



The evidence that evidence-based medicine omits

Brendan Clarke ^a, Donald Gillies ^a, Phyllis Illari ^a, Federica Russo ^{b,c}, Jon Williamson ^{d,*}

^a Department of Science and Technology Studies, University College London, UK

^b Center Leo Apostel, Vrije Universiteit Brussel, Belgium

^c Centre for Reasoning, University of Kent, UK

^d Department of Philosophy, University of Kent, UK

ARTICLE INFO

Available online 27 October 2012

Keywords:

Evidence-based medicine

Methods

Philosophy

Evidence hierarchy

RCT

Mechanism

Evidence

ABSTRACT

According to current hierarchies of evidence for EBM, evidence of correlation (e.g., from RCTs) is always more important than evidence of mechanisms when evaluating and establishing causal claims. We argue that evidence of mechanisms needs to be treated alongside evidence of correlation. This is for three reasons. First, correlation is always a fallible indicator of causation, subject in particular to the problem of confounding; evidence of mechanisms can in some cases be more important than evidence of correlation when assessing a causal claim. Second, evidence of mechanisms is often required in order to obtain evidence of correlation (for example, in order to set up and evaluate RCTs). Third, evidence of mechanisms is often required in order to generalise and apply causal claims.

While the EBM movement has been enormously successful in making explicit and critically examining one aspect of our evidential practice, i.e., evidence of correlation, we wish to extend this line of work to make explicit and critically examine a second aspect of our evidential practices: evidence of mechanisms.

© 2012 Elsevier Inc. All rights reserved.

All studies are fallible

The EBM movement views evidence of mechanisms as poor quality evidence. (Terminology: 'Evidence based medicine is the conscientious, explicit, and judicious use of current best evidence in making decisions about the care of individual patients' (Sackett et al., 1996). 'A mechanism for a phenomenon consists of entities and activities organised in such a way that they are responsible for the phenomenon' (Illari and Williamson, 2012, p. 120). Here evidence of mechanisms may include evidence obtained by laboratory studies or previous statistical studies and tends to be relayed by expert testimony, e.g., via scientific publications.) This dim view of mechanistic evidence is most obvious when one refers to the 2011 Levels of Evidence table issued by the Oxford Centre for Evidence Based Medicine (OCEBM, 2011), which places 'mechanism-based reasoning' at level 5 – the lowest level – of the hierarchy of evidence. (Here, 'Mechanistic reasoning is an inferential chain (or web) linking the intervention (such as HRT) with a patient-relevant outcome, via relevant mechanisms' (Howick, 2011, p. 929)). Earlier evidence hierarchies, although often less explicit, also tend to leave only one possible place for prior evidence of mechanisms: the bottom level. For Canadian Task Force (1979, p. 1195), for instance, levels I and II are occupied by statistical studies and everything else is relegated to the bottom level: 'III: Opinions of respected authorities, based on clinical experience, descriptive studies or reports of expert committees'.

Higher up these hierarchies of evidence come various kinds of statistical study, including controlled studies, randomised controlled studies and, at the apex, systematic reviews and meta-analyses. The thought is that for these statistical studies, the higher up the hierarchy the better the evidence copes with the problem of confounding. (This is the problem that an observed dependence between *A* and *B* may be attributable to variation in *B*'s other causes, rather than variation in *A*.) But it is important to note that no statistical study solves the problem of confounding. A randomised controlled experiment only yields treatment/non-treatment groups that are homogeneous with respect to the putative effect's other causes in the asymptotic limit, as the number of individuals assigned to each of these two groups goes to infinity. So, while the evidence hierarchies may correctly identify the relative merits of various kinds of statistical study for dealing with the problem of confounding, in absolute terms these studies are all very fallible. There is room, therefore, for other kinds of evidence to influence decision making, even when high quality RCT evidence is available.

