

This is a repository copy of *Reference Case Methods for Expert Elicitation in Healthcare Decision Making*.

White Rose Research Online URL for this paper:
<https://eprints.whiterose.ac.uk/175084/>

Version: Accepted Version

Article:

Bojke, Laura orcid.org/0000-0001-7921-9109, Soares, Marta O orcid.org/0000-0003-1579-8513, Claxton, Karl Philip orcid.org/0000-0003-2002-4694 et al. (7 more authors)
(Accepted: 2021) *Reference Case Methods for Expert Elicitation in Healthcare Decision Making*. *Medical Decision Making*. ISSN 1552-681X (In Press)

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

Medical Decision Making

Reference Case Methods for Expert Elicitation in Healthcare Decision Making

Journal:	<i>Medical Decision Making</i>
Manuscript ID	MDM-20-302.R4
Manuscript Type:	Original Research Article
APPLICATION AREAS:	GROUP DECISION MAKING, FORMULARY DECISION MAKING
DETAILED METHODOLOGY:	Delphi Method (consensus methods) < HEALTH SERVICE RESEARCH, Technology Assessment < HEALTH SERVICE RESEARCH, Provider Decision Making < DECISION AIDS--TOOLS

SCHOLARONE™
Manuscripts

Reference case methods for expert elicitation in healthcare decision making

Authors: Laura Bojke,¹ Marta Soares,¹ Karl Claxton,¹ Abigail Colson,² Aimée Fox,¹ Chris Jackson,³ Dina Jankovic,¹ Alec Morton,² Linda D Sharples,⁴ Andrea Taylor⁵

¹ Centre for Health Economics, University of York, York, UK

² The Department of Management Science, University of Strathclyde, Glasgow, UK

³ MRC Biostatistics Unit, University of Cambridge, Cambridge, UK

⁴ London School of Hygiene and Tropical Medicine, London, UK

⁵ Leeds University Business School, Leeds, UK

Corresponding Author:

Laura Bojke, PhD
Centre for Health Economics
University of York, York, UK
laura.bojke@york.ac.uk

This work was funded by a grant from the Medical Research Council (MRC): MR/N028511/1 “HEE: Developing a reference protocol for expert elicitation in health care decision making”

Background

The evidence used to inform healthcare decision-making (HCDM) is typically uncertain. In these situations, the experience of experts is essential to help decision makers reach a decision. Structured expert elicitation (referred to as elicitation) is a quantitative process to capture experts' beliefs. There is heterogeneity in the existing elicitation methodology used in HCDM and it is not clear if existing guidelines are appropriate for use in this context. In this paper we seek to establish reference case methods for elicitation to inform HCDM.

Methods

We collated the methods available for elicitation using reviews and critique. In addition, we conducted controlled experiments to test the accuracy of alternative methods. We determined the suitability of the methods choices for use in HCDM according to a pre-defined set of principles for elicitation in HCDM, which we have also generated. We determined reference case methods for elicitation in HCDM for Health Technology Assessment (HTA).

Results

In almost all methods choices available for elicitation, we found a lack of empirical evidence supporting recommendations. Despite this it is possible to define reference case methods for HTA. The reference methods include: a focus on gathering experts with substantive knowledge of the quantities being elicited as opposed to those trained in probability and statistics, eliciting quantities that the expert might observe directly, and individual elicitation of beliefs, rather than solely consensus methods. It is likely that there are additional considerations for decision makers in healthcare outside of HTA.

Conclusions

The reference case developed here allows the use of different methods, depending on the decision-making setting. Further applied examples of elicitation methods would be useful. Experimental evidence comparing methods should be generated.

Background

Evidence on health benefits and resource use associated with health interventions may be required to inform healthcare decision-making (HCDM), including assessments of cost-effectiveness.[1] In a model-based analysis, key parameters, such as treatment effects are not known precisely, due to sampling uncertainty. There are often other limitations in the evidence; for example, the licensing of cancer products may be based on evidence of progression-free survival rather than overall survival, or the evidence base may not be well developed (e.g. diagnostics, medical devices, early access to medicines or public health).

It is important that the uncertainty in this evidence is quantified. If not, any analysis using this evidence may give decision makers a misleading view of the consequences associated with their decision.[2] By quantifying uncertainty, it is also possible to assess the potential value of additional evidence.[3] In these situations, the experience of experts is useful and, in some cases, critical to reach a decision. To ensure accountability in the decision, expert judgements should be explicit and their inclusion in HCDM transparent. The process by which beliefs of experts can be quantified according to scientific principles has been called structured expert elicitation (hereafter 'elicitation').[4] When empirical evidence is unsuitable or does not exist, elicitation can provide point and interval estimates describing the state of knowledge for parameters required to make the decision. Where experimental evidence does not exist at all, the expert can utilise their knowledge of the parameter based on their observations, e.g. knowledge gained from clinical practice. Where the experimental evidence is unsuitable, for example in a different population, experts may be required to extrapolate from one population to another.

There is increasing interest in elicitation to inform HCDM, as new technologies are assessed progressively closer to their launch into the market. Elicitation may also be particularly valuable for early stage cost-effectiveness models, or for rare or emerging diseases, for which little or no evidence is available. A review of company submissions appraised by the National Institute for Health and Care Excellence (NICE) found that expert judgement is ubiquitous in company submissions (23/25).[5] In the context of cost-effectiveness analysis, a review of applied studies in decision modelling for cost-effectiveness analysis found heterogeneity in methodology used for elicitation, with little consideration of existing elicitation guidance reported.[8]

Elicitation has also been used widely in disciplines including weather forecasting and engineering.[6] Guidance that exists for elicitation in these contexts suggests several key issues to consider when designing, conducting, and analysing an elicitation exercise, with multiple methodological choices at each stage. The preferred methods are inconsistent across different guidance; examples include the use of group or individual level elicitation methods.[7]

In addition to discipline specific guidance, there are also published generic guidance documents. A number of these have been used in HCDM, the most notable being the Sheffield Elicitation Framework (SHELF)[9] and Cooke's classical method,[10]. Despite their use in HCDM, little is known about the suitability of methods proposed in these generic guidelines. Some of the methods recommended in generic guidance may not be suitable in HCDM, for example the elicitation of complex quantities or the use of more complex methods. The reasons for this include resource and time constraints in HCDM, the types of experts typically consulted; usually recruited for their subject knowledge rather than quantitative background, and the wide range of parameters required for elicitation. [11]

In this paper we describe the development of reference methods for expert elicitation to inform HCDM. Details are reported in full elsewhere.[12] The intention is for these reference methods to be used by a range of decision makers to generate their own guidance for expert elicitation, for

1
2
3 example across the globe and/or across different areas of HCDM. Here we describe the reviews
4 undertaken to compile methods available for expert elicitation, the approach used to critique the
5 different methods for expert elicitation and determine their suitability for use in HCDM, and finally
6 the set of reference methods that were produced. Given the infancy of expert elicitation in HCDM
7 and the lack of evidence to support many of the methods choices, we define these reference
8 methods for one aspect of HCDM, Health Technology Assessment (HTA). Thereafter we highlight the
9 complexities and challenges for HCDM outside of this setting.
10
11

12 **Methods**

13 We conducted systematic and non-systematic reviews of evidence to compile available methods
14 (described in **Reviews of evidence to compile methods for elicitation**). To generate reference
15 methods for HCDM we then developed resources to critique the identified methods for elicitation
16 (described in **Critiques of methods choices for elicitation**). Details are reported in full elsewhere.[12]
17 We summarise the evidence sources and critique approaches in the sections below (**Determining**
18 **appropriate methodological choices for elicitation in HCDM**). Figure 1 presents the broad structure
19 of the evidence sources and critique methods.
20
21

22 **Reviews of evidence to compile methods for elicitation**

23 We utilised a range of evidence sources. These sources are summarised below.
24
25

26 *Review of published guidelines:* We have undertaken a review of guidelines for elicitation published
27 in either the peer-reviewed or grey literature. The elicitation guidelines were systematically
28 reviewed according to the search strategy and inclusion criteria presented in the appendix (also
29 detailed elsewhere [12]). This review (not restricted to HCDM), included guidelines concerning
30 probabilistic judgements that offered guidance on more than one stage of the elicitation
31 process.[12] Information was extracted from these guidelines to create an overview of the
32 sequential stages of the elicitation process (design, conduct and analysis), the elements within each
33 of these stages and the choices involved in each of these elements (see the appendix for the
34 extraction template). For example, training and preparation is an element of elicitation, which
35 requires choices about whether and how to pilot the elicitation and what to cover when training
36 experts.
37
38

39 We determined where current advice conflicts or agrees across guidelines. Where the guidelines
40 agreed, we assumed this methodological choice represented best practice and accepted it as the
41 reference case method.
42
43

44 *Targeted searches:* A priori, we were aware of many elements of elicitation that were not discussed
45 in any depth in the existing guidelines. It was therefore not clear what methods choices were
46 available for these elements. To augment the existing published guidelines, we conducted semi-
47 structured searches to identify the full set of choices for these elements.[12] The searches also
48 aimed to identify any agreed 'best practice' for elicitation in these elements. Further details of the
49 methods of the targeted searches are reported elsewhere.[12] Specifically we conducted targeted
50 searches for five methods choices. The areas for targeted searches were chosen following
51 consultation with our project advisory group, and are as follows:
52

- 53 1) the selection of experts
- 54 2) the level of elicitation (individual or group)
- 55 3) the fitting of parametric models to elicited judgements and subsequent aggregation across
56 multiple experts
- 57 4) the assessment of the accuracy of expert judgements
- 58 5) the identification of cognitive heuristics and biases and methods to minimise the impact of these
59 on the elicitation.
60

2.2 Critique of methods choices for elicitation

In the anticipated absence of definitive statements of agreement from the review of guidelines and targeted searches, it was necessary to critique the methods choices identified in **Reviews of evidence to compile methods for elicitation**. We did this according to the principles for successful elicitation that we developed as part of this project (**Principles for elicitation in HCDM**). In addition, we augmented the critique using evidence on the choice of methods that might be suitable in HCDM from the applied studies review (**Review of applied studies**)[11], the constraints in using elicitation in HCDM that we identified (**Constraints in using elicitation in HCDM**) and conclusions drawn from experiments conducted as part of this project (**De novo experimental evidence**). Full details are presented elsewhere.[12]

Review of applied studies: We have previously published a review of cost-effectiveness studies that include elicitation [11]. The review considered the methods used and the specific challenges in conducting elicitation in this context. We identified twenty-one applied studies. Many authors expressed methodological uncertainty in justifying their choices. From the review, several aspects of the context area (HCDM) emerged as potentially important in determining methodological choices in elicitation. We used the findings from this review of applied studies to generate the core principles for elicitation in HCDM (**Principles for elicitation in HCDM**), and also to critique the methods choices from the review of guidelines.[11, 12]

Constraints in using elicitation in HCDM: In considering how reference methods for elicitation in HCDM might be used in practice, it is important to understand how different decision-making settings may influence the requirements for, and practicalities of, elicitation. We considered the potential practical constraints of using elicitation in HCDM at various levels of decision-making when generating the principles and assessing the applicability of methods identified in the review of published guidelines. A formal review of the challenges and constraints faced by different HCDM's was not possible. Instead this source of evidence drew on the observations and experiences of the authors and an advisory group convened as part of the project (see acknowledgements for details of this group). See [12] for further details.

De novo experimental evidence:

As part of this work, we generated evidence from randomised simulation experiments to compare method choices for elicitation.[12] Randomisation concerned the level of precision specified in the scenarios presented to participants for experiment 1, the distributions of subgroups for extrapolation for experiment 2 and the degree of discordance between individual and group summaries in experiment 3. Randomisation was undertaken to explore multiple scenarios within each of the experiments, whilst not over burdening participants. It also helped to standardise other aspects of the participants knowledge, such as their level of training in probability and statistics.

Specifically the experiments concerned: 1) different methods to encode experts judgements, the Variable Interval Method and Fixed Interval Method, 2) requiring experts to extrapolate from their individual knowledge to populations with different prevalence of a successful outcome, and 3) the use of Delphi type processes to understand how experts revise their estimates in the light of group summaries.[12] The experiments were conducted using the Shiny package for R.

To conduct these experiments, we used a simulated (virtual) learning process to standardise participants knowledge. This allowed us to compare the elicited probabilities directly to the distribution implied by the observed dataset, therefore giving a measurement of accuracy. Two main metrics were used: 1) bias in location (difference in the mean of the elicited distribution and the posterior distribution implied by the data) and 2) bias in uncertainty (the ratio of the standard

1
2
3 deviation of the elicited distribution to the standard deviation of the posterior distribution implied
4 by the data).

5
6 In the first experiment, each participant was shown observations from a stochastic simulation
7 model. The context was an abstract generic medical problem. The participants were asked to choose
8 between treatments with differing levels of effectiveness. Experiments 2 and 3 followed on from the
9 context specified in experiment 1 and 3. We report further details of the methods used in these
10 experiments elsewhere.[12] We use the results from these experiments to provide additional
11 information on the suitability of methods arising in these three areas, from the review of guidelines.
12
13

14 *Principles for elicitation in HCDM:* These principles are primarily informed by the findings of the
15 published review of the use and challenges in applying elicitation in cost-effectiveness modelling
16 **(Review of applied studies)**[11] and the review that identified constraints in using elicitation in
17 HCDM more generally **(Constraints in using elicitation in HCDM)**. [12] We also reflect the
18 requirements for elicitation reported in the guidance by Cooke (1991)[10]. These requirements
19 represent 'good practice' in elicitation generally and are widely referred to in the elicitation
20 literature.
21
22

23 We first drafted the principles and then amended these following a meeting with our advisory
24 group, convened to guide the project. The advisory group consisted of elicitation methodologists
25 and users (see acknowledgements for the list of advisory group members). We presented the
26 redrafted principles at the workshop described in **Determining appropriate methodological choices
27 for elicitation in HCDM.**
28
29

30 **Determining appropriate methodological choices for elicitation in HCDM**

31 We have made recommendations for each element within the stages for elicitation by assessing
32 which choices are supported by the principles or the evidence, and which could be left flexible
33 according to the specific elicitation context. We convened a stakeholder workshop, where we
34 presented our draft reference methods for elicitation in HCDM, for the purposes of gaining feedback
35 and establishing validity. We identified the relevant stakeholders as: HTA decision makers,
36 methodologists, industry representatives and commissioners. To gather stakeholders we reached
37 out to UK decision makers, including those from NICE, NHS England and Public Health England,
38 authors of the elicitation papers identified in our applied studies review[11] and key contacts in
39 industry and consultancy. Approximately 30 stakeholders attended the event.
40
41

42 We gathered opinion through presentation and discussion, followed by communication with specific
43 individuals who wished to speak about the topic outside of the meeting. Following feedback from
44 this workshop, we generated a set of final reference case methods. The workshop also sought to
45 identify challenges in using elicitation in different settings, for example where evidence is immature,
46 or where decisions concern orphan drugs. We documented these challenges along with the
47 examples of areas in which they arose.
48
49

50 **Results**

51 **Evidence to inform the set of choices**

52 *Review of published guidelines*
53
54
55
56
57
58
59
60

1
2
3 We identified sixteen unique guidelines (see the appendix for the search results and the full list of
4 included guidelines. See [12] for further details). Five are generic[9, 13-16] and eleven are domain
5 specific.[17-27] The guidelines include the widely cited European Food Safety Agency (EFSA)
6 guideline[25], Cooke's Classical Model[13] and SHELF.[9] Whilst some of these guidelines have been
7 used in HCDM, for example SHELF[9] and the IDEA protocol,[15] none were developed specifically
8 for this context, and none discuss their applicability to HCDM.
9

10
11 Details of the elements and methodological choices contained in existing guidelines are presented in
12 the appendix. In addition, we present, for each methodological choice, the level of agreement
13 between existing guidelines (see the appendix). There are relatively few methods choices for which
14 the existing guidelines entirely agree. Areas of agreement are: (1) the need to decompose
15 (breakdown) variables into several smaller, more observable quantities, (2) the number of experts
16 should be between 5 and 10, (3) the roles of experts within the elicitation task should be made
17 explicit, (4) there should be piloting of the task, (5) experts should provide rationales for their
18 judgements and (6) aggregation should be undertaken post elicitation. There are many
19 methodological choices for which guidelines have only partial agreement on the appropriate choice,
20 or else no agreement at all.
21
22

23 *Targeted searches*

24 In the five areas subjected to targeted searches, there is very little empirical evidence to support or
25 discount any specific choices, and none of the evidence that does exist focuses on HCDM.[12] Any
26 conclusions offered on these elements are generally anecdotal rather than empirically based. For
27 example, regarding the minimisation of bias, there is a suggestion that experts should not be asked
28 to express confidence intervals in a single stage process, as doing so results in participants focussing
29 on a narrow set of salient possibilities. Instead, lower bounds, upper bounds and median values
30 should be elicited separately.[9, 15, 19] Full details of the targeted searches results are reported
31 elsewhere.[12]
32
33

34 **Resources to critique methods choices**

35 The review of applied studies (**Review of applied studies**), the constraints that may have
36 implications for elicitation in HCDM (**Constraints in using elicitation in HCDM**) and the evidence
37 generated from the experiments (**De novo experimental evidence**) are reported in detail elsewhere,
38 and results are therefore not repeated here. [11, 12] Instead we describe the principles that were
39 generated and refer to evidence from Critique of methods choices for elicitation in the critique of
40 methodological choices section below (**Critique of methodological choices for elicitation in HCDM**).
41
42

43 We developed nine principles for judging the suitability of choices available for elicitation. These are
44 summarised below. Workshop participants agreed unanimously that these represented a complete
45 set of requirements for elicitation conducted in HCDM, with stakeholders only suggesting minor
46 changes to the wording.
47
48

49 *Principle 1. Ensure transparency in elicitation*

50 Systematic and transparent reporting of elicitation helps to improve the validity of the resulting
51 expert judgements, allows the elicitation to be peer-assessed, and supports others who use the
52 judgements in their own analysis. If there is insufficient space to describe the elicitation process in
53 the primary study report, separate details of the elicitation, ideally comprising an elicitation protocol
54 and results, should be fully documented.
55
56
57
58
59
60

1
2
3 *Principle 2. The elicitation must provide useful information for the decision problem*

4 The elicitation must be fit-for-purpose, in that it must provide information that is relevant to the
5 decision problem. If a decision model is employed by the analyst to synthesise evidence to
6 determine cost-effectiveness, then the quantities being elicited should be consistent with the
7 parameters and structure of the model. For example, suppose we believe that two model
8 parameters are likely to be correlated, such that a belief that one parameter is high implies belief
9 that the other one is high. In these circumstances an elicitation designed to inform these parameters
10 should give information about their correlation, e.g. by eliciting the second quantity conditionally on
11 the first. Multiple quantities must also be mathematically consistent, for example, probabilities of
12 mutually exclusive events should sum to one.
13
14

15
16 *Principle 3. Elicitation should aim for consistency, but respect the constraints of the decision-making*
17 *context*

18 There are different potential users of elicitation, from local level to national or international decision
19 makers, including reimbursement agencies and research funders. These different decision-making
20 entities have quite different capacities to conduct elicitation and incorporate it into their decision-
21 making processes. It is important that a degree of flexibility is retained in the reference case for
22 elicitation in HCDM, but sensitivity of results to the choices made should be explored.
23
24

25 *Principle 4. Elicitation should reflect uncertainty at the individual expert level*

26 Judgements elicited from experts need to reflect the imperfect knowledge they have. In elicitation,
27 experts are often required to provide both a point estimate of the quantity(s) of interest and an
28 assessment of their uncertainty in that estimate. An important concern is that, when reflecting on
29 their own experiences, experts may mistakenly report the extent of variability (e.g. between disease
30 outcomes for individuals) rather than uncertainty in knowledge (e.g. about the expected incidence
31 rate of the outcome).
32
33

34 *Principle 5. Elicitation should recognise and act on biases*

35 There are many biases and heuristics (cognitive shortcuts that individuals use when asked for
36 complex judgements) that apply to elicitation, including overconfidence/under confidence, over-
37 extremity (tendency to use the extremes when responding), discrimination (including prejudice or
38 stereotyping), or susceptibility to base rate neglect (a disregard of fundamental statistical reality). An
39 elicitation task should be designed and conducted using techniques that mitigate against heuristics
40 and other sources of bias, and appropriate training should be given to experts.[14]
41
42

43 *Principle 6. An elicitation task should be suitable for experts who possess substantive skills and who*
44 *are less likely be trained in probability and statistics*

45 Substantive experts in HCDM are often health professionals, who are unlikely to have had extensive
46 experience of elicitation, and unlikely to have developed the necessary normative skills, e.g. in
47 probability and statistics. Methods of elicitation employed in other areas may not be directly
48 suitable in HCDM unless there is additional training before use.
49
50

51 *Principle 7. The elicitation task should recognise where adaptive skills are required*

52 In some instances, adaptive skills may be relevant for elicitation in HCDM. For example, in early cost-
53 effectiveness modelling or early stage trial design, experts may not be familiar with the target
54 quantity for elicitation but are substantive experts in one or more related quantities (for example,
55 the quantity in a similar population to the target population). In this situation, knowledge of the
56 related quantity may need to be adapted.
57
58
59
60

1
2
3 *Principle 8. Elicitation should recognise, and act on, between-expert variation*

4 In the context of HCDM, between-expert variation is common. There may be genuine heterogeneity
5 in the populations that experts draw upon to formulate their judgements. In this case, it is desirable
6 to reflect this heterogeneity in the pooled distribution, whether through group consensus or
7 mathematical aggregation methods. It is also important to understand why between-expert
8 heterogeneity is present.
9

10
11 *Principle 9. Elicitation should promote high performance*

12 In HCDM, experts may be motivated to undertake the task to the best of their abilities because of
13 their interest in the topic area and for altruistic reasons. However, not all experts within an
14 elicitation may possess the same subject knowledge and there may be differences in normative (e.g.
15 probability and statistics) expertise. Where possible, an elicitation task should account for differing
16 levels of normative expertise and encourage experts with substantive knowledge to perform equally,
17 for example in providing unbiased estimates. As well as promoting high performance, an elicitation
18 may want to explore differences in expert performance.
19

20
21 **Critique of methodological choices for elicitation in HCDM**

22 The critique determined the suitability of elicitation methods according to their adherence to the
23 principles of elicitation for HCDM reported here (see the appendix for full details of which
24 principles are applied to which methods choices), the results from the experiments and the
25 constraints.[12]
26

27
28 *Selecting quantities (preparation and design)*

29 A key requirement is that the elicited information should be fit-for-purpose and describe an expert's
30 uncertainty regarding the quantity of interest. Experts in HCDM are often recruited due to their
31 subject expertise and may be less likely to have statistical expertise. To aid completion by experts in
32 HCDM, existing guidelines are consistent about the need to breakdown variables into simpler
33 quantities to elicit.
34

35
36 Despite the lack of empirical evidence to support this assertion, we believe that questions should be
37 posed in a manner consistent with how experts express their knowledge. As a result, elicitation tasks
38 should specify observable quantities, such as probabilities (expressed as proportions or frequencies)
39 and more complex quantities such as odds ratios or variances should be avoided. The use of
40 observable quantities may aid experts when they are required to extrapolate outside of their
41 knowledge base. The experiments we conducted concluded that such extrapolation is unlikely to
42 produce more biased judgements or more inaccurate expressions of an expert's uncertainty.
43

44
45 In some circumstances the quantities elicited may have a degree of dependency. In HCDM the aim
46 should be to ask about independent quantities where possible.[9, 13, 15, 17, 21-27] If this is not
47 possible, dependent quantities can be re-expressed in terms of independent quantities or
48 conditional quantities, or dependence methods can be used.[15]
49

50
51 *Methods to encode judgements (preparation and design)*

52 In general, existing guidelines suggest that both the Fixed Interval Method and the Variable Interval
53 Method can be used to encode judgements.[9, 13, 16, 17, 20, 23, 25-27] Because experts may be
54 recruited primarily because of their substantive skills, the suitability of alternative methods must
55 recognise differences in their normative (e.g. probability and statistics) skills. Evidence from our
56 experiments suggested that the Fixed Interval Method and the Variable Interval Method are equally
57 appropriate for HCDM in terms of providing accurate representations of an expert's uncertainty,
58 although there is some preference for the Fixed Interval Method, delivered using a "chips and bins"
59
60

1
2
3 approach.[12] Decision makers may therefore choose either, but should apply them consistently in
4 their setting.
5

6 *Selecting experts*

7 The existing guidelines and targeted searches suggest features to consider when selecting experts.
8 These include normative expertise, substantive expertise, and willingness to participate. The
9 constraints of conducting elicitation in HCDM may dictate that the selection process focuses on only
10 one or two key characteristics. It is worth noting that there may be a limited number of healthcare
11 professionals with the relevant substantive expertise and therefore more opportunistic methods for
12 recruitment may be required, such as peer nomination. In some instances, adaptive skills may be
13 required for an elicitation, particularly in the case of new and emerging technologies. It is not clear
14 what metrics can be used to determine an expert's level of adaptive skills.
15
16

17
18 Identifying an unbiased expert poses a challenge and indeed, an entirely unbiased expert may not
19 exist. The targeted searches showed that the elicitation should attempt to recruit experts who are
20 free from motivational biases by collecting disclosure of personal and financial interests and conflicts
21 of interest.[9, 14, 19, 22, 27] Additionally, efforts should be made to ensure that the sample of
22 experts contains a range of viewpoints, with the intention of "balancing out"[13, 15-17, 19, 20, 22].
23 This may dilute the effect of motivational biases.
24

25
26 Between-expert variation may exist, and methods used to select experts must attempt to capture
27 the range of plausible beliefs. Identification of experts through recommendations by peers, either
28 formally or informally, may generate a pool of experts that are all similar. Instead it may be
29 preferable to identify experts through research outputs, known experience or by using a profile
30 matrix. The elicitation can also seek diversity in background and a balance of different viewpoints.
31 Recruiting a larger number of experts may help to fulfil these criteria (5-10 experts are suggested by
32 the existing literature identified in the targeted searches[12]).
33
34

35 *Piloting and training*

36 All existing guidelines agree that an elicitation should be piloted on a smaller set of experts prior to
37 the actual task, with subsequent revision based on feedback and follow up of any issues that arise.
38 For example, priors that are incoherent may indicate the need to re-specify the quantities elicited or
39 the questions asked.
40

41 Training of experts is essential and should focus on enabling non-normative but substantive experts
42 to express their uncertain beliefs at an individual level. Training also plays a key role in minimising
43 biases. Although evidence in the context of HCDM is weak, there are some suggestions from the
44 literature that training can reduce the effect of anchoring and adjustment bias, confirmation bias
45 and overconfidence.
46
47

48 The training delivered to experts will be guided, in part, by the specific task, and include for example
49 the description of quantities, the description of the performance measurement and how to manage
50 dependence. The core elements of training are a description of what is required from the experts, an
51 outline of the elicitation process, an outline of the questions that will be asked, and example and
52 practice questions.[12]
53

54 *Level and conduct of elicitation*

55 Existing guidelines are inconsistent regarding whether elicitation should be individual, or group
56 based.[12] Group discussion can help experts with less clinical knowledge or probability and
57 statistics training. However, interaction between experts can also introduce biases, and the act of
58 striving for consensus can potentially eliminate important between expert variation. The constraints
59
60

1
2
3 apparent in HCDM, such as limited access to experienced experts and short timescales for decision-
4 making, may also discourage the use of consensus methods.[12] Additionally, there is no evidence
5 from our experiments that those who revise their judgements following group feedback have
6 different accuracy to those that who did not revise their judgements, which casts some doubts on
7 the benefits of the Delphi iteration process. For these reasons, we believe it is preferable to elicit
8 from experts individually.
9

10
11 When using individual elicitation there should still be possibilities for interaction between experts.
12 This should follow on from the individual elicitation where practically feasible and useful. For
13 consensus methods, again the elicitation should first conduct individual elicitation followed by the
14 group consensus stage. Feedback should follow the elicitation task, with graphical feedback
15 considered for experts unfamiliar with probability and statistics.
16

17
18 Many of the existing guidelines agree that face-to-face administration is preferred.[9, 13, 17-22, 24] It
19 is thought to promote good performance and maximise engagement with experts. Face-to-face
20 elicitation is necessary for some consensus methods; however, it is not necessary for aggregating
21 judgments mathematically. The constraints in HCDM are the biggest factors in driving the method
22 choice. If the task requires many experts, face to face elicitation may be prohibitively time consuming
23 and resource intensive.
24

25 *Aggregation, analysis and post-elicitation*

26 The existing guidelines agree that, following elicitation of judgements, there should be an
27 aggregation of the elicited information across experts. In the context of HCDM, however,
28 aggregation should not simply focus on reducing variability between experts; instead efforts should
29 also be made by the elicitation facilitator to understand the reasons for any variability. To generate
30 an aggregate summary, e.g. for use in a probabilistic decision model, it is necessary to fit a
31 probability distribution. For the purposes of using the elicitation results in further analysis, a smooth
32 fitted distribution is preferred to an empirical summary (without fitting). The choice of distribution
33 will depend on the quantity elicited. Parametric distributions (e.g. Normal, Beta or Gamma) may be
34 appropriate. The best-fitting distribution should be determined using standard statistical methods
35 (e.g. Ordinary least squares, generalized method of moments or Maximum Likelihood). Simple
36 mathematical rules for aggregation, like linear opinion pool with equal weights, are the most
37 commonly applied in HCDM and are straightforward to implement.
38
39

40
41 Documentation of the elicitation design, conduct and analysis is key to understanding the choices
42 made and the rationale for these. It is also important in assessing the validity of the distributions
43 elicited. Details should be reported in the body of a report if possible and as a separate appendix if
44 not. The documentation should include the justifications given by the experts, for their judgements.
45
46

47 *Managing biases*

48 There is very little in the existing guidelines on methods to minimise bias. The targeted search
49 conducted to identify methods to minimise cognitive heuristics and biases,[12] suggests that efforts
50 should be made to identify the likely biases given the type of experts who have been recruited.
51 Relevant strategies to minimise these should then be employed.[12] To mitigate biases, experts can
52 be told as part of their training about the likely sources of bias, and asked to be aware of these when
53 responding to questions. In addition, questions can be framed in a way to minimize bias and
54 ambiguity. This could include asking experts to first specify their plausible upper and lower bounds
55 and giving experts the opportunity to revise the information they provide.
56
57
58
59
60

Validation

Commonly discussed elements of validation include verifying that the elicitation captures what experts truly believe, and that the expressed probabilities reflect reality. Above all, validation should focus on the extent to which the elicited beliefs are fit for purpose for the intended task. This could be assessed by coherence and consistency with the intended HCDM it is to inform (i.e. an assessment of face validity). It is also important to understand how experts formulate their beliefs and why they present heterogeneous beliefs. An external review of the elicited priors, by experts not involved in the elicitation task, should be undertaken to assess validity.

Generating reference case methods for elicitation

The sparse evidence supporting the methodological options in elicitation means that, for many elements, there remains uncertainty about the most appropriate choices, and further research is necessary. The previous section lays out considerations that are required to generate reference case methods for elicitation in HCDM settings. These are also reported in detail elsewhere.[12] This critique helps to highlight the trade-offs required when developing context-specific methods, where we need to take into account not only accuracy, but context specific features, restrictions, and constraints.

Elicitation can inform HCDM in diverse settings, ranging from local-level prioritisation to strategic planning for emerging threats. It has, perhaps, been most frequently applied in national level reimbursement, price negotiation and clinical guideline development,[7] collectively referred to as HTA. We have developed an exemplar set of reference methods for elicitation in the HTA context (see *Table 1*).

