# Clinicians' guide to reading psychiatric literature: therapeutic trials and systematic reviews

## James Warner

Getting to grips with critical appraisal is not easy. To some extent, it is a metaphorical bandwagon, overloaded with dogma, misinformation and mathematics. Its introduction into the MRCPsych curriculum (Royal College of Psychiatrists, 2001) has left many trainers trying to catch up with their trainees' skills. This paper attempts to demystify some of the concepts, provide a balanced overview of the advantages and drawbacks of critical appraisal, and discuss ways of developing skills in critical appraisal.

We are currently in the middle of a shift in how medicine is practised. Medical knowledge is expanding at such a rate that it is impossible to keep pace with all advances. To keep up with the knowledge base in a single speciality requires several hours' reading each week. Not only do most clinicians have difficulty in finding the time to do this, it is also unlikely that any individual would be able to retain all the new information. The traditional top-down method of absorbing all relevant facts and selectively using them in a particular clinical scenario is being replaced by a bottom-up approach, advocated by proponents of evidence-based medicine (EBM), where up-to-date, good-quality information is sought to address a particular clinical situation (Sackett *et al*, 1996). This new approach requires new skills, time, resources and courage, as we often have to admit to patients that we do not know the answer to a particular question and need to look it up.

This paper is not about the whole of EBM. It is about how to appraise an article or paper that may inform clinical practice in relation to therapeutic trials. This is only one part of the EBM process.

## Reading medical literature

The temptation when reading any paper is to concentrate on the abstract, introduction and discussion, perhaps scan the results, and leave out the methods section. Some journals encourage this by putting the methods in a smaller font or a box. This is probably the wrong way to read a paper. The introduction is usually a brief review intended to set the study in context and outline the aims. It will be too short to be systematic or balanced. It is likely to take several minutes to read and the pay-back for this time commitment will be minimal. Many journals are now severely restricting the length of introductions to avoid over-long, redundant literature reviews.

The discussion will review the results, with the authors' interpretation. This may be interesting but it may not be balanced and it should be interpreted with caution. What is really important about any paper that describes a study is how much store the reader is willing to set by the results. This requires three separate judgements (Guyatt *et al*, 1993, 1994): How valid is the study? What are the results? Can I use the results with my patient or in my practice?

Referring to Boxes 1 and 2, which summarise the main points of critical appraisal for a single trial and a systematic review, may be helpful, both when reading this paper and when assessing published evidence.

James Warner is a senior lecturer in old age psychiatry at Imperial College School of Medicine (Department of Public Mental Health, Faculty of Medicine, Imperial College of Science, Technology and Medicine, Paterson Centre, South Wharf Road, London W2 1PD, UK; tel: 020 7886 1655; fax: 020 7886 1995; e-mail: j.warner@ic.ac.uk) and an honorary consultant in old age psychiatry at Brent Kensington Chelsea and Westminster Mental Health NHS Trust. His interests include dementia and sexual disorders in the elderly and he has extensive experience in teaching evidence-based practices.

## Box 1  Critical appraisal of a paper on treatment

### Are the results valid?

*Was assignment of patients randomised and was randomisation concealed?*
Is there an explicit statement about the method of randomisation? Was the method open to bias? Was it possible that investigators had control over allocation of of subjects? Computer-generated independent randomisation methods are best.

*Were all patients who entered the trial accounted for at the conclusion?*
All trials have subjects who drop out. If drop-outs are ignored, the distillate of completers may be very atypical. Results should usually be reported on an intention-to-treat (ITT) basis, where everyone who was randomised is included in the final analysis. Common methods of ITT include last observation carried forward (LOCF).

*Were patients, clinicians and study personnel blind to treatment allocation, and was blinding assessed?*
If subjects are not blind to their treatment, they may (unintentionally) give misleading answers to outcome questionnaires. Likewise, study personnel and clinicians may influence outcome if they are aware of treatment allocation. It is not enough to assume that initial blinding will be preserved. Subjects and clinicians alike can be de-blinded by adverse effects. One way to assess this is to ask all involved in the study to guess their allocation.

*Were the groups similar at the start of the trial?*
Often now there are no statistical comparisons on tables of baseline data. However, differences between groups at baseline can occur by chance or may indicate non-random allocation to groups. If differences are present, think how they may influence the results.