In particular, since statistical studies are fallible, strong evidence of mechanisms can sometimes override evidence gleaned by a statistical study that is high up in the hierarchy. This is especially clear when there is strong prior evidence that there is *no* mechanism linking the putative cause with the putative effect. In this case, the best remaining explanation is that the increase in the effect variable is due to confounding. Thus it can be reasonable to dismiss the claim that remote, retroactive intercessory prayer shortens length of stay in hospital, despite evidence from an RCT that yields a significant correlation between the two variables (Leibovici, 2001), on the grounds that current science holds no place for any mechanism that can

* Corresponding author.

E-mail address: j.williamson@kent.ac.uk (J. Williamson).

explain the putative effect in terms of the putative cause. Similarly for claims made in favour of precognition on the basis of a report of 9 experiments (Bem, 2011) – positive results which eventually turned out not to be replicable (Ritchie et al., 2012). Certain claims in favour of homoeopathy – including positive results from systematic review and meta-analysis (Cucherat et al., 2000) – can be treated analogously. While these examples are extreme, they clearly show that mechanistic evidence should not be confined to the bottom level of evidence hierarchies, but should, in certain cases, be considered alongside high-level statistical evidence.

Assessing RCTs

Let us now consider some simple historical reasons for our position, before returning to our philosophical argument in later sections. Historical examples suggest that the use of RCTs does not allow us to dispense with evidence of mechanisms, because evidence of mechanisms is needed to design and interpret RCTs. This is well illustrated by one of the first and most famous RCTs. This trial, carried out in the UK beginning in September 1947, was used to test whether streptomycin is an effective cure for pulmonary tuberculosis. The results of this RCT after 6 months were that, in the streptomycin group of 55, 28 (51%) had shown considerable improvement and only 4 (7%) had died, whereas, in the control group of 52, only 4 (8%) had shown considerable improvement whereas 14 (27%) had died (MRC, 1948, p. 771). Seemingly the RCT was an overwhelming success, but there was also some evidence that the bacilli responsible for the disease were developing resistance to streptomycin. This led the scientists conducting the trial to express caution, and to recommend that the observation of patients involved in the trial should be continued. This caution, based on evidence relating to the mechanism of the disease, proved to be amply justified. After 5 years, 32 of the 55 in the streptomycin group had died (58%), compared with 35 of the 52 in the control group (67%) (Florey, 1961, p. 133). The difference here is not statistically significant, and this showed that, over a longer period, streptomycin, on its own, was no better than existing treatments. If evidence of mechanisms had not been taken into account, the misleading impression that streptomycin on its own was an effective therapy would have been given, and this would have delayed the development of the first genuinely effective treatment, which was a combination of streptomycin with para-amino-salicylic acid (PAS).

As the streptomycin case demonstrates, evidence of mechanisms informs the design and interpretation of RCTs. While it is theoretically possible to conduct an RCT in the absence of evidence of mechanisms – as in the case of Leibovici (2001) – most clinical trials do evaluate interventions that are somewhat mechanistically understood. The same can be said for the manner in which this evaluation is performed – including the decision to clinically evaluate particular interventions; the way in which these interventions are carried out; and the measurement of the effects of these interventions. This makes evidence of mechanisms central to the business of conducting clinical trials. Note that this consideration of evidence of mechanisms does not mean that judgements of efficacy proceed on entirely mechanistic grounds (see e.g. Howick, 2011, p. 128).

Given that clinical trials typically seek to investigate novel interventions, the evidence of mechanisms upon which interventions and outcome measures rely is highly dynamic and rapidly changing. This can be demonstrated by estimating the age of such measures in contemporary clinical trials. Of the ten most-cited articles from the last five years of *The Lancet* (Scopus data April 21 2012), the age at publication of both interventions and outcome measures used was estimated. Of these ten articles, collectively cited 6132 times, three (Black et al., 2008; Daemen et al., 2007; Goldenberg et al., 2008) were identified as non-RCT publications, and excluded. From the remaining seven, two each dealing with HIV-AIDS (Bailey et al., 2007; Gray et al., 2007) and renal cell carcinoma (Escudier et al.,

2007; Motzer et al., 2008), and one each on HPV vaccination (Paavonen et al., 2007), vascular outcomes in diabetes (Patel et al., 2007), and breast cancer (Smith et al., 2007), a total of 35 intervention or outcome measures were identified (see Table 1, supplementary material). The age was estimated as per the methods discussed in the supplementary material. Where there was doubt about the introduction of a particular measure, the oldest recorded instance was used. The average age at the time of publication is 15 years, excluding those interventions or outcome measures thought to be older than 100 years, with the youngest intervention ranging between 1 and 23 years (mean = 10).