In summary, our reference case methods state that, in HTA, the elicitation should focus on gathering substantive expertise or experience. Elicitation skills can be developed during the training, which should focus primarily on avoiding bias and expressing uncertainty. In recruiting experts, conflicts of interest should be minimised and if necessary recorded. Experts external to the elicitation task should be included, i.e. not those involved in developing the task. Beliefs should be elicited face-to-face and from experts individually and then pooled. Between-expert heterogeneity should be explored explicitly. Simple observable quantities should be elicited where possible, with efforts made to capture dependence between quantities in a way that can be elicited reliably. Either the Variable Interval Method or the Fixed Interval Method can be used, with the choice depending on which best suits the type of expert and the elicitation task.

Whilst these reference methods are intended to reflect emerging 'best practice' in HTA, given the infancy of elicitation applied to HCDM, it is important to allow a degree of flexibility in the reference methods recommended here. A decision maker adopting this protocol could choose to specify methods for the reference case to ensure greater consistency across appraisals. In cases where non-reference case methods are employed, choices should be justified, and sensitivity analyses undertaken.

Elicitation may also be useful for decision makers outside of HTA, for example at a local level, or in the context of the appraisal of early technologies that have yet to progress through the regulatory process. In addition, there may be additional challenges in some HTA contexts, for example in the assessment of genomic treatments, or treatments for rare diseases. In such settings, a potential reference case should consider the additional issues summarised in the 3rd column of *Table 1*.

Conclusions

Elicitation can be a valuable method for HCDM, particularly to inform reimbursement decisions that are supported by model based economic evaluation. Elicitation provides the additional information needed to reach a decision when empirical evidence is lacking.

This paper describes work to generate reference case methods for elicitation for HCDM. We believe that the results will be useful for analysts and decision makers in HCDM. Elicitation conducted in this context to date has not used a common set of methods, and above all, has not consistently considered the implications of the methods choices made when designing and conducting an elicitation. To improve the accountability of HCDM, the procedure used to derive expert judgements should be transparent and documented.

The reference case methods presented here serve as a benchmark for good practice and reporting. Whilst consistency in the methods applied is desirable to ensure consistent evaluations, the lack of evidence on the performance of different methodological choices means we could not be prescriptive. This reference case is therefore, by virtue of the evidence used to support it, flexible in many choices. This may be a useful characteristic, as it is possible to apply the reference case across different settings within HTA. Deviations from the suggested methods should be justified and limitations discussed in the elicitation documentation. It may be useful to report the methods used in the applied elicitation using the reference case methods as a benchmark.

Here we illustrate the development of a reference case specific to the HTA setting. Different HCDM contexts have different constraints and requirements. Outside of HTA, there are key methodological choices which may involve additional or different considerations, for example as part of the commissioning process at a local level, or for early technologies that have yet to progress through the regulatory process. Moreover, in some circumstances, it may not be possible to conduct face-to-face elicitation. Group discussion may be needed to generate a distribution, where there is no practical experience of the quantity of interest.

The major limitation of this work lies in the evidence available from the wider literature, on which to base methods choices and determine their appropriateness. The lack of an agreed definition for accuracy of elicitation also limits the choice of 'best' methods. In many circumstances, expert beliefs are unobservable to the analyst, so that determining how well methods perform in enabling experts to express their beliefs is a complex task.

There are important areas warranting further research. These include strategies to recruit experts, methods for training experts to minimise bias, and methods for eliciting dependent quantities from non-normative experts. Application of the reference case in further studies, including in settings with a range of constraints, will generate valuable evidence regarding its applicability and value.

References

1. Bryan S, Williams I, and McIver S, *Seeing the NICE side of cost-effectiveness analysis: a qualitative investigation of the use of CEA in NICE technology appraisals*. Health Economics 2007. 16: p. 179-93.
2. Griffin SC, et al., *Dangerous omissions: the consequences of ignoring decision uncertainty*. Health Economics, 2011. 20: p. 212-24.
3. Claxton K, et al., *Probabilistic sensitivity analysis for NICE technology assessment: not an optional extra*. Health Economics, 2005. 14(4): p. 339-347.
4. O'Hagan A, et al., *Uncertain judgements: eliciting experts' probabilities*. 2006.
5. van Hest N, et al. *Trust the Experts? Acceptance of Expert Elicitation in the National Institute for Health and Care Excellence (NICE) Single Technology Appraisal (STA) Process*. in *ISPOR 22nd Annual European Congress*. 2019.
6. Babuscia, A and Cheung KM, *An approach to perform expert elicitation for engineering design risk analysis: methodology and experimental results*. Journal of the Royal Statistical Society Series a-Statistics in Society, 2014. 177(2): p. 475-497.
7. Ayyub BM., *Elicitation of expert opinions for uncertainty and risks*. 2001, Boca Raton, Fla.: CRC Press. 302 p.
8. Grigore B, et al., *Methods to Elicit Probability Distributions from Experts: A Systematic Review of Reported Practice in Health Technology Assessment*. PharmacoEconomics, 2013. 31: p. 991–1003.
9. Gosling JP, *SHELF: The Sheffield Elicitation Framework*. 2017: Springer.
10. Cooke, R.M., *Experts in Uncertainty: Opinion and Subjective Probability in Science*. Oxford University Press, 1991.
11. Soares MO, et al., *Experiences of structured elicitation for model based cost-effectiveness analyses*. Value in Health, 2018. 21(6): p. 715-723.
12. Bojke L, et al., *Developing a reference case for expert elicitation in healthcare decision making*, Health Technology Assessment, 2021 forthcoming. London, UK. Pre-print available at: <http://eprints.whiterose.ac.uk/151174/>
13. Cooke, R.M. and L.H.J. Goossens, *Procedures guide for structured expert judgement in accident consequence modelling*. Radiation Protection Dosimetry, 2000. 90(3): p. 303-309.
14. Garthwaite, PH., Kadane JB, and O'Hagan A, *Statistical Methods for Eliciting Probability Distributions*. Journal of the American Statistical Association, 2005. 100(470): p. 680-701.
15. Hemming V, et al., *A practical guide to structured expert elicitation using the IDEA protocol*. Methods in Ecology and Evolution, 2017. 12.
16. Meyer MA, Booker JM., *Eliciting and Analyzing Expert Judgment: A Practical Guide*. 2001, Society for Industrial and Applied Mathematics: Philadelphia, PA.
17. Keeney R. and von Winterfeldt D, *Eliciting Probabilities from Experts in Complex Technical Problems*. IEEE Transactions On Engineering Management, 1991. 38(3).
18. Kaplan S, *Expert information' versus 'expert opinion.'* Another approach to the problem of eliciting/combining/using expert knowledge in PRA. Reliability Engineering & System Safety, 1992. 35(1): p. 61-72.
19. Kotra, J., et al., *Branch Technical Position on the Use of Expert Elicitation in the High-Level Radioactive Waste Program*. Division of Waste Management, Office of Nuclear Material Safety and Safeguards, U.S. Nuclear Regulatory Commission, 1996.
20. Budnitz RJA, et al., *Recommendations for Probabilistic Seismic Hazard Analysis: Guidance on the Uncertainty and Use of Experts*. NUREG/CR-6372. Two volumes. 1997, Nuclear Regulatory Commission: Washington, D.C, U.S.
21. Walls L and Quigley J, *Building prior distributions to support Bayesian reliability growth modelling using expert judgement*. Reliability Engineering & System Safety, 2001. 74(2): p. 117-128.

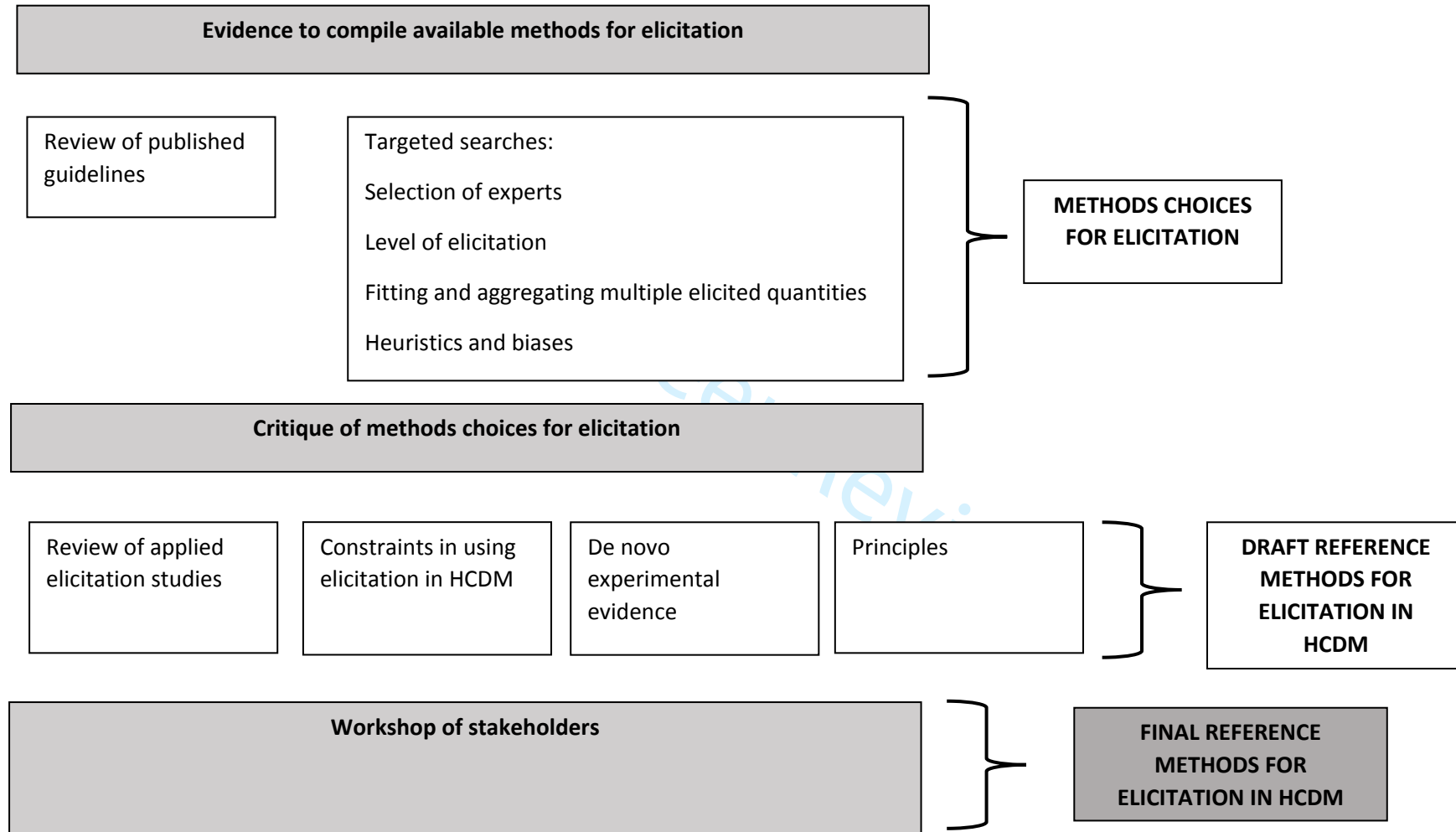
- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
22. U.S. Environmental Protection Agency, *Expert Elicitation Task Force White Paper*. Prepared for the U.S. Environmental Protection Agency by Members of the Expert Elicitation Task Force, Editor. 2009: Washington DC.
23. Choy, S., R. O'Leary, and K. Mengersen, *Elicitation by design in ecology: using expert opinion to inform priors for Bayesian statistical models*. *Ecology*, 2009. 90(1).
24. Knol, A.B., et al., *The use of expert elicitation in environmental health impact assessment: a seven step procedure*. *Environmental Health*, 2010. 9.
25. European Food Safety, A., *Guidance on Expert Knowledge Elicitation in Food and Feed Safety Risk Assessment*. *EFSA Journal*, 2014. 12(6): p. 3734-n/a.
26. Ashcroft, M., et al., *Expert Judgement*. Institute and Faculty of Actuaries' Solvency & Capital Management Working Party, 2015.
27. Tredger ERW, et al., *Bias, guess and expert judgement in actuarial work*. *British Actuarial Journal*, 2016. 21(3): p. 545-578.

Table 1: A reference case for HTA

Element	Reference methods suggested	Additional considerations outside of HTA
Selecting quantities	<ol style="list-style-type: none"> 1. Simple observable quantities should be elicited where possible; ratios or complex parameters such as regression coefficients should not be elicited directly. 2. Dependence between variables should be captured in elicitation. Expressing dependent variables in terms of independent variables is preferable when experts do not have strong normative skills. 3. Wording should be clear and quantities should be decomposed where this means a better fit with experts intuition. 	-
Methods to encode judgements	Both Variable Interval Methods or Fixed Interval Methods can be used. Decision makers should aim for consistency across applications.	Fixed Interval Methods may be more appropriate for less experts less familiar with elicitation or where face-to-face training is impossible.
Selecting experts	<ol style="list-style-type: none"> 1. Recruitment will be driven by the context, however the elicitation should pursue diversity, representing the full range of valid expert beliefs. Experts should be willing to participate. 2. Focus on gathering substantive expertise or experience. Normative skills (for example, in probability and statistics) can be developed during the training session as part of the elicitation. 3. Minimize and record conflicts of interest among the experts. Include experts external to the elicitation task, i.e. not those involved in developing the task. 4. At least 5 experts should be included in the elicitation. 	<ol style="list-style-type: none"> 1. Researchers may have limited access to sufficient experts, for example in rare diseases, therefore expert recruitment may be more challenging and rely on peer nomination. 2. Adaptive skills may be required for new technologies since indirect evidence may outweigh directly relevant evidence (e.g. childhood diseases may be informed by adult versions with some extrapolation and appropriate weighting).
Piloting & training	<ol style="list-style-type: none"> 1. Training is crucial and should focus on avoiding bias and expressing uncertainty. 2. Piloting should be undertaken. 	-
Level and conduct of elicitation	<ol style="list-style-type: none"> 1. Beliefs should be elicited from experts individually, even if a group interaction follows. 2. Although interaction between experts should be 	Group discussion may be needed to generate a distribution, for example in early technologies or when

	<p>structured through face-to-face sessions.</p> <p>3. Between-expert variation should be explored explicitly.</p> <p>4. Face-to-face where possible to allow a facilitator to deliver training to the expert.</p> <p>5. Feedback to experts should be given during the elicitation. Following feedback, experts should be given an opportunity to revise their distributions, either during or after an elicitation session.</p>	<p>eliciting more abstract/complex (non-observable) quantities cannot be avoided, for example those relating to service delivery, public health programmes or patient pathways.</p> <p>Practical constraints may dictate remote delivery of elicitation, for example through video conferencing.</p>
<p>Aggregation, analysis & post-elicitation validation</p>	<p>1. Probability distributions should be fitted to individually elicited judgements.</p> <p>2. Following fitting, a summary of the individual distributions should be obtained using linear pooling with equal weighting of experts.</p> <p>3. Any adjustments applied should be to improve coherence and consistency, and not to reduce variability. Internal and external review can be used to assess validity.</p> <p>4. Rationales for how the experts made their judgements should be collected and recorded post elicitation.</p> <p>5. All methodological choices for the elicitation must be documented and justified.</p>	<p>1. Pooling methods, other than linear pooling, may better accommodate expert heterogeneity. Further research is needed to explore which methods are most appropriate in these circumstances.</p> <p>2. Weighting may be preferable in some circumstances, for example where experts represent different disciplines or contribute different perspectives on the elicited quantities and therefore considerable heterogeneity is anticipated, but a single agreed consensus distribution is required. Weighting may be achieved implicitly through consensus or explicitly through performance weighting, although it is difficult to see how performance scores would be generated in this context.</p>

Figure 1: Evidence sources used to develop HCDM reference methods for elicitation



HCDM: Health Care Decision Making

1
2
3 **Title:** Developing a reference protocol for expert elicitation in healthcare decision making
4
5

6 **Authors:** Laura Bojke,¹ Marta Soares,¹ Karl Claxton,¹ Abigail Colson,² Aimée Fox,¹ Chris Jackson,³ Dina
7 Jankovic,¹ Alec Morton,² Linda Sharples,⁴ Andrea Taylor⁵
8
9

10
11 **Affiliations:**

12
13 ¹ Centre for Health Economics, University of York, York, UK

14
15 ² The Department of Management Science, University of Strathclyde, Glasgow, UK

16
17 ³ MRC Biostatistics Unit, University of Cambridge, Cambridge, UK

18
19 ⁴ London School of Hygiene and Tropical Medicine, London, UK

20
21 ⁵ Leeds University Business School, Leeds, UK
22

23
24 **Corresponding Author:**

25 Laura Bojke,

26 Centre for Health Economics, University of York, Heslington, York, YO10 5DD, UK

27 Laura.bojke@york.ac.uk

28
29 00441904321416
30
31

32
33 **Conflict of Interest:**

34
35 Competing interests: AM declares private consulting for AstraZeneca and Office of Health
36 Economics, not related to the current work. AM and AC were previously supported by a European
37 Union FP7 project which received in-kind support from pharmaceutical companies, including
38 Astellas, Roche and AstraZeneca. LB declares private consulting for Bresmed and Roboleo, not
39 related to the current work.
40
41
42
43
44

45 **Key words:**

46 Expert-elicitation, health, decision-making, transparency, uncertainty.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

Background

Many decisions in healthcare aim to maximise health, requiring judgements about interventions that may have higher health effects but potentially incur additional costs (cost-effectiveness framework). The evidence used to establish cost-effectiveness is typically uncertain and it is important that this uncertainty is characterised.

In situations where evidence is uncertain, the experience of experts is essential. The process by which the beliefs of experts can be formally collected in a quantitative manner is structured expert elicitation (SEE). There is heterogeneity in the existing methodology used in healthcare decision making (HCDM). A number of guidelines are available for SEE. It is not clear if any of these are appropriate for HCDM.

Objectives

The overall aim was to establish a protocol for SEE to inform HCDM. The objectives are:

- Provide clarity on methods for collecting and using experts judgements.
- Consider where alternative methodology may be required in particular contexts.
- Establish preferred approaches for elicitation on a range of parameters.
- Determine which elicitation methods allow experts to express uncertainty.
- Determine the usefulness of the reference protocol developed.

Methods

A mixed methods approach was used: systemic review, targeted searches, experimental work and narrative synthesis. A review of existing guidelines for SEE was conducted. This identified the approaches used in existing guidelines (the “choices”) and determined if dominant approaches exist. Targeted review searches were conducted for: selection of experts, level of elicitation, fitting and aggregation, assessing accuracy of judgements and heuristics and biases.

To sift through the available choices, a set of principles that underpin the use of SEE in HCDM was defined using evidence generated from the targeted searches, quantities to elicit, experimental evidence and consideration of constraints in HCDM. These principles, including fitness for purpose and reflecting individual expert uncertainty, were applied to the set of choices to establish a reference protocol. An applied evaluation of the developed reference protocol was also undertaken.

Results

For many elements of SEE, there was a lack of consistency across the existing guidelines. In almost all choices, there was a lack of empirical evidence supporting recommendations and in some circumstances the principles are unable to provide sufficient justification for discounting particular choices. It is possible to define reference methods for Health Technology Assessment (HTA). These include: a focus on gathering experts with substantive skills, eliciting observable quantities and individual elicitation of beliefs. Additional considerations are required for decision makers outside of HTA, for example at a local level, or for early technologies. Access to experts may be limited and in some circumstances group discussion may be needed to generate a distribution.

Limitations

The major limitation of the work conducted here, lies not in the methods employed in the current work, but the evidence available from the wider literature relating to how appropriate particular methodological choices are.

Conclusions

The reference protocol is flexible in many choices. This may be a useful characteristic, as it is possible to apply this reference protocol across different settings. Further applied studies, which use the choices specified in this reference protocol, are required.

(word count: 500)

Table of Contents

1		
2		
3		
4		
5		
6	Table of Contents.....	iii
7	List of Tables	iii
8	List of Figures	iii
9	Glossary	iii
10	Abbreviations.....	v
11	Plain English Summary	vii
12	Scientific Summary.....	viii
13	Chapter 1 Background	1
14	Chapter 2 Good practice in SEE: Learning from the available guidance.....	8
15	Section 2.1 Introduction	8
16	Section 2.2 Methods.....	8
17	Section 2.3 Included SEE guidelines.....	9
18	Section 2.4 Analysis of the elicitation process.....	9
19	Section 2.5 Identifying elicitation variables.....	23
20	Section 2.5.1 What quantities to elicit	23
21	Section 2.5.2 Encoding judgements.....	24
22	Section 2.6 Identifying and selecting experts.....	25
23	Section 2.7 Training and preparation	26
24	Section 2.8 Conducting the elicitation.....	26
25	Section 2.8.1 Mode and level of elicitation	26
26	Section 2.8.2 Feedback and revision	27
27	Section 2.8.3 Interaction	27
28	Section 2.8.4 Rationales	28
29	Section 2.9 Post-elicitation.....	28
30	Section 2.9.1 Aggregation.....	28
31	Section 2.9.2 Fit to distribution	30
32	Section 2.9.3 Other post-elicitation components.....	30
33	Section 2.10 Managing heuristics and biases	31
34	Section 2.11 Considering the validity of the process and results	31
35	Section 2.12 Conclusions	32
36	Chapter 3 Expert elicitation in different decision making contexts.....	33
37	Section 3.1 Introduction	33
38	Section 3.2 Levels of decision making	33
39	Section 3.2.1 Individual practitioners and local 'population level' decision makers	34

1		
2		
3	Section 3.2.2 National decision makers.....	36
4		
5	Section 3.2.3 Research commissioners.....	39
6		
7	Section 3.3 Conclusions	40
8	Chapter 4 Challenges in structured elicitation in HCDM	42
9		
10	Section 4.1 Introduction	42
11		
12	Section 4.2 Methods.....	42
13		
14	Section 4.3 Aspects related to the design of the SEE	42
15		
16	Section 4.4 Experiences with the conduct of the exercise	44
17		
18	Section 4.5 Experiences with the analyses and interpretation of elicited evidence	44
19		
20	Section 4.6 Conclusions	45
21		
22	Section 4.7 Discussion	45
23	Chapter 5 Reviewing the evidence: expert selection, level of elicitation, fitting and pooling	47
24		
25	Section 5.1 Introduction	47
26		
27	Section 5.2 Identifying literature	47
28		
29	Section 5.3 Selection of experts	48
30		
31	Section 5.4 Level of elicitation (individual versus group)	50
32		
33	Section 5.5 Fitting & pooling	54
34		
35	Section 5.5.1 Distribution choice and fitting	54
36		
37	Section 5.5.2 Mathematical aggregation.....	56
38		
39	Section 5.6 Assessing the expected accuracy of experts judgements.....	57
40		
41	Section 5.7 Conclusion.....	59
42	Chapter 6 Reviewing the evidence: Heuristics and biases	62
43		
44	Section 6.1 Introduction	62
45		
46	Section 6.2 Cognitive and motivational biases	63
47		
48	Section 6.2.1 Cognitive biases	63
49		
50	Section 6.2.2 Motivational biases.....	64
51		
52	Section 6.2.3 Overconfidence bias	64
53		
54	Section 6.3 Addressing psychological biases in in SEE.....	64
55		
56	Section 6.4 Technical bias reduction strategies.....	65
57		
58	Section 6.5 Behavioural bias reduction strategies with consistent support.....	65
59		
60	Section 6.5.1 Consider more information.....	65
	Section 6.5.2 Feedback	66
	Section 6.5.3 Avoid unnecessary anchors	67
	Section 6.5.4 Reduce bias through expert selection	67
	Section 6.6 Behavioural bias reduction techniques with conflicting evidence.....	68

1
2
3 Section 6.6.1 Bias warnings and training68
4 Section 6.6.2 Fixed value versus fixed probability methods.....68
5 Section 6.6.3 Face-to-face versus online elicitation69
6
7
8 Section 6.7 Conclusions69
9
10 Chapter 7 Quantities to elicit.....70
11 Section 7.1 Introduction70
12 Section 7.2 Overview of probability-, rate- and hazard-type parameters71
13 Section 7.2.1 Simple probability and conditional probability parameters71
14 Section 7.2.2 Transition probability parameters in discrete-time STMs71
15 Section 7.2.3 Time to event and survival.....72
16 Section 7.2.4 Continuous time-to-event decision models.....73
17 Section 7.2.5 Repeated event rates.....74
18 Section 7.3 Eliciting probability, rate and hazard-related parameters.....74
19 Section 7.4 Current practices in elicitation.....75
20 Section 7.4.1 Identification of examples75
21 Section 7.4.2 Simple and conditional probability or odds76
22 Section 7.4.3 Transition probability parameters in discrete-time STMs77
23 Section 7.4.4 Time to event and survival.....78
24 Section 7.4.5 Treatment effects on time-to-event distributions.....79
25 Section 7.4.6 Elicitation of hazards or of parameters of a time to event distribution directly81
26 Section 7.5 Steps and considerations in defining the quantities to elicit81
27 Section 7.6 Discussion83
28 Chapter 8 Three methodological experiments on the elicitation of subjective probabilistic belief.....85
29 Section 8.1 Introduction85
30 Section 8.2 General approach to the experiments.....86
31 Section 8.3 Methods.....87
32 Section 8.3.1 Overview of the experimental approach87
33 Section 8.3.2 Experiment 1: comparing different methods of elicitation89
34 Section 8.3.3 Experiment 2: Are individuals' able to 'extrapolate' from their knowledge-base? 90
35 Section 8.3.4 Experiments 3: To understand how individuals review their own probabilistic
36 assessments when presented with Delphi-type summaries.92
37 Section 8.4 Methods of analyses94
38 Section 8.4.1 Outcomes and metrics used94
39 Section 8.4.2 Methods.....96
40 Section 8.5 Results.....97
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1		
2		
3	Section 8.5.1 Description of the sample recruited	97
4		
5	Section 8.5.2 Experiment 1.....	98
6		
7	Section 8.5.3 Experiment 2.....	101
8		
9	Section 8.5.4 Experiment 3.1.....	102
10		
11	Section 8.5.5 Experiment 3.2.....	105
12	Section 8.6 Conclusions	107
13		
14	Section 8.6.1 Key findings from the experiments.....	107
15		
16	Section 8.6.2 Limitations and suggestions for future research (using the same experimental approach).....	109
17		
18	Section 8.6.3 Experimental approach vs. using almanac quantities.....	110
19	Chapter 9 Consideration of the methodological choices emerging from existing guidelines	111
20		
21	Section 9.1 Introduction	111
22		
23	Section 9.2 Principles underpinning the use of SEE to inform HCDM	111
24		
25	Section 9.2.1 Principle 1. Transparency.....	112
26		
27	Section 9.2.2 Principle 2. Fitness for purpose.....	113
28		
29	Section 9.2.3 Principle 3. Consistency, but respecting constraints of the decision making context	113
30		
31	Section 9.2.4 Principle 4. Reflecting uncertainty at the individual expert level	113
32		
33	Section 9.2.5 Principle 5. Recognising and acting on biases.....	113
34		
35	Section 9.2.6 Principle 6. Suitability for experts who possess substantive skills, who are less likely to be normative.....	114
36		
37	Section 9.2.7 Principle 7. Recognising where adaptive skills are required.....	114
38		
39	Section 9.2.8 Principle 8. Recognising and act on between-expert variation.....	114
40		
41	Section 9.2.9 Principle 9. Promoting high performance.....	115
42	Section 9.3 How do SEE elements and methodological choices reflect the principles underpinning healthcare?.....	115
43		
44	Section 9.3.1 Selecting quantities (preparation and design).....	115
45		
46	Section 9.3.2 Methods to encode judgements (preparation and design).....	117
47		
48	Section 9.3.3 Selecting experts.....	118
49		
50	Section 9.3.4 Pilot exercise.....	119
51		
52	Section 9.3.5 Training and preparation for experts.....	119
53		
54	Section 9.3.6 Level of elicitation (elicitation).....	120
55		
56	Section 9.3.7 Mode of Administration (elicitation)	121
57		
58	Section 9.3.8 Feedback to experts and revision (elicitation).....	121
59		
60	Section 9.3.9 Opportunity for interaction (elicitation)	122
	Section 9.3.10 Feedback from experts on process (elicitation).....	123

1
2
3 Section 9.3.11 Rationales (elicitation)123
4 Section 9.3.12 If/how to aggregate (aggregation, analysis and post-elicitation)124
5 Section 9.3.13 Fit to distribution (aggregation, analysis and post-elicitation)125
6 Section 9.3.14 Adjusting judgements (aggregation, analysis and post-elicitation)126
7 Section 9.3.15 Documentation (aggregation, analysis and post-elicitation)126
8 Section 9.3.16 Managing biases127
9 Section 9.3.17 Validation127
10 Section 9.4 Conclusions128
11 Chapter 10 Reference protocol for expert elicitation in healthcare130
12 Section 10.1 Evidence in support of a reference protocol for HCDM130
13 Section 10.2 How the evidence is used to generate a reference protocol for SEE in HCDM131
14 Section 10.3 Reference protocol for SEE in HTA.....131
15 Section 10.4 Important considerations for decision makers outside of the HTA setting133
16 Section 10.5 Conclusions134
17 Chapter 11 Applied evaluation of developed reference protocol136
18 Section 11.1 Background136
19 Section 11.2 The evaluation topic136
20 Section 11.2.1 Diagnosis of asthma137
21 Section 11.2.2 Diagnostic model developed.....137
22 Section 11.3 Description of elicited parameters138
23 Section 11.3.1 Time until correct diagnosis.....138
24 Section 11.4 Application of developed reference protocol139
25 Section 11.4.1 Selecting the quantities (preparation and design stage)140
26 Section 11.4.2 Methods to encode judgements (preparation and design stage)141
27 Section 11.4.3 Validation (preparation and design stage).....141
28 Section 11.4.4 Selecting experts (preparation and design stage).....141
29 Section 11.4.5 Pilot exercise (preparation and design stage).....142
30 Section 11.4.6 Training and preparation for experts (preparation and design stage).....142
31 Section 11.4.7 Level of elicitation (elicitation stage)142
32 Section 11.4.8 Mode of administration (elicitation stage)142
33 Section 11.4.9 Feedback to experts and revision (elicitation stage)143
34 Section 11.4.10 Opportunity for interaction (elicitation stage)143
35 Section 11.4.11 Feedback from experts on process (elicitation stage)143
36 Section 11.4.12 If/how to aggregate (aggregation, analysis and post-elicitation)144
37 Section 11.4.13 Fit to distribution (aggregation, analysis and post-elicitation)144
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1	
2	
3	Section 11.4.14 Data Protection and Anonymity (aggregation, analysis and post-elicitation) ..144
4	
5	Section 11.5 Results.....144
6	
7	Section 11.5.1 Elicitation Results.....146
8	
9	Section 11.6 Feedback from experts on the process (including rationales and validation)146
10	
11	Section 11.6.1 Observability of the quantity asked146
12	
13	Section 11.6.2 Completion of exercise146
14	
15	Section 11.6.3 Wording of the question147
16	
17	Section 11.6.4 Value of interaction147
18	
19	Section 11.6.5 Value of facilitator147
20	
21	Section 11.6.6 Experts rationales147
22	
23	Section 11.6.7 General feedback from experts148
24	
25	Section 11.6.8 Practicality of conducting the SEE process148
26	
27	Section 11.7 Conclusion.....148
28	
29	Chapter 12 Discussion & conclusions150
30	
31	Section 12.1 Conclusions on evidence generated150
32	
33	Section 12.2 Key considerations for using the reference protocol in HCDM152
34	
35	Section 12.3 Key areas for further research154
36	
37	Section 12.4 Limitations of the work conducted156
38	
39	Chapter 13 Acknowledgements158
40	
41	Chapter 14 Contributions of authors.....159
42	
43	Chapter 15 Patient and public involvement.....160
44	
45	Chapter 16 References161
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	

List of Tables

1
2
3
4
5 Table 1 Summary of the elicitation elements, components, and choices described in SEE guidelines 10
6 Table 2 Level of agreement on recommendations and choices in SEE guidelines18
7 Table 3: Summaries of selected survival distributions73
8 Table 4: List of published examples of eliciting probability, rate and time to-event parameters76
9 Table 5 : Wording of questions in experiment 291
10 Table 6 Sample characteristics97
11 Table 6 Response to question 1 about the ease of completion.100
12 Table 7 Response to question 2 about method preference.100
13 Table 8 Proportion of participants who revised their priors.103
14 Table 9 Accuracy of initial priors compared between participants who did and did not revise their
15 priors: mean (SD) [median and (interquartile range)] of accuracy metric over participants103
16 Table 10 Proportion of participants who revised their priors.105
17 Table 11 Results of experiment 3.2: Outcomes in participants who did and did not revise their priors.
18106
19 Table 12 Key principles of SEE in HCDM111
20 Table 13 A reference case for HTA132
21 Table 14 Additional issues in generating a reference case outside of HTA133
22 Table 15 Application of reference case: Parameters elicited in DAR¹⁹⁶138
23 Table 16 Application of reference case: Summary of experts recruited145
24 Table 17 Application of reference case: Elicitation Results145
25 Table 18 Areas for further research on SEE in HCDM154
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

List of Figures

1	
2	
3	
4	
5	Figure 1 General schematic for SEE4
6	Figure 2 Summary of project activities5
7	Figure 3 The elicitation process10
8	Figure 4 Levels of expert elicitation*51
9	Figure 5 Elicited points on a cumulative distribution function (CDF) and alternative fitted
10	distributions.....55
11	Figure 6 Example of a transition diagram for a state-transition model (STM)72
12	Figure 7 Example of eliciting conditional probabilities for a decision tree.....77
13	Figure 8 Example of eliciting transition probabilities for a novel setting in a discrete-time STM78
14	Figure 9 Snapshot of the R shiny app (1)87
15	Figure 10 Snapshot of the R shiny app (2)87
16	Figure 11 Snapshot of the R shiny app (4)91
17	Figure 12 Snapshot of the R shiny app (5)91
18	Figure 13 Illustrative example of the scenarios evaluated in experiment 3.193
19	Figure 14 Illustration of the scenarios defined for experiment 3.294
20	Figure 15 Distribution of bias, lnSDR and lnKL. (SDR= SD elicited / SD true).....98
21	Figure 16 Within-participant comparison of accuracy when different elicitation methods were used,
22	for each precision scenario.....100
23	Figure 17 Within-participant comparison of accuracy with and without extrapolation, for different
24	levels of extrapolation.102
25	Figure 18 Within-participant comparison of accuracy to initial priors and extent of revision, for
26	different types of group summaries.104
27	Figure 19 Within-participant comparison of accuracy and extent of revision (using selected outcome
28	metrics), for different types of group summaries.106
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	

1
2
3 **List of supplementary material**
4
5
6

7 Supplementary material 1: Review of SEE guidelines

8
9 Supplementary material 2: Design of the experiments

10
11 Supplementary material 3: Analysis of the experiments

12
13 Supplementary material 4: Applying choices to the principles for HCDM

14
15 Supplementary material 5: The evaluation of the protocol
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

DO NOT COPY
For Peer Review

Glossary

Adaptive expertise

The expert's ability to adapt their knowledge to new situations for which they do not have prior experience, such as entirely new medical interventions or patients with different characteristics.

