

# Essential statistics

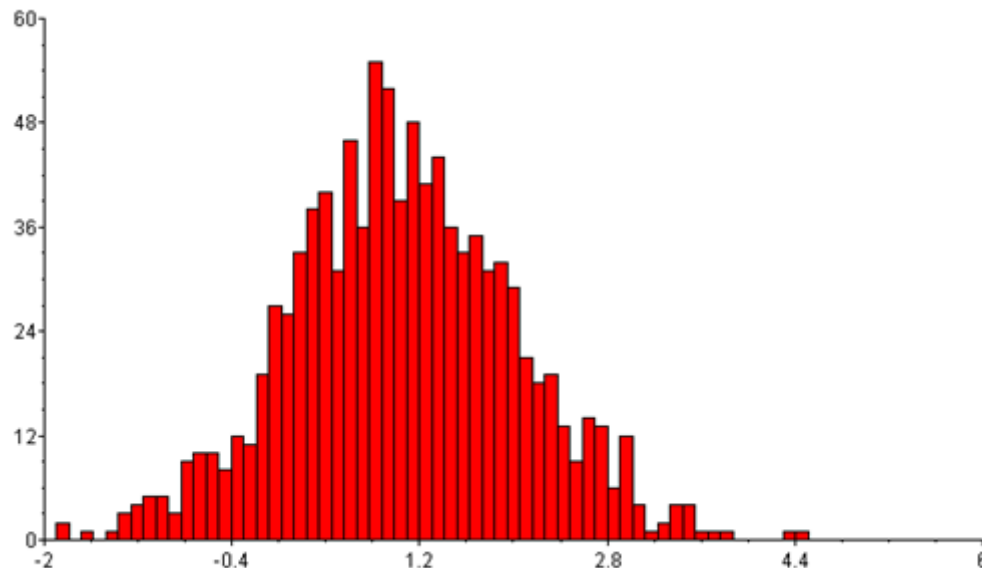
Jason Smith MD DMI FRCS(Gen.Surg)

Paris Tekkis MD FRCS(Gen.Surg)

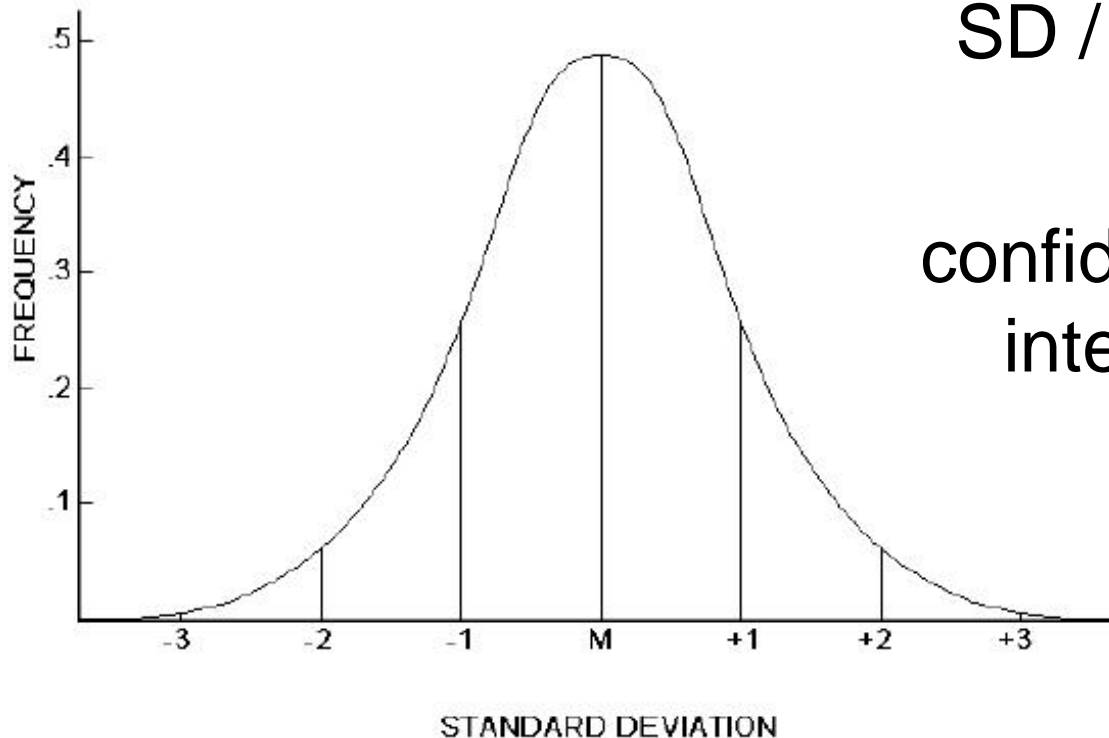


# Types of data (variables)

- Categorical - male/female
- Ordered – mild, moderate, severe
- Numerical - age, BP
- Dependent, independent



# Are the data normally distributed? parametric analysis



Mean

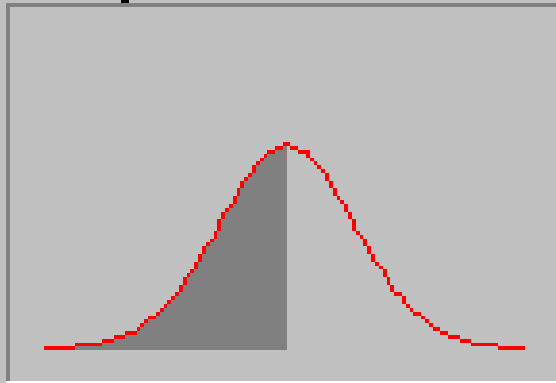
SD / SEM

95%  
confidence  
intervals

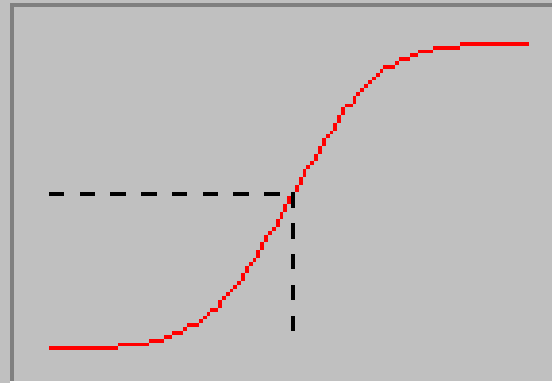


# normally distributed data

Density Function:



Distribution Function:



$$z = 0.00$$

$$p = .50$$

- Tests of normality (Kolmogorov-Smirnov test, or the Shapiro-Wilks' W test)
- Histogram (best)



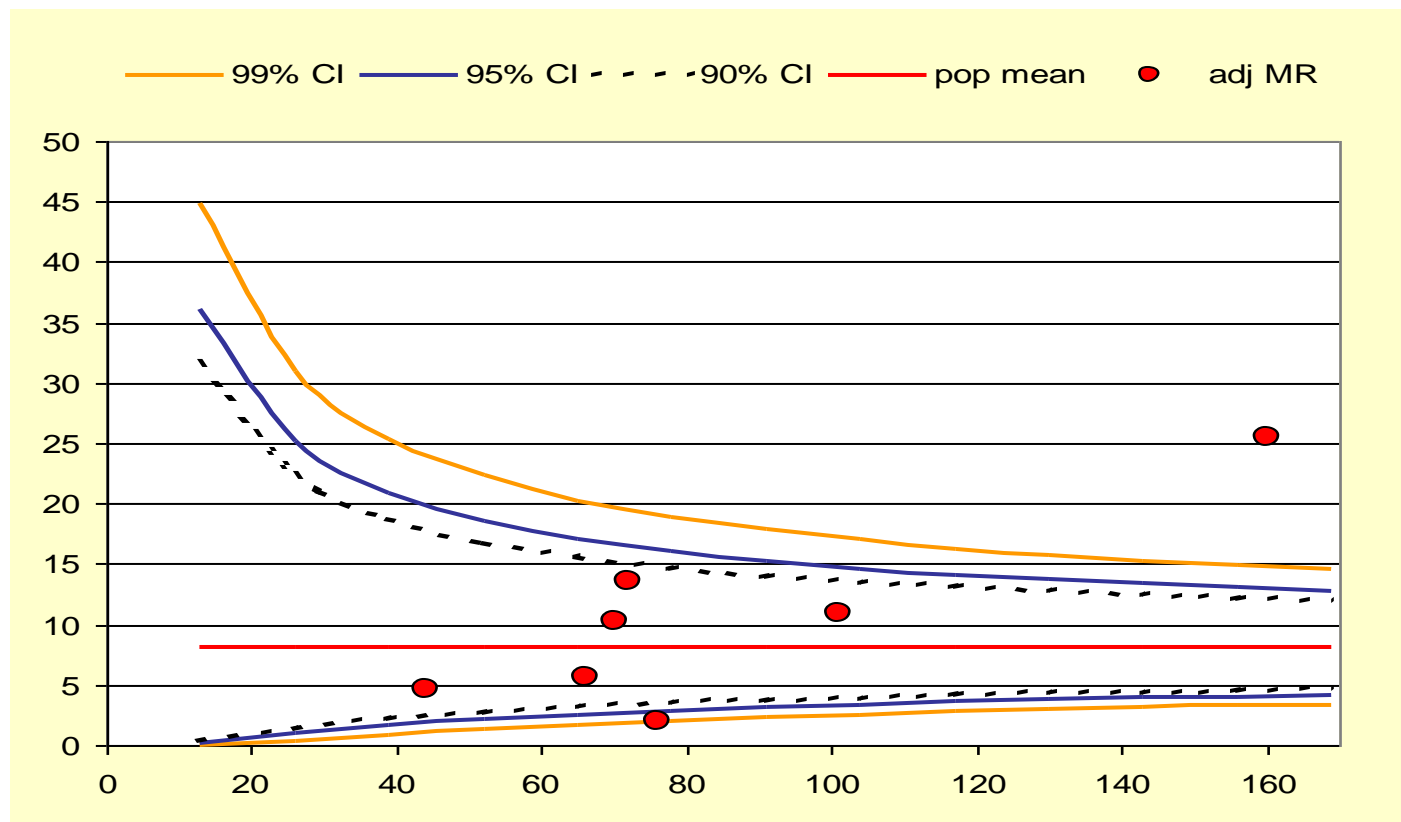
# standard deviation

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

- Is the square root of the variance!
- Simply put it is a description of how far an individual observation is away from the mean



# confidence intervals



- a range of values where the "true" (population) value of a data point can be expected to be located
- You can specify certainty

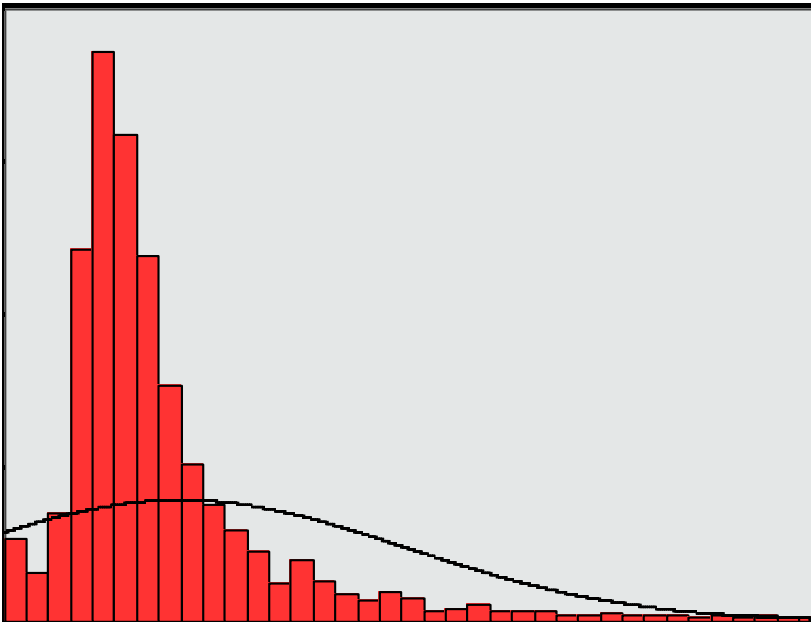


# Other distributions

- **Binomial** – only 2 possibilities (eg blood grp B and not grp B). As sample size increases, approximates normal
- **Poisson** – number of occurrences of an event, eg daily number of notifications of new cancer to a registry. Again approximates to normal as size increases
- **Others include** (Bernoulli Distribution Beta Distribution Cauchy Distribution Chi-square Distribution Exponential Distribution Extreme Value Distribution F Distribution Gamma Distribution Geometric Distribution Gompertz Distribution Laplace Distribution Logistic Distribution Log-normal Distribution Pareto Distribution Rayleigh Distribution Rectangular Distribution Student's t Distribution Weibull Distribution)



