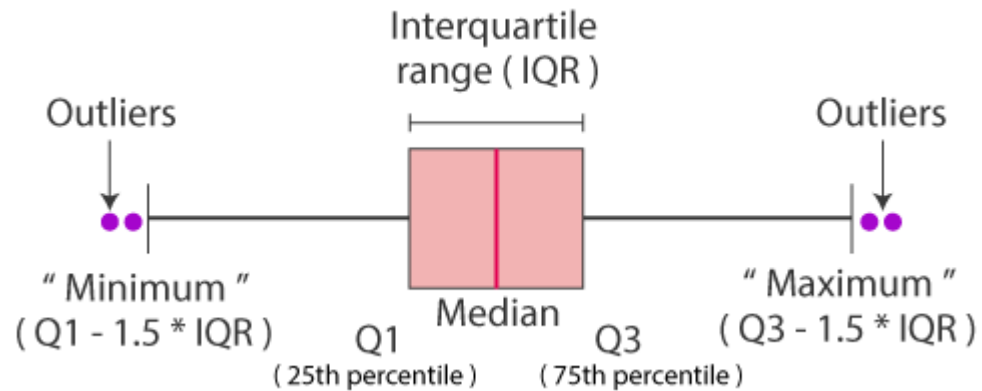
Mode – the most common value in your dataset



Diva Jain

**Always graph your data!**
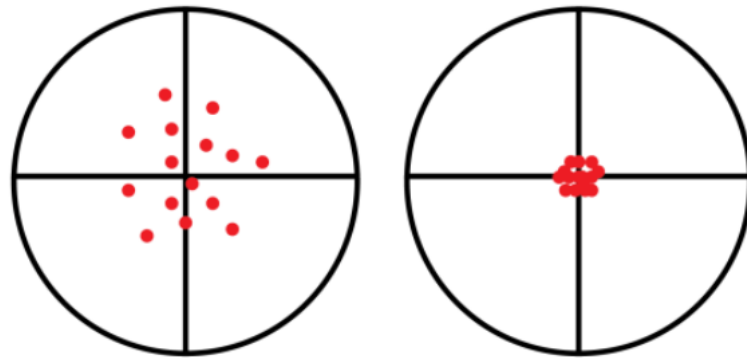
Box plots



One way to determine outliers



**Descriptors**
- Range (w/wo outliers)
- IQR
- 95% confidence intervals

**Mean – the "traditional" average (sum/n)**

How accurate have we been in measuring the mean of the population?



**Standard deviation** – how much do the underlined <u>individual measurements</u> differ from the mean value?

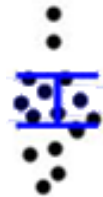Low SD gives us more confidence in our assessment of the mean value

Standard deviation vs Standard error **of the mean**

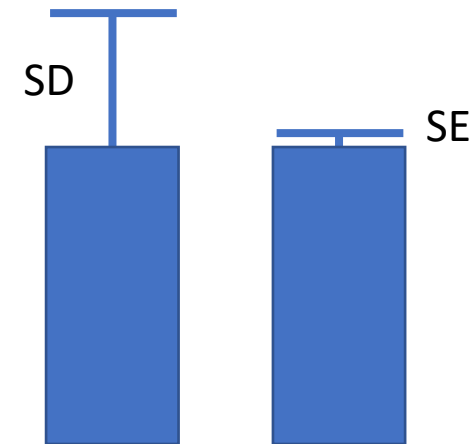$$SD = \sqrt{\frac{\sum |x - \mu|^2}{N}}$$