The sheer novelty of these critical parts of trial construction indicate that, far from being background or common knowledge, the evidence used to build and interpret trials changes rapidly. Given that a central principle of EBM practice is the “...conscientious, explicit, and judicious use of current best evidence...” (Sackett et al., 1996), we suggest that evidence of mechanisms should therefore be subject to the same process of systematic critical appraisal as evidence gleaned from trials themselves.

External validity

Even if we grant the soundness of an RCT, a question remains about its external validity. There is no a priori reason why the results of an RCT should be straightforwardly applicable to another population. This concerns medical treatments as well as policy actions. This problem is thoroughly discussed by Victora et al. (2004), where the authors point to several issues that hinder the external validity of RCTs. In particular, the authors dispute that the internal validity of an RCT also ensures its generalisability. The assumption that it does follows, Victora et al. explain, from the assumption of ‘universal biological response’. Victora et al. (2004) challenge this view and argue that although this assumption might well be hold for “interventions with short causal pathways”, it is certainly not the case for “interventions involving long, complex causal pathways, or in large-scale evaluations where these pathways can be affected by numerous characteristics of the population, health system, or environment”, such as policy interventions. In fact, there might be two threats to successful extrapolation in the case of policy: one is “behavioural effect modification” and the other is “biological effect modification” (i.e., respectively, “differences in the actual dose of the intervention delivered to the target population” and “differences in the dose–response relationship between the intervention and the impact indicator”).

Cartwright (2011) makes a similar point and illustrates it with the example of the ‘Bangladesh Integrated Nutrition Policy (BINP)’, a programme that largely failed to have an impact on child nutrition, although a very similar programme proved highly successful in Tamil Nadu (TINP – the Indian Tamil Nadu Integration Project). Cartwright makes the point that policy makers neglected the different social structure of the populations to which they applied the programme, and this explains the success in one case and the failure in another case. Social structures can in fact be understood in mechanistic terms too (see e.g. Demeulenaere, 2011). Evidence of mechanisms helps assess the external validity of an RCT (or indeed of any study) because it adds precious knowledge about the similarities between the test and target populations. This point has been forcefully argued for by Steel (2008). ‘Mechanism-based’ external validity inferences are a significant step forward with respect to the Cook and Campbell tradition (Cook and Campbell, 1979) that connects validity merely to the representativeness of the sample and to the possibility of replicating the study.

The problem of inferring from the population to the single case

There is also another sense in which external validity poses a problem. Above, we discussed the inference from one population to another population. Here, the issue concerns the inference from the

population (studied in the RCT) to a *particular patient*. While it is a merit of the evidence-based movement to have fostered protocols for treatment in order to ensure standardisation and comparability, there is no a priori guarantee that an individual patient will be similar enough to the average individual of the RCT and that, consequently, s/he will respond to the treatment in the same way. In such cases, considerations to do with single-case individual responses will be vital to support a claim that the same treatment will work in the single case.

This kind of consideration is particularly vital in treatments for diseases where a variety of distinct causal mechanisms produce clinically similar effects. In the case of breast cancer, tumours may be distinguished by the kinds of receptors they express, and this classification is predicated on the different mechanisms at work in these tumours. Similarly, melanoma classifications now often include consideration of particular genetic mutations (Clarke, 2011). Both these cases are motivated by therapeutic considerations: statistical evidence suggests that differently constituted tumours respond very differently to particular treatments. Thus one needs to know which mechanisms, or features of mechanisms, are instantiated in the particular patient. Again, statistical evidence works better to “make decisions about the care of individual patients” (Sackett et al., 1996) when integrated with evidence of mechanisms.