# skewed distribution? non-parametric analysis

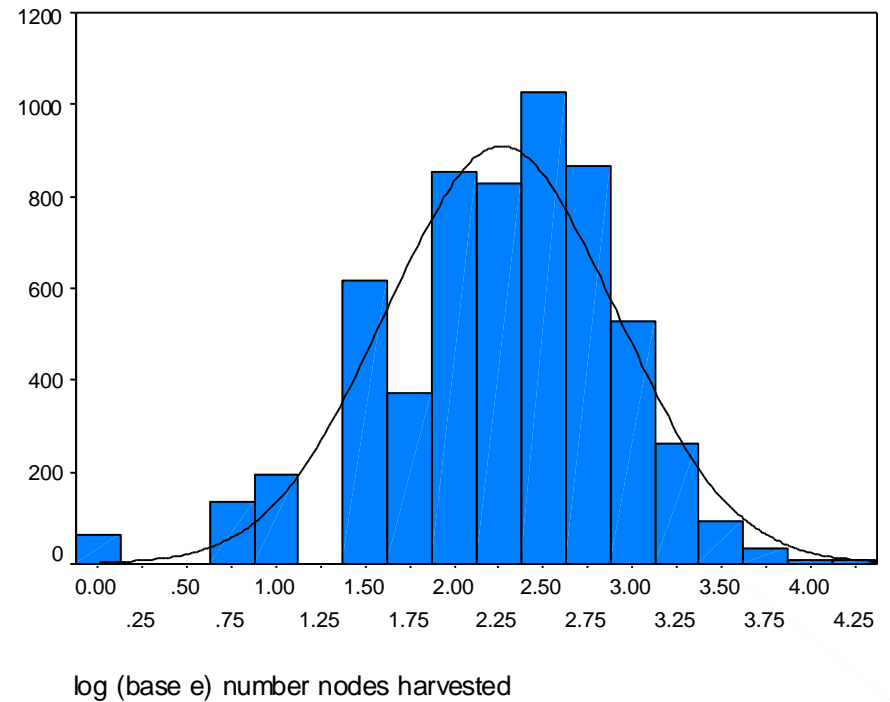
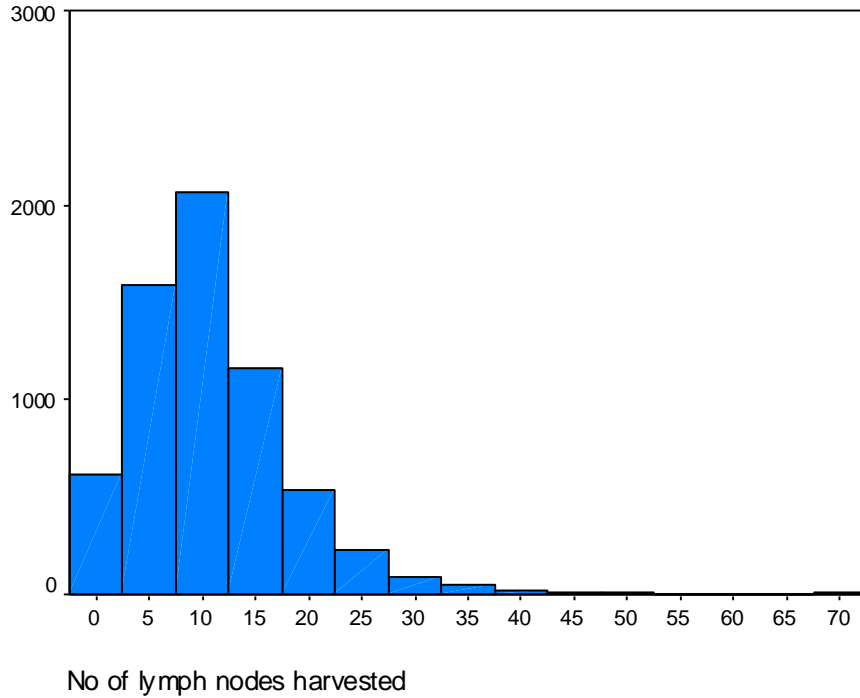


- Mode
- Median
- Geometric mean?