SD – how well do my **individual** measurements support **my observed mean** value?
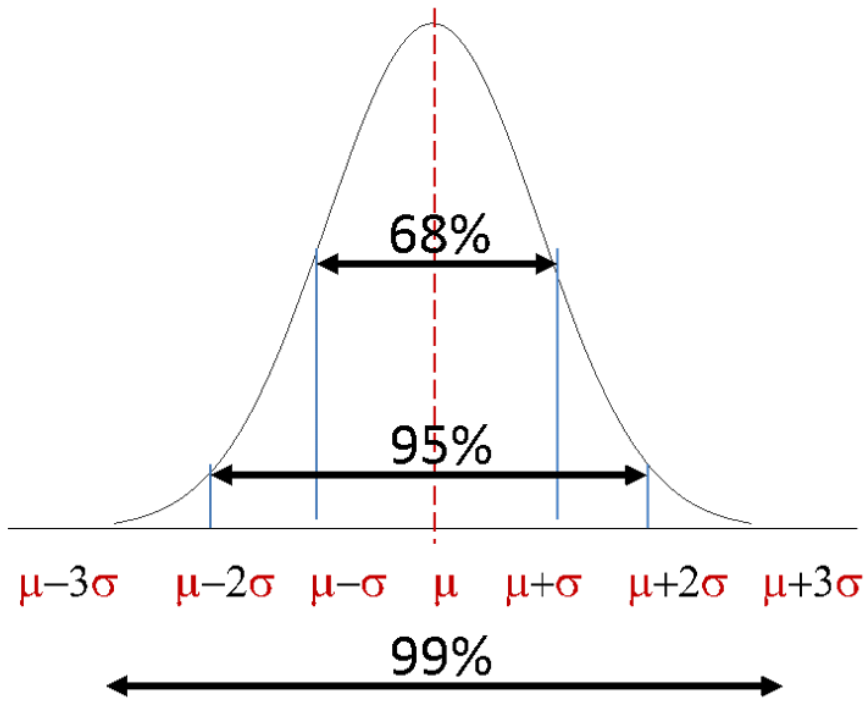
$$SE = \frac{SD}{\sqrt{n}}$$

SE – how well do my **repeat** measurements support the **actual population mean** value?

SD

SE

Normal distribution = Gaussian distribution



Mean determines the peak of the curve
SD determines the shape of the curve

95% of the data points fall within +/- 2SD from the mean

100% - 95% = 5% = 0.05

This is where the significance threshold comes from…

Parametric testing requires data to be normally distributed (ish)

- Equal variance between groups

- Groups are independent measures

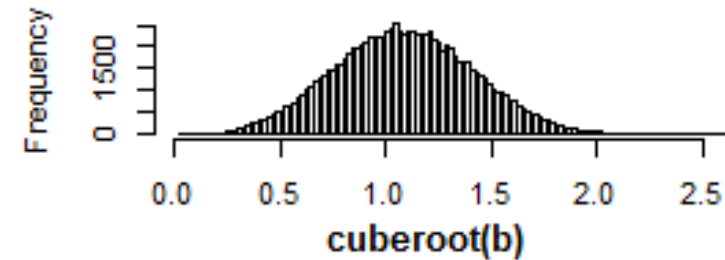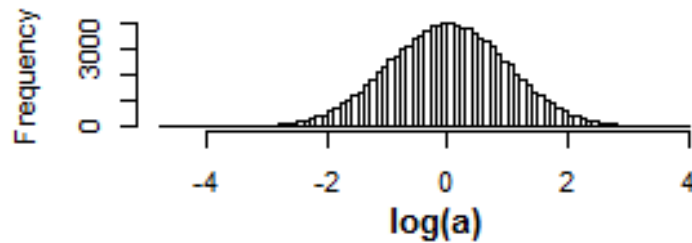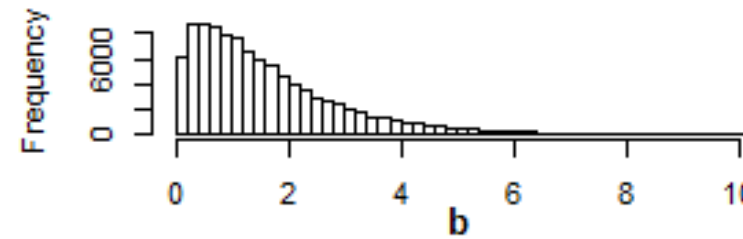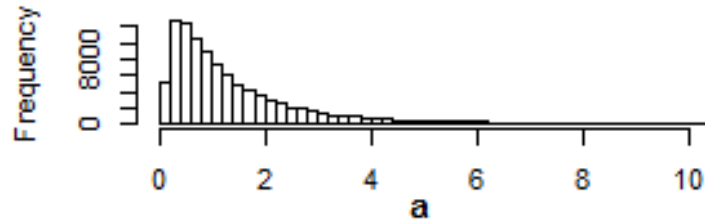- No distribution altering outliers

