

BRITISH JOURNAL OF
**HOSPITAL
MEDICINE****MMC**
Modernising Medical Careers**MODERNISING
MEDICAL CAREERS****Basic statistics: a guide
for the foundation year doctor** M146*Walter Andreatta, Alka Kothari, Deepali Trivedi,
Rachel Hooke***Applied anatomy of
cricothyrotomy and tracheostomy** M148*Harold Ellis***The unwell patient on
haemodialysis: what you need to
know on an acute medical take** M152*WT Hinchliffe, JS Murray, NS Kanagasundaram***Understanding the metabolic
response to trauma** M156*Kaji Sritharan, Hilary Thompson***Good medical records:
a guide for the foundation
year doctor** M159*Kiran Singiseti, Helen Staley***So you want to be ...
a medical oncologist** M160*Martin Gore***IN NEXT MONTH'S
MMC SUPPLEMENT****The safe referral and transfer of
patients to major trauma centres****Preoperative assessment****The emergency airway****So you want to be...
a sport and exercise
medicine physician**

Basic statistics: a guide for the foundation year doctor

Introduction

Statistics encompasses the methods of collecting, summarizing, analysing and drawing conclusions from data (Petrie and Sabin, 2005; Greenhalgh, 2006).

Types of data

Categorical (qualitative) data

Each individual value belongs to one of a number of distinct categories of the variable:

- Nominal data – data one can name. They can be simply counted and not measured, such as male or female, pregnant or non-pregnant
- Ordinal data – the categories are ordered in some way, for instance, grading degree of pain, disease staging and various others.

Numerical (quantitative) data

Each variable takes some numerical value:

- Discrete data – the variable can only take certain whole numerical values, such as number of episodes of illness, number of babies born
- Continuous data – there are no limitation on the value that the variable can take, for example, height or weight.

Derived data

Various types of data may be encountered in the medical field, such as percentages, ratios or quotients, rates or scores.

Summarizing data

Diagrams or graphical presentations can provide simple and quick summary pictures, spotting outliers and trends before any formal analysis is performed

Dr Walter Andreatta is Foundation Year 2 Doctor, Ms Alka Kothari is Research Nurse, Miss Deepali Trivedi is Ophthalmologist, Sandwell and West Birmingham Hospitals NHS Trust, Birmingham B18 7QH and Dr Rachel Hooke is Working Time Directive (WTD) Implementation Manager, Airedale NHS Trust, Steeton, Keighley, West Yorkshire

Correspondence to: Miss D Trivedi

(Campbell and Machin, 1999; Petrie and Sabin, 2005).

Frequency distribution of a variable can be displayed graphically using:

- Bar or column chart
- Pie chart
- Histogram
- Dot plot
- Stem-and-leaf plot
- Box plot
- Clustered or segmented bar chart, which is useful for two variables
- Scatter diagram, which is useful for two variables.

Data are summarized to condense the information by providing measures that describe the important characteristics of the data. There are mainly two measures that summarize continuous variables:

Measure of location or average

- Mean – derived by adding up all the values in a data set and dividing the sum by the number of values in the data set. Although it uses up all the data values, the result could be distorted by skewed data or outliers
- Median – calculated by arranging the data in order of magnitude in ascending order, and the middle value is the median. Although the result is not distorted by outliers, it may ignore most of the information
- Mode – the value that occurs most frequently in a data set. This method may ignore most of the information but may be useful when the data are skewed.

Measure of spread or variability

- Range – the difference between the largest and smallest observations in the data set. This method uses only two values from the data and the results can be easily distorted by outliers.

The interquartile range is the numerical difference between the values within the first and the third quarter of the data when arranged in ascending order of magnitude. This method encompasses the middle 50% of all

values in the distribution. It is unaffected by outliers and appropriate for skewed data.

- Standard deviation – determines how deviated or spread out the data is from its mean. If the data values are close to the mean, the standard deviation is small, whereas if the data values are far from the mean, the standard deviation is large.