*Apart from the intervention, were groups treated equally?*
Apart from the intervention being assessed, subjects should have similar treatment. Look especially for different adjunctive treatments between groups, and for different management or investigations by the study team.

*Did the study have adequate power?*
Was there a pre-trial power calculation? This is especially important if the study did not show significant differences between groups. It is important to be able to judge whether no difference was found because there truly was none, or whether the study was too small (type II error).

### What are the results?

*How large is the treatment effect?*
What is the number needed to treat (NNT)? This is a four-stage process:
(1) calculate the proportion of subjects improving in the experimental group (the experimental event rate, EER) – this may be given or be derivable;
(2) calculate the proportion of subjects improving in the control group (control event rate, CER);
(3) subtract (1) from (2) to give the net effect that treatment confers over the control (absolute risk reduction, ARR);
(4) divide 1 by the ARR (the reciprocal) to give the NNT.

*How precise is the treatment effect?*
What are the confidence intervals? The following equation will give the 95% confidence limits of the ARR.

$$ARR \pm 1.96 \times \sqrt{\frac{CER \times (1 - CER)}{\text{no. of control points}} + \frac{EER \times (1 - EER)}{\text{no. of experimental points}}}$$

The reciprocal of these numbers gives the confidence limits for the NNT.

### Will the results help me in caring for my patient?

*Are my patients similar to those in the trial?*
Consider exclusion factors of the study being appraised (medication, general health, etc.) and the population from which the sample was drawn (e.g. are they the same age and ethnicity as your patient?).

*Were all clinically relevant outcomes considered?*
Do the authors report adverse events as well as benefits.

*Are the benefits worth the harms and costs?*
This is a matter of judgement. Think not only of monetary cost, but of the risk–benefit ratio of the intervention.

*What does my patient think?*
A most important question. The clinician is now in a position to take the results to a patient and discuss the evidence.

### Conclusion
*Where on the continuum between very good and very bad does the study lie? Decide for yourself.*

Very good ⊢————————————————————————————————⊣ Very bad

**Box 2  Critical appraisal of a systematic review**

**Are the results valid?**

| | |
|---|---|
| *Does the review address a focused, relevant question?* | Systematic reviews that attempt to address a broad question are unlikely to be very helpful, because the gamut of evidence is too wide, and the presentation will probably be confusing. |
| *Does it include a methods section that describes finding and including all the relevant trials?* | A good systematic review should identify all relevant research literature, whether published or not. Look for the databases searched (Medline alone is insufficient), the contacts with researchers in the field and drug companies for unpublished data, and cross-checking references in reviews. Details of the search terms (headings and keywords) should be provided and should appear comprehensive. |
| *Does it include a methods section that describes assessing their individual validity?* | After all relevant studies have been identified, many are then likely to be excluded because of poor study methods, lack of usable data, etc. The decision to exclude studies should be based on specified criteria, and made by two authors blind to each other's decisions. Studies are often given a score depending on validity factors such as randomisation, allocation concealment and blinding. |
| *Were the results consistent from study to study?* | Heterogeneity is when results differ significantly between studies. It may be due to differences in methods, study populations, etc., or to a statistical quirk. If heterogeneity is present, different statistical techniques should be used for the meta-analysis and the authors should attempt to explain its presence. |

**What are the results?**

| | |
|---|---|
| *What is the odds ratio or relative risk or difference between means?* | Results of continuous data may be expressed as a difference between means weighted for the size of the study (weighted mean difference, WMD) or for categorical variables as an odds ratio (OR) or relative risk (RR). Results are frequently displayed on a Forest plot ('blobbogram') that shows the WMD, RR or OR and confidence interval (CI) or each study, and the result and CI for the meta-analysis (sometimes known as the pooled effect). |
| *What are the confidence intervals?* | If the 95% CI of a WMD crosses 0 or, for an OR or RR, crosses 1, then the result is not significant at the 5% level. On Forest plots the vertical anchor-line represents no effect; if a CI crosses this, the result is not statistically significant. |

**Is the evidence applicable to my patient?**

| | |
|---|---|
| *Is your patient similar to those described in the review?* | This information is often not given in sufficient detail, because each study that is included in the review is only briefly described. |
| *Are the benefits worth the harms and costs?* | A systematic review, especially in the Cochrane database, will often provide large amounts of data on adverse effects that can help the clinician and patient decide which way the risk–benefit equation tips. |
| *What are my patient's preferences?* | It may be worth putting the evidence to the patient (or carer) and asking his or her opinion. |