**t test**

**ANCOVA**

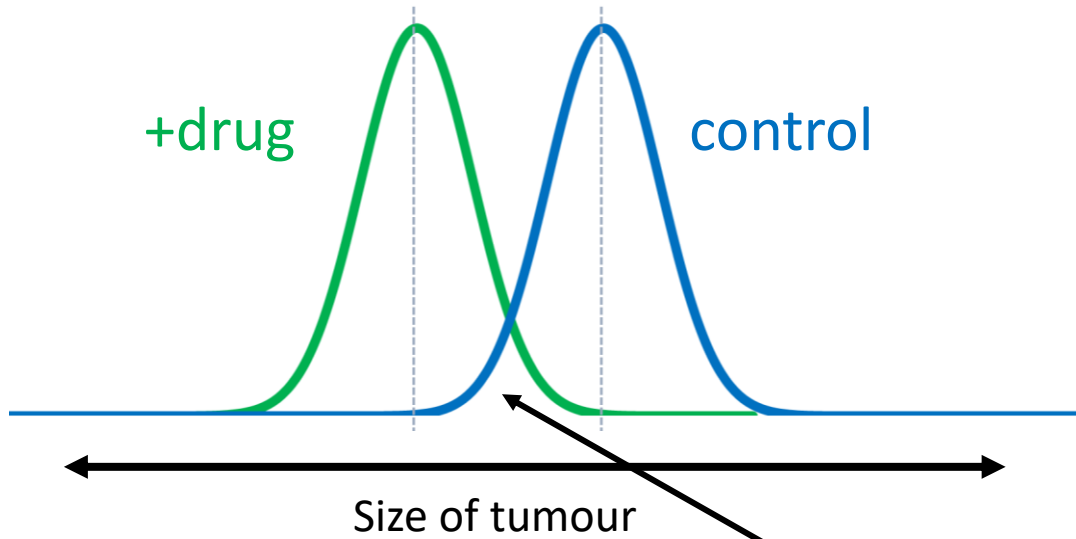**Pearson correlation**

**ANOVA**

Data transformations are legit!



Just be careful when:
- reporting results (is it understandable?)
- considering outliers
- drawing error bars

**NON PARAMETRIC TESTING DOES EXIST!!**

# t test logic



+drug    control

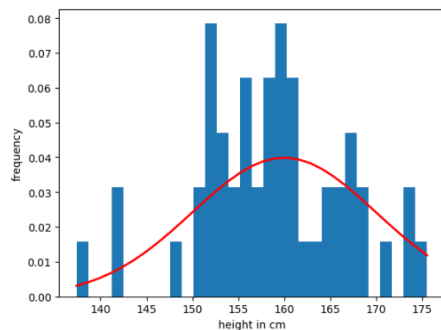Size of tumour

Has the tumour significantly decreased in size with addition of the drug?