Aleatory uncertainty

Uncertainty due to randomness. Inherently irreducible and unpredictable in nature.

Behavioural aggregation

The process of grouping together individual experts to generate an overall, **consensus** aggregate distribution

Beta-binomial

The binomial distribution in which the probability of success at each trial is fixed but randomly drawn from a beta distribution.

Biases

Systematic errors in the processes that people use to make judgements. The biases most commonly referred to in SEE are cognitive and motivational biases.

Bisection method

A method of elicitation that can be used for any type of continuous univariate distribution. The expert is asked to divide the interval into two equally likely intervals, for example the interval containing the minimum and the median.

Calibration

A process to determine the accuracy of a probabilistic judgement. Estimates of expected calibration can be used to weight expert judgements.

Central tendency

A measure of the "centre" or typical value for a probability distribution. Examples are the mean, median and the mode.

Chips and bins

A graphical representation of the **fixed interval method** for elicitation Experts are asked to give an interval for an uncertain quantity and then place 'chips' in 'bins' which divide this interval. The chips represent the weight of their belief.

Choice(s)

See Elements

Consensus

The act of reaching agreement. In SEE this refers to agreement between experts on a distribution.

Credible range

A summary of a probability distribution for an uncertain quantity, describing an interval within which the quantity falls with a particular probability.

Decomposition

See **disaggregation**

Delphi method

A structured communication technique based on several rounds of questionnaires, feedback and revision. The modified Delphi (EFSA) is a form of Delphi method used to elicit an uncertain quantity.

Dependence

A statistical relationship between two random variables.

Domain-specific

Relating to a particular discipline, e.g. health, education, engineering, or subdivision, e.g. oncology, geriatrics.

Disaggregation

Dividing into constituent parts. In SEEE this refers to **decomposition** of a complex quantity into less complex, **observable quantities**.

Element(s)

The SEE process comprises numerous elements which encompass several possible components for which choices need to be made. The selection of experts is an example of a SEE element for which the analyst needs to make choices regarding the different components such as how many expert to include, how to recruit the experts and the type of expertise the expert should possess and so on.

Epistemic uncertainty

Uncertainty which arises primarily from limited or imperfect knowledge. It is, in principle, reducible by obtaining more or better information.

Expert

The individual (s) from whom subjective beliefs are sought. Experts may be defined as such on the basis of their **substantive, normative or adaptive expertise**.

Facilitator

An unbiased, impartial individual that works with experts to obtain their **subjective beliefs**. This may involve coordinating active discussion between experts to achieve a consensus, or less hands-on guidance to enable individuals to provide their own judgements.

Fit for purpose

Here 'fit for purpose' relates to the suitability of a technique or results for its designated role or purpose. In the context of HCDM this will often involve future statistical analysis or modelling.

Fixed interval method

A method of elicitation in which experts are presented with an interval and asked to assess the probability that the quantity will fall into that interval.

Frequencies

The observed number of successes or failures out of a finite number of trials.

Hazard

The probability of transition in a short time interval divided by the length of the interval, in the limit, as this length becomes shorter.

Kullback-Liebler

A measure of the difference between two distributions. In the context of elicitation, this is a measure of the information lost when the true distribution is approximated by the elicited distribution¹³.

Level of elicitation

Describes either an elicitation conducted at the individual level (which may be followed by **mathematical aggregation**) or at the group level (**behavioural aggregation**).

Heuristics

Mental shortcuts that ease the cognitive load of making a decision, e.g. rule of thumb, an educated guess, an intuitive judgment, a “guesstimate”, profiling, or common sense.

Linear opinion pooling

A mechanistic rule for combining probabilities or distributions elicited from two or more sources into a single probability or distribution. In linear opinion pooling the single distribution is calculated as the unweighted linear average of individual distributions.

Mathematical aggregation

Combining the beliefs of individual experts using a mathematical rule, such as **linear opinion pooling**.

Model based economic evaluation (MBEE)

An evaluation of cost-effectiveness which employs some form of decision model or statistical model.

Normative expertise

The expert's ability to accurately assess and clearly communicate their beliefs in probabilistic form.

Observable quantities

A quantity that could be estimated as a simple function of observed data, if that data were available. For example, the probability of an outcome, which can be estimated by the observed frequency of the outcome. This contrasts with composite quantities such as odds ratios which are complex functions of observable data and generally more difficult for an expert to conceptualise.

Parameters

A variable within an analysis, for example a **MBEE** or a regression model.

Precision

A measure of statistical variability or statistical bias. Repeated measures are said to be precise if the values are close together. Calculated as the reciprocal of the variance.

Probabilistic

Relating to probabilities. Involving quantities whose values are uncertain, or which may take multiple possible values. Probabilistic sensitivity analysis is used to explore the consequences of uncertainty in **parameter** inputs in **MBEE**.

Probability

1
2
3 A measure of the likelihood/chance of occurrence of a particular event. Can take values between 0
4 and 100%.
5

6
7 **Proportion**

8 The number of something in comparison to the whole, e.g. the proportion of females in the
9 population. Can take values between 0 and 100%.
10

11
12 **Quantiles**

13 Points on a probability distribution which divide it into continuous intervals. The 50th quantile is
14 known as the *median*.
15

16
17 **Relative risk**

18 The ratio of the probability of an event occurring in the exposed or treated group versus the
19 probability of the event occurring in the non-exposed group.
20

21
22 **Structure Expert Elicitation (SEE)**

23 Refers to a formal documented process by which experts beliefs (priors) are obtained in a
24 quantitative form.
25

26
27 **Subjective beliefs**

28 An individual's own beliefs/opinions about an uncertain quantity, which may be expressed as a
29 distribution. If this is prior to data collection/availability, this is called a "prior" distribution.
30

31
32 **Substantive expertise**

33 An expert is said to be a substantive expert if they possess skills/knowledge pertaining to a particular
34 domain or subject within that domain.
35

36
37 **Survival function**

38 Defined as the probability that the time T that an event (e.g. death) occurs is greater than t, $P[T>t]$. It
39 can be defined as $S(t) = \exp\{-\int_0^t h(u)du\} = \exp\{H(t)\}$, where $\int_0^t h(u)$ is the cumulative
40 hazard function, H(t). Survival can alternatively be described using the probability density function
41 for the survival times, f(t), using the following relationship: $S(t) = 1 - F(t) = 1 - \int f(t)dt$, where $\int f(t)dt$
42 is the cumulative distribution function, F(t). The hazard and the probability density functions can also
43 be used together to determine survival $S(t) = f(t)/h(t)$.
44

45
46 **Variance**

47 A measure of the spread of a random variable.
48

49
50 **Variable interval method**

51 A method to elicit a distribution. The expert is asked to express the quartiles or credible intervals of
52 a distribution, e.g. tertiles are used in the **bisection method**.
53

54
55 **Validity**

56 Generally refers to the quality of making logical sense. In SEE, validity can mean that the exercise
57 captured what the experts believe, or that the expressed quantities correspond to reality, or are
58 consistent with the laws of probability, or are internally coherent.
59
60

Abbreviations

HTA	Health Technology Assessment
NICE	National Institute for Health and Clinical Excellence
RCT	Randomised Control Trial
SEE	Structured Expert Elicitation
SHELF	Sheffield Elicitation Framework
HCDM	Healthcare decision-making
DAR	Diagnostic Assessment Report
EFSA	European Food Safety Authority
EPA	Environmental Protection Agency
IFA	Institute and Faculty of Actuaries
NRC	Nuclear Regulation Commission
GP	General Practitioner
CCG	Clinical Commissioning Group
LA	Local Authority
DHSC	Department of Health and Social Care
MRC	Medical Research Council
PHE	Public Health England
IFR	Individual Funding Request
SMC	Scottish Medicines Consortium
AMWSG	All Wales Medicines Strategy Group
HPA	Health Protection Agency
PHP	Public Health Programme
STA	Single Technology Appraisal
NIHR	National Institute for Health Research
MBEE	Model Based Economic Evaluation
FIM	Fixed Interval Method
VIM	Variable Interval Method
PEGs	Prior Elicitation Graphical Software
GEM	Generalised Expertise Measurement
E-SQ	Expert-Selection Questionnaire
IDEA	Investigate, Discuss, Estimate, Aggregate
CDF	Cumulative Distribution Function
STM	State Transition Model

1		
2		
3	DES	Discrete Event Simulation
4		
5	HRQoL	Health Related Quality of Life
6		
7	TPM	Transition Probability Matrix
8		
9	Crls	Credible Intervals
10		
11	KL	Kulback Leibler
12	MCMC	Markov Chain Monte Carlo
13	SDR	Standard Deviation Ratio
14		
15	InSDR	Log of Standard Deviation Ratio
16		
17	InKL	Log of Kulback Leibler
18	SCHARR	Sheffield School of Health and Related Research
19		
20	FeNO	Fractional Exhaled Nitric Oxide
21		
22	BTS	British Thoracic Society
23		
24	SIGN	Scottish Intercollegiate Guidelines Network
25	FTE	Full-Time Equivalent
26		
27		
28		
29		
30		
31		
32		
33		
34		
35		
36		
37		
38		
39		
40		
41		
42		
43		
44		
45		
46		
47		
48		
49		
50		
51		
52		
53		
54		
55		
56		
57		
58		
59		
60		

Plain English Summary

Background

Decisions in health care aim to maximise health, requiring judgements about treatments. The evidence used to make these judgement is typically uncertain.

In these situations, the experience of experts is essential. Structured expert elicitation (SEE) collects beliefs from experts. There are different guidelines available for SEE, however it is not clear if any of these be can be used in healthcare decision making (HCDM), for example in considering if a treatment should be made available in the NHS. This project aimed to develop a guidance for SEE to inform HCDM.

Methods

Reviews and experimental techniques were used to gather a list of methods to conduct SEE. The suitability of these choices in HCDM was then determined by comparing these with a set of standards that support the use of SEE in HCDM was generated.

Results

Different guidelines prefer different approaches to conduct SEE. There is a lack of evidence available to determine which of these methods is most appropriate across the whole of HCDM.

It is possible to define reference protocol methods that could be used in a particular type of HCDM, Health Technology Assessment (HTA). This includes: gathering experts with knowledge of the clinical area, asking experts about things that they observe in clinical practice and asking experts individually for their beliefs. For decision makers working outside of HTA, for example at a local level, or for treatments that are not yet available to patients, these choices may not be appropriate.

Conclusions

This flexibility of this guidance is a useful feature. It is possible for different decision makers in healthcare to interpret the reference protocol for their own circumstances.

Scientific Summary

Background

At the forefront of decisions in healthcare is the aim of maximising health, requiring judgements about interventions that may have higher health effects but potentially incur additional costs. The evidence used to establish cost-effectiveness is typically uncertain, for example the evidence may not be on 'final' outcomes (e.g. cancer products licensed on evidence of progression-free survival), or the evidence base may not be well developed (e.g. in diagnostics, medical devices, early access to medicines scheme). It is important that the uncertainty in this evidence is characterised. If not, any analysis using this evidence may give decision makers a misleading view of the risks associated with their decision.

In situations where evidence is subject to uncertainty, the experience of experts may be necessary essential. To ensure accountability in the decision, these expert judgements should be made explicit and incorporated transparently into the decision making process. The process by which the beliefs of experts can be formally collected in a quantitative manner is structured expert elicitation (SEE). If conducted in an appropriate manner, SEE can characterise uncertainties associated with the cost-effectiveness of competing interventions and assess the value of further evidence. This may be the approach best suited to a transparent decision making process.

There is an increasing interest in SEE as new technologies are assessed progressively closer to their launch on the market. SEE is also valuable for 'early modelling' of new interventions or unknown diseases for which little or no evidence is available. A review of applied studies in health care decision making (HCDM) found heterogeneity in methodology used and a lack of consideration for any existing guidance on the topic.¹

No standard guidelines exist to conduct expert elicitation in health technology assessments (HTA) but there are a number of generic guidance documents, some of which have been used in HTA. The most notable of these are the Sheffield Elicitation Framework (SHELF) and Cooke's classical method. It is not clear if any of the existing guidelines, generic and domain specific are appropriate for us in HCDM.

Objectives

The overall aim of this report was to establish a reference protocol or guideline for the elicitation of experts judgements to inform HCDM. To achieve this overall aim, the report focussed on the following objectives:

1. Providing clarity on the methods for collecting and using experts judgements within an assessment of cost-effectiveness.
2. Exploring where alternative methodology may be required in particular context/constraints e.g. time.
3. Establishing preferred approaches for elicitation for a range of parameters and a range of decision-making contexts.
4. Determining which elicitation methods allow experts express parameter uncertainty as opposed to variability.
5. Determining the applicability and usefulness of the reference protocol developed within a case-study application.

For (4), statistical experiments were conducted. The aim of these experiments was threefold: (1) to evaluate alternative methods of elicitation and how they perform in representing parameter uncertainty; (2) to explore individuals' ability to extrapolate from their knowledge base; and (3) to explore how individuals revise their answers when presented with group summaries.

Methods

To achieve these objectives, a mixed methods approach was used, combining formal systematic review, targeted searches, experimental work and narrative synthesis. Specifically, first a systematic review of existing guidelines for formal elicitation, published in either the peer-reviewed or grey literature, was conducted. This identified the approaches used in existing guidelines (the "choices") and determined if dominant approaches evolve. Less formal, targeted searches were also conducted to determine the state of the evidence on choices relating to: the selection of experts, level of elicitation, fitting and aggregation, assessing expected accuracy of experts judgements and heuristics and biases. The advantages and disadvantages of each available choice for these elements were extracted from the papers and potential constraints to their application in HCDM determined.

HCDM is not a homogenous domain, as different decision makers face different constraints and this may have implications for expert elicitation methodology. The contexts in which SEE in HCDM may

1
2
3 be conducted, is therefore discussed in detail and conclusions made regarding the use of a reference
4 protocol for SEE. Alongside this, a systematic review of SEE applications in cost-effectiveness
5 modelling was undertaken. This details the challenges that were reported by the authors conducting
6 these analyses. When available, the basis for the methodological choices made in each application is
7 extracted. This also provided a view of the current scope of the landscape with regards to applied
8 SEE in HCDM.
9
10
11
12

13
14
15 When designing a SEE, deciding what quantities to elicit is a major challenge. There is no guidance
16 covering the spectrum of quantities that may be appropriate to elicit to inform HCDM, including
17 measures of treatment effects and baseline event rates. To address this lack of guidance, a review
18 was undertaken of alternative quantities that can be elicited to inform the probability-related or
19 time to event-related parameters commonly used in HCDM.
20
21
22
23

24
25 The statistical experiments, conducted to explore multiple uncertainties in SEE methodology, utilised
26 a simulated learning process (e.g. Wang 2002). Individual's knowledge was determined by recorded
27 observations. The 'dataset' observed then determines participants' belief about the quantity of
28 interest, from which accuracy can be measured. This approach allows the conditions of the
29 experiment to be defined, e.g. equal vs different knowledge-base, and isolate potential
30 determinants, e.g. precision. Participants were shown random observations from a statistical model
31 representing an abstract medical problem. Following this they were asked to express their beliefs'
32 regarding treatment effectiveness. All participants (72) were students at the University of York, the
33 large majority of which were undergoing clinical training. The exercises was delivered face-to-face
34 and financial incentives were offered according to accuracy. The experiments measured:
35
36
37
38
39
40
41

- 42 • Bias: difference in the means of the true and elicited (and fitted) distributions;
- 43 • Uncertainty: ratio of the SDs of the two distributions.
- 44 • Kullback-Leibler divergence (KL): information lost when one distribution is approximated by
45 another¹³.
- 46 • Participants preference for alternative methods.
47
48
49
50
51

52 Given the full range of evidence generated, on which to base a reference protocol for SEE in HCDM,
53 it was necessary to use this evidence to generate a set of principles that underpin the use of expert
54 elicitation in HCDM. Available choices, from the review of guidelines, are considered in light of these
55 principles and any empirical evidence available to support the choices. This informs the reference
56 protocol, by discounting or supporting particular choices.
57
58
59
60

1
2
3 The work also included an applied evaluation of the developed reference protocol. This uses an
4 existing cost-effectiveness model, in which SEE was used to generate initial estimates of uncertain
5 parameters. In addition to demonstrating the usefulness of the reference protocol in navigating the
6 SEE process, the practicality of SEE is determined using narrative feedback from experts and by
7 generating estimates of resources required to design and conduct the SEE.
8
9
10
11
12

13 Finally, a dissemination workshop was convened, which explored the usefulness, and challenges in
14 using SEE in HCDM. It was also used to refine, using discussion, a set of recommendations for further
15 research.
16
17
18
19

20 *Results*

21 A comprehensive list of elements and choices for SEE was developed by reviewing existing protocols
22 (work package 1). This covered the design, implementation and analysis stages of SEE. The review
23 showed that for many elements of the SEE, there was a lack of consistency across the existing
24 guidelines. Targeted searches also revealed that the majority of choices are not supported by any
25 empirical evidence, both specific to HCDM and more generally.
26
27
28
29
30

31 Empirical evidence, generated by the experiments conducted here (work-packages 2 and 3),
32 determined that there is little difference between variable interval methods (VIM) and fixed interval
33 methods (FIM) methods to encode judgements, in terms of procedural performance. Therefore a
34 decision maker can consider either of these choices suitable. This experiment also determined that
35 participants did not adjust uncertainty levels sufficiently to reflect differences in the underlying
36 heterogeneity of the populations; in particular, uncertainty was consistently underestimated in the
37 case of high heterogeneity. This case is frequently encountered in healthcare settings. The
38 experiments also sought to explore extrapolation beyond data observed and updating of priors after
39 presentation of group summaries – issues which feed into multiple choices for SEE. It was difficult to
40 form definitive conclusions given that the experiments were underpowered for these elements. The
41 experiments did provide some evidence that experts changed their estimates in a rational way when
42 provided with estimates from others, suggesting that group discussion or feedback may be useful.
43 Extrapolation outside of the observed sample does not seem to affect accuracy, suggesting that it is
44 reasonable to ask experts about patients and practices in which they do not have direct clinical
45 experience, or for whom there is no relevant literature.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 In order to sift through the available choices, a set of principles that underpin the use of SEE in HCDM
4 was defined using evidence generated from targeted searches, experimental evidence on methods to
5 encode judgements and consideration of the constraints on the decision making processes in health
6
7
8 (work package 1). These nine principles are:

- 10 1. Transparency
- 11 2. Fitness for purpose
- 12 3. Consistency, but respecting constraints of the decision making context
- 13 4. Reflecting uncertainty at the individual expert level
- 14 5. Recognising and acting on biases
- 15 6. Suitability for substantive experts who are less likely to be normative
- 16 7. Recognising where adaptive skills are required
- 17 8. Recognising between-expert variation
- 18 9. Promoting high performance

19
20
21
22
23
24
25
26
27 Not all principles for SEE in HCDM were relevant for all elements. The most relevant principles for
28 each element and components within SEE were considered.

29
30
31
32 In almost all choices, there is a lack of empirical evidence and in some circumstances the principles
33 are unable to provide sufficient justification for discounting particular choices (work package 1). It
34 is, however, possible to define reference methods that could be used in a more narrowly defined
35 area of HCDM, namely HTA. These include:

- 36
37
38
39
40 • Focus on gathering substantive expertise or experience. Normative skills can be developed
41 during the training session as part of the SEE.
- 42
43 • Simple observable quantities should be elicited where possible; ratios or complex
44 parameters such as regression coefficients should not be elicited directly.
- 45
46 • Minimize and record conflicts of interest among the experts. Include experts external to the
47 SEE task, i.e. not those involved in developing the task.
- 48
49 • Dependence between variables should be captured in SEE. Expressing dependent variables
50 in terms of independent variables is preferable when experts do not have strong normative
51 skills.
- 52
53 • Use of either VIM or FIM work well, however decision makers should aim for consistency
54 across applications.
- 55
56 • Beliefs should be elicited from experts individually, even if a group interaction follows.
- 57
58 • Between-expert variation should be explored explicitly.
- 59
60

- Following fitting, a summary of the individual distributions should be obtained using linear pooling.
- Interaction should be face-to-face where possible to allow a facilitator to deliver training to the expert.
- Training is crucial and should focus on avoiding bias and expressing uncertainty.
- All methodological choices for the SEE must be documented and justified.

Additional considerations are required for decision makers outside of HTA, for example at a local level, or for early technologies that have yet to progress through the regulatory process. Access to experts may be limited and in some circumstances, group discussion may be needed to generate a distribution.

The application of the case-study, a diagnostic model for asthma, explored practical issues. This highlighted sufficient information needs to be presented to the experts. The level of information presented to the experts and the wording of this information is paramount in ensuring that the quantity of interest is observable to the expert. When deciding on the information to provide to experts it may be useful to consult existing policies. With regards to time constraints, the applied evaluation was undertaken over a 7 month period and involved 3 analysts in varying proportions. Overall this equated to 5 months full time equivalent researcher time.

Limitations

The major limitation of the work conducted here, lies not in the methods employed, but the evidence available from the wider literature, on which to base the set of choices and determine how appropriate these are. Concluding on the suitability of the choices available from the existing guidelines is challenging due to the lack of empirical evidence to support specific choices. Instead, it was necessary to develop principles for SEE in HCDM, using the sources of evidence as described above and published guidelines for good SEE. Using only the principles, in the absence of empirical evidence, meant that it was not always possible to give definitive conclusions on choices.

Areas for further research

In considering the appropriateness of choices for SEE in HCDM and exploring how these choices may be affected by the context in which the SEE is applied, there are several areas in which further research is required before definitive statements can be made regarding their appropriateness for a reference protocol. Researchable questions in these areas include the following:

- Which methods for expert recruitment are most practical and what are the challenges?
- What training strategies can be used to minimise bias?
- Which methods for eliciting dependent quantities work best for non-normative experts?
- Which consensus approach works best in HCDM in practice and for which types of quantities and decision makers?
- Should individual priors be combined when there is significant expert variation? If so, how?

At the dissemination workshop, participants were asked to discuss areas for further research, specifically considering what decision makers in HCDM may require when determining a reference protocol for SEE for use within their setting. Participants were not asked to define which research topics are highest priority for their setting. Selecting experts, minimising bias, adaptation to specific setting in which SEE may be applied, for example choosing individual or group elicitation, appropriate wording of questions, methods for multivariate elicitation and what information should be presented to the experts to help them formulate their beliefs. Some of these topics would benefit from empirical research and others may be resolved through application of the proposed reference protocol to HCDM, including in settings with a range of constraints.

Conclusions

SEE can offer opportunities in HCDM, particularly reimbursement decisions supported by model based economic evaluation (MBEE). SEE allows the uncertainty in the evidence used to populate these models to be characterised, or where evidence is completely lacking, provides additional information needed to reach a decision.

The work described in this report has attempted to generate evidence which is useful for analysts and decision makers in HCDM. SEE conducted in this context to date, has not used a set of consistent methods, and above all, has not considered the implications of the choices made when designing and conducting a SEE. To improve the accountability of HCDM the procedure used to derive expert judgements should be transparent.

The reference protocol presented here is intended to serve as a guide to good practice and reporting, and is flexible in many choices rather than being prescriptive regarding methods. It can therefore be thought of as a reference guide. This was necessary due to the lack of empirical data specific to HCDM and more generally in SEE. This may be a useful characteristic, as it is possible to apply this reference protocol across different settings.

1
2
3 *Funding details*
4

5 This work was funded by the Medical Research Council. HEE: Developing a reference protocol for
6 expert elicitation in health care decision making. Bojke L. 2017. Reference: MR/N028511/1.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

DO NOT COPY
For Peer Review

Chapter 1 Background

In the UK, decisions about the use of health care interventions are made by various NHS organisations as well as the immediate beneficiaries, namely patients. In England, these NHS organisations include the National Institute for Health and Care Excellence (NICE), NHS England and Public Health England. At the forefront of these decisions is the aim of maximising health, calling for judgements about the interventions that are expected to lead to higher health effects. Where resources are limited, additional costs incurred will impact on access to care for other patients, and health foregone in this way should also be taken into account (a cost-effectiveness framework).²

Although Randomized Control Trials (RCTs) have been described as the principle source of evidence for such decision-making, RCTs have considerable limitations including a lack of external validity, short study periods to assess long-term treatment effect and invalid generalizations of findings outside the study group.³⁻⁵ Also RCTs are not possible or ethical in some situations.

These limitations also impact the use of RCTs for urgent health issues for which decisions need to be made promptly on the basis of limited and often imperfect available data.⁶ Health Technology Assessments (HTAs) traditionally use decision-modelling methods that gather different forms of evidence by defining mathematical relationships between a varied set of input parameters, in a way that describes aspects of the history of the disease of interest and the impact of the intervention.

