

A large graphic consisting of three overlapping, curved leaf-like shapes. The leftmost shape is a dark green, the middle one is a light green, and the rightmost one is a yellowish-green. They all curve towards the right.

WEO Research Workshop

Seoul, 17 November 2018





Biostatistician's role in study design

Dr Sunny H Wong

MBChB, DPhil, FRCPEd, FRCPath, FHKCP, FHKAM

Assistant Professor
Institute of Digestive Disease
The Chinese University of Hong Kong



Breath of biostatistics

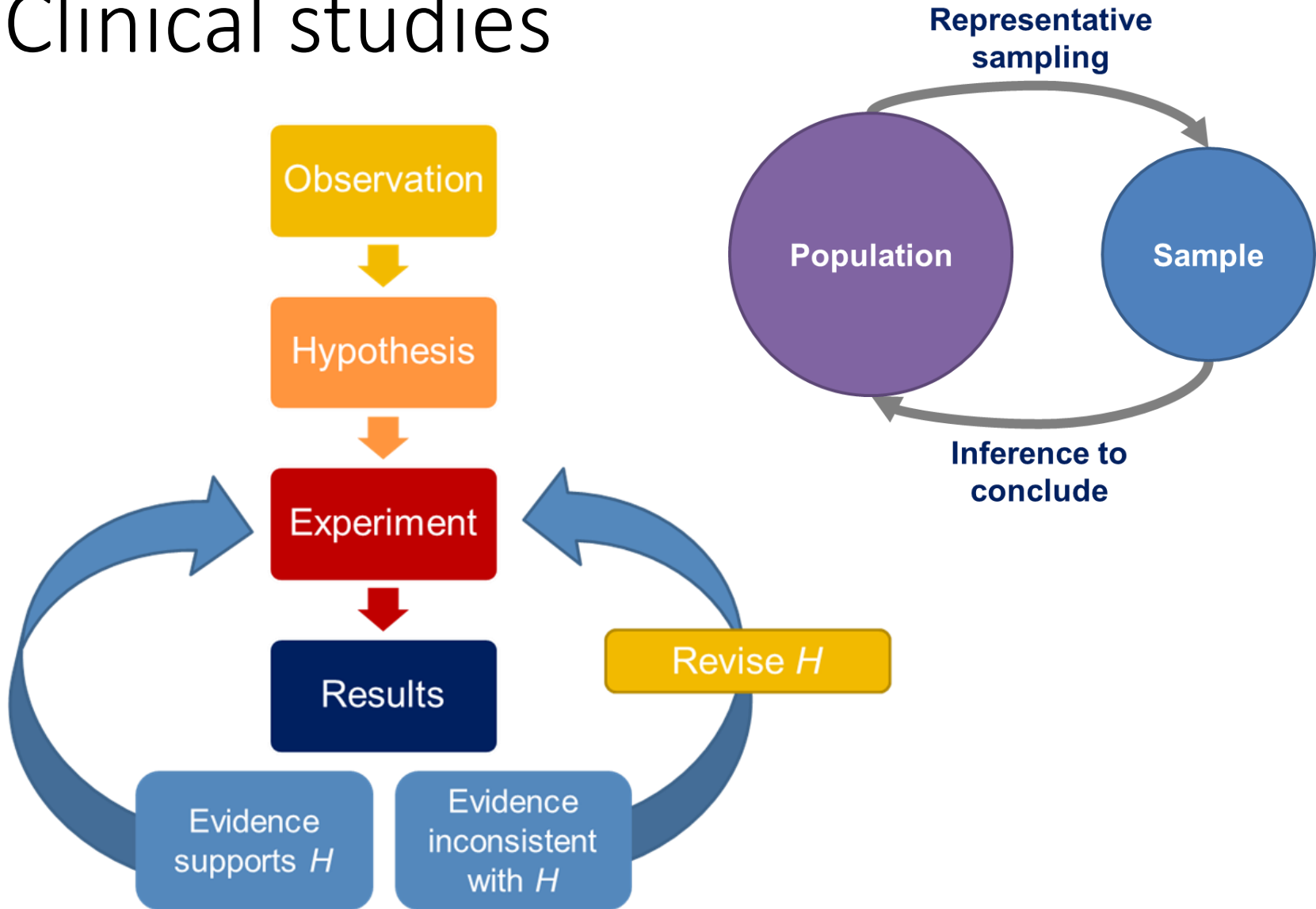
- Descriptive statistics (data types, central tendency, dispersion, exploratory data analysis)
- Probability distributions & confidence intervals
- Hypothesis testing (null hypothesis, type I and II errors, sample size, power)
- Inferential statistics (t-test, chi-square, trend test, Fisher's test, log rank test, comparative data analysis)
- Correlation & regressions
- Multiple comparisons & corrections
- Survival analysis
- Meta-analysis
- Bayesian statistics
- Others (diagnosis, public health, bioinformatics)

Outline

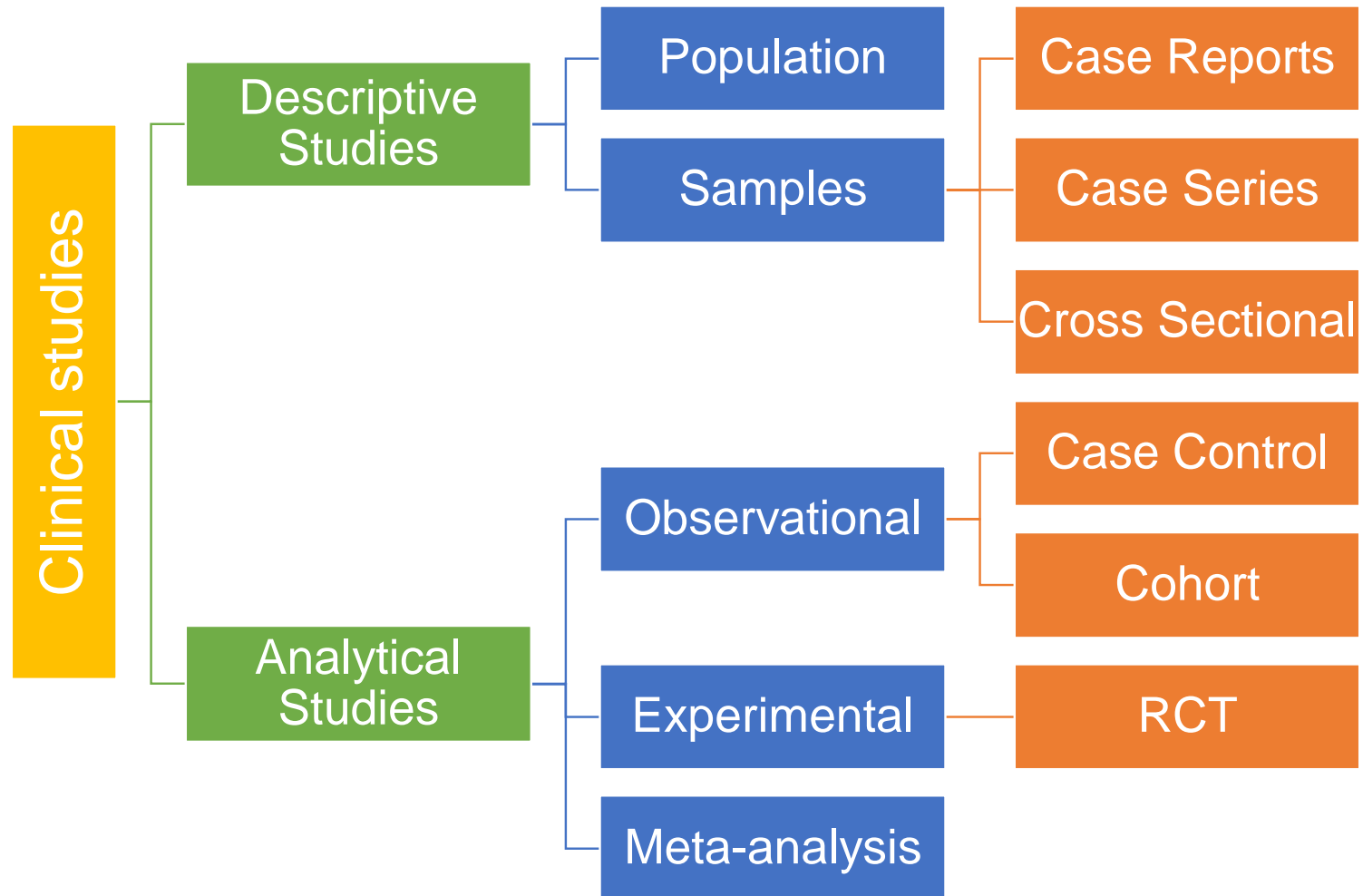
1. Study design and power
2. Descriptive statistics
3. Inferential statistics
4. Software and tips