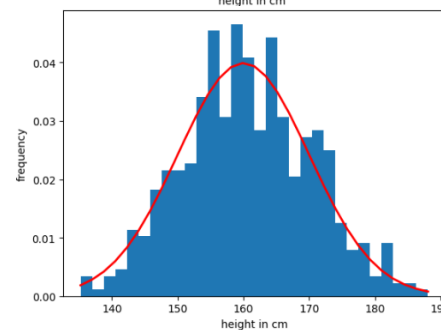What is the overlap here, what proportion of the data could be found in <u>both</u> distributions?
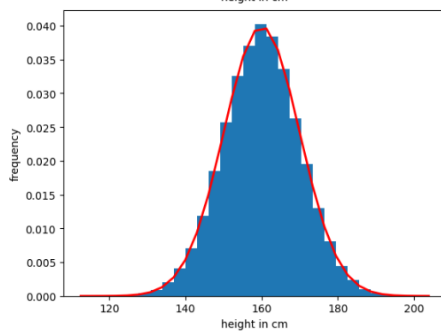
1) Small n vs big n

e.g. height dataset - mean 160cm, SD 10cm

n=50

n=500

n=50,000

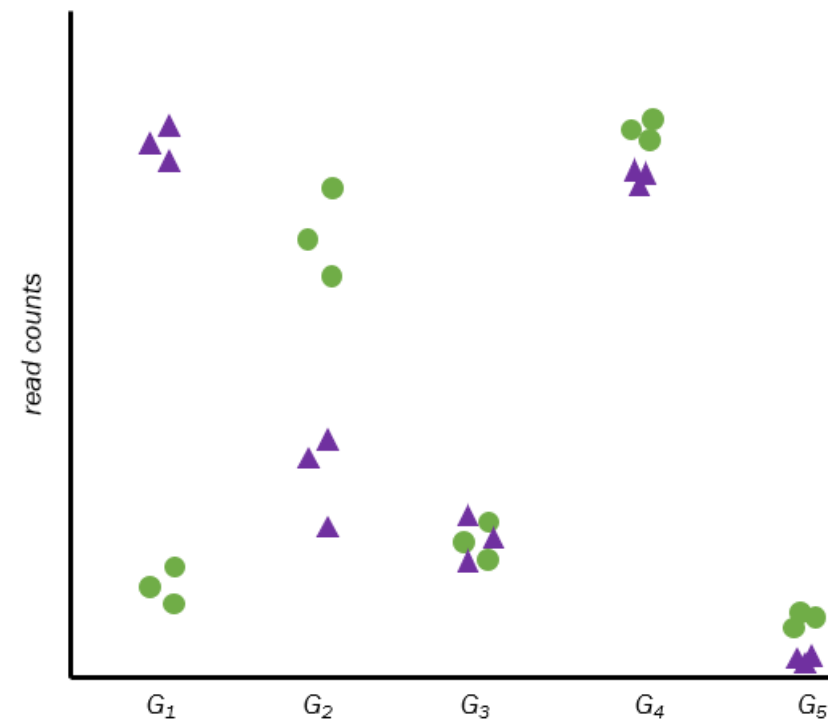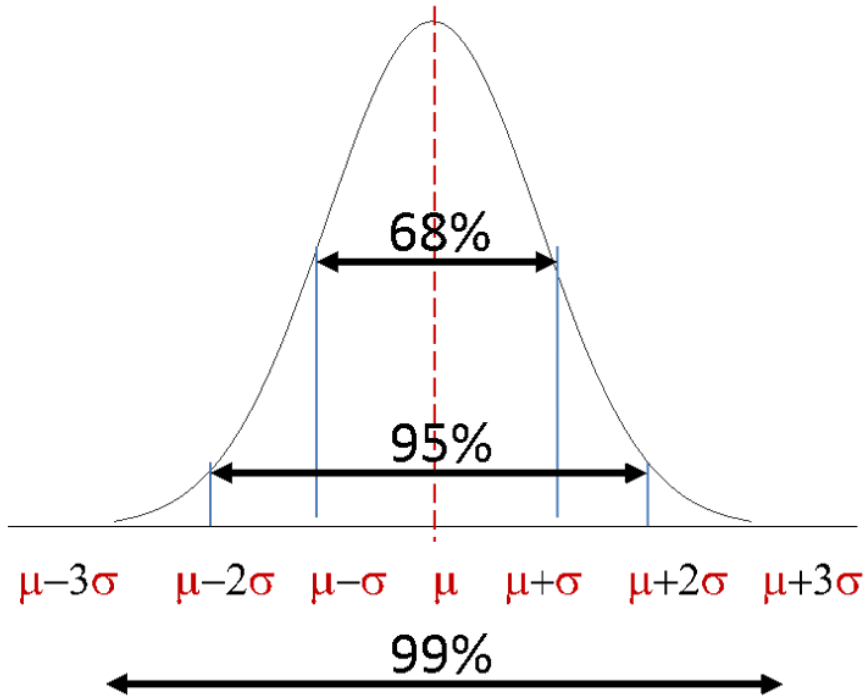2) Biological vs Statistical significance