Uncertainty in the evidence is pervasive in cost-effectiveness modelling and the analysis may be biased if uncertainty in the model inputs is not reflected. Uncertainty can be distinguished as epistemic or aleatory.^{7,8} Aleatory uncertainty arises due to randomness i.e. unpredictable variation in a process and expert knowledge cannot reduce this type of uncertainty.⁷ Therefore, it is sometimes referred to as irreducible uncertainty. Epistemic uncertainty is due to imperfect knowledge and it can be reduced with sufficient study and therefore, expert judgement may be useful in its reduction.⁷ Additional evidence can reduce uncertainty and provide a more precise estimate of cost-effectiveness. By quantifying uncertainty, it is possible to assess the potential value of additional evidence, inform the types of evidence that might be needed, and consider restricted use until the additional evidence becomes available.⁹

In some situations, several input parameters in the decision model may have only limited empirical data. For example, the evidence may not be on 'final' outcomes (e.g. cancer products licensed on evidence of progression-free survival), or the evidence base may not be well developed (e.g. in the

1
2
3 areas of diagnostics, medical devices, early access to medicines scheme, or public health). In these
4 situations, judgements are required for a decision to be reached regarding that parameter. To
5 ensure accountability in the decision, these judgements should be made explicit and incorporated
6 transparently into the decision making process, an inherently Bayesian view on decision-making.
7 Formal methods to quantify prior beliefs in the form of experts judgements exist, and are termed
8 Structure Expert Elicitation (SEE) methods.⁸
9

10
11
12
13
14
15 Structural Expert Elicitation (SEE) is a process that allows experts to express their beliefs in a
16 statistical, quantitative form. If conducted in an appropriate manner, SEE is the best approach to
17 characterise uncertainties associated with the cost-effectiveness of competing interventions and
18 assess the value of further evidence. SEE methods have been used in disciplines including weather
19 forecasting and reliability analysis within engineering¹⁰ but the research findings in these disciplines
20 are often interpreted as contradictory, in particular the appropriateness of generating consensus
21 amongst experts.¹¹ In terms of SEE in healthcare, NICE uses expert judgement across all guidance-
22 making programmes but, expert elicitation (compared to expert opinion) is used less frequently.¹²
23 Existing timelines and consequent time constraints are reported as the common obstacles when
24 conducting expert elicitation in healthcare.¹²
25
26
27
28
29
30
31
32

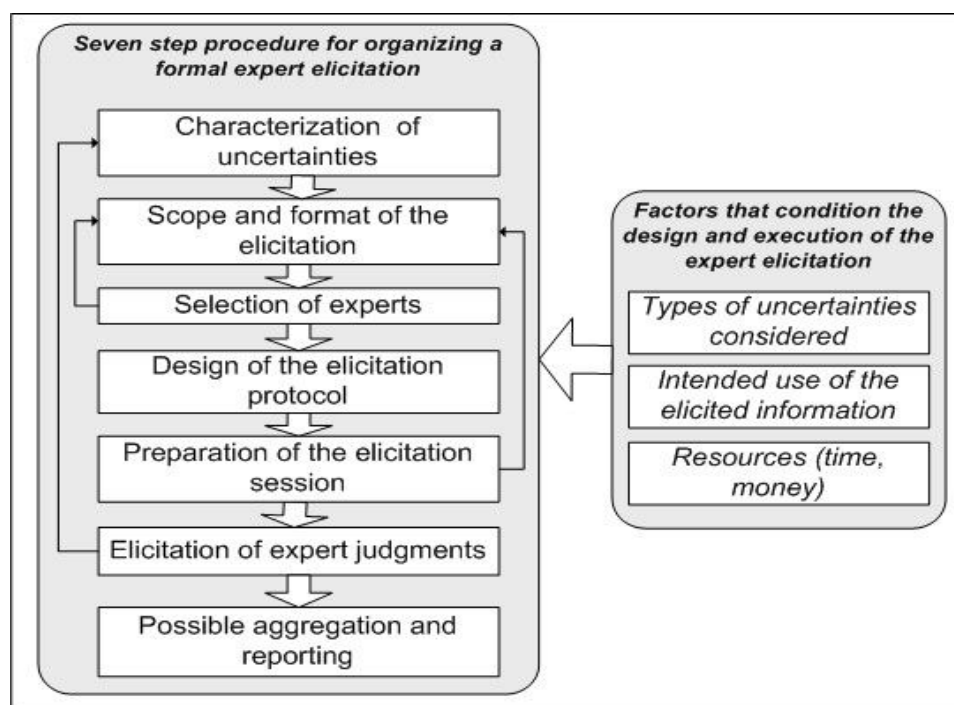
33 There is an increasing interest in SEE, as HTAs are conducted progressively closer to the launch of the
34 intervention of interest.¹³ SEE is also essential for 'early modelling' of new interventions or unknown
35 diseases for which little or no evidence is available.
36
37
38
39

40 No standard guidelines exist to conduct expert elicitation in HTA but there are a number of generic
41 guidance, some of which have been used in HTA.^{14, 15} The most notable of these is the Sheffield
42 Elicitation Framework (SHELF)¹⁵. This is a package of documents, templates and software for eliciting
43 probability distributions. The method begins by eliciting judgements from each expert individually
44 and then elicits a single probability distribution from the group of experts. Cooke's classical method
45 is another generic technique that has been applied in HTA. This method primarily focuses on the
46 synthesis of multiple experts beliefs. Patients are scored based on their performance on calibration
47 questions (questions for which experts do not know true values) and their assessments are weighted
48 according to their scores.¹⁶ The third generic guidance applied in HTA is the DELPHI method. This is
49 an iterative survey that provides feedback from the experts over successive rounds providing an
50 opportunity for consensus as experts review their opinions based on new information from their
51 peers.¹⁷
52
53
54
55
56
57
58
59
60

1
2
3
4
5 While generic processes have been applied in HTA (*Figure 1*), there is an absence of a published
6 guidance that is specific to HTA. Certain elements of the generic guidance may not be appropriate in
7 a HTA context due to resource and time constraints that are inherent in HTA.
8
9

10
11 At present, an analyst needs to be aware of a number of key issues to consider when designing,
12 conducting and analysing an elicitation exercise. In terms of the design, the analyst must decide
13 what quantities to elicit. This will largely be informed by the requirements of the decision model. As
14 a rule, experts should be asked to express their beliefs about observable quantities such as
15 probabilities rather than unobservable quantities i.e. moments of a distribution or covariates. Once
16 the quantities have been chosen, the next choice will be based on which method(s) will be employed
17 to express the parameters. Possible methods include fixed or variable interval methods. The analyst
18 must then choose which experts should be recruited to elicit these judgements. Once the beliefs
19 have been elicited, a decision must be made on how to synthesis the beliefs.
20
21
22
23
24
25
26
27

28 There is heterogeneity in the existing methodology used in HTA. Given the lack of guidance, there is
29 a need to develop a standard set of principles to guide the design and conduct of expert elicitation in
30 HTA. It is essential that the elicited information represent how uncertain experts are about the
31 current state of knowledge regarding a parameter of interest. There is a need to reflect the range of
32 reasonable judgements that may be expressed across experts (between-expert variation) and
33 determine how decision makers use these elicited judgements in the decision making process.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



28 **Figure 1 General schematic for SEE**

29
30
31 The overall aim of this report was to establish a reference protocol or guideline for the elicitation of
32 experts judgements to inform healthcare decision-making (HCDM). To achieve this overall aim, the
33 report will focussed on the following objectives:

- 34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
- Providing clarity on the methods for collecting and using experts judgements within an assessment of cost-effectiveness.
 - Demonstrating where alternative methodology may be required in particular context/constraints e.g. time.
 - Establishing preferred approaches for elicitation for a range of parameters and a range of decision-making contexts.
 - Determining which elicitation methods allow experts express parameter uncertainty as opposed to variability.
 - Determining the applicability and usefulness of the reference protocol developed within a case study application.

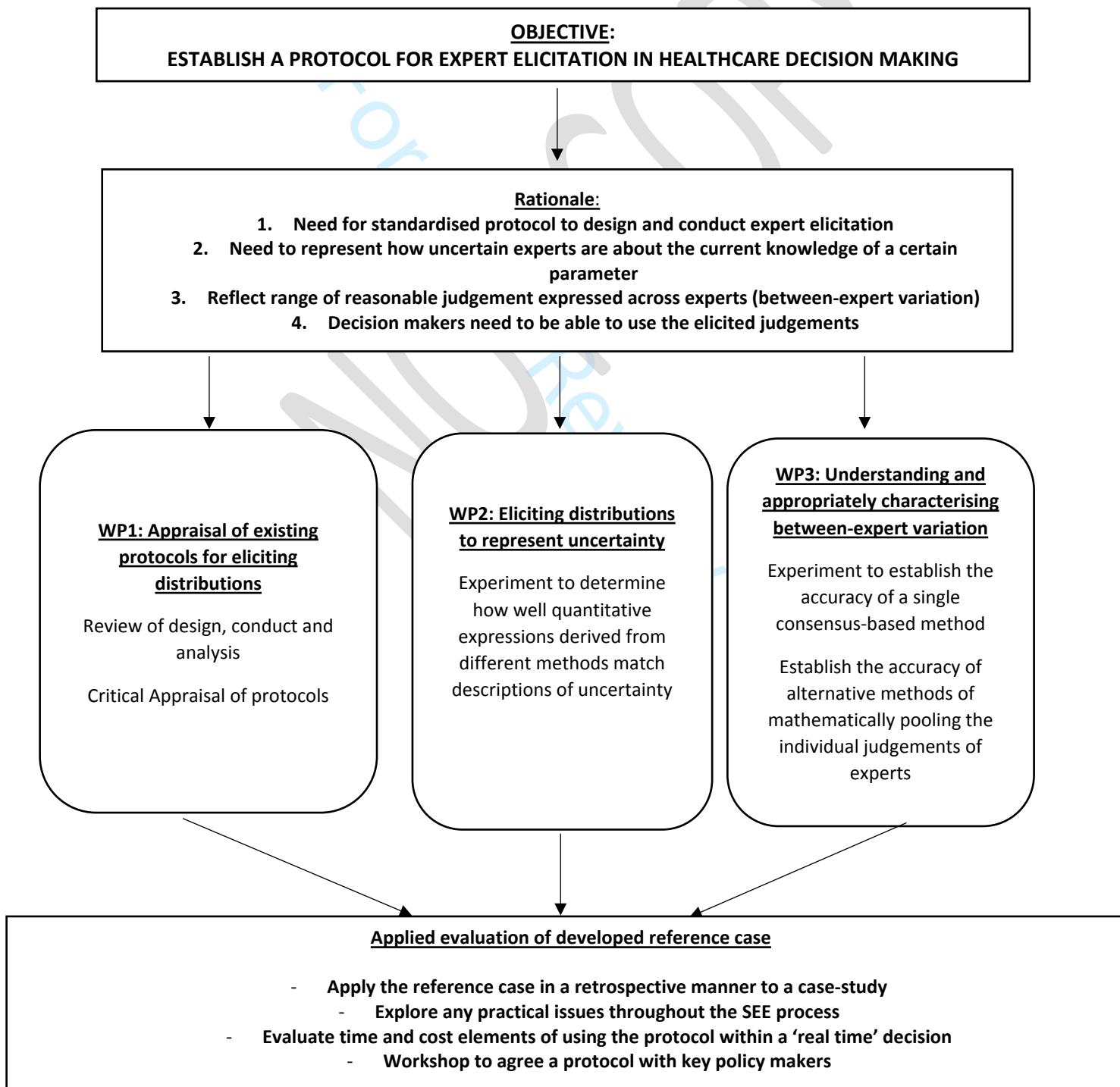
55
56
57
58
59
60

The initial research protocol outlined two additional objectives: to establish the accuracy of consensus-based methods in generating representations of uncertainty, and to establish the accuracy of alternative methods of mathematically pooling the individual judgements of experts. The

objectives were subsequently refined to explore individual factors that can affect the accuracy of consensus-based methods, in particular: to explore individuals' ability to extrapolate from their knowledge base, and to explore how individuals revise their answers when presented with group summaries. Further details on the reason for these deviations is provided in Chapter 8 .

To achieve these objectives, the activities of this project were split into three work packages and an evaluation. The activities of the project are summarised below in Figure 2.

Figure 2 Summary of project activities



1
2
3 Specifically the remaining chapters in this report provide the following:
4
5

6 *Chapter 2* reviews existing guidelines for formal elicitation (SEE). This review identifies the
7 approaches used in existing guidelines and aims to identify whether dominant approaches evolve in
8 terms of the choices that need to be made in the elicitation process.
9
10

11
12
13 In light of this review, *Chapter 3* considers contexts for structured elicitation in HCDM. Different
14 contexts may influence the requirements and feasibilities of expert elicitation and so *Chapter 3*
15 discusses this in detail and identifies the potential constraints in decision-making in healthcare and
16 discusses the implications for expert elicitation methodology.
17
18
19

20
21
22 *Chapter 4* is a review of SEE applications in cost-effectiveness modelling. The chapter summarises
23 the basis for the methodological choices made in each application and details the challenges that
24 were reported by the authors.
25
26
27

28
29
30 *Chapter 5* reviews the evidence on the potential choices that are available for different components
31 of the elicitation process. This focuses on the elements: selection of experts, level of elicitation,
32 fitting and aggregation and adjusting judgements. This chapter discusses the advantages and
33 disadvantages of each available choice and identifies any potential constraints to their application in
34 cost-effectiveness analyses.
35
36
37

38
39
40 Heuristics and biases are concerns that are predominant across all elements in SEE and SEE should
41 be conducted in such a way that minimises these errors. *Chapter 6* reviews the existing evidence on
42 heuristics, biases and de-biasing techniques that are of most relevance to HCDM.
43
44

45
46
47 *Chapter 7* discusses what quantities to elicit. This chapter provides a list of alternative quantities that
48 can be elicited to inform certain types of parameters that are commonly used in healthcare. This is
49 particularly relevant in cost-effectiveness analyses, as parameters are often complex constructs,
50 such as relative treatment effects or time to events, which experts will not directly observe in
51 practice. *Chapter 7* compiles a list of alternative quantities that may be elicited to inform specific
52 parameters.
53
54
55

56
57
58 *Chapter 8* provides the experimental plan for experiments that were conducted as part of this
59 research. The aim of these experiments was threefold: (1) to evaluate alternative methods of
60

1
2
3 elicitation and how they perform in representing parameter uncertainty; (2) to explore individuals'
4 ability to extrapolate from their knowledge base; and (3) to explore how individuals revise their
5 answers when presented with group summaries. The results and interpretation of these
6 experiments is then presented.
7
8
9

10
11 *Chapter 9* discusses the methodological choices for each of the different components of SEE: design,
12 conduct and analysis. Managing biases and validity assessment are then considered as overarching
13 concerns for throughout the SEE process. In order to conclude on their suitability for HCDM, *Chapter*
14 *9* first presents a set of principles that underpin the use of expert elicitation in HCDM. Available
15 choices are considered in light of these principles and any empirical evidence available to support
16 the choices.
17
18
19
20
21

22
23 Chapters 2 to 9 are then used to generate a reference protocol for HCDM (*Chapter 10*). This presents
24 the choices that are supported by the principles for HCDM, and/or empirical evidence in this
25 domain. Given the paucity of empirical evidence relating to HCDM, it was necessary to define this for
26 a specific type of HCDM, HTA. Considerations when using the reference protocol outside of this
27 context are also presented.
28
29
30
31

32
33 *Chapter 11* describes the applied evaluation of the developed reference protocol. This uses an
34 existing cost-effectiveness model, in which SEE was used to generate initial estimates of uncertain
35 parameters. In addition to demonstrating the usefulness of the reference protocol in navigating the
36 SEE process, the practicality of SEE is determined using narrative feedback from experts and by
37 generating estimates of resources required to design and conduct the SEE.
38
39
40
41

42
43 The report closes with discussion and conclusions (*Chapter 12*) based on the findings of this
44 research. The feedback from a dissemination workshop, exploring the usefulness and challenges in
45 using SEE in HCDM are reported. The limitations of the research and areas of further research are
46 also discussed here.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Chapter 2 Good practice in SEE: Learning from the available guidance

Section 2.1 Introduction

Over the last few decades, SEE has been used in areas such as natural hazards, environmental management, food safety, healthcare, security and counterterrorism, economic and geopolitical forecasting, and risk and reliability analysis. All of these areas require consequential decisions be taken in the face of significant uncertainty about future events or scientific knowledge.

How judgements are elicited is critical to the quality of the resulting judgements and hence the ultimate decisions and policies. Methods for SEE should be suitable for specific contexts and understood by content experts to be useful to decision makers. Example applications and recommended practices do exist in certain fields, but the specifics vary.

In developing a reference protocol for SEE specific to the needs of HCDM, the methodological recommendations and choices, which exist in other fields, need to be understood. This chapter surveyed existing best practices for SEE as reflected in published elicitation guidance to identify areas of consensus, places where no consensus exists, and other gaps. Identifying areas of commonality across current guidance can support elicitation practice in areas that lack context-specific guidance such as HCDM. The recommendations and choices for the SEE process identified in this chapter are further explored in *Chapter 5*, *Chapter 6*, *Chapter 7* and *Chapter 8*, and their suitability for HCDM is considered in *Chapter 9*.

Section 2.2 Methods

To identify areas of agreement and disagreement in elicitation practice, both domain-specific and generic elicitation guidelines were systematically reviewed according to the search strategy and screening process detailed in **Error! Reference source not found..** A *SEE guideline* is defined as a document, either peer-reviewed or in the grey literature, that advises on the design, preparation, conduct, and analysis of a structured elicitation exercise. The review focussed on SEE guidelines rather than applications to determine a full list of the possible methodological options rather than relying on the partial reporting available in applications.

To constrain the scope of this review, guidelines needed to concern explicitly probabilistic judgements and offer guidance on more than one stage of the elicitation process. Literature relating to only one element of elicitation is considered in the targeted searches discussed in *Chapter 5* and *Chapter 6*. When the same or similar author lists published multiple guidance documents making

1
2
3 similar recommendations, only one version was included. An extraction template was used to collect
4 information from each guideline. The extracted data was analysed to create an overview of all the
5 stages, elements, and choices involved in an elicitation and to understand where current advice
6 across guidelines conflicts or agrees. Where the guidelines agree, we assume this represented best
7 practice which can be taken forward within the HCDM context, as applicable. Where the guidelines
8 disagree, we sought additional evidence to support the development of a reference protocol for
9 HCDM (*Chapters 3-8*). .

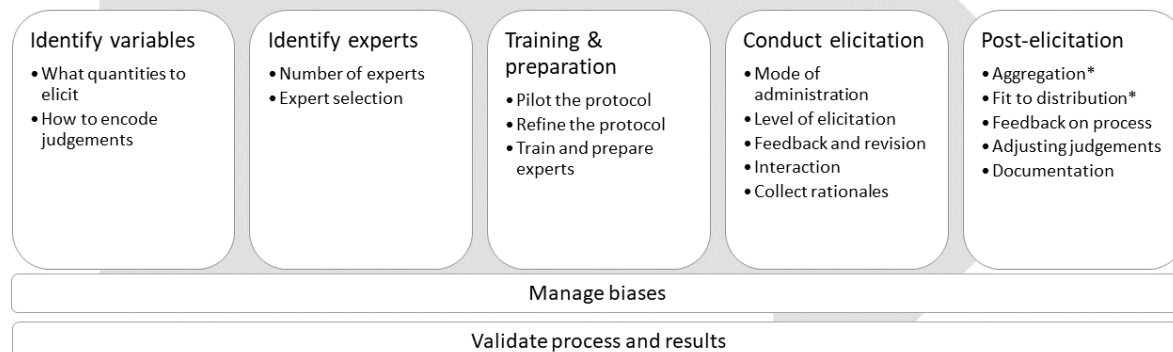
16 **Section 2.3 Included SEE guidelines**

17 The searches identified 16 unique SEE guidelines (**Error! Reference source not found.**, *Table 2*). Five
18 of the guidelines are generic and aim to inform practice across disciplines, and 11 focus on specific
19 domains. Six of the domain-specific guidelines are agency white papers or agency-sponsored peer-
20 reviewed articles and are tailored to the specific decision-making processes the agencies govern.
21 Agencies issuing guidelines include the European Food Safety Authority (EFSA), U.S. Environmental
22 Protection Agency (EPA), the Institute and Faculty of Actuaries (IFA), and the U.S. Nuclear Regulatory
23 Commission (NRC). Both the IFA and NRC have published two distinct guidelines. The 10 guidelines
24 not connected to agencies are based on reviews of existing evidence and practice about elicitation
25 methods (two guidelines), reflections on personal experience and practice (three guidelines), or
26 combinations of review and reflection (five guidelines).

27 Two of the agency SEE guidelines were included with caveats. First, the EFSA guideline covers three
28 distinct elicitation methods, but the Classical Model and SHELF are presented in other guidelines, so
29 only the portions of the EFSA document related to the EFSA Delphi method are included in this
30 review.¹⁷ Second, the EPA guideline is a White Paper released for public review that was not
31 intended to be the final agency report on the subject.¹⁸ However, a final version was never released,
32 and thus the document is widely cited in elicitation literature and has served as a de facto guideline
33 as nothing has superseded it.

36 **Section 2.4 Analysis of the elicitation process**

37 Although the characterization of the process, including the number and categorization of steps,
38 differed among the 16 guidelines, the underlying elicitation process described, depicted in *Figure 3*,
39 was remarkably similar.



* Note: These steps are described as post-elicitation in some guidelines.

1

Figure 3 The elicitation process

At each step of the elicitation process, analysts are faced with a variety of methodological choices. *Table 1* provides the full list of choices described in the 16 guidelines, and *Table 2* summarises the level of agreement in the recommendations and choices discussed for each element. The following sections discuss the variety of methodological recommendations for each stage made across the guidelines. See *Supplementary material 1, Tables 4-15* for further detail.

Table 1 Summary of the elicitation elements, components, and choices described in SEE guidelines

Element	Component	Choices
Identifying elicitation variables		
What quantities to elicit	Type of parameter	<ol style="list-style-type: none"> Elicit observable quantities Elicit required model parameters directly
	Type of quantity	<ol style="list-style-type: none"> Proportions Frequencies Probabilities Odds ratios
	Selection criteria	<ol style="list-style-type: none"> Define selection criteria (probabilities, consequences, constraints, etc) Minimal assessment of each possible uncertain parameter and sensitivity analysis to see which uncertain parameters have the biggest impact
	Principles for describing	<ol style="list-style-type: none"> Ask clear and well-defined questions

	quantities	<ol style="list-style-type: none"> 2. Ask questions in a manner consistent with how experts express their knowledge 3. Uncertainty in the elicited variables should impact the model and/or decision 4. Use neutral wording
	Decomposition/disaggregation	<ol style="list-style-type: none"> 1. Decompose variables of interest to aid experts in the elicitation task 2. Don't decompose variables for the experts
	Handling dependence	<ol style="list-style-type: none"> 1. Express dependent variables in terms of independent variables 2. Use conditional probabilities 3. Use other dependence elicitation methods
Encoding judgements	General approach	<ol style="list-style-type: none"> 1. Fixed interval method: <ul style="list-style-type: none"> • Roulette or chips and bins method • Ask for % that falls within a specific range 2. Variable interval method: <ul style="list-style-type: none"> • Quantiles (Quartiles, Tertiles, 5%, 95% & median, 17%, 83% & median, other) • Bisection • Plausible probabilities (Lowest plausible probability, Highest plausible probability, Best guess for the probability) • Plausible quantities (Upper and lower plausible bounds, best guess, degree of belief) • NUSAP (Numeral, Unit, Spread, Assessment, Pedigree) 3. Hybrid fixed/variable interval methods 4. Summary statistics, moments, measures of central tendency 5. Elicit evidence, not parameter values, and analyst/facilitator defines probability distribution that reflects the body of evidence 6. Other
	Use of visual aids	<ol style="list-style-type: none"> 1. Use to aid elicitation task 2. Do not use
Identifying and selecting experts		
Number of experts	Number of experts	<ol style="list-style-type: none"> 1. Depends on application 1. Options mentioned in different guidelines: about 10; about 5 specialists

		and 2-3 generalists; 10-20; 6-e12; at least 4; 8 a "rule of thumb"; 5-9
Selecting experts	Roles within SEE	<ol style="list-style-type: none"> 1. Facilitator (assessor, analyst, coordinator): prepare and conduct elicitation 2. Expert (technical expert, specialist, subject-matter expert): provide judgments (and/or evidence) 3. Generalists: may provide judgments, advise on design, or help with the elicitation
	Desired characteristics for those providing judgements	<ol style="list-style-type: none"> 1. Normative expertise 2. Substantive expertise 3. Willingness (interest and availability) to participate 4. Ability to understand questions 5. Ability to apply skills 6. Notability
	Identification procedure	<ol style="list-style-type: none"> 1. Recommendations by peers, either formally or informally 2. Research outputs 3. Known experience 4. RFP to seek out experts 5. Profile matrix to identify types of expertise required
	Selection procedure	<ol style="list-style-type: none"> 1. Disclosure of personal and financial interests 2. Pursue diversity in opinions, specialisation, area, institution, etc. 3. Pursue diversity in age, gender, culture 4. Formal selection criteria developed and applied 5. Send potential experts a questionnaire 6. Review CVs of possible experts and have a committee select accordingly 7. Match possible experts against profile matrix
	Possible selection criteria	<ol style="list-style-type: none"> 1. Reputation 2. Experience and qualifications 3. Publication history 4. Diversity in background 5. Conflicts of interest 6. Awards 7. Balancing different viewpoints and managing group dynamics 8. Peer assessment (such as GEM) 9. Convenience 10. Balance of internal and external experts (e.g., include at least 2 external experts)

Training and preparation		
Pilot the protocol	Pilot exercise	<ol style="list-style-type: none"> 1. Pilot 2. No mention of pilot
Training and preparation for experts	What to cover in training	<ol style="list-style-type: none"> 1. Probability, including subjective probability, and related concepts 2. Motivation for elicitation 3. Description of what is required from experts 4. How results will be used 5. Elicitation questions 6. Example and practice questions 7. Review of potential biases 8. Relevant background information, data, and sources 9. Review assumptions and definitions used in the elicitation 10. Description of performance assessment (if relevant) 11. Introduction to dependence (if relevant)
Conducting the elicitation		
Mode of administration	Location	<ol style="list-style-type: none"> 1. Face-to-face <ul style="list-style-type: none"> • 1-on-1 • Group • Plenary 1. Remote (web, mail, email, phone, video conference, etc)
Level of elicitation	Level of elicitation	<ol style="list-style-type: none"> 2. Individual 3. Group 4. Combination (individual assessment followed by group discussion and assessment)
Feedback and revision	Type of feedback	<ol style="list-style-type: none"> 1. Graphical feedback 2. Fitted distributions 3. Written description of the expert's rationale 4. Rationales from other experts 5. Data collected in the future 6. Discussion of elicited values 7. The expert's performance scores 8. Result of using elicited values in the model 9. Decision resulting from the expert judgment 10. Draft elicitation report
	What to feed back	<ol style="list-style-type: none"> 1. The individual's judgments 2. Aggregated group judgments 3. Judgments from other individual experts
	Opportunity for revision	<ol style="list-style-type: none"> 1. Iterate elicitation/feedback rounds

		<ol style="list-style-type: none"> Update after future data is collected Update for revisions/clarifications after circulating draft elicitation report
Interaction	Opportunity for interaction	<ol style="list-style-type: none"> No interaction Group discussion prior to individual elicitation Group discussion and group elicitation Group discussion following individual elicitation (with opportunity for revision) Remote, anonymized interaction
Rationales	Rationales	<ol style="list-style-type: none"> Collect/record rationales from experts (about how they made their judgments) Collect/record rationales from decision makers (about how they used the expert judgments)
Aggregation	Aggregation	<ol style="list-style-type: none"> Aggregate Don't aggregate <ul style="list-style-type: none"> Analyst provides a distribution that captures knowledge from all experts (the Kaplan approach) Only use individual distributions
	Aggregation approach	<ol style="list-style-type: none"> Mathematical <ul style="list-style-type: none"> Opinion pool: equal-weighting, performance-based weighting (with seed questions), analyst-defined weighting (based on rationales, expert qualifications, or other criteria) Bayesian aggregation Behavioural Combination Other
Fit to distribution	Fit	<ol style="list-style-type: none"> Fit to parametric distribution Use non-parametric approaches Don't fit at all
	Distribution	<ol style="list-style-type: none"> Uniform Triangular Uniform over elicited intervals Normal/beta/other parametric distribution
	Fitting method	<ol style="list-style-type: none"> Minimum least squares Method of moments Other
Post-elicitation		
Feedback on process	Feedback from experts on process	<ol style="list-style-type: none"> Get feedback on the procedure if future data collection contradicts elicitation results

		<ol style="list-style-type: none"> 2. Ask experts to appraise the elicitation exercise after completing it
<p>Adjusting judgements</p>	<p>Methods for adjusting judgments</p>	<ol style="list-style-type: none"> 1. Do not adjust experts' assessments 2. Possible adjustments <ul style="list-style-type: none"> • Calibrate experts' assessments • Adjust to improve coherence (described by Lindley et al. (1979)) • Small adjustments allowed, if they are fed back to the experts • Drop an expert from the panel
<p>Documentation</p>	<p>What to include</p>	<ol style="list-style-type: none"> 1. Elicitation questions 2. Responses from individual experts (if elicited) 3. Description of process and assumptions for fitting a distribution 4. Discussion of elicitation procedure (and justification for choices made) 5. Rationales 6. Evidence related to elicited quantities 7. Aggregated judgements and/or consensus curves 8. Discussion of use/impact of elicitation results 9. Recording of session(s) 10. List of experts 11. Definitions and assumptions 12. The process for updating judgments
<p>Managing heuristics and biases</p>		
<p>Managing heuristics and biases</p>	<p>Biases relevant for SEE</p>	<ol style="list-style-type: none"> 1. Cognitive biases <ul style="list-style-type: none"> • Overconfidence • Representativeness • Availability • Anchoring and adjustment • Conservatism • "law of small numbers" • Hindsight bias • Discrepancy between expert's beliefs and responses (conscious or unconscious) • Location errors • Tacit assumptions • Inconsistency 2. Motivational biases <ul style="list-style-type: none"> • Management bias • Expert bias • Social pressure • Group think • Impression management • Wishful thinking

		<ul style="list-style-type: none"> • Misinterpretation • Misrepresentation
	Bias elimination or reduction strategies	<ol style="list-style-type: none"> 1. Give experts practice and feedback 2. Identify biases through discussion with experts 3. Provide training on biases 4. Frame questions to minimize biases and ambiguity 5. Provide relevant background evidence 6. Ask for upper/lower bounds first 7. Ask experts to specify the credible interval they have provided 8. Minimize and record conflicts of interest among the experts 9. Require the experts address conflicting information 10. Collect rationales from experts 11. Report anonymous results 12. Anticipate likely biases 13. Ask experts about evidence, not the probability 14. Avoid numbers in questions
Considering the validity of the process and results		
Validation	Characteristics of validity and supporting actions	<ol style="list-style-type: none"> 1. Faithfully capturing experts' beliefs <ul style="list-style-type: none"> • Provide feedback (graphical feedback often mentioned) • Calibration could be a pragmatic proxy • Test that the question is understood 2. Fitness for purpose 3. Calibration <ul style="list-style-type: none"> • Ask questions with realizations (i.e., seed questions) that allow calibration to be tested 4. Calibration and informativeness scoring on seed questions (i.e., the Classical Model) <ul style="list-style-type: none"> • Score experts according to calibration and informativeness • Use scores as basis for performance-based weights (related to Aggregation choices) • Score both individual experts and combinations of experts 5. Coherence <ul style="list-style-type: none"> • Ask for sets of probabilities that allow coherence to be tested • Overfitting (asking for one more summary than is needed)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

		<ul style="list-style-type: none"> • Ask for rationales from experts <p>6. Consistency</p> <ul style="list-style-type: none"> • Ask for rationales from experts (and check for inconsistencies) • Provide feedback • Derive/give feedback on density function during elicitation • Multiply/integrate decompositions during elicitation • Use different elicitation methods and compare results <p>7. Internal peer review of process and/or results</p> <p>8. External peer review of process and/or results</p>
--	--	---

DO NOT COPY
 For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Table 2 Level of agreement on recommendations and choices in SEE guidelines

Element	Component	Agreement level	Explanation
Identifying elicitation variables			
What quantities to elicit	Type of parameter	Some disagreement	Guidelines agree that observable quantities are preferred, but disagree on whether directly eliciting model parameters is an acceptable choice.
	Type of quantity	Disagreement	Guidelines offer conflicting recommendations on whether eliciting probabilities (compared with other uncertain quantities) is an acceptable choice.
	Selection criteria	Some agreement	Fewer than five guidelines discuss this, but they agree selection criteria should be defined.
	Principles for describing quantities	Some agreement	Some guidelines describe slightly different principles (e.g., asking clear questions, ensuring uncertainty on elicited parameters impacts the final decision or model), but they do not conflict.
	Decomposition	Agreement	The guidelines that discuss decomposing the variables of interest all agree it should be a choice.
	Handling dependence	Some agreement	The guidelines that discuss dependence agree it should be avoided if possible or addressed separately,

			but they discuss a range of methods for considering dependence.
Encoding judgements	General approach	Disagreement	Guidelines recommend and discuss different, conflicting methods for encoding judgements.
	Use of visual aids	Some agreement	Fewer than five guidelines discuss this, but they agree visual aids can be a useful choice.
Identifying and selecting experts			
Number of experts	Number of experts	Agreement	The experts agree that multiple experts are important, with most guidelines recommending around 5-10 experts.
Selecting experts	Roles within SEE	Agreement	The guidelines are very consistent in their description of the roles involved with elicitation.
	Desired characteristics for those provide judgements	Some agreement	Characteristics discussed in the guidelines are largely consistent, aside from differing views on if normative expertise is a requirement or just desired.
	Identification procedure	Some agreement	Recommendations differ but do not conflict across the guidelines. Agency guidelines tend to offer more detail.
	Selection procedure	Some agreement	Recommendations differ but do not conflict across the guidelines. Agency guidelines tend to offer more detail.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

	Possible selection criteria	Some agreement	Recommendations differ but do not conflict across the guidelines.
Training and preparation			
Pilot the protocol	Pilot exercise	Agreement	Almost all guidelines recommend conducting a pilot exercise.
Training and preparation for experts	What to cover in training	Some agreement	The lists of what should be included in training vary across guidelines but do not conflict.
Conducting the elicitation			
Mode of administration	Location	Some agreement	Most guidelines agree that face-to-face administration is preferred, though remote options may be pragmatically useful alternative in some situations.
Level of elicitation	Level of elicitation	Disagreement	Guidelines recommend and discuss conflicting levels of elicitation.
Feedback and revision	Type of feedback	Some agreement	Recommendations differ but do not conflict across the guidelines.
	What to feed back	Some agreement	Recommendations differ but do not conflict across the guidelines.
	Opportunity for revision	Some agreement	Guidelines either recommend revision take place following an elicitation (as part of an iterative process or immediately following the elicitation) or further in

			the future, following a draft report or additional data collection.
Interaction	Opportunity for interaction	Disagreement	Guidelines offer conflicting recommendations about when and how to facilitate interaction between the experts.
Rationales	Rationales	Agreement	Almost all guidelines recommend collecting expert rationales in some form.
Post-elicitation			
Aggregation	Aggregation	Agreement	All guidelines discuss aggregation as a recommendation or valid choice.
	Aggregation approach	Disagreement	Guidelines offer conflicting recommendations on the approach and method to aggregate judgements.
Fit to distribution	Fit	Some disagreement	The guidelines make few recommendations, but their choices differ.
	Distribution	Some agreement	Fewer than five guidelines discuss this, but they generally agree that many parametric distributions could be chosen.
	Fitting method	Some agreement	Fewer than five guidelines discuss this, but they generally agree that choices include minimum least squares and method of moments.
Feedback on process	Feedback from experts on process	Some agreement	Fewer than five guidelines discuss this, and they recommend complementary approaches.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Adjusting judgements	Methods for adjusting judgements	Some disagreement	Fewer than five guidelines discuss this, but they offer different perspectives.
Documentation	What to include	Some agreement	The lists of what should be included in final documentation vary across guidelines but do not conflict.
Managing heuristics and biases			
Managing heuristics and biases	Biases relevant for SEE	Some agreement	The lists of potential biases vary across guidelines but do not conflict.
	Bias elimination or reduction strategies	Some agreement	The list of possible strategies vary across guidelines but do not conflict.
Considering the validity of the process and results			
Validation	Characteristics/measures	Disagreement	The guidelines differ in their definitions of validity and discussion of how the concept can be operationalised in an elicitation.