**Conclusion**

*Where on the continuum between very good and very bad does the systematic review lie? Decide for yourself.*

Very good ⊢─────────────────────────────────────────────────┤ Very bad

# Critical appraisal

Two things are important when assessing a piece of evidence. First, by reading the methods section, it is possible to appraise the way the study was conducted and judge whether the study is applicable to a particular clinical setting. The methods should provide sufficient information to make judgements about the process of randomisation, concealment of treatment, completeness of follow-up, how missing data were accounted for and the representativeness of the sample in relation to patients in a particular clinical setting (Table 1). Second, it is important to assess how significant the results of the study are in clinical as well as statistical terms. In other words, are the results sufficiently impressive to alter clinical practice? Different types of study best address different clinical questions and these are outlined in Table 1. Here I concentrate on evidence relating to treatment (which is the subject of the bulk of

| Table 1  Types of study method, themes of evidence and statistical concepts | | |
|---|---|---|
| Study method | Theme of evidence | Statistical concepts |
| Systematic review | Therapy<br>Diagnosis<br>Prognosis | Meta-analysis, pooled-effect size (can be odds ratio, relative risk or mean difference) |
| Randomised controlled trial | Therapy<br>Harm (as derivative) | Control event rate, experimental event rate, absolute risk reduction, number needed to treat |
| Cross-sectional survey | Diagnosis<br>Epidemiology | Sensitivity, specificity, positive predictive value, negative predictive value, likelihood ratio |
| Cohort study | Prognosis<br>Aetiology<br>Harm<br>Epidemiology | Incidence, prevalence, relative risk, odds ratio, standardised mortality ratio |
| Case–control study | Prognosis<br>Aetiology<br>Harm<br>Epidemiology | Odds ratio |

'evidence' in the literature) and focus on single, randomised controlled trials and systematic overviews.

## The single randomised trial

### Validity

The randomised controlled trial (RCT), when subjects entering a trial are randomly allocated to one of two (or more) groups, is currently regarded as the cornerstone of evidence about whether a treatment is useful. However, other types of evidence that are considered below the RCT in the evidence hierarchy may still provide useful information.

The main purpose of randomisation is to ensure comparability of groups and reduce bias. If two groups differ with respect to a socio-demographic or disease characteristic at the start of a trial, that characteristic may influence the result. For example, a study investigating the efficacy of an antidepressant may show spurious results if one group was more depressed than the other at the start of the study. Several potential sources of bias (systematic error), such as selecting subjects for one or other treatment, should also be reduced by randomisation and, wherever possible, concealment of allocation. Having a control group (which should be managed identically to the intervention group in all respects except the treatment under scrutiny) enables the net effect of the intervention to be estimated, as any other changes that may influence the result (effect of time, participation in the study, etc.) should occur equally across both groups.

It is highly unlikely that any study will collect a full set of data on all subjects who were randomised. Patients may drop out of a study for a variety of reasons (attrition). For example, they may get better, move away, die, stop because they cannot tolerate the drug, or lose patience with the study. These missing subjects should be properly accounted for in the final result. Failure to do so may introduce bias. For example, if 20% of patients in a therapy trial drop out because they get worse and these subjects are excluded from the analysis, the efficacy of the drug will be overestimated. There is no ideal way of accounting for subjects who go missing from a trial, but a common method is 'last observation carried forward' (LOCF), where the last set of results obtained on a subject are used as if these were the results at the end of the trial.

When critically appraising the validity of a study, the emphasis should be on the method of randomisation, the quality of concealment (blinding) and the way attrition is accounted for. Unfortunately, owing to constraints of journal space, poor authorship or poor editorial practices, the information needed to make these judgements is often missing. A simple statement such as "subjects were randomised" will provide no information as to whether the process was robust, for example using a computer-generated code, or flawed, such as a roll of a dice. The reporting of clinical trials is improving and less guesswork is required in more recent papers.

### Results

Many papers comparing a treatment and placebo, or two alternative treatments, present the results as a difference in means between the two groups, for example, a difference between groups in mean scores on the Brief Psychiatric Rating Scale (Overall & Gorham, 1962). This can be misleading. First, if the study is large enough, even very small differences

of any outcome between groups will be statistically significant. Furthermore, how often do clinicians assess a patient's response by changes in a scale or score?