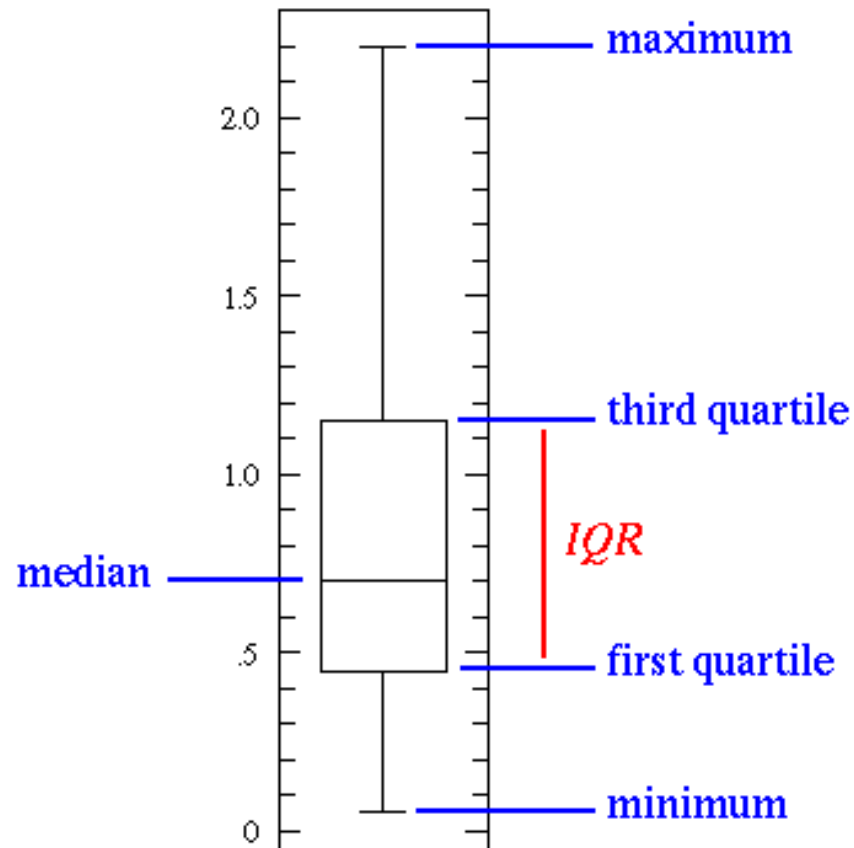
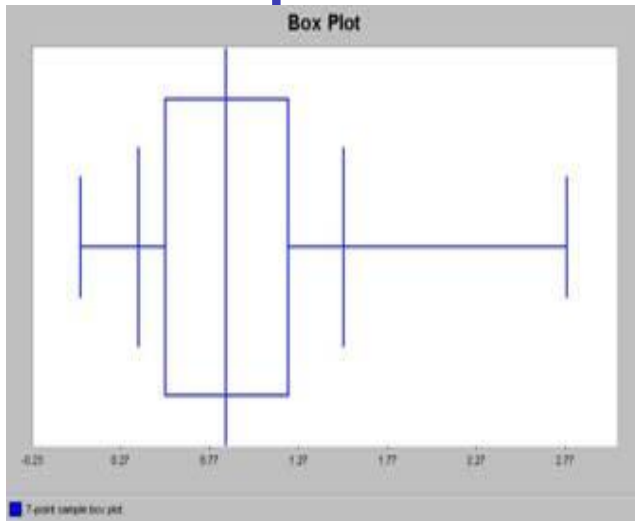




# geometric mean



# non-parametric analysis: box-plots



- Median
- Inter-quartile range
- Range



**Operation A has a lower  
30-day mortality than  
operation B**

**$P < 0.05$**



# what is significance?

*the probability that an observed outcome of an experiment or trial is due to chance alone*

- What are P values
- Setting acceptable levels (2SD)
- What does  $P < 0.05$  mean?
- Arbitrary assignment



# sample size

- Few observations = few combinations = high chance of occurring by chance
- Magnitude of difference
- Power analysis - how large a sample is needed to enable statistical judgments that are accurate and reliable



# hypothesis testing

- Research hypothesis: Do diabetic patients have a high blood pressure?
- Null hypothesis: the mean blood pressure of diabetic and non-diabetic patients is the same
- Alternative hypothesis?



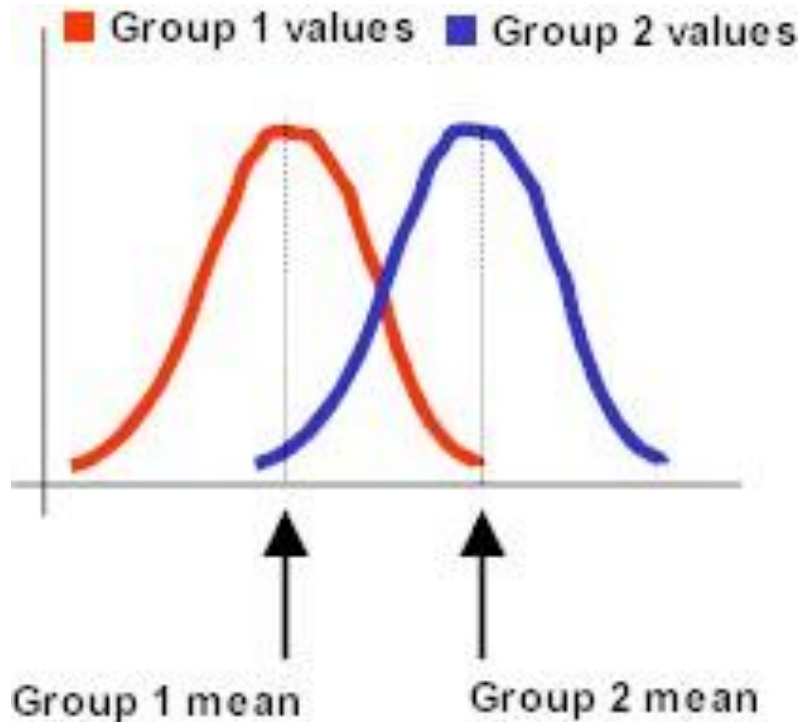
# Comparing two samples

- Compare the average and spread
- Assume that both samples are taken from the same population (null hypothesis)
- Tests the likelihood that they are from the same population
- Arbitrary likelihood ( $p$ ) less than 5% ( $p < 0.05$ ) taken as statistically significant ( $> 2SD$ )



# Parametric two-sample tests

Null-hypothesis



Related samples  
(before and after)