Study design and power

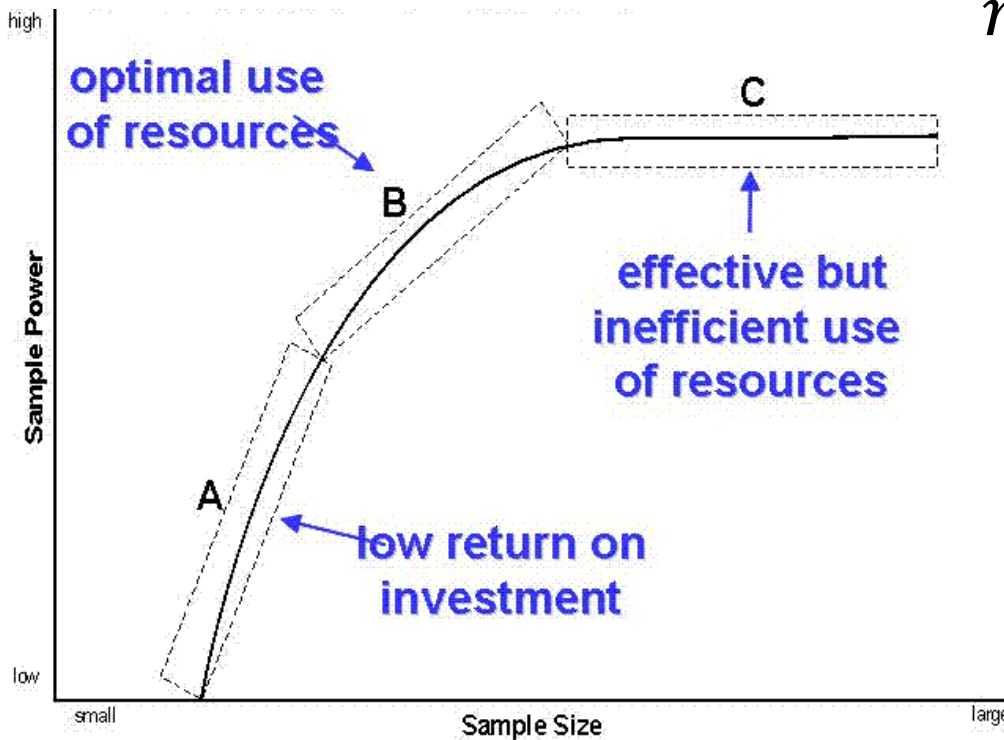
Clinical studies



Types of clinical studies



Power calculation



$$n = \frac{2\sigma^2(Z_\beta + Z_{\alpha/2})^2}{d^2}$$

n = sample size
 σ = standard deviation
 β = power
 α = level of significance
 d = effect size

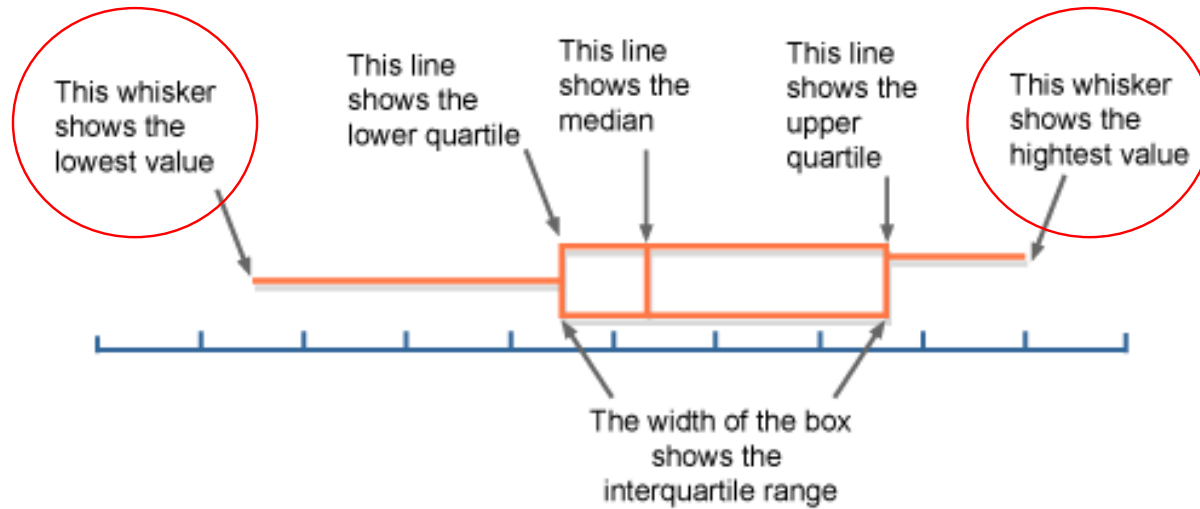
Loss to follow-up

Descriptive statistics

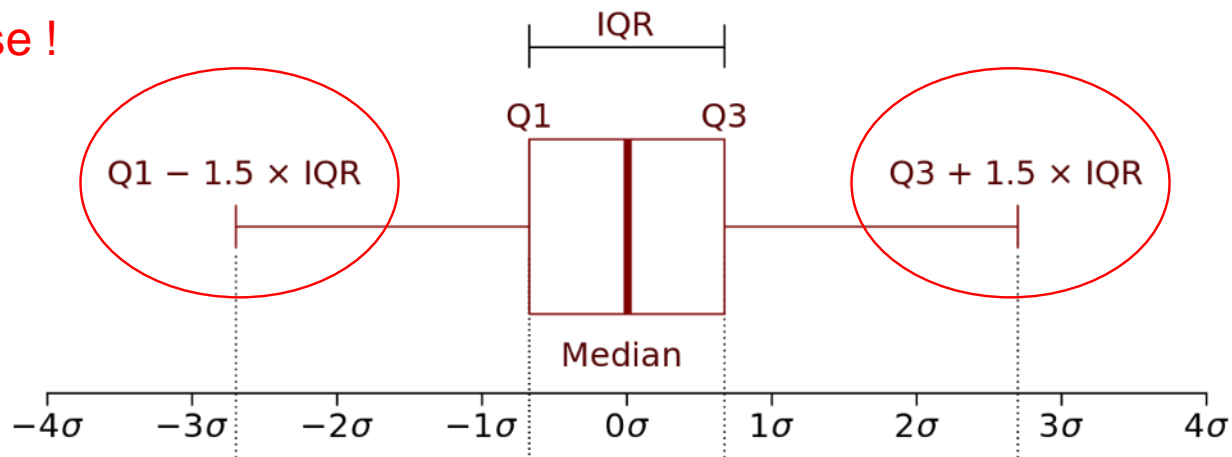
Types of variables

- Nominal - categories (e.g. gender, ethnicity)
- Ordinal - categories with a trend (e.g. cancer stage, grade)
- Numerical / scalar - quantitative
 - Continuous scale (e.g. age, height)
 - Discrete scale (e.g. number of polyp)

Box-and-whisker plot



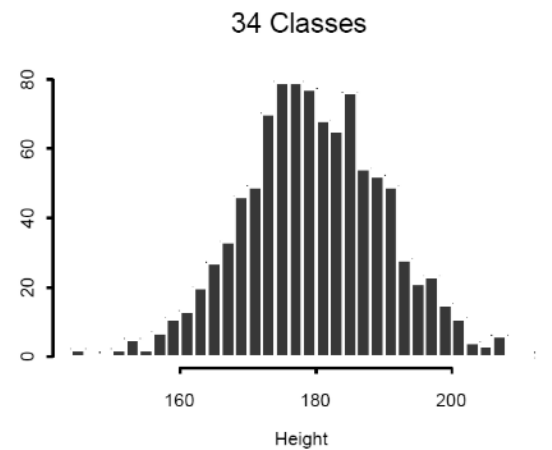
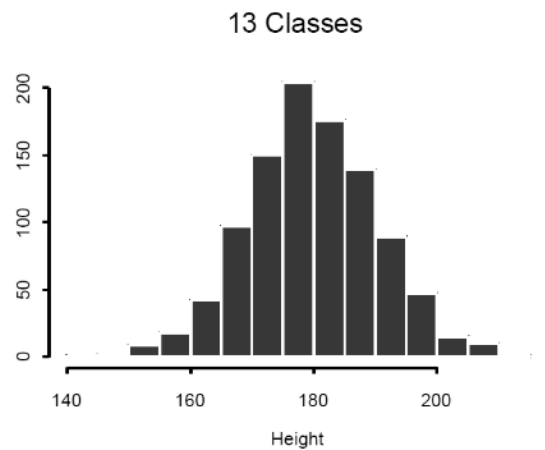
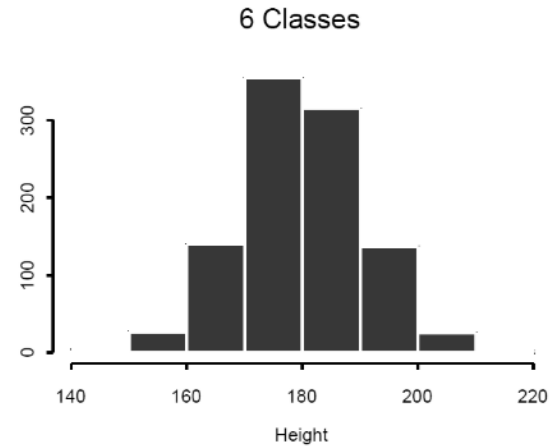
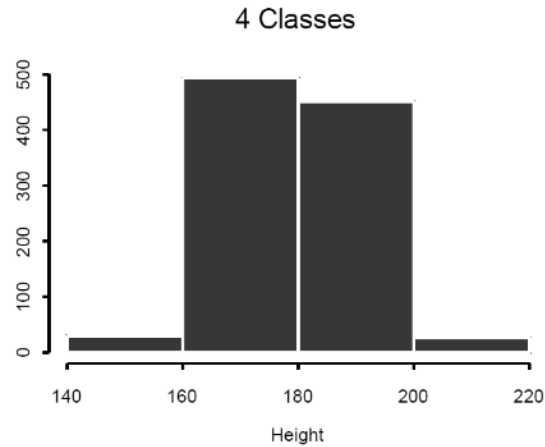
Mind these !