It is often more useful to present data that are clinically meaningful, such as the number of subjects who got better, or left hospital, or did not relapse or some other categorical outcome. Categorical outcomes can be used to calculate how many individuals need to receive treatment rather than placebo in order that one additional person achieves the desired outcome: the number needed to treat (NNT). This is a useful method of deciding how clinically useful a treatment is. Authors rarely report results in terms of the NNT. The percentage of subjects achieving the desired outcome in each group is often provided, and calculating the NNT from these is a simple process; there is no complicated mathematics involved. Any results expressed as proportions or percentages can be used to calculate NNTs.

For example, a trial may report that 40% of subjects in the treatment arm 'improved', compared with 20% in the control group. The difference between the two groups is 20% (i.e. this is the net improvement conferred by the treatment). The net effect of the treatment (over the control) was to improve one person by 20%, therefore five people need to be treated to make one whole person 'improved' (i.e. the NNT is 5). The easy way to do this is to divide 100 by the difference between groups (or divide 1 by the difference, if the results are expressed as proportions). Numbers needed to treat can vary between 1 (very good) and infinity. The value of an NNT depends on the effort and cost of achieving the outcome weighed against the relative importance of that outcome. Numbers needed to treat can also used to compare the relative effectiveness of two similar treatments. The 95% confidence intervals (when provided, or if they can be calculated) will provide a range of values within which the true value of the NNT is likely to lie.

### Applicability

One of the biggest drawbacks of the RCT is whether the results translate to clinical practice. Efficacy is the measure of how useful a treatment is under ideal trial conditions, effectiveness is how well a treatment works in routine clinical settings.

For example, it is likely that only a small proportion of patients with depression will be prepared to take part in a 12-week controlled trial. Furthermore, all trials exclude certain people, often the very severely ill or suicidal and those with dual diagnosis. The distillate after these filters is usually highly selected and may behave very differently from the average patient. Some feeling of applicability can be derived by comparing the demographic and disease characteristics of the trial subjects with patients in a clinical setting, but no trial will replicate a true clinical setting, warts and all.

## Systematic reviews

A systematic review attempts to identify all evidence (published and unpublished) on a particular topic and appraise the validity of the evidence using 'quality filters'. Often the evidence from the individual studies that have passed the validity test is then combined into a single unified result: a meta-analysis. So, instead of reading about 10 trials each with 100 subjects, the meta-analysis will report the combined result for 1000 subjects. Increasing the number of subjects should increase the precision of the final result because there is greater confidence that the results reflect the 'true' value. However, meta-analyses have been criticised because the results can mislead, owing to publication bias (studies with negative results are less likely to be published, or are published in obscure journals), and are guilty of the dichotomous approach to evidence (addressed below). Systematic reviews and the meta-analyses contained therein can be helpful, but they should be viewed in the same way as any other source of evidence and be critically evaluated.

A particularly good source of systematic reviews is the Cochrane database. This has comprehensive systematic reviews on many subjects, updated four times per year. It is easy to search and very user-friendly. Many libraries have Cochrane reviews on disk or they can be accessed on-line through a variety of sources, including some local health authority websites (http://www.doh.gov.uk) or through Doctors.net.uk (http://www.doctors.org.uk), which is free to registered doctors and medical students. Systematic reviews and meta-analyses are not available just for studies on therapy. Increasingly, they are done for diagnostic and prognostic data.

### Validity

The critical appraisal of a systematic review focuses on the quality of the review, not the quality of the component studies (Box 2). The main points are whether the authors of the review identified all relevant trial data worldwide and adequately assessed each study they did identify. Systematic reviews should include an explicit statement about which databases the authors searched, what search terms they used and how they identified unpublished data (by contact with researchers in the field, pharmaceutical companies and the clinical trials register). Each study identified for the review should be independently scrutinised by two authors to decide which are of sufficient quality to be included in the

meta-analysis. The criteria for assessing the quality of individual studies vary between reviews, but usually include evaluation of randomisation and blinding. Studies considered for a systematic review are either included or not, and the quality of individual studies has no bearing on the weight attached to their results once they are included. This dichotomous approach may result in a loss of useful evidence in what the reviewers perceive to be substandard trials and may introduce bias to the meta-analysis.

Another issue with systematic reviews is that of heterogeneity. Heterogeneity is when the results of one or more studies in a meta-analysis differ significantly from the others. This may be caused by variation between studies, such as differences in the methods or study population, or may occur by chance. Authors of the systematic review should identify whether heterogeneity is present and seek to explain it. The presence of heterogeneity should also result in a different statistical approach to the meta-analysis.