Section 2.5 Identifying elicitation variables

Section 2.5.1 What quantities to elicit

SEE is often done in areas with many relevant uncertainties, and a decision has to be made about what will be elicited. Only 1 of the 16 guidelines¹⁹ does not provide advice on selecting what quantities to elicit. Recommendations and choices from the other guidelines are summarised in **Error! Reference source not found.**, Table 3.

Five guidelines recommend that elicited variables should be limited to quantities that are, at least in principle, observable. This includes probabilities, which can be conceptualised as frequencies of an event in a sample of data (even if such data may in practice not be directly available to the expert). However, three guidelines argue that elicited quantities can be “unobservable” model parameters, such as odds ratios, provided they are well-defined and understood by the participating experts. Parameters are here described as “unobservable” if they are complex functions of observable data, such as odds ratios. The guidelines list many types of quantities or parameters that can be elicited, including physical quantities, proportions, frequencies, probabilities, and odds ratios. They give few recommendations, though, aside from Cooke and Goossens²⁰ recommending experts not be asked about uncertainty about probabilities, but that questions be reframed as uncertainty about frequencies in a large population. Choy et al. also recommends against eliciting probabilities directly²¹, but two other guidelines list it as a possible choice. *Chapter 7* further considers the possible types of quantities relevant for HCDM.

Three of the guidelines recommend formal processes for selecting what to elicit, and several guidelines describe principles the elicited quantities should adhere to. Principles discussed include that questions should be clear and well-defined with neutral wording, asked in a manner consistent with how experts express their knowledge, and only be elicited when the uncertainty impacts on the final model and/or decision.

Some SEE guidelines describe two issues related to the quantities to elicit: disaggregation and dependence. Five guidelines suggest disaggregating or decomposing a variable makes the questions clearer and the elicitation easier for experts. Five guidelines also discuss the importance of considering dependence between variables. When dependence is discussed, guidelines recommend reframing dependent items in terms of independent variables wherever possible. If dependence cannot be avoided, the elicitation task will be more complicated, but they recommend assessing conditional scenarios or using other elicitation framing and related techniques to estimate dependence.

Section 2.5.2 Encoding judgements

In addition to choosing what questions to put to experts in an elicitation, analysts must also choose how questions will be put to experts. That is, how will experts be asked to assess their uncertainty about the unknown quantities?

Three guidelines--all agency documents--either do not discuss methods for encoding judgements at all^{22, 23} or offer no advice (i.e., neither recommendations nor a list of choices) on the matter.²⁴ **Error! Reference source not found.**, Table 4 summarises the recommendations and choices described by the other 13 guidelines.

Most approaches can be classified as either fixed or variable interval. Fixed interval techniques (discussed in 6 of the 16 guidelines) present experts with a specific set of ranges, and the experts provide the probability the quantify falls within that range. A popular fixed interval technique is the roulette or "chips and bins" method, in which experts construct histograms that represent their beliefs. Variable interval methods (recommended by five guidelines and discussed in another five), in contrast, give the experts set probabilities and ask for the corresponding values. Popular variable interval methods include the bisection and other quantile techniques. These methods are described further in *Chapter 8*.

Two guidelines recommend methods that cannot be classified as either fixed interval or variable interval. The Investigate, Discuss, Estimate and Aggregate (IDEA) protocol utilizes a combination approach, asking experts to provide a minimum, maximum, and best guess for each quantity as well as a "degree of belief" that reflects the probability the true value falls between the minimum and maximum. Experts may all provide assessments for different credible ranges, and the analyst standardizes them to an 80% or 90% credible interval using linear extrapolation.²⁵

Kaplan's method takes a very different approach.¹⁹ Rather than asking experts to encode their beliefs in a way that can be transformed or interpreted as a probability distribution, the method requires that experts only discuss evidence related to the quantity of interest, and then a facilitator creates a probability distribution that reflects the existing evidence and uncertainty.

In addition to the core encoding method, three guidelines also discuss that physical or visual aids can be used by the elicitor(s) to assist with the encoding process.^{18, 26, 27}

1
2
3
4
5 Despite the variety of encoding methods discussed, none of the guidelines present empirical or
6 anecdotal evidence or other justification for their recommendations or choices. *Chapter 8* provides
7 new evidence relating to the choice of encoding method.
8
9

10 11 12 **Section 2.6 Identifying and selecting experts**

13 Recommendations and choices related to identifying and selecting experts are summarised in **Error!**
14 **Reference source not found.**, *Tables 5-6*. Only one guideline does not discuss the number of experts
15 to include in an elicitation.²³ The others either explicitly recommend or imply that judgements will
16 be elicited from multiple experts. The range of how many experts should be included spans from 4
17 experts²⁰ to 20²⁵. The EPA White Paper is the only guideline that gives considerations beyond
18 practical concerns for how many experts to include in an exercise. It observes that if opinions vary
19 widely among experts, more experts may be needed. On the other hand, if the experts in a field are
20 highly dependent (based on similar training or experiences, for example), adding more experts has
21 limited value.¹⁸ The risk of dependence between experts is only discussed in three other
22 guidelines.^{22, 26, 28}
23
24
25
26
27
28
29
30

31
32 Most guidelines do not address how many facilitators or analysts should be involved in an elicitation.
33 The few that do state that 2 or 2-3 facilitators is ideal, with the facilitators having different
34 backgrounds or managing different tasks during the elicitation.^{18, 20, 24, 29, 30}
35
36
37

38 Identifying and selecting experts is discussed in all but three guidelines.^{19, 21, 28} Recommendations
39 from the other 13 guidelines overlap considerably. Common criteria relate to reputation in the field,
40 relevant experience, the number and quality of publications, and the expert's willingness and
41 availability to participate. Normative expertise is listed as desired by five guidelines, but three
42 specify that it is not a requirement.
43
44
45
46
47

48 Five guidelines recommend all potential experts disclose a list of their personal and financial
49 interests, often noting that interests should be recorded but will not automatically disqualify an
50 expert from participating, as that may impose too extreme a limit on the pool of possible experts.
51 Eight guidelines recommend the group of experts included in an elicitation should reflect the
52 diversity of opinions and range of fields relevant to the elicitation topic. The agency guidelines tend
53 to provide more details on identifying and selecting experts, with four describing optional
54 procedures producing a longlist of possible experts that are then winnowed down based on agreed
55
56
57
58
59
60

1
2
3 upon selection criteria. While many guidelines suggest identifying experts through peer nomination,
4 Meyer and Booker³¹ cautions that this process can, if not well-managed, lead to issues related to
5 experts only nominating other people with similar views. *Chapter 5* considers the broader literature
6 on selecting and identifying experts.
7
8
9

10 11 12 **Section 2.7 Training and preparation**

13 Recommendations and choices related to identifying and selecting experts are summarised in **Error!**
14 **Reference source not found.**, *Table 7*. Nine guidelines either explicitly recommended piloting the
15 elicitation protocol with a subject-matter expert not participating in the exercise or imply that
16 piloting will be done. The remaining seven guidelines did not discuss piloting.
17
18
19

20
21 Only one guideline offered training as a choice¹⁹; the other 15 all requiring at least some form of
22 training. Recommendations and suggestions for what should be included in expert training are
23 largely consistent across the guidelines and cover issues related to elicitation generally and the
24 subject-matter at hand specifically. Commonly recommended aspects of training include: an
25 introduction to probability and uncertainty, an overview of the elicitation process, an introduction to
26 heuristics and biases, the aim and motivation for the elicitation, information on how elicitation will
27 be used, relevant background information, and details of any assumptions or definitions used in the
28 elicitation. Four guidelines recommend using practice questions to ensure experts understand the
29 elicitation process.
30
31
32
33
34
35
36
37

38 Most guidelines do not discuss what, if any, training should be provided to the elicitation
39 facilitator(s) or other roles involved in conducting an elicitation. Five guidelines, including four
40 generic guidelines, provide material meant to assist the facilitator, including sample text and
41 forms.^{17, 20, 25, 29, 31}
42
43
44
45
46

47 **Section 2.8 Conducting the elicitation**

48 *Section 2.8.1 Mode and level of elicitation*

49 Recommendations and choices about the mode of administration and level of elicitation (group or
50 individual) are summarised in **Error! Reference source not found.**, *Table 8*.
51
52
53
54

55 Elicitations can be conducted in-person, in either individual interviews or group workshops, or
56 remotely, via the internet, email, mail, phone, video conferencing, or other means. Nine guidelines
57 recommend in-person elicitation, and only one recommends remote elicitation. Eight guidelines list
58
59
60

1
2
3 remote elicitation as a choice, recognizing that it may be logistically easier to arrange than an in-
4 person elicitation.
5
6
7

8 The mode of administration may be governed by whether a method elicits judgements from
9 individual experts (i.e., each expert provides an individual assessment) or groups (i.e., a group of
10 experts provides a single assessment). Of the 16 guidelines, only Choy et al.²¹ does not discuss the
11 level of elicitation. Group-level elicitation is only recommended by Kaplan,¹⁹ who recommends a
12 process where experts discuss the evidence relevant to an elicitation variable and then the facilitator
13 proposes a probability distribution that matches the input provided by all of the experts. Individual-
14 level elicitation is recommended by five guidelines, and two guidelines recommend a combination
15 approach, wherein individual assessments are elicited first and then the group works to provide a
16 communal assessment that reflects the diversity of opinion in the group. *Chapter 5* provides more
17 detail on individual- versus group-level elicitation.
18
19
20
21
22
23
24
25

26 27 *Section 2.8.2 Feedback and revision*

28 All but one guideline³¹ discussed the importance of feedback and revision, but three did not provide
29 information on how it should be done.^{23, 26, 27} The other guidelines discuss a range of possible
30 feedback methods, which can provide information on an individual's judgements, the aggregated
31 group judgements, or a summary of what the other experts provided. Recommendations and
32 choices about the mode of administration and level of elicitation are summarised in the **Error!**
33 **Reference source not found.**, *Table 9*.
34
35
36
37
38
39

40 Knol et al²⁷ is the only guideline that warned of a possible negative impact of feedback and revision,
41 cautioning that it can cause unwanted regression to the mean in the experts revised assessments.
42 None of the guidelines recommend against providing feedback and opportunities for revision in any
43 form. The feedback of group summary judgements is investigated in *Chapter 8*.
44
45
46
47
48

49 *Section 2.8.3 Interaction*

50 Recommendations and choices regarding interaction and rationales are summarised in the **Error!**
51 **Reference source not found.**, *Table 10*. Three guidelines did not explicitly discuss interaction
52 between the experts.^{20, 21, 32} While no guidelines recommended avoiding interaction, seven say no
53 interaction is a possible choice. Interaction is closely related to level of elicitation, with guidelines
54 recommending group discussion prior to individual elicitation, group discussion prior and during a
55 group elicitation, and group discussion following an individual elicitation. One guideline
56
57
58
59
60

1
2
3 recommended interaction be limited to a remote, anonymous, facilitated process.¹⁷ Other guidelines
4 also described these options as choices.
5
6
7

8 Although the guidelines disagreed about if and how interaction should be managed in an elicitation,
9 many do present more justification for the recommendations or choices around interaction than
10 they do for other methodological choices. The benefits of interaction between experts is that it
11 minimizes the differences assessments that are due to different information or interpretation²⁷ and
12 allows analysts to explore correlation between experts.²⁸ The drawbacks, though, are that it can
13 allow strong personalities to carry too much weight²⁶⁻²⁸, the experts may feel pressure to reach a
14 consensus²⁶, there may be risk of confrontation²⁸, and interaction can encourage groupthink,
15 resulting in the experts being overconfident.²³ Practical considerations can also guide the choice of if
16 and how to include interaction, as individual interviews may take more time, but a group workshop
17 may be more expensive.²⁷ These issues are further discussed in *Chapter 5*.
18
19
20
21
22
23
24
25

26 *Section 2.8.4 Rationales*

27 Only one guideline presented collecting the experts rationales during an elicitation as a choice rather
28 than a recommendation.³¹ The other 15 guidelines all recommend collecting rationales because they
29 help analysts and decision makers understand what an answer is based on^{23, 26, 28}, provide a check of
30 the internal consistency of an expert's responses²⁶, record any assumptions³⁰, and may help limit
31 biases.²¹ The information collected in rationales can also be useful for peer review or for future
32 updating of the judgements.²³
33
34
35
36
37
38
39

40 One guideline also recommended collecting rationales from the decision maker about how they use
41 the expert judgement results.³²
42
43
44

45 **Section 2.9 Post-elicitation**

46 *Section 2.9.1 Aggregation*

47 Even when eliciting judgements from multiple experts, it can be important to have a single
48 distribution that reflects the beliefs of the experts that can be used in modelling. Recommendations
49 and choices on aggregation methods are summarised in ***Error! Reference source not found., Table***
50 ***11***. Five guidelines presented aggregation as a choice, but the remaining 11 recommended
51 aggregation always be done.
52
53
54
55
56
57
58
59
60

1
2
3 Aggregation can be behavioural or mathematical. In behavioural aggregation, experts interact with
4 the goal of producing a single, consensus distribution. Mathematical aggregation involves the
5 facilitator(s) eliciting individual assessments from the experts and then combining them into a single
6 distribution through a mathematical process. Two guidelines recommend behavioural aggregation.
7
8 Kaplan¹⁹ recommends a process that includes group-level elicitation and behavioural aggregation:
9 the experts discuss the evidence relevant to an elicitation variable, the facilitator suggests a
10 probability distribution that reflects the diversity of evidence on the subject, and then the process
11 concludes when there is consensus from the experts about the proposed distribution. The SHELF
12 method recommends an initial round of individual-level elicitations followed by expert discussion
13 designed to produce a single distribution that represents how a "rational independent observer"
14 would summarise the range of expert opinions.²⁹

23 Four of the guidelines recommended variations on mathematical aggregation. Three guidelines
24 recommended combining expert judgements in a linear opinion pool that equally weights all of the
25 experts. Cooke and Goossens²⁰ is the only guideline that recommends mathematical aggregation
26 with differential weights for the experts. They suggested a method whereby the experts are scored
27 and weighted according to their performance assessing a set of seed questions, which are items that
28 are unknown to the expert but known to the facilitator.

35 Budnitz et al.³³ recommend a unique approach wherein the analysts determine the aggregation
36 method during an elicitation, based on an evaluation of how the process is unfolding and
37 determining what is most appropriate. They recommend a behavioural aggregation-based consensus
38 is the best choice, but believe it is not appropriate in all situations. The analysts can also decide to
39 use mathematical aggregation with equal weights or analyst-determined weights or a process similar
40 to that recommended by Kaplan,¹⁹ in which the analysts supply a distribution they believe captures
41 the discussion and evidence presented by the experts.

48 Like interaction, several of the guidelines give more background to help guide an analyst in his or her
49 choice of method. The main drawback of aggregation, according to Tredger et al.,²³ is that it can lead
50 to a result that no one believes. Two guidelines warn that the expert selection is of increased
51 importance if an elicitation will use mathematical aggregation with an opinion pool, particularly
52 equal weights, as increasing the number of experts with similar beliefs will result in those beliefs
53 having more influence in the final, aggregated distribution.^{24, 26} Garthwaite et al.²⁶ also suggest
54 opinion pools may be problematic as they result does not represent any one person or group's
55
56
57
58
59
60

1
2
3 opinion, but Bayesian weighting requires a lot of information on the decision maker's views of the
4 experts opinions. Finally, several guidelines discuss that the possible issues around behavioural
5 aggregation are linked to the challenge of properly managing group interactions, the topic discussed
6 next. The broader literature on aggregation is discussed in *Chapter 5*.
7
8
9

10 11 *Section 2.9.2 Fit to distribution*

12 Recommendations and choices on fitting to distribution are summarised in the **Error! Reference**
13 **source not found.**, *Table 12*. Analysts can fit the elicited data to a probability distribution either as
14 part of the elicitation or during post-elicitation analysis of the data. Possible choices, discussed in
15 about half of the guidelines, include fitting to a parametric distribution, using non-parametric
16 approaches, or just use the information directly elicited from the experts.
17
18
19
20
21
22

23 None of the guidelines recommended specific distributions be used in fitting, but they say the
24 analysts should choose based on the nature of the elicited quantity and the information provided by
25 the experts. Cooke and Goossens²⁰ describe probabilistic inversion, a method that can be done if
26 the observable elicited variable needs to be transformed into a distribution on an unobservable
27 model parameter. *Chapter 5* explores issues of fitting judgements to distributions in more detail.
28
29
30
31
32

33 *Section 2.9.3 Other post-elicitation components*

34 Recommendations and choices related to the other post-elicitation components are summarised in
35 the **Error! Reference source not found.**, *Table 13*. Only 2 guidelines discussed obtaining feedback
36 from the experts on the elicitation process. Walls and Quigley recommended analysts ask experts
37 what could have been done differently if new data is later collected that differs from the experts
38 judgements.²⁸ The EFSA Delphi recommended that analysts give experts a questionnaire with the
39 opportunity to provide general comments on the elicitation questions and process.
40
41
42
43
44
45
46

47 None of the guidelines recommended that analysts should adjust experts assessments, but five
48 describe related choices, such as manually adjusting assessments, dropping an expert from the
49 panel, or adjusting assessments to be more accurate, which is recommended against by two
50 guidelines.
51
52
53

54 Documenting the elicitation process and results is the only elicitation element discussed by all 16
55 guidelines. Although the specific recommendations regarding what to include in the final
56 documentation varies across the guidelines, they do not conflict. The guidelines typically
57
58
59
60

1
2
3 recommend documentation include the elicitation questions, experts individual (if elicited) and
4 aggregated responses, experts rationales, and a detailed description of the procedures and design of
5 the elicitation, including the reasoning behind any methodological decision. Many of the agency
6 guidelines are more prescriptive about what documentation should entail, and some provide
7 detailed templates.^{17, 18, 32}

13 **Section 2.10 Managing heuristics and biases**

15 Expert judgements are affected by a variety of heuristics and biases.^{34, 35} Morgan³⁶ argues that these
16 biases cannot be completely eliminated, but that the elicitation process is designed to minimize their
17 influence on the results. The 16 reviewed guidelines discussed 11 different cognitive biases and 8
18 motivational biases that can affect an elicitation. A list of the biases discussed and possible actions to
19 minimize them can be found in **Error! Reference source not found.**, Table 14.

25 Most of the bias-reducing actions mentioned by SEE guidelines are only discussed in one or two, but
26 the actions do not conflict with one another. The most frequently recommended actions are to
27 frame questions in a way that minimizes biases (discussed in five guidelines) and to ask for the upper
28 and lower bound first, to avoid anchoring (discussed in three guidelines). While most guidelines offer
29 some recommendations for mitigating and managing biases, they present little to no empirical
30 evidence to support that their recommended actions have the intended effect. The broader
31 literature on heuristics and biases is reviewed in Chapter 6.

39 **Section 2.11 Considering the validity of the process and results**

41 Four guidelines do not discuss how to ensure the validity of elicited results^{17, 19, 29, 31}, and the other 12
42 guidelines present a range of perspectives on what is meant by validity, summarised in **Error!**
43 **Reference source not found.**, Table 15. Validity can mean that the exercise captured what the
44 experts believe (even if that is later proven false).²⁶ It can also refer to whether the expressed
45 quantities correspond to reality^{20, 25, 26, 28}, are consistent with the laws of probability^{26, 28}, or are
46 internally consistent.^{22, 27} Some guidelines—all agency documents—also view validity as mostly
47 concerned with the process rather than the results and suggest an elicitation is valid if it has been
48 subjected to peer review.^{18, 24, 32} Recommendations and choices for handling validity differ across the
49 guidelines and can involve actions at any stage of the elicitation process, depending on what
50 definition of validity the guideline seeks to achieve.

Section 2.12 Conclusions

The SEE guideline review reveals a developing body of work designed to guide elicitation practice. Although the guidelines evolved separately in different fields, they largely agree on issues around what quantities to elicit, expert selection, the importance of piloting the exercise and training experts, face-to-face elicitation being preferable to remote modes, the importance of collecting rationales from the experts alongside the quantitative assessments, fitting assessments to distributions, the key role documentation plays in supporting and communicating an elicitation exercise, and how to manage heuristics and biases. The guidelines recommend different approaches for encoding judgements, using individual- or group-level elicitation, aggregating judgements, and managing interaction between the experts. Although the guidelines agree that validation is important, they disagree on what actions an analyst can take to encourage or demonstrate validity. Finally, some areas seem under-discussed. Dependence between questions, for example, is a complicated issue that could be critically important when interpreting elicitation results, but little guidance exists on the topic.

The elicitation choices identified in this review are further considered in *Chapters Chapter 5, Chapter 6, Chapter 7 and Chapter 8*, and their suitability for use in the HCDM context is evaluated in *Chapter 9*.

Chapter 3 Expert elicitation in different decision making contexts

Section 3.1 Introduction

Challenges in the conduct of SEE in HCDM are discussed in *Chapter 4*. The challenges identified in the applied examples were largely practical and related to the design of the SEE for that particular task. There are, however, much broader challenges and opportunities, which relate to the decision making context in which SEE is applied. These issues are discussed in this chapter.

The specificities of the context in which expert elicitation is conducted should be distinct from the principles and methods employed. That is, best practice should always be regarded as an appropriate starting point regardless of the context. Doing so, however, may ignore many important factors which influence the choice of method employed for an elicitation. A reference protocol which does not at least consider context specific constraints, is unlikely to be used widely or may be only restricted to a sub-set of decision makers, perhaps those operating at a national level. In considering how a reference protocol for expert elicitation in HCDM might be utilised in practice, it is important to understand how different decision making contexts may influence the requirements for and practicalities of, expert elicitation. In particular there may be practical constraints in certain contexts that imply the use of a second best methodology. Some of these issues are explored in the evaluation (see *Chapter 11*), however this chapter considers the range of decision making contexts more generally, highlights the potential constraints and the implications for SEE methodology. Given the lack of experience with SEE in formal decision making processes, a formally review of the challenges and constraints faced by different HCDM's is unlikely to be informative. Instead this chapter is intended as a discussion, rather than a formal review. It draws on observations and experiences of the project team and the wider advisory group.

Section 3.2 Levels of decision making

In England reimbursement decision making bodies can be described at three levels, implying the population they serve and the jurisdiction for their decision making activities.³⁷ These are:

1. Individual practitioners, such as General Practitioners (GPs), secondary care clinicians and local decision makers, such as Clinical Commissioning Groups (CCGs), Local Authorities (LAs) and hospital trusts
2. National decision makers, such as NICE, Department of Health and Social Care (DHSC), NHS England, Public Health England

- 1
 - 2
 - 3
 - 4
 - 5
 - 6
 - 7
 - 8
 - 9
 - 10
 - 11
 - 12
 - 13
 - 14
 - 15
 - 16
 - 17
 - 18
 - 19
 - 20
 - 21
 - 22
 - 23
 - 24
 - 25
 - 26
 - 27
 - 28
 - 29
 - 30
 - 31
 - 32
 - 33
 - 34
 - 35
 - 36
 - 37
 - 38
 - 39
 - 40
 - 41
 - 42
 - 43
 - 44
 - 45
 - 46
 - 47
 - 48
 - 49
 - 50
 - 51
 - 52
 - 53
 - 54
 - 55
 - 56
 - 57
 - 58
 - 59
 - 60
3. Research commissioners, including organisations such as NIHR and the Medical Research Council (MRC) but also industry sponsors of research.

Reimbursement bodies range from local practitioners and commissioners to national decision makers (2). Here individual and local decision makers (1) are grouped together as many of the constraints are relevant in both contexts. In addition, there are multiple organisations which commission research (3), potentially including SEE, these can also be regarded as decision makers.

Section 3.2.1 Individual practitioners and local 'population level' decision makers

Section 3.2.1.1 Features

There are a number of decisions that are made on an individual practitioner – patient level in the NHS and other health care systems. These usually concern a patient's course of treatment and the most effective, and sometimes cost-effective, choice given the particular circumstances. Such decisions are made in a primary care setting, usually involving a GP, and a secondary care setting, usually involving a consultant or other medical specialist. Such decision makers may also make choices for groups of patients, for example in deciding which device to purchase within a hospital or organising surgical lists.

HCDM occurs at a population level in several forms. In England within primary care, a CCG (see below) supports decision making between GPs and their patients through local guidance such as the referral support system (RSS)[see ³⁸ for an example]. This may also extend to services offered within secondary care, e.g. referrals for further testing/investigations. For both primary and secondary care there may be relevant guidance produced by NICE to support decision making. Individual practitioners and secondary care clinicians are also influenced by their professional bodies/councils. NHS England, as well as commissioning for primary care, also produces key strategic guidance for CCGs to support them to fulfil their duties to their respective populations. So, although individual health care professionals in the NHS make decisions about individual patients, this is very much governed by the organisations that are intended to harmonise provision of services across England and encourage best and most cost-effective provision of services.

The NHS also works closely with social services, and individual practitioners in this respect include social workers, residential care homes and carers. Within public health, services are delivered across the NHS, social care and local authorities (LAs) and many are supported by the work of Public Health England (PHE). Public health practitioners, including nutritionists, smoking cessation coordinators,

1
2
3 and teenage pregnancy coordinators, are again employed by the NHS, PHE, NHS England and
4 individual LAs and CCGs, and as such work within their codes of practice and adhere to appropriate
5 guidance regarding provision of services. LA Public Health is also accountable to Public Health
6 England.³⁷
7
8
9

10 11 12 *Section 3.2.1.2 Constraints*

13 The particular constraints in these context relate to the degree of autonomy that individual
14 practitioners/CCGs have in making decisions regarding individual patients or groups of patients
15 within their jurisdiction. Since the abolition of GP fund holding (in 1997/98)³⁹ and subsequent
16 changes to commissioning before the 2014 care act, individual practitioners are more constrained
17 with regards to patient level decision making.
18
19
20
21
22

23 In both CCGs and LAs there are significant budget constraints. Whilst the average CCG's budget grew
24 by 3.4% in 2016/2017⁴⁰ there are a number of new pressures on CCGs which require them to
25 cutback or reorganise local services. These include cutbacks to public health and social care funding.
26 Post the first two years of the move of public health to LAs there are services which LAs may be
27 forced to reduce investment in, some of which have implications for public health.⁴¹
28
29
30
31
32

33 CCGs and LA are also constrained by budget cycles, which are typically 1-3 years. There may be an
34 incentive to replace activities that cannot prove 'value' within these time frames with those that
35 have a higher immediate payoff, for example less investment on prevention.
36
37
38
39

40 Both CCGs and LAs face multiple competing demands for money and resources. There are big
41 differences across regions with regards to commissioning of and participation in research. Some
42 CCGs and LAs work with health economists and are therefore directly involved in either
43 commissioning participating in or understanding the results of cost-effectiveness evidence and the
44 implications for their population. Others do not have access to such resources.
45
46
47
48
49

50 51 *Section 3.2.1.3 Implications for expert elicitation in this context*

52 We are not aware of any examples where formal SEE has been used to support decision making of
53 the system, although of course, judgement is used routinely in the clinical and management settings
54 every day. (This does not preclude that such elicitations have not been done, of course, but if so,
55 they have not been documented.) Indeed, it might initially seem practically unfeasible to use SEE to
56
57
58
59
60

1
2
3 support decision-making at the individual practitioner level, particularly given that many decisions
4 are made on a national level with implementation at a local level. However, there are still a number
5 of decisions that can be made by individual practitioners or groups of practitioners/local
6 commissioners, many of which may rely on assumptions and opinion rather than experimental data.
7
8 For example in considering individual cases and episodes of care, such as for procedures and services
9 not routinely funded by the NHS, an individual funding request (IFR) panel, will consider specific
10 cases for reimbursement, for example cosmetic services.⁴² Those conducting SEE to inform other
11 decision making process may also rely on individual practitioners to act as experts. In some
12 circumstances the SEE may be required to consider parameters for a specific patents rather than at a
13 population level, for example in the IFR process. This can have implications for how a SEE is
14 designed, specifically elicitation of uncertain and communication with experts on how to express
15 their uncertainty.
16
17
18
19
20
21
22

23
24
25 SEE undertaken in this context must also adapt to the practical constraints, in particular it may not
26 be possible to invest significant amounts of time and resource into SEE and availability of experts to
27 inform often practice level decision making maybe be limited. Such experts are unlikely to possess
28 any normative skills or have any experience with SEE. Group based SEE may be a challenge in this
29 context, as well as individual SEE which requires face-to-face interaction. It may be necessary to
30 trade off recruiting large numbers of experts for face-to-face SEE with obtaining larger numbers
31 through a remote SEE.
32
33
34
35
36
37
38

39 *Section 3.2.2 National decision makers*

40 *Section 3.2.2.1 Features*

41 In England the DHSC governs health and social care matters and has responsibility for some
42 elements not covered separately by Scottish, Welsh or Northern Irish governments.³⁷ The DSHC
43 itself takes responsibility for a number of services and activities provided by the NHS, and are also
44 supported by a number of agencies and public bodies. The DSHC provides a mandate to NHS England
45 to help guide its decisions regarding the allocation of resources, commission specialist services and
46 its strategic direction. NHS England oversees commissioning aided by four regional offices. It has
47 responsibility for commissioning contracts for GPs, pharmacists, and dentists and supporting CCGs in
48 their commissioning roles.
49
50
51
52
53
54
55
56

57 There are a number of special health authorities and other bodies which are either part of the NHS
58 or closely associated with it. They include the NICE and the Prescription Pricing Authority. These
59
60

1
2
3 organisations are either accountable to the Secretary of State, or have formal agreements with the
4 DHSC. In general they provide national services. NICE was set up in 1999 as a special health
5 authority.⁴³ Officially NICE only has jurisdiction in England and is supported in considering its
6 guidance in Scotland and Wales by the Scottish Medicines Consortium (SMC) and the
7 All Wales Medicines Strategy Group (AWMSG) respectively. NICE provides guidance on a range of
8 healthcare products and services, including: pharmaceuticals, diagnostics, medical devices and
9 public health interventions. In compiling evidence to generate these guidance it often relies on the
10 use of expert opinion in some form. A review of practices relating to the use of evidence elicited
11 from experts across NICE guidance-making programmes, was recently published⁴⁴ and concluded
12 that “NICE uses expert judgement across all its guidance-making programmes, but its uses vary
13 considerably”. In addition it agreed that “there is no currently available tool for expert elicitation
14 suitable for use by NICE”.