Statistical or clinical significance

Statistical analysis is concerned not only with summarizing data but also with investigating relationships. In any randomized control trial, the difference found between the groups might be statistically and/or clinically significant. For instance, one treatment might be superior to the other statistically (statistical significance) but might not show much clinical difference (clinical significance). On the other hand, one treatment might show obvious clinical difference but no statistical significance, particularly if the study is too small to demonstrate this effect.

There are various methods that might help to establish the significance of a study (Campbell and Machin, 1999; Petrie and Sabin, 2005), which include:

Confidence interval

The range of values within which one is usually 95% confident that the true population parameter lies. The upper and lower limits of the confidence interval provide a means of assessing whether the results are important. In other words, a wide confidence interval indicates that the estimate is imprecise and a narrower range determines a more accurate estimate. The width of the confidence interval depends on the sample size and variability of the data.

P value

In a study, the null hypothesis means that there is no effect or no difference between the two parameters being studied. The *P* value is an alternative hypothesis which rejects the null hypothesis. In other words, the *P* value determines how likely it is that a particular result in a study has occurred by chance alone, assuming that the null hypothesis is true and there is no difference between the outcomes of the parameters measured. Conventionally, if the *P* value

is less than 0.05, it is considered that there is significant evidence to reject the null hypothesis, as there is only a small chance of the result occurring if the null hypothesis is true. Thus, the null hypothesis is rejected and it is stated that the results are significant at the 5% level. In contrast, if the *P* value is greater than 0.05, there is not enough evidence to reject the null hypothesis and the results are insignificant.

Sample size and power

The sample size is the number of people that are included in a study. An appropriate sample size should be chosen that will provide evidence to support or reject the null hypothesis.

The power of a study is a measure to find a certain size of difference between the group being compared, assuming such a difference does exist. The power of a study gives a good chance of detecting a clinically relevant effect, if one exists.

Relative risk

The risk of an event is the probability that an event will occur within a stated time period compared with the total number of possible events. Relative risk tells you how many times more likely it is that an event will occur in the treatment group relative to the control group:

- Relative risk = 1 → no effect
- Relative risk <1 → the treatment decreases the risk of the outcome
- Relative risk >1 → the treatment increases the risk of the outcome.

Absolute risk reduction gives an indication of the baseline risk and treatment effect. Absolute risk reduction is the amount by which a treatment reduces the risk of an event.

Relative risk reduction indicates the reduction in the rate of the outcome in the treatment group relative to that in the control group.

Odds

The odds of an event occurring is another way of describing the chance of it happening. The odds ratio is calculated by dividing the odds of an outcome in one group by the odds of an outcome in another group.

Number needed to treat

This represents the number of patients who need to be treated with the experimental therapy in order to prevent one adverse event. The reciprocal of absolute risk reduction gives the number of patients needed to be treated for a defined period in order to prevent one unwanted outcome.

There are several readily available computer packages such as SPSS (Statistical Package for the Social Sciences) which help in analysis and calculation of small or large data.

Conclusions

Statistics are useful for analysing and interpreting data. Foundation doctors should be familiar with the different methods which are used. **BJHM**

Conflict of interest: Dr Hooke has worked in both management and medicine. Her views are her own and do not necessarily reflect those of her employer or any other organization that she is associated with.

- Campbell MJ, Machin D (1999) *Medical Statistics: A Commonsense Approach*. 3rd edn. John Wiley & Sons Ltd, Chichester
- Greenhalgh T (2006) *How to Read a Paper: the Basics of Evidence-based Medicine*. 3rd edn. Blackwell Publishing Ltd, Oxford
- Petrie A, Sabin C (2005) *Medical Statistics at a Glance*. 2nd edn. Blackwell Publishing Ltd, Oxford

Further reading

- Anonymous (2009) Understanding statistical terms: 1. *Drug Ther Bull* 47: 22–4
- Anonymous (2009) Understanding statistical terms: 2. *Drug Ther Bull* 47: 35–6
- Anonymous (2009) Understanding statistical terms: 3. *Drug Ther Bull* 47: 59–60

KEY POINTS

- Foundation doctors need to understand basic statistical methods.
- Statistics help to interpret data.
- Data can be qualitative, quantitative or derived.
- There are different ways of presenting data.
- The *P* value helps determine if results are statistically significant.