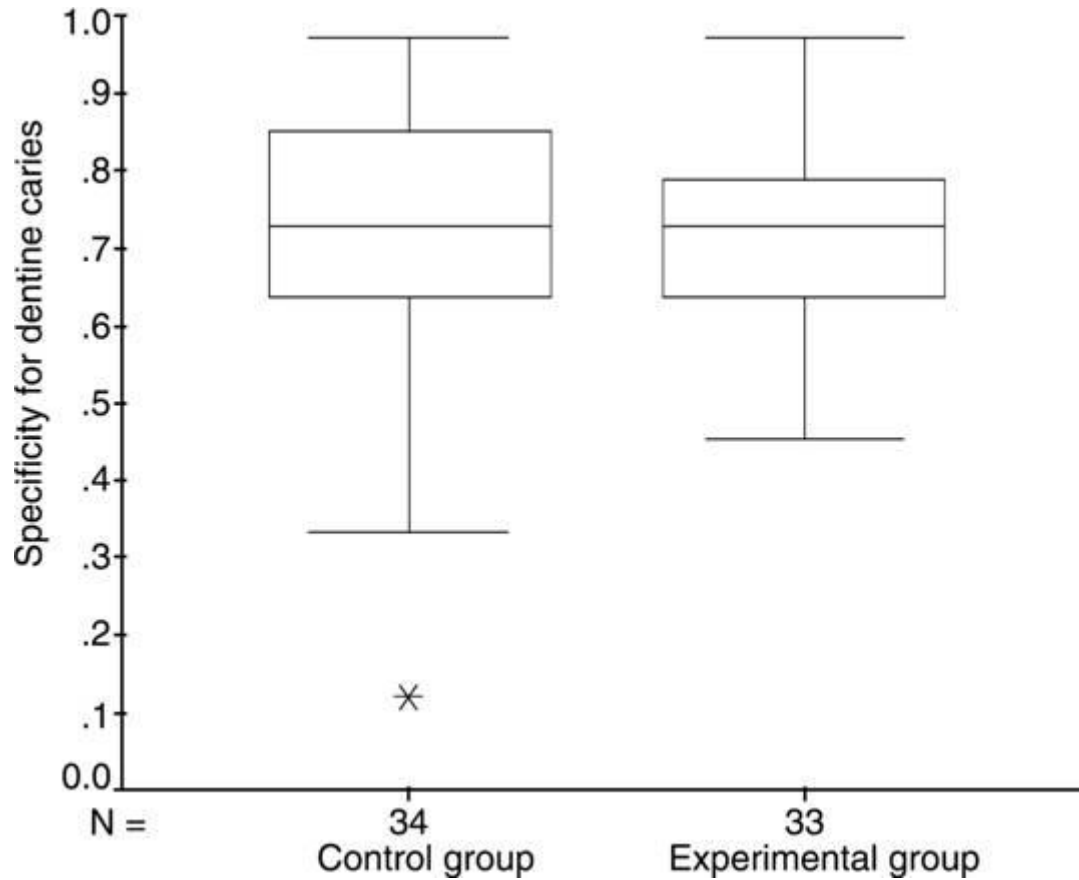
Paired t-test

non-related  
samples (male vs  
female) two-  
sample t-test





# Non-parametric – two-box plots



# Two sample comparisons

	Parametric	Non-parametric
Two related samples	Paired t-test	Wilcoxon rank sum test
Two independent	Two-sample t-test	Mann-Whitney U test



# Statistical Errors

- Type II
  - A difference is present between the samples but your methods failed to show it (small sample size)
- Type I
  - There are no differences between the samples but the your methods was such that you showed a difference (repeated testing)



# Comparative tests – categorical variables

- 2 x2

Is the data large enough\*?