15
16 Working alongside NICE on public health issues are PHE, which was formed in 2013 and took over
17 the role of a number of other health bodies, including the Health Protection Agency (HPA).⁴⁵ PHE
18 generates and interprets evidence and therefore there is potential for it to utilise SEE. Like the Public
19 Health Programme at NICE, the evidence base it considers is more likely to be low quality and/or
20 sparse, and therefore the opportunities for SEE may be significant

21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 *Section 3.2.2.2 Constraints*

37 The likes of the DSHC, NICE and PHE are required to make decisions about reimbursement, best
38 practice and access across the whole of their population. Therefore decisions have to be relevant
39 across different, perhaps heterogeneous, populations.

40
41 The separation of research commissioning and reimbursement can also generate complexities.
42 Decisions may be reached on the basis that further data collection may be required, however some
43 national decision makers do not commission their own research and therefore cannot ensure that
44 data collection takes place and/or addresses the uncertainties identified.

45
46 As with more regional decision making, national decision making is also subject to the constraints of
47 time and resources. Although not necessarily as constrained as local commissioning cycles dictate,
48 national decision makers do still have to generate guidance within acceptable timescales. The
49 process of generating guidance through the NICE single technology appraisal process (STA)⁴⁵ can
50 take around 6-months including committee meetings. Despite the fairly rapid timescales, formal
51
52
53
54
55
56
57
58
59
60

1
2
3 decision making processes, particularly those which imply mandatory implementation of guidance,
4 such as the NICE technology appraisals process, require full accountability for decisions reached. The
5 need to make decisions in a timely manner therefore cannot compromise the quality of the
6 deliberations used to make these decisions, including any evidence generation that contributes
7 towards this.
8
9
10

11 12 13 14 *Section 3.2.2.3 Implications for expert elicitation in this context*

15
16
17 Historically, SEE has been commissioned to support policy challenges. For example, policy on
18 surgical equipment sterilisation to protect against the risk of new variant Creutzfeld-Jakob disease
19 (nvCJD) prion transfer has been informed by SEE in the wake of the Bovine Spongiform
20 Encephalopathy (BSE) crisis in the UK.⁴⁶ More recently, the European Commission commissioned SEE
21 studies of the future antibiotic resistance (AMR) rates in four European Countries to inform policy
22 and the UK Department of Health and Social Care is currently commissioning additional UK-focussed
23 work in this area.⁴⁷ Across national decision makers, the quality of evidence to inform decisions is
24 quite heterogeneous. This can be at various stages of maturity and in some areas, for example public
25 health, evidence may not be particularly robust. SEE could be useful to help inform decisions in
26 these situations, although it is likely that some of the parameters required may also be difficult for
27 experts to make judgements about, for example population uptake of a screening programme.
28
29
30
31
32
33
34
35

36
37 Indeed many examples of SEE conducted in the area of HCDM have been undertaken to inform
38 national decision making organisations such as NICE (see *Chapter 4*). As a result there is a degree of
39 familiarity with the approaches used and an acceptance of its limitations. NICE only makes brief
40 reference to the use of expert opinion to generate evidence in its guide to the methods of
41 technology appraisal.⁴⁵ NICE do not suggest a preferred methodology for this and they have not used
42 any consistent criteria to judge SEE submitted as part of any appraisal process.
43
44
45
46
47
48

49 It is true that decision makers have differing capacities to undertake SEE, specifically in reference to
50 resourcing of SEE. Evidence generation does not constitute a significant proportion of the remit for
51 some decision makers. Therefore, similarly to the use of SEE in local decision making, SEE
52 undertaken in this context must adapt to the practical constraints. Timescales for evaluation are
53 often tight and there are implications of any delay in approving a technology or service. Although
54 SEE takes significantly less time than many other forms of empirical evidence to collect, if conducted
55 appropriately the time resource can still be unachievable in some instances. Political cycles can
56
57
58
59
60

1
2
3 generate promises around improving efficiency and accesses to NHS services. Tight turnaround for
4 evidence to support these promises can negate the ability to undertake SEE and in this instance less
5 formal approaches to filling data gaps may be employed.
6
7
8
9

10 In terms of specifics, as discussed above, decisions may have to be relevant for potentially
11 heterogeneous populations. Eliciting uncertainty around a measure of central tendency across a
12 heterogeneous population, can be a challenge for experts. Rather than eliciting across the entire
13 population it may be advantageous to express quantities for multiple patient types, which will
14 increase the size of the SEE task.
15
16
17
18
19

20 *Section 3.2.3 Research commissioners*

21 *Section 3.2.3.1 Features*

22 In addition to those discussed above, there are also other decision makers not concerned with
23 reimbursement, such as HTA, NIHR more generally and the MRC. These bodies commission research
24 and use expert opinion in cost-effectiveness analyses and therefore any guidance on appropriate
25 design and conduct of SEE would have implications for their practices. Industry can also commission
26 research as part of the license and reimbursement processes.
27
28
29
30
31
32

33 Such decision makers, typically do not fund interventions *per se*. but instead commission
34 effectiveness and cost-effectiveness research across their respective areas of interest. Many of these
35 could potentially use SEE to help inform their decisions regarding which research to fund and the
36 specific form that this research might take. One example is the use of SEE in determining sample size
37 calculations for clinical studies.⁴⁸ Here the SHELF framework¹⁵ has been used to generate prior
38 beliefs to aid clinical study design, specifically on the probability of success (assurance parameter).
39
40
41
42
43
44
45

46 *Section 3.2.3.2 Constraints*

47 The scale of the commissioning of research varies across funders and within their programmes of
48 work. Some funders are constrained to commission research with a specific area, for example clinical
49 speciality, whereas other, such as the National Institute for Health Research (NIHR) and the MRC
50 commission across a range of topic areas. SEE used outside of the context of a decision making
51 (reimbursement) process may not be subject to the same constraints in terms of time or resources,
52 however *Chapter 4* does not identify any applied examples where SEE has been the sole purpose of
53 the research, instead SEE is likely to constitute only a small proportion of the research funding.
54
55
56
57
58
59
60

Section 3.2.3.3 Implications for expert elicitation in this context

As the rationale for the research is to reduce uncertainty, and as research priorities are inevitably contentious there seems to be a very strong case for using SEE in this context to focus research. For example, Dallow et al discuss several examples of the use of expert elicitation at GlaxoSmithcline to inform trial design and the management of the company's research portfolio.⁴⁸ In a similar vein, Walley et al describe a case study of Pfizer where elicitation was used.⁴⁹ Given that research commissioners tend to focus on particular specialities, for example clinical areas, it may also be possible to generate a level of expertise to undertake SEE, both in terms of the analyst and the experts. Where it has been used in the clinical trial setting, to inform sample size calculations, an expert panel has been established to speed up the generation of experts priors in subsequent SEE's.

The lack of consistency between research commissioners presents a challenge for the application of SEE in this context. Not all commission cost-effectiveness studies and there is also diversity in topics which may have implications for the way in which SEE is conducted. Public Health and complex interventions, for example vaccinations programs or other non-pharmacological interventions such as service changes, may imply different methods for SEE compared to medicines (see *Chapter 10*).

Section 3.3 Conclusions

SEE can, in principle, be applied in many different settings, across a range of types of decision makers. In practice, to date, its application has largely been restricted to informing national level HTA decisions and for the purposes of generating evidence as part of larger research projects (see *Chapter 4*). The lack of SEE at an individual practitioner and local population level is likely to be driven by resource and time constraints and the fact that constant changes to policy making at local/national level can also shift the focus on a frequent basis. One solution is to move away from the use of SEE as an 'addition to the analysts' toolkit' and instead as a substitutive for other forms of evidence, for example, a systematic review or modelling exercise. This is likely to be a challenge in systems that have relied heavily on such forms of evidence to inform decision making, but may be more feasible in local decision-making settings.

Guidance on appropriate conduct of SEE in HCDM is likely to be useful in all the contexts discussed, however time constraints and lack of capacity to conduct such exercises is likely to remain a challenge, where SEE is forced to fit into existing processes. For this reason SEE is most likely to gain traction in national and multinational settings, where a capacity for such activities can be generated, simply through economies of scale.

1
2
3 Even within national and multi-national decisions making process, there are likely to be different
4 challenges in conducting SEE and some of these may imply that methodological choices need to be
5 adapted to suit that particular application. Such issues are discussed in *Chapter 10*.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Chapter 4 Challenges in structured elicitation in HCDM

Section 4.1 Introduction

Reimbursement decisions in health are often supported by model based economic evaluation (MBEE).⁵⁰ There may be circumstances in which SEE is required to address data limitations in MBEEs, such as short time horizon or missing entirely.

A review of applications in this area, published in 2013,⁴⁴ identified only a small number (14) of studies reporting the use of SEE. This review did not seek to determine the reasons for heterogeneity of approach, nor did it look at the challenges faced when conducting SEE to support MBEE in health, and inform directions for future research. In pursuit of further clarity, instead focusses on summarising the basis for methodological choices made in each application (design, conduct, and analysis) and the difficulties and challenges reported by the authors. Further details of this review are reported elsewhere¹ and so only a summary has been presented in this chapter.

Section 4.2 Methods

To identify applications of SEE, the 2013 review⁴⁴ was updated (identifying studies up to April 11, 2017). Further details on the methods of the search are given elsewhere.¹ Studies were included only if they contained an SEE to elicit uncertain parameters (in the form of a distribution) to inform MBEE in health.

The methods used in each application were extracted along with the criteria used to support methodological and practical choices and any issues or challenges discussed in the text. Issues and challenges were extracted using an open field, and then categorised and grouped for reporting.

Section 4.3 Aspects related to the design of the SEE

In existing applications, experts beliefs were sought for only a few parameters of a decision model, often not elicited directly but calculated from one or more alternative elicited quantities. Quantities included: event probabilities, relative effectiveness, time to event and diagnostic accuracy (see Soares, et al¹ for full details). The choice of which quantities to elicit was based on a number of criteria. The first was appropriateness for experts, specifically that parameters in decision models can be complex and may not be directly observable by experts. Second are statistical concerns. The quantities elicited should be fit-for-purpose for further analysis, for example allowing elicited evidence to be combined with any existing empirical evidence, statistically coherent and reflect any dependencies between the quantities elicited. Finally the burden to experts need to be considered.

1
2
3 Burden can be reduced by, for example, limiting the number of target parameters to elicit or eliciting
4 homogeneous quantities throughout the exercise.
5
6
7

8 Almost exclusively, applications have recruited health care professionals, based on the following
9 criteria: recognition by peers,⁵¹ specialist knowledge or clinical experience,^{52, 51, 53, 54, 55, 56, 57} based in
10 the relevant jurisdiction,^{52, 51, 54,55} research experience,^{51, 56, 57} and lack of involvement in product
11 development.⁵³ A number of authors^{52, 58,59} recognised that health care professionals are unlikely to
12 have knowledge of elicitation and may have only sparse quantitative skills. This has driven choices
13 made in designing and conducting the SEE, such as training needs, method of elicitation and
14 definition of the quantities to elicit.⁵²
15
16
17
18
19

20
21 Many of the applications have included a varied sample of experts by recruiting them from a range
22 of relevant specialties,^{51, 60, 61} clinical settings^{52, 51,61} and geographical areas/countries^{51, 57} to capture
23 heterogeneity in beliefs (reflecting underlying heterogeneity in patient populations), and avoid
24 dependency between experts.⁵¹ The potential for bias in expert opinion was recognised in some
25 SEE^{52, 56} with reported attempts to minimize bias in the design.⁶² Two applications make explicit
26 efforts to avoid recruiting experts that may have motivational biases.^{53,61} Two studies provided
27 information on cognitive biases in the training session.^{52,59}
28
29
30
31
32
33

34
35 In eliciting uncertainty, applications have typically used either the FIM^{52,51, 63, 60,54, 55, 59, 64, 62, 65} or
36 VIM.^{53,58,66, 67,56, 57, 68} Choices were justified on the basis of: pilot exercises designed for the purpose
37 (see below), generic methods research, previous use in MBEE, and claims of lower burden or
38 intuitiveness for experts.
39
40
41
42
43

44 Fourteen studies elicit individually from experts and aggregate mathematically, three aimed to
45 achieve consensus amongst experts^{69,66, 67} and three others did not explicitly report the method of
46 aggregation used.^{70, 61,64} None of the three studies using consensus was explicit about the reasons
47 for choosing consensus or the process of achieving it. Authors justify the choice of mathematical
48 aggregation based on the desirability to reflect variation within and between experts⁶⁰, because
49 consensus is known to lead to overconfident results (i.e., narrow distributions)⁵¹ and because it
50 raises practical difficulties of convening experts and providing experienced facilitation. With regards
51 to weighting in the mathematical approach, most of the applications reviewed claim insufficient
52 justification for generating differential weights^{52, 51} and lack of clarity on how to appropriately
53 generate the weights^{52, 53, 62} and hence apply equal weighting. Five studies, however, explored
54
55
56
57
58
59
60

1
2
3 unequal weighting, either based on responses to seed questions^{52, 54,57, 59} (performance-based
4 weighting) or using the clinical background of experts(objective weighting).⁵³
5
6
7

8 **Section 4.4 Experiences with the conduct of the exercise**

9

10 No studies reported major challenges in the conduct of the SEE, despite the complexity of the task.
11 Some conducted a group based session, which were typically face-to-face, although studies^{51, 60, 58, 62}
12 departed from this format due to time constraints, geographical limitations and availability of
13 experts. Mathematical exercises adopted a mix of formats, ranging from individual interviews to
14 remote completion via email. Administration, where details were specified was via bespoke tools
15 using Excel,^{52, 51,54, 55,59, 62} paper questionnaires, a generic elicitation package(the Sheffield Elicitation
16 Framework or SHELF)^{61, 67} and a software package for the elicitation of dependency (Prior Elicitation
17 Graphical Software, PEGS).⁵⁸ Some exercises were explicit about piloting the tool to ensure clear
18 wording of the questions,^{52,53, 55,56} and most offered opportunities for revision and/or graphical
19 feedback.
20
21
22
23
24
25
26
27

28 Five applications were explicit about training of experts^{51, 60, 54, 59} covering: overview of the project
29 and of the role of elicitation,^{52, 60, 53, 54, 59} quantities required and definitions,^{52,60, 54, 59} explanation
30 and expression of uncertainty;^{51, 53} consideration of potential biases,^{52, 54, 59} use of the elicitation
31 instrument⁵² and delivery of practice exercises.^{51, 60,54,59} Studies that implemented elicitation
32 remotely generally included some form of instructions, although none reported these in detail.
33
34
35
36
37
38

39 **Section 4.5 Experiences with the analyses and interpretation of elicited evidence**

40

41 Studies did not report details sufficiently, however validity was assessed according to: missingness,
42 validity checks and self-reported face-validity. Some applications requested feedback from experts
43 on: the ease of completion of the SEE,^{52, 51, 59,62} the basis for experts answers (to reveal the sources
44 of evidence considered by the experts and their level of knowledge),⁵¹ or on self-reported face
45 validity.^{52, 51, 60, 59}
46
47
48
49
50

51 Of the 14 studies that used a mathematical approach to aggregation, one did not generate a group
52 estimate and instead used the responses of each expert individually.⁵⁸ The majority linearly pooled,
53 by averaging individual distributions. In order to pool, some applications were not explicit about
54 how prior distributions were derived from elicited summaries. Those that were explicit used
55 parametric distributions, with the choice of distribution either not justified or based on general
56
57
58
59
60

1
2
3 MBEE literature on distribution choice for probabilistic sensitivity analyses.⁵¹ To fit the distribution
4 some applications used software^{63, 53} and others a specific fitting method, such as maximum
5 likelihood fitting.
6
7

8
9
10 In some applications elicited evidence was used directly as input to a cost-effectiveness model.^{51, 60,}
11 ^{58, 69, 71, 59} Where external evidence existed on elicited parameters, some authors present both
12 sources separately using scenarios,^{60, 70} while others combine them using Bayesian updating.^{52,70,67}
13
14 Three authors^{63, 58, 55} explored use of individual experts beliefs and found that results and associated
15 allocation decisions varied between experts.
16
17

18 19 20 **Section 4.6 Conclusions**

21 This critical review demonstrates that reporting is poor (as also identified elsewhere),⁷² and there is
22 a lack of consensus on methodology. A number of principles from the elicitation literature are
23 expected to generalise to the MBEE setting, such as the need for piloting and training; however, for
24 many other areas of SEE, it is not clear that methods used in other disciplines translate to HTA.
25
26
27

28 29 30 **Section 4.7 Discussion**

31 The review highlights a number of specificities/constraints that can shape the development of
32 guidance and target future research efforts in this area. Firstly, there exists important between-
33 expert variation. In other disciplines, variation is generally linked to different levels of bias and hence
34 regarded as undesirable, warranting the use of strategies to reduce or discourage variation, such as
35 consensus methods. The majority of applications in MBEE, however, expect wide variation in the
36 beliefs of multiple experts due to genuine heterogeneity in the populations' experts draw upon.
37
38
39

40
41
42
43 Secondly, substantive experts in HTA are health professionals who may not be trained in
44 quantitative subjects, unlike other areas of science in which elicitation is used such as engineering or
45 meteorology. Further research on SEE should consider the appropriateness of alternative methods
46 of elicitation (e.g. chips and bins, or bisection method) for the potentially less normative experts, or
47 on how to facilitate the elicitation of complex parameters, including dependency. Some of the
48 applied examples seek for assurance on the validity of the particular exercise. It is, however, not
49 clear how such an assessment should proceed. Examples have used self-reported face-validity
50 assessments, sensitivity analyses, and performance weighting (calibration). Particularly for
51 performance weighting, despite a growing (generic) literature discussing the validity of this approach
52
53
54
55
56
57
58
59
60

1
2
3 (see for example),^{73,74,75} the applied literature struggles with supporting the methodological choices
4 that need to be made.
5
6
7

8
9 Finally, while it is generally agreed that SEE should be designed and conducted in a way that
10 minimises the use of heuristics and other sources of bias, there is little integration in the applied
11 literature of the findings from behavioural research. A recent review placing special emphasis on
12 debiasing techniques⁷⁶ is a helpful resource to be reflected in future research.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Chapter 5 Reviewing the evidence: expert selection, level of elicitation, fitting and pooling

Section 5.1 Introduction

Each element of a SEE process encompasses several possible components for which choices need to be made, with each choice successively impacting on the next. If inappropriate choices are made for some component, the process can provide inaccurate or misleading judgements. Despite this risk, there is little or no empirical evidence which compares many of the alternative choices available when designing, conducting and analysing a SEE.

The aim of this chapter was to review the SEE literature to identify the possible choices and related evidence available in four components of the SEE process. Following discussion with the advisory group, the following four were chosen:

1. Selection of experts
2. Level of elicitation
3. Fitting and aggregation
4. Assessing the expected accuracy of experts judgements

Section 5.2 Identifying literature

The use of SEE in HCDM and cost-effectiveness modelling is still evolving. Considering this, literature from both a cost-effectiveness context and other disciplines are reviewed in this chapter. While this chapter is not a systematic review, a semi-structured approach was employed to identify the relevant SEE literature and to summarise existing evidence. The targeted searches for the SEE elements listed above were conducted using the same semi-structured approach.

Literature was initially searched by reading a selection of well-known books and papers addressing SEE.^{8, 14, 31, 77} This approach defined each component and identified the associated choices available for each one. Following this, more recent literature on SEE was then explored to investigate the availability of any recommended choices based on sound principles or evidence for the identified components.^{33, 44, 58, 59} The guidelines reviewed in *Chapter 2* were also included if they provided evidence. Papers that were specific to the four components were also included in the review.) For applications of elicitation in the context of HCDM, research from Leal et al (2007),⁵¹ Bojke et al (2010)⁵⁴ and Soares et al (2011)⁵² were consulted, along with a review of elicitation methods in cost-effectiveness analysis.¹

Each section describes the requirements for each part of the process, identifies the choices available and presents any evidence or principles to inform the choices. In *Chapter 9*, choices recommended in this chapter are considered against a set of principles underpinning elicitation in HCDM.

Section 5.3 Selection of experts

An 'expert' is defined as someone who has great knowledge of the subject domain⁸ and who is competent in the practical application of this knowledge.

The SEE literature recognises several choices regarding to expert selection, relating to;

1. *Finding experts with the relevant skills.* The literature recommends that an expert possesses substantive skills in the target domain and normative skills regarding the expression of uncertain probabilities.^{44, 59} Ideally the expert will have specific knowledge in the field, together with a broad perspective.⁵⁹ Where evidence is less developed, the expert will require adaptive skills to elicit judgements.⁸
2. *Measuring competencies in these skills.* Substantive skills can be identified and measured using social indicators of expertise or peer- and self-assessment tools, such as the Generalised Expertise Measurement (GEM)¹⁴ and the Expert-Selection Questionnaire (E-SQ),⁹² respectively. As normative expertise is more generic, the quality of probability judgement can be measured against seed questions.^{14, 54, 59}
3. *Identifying experts.* The elicitation literature reports the following criteria for sourcing possible experts: recognition by peers,⁵¹ specialist knowledge or clinical experience,⁵¹⁻⁵⁷ current work in the target domain,⁵¹⁻⁵⁴ research output^{51, 56, 57} and lack of involvement in the product of interest.⁵³
4. *Recruiting experts.* Literature reports recruiting experts with a formal nomination process to develop a long-list of potential experts.^{20, 22, 93} More recent literature describes the use of a profile matrix listing the essential and desirable characteristics that are required.⁹⁴
5. *Identifying the optimal number of experts.* The analyst needs to decide whether to recruit multiple experts, and if so, how many, or one 'top' expert.

In the literature, substantive and normative skills are the most frequently referenced skills while the role of adaptive skills is not addressed as often. We think that adaptive skills can be of particular

1
2
3 importance in a HCDM context where experts may need to elicit judgements on new or emerging
4 technologies in which they do not have a great deal of experience.
5
6
7

8 There are concerns in the literature regarding the methods by which experts skills can be measured,
9 with one source branding these measures as subjective, particularly when referring to social
10 indicators of expertise.⁸⁰ Peer and self-assessment tools have also been queried in terms of how
11 accurately they measure substantive expertise.^{52, 57, 59} When measuring normative skills, selecting
12 appropriate seed questions is difficult to assess (see Section 5.6). Until the relative advantages of the
13 different methods for identifying and measuring substantive and normative skills are better
14 understood, no technique in particular can be recommended.
15
16
17
18
19

20
21 The SEE literature suggests the recruitment of experts will be largely influenced by the target
22 domain. In the healthcare literature, recruited experts are largely clinicians or professionals
23 practicing in the target domain.^{44, 51, 52, 54, 59, 95} The recruitment strategies summarised above come
24 from outside a HCDM context. Consequently, their applicability to elicitation in healthcare may not
25 be appropriate.
26
27
28
29

30
31 There are not consistent recommendations on sourcing and recruiting experts, possibly because
32 strategies are domain-dependent making it difficult to find a strategy that is appropriate across
33 different contexts. Some of the processes of expert recruitment reported in literature, such as
34 research output, can be interpreted as social indicators of expertise. These indicators have been
35 subject to critique (Burgman reference) and should be treated with caution if used for expert
36 selection.
37
38
39
40
41

42
43 Many reported influential factors can impact the optimal number of experts to recruit, such as the
44 number of experts available with the relevant expertise,⁵¹ time and budget constraints and mode of
45 administration of the SEE process.^{56, 59} Smaller samples are recommended for face-to-face modes of
46 administration. The SHELF and Classical methods are optimally conducted with between five to ten
47 experts as there are diminishing returns to accuracy improvement with more experts.⁹⁴ In contrast,
48 larger sample sizes are possible using remote methods of administration, as done in the Delphi
49 method.⁹⁶ When determining the optimal number of experts, heterogeneity amongst the experts
50 also needs to be accounted for. It is logical to think that in a HCDM context, particularly for a new
51 intervention or unknown condition, the pool of available experts will be limited. Despite the fact that
52
53
54
55
56
57
58
59
60

1
2
3 SEE is less costly than primary data collection, the financial and time resources that are available for
4 the design, conduct and analysis will dictate the number of experts that should be recruited.
5
6
7

8 Selecting the optimal number of experts is one of the few components of expert recruitment that
9 have been studied empirically. Budescu and Chen (2015)³³ assessed the benefits of adding additional
10 experts, concluding that the best performance is found using between three to 16 experts, with
11 around six being optimal. This assumes that all experts perform to the best of their ability. If there is
12 a redundancy in expertise, the number of experts will need to be greater than six.³³
13
14
15
16
17

18 **Section 5.4 Level of elicitation (individual versus group)**

19 Generally, judgements from multiple experts will be sought in the SEE process. These judgements
20 can be elicited individually from experts or from a group of experts.
21
22
23

24
25 When choosing the level of elicitation, the SEE literature suggests the analyst will need to consider
26 the following choices and their associated principles;
27
28
29

- 30 1. The expert may provide their own judgement individually without interacting with other
31 experts and this is referred to as *individual* level elicitation.¹⁷
32
33
- 34 2. Alternatively, experts are encouraged to interact with one another in a group to discuss the
35 uncertain quantity until they achieve consensus.⁸ This is called *group* level elicitation.
36
37
- 38 3. A third approach uses a *combination of individual and group level* elicitation where experts
39 first provide judgements individually and then engage in a facilitated face-to-face group
40 discussion until they reach a consensus(SHELF method).¹⁵ This approach can also be
41 conducted remotely using an iterative survey with several rounds of elicitation where each
42 expert has access to the opinion of others through a highly restricted level of interaction
43 (Delphi method).^{17, 97} On receiving the new information from peers, experts are given the
44 opportunity to reach consensus using this remote method.
45
46
47
48
49
50
51

52 Figure 4 presents these different levels of expert elicitation describing their level of interaction,
53 consensus and how uncertainty is quantified.
54
55
56
57
58
59
60

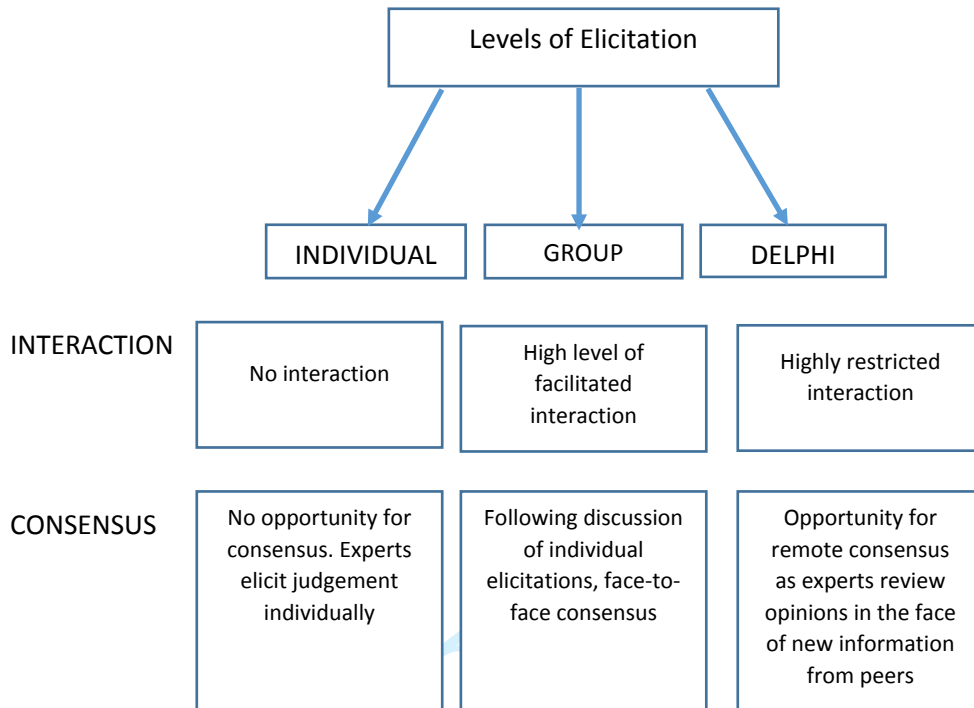


Figure 4 Levels of expert elicitation*

*Information taken from EFSA: Guidance on Expert Knowledge Elicitation in Food and Feed Safety Risk Assessment.¹⁷

SEE literature suggests that individual level elicitation is the most effective approach when eliciting expert beliefs,^{8, 77} as it reduces the risk of bias due to experts influencing each other⁶⁰ and promotes accountability to the decision maker. Individual level elicitation is also recommended for its transparency, in particular where expert judgements are available for peer review.²⁷ Despite these advantages, one of the concerns associated with this approach is that experts may not feel confident in expressing uncertainty individually compared to a group situation.⁵³ The literature reports that experts involved in group level elicitation are more confident about their decisions compared to experts at individual level elicitation.⁸¹

Group level elicitation provides a substantial exchange of information amongst experts^{51, 92} and consequently, it is expected that more informed experts will have greater influence in the group.⁸¹ When aiming to achieve consensus face-to-face, a facilitator is required. The role of the facilitator is to engage the entire group and protect the process from becoming dominated by a subset of experts.⁸ If group elicitation is not monitored by an experienced facilitator, the interaction may pressure experts into reaching a consensus, and some experts may suppress their opposing beliefs.^{51, 58} Another concern associated with group level elicitation is that convening experts from various

1
2
3 geographical areas at a time and place suitable for all can prove to be problematic.⁵¹ Depending on
4 appropriateness, the Delphi method can help overcome this issue as this method is conducted
5 remotely using web-based survey tools. This approach is adequate when little but crucial
6 information is required and a quick response is needed.¹⁷
7
8
9

10
11 In the healthcare literature, a 2013 review by Grigore reports individual elicitation is the dominant
12 approach with 13 of the 14 included studies reporting individual level elicitation.^{51, 52, 54, 55, 58, 64, 53, 56,}
13 ^{60, 62, 66, 68, 70} Only two of these studies justify why the particular level was selected. Bojke et al (2010)
14 report choosing individual level elicitation rather than group level to capture uncertainty within and
15 between experts.^{19, 54} Leal (2007) reports choosing it based on guidelines for HDCM developed by
16 Philips et al in 2004.^{51, 77}
17
18
19
20
21
22

23 As discussed in the introduction of this section, choices made for each component can have an effect
24 on subsequent components of the SEE process, in particular, the chosen level of elicitation is related
25 to the mode of administration and the method of aggregation used in the post-elicitation phase.
26
27
28
29