Normal distribution = Gaussian distribution



68%

95%

$\mu-3\sigma$    $\mu-2\sigma$   $\mu-\sigma$   $\mu$   $\mu+\sigma$   $\mu+2\sigma$   $\mu+3\sigma$

99%



5%

2.5%     2.5%

0   1.645       -1.96   0   1.96

(a) One-tailed test      (b) Two-tailed test

Null hypothesis is typically that "there is no difference"
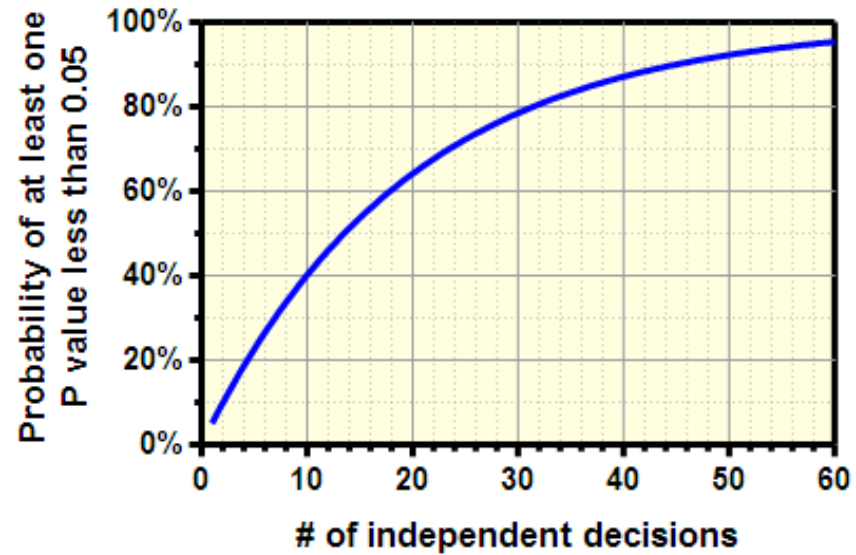
Only use a one-tailed test when there is no possibility of there being the alternate direction of response
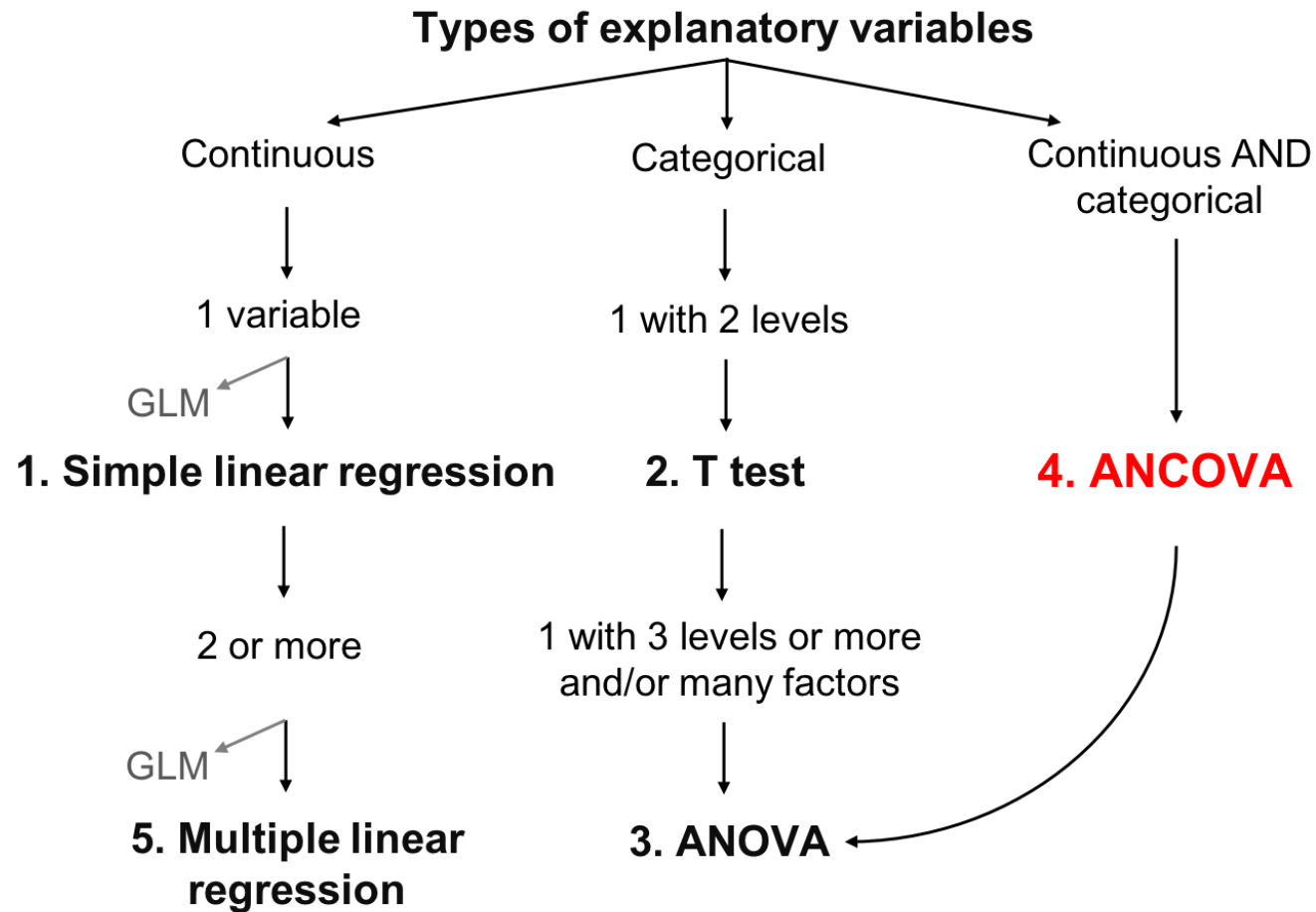
Fishing exercises vs the art of data science