### Results

Results of meta-analyses are often expressed graphically (Forest plots), on which the results and confidence intervals of individual studies are plotted together with a combined effect (Lewis & Clarke, 2001). Results are expressed as differences in means for continuous variables (such as questionnaire scores) or odds ratios for categorical results (such as better/not better). The pooled effect size (i.e. combined result of the included trials) is represented as a lozenge-shaped bar at the bottom of the graph. The vertical line (at 0 for differences in means, 1 for odds ratios) indicates no effect. If the horizontal confidence interval bars cross this line, the results are not statistically significant.

### Applicability

As with single therapeutic trials, the main questions here are whether the samples in the studies included in the meta-analysis are similar to a patient or group of patients, and whether the magnitude of the effect, weighed against the costs (human and financial), means that the treatment is worth a try.

# Developing critical appraisal skills

Reading this paper is insufficient preparation for how to read the literature in an efficient and effective way. However, there are strategies and resources that will help develop and maintain these skills, and some of these are shown in Box 3.

Even if none of these methods is available, simply concentrating on the validity of a paper's study methods and assessing the results for clinical relevance will enhance skills over time. Books by Sackett *et al* (2000) and Lawrie *et al* (2000) give sound guidance on critical appraisal, and a list of useful websites is given in Box 4.

# EBM tyranny

Judgement is an essential part of the EBM process. Some advocates of EBM, who seem to want to tyrannise rather than inform, suggest that all medical studies should be dichotomised into 'good' or 'bad', with the implication that bad studies should be ignored. This would result in ignoring many medical studies. This is not what EBM is about. All medical studies lie on a continuum between the very good and very bad. Where a study is on this continuum is a matter of judgement. Almost

---

**Box 3  Strategies and resources for developing and maintaining skills in critical reading**

*Courses*  May last up to 1 week and can provide a comprehensive introduction to evidence-based medicine (see the Centre for Evidence Based Mental Health's website in Box 4)

*Local evidence-based journal clubs*  A useful way to have regular practice (Warner & King, 1997). Each week the journal presentation could be a structured critical appraisal rather than an unstructured critique. This is particularly attractive and salient to trainees preparing for the MRCPsych examination.

*Virtual journal clubs* Every month each club member scans an allocated journal, selects one article of interest and appraises it. The appraisal is then distributed to others in the group. Thus, in a club of five members, each will receive four other critically appraised articles each month. This method could work well using e-mail.

*Introduction of evidence-based practices to clinical services* Each consultant team could set, for example, one critical appraisal exercise each week involving trainees and other disciplines.

all studies will be useful to some degree and the judgement comes in deciding just how 'good' or 'bad' a particular study is, and consequently how much that study will inform practice. It is worth noting that the good reputation of the journal or of the authors is no guarantee of quality and that studies funded by the pharmaceutical industry do tend to report better results than independent researchers.

Another example of EBM tyranny is the devaluing of evidence of therapeutic efficacy that is not from an RCT. It is a fallacy that only RCTs can yield useful information on therapy. For example, faced with a particular clinical question, the only evidence available may be from a series of case reports. This is better than no evidence and should be given some credence. In open-label studies without blinding or randomisation bias may be a significant problem, but the evidence presented should not be disregarded: it can be interpreted in the light of potential flaws. Even RCTs may not be perfect, and you should consider the validity, results and applicability of each piece of evidence found.

## EBM nihilism

There are two main sources of nihilism that most EBM practitioners (except the evangelists) will have to negotiate. First, there will be times when it is not possible to find a piece of evidence relating to a question, and second, there will be many times when the evidence found is towards the 'very poor' end of the spectrum owing to problems with validity or applicability. Recognising that the evidence base for

---

**Box 4  Useful websites**

*http://www.cebmh.com* Centre for Evidence-Based Mental Health. Useful for critical appraisal tools, information about courses and links to the National Electronic Library for Mental Health

*http://www.bma.org.uk* Free on-line library and search facilities for BMA members

*http://www.doctors.org.uk* Free access to Cochrane database, search facilities and other useful resources, including *Bandolier*

*http://www.omni.ac.uk* Organised Medical Networked Information. A very useful site that helps sift out poor-quality information from internet searches

---

much of clinical practice is inadequate can be very therapeutic. It is also helpful to bear in mind that evidence does not necessarily 'trump' clinical experience and the absence of evidence is not the same as evidence of absence.