Integration of evidence

What we really need is to use the totality of evidence available to us. When we must use fallible sources of evidence – and all sources of scientific evidence are fallible – it is better to look for independent converging sources of evidence, as a single good source of evidence will fail significantly often (Wimsatt, 2007). Evidence of correlations obtained from RCTs or observational studies and evidence of mechanisms are independent sources of evidence that are usefully complementary. We have shown that evidence of mechanisms supplements evidence of correlation in designing and assessing RCTs, and in inferring from population to population, and from a population to the single-case. A serious problem with evidence of correlation is the problem of confounding: e.g., when a correlation between variables *A* and *B* may be the result of a common cause of *A* and *B*. Tracing a mechanism from *A* to *B* helps alleviate that worry by offering a direct connection to account for the correlation.

The parallel problem is that evidence of a mechanism does not on its own establish an average causal effect between *A* and *B*. Evidence of one mechanism linking *A* and *B* cannot establish that there aren't other mechanisms linking *A* and *B*, which may balance out, or *mask*, the effect of the known mechanism. But evidence of correlation between *A* and *B* is exactly what is needed to address this masking problem. The best evidence that *A* causes *B* is evidence of a mechanism linking *A* and *B*, where the expected effect size between *A* and *B* is commensurate with the effect size observed in RCTs (if possible) or observational studies seeking a correlation between *A* and *B*. Evidence of mechanisms and evidence of correlation are complementary: each addresses the primary weakness of the other. What we advocate is a pragmatic evidential pluralism, which uses the totality of available evidence.

The problem that we have identified is not that mechanistic evidence is being ignored. Mechanistic evidence is being used to eliminate confounding, to set-up and interpret RCTs, and to extrapolate from one population to another. It is clear from informal discussions with researchers and those charged with approving drugs that mechanistic evidence is being used – often tacitly – alongside statistical evidence in order to establish causal claims. But all this happens *despite* the protocols offered by evidence hierarchies, which urge that, when good statistical evidence is available, it should be considered to the exclusion of other forms of evidence. Evidence hierarchies need revising to ensure that complementary forms of evidence are

treated as complementary, and that evidence of mechanisms, currently treated implicitly, is examined explicitly.

Conflict of interest statement

The authors declare that there are no conflicts of interests.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <http://dx.doi.org/10.1016/j.ypmed.2012.10.020>.