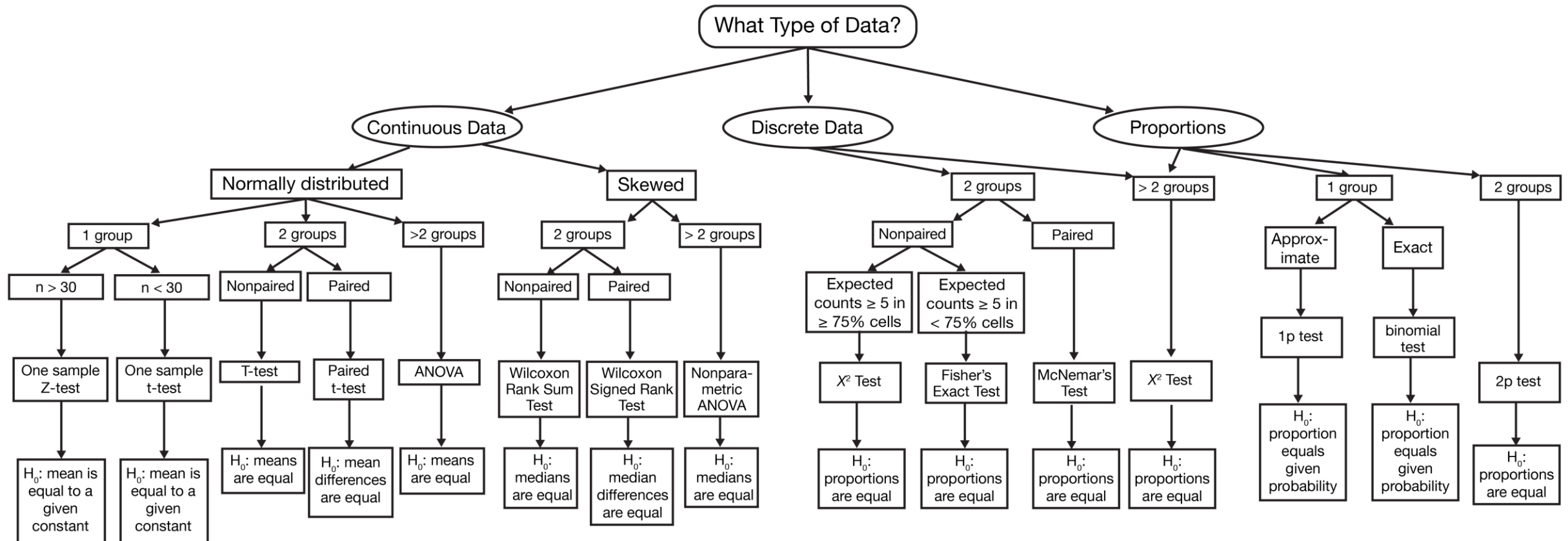Increasing chance of a false positive when you do lots of tests