Bar chart

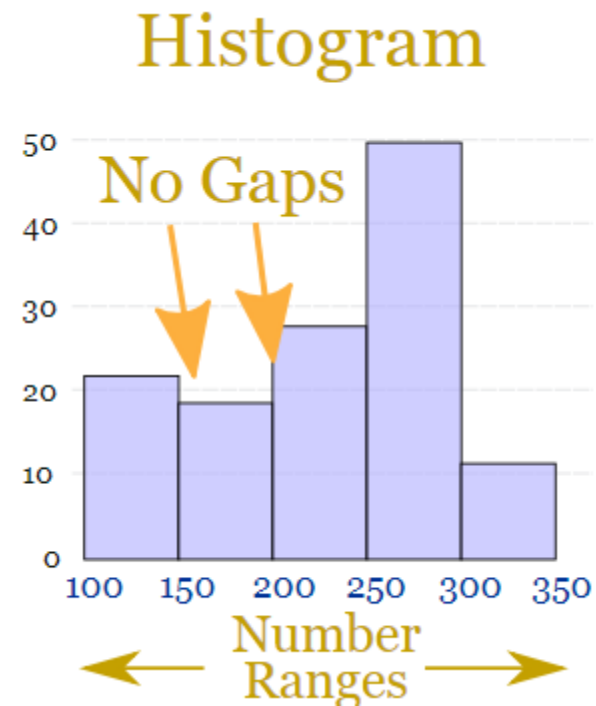
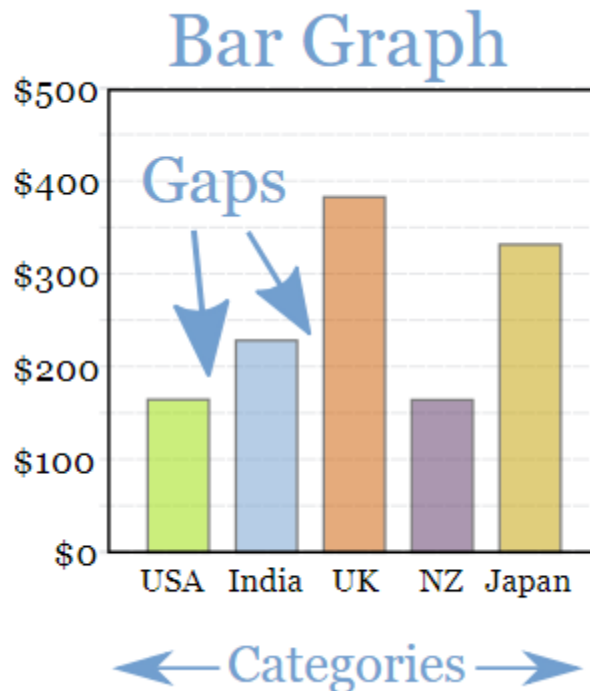
- Height is the mean
- How about error bars?
 - Standard deviation (SD)
 - Standard error of mean (SE / SEM)
 - Confidence interval (95% CI)

Histogram



Q: What are the differences between bar chart and histogram? (P.T.O.)

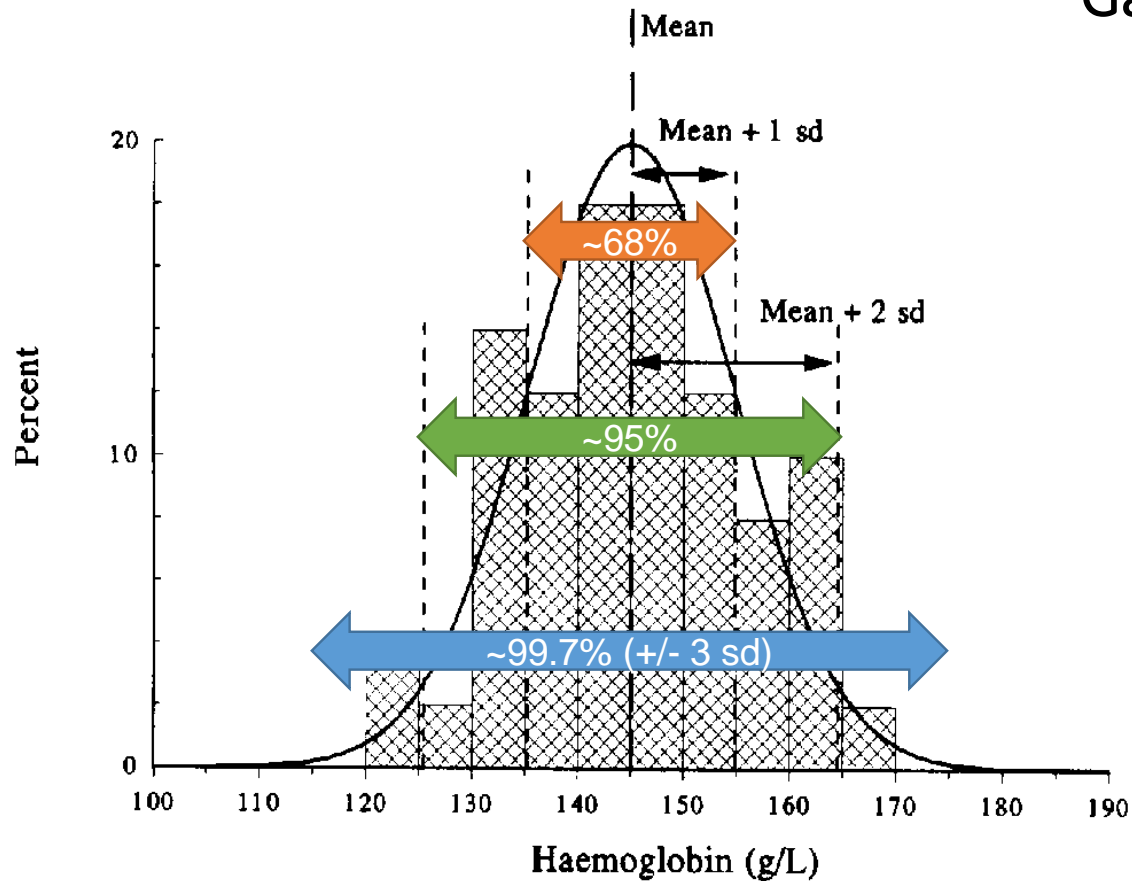
Bar chart and histogram



Bonus Q: what is a Pareto chart?

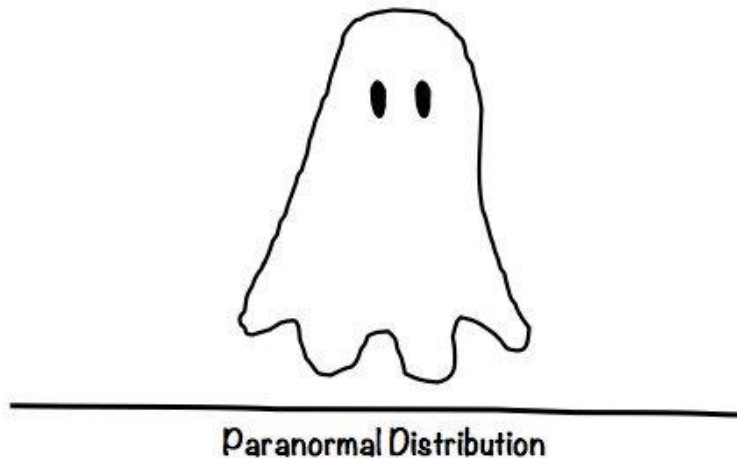
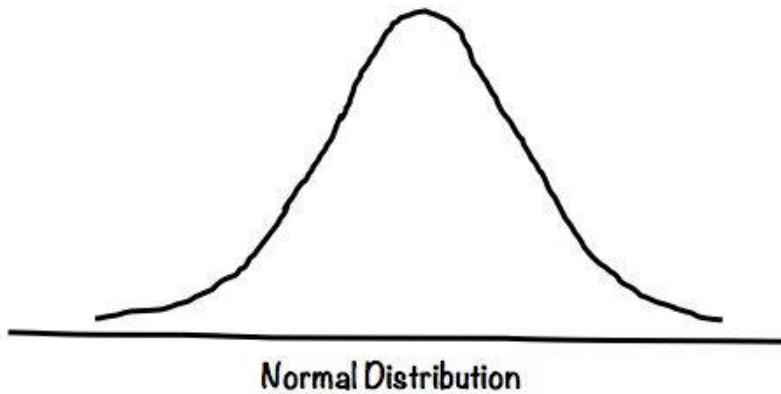
Normal distribution

Gaussian distribution
(z-statistics)



de Moivre,
Gauss & Laplace

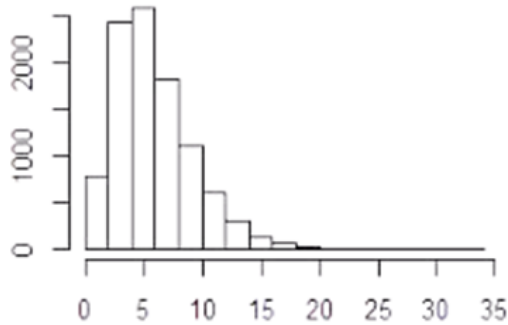
Normality



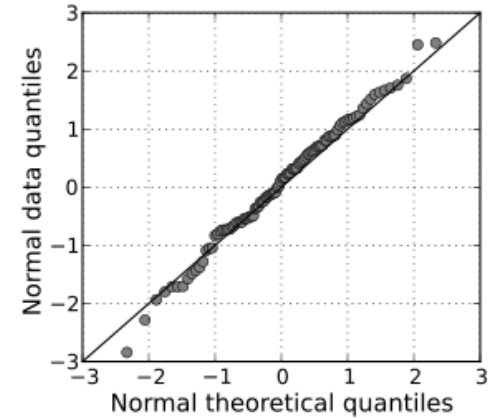
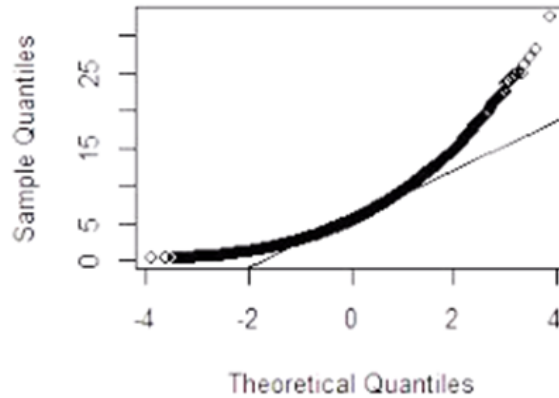
- Histogram inspection
 - $n > 30$
 - Fitting shape
- Quantile-quantile plot
- Formal tests
 - Komolgorov-Smirnov test
 - Shapiro Wilk test