References

- Bailey, R.C., Moses, S., Parker, C.B., et al., 2007. Male circumcision for HIV prevention in young men in Kisumu, Kenya: a randomised controlled trial. *Lancet* 369, 643–656.
- Bem, D.J., 2011. Feeling the future: experimental evidence for anomalous retroactive influences on cognition and affect. *J. Pers. Soc. Psychol.* 100, 407–425.
- Black, R.E., Allen, L.H., Bhutta, Z.A., et al., 2008. Maternal and child undernutrition: global and regional exposures and health consequences. *Lancet* 371, 243–260.
- Canadian Task Force on the Periodic Health Examination, 1979. The periodic health examination. *Can. Med. Assoc. J.* 121, 1193–1254.
- Cartwright, N.D., 2011. Evidence, external validity and explanatory relevance. In: Morgan, G. (Ed.), *The Philosophy of Science Matters: The Philosophy of Peter Achinstein*. Oxford University Press, New York, pp. 15–28.
- Clarke, B., 2011. Causation and melanoma classification. *Theor. Med. Bioeth.* 32, 19–32.
- Cook, T., Campbell, D., 1979. *Quasi-experimentation: Design and Analysis Issues for Field Settings*. Rand MacNally, Chicago.
- Cucherat, M., Haugh, M.C., Gooch, M., Boissel, J.-P., for the HMRAG group, 2000. Evidence of clinical efficacy of homeopathy: a meta-analysis of clinical trials. *Eur. J. Clin. Pharmacol.* 56, 27–33.
- Daemen, J., Wenaweser, P., Tsuchida, K., et al., 2007. Early and late coronary stent thrombosis of sirolimus-eluting and paclitaxel-eluting stents in routine clinical practice: data from a large two-institutional cohort study. *Lancet* 369, 667–678.
- Demeulenaere, P. (Ed.), 2011. *Analytical Sociology and Social Mechanisms*. Cambridge University Press, Cambridge.
- Escudier, B., Pluzanska, A., Koralewski, P., et al., 2007. Bevacizumab plus interferon alfa-2a for treatment of metastatic renal cell carcinoma: a randomised, double-blind phase III trial. *Lancet* 370, 2103–2111.
- Florey, M.E., 1961. *The Clinical Application of Antibiotics. Volume II Streptomycin and Other Antibiotics Active Against Tuberculosis*. Oxford University Press, London.
- Goldenberg, R.L., Culhane, J.F., Iams, J.D., Romero, R., 2008. Epidemiology and causes of preterm birth. *Lancet* 371, 75–84.
- Gray, R.H., Kigozi, G., Serwadda, D., et al., 2007. Male circumcision for HIV prevention in men in Rakai, Uganda: a randomised trial. *Lancet* 369, 657–666.
- Howick, J., 2011. Exposing the vanities—and a qualified defense—of mechanistic reasoning in health care decision making. *Philos. Sci.* 78, 926–940.
- Illari, P.M., Williamson, J., 2012. What is a mechanism? Thinking about mechanisms across the sciences. *Eur. J. Philos. Sci.* 2, 119–135.
- Leibovici, L., 2001. Effects of remote, retroactive intercessory prayer on outcomes in patients with bloodstream infection: randomised controlled trial. *BMJ* 323, 1450–1451.
- Motzer, R.J., Escudier, B., Oudard, S., et al., 2008. Efficacy of everolimus in advanced renal cell carcinoma: a double-blind, randomised, placebo-controlled phase III trial. *Lancet* 372, 449–456.
- MRC, 1948. Streptomycin treatment of pulmonary tuberculosis. *BMJ* 2, 769–782.
- OCEBM Levels of Evidence Working Group, 2011. *The Oxford 2011 Levels of Evidence*. Oxford Centre for Evidence-Based Medicine. (Last access 24-09-2012 <http://www.cebm.net/index.aspx?o=5653>).
- Paavonen, J., Jenkins, D., Bosch, F.X., et al., 2007. Efficacy of a prophylactic adjuvanted bivalent L1 virus-like-particle vaccine against infection with human papillomavirus types 16 and 18 in young women: an interim analysis of a phase III double-blind, randomised controlled trial. *Lancet* 369, 2161–2170.
- Patel, A., ADVANCE Collaborative Group, MacMahon, S., et al., 2007. Effects of a fixed combination of perindopril and indapamide on macrovascular and microvascular outcomes in patients with type 2 diabetes mellitus (the ADVANCE Trial): a randomised controlled trial. *Lancet* 370, 829–840.
- Ritchie, S.J., Wiseman, R., French, C.C., 2012. Failing the future: three unsuccessful attempts to replicate Bem's 'retroactive facilitation of recall' effect. *PLoS One* 7, e33423. <http://dx.doi.org/10.1371/journal.pone.0033423>.
- Sackett, D.L., Rosenberg, W.M., Gray, J.A., Haynes, R.B., Richardson, W.S., 1996. Evidence based medicine: what it is and what it isn't. *BMJ* 312, 71–72.
- Smith, I., Procter, M., Gelber, R.D., et al., 2007. 2-Year follow-up of trastuzumab after adjuvant chemotherapy in HER2-positive breast cancer: a randomised controlled trial. *Lancet* 369, 29–36.
- Steel, D., 2008. *Across the Boundaries. Extrapolation in Biology and Social Science*. Oxford University Press, Oxford.
- Victora, C.G., Habicht, J.P., Bryce, J., 2004. Evidence-based public health: moving beyond randomized trials. *Am. J. Public Health* 94, 400–405.
- Wimsatt, W.C., 2007. *Re-engineering Philosophy for Limited Beings: Piecewise Approximations to Reality*. Harvard University Press, Harvard.