Yes- Chi-squared test

No – Fisher exact test

	Dead	Alive
Rx A	25	5
Rx B	3	23

\* USE Yates correction for total <100 or any cell <10



# Comparative tests – categorical variables

	Dead	Alive
ASA I	2	98
ASA II	5	55
ASA III	10	30

- 3 x2 ... 4x2 ... 5x2

Is there natural ordering

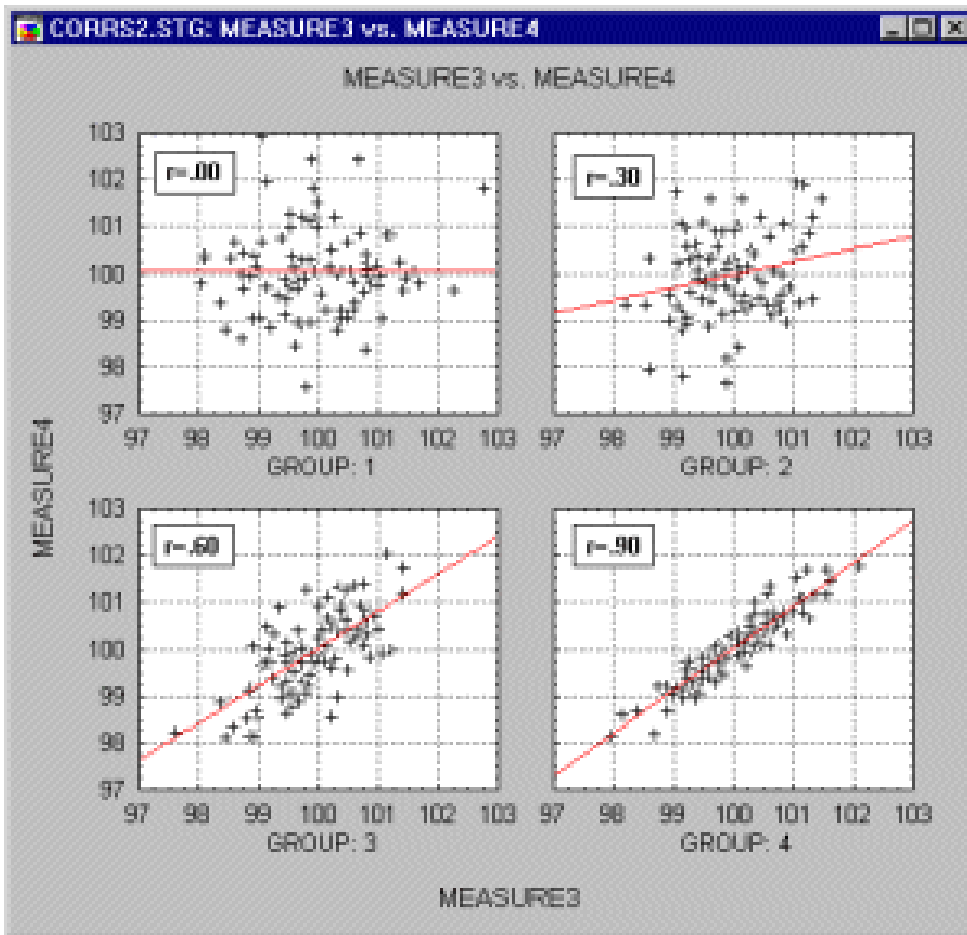
Yes – Chi-squared test for trend

No –  $r \times c$  Chi-squared test



# correlation

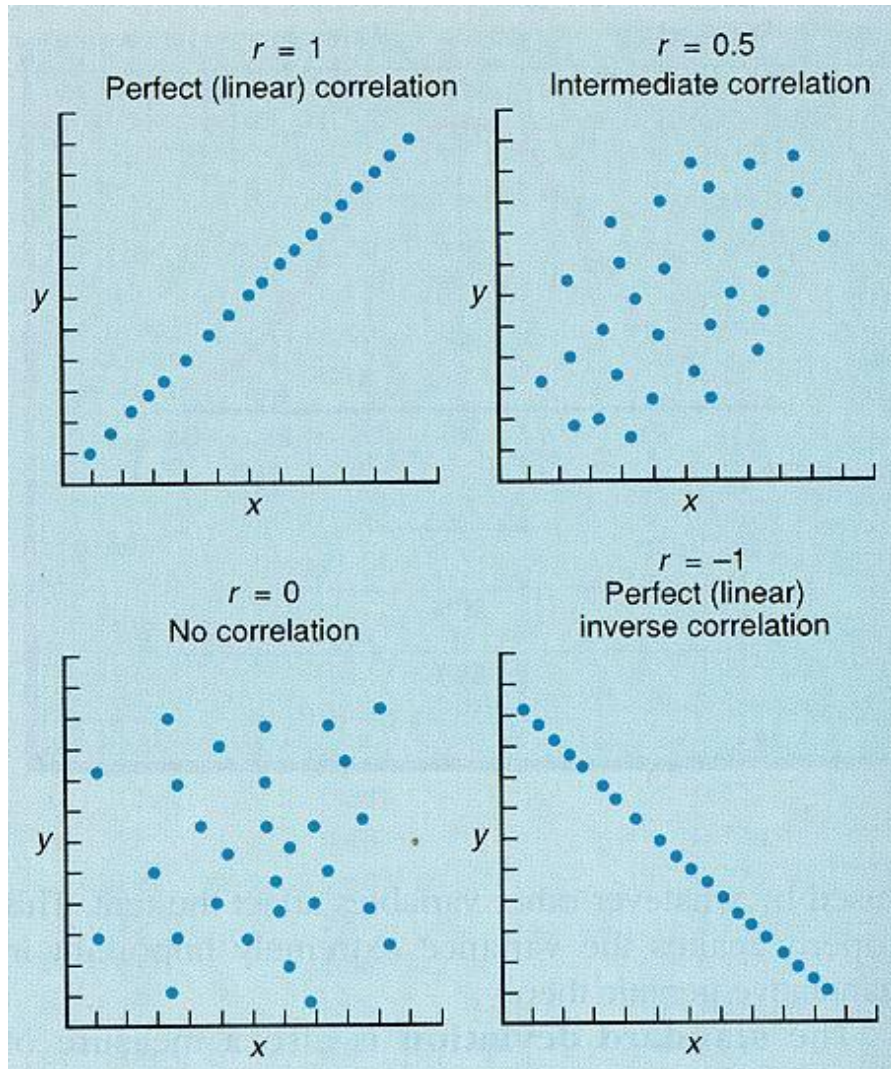
- a measure of the relation between two or more variables. Correlation coefficients can range from -1.00 to +1.00
- Regression lines
- Significance depends on sample size



# Correlation

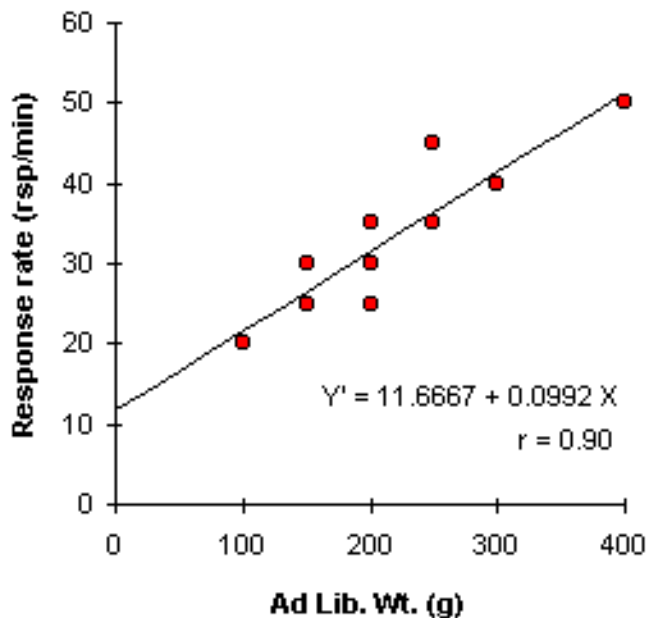
Pearsons  
correlation  
coefficient – linear