QQ plot

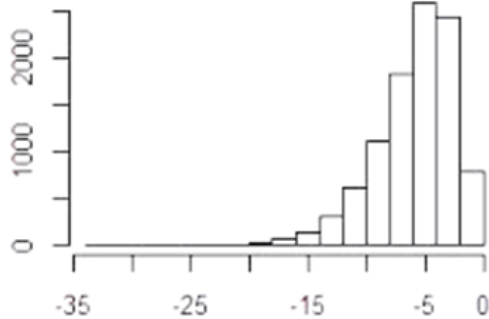
Positive skew



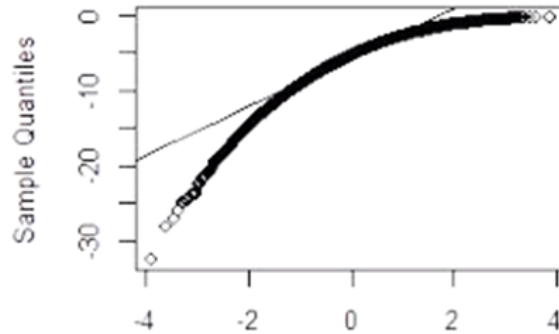
Normal Q-Q Plot



Negative skew



Normal Q-Q Plot

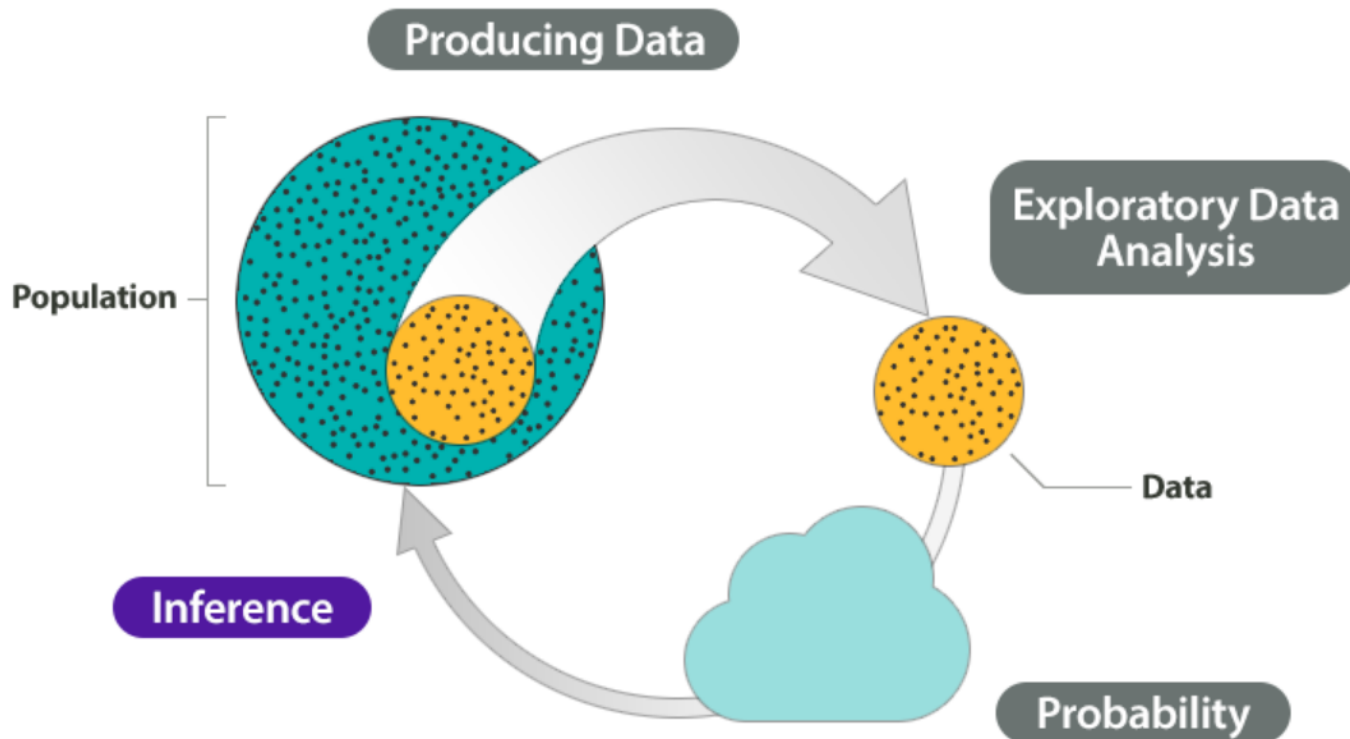


Q: What is the hospital length-of-stay distribution? (right skewed distribution)

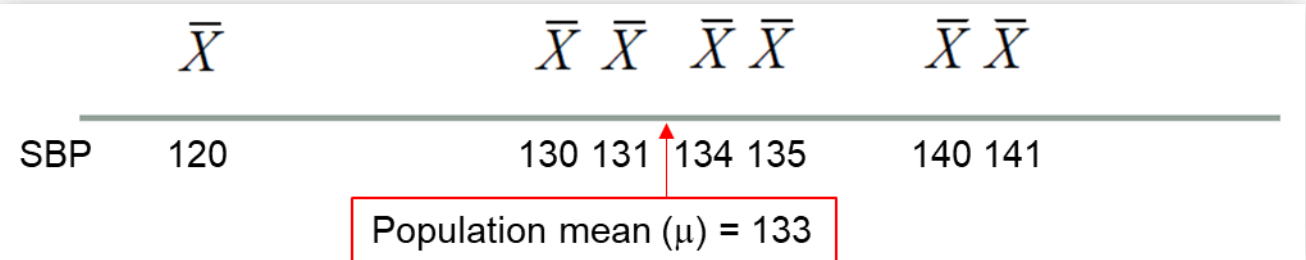
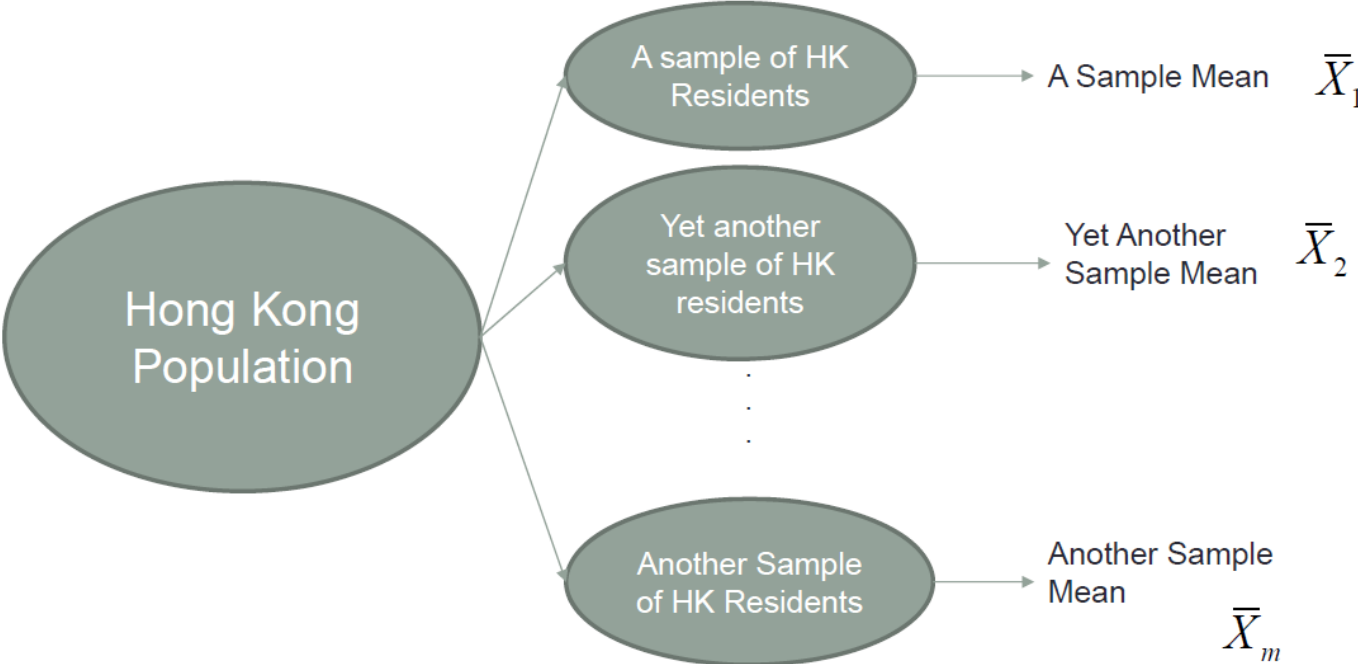
Inferential statistics