Bonferroni correction
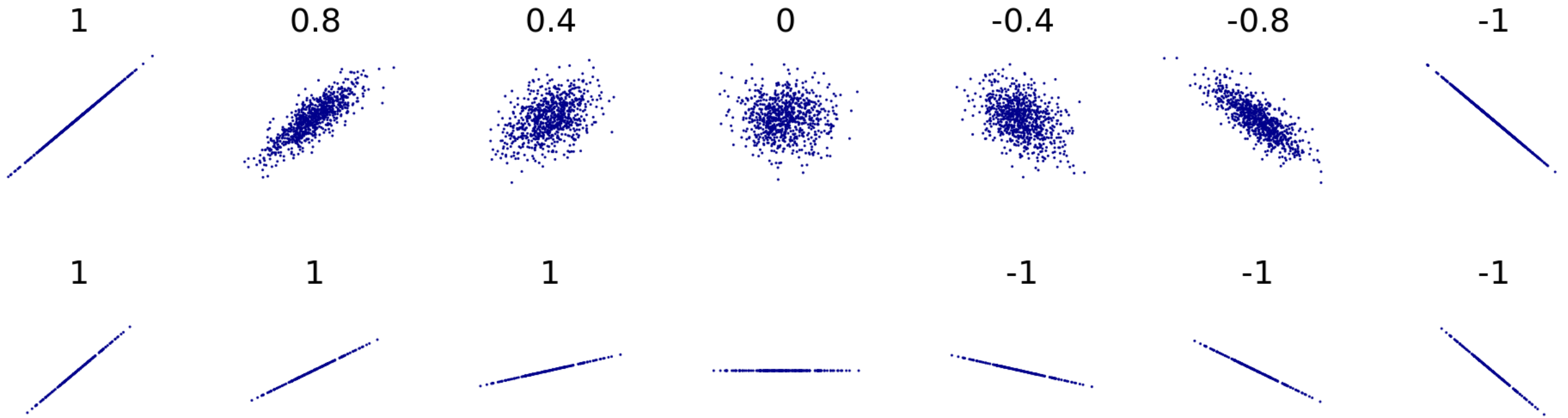Benjamini-Hochberg
Etc etc

