



Form of evidence	Examples of quality-assessment tools
Types of evidence for which quality-assessment tools exist	
Data analytics	<p>ROBINS-I (riskofbias.info) for observational studies, such as those that examine associations between selected factors (including confounding – here the observed relationship between a factor and an outcome, differs from their relationship because of one or more additional factors that are not accounted for)</p> <p>selection of participants included</p> <p>classification of interventions</p> <p>definitions from intended interventions</p> <p>missing data</p> <p>measurement of outcomes</p> <p>selection of the reported results</p>
Evaluation	<p>Risk of Bias (RoB) 2 (riskofbias.info) for randomised-controlled trials, here the risk of confounding is less, but here is a risk of bias from some (albeit fewer) of the same sources as above:</p> <p>randomisation process</p> <p>definitions from the intended interventions</p> <p>missing (outcome) data</p> <p>measurement of outcomes</p> <p>selection of the reported results</p>
Behavioural/implementation research	See other rows for the relevant types of studies or syntheses
Qualitative research	JBI critical appraisal checklist for qualitative research (bit.ly/31SibI), here different considerations come into play, such as: <p>congruence between the research methodology and the research question, data collection methods, data representation and analysis, and research interpretation, as well as between the stated philosophical perspective and the methodology</p> <p>reflectiveness on the part of the researcher, such as statement's locating the researcher critically and theoretically, and addressing the researcher's influence on the research and vice versa</p> <p>representation of study participants' and their voices</p> <p>flow of conclusions from the analysis and interpretation of the data</p>
Evidence synthesis	<p>See above for the relevant types of studies considered in the evidence synthesis</p> <p>A MeAS remen Tool to Assess Systematic Reviews (AMSTAR; amstar.ca) for the quality of the evidence synthesis, here the risk of bias can arise from:</p> <p>identification of all potentially relevant studies through a comprehensive search of both published and grey literature and through language restrictions</p> <p>selection of all studies addressing the research question using explicit criteria about study designs and about participants, interventions/factors, comparisons and outcomes, and at least three reviewers applying the criteria</p> <p>quality appraisal of and data extraction from all included studies</p> <p>synthesis of findings from all included studies</p> <p>Note that there are two versions of AMSTAR: 1) the original version that can be applied across all types of syntheses, albeit to some criteria removed from both the numerical and denominational; 2) a second version of AMSTAR that is more specifically relevant to syntheses of randomised-controlled trials</p> <p>Grading of Recommendations, Assessment, Development and Evaluation (GRADE; bit.ly/3C9pMrx) for the certainty of evidence for the outcomes of an intervention, that:</p> <p>certainty rated down because of risk of bias (the evidence from randomised-controlled trials rating a high certainty and evidence from observational studies rating a low quality and then being adjusted based on RoB2 or ROBINS-I), imprecision (e.g., one or too small studies), inconsistency (e.g., studies showing different findings), indirectness (e.g., surrogate measures used or atypical settings studied), and publication bias (e.g., more common observational studies because of the lack of standard registries or hind-sighted studies because of the commercial incentive to publish positive studies)</p> <p>certainty rated up for large magnitude of effect, dose-response gradient, and when all residual confounding could decrease the magnitude of effect</p> <p>GRADE CERQual (cerqual.org) for the certainty of evidence for the quality of representation of a phenomenon of interest. That:</p> <p>certainty rated down because of concerns about methodological limitations (because problems in the studies were designed or reported were identified using a critical-appraisal tool like the JBI one above), relevance (because the context in which the primary studies were conducted are substantially different from the context of the synthesis question), coherence (because some of the data contradicted the findings or are ambiguous), and adequacy (because the data are not sufficiently rich or only come from a small number of studies or participants)</p> <p>A <i>grey area</i> exists at the <i>ali-a-e</i> stage.</p>



Technology assessment / cost-effectiveness analysis

International Network of Agencies for Health Technology Assessment (INAHTA) checklist (bit.ly/2YJVMVK) for the quality of technology assessments. Out of the 14 questions addressing the approach to assessing the evidence (which prompts similar to AMSTAR) and another question addressing the assessment as compared to other approaches, one question addresses the local meaning (national or sub-national costing data), and consideration of local legal, ethical and social implications.

Drummond checklist of cost-effectiveness analyses (bit.ly/3FbnB8R), and for economic evaluations more generally, the questions about design, data collection, and the analysis and interpretation of results.

Philips checklist for cost-effectiveness analyses that include a decision-analytic modeling component (bit.ly/3FcWBGc). The questions about the structure of the model (e.g., explicit rationale, justified assumptions and appropriate time horizon), the data used (e.g., baseline probabilities from observational studies, treatment effects from randomised-controlled trials, and assessments of responses of uncertainty), namely the structure of the model, the methodological steps followed, the heterogeneity in the population studied, and the parameters used, and the consistency (internal and external). There is also the complementary TRUST tool to assess uncertainties in decision-analytic models (bit.ly/3quFSKp).



Guidelines

AGREE II tool (bit.ly/30qyFAb) for assessing the development, reporting and evaluation (or quality appraisal) of guidelines, which uses 23 items grouped into six domains, each of which is scored independently:

- scope and purpose described
- stakeholder (client/patient and professional) involvement
- rigor of development (the evidence syntheses used as an input, a robust recommendations-development process, and recommendations linked to the supporting evidence)
- clarification of presentation
- applicability
- editorial independence (in relation to funder and panel members' conflicts of interest)

Grading of Recommendations, Assessments, Developments and Evaluations (GRADE; bit.ly/3C9pMrx) for assessing the strength of recommendations, which uses four key considerations:

balance between desirable and undesirable outcomes (trade-offs), taking into account estimates of the magnitude of effects on desirable and undesirable outcomes, and the importance of those outcomes (estimated clinical values and preferences), confidence in the magnitude of estimates of effects of interventions on important outcomes (see GRADE in a previous section), confidence in values and preferences and their variability, resources.

Types of evidence for which quality-assessment tools don't yet exist



Modeling

No ideal accepted tool exists for most types of models, however, there are some general questions that can be asked about models (much like those listed as part of the Philips checklist above), such as:

- structure of the model (e.g., explicit rationale, justified assumptions, and appropriate time horizon)
- data used (e.g., baseline probabilities from observational studies, treatment effects from a range of sources*, and assessments of responses of uncertainty), namely the structure of the model, the methodological steps followed, the heterogeneity in the population studied, and the parameters used
- consistency (internal and external)
- availability of the software or tool that can be assessed by others

*One of the challenges in COVID-19 is that standard designs typically used to capture intervention effects, such as randomised-controlled trials, are ethically or logistically difficult and/or took time to complete, so other standard designs needed to be used and expert opinion needed to be sought (and there are approaches that enable this to be done in a way that is systematic and transparent, such as SHELF, see bit.ly/30nteC4).

Approaches used with certain types of evidence for which quality-assessment tools don't yet exist



Artificial intelligence

No ideal accepted tool exists.