Sampling

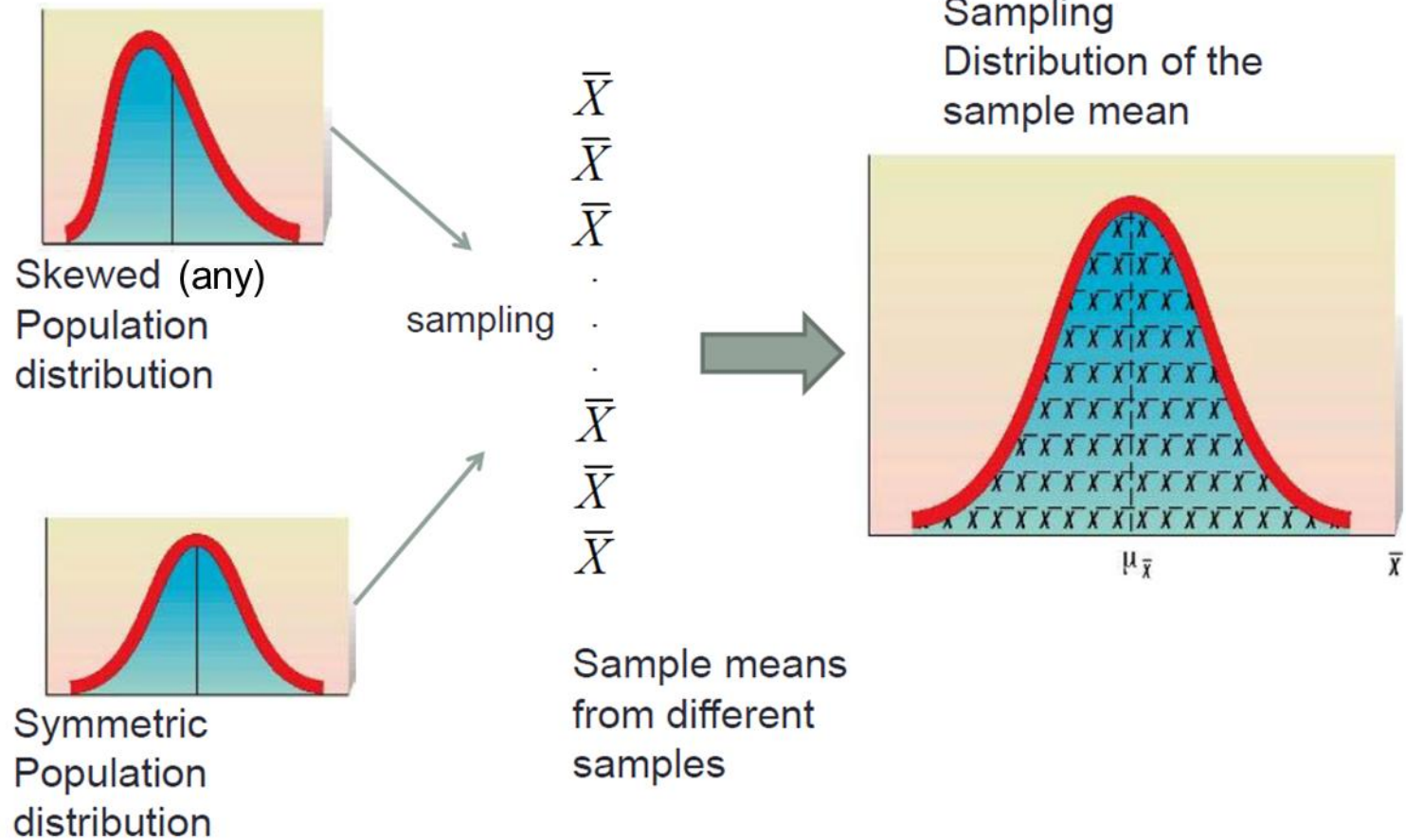
e.g. SBP of the Hong Kong population



Sampling distribution



Central limit theorem

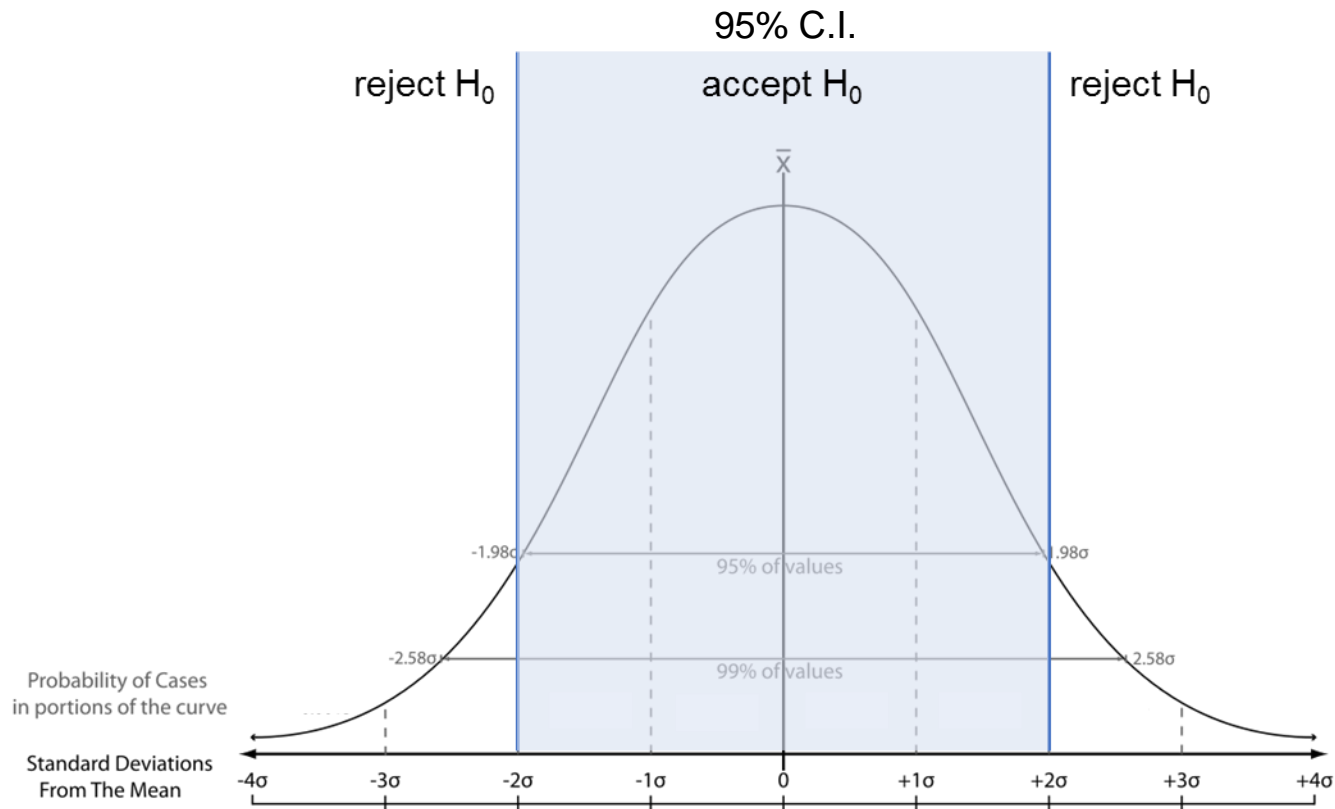


Distribution of the means

$$E(\bar{X}) = \mu$$

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

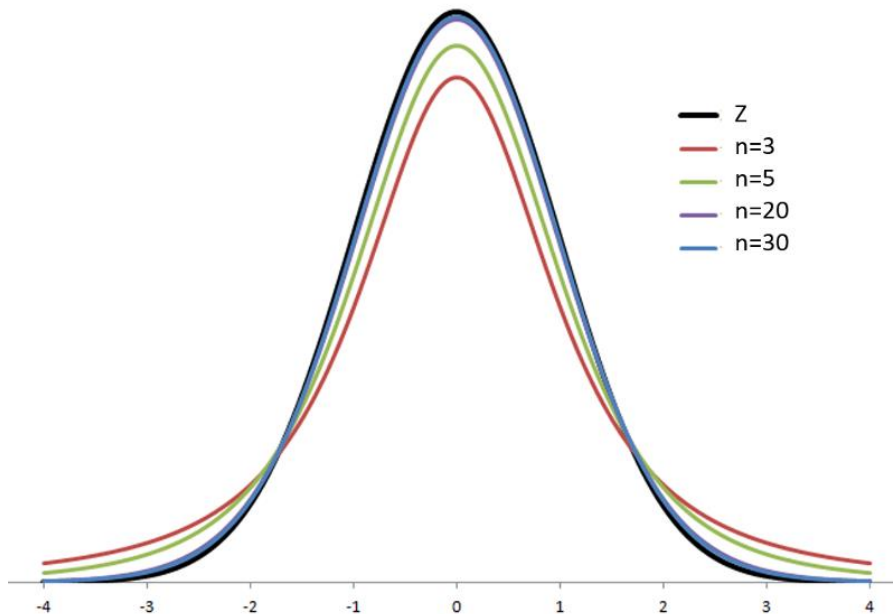
$$se = \frac{\sigma}{\sqrt{n}}$$



Student's t distribution

- Described by William Sealy Gosset

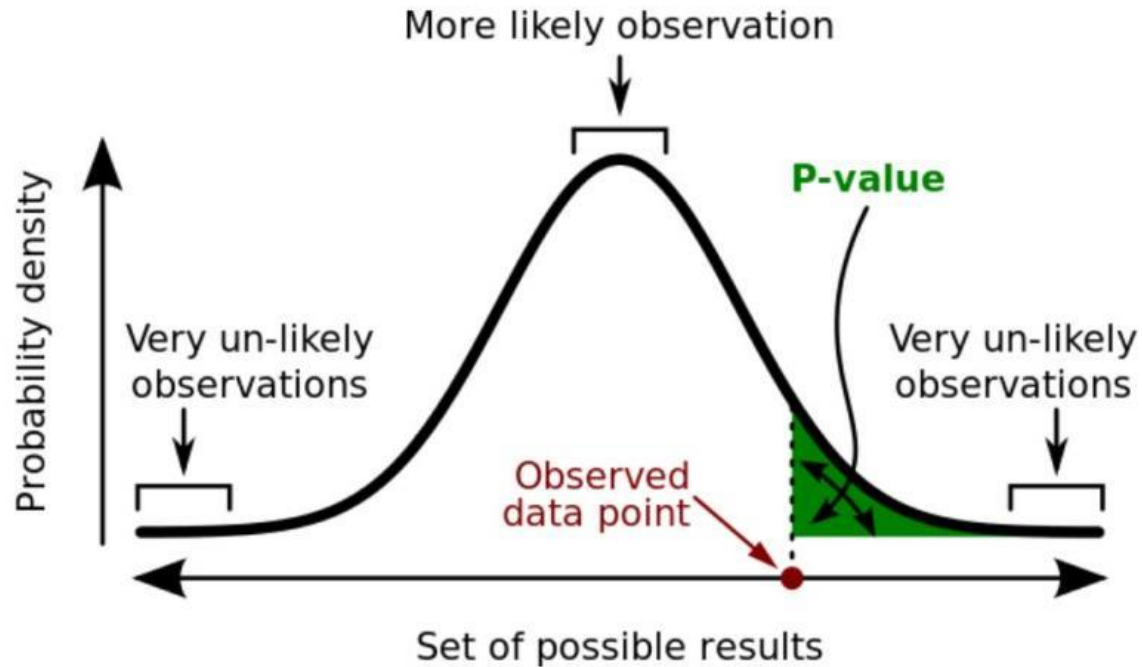
Resemble normal distribution if
sample size is large ($n > 30$)