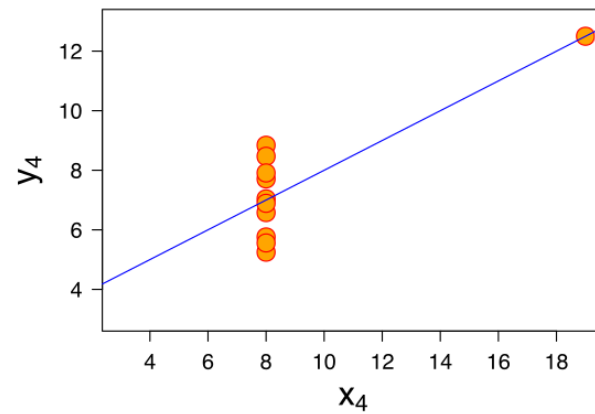t test vs ANOVA vs ANCOVA (and what is post hoc testing?!)
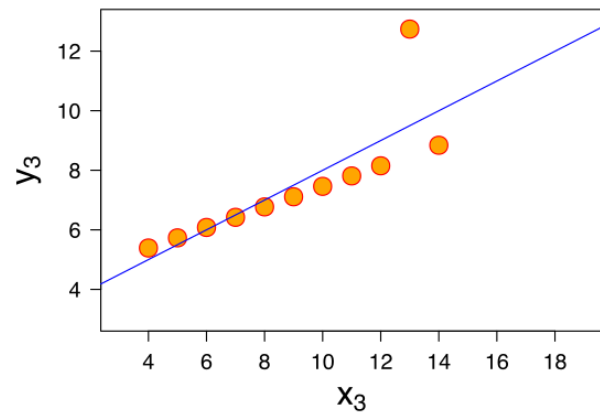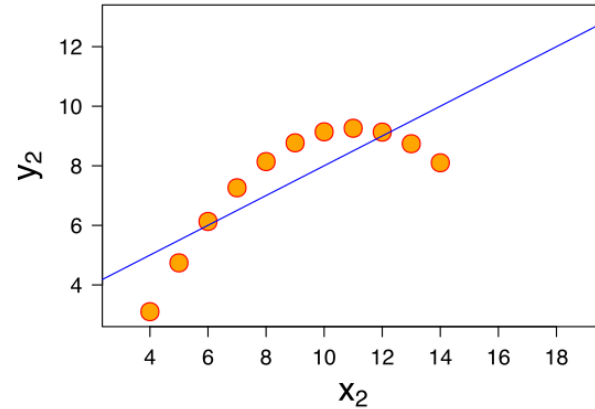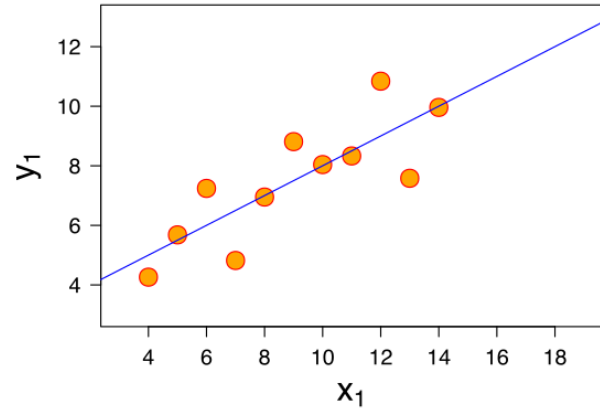
Flow chart: which test statistic should you use?

Cautionary tale – Anscombe's quartet



Same considerations for parametric vs non-parametric
- often Pearson vs Spearman

https://www.socscistatistics.com/

**How often will you be doing stats like this?**

This presentation, including a recording, will (shortly) be on my website:

https://tinyurl.com/BMS-Stats-May2023

Materials for previous "Introduction to cBioPortal" course:

https://tinyurl.com/Intro-cBioPortal-Jan2023

Elixir Research Data Management "Bites" on sequencing data:

https://tinyurl.com/RDM-Seq-Videos