Spearman rank  
correlation -  
non-linear

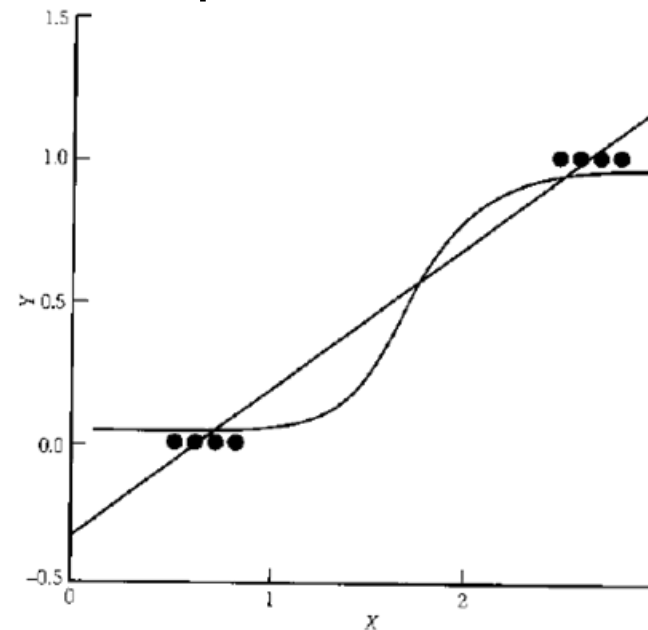


# Regression

- Linear: continuous dependent variable



- Logistic – binary dependent variable





# Univariate vs multivariate analysis

## Mortality in colorectal cancer by age

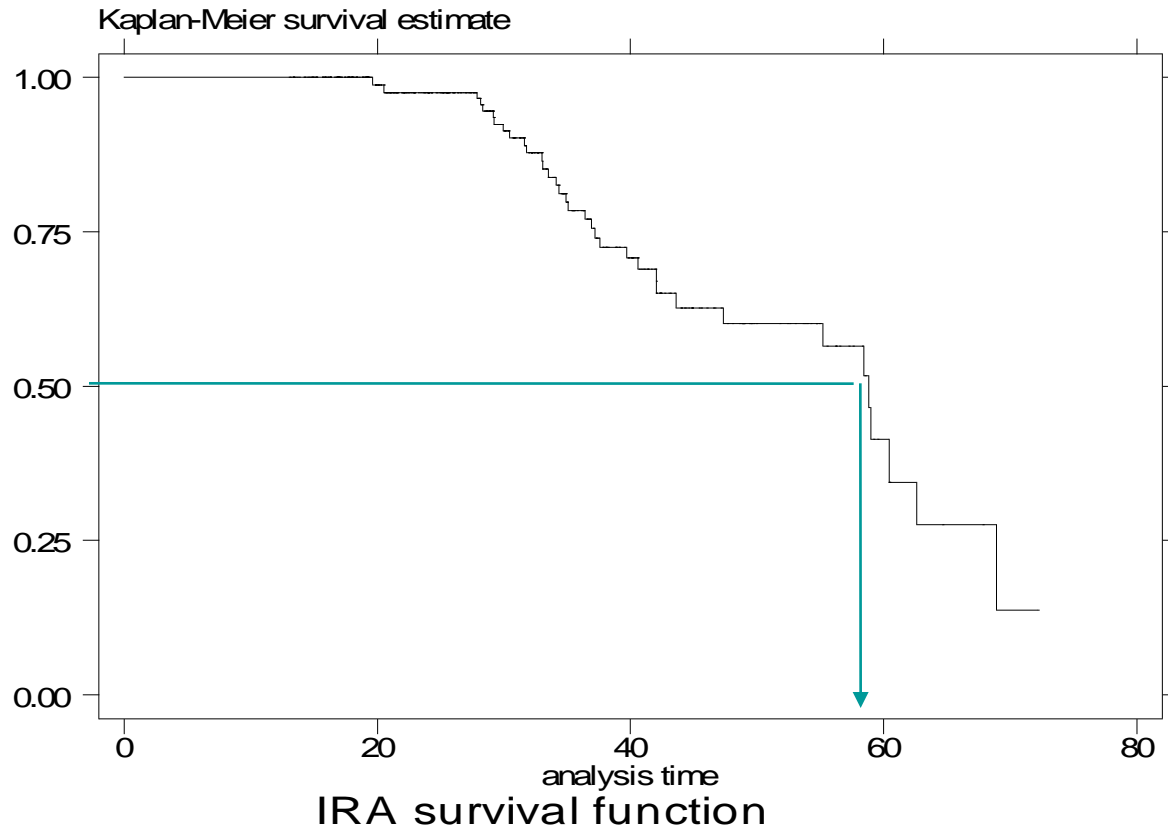
	Coefficient	P-value	Odds ratio	95% CI for OR
<i>Age (10-year<sup>-1</sup> increase)</i>	0.564	<0.001	1.757	1.600-1.931

## Mortality in colorectal cancer by age and mode of surgery

	Coefficient	P-value	Odds ratio	95% CI
<i>Age (10-year<sup>-1</sup> increase)</i>	0.570	<0.001	1.768	1.606-1.947
<i>Urgent</i>	0.947	<0.001	2.578	2.094-3.173
<i>Emergency</i>	1.446	<0.001	4.246	2.978-6.054