Hypothesis testing

- Confirmatory data analysis to determine the probability that a given hypothesis is true
- Null hypothesis ' H_0 ': statement of no differences or association between variables
- Alternative hypothesis ' H_1 ': statement of differences or association between variables
- Type I (alpha) and Type II (beta) error

P-value

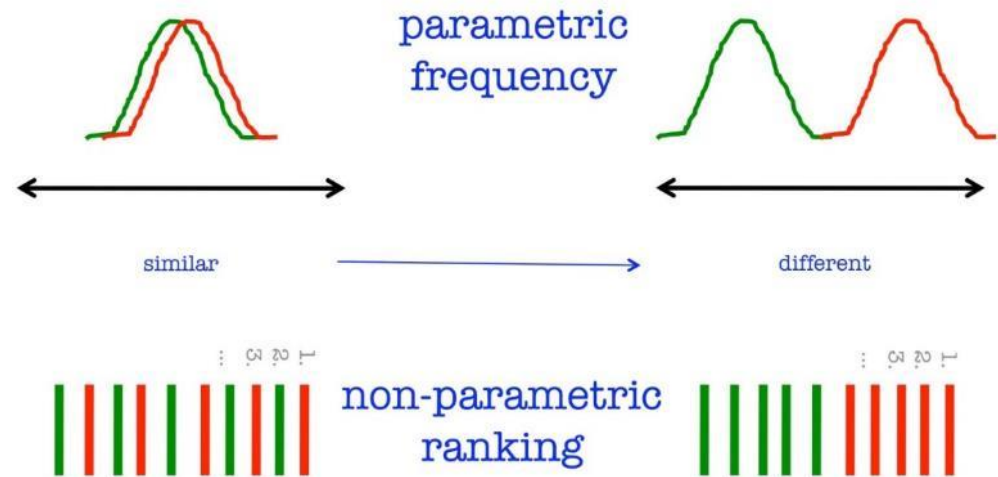


A **p-value** (shaded green area) is the probability of an observed (or more extreme) result assuming that the null hypothesis is true.

Probability of mistakenly rejecting the null hypothesis (α)

Parametric tests

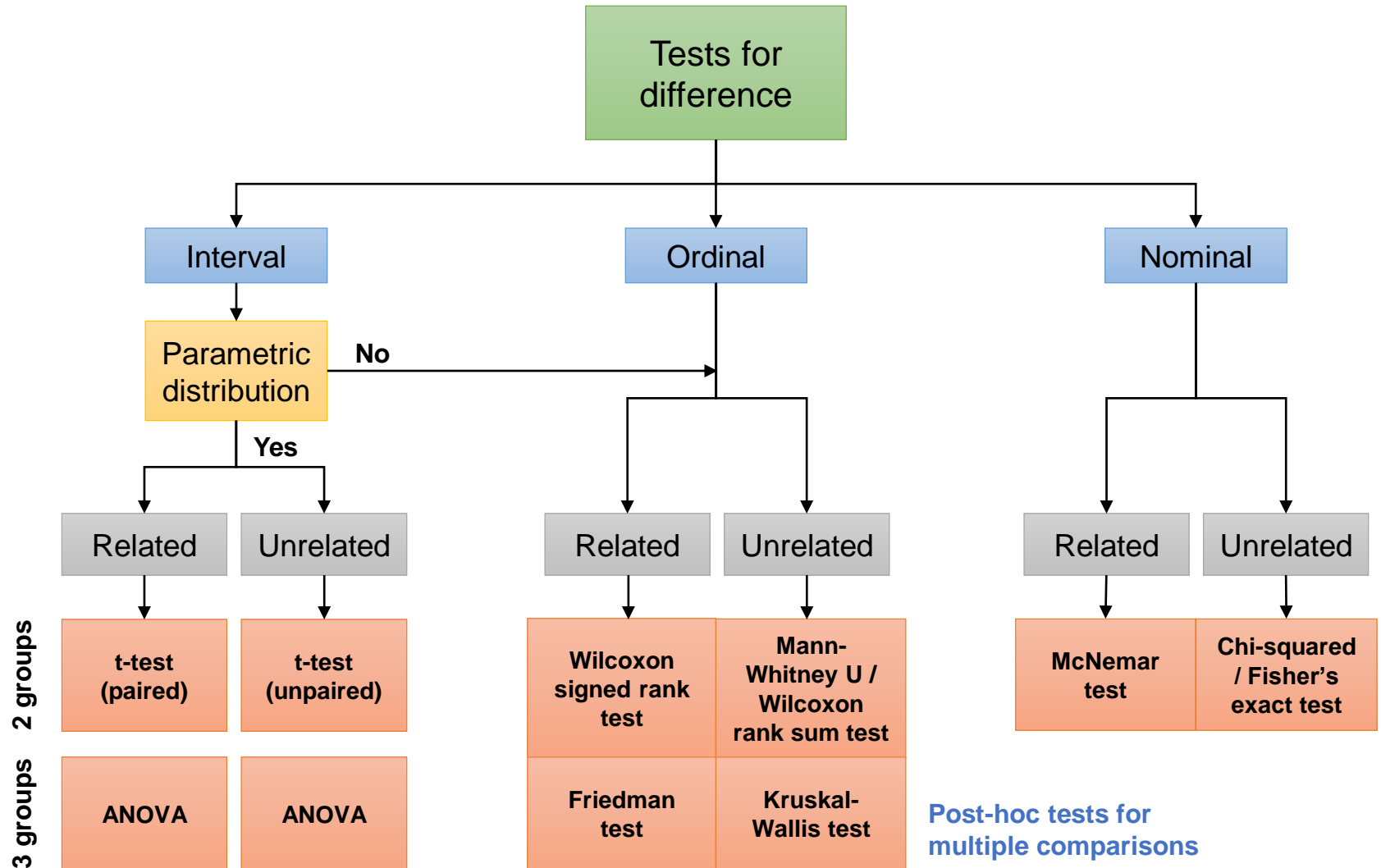
- Parametric tests assume data comes from a population



- With known probability distribution (e.g. normal)
- Based on a fixed set of parameters (e.g. mean, SD)

Non-parametric tests usually less powerful (values discarded)

Comparative tests (interval)



A walk-through of rank sum test

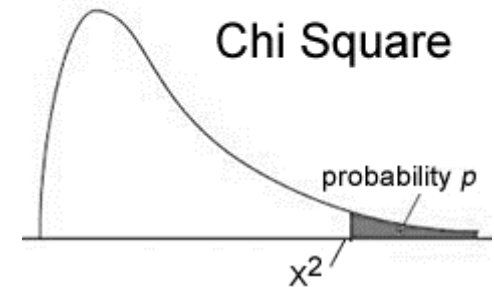
Group A	Group B	Rank A	Rank B		A	B
87	71	19	9	Total rank	127	148
72	42	10	1	Median	74	75.5
94	69	22	8	n	11	12
49	97	2	23			
56	78	4	14.5	U(A)	71	
88	84	20	17	U(B)	61	
74	57	12	5	U statistic	61	
61	64	6	7	P-value	0.76	
80	78	16	14.5			
52	73	3	11			
75	85	13	18			
	91		21			

Comparative tests (categorical)