30 In terms of mode of administration, individual level elicitation is relatively versatile. Various
31 platforms are reported when using the this approach; face-to-face^{53, 60, 64, 68} computer-based⁵⁶ and
32 email.⁵¹ Given these possibilities it is not surprising that the literature indicates individual level of
33 elicitation can be conducted face-to-face or remotely, and may be facilitated^{26, 52, 60} or non-
34 facilitated.^{51, 54, 56} In contrast, group level elicitation is not as flexible. At the group level, experts
35 must be brought together, usually in one location²⁶ and the elicitation must be guided by an
36 experienced facilitator. Thus, this approach can be time consuming and costly. Once the experts are
37 organized, two modes of administration are reported in the literature for group level elicitation;
38 paper-based^{57, 60} and excel-based.⁵²
39
40
41
42
43
44
45
46

47 As discussed in *Chapter 2*, the SEE process may strive to achieve a unique distribution that reflects
48 the beliefs of all experts. This process is described as either “behavioural” or “mathematical”
49 aggregation, and is related to the level of elicitation^{8, 8, 81}. Behavioural aggregation relies on
50 interaction between the experts to create a single distribution that reflects either the experts
51 consensus beliefs or how an independent, rational observer would summarise the collective
52 opinions of the experts. Mathematical aggregation, in contrast, combines individual assessments
53 from several experts (who may or may not have interacted) into a combined group assessment
54
55
56
57
58
59
60

1
2
3 based on an algorithm or mathematical process. Methods of mathematical aggregation are
4 discussed in Section 5.5.2.
5
6
7

8 Behavioural aggregation can involve either the elicitation of a single group distribution, or elicitation
9 of individual judgements followed by interaction to produce a consensus distribution (O'Hagan et al.
10 2006). Both Clemen and Winkler (1999) and O'Hagan et al. (2006) recognize that it may not always
11 be desirable to strive for a consensus distribution.^{8, 81} Both suggest that in such cases, mathematical
12 aggregation can be applied at the end of group interaction, thus combining behavioural and
13 mathematical approaches. The EFSA Delphi method and IDEA protocol are both examples of this.^{17, 25}
14 Clemen and Winkler (1999)⁸¹ state that the benefit from group interaction is the sharing of
15 information and not the forced consensus. They also believe that individual probability assessments
16 are useful for understanding the range of expert opinion and conducting sensitivity analysis, thus
17 supporting combined aggregation approaches or behavioural approaches that also involve individual
18 elicitation (such as the SHELF method.²⁹
19
20
21
22
23
24
25
26
27
28

29 The empirical evidence on the merits of individual and group level elicitation is dated and may not
30 be entirely relevant to HCDM. When solving problems that require 'originality and insight', Fogel
31 (1967) commends the use of an interactive group.⁸² This is supported by Seaver (1976) who
32 performed a comparative study of the three methods; individual elicitation, group interaction and
33 the Delphi method and reports that the interactive group produced a larger number of ideas
34 compared to the Delphi method.⁸³ Stael Von Holstein (1971) report that results using the individual
35 level approach are judged to be poorer than results from an interactive group.⁸⁴ However, Meyer
36 and Booker (1991)⁹², emphasise that these studies were not applying this method for deep problem-
37 solving, the type of problem-solving for which it is most suited.
38
39
40
41
42
43
44

45 While there is a lack of empirical evidence available comparing individual and group level elicitation
46 in healthcare, individual level elicitation is the most commonly adopted approach in the healthcare
47 literature and the most recently recommended choice in guidelines for decision-making in
48 healthcare.⁷⁷ The cost of using group-level elicitation will depend on the context, crucially, how
49 physically dispersed the experts are. There may be a particular case for using group level elicitation:
50 when the problem structure is unclear and there is a need for experts to develop a consensus
51 problem structure and to specify the elicitation questions; when experts have distinctively different
52 disciplinary backgrounds and knowledge bases (e.g. practising clinicians versus epidemiologists) and
53 so require discussion to assess each other's claims to expertise; or when it is expected that the
54
55
56
57
58
59
60

1
2
3 experts will work together repeatedly and so group-level elicitation may also serve a team-building
4 purpose.
5
6
7

8 **Section 5.5 Fitting & pooling**

9
10 This section concerns the translation of the information elicited from one or more experts into a
11 probability distribution representing the evidence on an uncertain quantity to inform decision-
12 making. The following choices can arise:
13
14

- 15
16 1. To obtain a probability distribution from a single expert's belief (or a behavioural
17 aggregation process), should we pre-specify the form of the distribution and elicit its
18 parameters directly, or elicit characteristics of an unspecified distribution, such as quantiles?
19
20 a. If a distribution is pre-specified, what distribution should this be, and how should its
21 parameters be elicited?
22
23 b. If characteristics of an unspecified distribution have been elicited, what distribution
24 should be fitted to the elicited information, and how?
25
26 2. If individual views are elicited from multiple experts, and we choose to mathematically
27 aggregate them into a single distribution representing the overall spread of views, how
28 should this be done?
29
30
31
32
33
34

35 *Section 5.5.1 Distribution choice and fitting*

36
37 Methods of eliciting parameters of common distributions have been described, e.g. by Winkler
38 (1967) and O'Hagan et al (2006), for example, the Beta distribution for probabilities.^{98, 78} Only a small
39 number of these methods have been evaluated and compared, and O'Hagan et al (2006)
40 recommended much more work before advocating one particular method for use in practice.⁷⁸
41
42
43
44

45 The disadvantage of pre-specifying a distribution is that it may not fit the expert's belief. Instead of
46 eliciting parameters of pre-specified distributions, we may elicit characteristics of an unspecified
47 distribution. Typically the expert is asked either for the quantiles that contain a given probability
48 mass (e.g. median and credible intervals, as in the "bisection" method), or the probability masses
49 that lie within a given set of quantiles (e.g. the "chips and bins" method) (see *Chapter 8*). Either
50 way, the elicited data consist of a set of points on a cumulative distribution function (CDF) (*Figure 5*).
51 To obtain a fully specified distribution from these, the points could simply be interpolated, as
52 described by O'Hagan et al.⁷⁸ More commonly, a parametric family of distributions is specified,
53 followed by identifying the parameters which best fit the elicited data. The advantage compared to
54
55
56
57
58
59
60

interpolation is that the fitted CDF is not necessarily assumed to pass through the elicited points, acknowledging that the expert may not be fully confident in the precise values that they provide. The “best fitting” parameters can be determined by numerical methods, such as least squares, as in SHELF,⁹⁹ Leal et al (2007)⁵¹ and Thall et al (2017),⁸⁵ which is justified by maximum likelihood principles. The red line in *Figure 5* shows the CDF of the best-fitting Beta distribution determined by this method. A related approach is the method of moments, an approximation to maximum likelihood, used by Bojke et al (2010)⁵⁴ and Soares et al (2011)⁵² for two-parameter distributions.

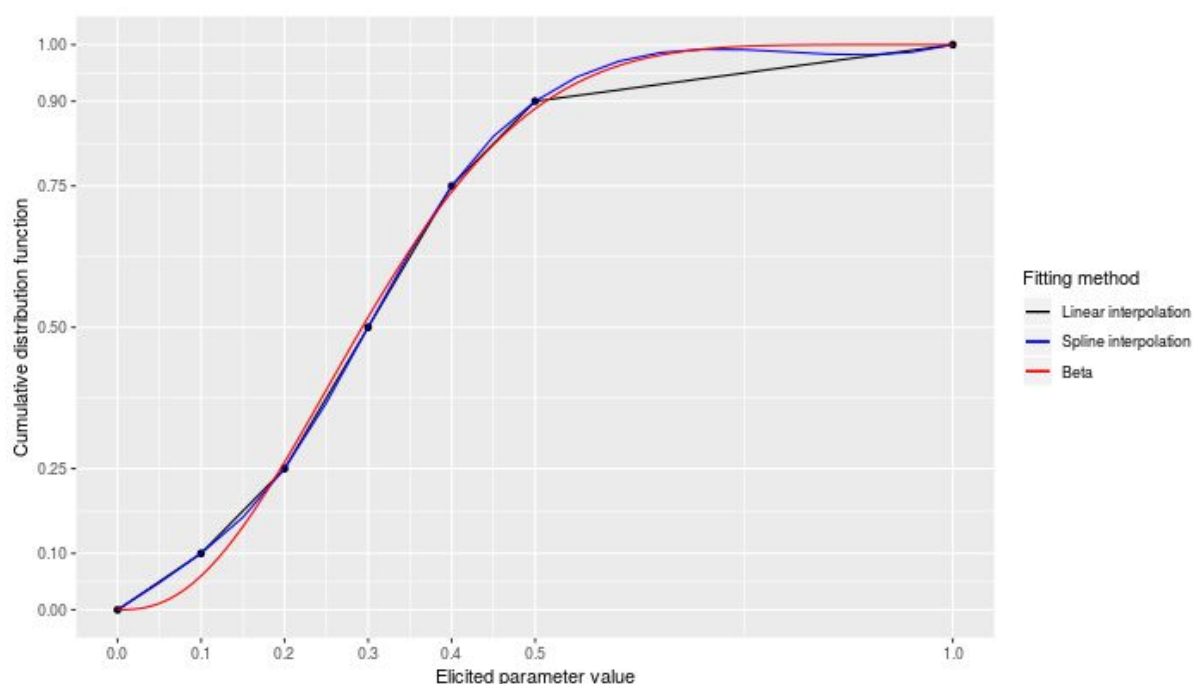


Figure 5 Elicited points on a cumulative distribution function (CDF) and alternative fitted distributions

Standard conjugate families (such as the Beta or Normal) can be combined easily with future observed data using Bayesian inference. As an alternative to standard families, which may not fit the elicited data well, Bornkamp and Ickstadt (2009) proposed to fit a penalized spline function.⁸⁶ This method was implemented in their R package “SEL”. O’Hagan et al (2006) recommend a process of “feedback”, where the fitted distribution is presented to the expert with the opportunity to revise it to better reflect their beliefs.⁷⁸ However in practice, there may not be sufficiently many elicited points on the CDF to identify distributions that fit better than standard ones.

As well as the expert not being fully confident in the precise values they provide, there may be multiple distributions which fit the elicited data equally well.¹⁰⁰ Bayesian nonparametric methods to

1
2
3 handle both these forms of uncertainty have been developed, see Oakley and O'Hagan (2007),
4 Gosling et al. (2007), Moala and O'Hagan (2010) and Daneshkahr et al. (2017).^{101, 102, 87, 88} These
5 methods generally require Markov Chain Monte Carlo (MCMC) simulation to implement, and, as far
6 as we are aware, there is no accessible software to implement them. Such methods tend to be
7 computationally intensive, which may not allow the fitted distribution to be instantly "fed back" to
8 the expert during an elicitation session.
9
10
11
12

13 14 15 *Section 5.5.2 Mathematical aggregation*

16 Mathematical aggregation methods fall into two general approaches: Bayesian combination (or
17 "supra-Bayesian") methods, and "axiomatic" (or "opinion pooling") methods.
18
19

20
21
22 In Bayesian combination, the decision maker or modeller treats each expert assessment as new data
23 and uses it to update his/her own distribution for the unknown parameter.^{103, 104} The resulting
24 distribution thus represents the beliefs of the modeller given the elicited data.⁷⁸ This is difficult to
25 apply in practice, though, due to the detailed information required on biases in and dependences
26 between the experts assessments.^{105, 81, 78} See, e.g. Lipscomb et al. (1997), Albert et al (2012), West
27 and Crosse (1992), Gelfand, Mallick and Dey (1995) for some examples.^{106, 107, 108,109}
28
29
30
31

32
33
34 In opinion pooling methods, the aggregated distribution is an average of the distributions from each
35 expert. "Linear pooling" uses an arithmetic average and was the most common aggregation
36 approach used in our review of elicitation in HCDM.¹ The alternative "log pool" uses a geometric
37 average, and harmonic averaging has also been used.^{110,111} Past work has shown that there is no
38 mathematical formula that can simultaneously satisfy a number of potentially desirable criteria,^{112, 81}
39 so there is no obvious justification for one combination rule over another. O'Hagan (2006) observe
40 that log pooling discounts values found implausible by at least one expert, leading to a distribution
41 concentrated on areas of agreement, while a linear pool encompasses all values that any expert
42 finds plausible, leading to a broader distribution. Harmonic averaging gives an even more
43 concentrated result than log pooling.⁷⁸ Hammitt et al (2013) performed simulation studies to
44 compare linear pooling with other mathematical aggregation methods, in situations where the
45 experts beliefs were generated by a known mechanism, concluding that linear pooling performed
46 worst, but it is unclear whether this mechanism holds generally.¹¹³ In HCDM, we would argue that
47 the broader distribution from linear pooling is preferable, as it acknowledges the uncertainty arising
48 from between-expert variation. This could motivate further research to obtain observed data on the
49 uncertain quantity, and strengthen the evidence base for decision-making.
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5 The weights applied to each expert's belief in an opinion pool are commonly chosen to be equal.
6 Alternatively they could represent an estimate of expected accuracy, determined using seed
7 questions (discussed in the final section of this chapter), prior assessments of the expert's
8 background^{80, 114, 115} or agreement of the expert's elicited data with subsequently-observed data on
9 the quantity of interest,¹¹⁵ though observed data would not generally be available in a health care
10 context. . More technically advanced methods that use only the elicited data to form weights are
11 presented by, e.g., Ranjan and Gneiting (2010), Rufo et al (2012) and Hora and Kardes (2015).^{116,}
12 ^{117,118} Essentially, these adjust the simple linear or log pools to give a better expected balance of
13 overall bias and over / under-confidence.
14
15
16
17
18
19

20
21 Meta-analysis methods have also been considered for pooling expert beliefs^{68, 54} but have been
22 argued to be inappropriate,⁵⁴ since they assume each expert's view is fully based on evidence that
23 no other experts have seen.
24
25
26
27

28 **Section 5.6 Assessing the expected accuracy of experts judgements**

29 This section focuses on estimating the expected accuracy of the judgements elicited from an expert
30 or group of experts, compared to the truth. The accuracy of a probabilistic judgement is often
31 referred to as "calibration".^{89, 14, 78} The SEE literature recognises a number of different choices in
32 this area, relating to:
33
34
35
36
37

- 38 1. *Measuring calibration using scoring rules.* When experts answer questions about quantities
39 that have "realizations" (that is, known answers), called "calibration" or "seed" questions,
40 the accuracy of their elicited judgements can be estimated by scoring rules that compare the
41 elicited assessments to the realizations, a practice first proposed by Winkler and Murphy
42 (1968).¹¹⁹ The most common strictly proper scoring rule used in SEE is that of Cooke's
43 Classical Model, which has been used to elicit judgements and measure calibration in over
44 100 expert panels.¹⁶ More detail on this method is available elsewhere (e.g., Cooke 1991;
45 Cooke and Goossens 2000; Quigley et al. 2018).^{14, 20, 120}
46
47
48
49
50
51
52
53
- 54 2. *Using the scoring rule to create and/or evaluate combinations of experts.* In the Classical
55 Model, scoring rules are also the mechanism for creating performance-based weights for
56 mathematically aggregating expert assessments. The Classical Model's scoring rule can also
57 enable the evaluation of combinations of the experts.¹²⁰ In practice, this is typically done by
58
59
60

1
2
3 comparing scores from the performance-weight linear combination of experts to the equal-
4 weight linear combination to see which has better performance. This is done both within a
5 study, to inform the choice of using performance- or equal-weights in the final reported
6 results,⁹⁰ and across studies, to evaluate the method more broadly (Cooke and Goossens
7 2008¹²¹; Colson and Cooke 2017⁷³; Quigley et al. 2018¹²⁰).

- 8
9
10
11
12
13 3. *Deciding how many seed questions should be used and how relevant seed questions are*
14 *identified.* Multiple seed questions are needed to assess the accuracy of elicited probability
15 distributions,⁷⁸ since poor calibration based on one seed question may indicate bad luck
16 rather than bad performance. Identifying appropriate seed questions is a challenge because
17 the questions must be closely related to the target questions but unknown to the experts
18 participating in the elicitation. Seed questions are used to assess an expert's skill in
19 quantifying uncertainty, so they should not just be a test of the expert's ability to recall
20 established facts or familiar quantities.
21
22
23
24
25
26
27

28 Scoring rules for sets of variables, rather than individual variables, have been argued to be
29 preferable, as the latter does not depend on the distribution of the realization.¹²² The Brier score,
30 for example, is a scoring rule for individual variables that was recently used to score expert forecasts
31 of geopolitical events.^{123,124} Cooke (2014)¹²², however, provides simple counterexamples that
32 demonstrate the issues with this approach to scoring. Strictly proper scoring rules are rules in which
33 an expert maximizes their score by stating their true beliefs. As the objective of an elicitation is to
34 capture the beliefs of the experts, it is critical that, if scoring rules are used to measure calibration,
35 they must be strictly proper.¹²² An improper scoring rule may reward an expert for providing
36 assessments more extreme than their real beliefs.
37
38
39
40
41
42
43
44

45 Two studies compare scores from the performance-weight and equal-weight combinations of
46 experts in 78 total applications of the Classical Model, and they find that the performance-weight
47 combination is consistently both more informative and more statistically accurate than the equal-
48 weight combination.^{121,16} These studies are based on in-sample comparisons, where the same set of
49 questions are used both to calculate the expert weights and evaluate the performance of the
50 method. Out-of-sample validation, in contrast, would estimate performance using external data.
51 However, data rarely become available on elicited quantities of interest (which is why elicitation is
52 needed). An alternative approach is cross-validation, where the set of seed questions are divided
53 into a training set and a test set. Expert scores are calculated based on the training set and then the
54
55
56
57
58
59
60

1
2
3 performance-weight combination is evaluated on its performance on the test set. The most recent
4 and extensive cross validation study of the Classical Model done to date found that the
5 performance-weight combination of experts outperforms the equal-weight combination in 26 of the
6 33 studies.⁷³ An evaluation of separate data, based on expert forecasts of the probability of various
7 geopolitical events, also concluded that the accuracy of an expert's assessments can be predicted by
8 past performance on related questions, supporting the use of performance-based expert weighting
9 (Hanea et al. 2018).¹²¹
10
11
12
13
14
15

16 In the Classical Model, because the scoring rule is asymptotically proper, there is no theoretical basis
17 for the number of seed questions required, but at least ten seeds is the recommended rule of
18 thumb.^{14, 20} Simulations of expert scores show that using ten seed questions allows an analyst to
19 distinguish between a well-calibrated and a slightly overconfident expert.¹²⁰ Another paper argues
20 that significantly more seed questions are needed, but it incorrectly understands the Classical
21 Model's scoring to be for the purpose of hypothesis testing, rather than for discriminating between
22 experts⁹⁴
23
24
25
26
27
28
29

30 Seed questions commonly come from four sources: future measurements, unpublished
31 measurements, unfamiliar information from standard datasets, or combining or comparing different
32 datasets.¹²⁰ They discuss examples of each of these strategies from past applications of the Classical
33 Model.
34
35
36
37

38 Although seed questions should be related to the subject of the elicitation, there is no clear test to
39 measure if a question is "close enough" to the target questions. In practice, Classical Model
40 elicitation focus on ensuring that the link between seed and target questions is strong enough that
41 the problem owner, experts, and knowledgeable reviewers accept the resulting unequal weights of
42 similarly qualified and knowledgeable experts. The Classical Model also recommends specific
43 sensitivity analysis to identify if any seed questions have a large impact on the results.²⁰
44
45
46
47
48
49

50 **Section 5.7 Conclusion**

51 Different methods for selecting and recruiting experts are recommended in the SEE literature with
52 very little empirical comparison. In terms of the skills that experts should possess, we believe
53 adaptive skills are important in SEE in HCDM, given the potentially novel nature of some of the
54 technologies the expert may make their judgements on. Yet there is a lack of acknowledgment of
55 this skill in the existing literature. This is explored further in a HCDM context in *Chapter 8*. While
56
57
58
59
60

1
2
3 some recommendations, such as the construction of a Profile Matrix, appear very useful, including
4 such an exercise within the conduction of an expert elicitation process in HCDM may not be feasible
5 due to time constraints on the overall project for which the elicitation is being conducted. With few
6
7
8 conclusive recommendations, the findings of this targeted search suggest further analysis is required
9
10 relating to the selection of experts in expert elicitation.
11

12
13 Individual level elicitation is conducted without interaction between experts, while group elicitation
14 requires experts to interact to discuss the uncertain quantity. The relative merits of individual and
15 group-level elicitation for HCDM are unclear, yet there are guidelines recommending individual level
16 elicitation as more appropriate given the complex nature of quantifying uncertainty (likely to apply
17 to HCDM) and the importance of seeing and understanding differences between experts. Despite
18 these recommendations, the findings of this targeted search suggest that further research is
19
20
21 required in this area.
22
23
24
25

26
27 Elicited data consisting of points on a CDF should be converted into a smooth distribution
28 representing the assumed state of belief. To represent one expert's belief, numerically fitting
29 standard distributions such as the Beta to the elicited data will often be sufficient. If the number of
30 elicited points is small, then more elaborate models would be difficult to identify. Standard
31 distributions are simpler to implement, but more complex approaches are worth considering,
32 particularly if they can be shown to give a better fit to the elicited data. Spline regression models
33 can be fitted instantly in general software, however more experience of these is needed. More
34 complex Bayesian non-parametric approaches can better represent uncertainty about the full belief
35 distribution given the elicited data, however more guidance and accessible software are required
36 before these can be recommended for routine use in HCDM.
37
38
39
40
41
42
43
44

45 To mathematically aggregate elicited data from multiple experts, linear pooling is simple to
46 implement, and allows all experts views to be considered by the decision maker. While possibly
47 difficult to interpret, as the final distribution does not represent any one person's or group's
48 beliefs,⁷⁸ it gives a conservative estimate of the extent of uncertainty, which can motivate future
49 research (such as clinical trials) to obtain more evidence to support decision making. More guidance
50 and accessible software would be needed before recommending more advanced pooling methods
51 for use in HCDM.
52
53
54
55
56
57
58
59
60

1
2
3 The practice of weighting experts judgements according to estimates of their expected performance
4 or calibration, particularly as implemented in the Classical Model, has been widely applied and
5 studied, and has been found to improve the accuracy of aggregated expert assessments (Cooke and
6 Goossens 2008; Colson and Cooke 2017, 2018).^{121, 16, 73} However, the method has been largely
7 underexplored in HCDM. After reviewing past HCDM-related elicitations, Soares et al. (2018)
8 conclude that further research is needed to support the use of performance-based weighting in this
9 area.¹ Past elicitations in HCDM that have tested the use of performance-based weighting used four
10 or fewer seed questions to evaluate the approach. Applications are needed that use the
11 recommended 10 or more seed questions to better evaluate if performance-based weighting in this
12 domain has the same benefits that have been identified across other application areas. Finally,
13 future analysis of Classical Model applications would be beneficial to identify strategies for
14 identifying and testing seed questions that have been useful specifically in applications with
15 heterogeneous experts from a variety of disciplines and fields.
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Chapter 6 Reviewing the evidence: Heuristics and biases

Section 6.1 Introduction

Formal models of judgement and decision making hold that judgements of probability and utility should be assessed using all of the information available to the decision maker, with the application of appropriate statistical rules.¹²⁵ However, humans are not perfect information processors. The amount of information processed can be affected by time pressure, limitations in cognitive capacity, lack of motivation and personal desire for a particular outcome. When it comes to probabilistic reasoning specifically failure to recognise when a statistical rule should be applied and unfamiliarity with the processes for making statistical inferences, all mean that probability judgements do not always conform to normative rules.¹²⁶ Experts, being human, are not immune from this. Indeed, even amongst highly educated populations, awareness of how to make simple statistical inferences can be limited.¹²⁷ In the context of HCDM, those practitioners with the greatest relevant knowledge and expertise (e.g. nurses, physiotherapists), may not necessarily have a high level of training in statistics or experience with elicitation.

Humans often make judgements using simple rules of thumb (or “heuristics”).^{126, 128} These strategies are usually effective in appropriately guiding judgement¹²⁹, especially amongst experts who have a large base of experience and knowledge to draw on.¹³⁰ However, in some contexts they can lead to systematic errors, known as “biases”. SEE should seek to elicit probability judgements in a way that minimises the effect of these systematic errors. This is increasingly recognised in the literature on HCDM, where SEE can be used to inform health policy and treatment recommendations.^{13, 44, 59, 97, 131} However, while heuristics, biases and strategies for bias reduction have been widely studied in the broader risk, judgement and decision making literature, there is a dearth of evidence for HCDM and what does exist has not been summarised in this context.

This chapter reviewed evidence relating to the psychological biases of greatest relevance to SEE for HCDM, specifically evidence on how these can be minimised. First outlining key cognitive and motivational biases that have the potential to negatively impact on the quality of expert elicitation for HCDM are outlined (see Section 6.2), and then potential strategies for addressing them (see Section 6.3) through technical measures (see Section 6.4) and behavioural bias reduction techniques (. Reflecting the fact that some behavioural bias reduction techniques have a large amount of evidence to support them, while others are more tentative, techniques are categorised according to

1
2
3 those for which a high degree of consensus exists and those for which evidence is lacking or
4 conflicted. Finally, the key recommendations are summarised in Section 6.7.
5
6
7

8 **Section 6.2 Cognitive and motivational biases**

9 A distinction may be drawn between cognitive biases, which result from how information is
10 processed, and motivational biases that come about due to preferences for particular outcomes.⁷⁶
11
12
13 ¹³² Both have been implicated in systematic overconfidence, which poses a threat to calibration in
14 SEE.
15
16
17

18 *Section 6.2.1 Cognitive biases*

19 Cognitive biases arise when decision makers do not process the full range of information available to
20 them. This may result from limitations in cognitive capacity, time pressure, or a lack of motivation to
21 expend cognitive effort on a task. They may also arise due to decision makers lacking the normative
22 skill to make appropriate probabilistic inferences. In the context of SEE, cognitive biases of particular
23 importance include availability and anchoring and insufficient adjustment; firstly, because they are
24 both implicated in overconfidence, which leads to the systematic underestimation of uncertainty in
25 probability judgments. Secondly, because unlike biases that may result from deficits in substantive
26 knowledge of a subject area, or from a lack of knowledge about how to reason with statistical
27 information, both have the potential to affect expert judgement^{76, 133}.
28
29
30
31
32
33
34
35

36 In making probabilistic judgements people may rely on how easily examples of an outcome come to
37 mind as a guide to how likely it is (the availability heuristic).¹³⁴ While this is often a good guide to
38 frequency, it means that probability judgements can easily be distorted by very recent or very
39 prominent events.¹³⁵ For instance, a clinician may focus on particularly memorable examples of
40 treatment success or failure when making probability judgements, neglecting instances that come
41 less readily to mind. Availability bias has been linked to the systematic underestimation of
42 uncertainty.¹³⁶ Anchoring and insufficient adjustment occurs when people fix (“anchor”) on an initial
43 value, and fail to sufficiently adjust their estimates away from it to provide an accurate judgement.
44 For example in judging the success of an intervention, a clinician may “anchor” on a value provided
45 by a source that they know to be flawed (e.g. a poor quality empirical study), and fail to sufficiently
46 adjust their own experienced-based estimate from this point, despite being aware of the flaws and
47 adjusting in the right direction.¹²⁸ Anchoring has proved challenging to debias, with even arbitrary
48 and irrelevant values being found to affect judgement (see Kahneman¹²⁶ for overview). experts In
49 SEE this can decrease accuracy in judgements of location and central tendency (e.g. mean, median).
50
51
52
53
54
55
56
57
58
59
60

Section 6.2.2 Motivational biases

Motivational biases, sometimes referred to as “self-serving” biases, result from being invested in a specific outcome (e.g. a particular treatment being successful) (see Bazerman, 2008¹³² for discussion). In situations where individuals are aware of potential conflicts of interest and strive to make objective and honest judgements, motivational biases can still distort judgements through rendering some information and experiences more salient (cognitively ‘available’) and easier to recall than others. Confirmation bias for instance, leads individuals to focus on information that is consistent with their existing beliefs and preferences, and subject it to a less critical appraisal than inconsistent information. Desirability bias (also referred to as “optimistic bias” or “wishful thinking”), leads people to overestimate the likelihood of positive outcomes. Undesirability bias meanwhile leads to an overestimation of the likelihood of negative outcomes and worst case scenarios (e.g. due to a focus on taking a precautionary approach). These biases result from motivated reasoning rather than a lack of knowledge or expertise.^{76, 132} Hence, they have the potential to adversely affect the outcomes of structured expert elicitation. In HCDM, where those with greatest knowledge of a particular treatment or procedure, may be those most invested

Section 6.2.3 Overconfidence bias

As a consequence of limiting the amount of information considered by decision makers, both availability¹³⁶ and confirmation bias¹³⁸, may lead to the uncertainty surrounding future outcomes being underestimated. This is known as “overconfidence bias”. It leads to interval judgements and probability distributions that are too narrow (e.g. estimates of 80% confidence intervals containing fewer than 50% of subsequent realisations). Overconfidence is prevalent amongst experts as well as novices^{36, 139}, making it an important consideration for any form of SEE

Section 6.3 Addressing psychological biases in in SEE

Strategies for reducing psychological biases could be said to fall into three categories: technical (e.g. using formal statistical procedures to correct for systematic errors in judgement), directly changing individual behaviour and perceptions (e.g. through training, incentives, feedback) and changing the structure of the judgement or decision task (e.g. how questions are asked).^{140, 141} In practice however they represent two fundamental approaches: post-hoc statistical techniques to make corrections after the fact, most notably through calibration (discussed in *Chapter 5*) (technical); and interventions to change judgement and behaviour (behavioural).

In reviewing approaches for reducing psychological bias (or ‘debiasing’), we restricted our search to studies that provide empirical evidence for the efficacy of bias reduction in the context of SEE. For

1
2
3 this reason, we have excluded papers that suggest approaches, but do not present empirical
4 evidence to support them. We also exclude studies that focus on biases in decision from description
5 (i.e. where choices can be made through analysis of a complete information set) rather than elicited
6 judgements. Relevant papers that did not appear in the searches, but that were cited in the papers
7 identified were examined and included where appropriate. A potential weakness of this approach is
8 that bias reduction techniques that are relevant to SEE, but that do not mention expert elicitation
9 directly, may have been missed if they were not cited in other papers identified through the search.
10 However, a full review of the heuristics and biases literature, which often focusses on novice rather
11 than expert judgement, is beyond the scope of this targeted search.
12
13
14
15
16
17
18
19

20 **Section 6.4 Technical bias reduction strategies**

21 These are commonly discussed with respect to overconfidence. These can involve statistical bias
22 correction and the weighting of experts based on their performance on seed questions, as is the
23 case in Cooke's Classic Model.¹⁴²⁻¹⁴⁴ These approaches do not require interventions at the individual
24 or task level, as the procedures are applied post-hoc. However, they do rely on the availability of
25 appropriate seed question from which the level of experts propensity to overconfidence can
26 measured.¹⁴⁵ This may be relatively easy in contexts where past realisations of the same or similar
27 target variables are available (e.g. probabilistic weather forecasting). In HCDM however it could
28 prove challenging to implement, as contextually similar seed variables with appropriate realisations
29 are not always readily available. Likewise, HTA brings together diverse sets of experts who have
30 specialist knowledge of specific treatments, interventions or procedures. They are not therefore
31 guaranteed to have similar expertise on the subject of seed questions.⁵²
32
33
34
35
36
37
38
39
40
41

42 **Section 6.5 Behavioural bias reduction strategies with consistent support**

43 Given the challenges in applying technical approaches to bias reduction outlined above, it is
44 important for those implementing SEE in the context of HCDM to consider behavioural approaches.
45 In this section we outline bias reduction strategies for which there is consistent empirical support. In
46 the next section we briefly discuss debiasing approaches for which there is conflicting evidence.
47
48
49
50

51 *Section 6.5.1 Consider more information*

52 It has been found that individuals with a greater predisposition towards open minded thinking
53 demonstrate better calibration on judgement tasks.¹⁴⁶ Increasing the amount of information
54 considered by participants may therefore be effective in countering these biases. Behavioural bias
55 reduction techniques that prompt experts to consider more information (increasing the range of
56
57
58
59
60

possibilities considered) have perhaps been the most frequently tested in the context of expert judgement.