# Survival analysis

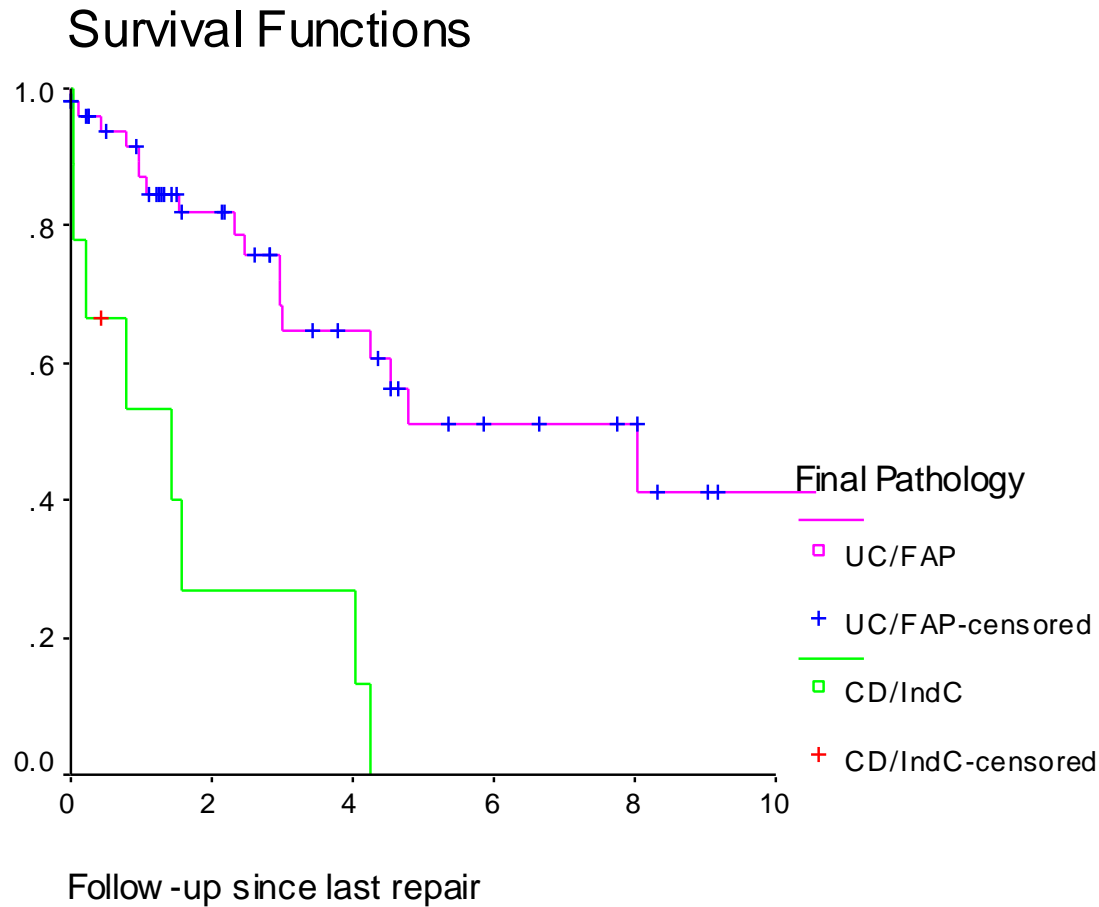


Censored data:      left (late start)  
                             right (lost to follow up, death)

Median survival (95% confidence intervals)



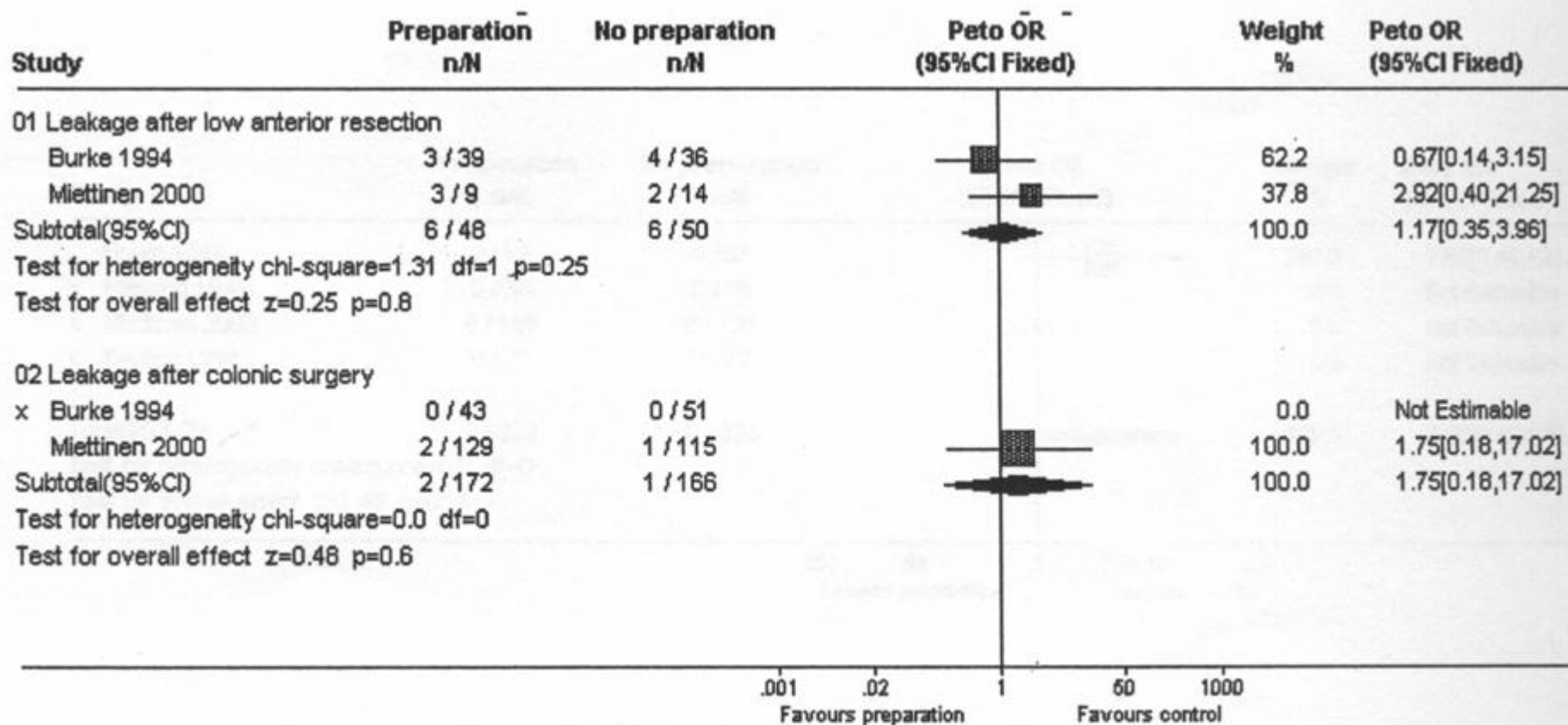
# Survival analysis – two sample data



Log-rank test 17.56, 1df,  $p < 0.001$



# Meta-analysis !



# Levels of evidence

- Level 1 – Randomised controlled trials
  - 1a Meta-analysis of RCTs
  - 1b RCT
- Level 2 - Non-randomised trials
  - 2a Control study no randomization
  - 2b Quasi-experimental study
- Level 3 – Descriptive study – case-control study
- Level 4 – Reports, individual opinion