- Examine differences between *observed* and *expected* counts

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

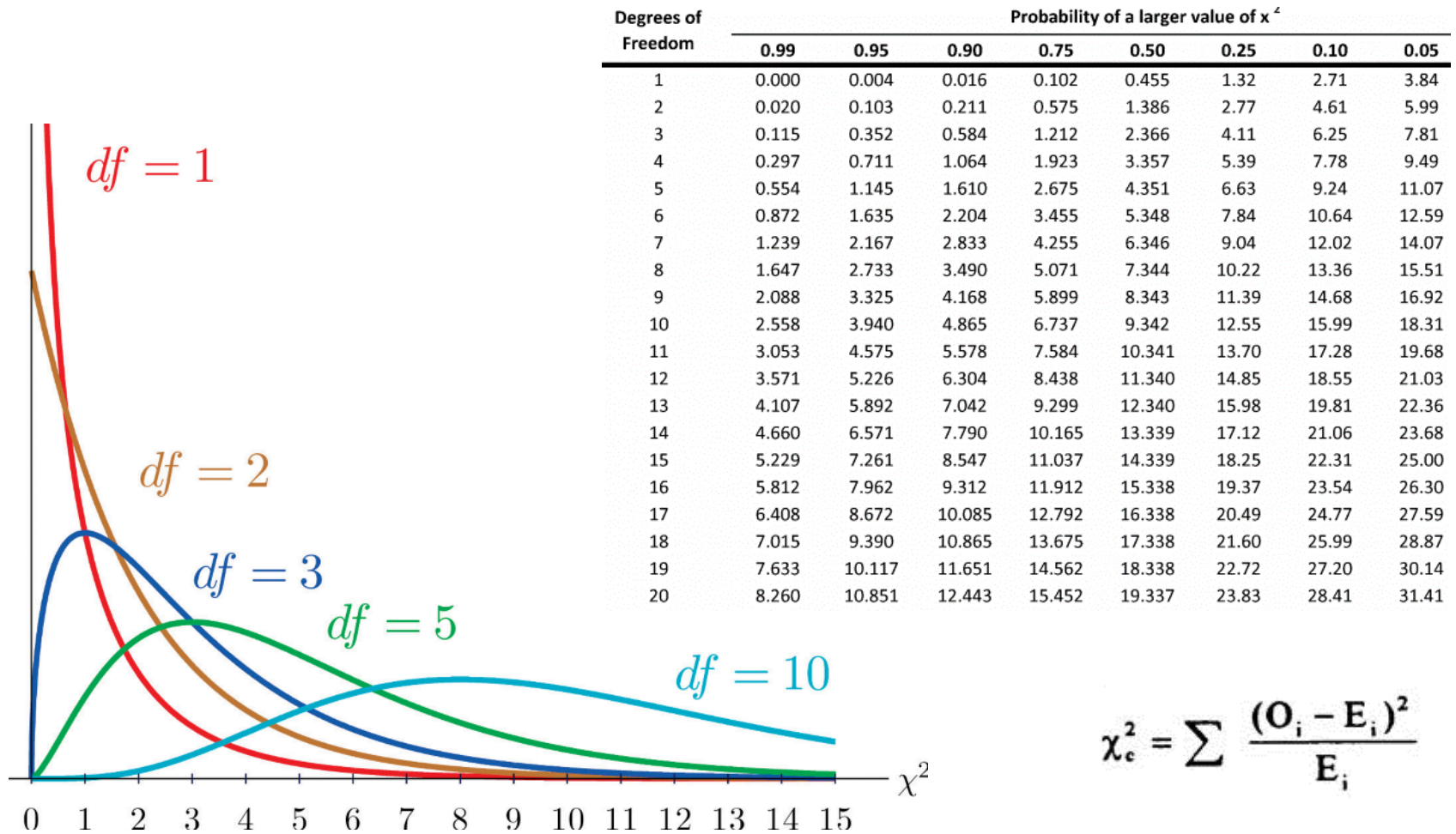
- Two assumptions
 - Independence of observations
 - Count of all cells >5



- Degree of freedom
 - $df = (R-1) \times (C-1)$
 - No of free variables

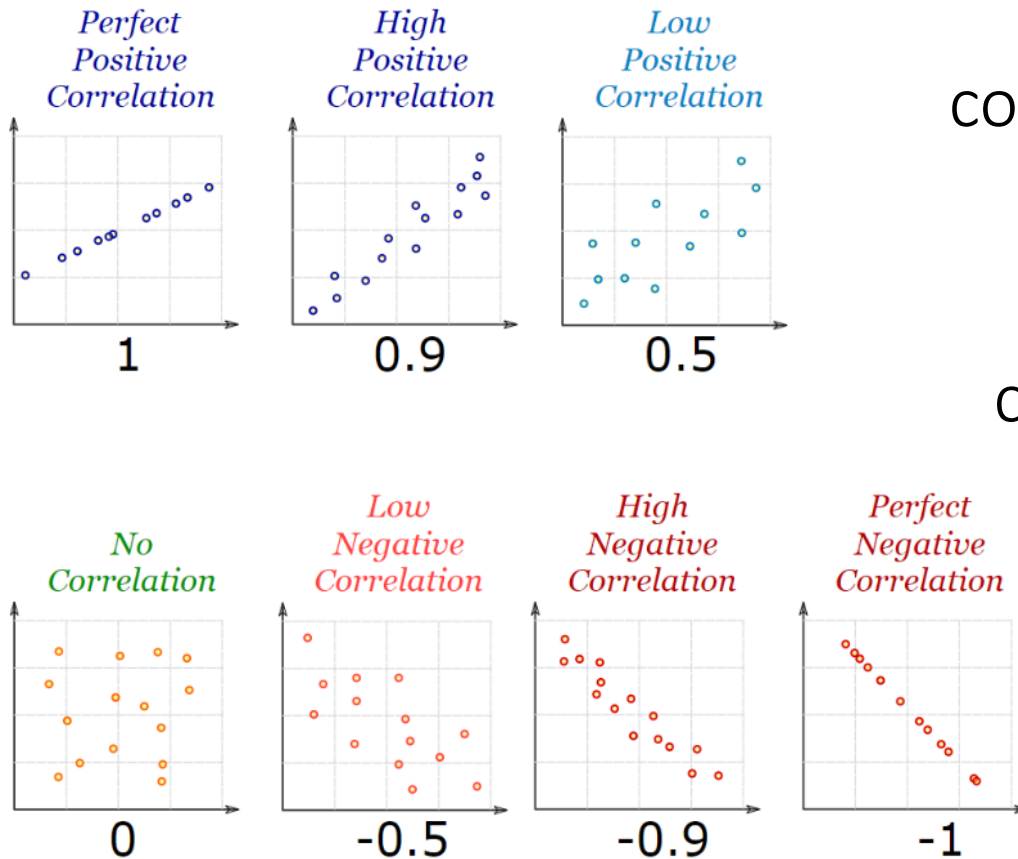
	Drug A	Drug B	
Cured	20	10	30
Not cured	12	22	34
	32	32	64

χ^2 distribution and df



$$\chi_c^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Correlation



Pearson's correlation coefficient (parametric)

Spearman's rank correlation coefficient (non-parametric)

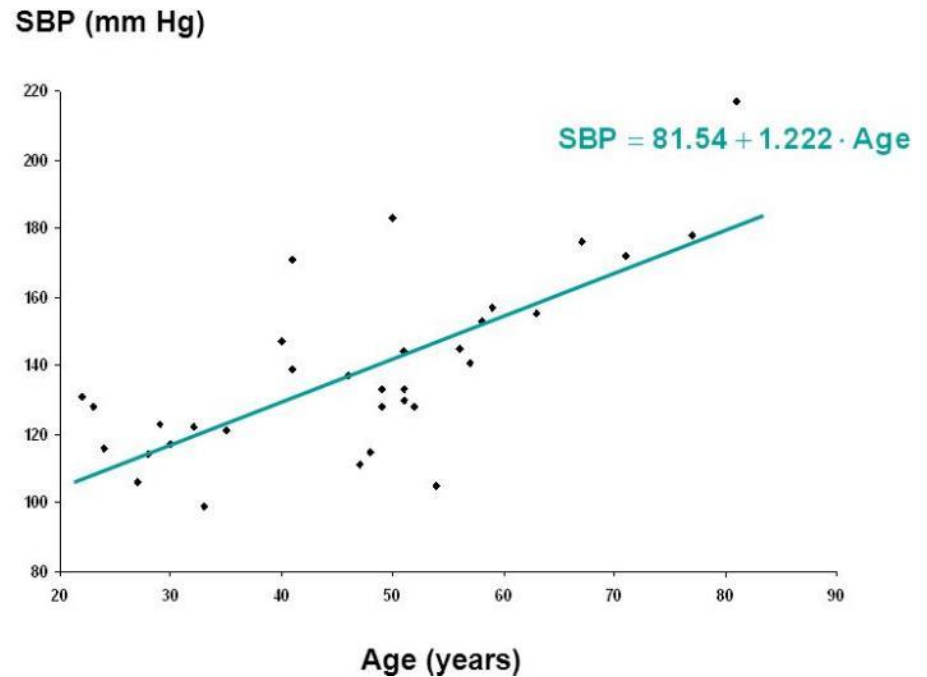
Regression

- Estimating relationships between variables (between dependent and independent variables)