## The future

There are already several journal sources of critically appraised evidence, including *Clinical Evidence*, *Evidence-Based Mental Health* and *Bandolier*. With the increase of medical knowledge, improvements in study design and reporting of results, and better information technology, there is likely to be a move towards the use of instant, high-quality information to guide medical practice. Evidence-based technologists may become integral to ward rounds and community meetings. Armed with a laptop computer with internet access, it is now possible to search, identify and download full-text articles from thousands of medical journals. These can then be critically appraised and an answer to a clinical question provided within an hour or so. EBM is not a fad, and the need to embrace it, master the basic skills and use it sensibly is becoming essential.

## Conclusions

This paper is a guide to effective reading of the medical literature. It will take practice and dedication for the skills to become second nature. Remember, research is very hard to do and very easy to criticise, so be gentle with your criticisms.

## References

Guyatt, G., Sackett, D. & Cook, D., for the Evidence-Based Working Group (1993) Users' guide to the medical literature. II. How to use an article about therapy or prevention. A: Are the results valid? *Journal of the American Medical Association*, **270**, 2598–2601.

—, — & —, for the Evidence-Based Working Group (1994) Users' guide to the medical literature. II. How to use an article about therapy or prevention. B: What were the results and will they help me in caring for my patients? *Journal of the American Medical Association*, **271**, 59–63.

Lawrie, S. M., McIntosh, A. M. & Rao, S. (2000) *Critical Appraisal for Psychiatry*. London: Churchill Livingstone.

Lewis, S. & Clarke, M. (2001) Forest plots: trying to see the wood and the trees. *BMJ*, **322**, 1479–1480.

Overall, J. E. & Gorham, D. R. (1962) The Brief Psychiatric Rating Scale. *Psychological Reports*, **10**, 799–812.

Royal College of Psychiatrists (2001) *Curriculum for Basic Specialist Training and the MRCPsych Examination* (Council Report CR95). London: Royal College of Psychiatrists.

Sackett, D. L., Rosenberg, W. M. C., Muir Gray, J. A. , *et al* (1996) Evidence-based medicine: what it is and what it isn't. *BMJ*, **312**, 71–72.

—, Straus, S., Richardson, W. S., *et al* (2000) *Evidence-Based Medicine* (2nd edn). London: Churchill Livingstone.

Warner, J. P. & King, M. (1997) Evidence-based medicine and the journal club: a cross-sectional survey of particiants' views. *Psychiatric Bulletin*, **21**, 532–534.

# Multiple choice questions

1. An RCT reported remission rates for depression of 65% for fluoxetine and 60% for imipramine; the NNT is therefore:
   a  60
   b  20
   c  15
   d  10
   e  5.

2. Important factors when assessing the validity of an RCT include:
   a  the patient's opinion of the treatment
   b  the method of randomisation
   c  the adverse-effect profile of the drug
   d  the power of the study
   e  how missing subjects were accounted for.

3. The following statements are true:
   a  a meta-analysis is always better evidence than a single randomised trial
   b  studies with large numbers tend to have wider confidence intervals
   c  Forest plots are a graphical representation of the studies in a systematic review
   d  in a meta-analysis, if the 95% CIs of a pooled weighted mean difference are –2.3 to –0.8 the results are statistically significant
   e  *P* values are a good indicator of clinical significance.

4. When making decisions about a patient's drug treatment, the following should be considered whenever possible:
   a  the patient's opinion
   b  the doctor's experience of using the drug
   c  the risk–benefit ratio of the drug
   d  the validity of the evidence for the treatment
   e  the effect size of the treatment.

5. The following evidence types are matched with a suitable study method:
   a  diagnostic test: case–control study
   b  prognosis: prospective cohort study
   c  therapy: cross-sectional survey
   d  aetiology: case–control study
   e  epidemiology: cross-sectional survey.

**MCQ answers**

| 1 | | 2 | | 3 | | 4 | | 5 | |
|---|---|---|---|---|---|---|---|---|---|
| a | F | a | F | a | F | a | T | a | F |
| b | T | b | T | b | F | b | T | b | T |
| c | F | c | F | c | T | c | T | c | F |
| d | F | d | T | d | T | d | T | d | T |
| e | F | e | T | e | F | e | T | e | T |