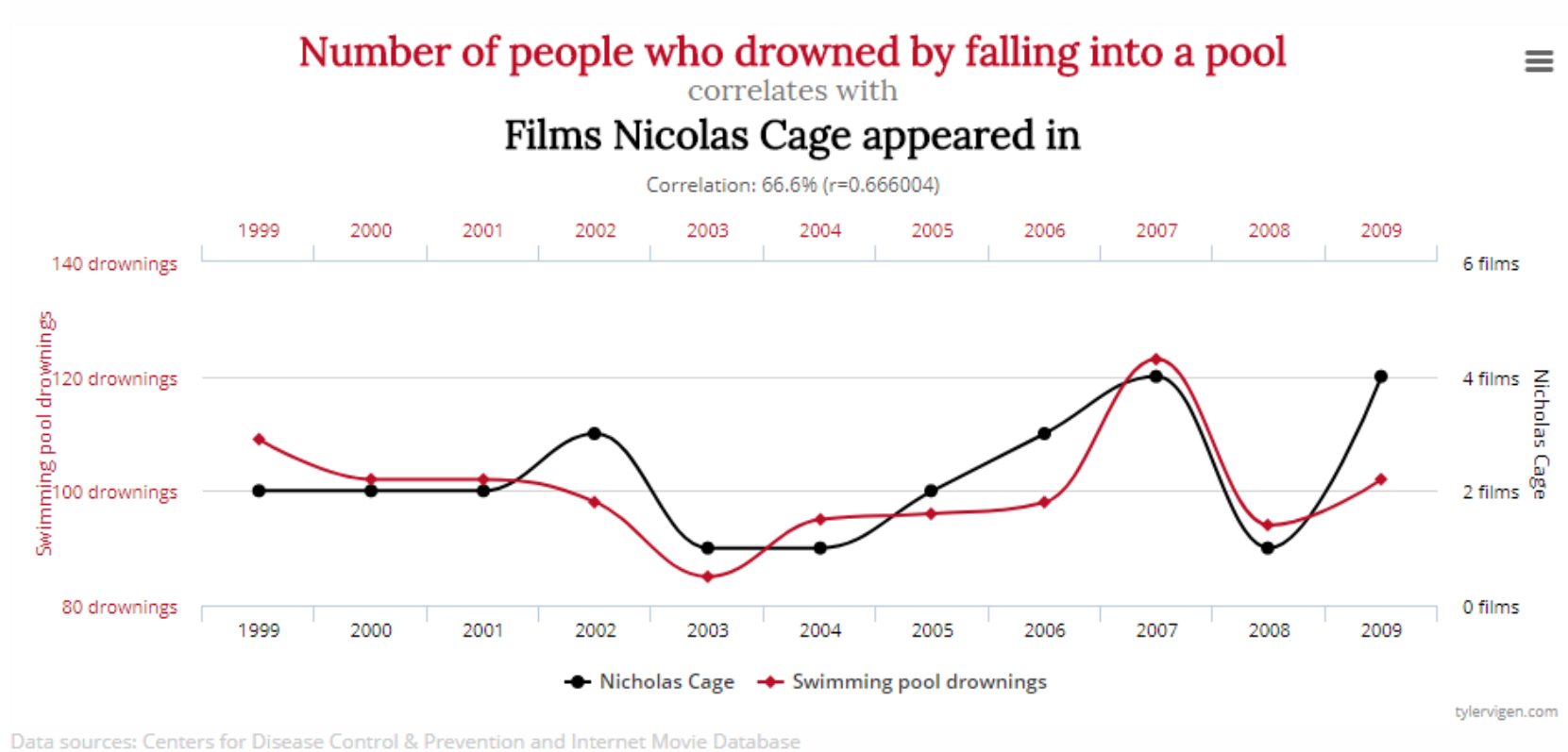
- Logistic regression

- Linear regression

$$\hat{Y} = bX + a$$



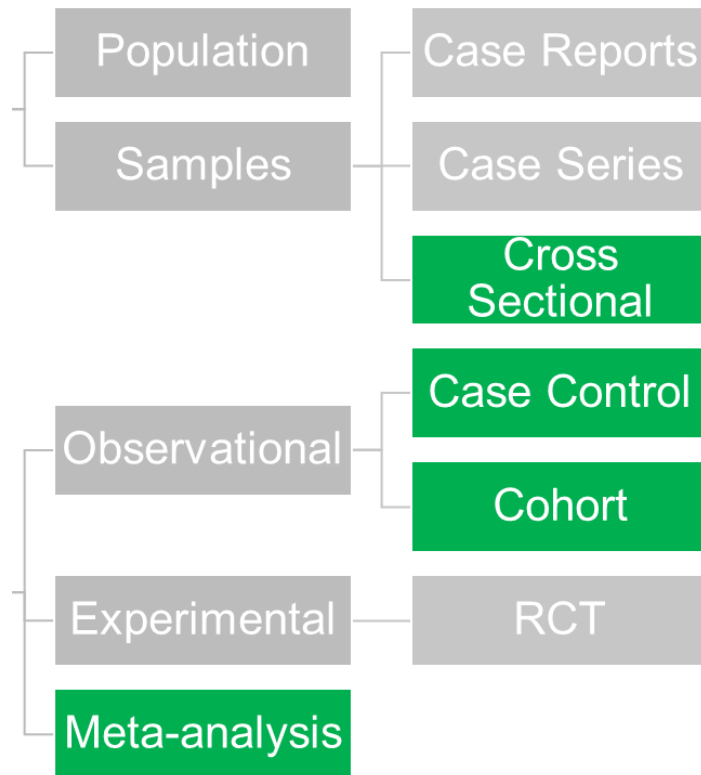
Correlation/association \neq causality



Q: What are the differences between correlation and regression?

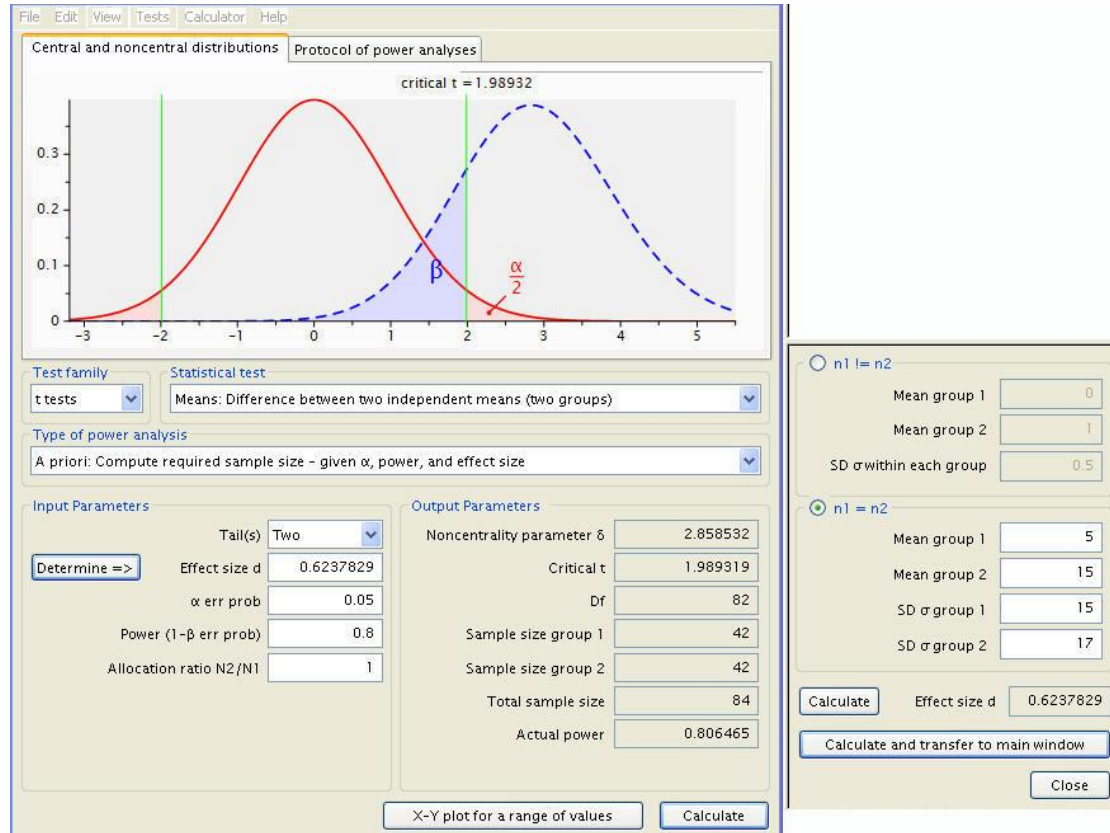
Tips and software

Study design



- Most important is your **research question**
- Consider
 - Time
 - Effort
 - Infrastructure
 - Clinical ethics, governance and compliance

Power calculators



GPower: <http://www.gpower.hhu.de/>

CUHK CCRB: <http://www2.ccrb.cuhk.edu.hk/web/>

Clin Calc: <http://clincalc.com/stats/samplesize.aspx>

Statistical software



R Project for Statistical Computing

```
summaryStats.R x HollywoodMovies2011 x
Source on Save Run Source
1 # loading the data set
2 HollywoodMovies2011 <- read.csv("HollywoodMovies2011.csv")
3
4 # Calculating summary stats for all columns of a data set
5 summary(HollywoodMovies2011)
6
7 # Calculating summary stats for an individual column
8 summary(HollywoodMovies2011$RottenTomatoes)
9
10 # Calculating specific summary statistics
11
12 # mean
13 mean(HollywoodMovies2011$RottenTomatoes)
14 mean(HollywoodMovies2011$RottenTomatoes, na.rm = TRUE)
15
16 # median
17 median(HollywoodMovies2011$RottenTomatoes, na.rm = TRUE)
18
19 # standard deviation
20 (Top Level) >
```

Console ~/Documents/m107/data/

Genre	TheatersOpenWeek	BOAverageOpenWeek	DomesticGross	ForeignGross	WorldGross
Action :32	Min. : 3	Min. : 1513	Min. : 0.02	Min. : 0.24	Min. : 0.025
Comedy :27	1st Qu.:2550	1st Qu.: 3779	1st Qu.: 19.03	1st Qu.: 14.25	1st Qu.: 30.706
Drama :21	Median :2995	Median : 5686	Median : 37.35	Median : 47.00	Median : 76.659
Horror :17	Mean :2828	Mean : 8339	Mean : 63.22	Mean : 96.92	Mean : 150.742
Thriller :13	3rd Qu.:3400	3rd Qu.: 8923	3rd Qu.: 80.46	3rd Qu.:102.00	3rd Qu.: 173.691
Animation:12	Max. :4375	Max. :93230	Max. :381.01	Max. :947.10	Max. :1328.111
(Other) :14	NA's :16	NA's :16	NA's :2	NA's :15	NA's :2

Budget	Profitability	OpeningWeekend
Min. : 0.20	Min. : 0.000	Min. : 0.00
1st Qu.: 20.25	1st Qu.: 1.065	1st Qu.: 7.71
Median : 36.50	Median : 2.199	Median : 13.10
Mean : 53.48	Mean : 3.315	Mean : 20.34
3rd Qu.: 70.00	3rd Qu.: 3.667	3rd Qu.: 25.00
Max. :250.00	Max. :64.673	Max. :169.19
NA's :2	NA's :2	NA's :3

>

Statistical software



SPSS



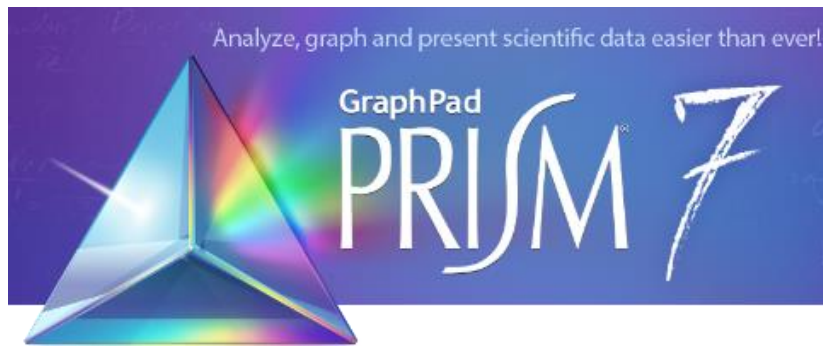
SAS



Stata



Minitab



Questions?