Early research with student samples failed to find added value from instructing groups of participants to consider why their estimates may be wrong, or appointing one member to be a “devil’s advocate”.¹⁴⁷ However, more structured approaches have had far greater success.^{138, 148, 149} Soll and Klayman¹³⁸ found that asking student participants to separately give lowest plausible estimates, highest plausible estimates and median estimates for almanac question with which students were likely to have some familiarity, led to lower levels of overconfidence than simply asking for a single 80% confidence interval. It was suggested that making people consider lowest, highest and median estimates sequentially focusses attention on a wider range of possibilities than asking for a single range (e.g. forcing participants to think of reasons why a value might be below (or above) a specific value. Building on this, Haran et al.¹⁴⁸ found that further increasing the number of considerations by asking participants to make judgements about the likelihood of different local seasonal temperature intervals reduced overconfidence. Adding a fourth step to the procedure suggested by Soll and Klayman¹³⁸, Speirs-Bridge et al.¹⁴⁹ found that ranges were widened further when participants (epidemiologists and ecologists) were asked how likely it was that the “true” value would fall within their specified range, and allowed to revise their estimates accordingly. This is consistent with research suggesting that people may be better at evaluating confidence intervals than providing^{150, 151}. More recently, Ferreti et al.¹⁵² noted reductions in overconfidence when environmental science students instructed to (a) actively think of reasons why their initial highest and lowest estimates of sea level rise may be incorrect; and (b) consider their willingness to place hypothetical bets on elicited confidence intervals.

Together, these studies provide strong evidence that structuring tasks in a way that increases consideration of a wider range of possibilities can reduce bias and improve calibration. They demonstrate that confidence intervals should not be elicited as a single stage process. Lower and upper bounds should be elicited individually^{138, 149}, or multiple smaller intervals should be considered individually.¹⁴⁸ Likewise, they show that participants should be given the opportunity to evaluate and adjust their confidence intervals.

Section 6.5.2 Feedback

There is extensive evidence that receiving repeated feedback on one’s judgements both improves accuracy and reduces overconfidence.^{36, 126, 145} Experts, such as weather forecasters, who receive direct and timely feedback on the accuracy of their judgements tend to be well calibrated in their

1
2
3 domain of expertise¹⁵³, although this does not result in a domain general improvement.¹⁴¹ One
4 suggestion for reducing the overconfidence bias in expert elicitation is to provide feedback on a set
5 of practice questions.¹⁵⁴ A challenge in doing this is the fact that domain specific seed variables may
6 be more readily available in some contexts than others (e.g. past realisations in forecasting tasks).
7 Hence, while this approach may be broadly effective in improving the calibration of expert
8 judgement, it could be difficult to implement in some HTA contexts where identifying appropriate
9 seed questions that a diverse set of experts will be familiar with could be challenging. Nonetheless,
10 in cases where these are available, the existing evidence suggests that providing feedback seed
11 questions can reduce overconfidence.
12
13
14
15
16
17
18
19

20 *Section 6.5.3 Avoid unnecessary anchors*

21 Ensuring that elicitation materials do not contain unnecessary anchor values is a “common sense”
22 approach to reducing biases cause by anchoring and insufficient adjustment.⁷⁶ For instance,
23 elicitation tools should not feature pre-set values that participants are then asked to adjust to match
24 their views. However, it may not always be possible to eliminate anchors entirely. In the case of
25 ‘carry over’ effects, for example, experts may use their own judgement on a previous question as an
26 anchor.¹⁵⁵ While there is some evidence to suggest that self-generated median anchors do not
27 threaten accuracy and calibration to the same extent as those that are externally imposed,^{138, 156}
28 Morgan³⁶ advises that measures of central tendency (i.e. the median) only be elicited after lower
29 and upper bounds have been estimated. Hence, while it may not be possible to eliminate all
30 potential anchor values in an elicitation task, a clear recommendation to avoid unnecessary anchors
31 can be made. Likewise, when eliciting confidence intervals, eliciting lower and upper bounds before
32 the median, may reduce the tendency to anchor on the median value.
33
34
35
36
37
38
39
40
41
42

43 *Section 6.5.4 Reduce bias through expert selection*

44 Addressing biases through expert selection means that experts are included or excluded based on
45 their potential susceptibility to bias (see *Chapter 5*). As noted above, motivational biases such as
46 desirability bias and confirmation bias are difficult to eliminate. Restricting participation to those
47 without any conflicts of interest is therefore one recommended approach to reducing motivational⁷⁶
48 biases. In HCDM this may be challenging, as those with the greatest knowledge about a particular
49 treatment or technology may also be those with the greatest vested interest in the elicitation’s
50 outcome.⁴⁴ Rejecting those with any conflict of interest or strong opinions may eliminate those with
51 greatest relevant knowledge. In such cases an alternative strategy is to ensure that a range of
52 viewpoints are represented in the sample, with the intention of “balancing out” or at least diluting
53 the effect of motivational biases.⁷⁶
54
55
56
57
58
59
60

Section 6.6 Behavioural bias reduction techniques with conflicting evidence

Section 6.6.1 Bias warnings and training

Within the HCDM literature that considers heuristics and biases, training is the most commonly referenced approach to behavioural debiasing.⁹⁷ Simply warning experts not to be biased (e.g. by stating that many people make their confidence intervals too narrow) is largely ineffective.^{147, 156, 157} However, in-depth training on the nature of biases, and strategies for avoiding them has been found to be more effective. Where biases occur due to experts not being familiar with rules for using and expressing probabilities, training on how to do so can reduce errors.¹⁵⁸ Likewise, educating participants about biases and explicitly outlining strategies for combating them (i.e. through systematically considering more information) reduced overconfidence in a study of petroleum engineering students.¹⁵⁹ However, this education programme was not effective for reducing anchoring, possibly because the student sample lacked the substantive knowledge of the field to give a more accurate value. Nonetheless, in a study with a general population sample¹⁶⁰ found evidence that interactive training interventions, explaining what anchoring and confirmation bias were, reduced instances of these biases on post-intervention tests relative to pre-intervention tests. These tests comprised tasks from the wider literature found to elicit the psychological biases¹⁶¹⁻¹⁶³ Hence, while the available evidence on the effectiveness of warnings and training for reducing psychological biases is not always consistent, it does provide an indication of the conditions under which bias avoidance training may be effective. Firstly, it must go beyond simple warnings and admonitions not to be bias and explain the causes and consequences of biases. Second, it should provide instruction as to how to avoid bias (e.g. consider why upper and lower bounds may be incorrect). Third, it can only be useful if participants have the substantive expertise to produce accurate responses.

Section 6.6.2 Fixed value versus fixed probability methods

A small number of studies have examined whether fixed value (where one must allocate probabilities to potential values of a target variable) or fixed probability (where one allocates values of the target variable to probabilities) affect overconfidence. In eliciting cumulative probability judgements from students regarding forecast variables with which they were expected to have some familiarity (i.e. local temperature and the Dow Jones), Abbas et al¹⁶⁴ found less evidence of overconfidence using the fixed value method. However, Ferreti et al.¹⁵² found that this resulted in relatively little improvement in performance. Hence, while there is some evidence that fixed value approaches may reduce overconfidence, this is limited.

Section 6.6.3 Face-to-face versus online elicitation

In one recent study it was found that face-to-face elicitation of energy demand with sectoral experts led to lower overconfidence than online elicitation.¹⁶⁵ However, this finding was not replicated in a recent comparison of face-to-face and online SEE.⁵⁹

Section 6.7 Conclusions

The objective of this review has been to synthesise existing knowledge on the effectiveness of different behavioural bias reduction techniques for expert elicitation, focussing specifically on their potential usefulness in the context of HCDM. While the efficacy of some of these approaches remains under-tested, the following five recommendations are supported based on the available evidence.

- Confidence intervals should not be elicited as a single stage process, as doing so leads participants to focus on a narrow set of salient possibilities. Instead, lower bounds, upper bounds and median values should be elicited separately. Eliciting lower and upper bounds before median values may also prevent participants from anchoring on median values.
- Participants should be allowed to evaluate and revise their confidence intervals or probability distributions.
- In selecting experts those with pronounced conflicts of interest should be excluded. However, excluding all participants who may have strong feelings or vested interests in the outcome, may result in the exclusion of those individuals with the greatest expertise in the subject. Hence, it is important to ensure that different viewpoints will be represented.
- Where suitable seed questions are available these may be useful in providing practice feedback to participants on their performance, and thus reduce overconfidence. However, care should be taken to ensure that all participants will have familiarity with the topic of these seed questions.
- Bias training may reduce biases, but only if this goes beyond simple warnings, and explains what bias is and provides strategies for avoiding it.

Chapter 7 Quantities to elicit

Section 7.1 Introduction

Health care decision making (HCDM) is underpinned by (i) evidence of clinical- and cost-effectiveness from randomised trials, to support regulatory approval of drugs, and (ii) decision modelling based on clinical and epidemiological evidence, to support reimbursement decisions. *Chapter 4* provided a review of published applications in this area in order to determine the reasons for methodological choices made in published scientific literature (design, conduct, and analysis) and the challenges faced by the authors when conducting SEE. This chapter discusses different choices available for the specific quantities to elicit in SEE, particularly those related to simple and conditional probabilities of events, as well as parameters to inform survival rates and other time-to-event variables.

Recommendations are made based on statistical theory underlying commonly adopted models and a series of targeted reviews of literature reporting current SEE practice.

Whilst data collected from trials typically aims to inform inference on a single probability, rate and hazard-related parameter, decision models¹⁶⁶ combine a number of these to describe how different courses of action (for example treatments) affect patients' progression through disease stages. Such models typically belong to one of three types.¹⁶⁷ Decision trees, defined using simple and conditional probabilities, describe a set of possible pathways each assigned a probability that is influenced by the treatment being considered. Discrete-time state-transition models (STMs), such as discrete-time Markov chains, define the disease process using a finite set of health states, known to have distinct health and cost implications, and patients transit between states through time. The speed of progression is defined using a set of transition probabilities. Decision models can also be defined in continuous time models, and can be STMs¹⁶⁸ or discrete-event models^{169, 170} (DES). STMs in continuous time are defined using transition rates. DES models use a number of events and use survival distributions to determine the time between events. For alternative treatments, STM and DES models determine the time spent in the different health states. To evaluate differences in lifetime quality-adjusted life years and costs (i.e. cost-effectiveness) from these models, costs and the health related quality of life (HRQoL) are attributed to time spent in each health state (or between events). In decision trees, the cost and HRQoL of each pathway is weighted by its probability.

To inform decision making, either based on single parameters, such as in clinical trials, or based on multiple parameters as in decision modelling, empirical and/or elicited evidence may be used.

Quantitative expression of 'individuals' beliefs regarding a parameter (or parameters) of interest, should be expressed as probability distributions. This chapter gives examples of alternative

1
2
3 quantities that can be elicited to inform the probability-related or time to event-related parameters
4 commonly used in HCDM.
5
6
7

8 **Section 7.2 Overview of probability-, rate- and hazard-type parameters**

9
10 First the main parameters of interest for HCDM and the relationships between them are described.
11
12

13 *Section 7.2.1 Simple probability and conditional probability parameters*

14
15 The probability of a discrete event that an individual may experience once, for example incidence of
16 a disease in a specified time interval and post-operative mortality, can be represented by a single
17 parameter $\pi = p(E)$. Probabilities may be altered by (or associated with) a particular attribute, e.g.
18 treatment, another event or a particular characteristic of the individual. In a conditional probability,
19 the event D is conditioned on the specific value the attribute takes. Conditional probabilities arise,
20 for example, in diagnostics where the sensitivity of a test reflects the probability of testing positive
21 conditional on having the disease, or in logistic regression analyses (say) where the coefficients
22 represent how the outcome of interest is affected by certain attributes, such as age and disease
23 severity.
24
25
26
27
28
29
30
31

32 Where the event of interest has more than two levels, a set of probability parameters is relevant,
33 which are constrained to sum to one given the fact that the categories are mutually exclusive. The
34 probabilities of potential events can be modelled using a Multinomial distribution or, alternatively,
35 expressed as conditional Binomials.
36
37
38
39
40

41 *Section 7.2.2 Transition probability parameters in discrete-time STMs*

42 STMs define a set of health states and describe transitions between them (for example, alive to
43 dead) over a prolonged time horizon. In discrete time STMs, such as Markov chains, the total follow-
44 up period is divided into a number of short consecutive time intervals (cycles). The speed at which
45 transitions between the health states occur is governed by probabilities of the occurrence of the
46 transitions in a particular cycle, termed transition probabilities.
47
48
49
50
51

52 An important feature of discrete-time STMs is that, in any cycle, they typically consider transitions
53 from the current state to one of several health states or competing events. Moreover, individuals
54 may re-enter previously occupied health states. Consider the simple example of a homogeneous
55 (through time) three state-transition model in *Figure 6* **Figure 6**. From state A individuals may transit
56 to either health state B or die (state Death) – these are, in this context, competing events. From
57
58
59
60

health state B individuals are allowed to move back to state A (i.e. backwards transitions are allowed). Death is an absorbing state since once entered, it cannot be exited. To evaluate this discrete-time STM, a transition probability matrix, or TPM (as illustrated in *Figure 6*) needs to be defined, where each row of the matrix must sum to one, so that that all individuals in a particular state at a cycle c are allocated to the allowed states at cycle $c+1$. Transition probabilities may not change over time or, alternatively, cycle-dependent transition probabilities may be specified.

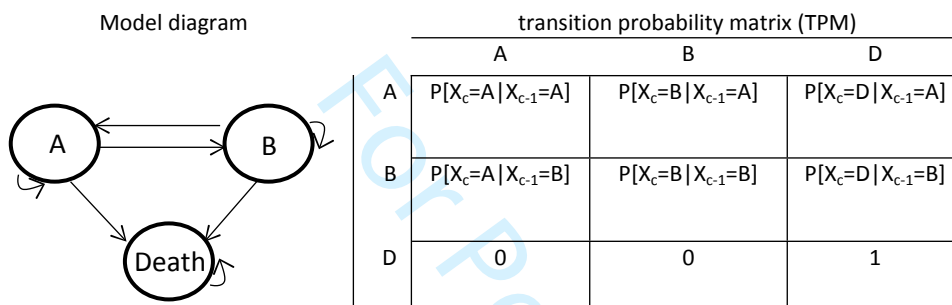


Figure 6 Example of a transition diagram for a state-transition model (STM)

Section 7.2.3 Time to event and survival

Survival and time to event outcomes are defined in continuous time typically using a hazard function, $h(t)$, representing the rate of the event which can take any value above zero. The hazard function can also be used to calculate the survival function, $S(t)$. The mean or expected value of T is the area under the survival curve and can be derived by integrating the survival function.

A number of parametric statistical models can be used to specify the time-to-event distribution, the simplest being the Exponential, which assumes the hazard is constant for all times, $h(t)=\lambda$. In this case, λ can be interpreted as the event rate per unit time. The mean and median times to an event occurring in the Exponential model are respectively $1/\lambda$ and $\ln(2)/\lambda$, where $\ln(2)$ is the natural logarithm of 2, (approximately equal to 0.69). $S(t)$ for the exponential model is $\exp(-\lambda*t)$, where t is the relevant time-frame for the survival estimate (*Table 3*). In many circumstances, a constant hazard is unrealistic, and more complex parametric models than the Exponential are required – examples are the Weibull, Gompertz, Log-Logistic, Log-Normal and the Generalised Gamma, for details see, for example, reference ¹⁷¹. *Table 3* describes key functions and summaries for parametric models commonly used.

Table 3: Summaries of selected survival distributions

Distribution	Hazard function	Survival function P[T>t]		Summaries of the time to event distribution		
				Mean	Median	Variance
Exponential	λ	$\exp(-\lambda * t)$	$\lambda > 0$	$1/\lambda$	$\ln(2)/\lambda$	$1/\lambda^2$
Weibull	$\rho\lambda t^{\rho-1}$	$\exp(-(\lambda t)^\rho)$	$\lambda > 0, \rho > 0$	$\Gamma(1+1/\rho)/\lambda$	$(\ln(2))^{1/\rho}/\lambda$	$\{\Gamma(1+2/\rho) - \Gamma(1+1/\rho)^2\} / \lambda^2$
Gompertz	$\exp\{\alpha + \beta t\}$	$\exp(-\varepsilon(e^{\beta t}-1))$	$\varepsilon > 0, b > 0$	Requires integration	$(1/b) * \ln[(1/\varepsilon) * \ln(1/2)+1]$	Complex

Γ is the Gamma function

Use of continuous time to event parameters in discrete-time STMs

Despite many analyses using a discrete time model, hazard functions from continuous time to event models are often used to derive their parameters. The simplest way to consider a time-varying hazard is to use a value for the hazard that changes between cycles but is constant within each cycle. Transition probabilities for different cycles are then derived to define a discrete time Markov chain that approximates the continuous time estimates using the following relationship:

$$P[t < T \leq t+c] = P[T \leq t+c | T > t] = 1 - S(t+c)/S(t).$$

Section 7.2.4 Continuous time-to-event decision models

Two types of continuous-time models are currently used more often in health care cost-effectiveness analyses: continuous-time STMs and discrete event simulation models. For both types, the clock reset model is adopted¹⁷², where all transition probabilities are expressed in terms of time spent in the current state. Continuous-time STMs are typically represented by state transition diagrams of the type exemplified in *Figure 6*, however these models are defined by a matrix of transition intensities which are derivatives of transition probabilities with respect to time (at time zero) and may vary with the length of time spent in the currently-occupied state or time spent in the study overall. The discrete time-evolution of the probability of transiting between health states can be evaluated using a system of (partial) differential equations defined on the transition intensities called the Kolmogorov equations.¹⁷³

Discrete-event models are informed by the times at which discrete *events* occur. These models define (through a parametric time-to-event model) the hazard of exiting each particular health state (independently of where to), which determines the mean time spent in that particular health state. Separately, additional probability parameters define the arrival state. Discrete-event models are typically evaluated by Monte Carlo simulation.¹⁷⁰

Section 7.2.5 Repeated event rates

Some models may represent events that occur multiple times, for example, rejection episodes in transplant patients, infections or other exacerbations in chronic lung diseases, hypoglycaemic episodes in diabetics, asthma attacks, or seizures in epileptics. These may be represented by a process that counts the number of events over a given time-interval, called a Poisson process governed by a rate parameter, λ , representing the number of events occurring per unit of time (independently of time, in a homogenous process). The time taken between consecutive pairs of events of a homogeneous process is Exponentially distributed with a parameter λ . The common situation in which event rates differ between individuals may be modelled using a negative binomial distribution, which requires that an additional parameter representing inter-individual variation.

Repeated events could reasonably be represented in a discrete-time STM if the time intervals are sufficiently short,¹⁷⁴ or one could construct an STM where transition to states that represent these events can occur repeatedly. The probability of an event per unit of time, or the number of events occurring by time period (which may vary with time) can both be considered in informing these quantities.

Section 7.3 Eliciting probability, rate and hazard-related parameters

The relationships described in Section 7.2 demonstrate that target parameters for elicitation (for example transition matrices in STM or time to event distributions) may be described using a number of alternative quantities. This is because the underlying phenomena described, of progression of a disease process, inherently involves multiple events happening in continuous time and models that do not explicitly consider time (decision trees) or those in discrete-time constitute simplifications of a more complex underlying process. Whilst models simplify reality, this does not mean that the evidence (empirical or elicited) that informs these models cannot be gathered in a way that more closely reflects the underlying processes. For example, a discrete-time STM can be based on inference obtained from empirical time to event data modelled using a particular (continuous time) survival function. Conversely, simpler evidence on a probability (for example, reflecting one point in the survival function) can also be used (alongside assumptions or other evidence) to inform continuous time models that use the entire survival function. To inform one or more parameters of interest, an elicitation need not be restricted to directly eliciting the model parameters; instead, a range of related quantities can be considered.

When eliciting survival functions, for example, if a constant hazard can be assumed, then prior elicitation can be achieved via the mean time to the event. Alternatively, the median time to the event may be more intuitive, since experts may have a clearer picture of the time when the first half of the individuals have experienced the event compared to the second half, some of whom may have very long times to the event. Alternatively, the proportion π of individuals who experience the event by a particular time t can be elicited. In this case, the parameter from the exponential can be calculated as,

$$\lambda = \frac{-\ln(1-\pi)}{t},$$

Which, in turn, can be used to calculate the full survival function $S(t)$. If a more complex survival function reflecting a time dependent hazard is used, its parameters can be elicited from experts, either directly, or more likely indirectly from survival or conditional survival estimates for two or more time points, though there may be multiple ways of accomplishing this.

Section 7.4 Current practices in elicitation

Section 7.4.1 Identification of examples

The aim here is to find examples for illustration, rather than provide a comprehensive review. To identify examples, we initially drew on three reviews that had been used by the authors in recent overviews of SEE for clinical trials or decision modelling.

1. A previous review of studies that included SEE to elicit distributions of parameters included in health care decision models was updated¹. This review excluded studies that generated utility estimates for health states using preference elicitation.
2. A review of experimental studies that used Bayesian survival analyses, most of which were trials in cancer patients and described randomised trials, whilst half were in diseases described as rare¹⁷⁵. Only one of the 28 applications elicited priors from experts¹⁷⁶.

A systematic review of 33 studies that used elicitation of prior beliefs for Bayesian analysis was reassessed in order to extract information on the validity, reliability and responsiveness of quantities that were elicited¹⁷⁷. In order to identify a wide range of applications, other targeted (non-systematic) searches were undertaken. These targeted searches sought Bayesian studies including elicitation of simple and conditional probabilities, survival rates, and other time-to-event variables, as well as reviewing recent technology assessment reports that stated that Bayesian methods, including SEE, were used. The references for all the applications considered in this paper are listed in

Table 4 by type of quantity elicited. Since many reports lack clarity in how the quantities elicited related to the target parameters of interest, the next section discussed is structured according to the type of quantities elicited.

Table 4: List of published examples of eliciting probability, rate and time to-event parameters

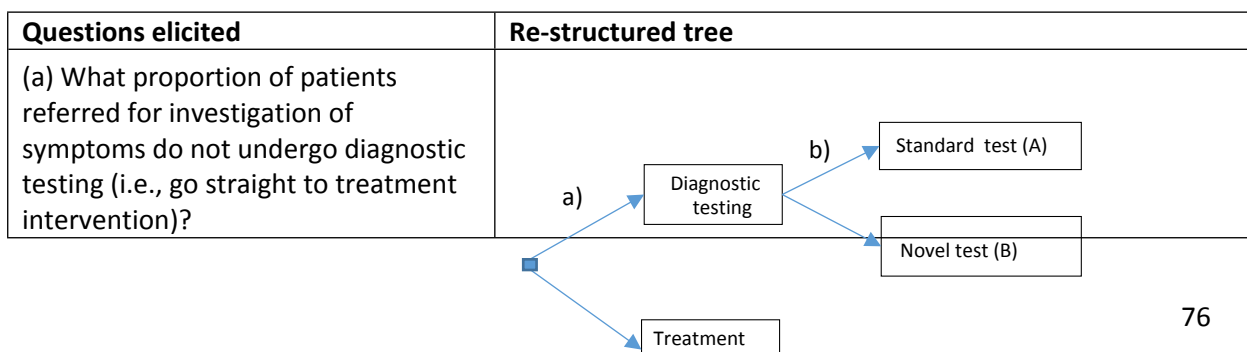
Quantities elicited	Examples
Simple and conditional probability or odds (see Section 7.2.1)	Decision modelling healthcare literature: ^{52, 53, 56, 58-61, 67, 69, 70} Broader health literature: ¹⁷⁸⁻¹⁸⁰
Transition probabilities of an STM (see Section 7.2.2)	Decision modelling healthcare literature: ^{52, 63, 181, 182}
Time to event and survival (see Section 7.2.3)	Decision modelling healthcare literature: ^{51, 52, 57, 58, 64, 65, 69, 70} Broader health literature: ^{48, 176, 183-190}
Hazards or parameters of the hazard function (see Section 7.2.4)	Outside Health: ¹⁹¹
Repeated event rates	No studies found

Section 7.4.2 Simple and conditional probability or odds

In decision modelling, the elicitation of simple probabilities is relatively common^{52, 56, 59-61, 67, 69, 70}, for example, proportion of individuals susceptible to clinical infection⁶⁶, perioperative mortality⁶⁹ or prevalence of cervical cancer recurrence.⁶⁰ Some of these papers elicit independently for different subgroups^{56, 60, 66, 67} or for different durations of treatment.⁶¹ One example of the elicitation of a prevalence was also found in the broader health literature.¹⁷⁸

Applications in decision modelling also elicit probabilities conditional on other events, for example to inform decision trees. An example is provided in Garthwaite *et al*⁵⁸, where an elicitation exercise was designed to inform a decision tree in which more than two outcomes in a single branch were possible.

The authors used conditional probabilities to re-structure the tree and elicited a set of conditional Binomials. A simple extension of the basic structure is presented in Figure 7.



(b) What proportion of the patients referred for testing undergo novel test B, rather than standard test A?	
---	--

Figure 7 Example of eliciting conditional probabilities for a decision tree

The same authors⁵⁸ also considered the need to more formally elicit how the probabilities depend on covariate values. The approach to the elicitation of dependencies was based on conditional probabilities: experts were asked about the quantity of interest by conditioning on a set of values of the covariate(s). These assessments were then analytically transformed to determine regression coefficients using a generalized piecewise-linear model. Specially developed software based on graphical displays was used (Prior Elicitation Graphical Software¹⁶).

Two examples^{53, 60} elicited probabilities related to diagnostic accuracy parameters. Despite requiring sensitivity and specificity, both studies elicited probabilities conditional on test results. One⁶⁰ elicited the proportion of false positives and false negatives (independently), and the other⁵³ the proportion of true positives and true negatives.

Treatment effects on probabilities or odds

A few examples were found in the broader health literature eliciting absolute difference in probabilities,^{184, 188} or ratios of probabilities (i.e. relative risks).^{179, 180}

Section 7.4.3 Transition probability parameters in discrete-time STMs

Soares et al⁵² elicited several quantities assuming conditional independence, but one strategy used in this study extended such the approach to ensure that elicited quantities were consistent with existing relevant data. In their applied example, a particular health state representing healing could be achieved in two ways – either by a wound being left open to heal or via further surgery to close the wound edges. The existing data did not distinguish between these healing types. In order to delineate experts beliefs regarding surgery from beliefs about healing, the probability of closure surgery conditional on healing outcomes was elicited. Denoting the closure surgery event as S and healing as H, the authors elicited: the probability of patients having had surgery given that they healed, $P[S|H]$; and the probability of healing in patients who received surgery, $P[H|S]$. By knowing the unconditional probability of healing, $P[H]$, and applying Bayes' theorem, the probability of receiving closure surgery was calculated as: $P[S] = (P[S|H]*P[H])/P[H|S]$.

Wilson et al¹⁸¹ elicited a total of 12 STM transitions on the progression of undiagnosed melanoma between cancer stages or death. The authors considered each row of the TPM as parameters of a Multinomial distribution. Experts were asked to distribute a cohort of 100 patients according to the stages they would be in six months later. These values were described as medians, and the software restricted the values introduced to sum to 100. The participants were then asked to elicit the lower and upper limits of 95% credibility intervals (CrIs) for each stage; for these no restriction was imposed on the values provided. The participants elicited in the same way for all other starting health states.

Cao et al⁶³ took a similar approach to Wilson et al¹⁸¹ but elicited membership of each health state at a particular point in time. Experts were initially presented with a diagram of the model with relevant empirically derived numbers for standard of care, and were then asked to revise these for a new care setting, as exemplified in *Figure 8*. Cao et al⁶³ elicit in 'discrete-time' but use the elicited quantities to inform a continuous-time model with the same structure. They argue that transition rates are complex and not-observable quantities, and hence did not elicit these directly.

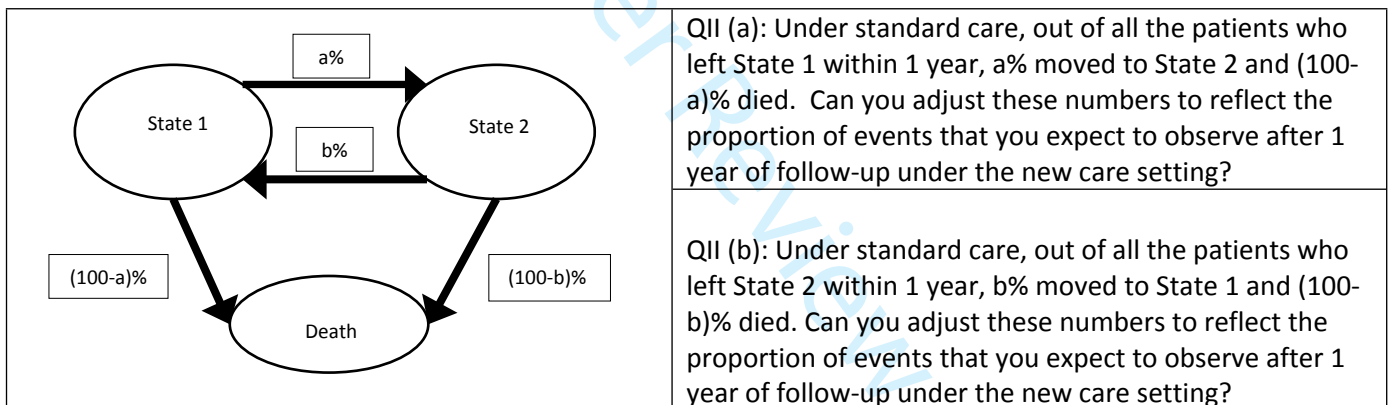


Figure 8 Example of eliciting transition probabilities for a novel setting in a discrete-time STM

Vargas et al¹⁸² also inform the transition probabilities of an STM, but the publication provides little detail on how these quantities were elicited.

Section 7.4.4 Time to event and survival

Some studies elicited summaries of time to event distributions, particularly median⁶⁹ and mean survival.⁵⁷ Survival functions have also been elicited to inform transition probabilities in a discrete-time STM.^{51, 52, 64, 65} Some studies elicited event probabilities at a single time point.^{51, 64, 65} For example, Leal et al 2007⁵¹ asked: "If 100 hypertrophic cardiomyopathy patients were stratified as low/medium risk at the age of 18, how many would be classified as high risk at age 50". Some

1
2
3 applications elicited multiple event probabilities without eliciting dependency between them.

4 Examples are Poncet et al,⁶⁵ that elicited separately for different subgroups of patients, and Speight
5 et al,⁶⁴ that elicited probabilities of sequential cancer progression (e.g. pre-cancer to stage I, stage I
6 to stage II, and so on).
7
8
9

10
11 To explore time-dependency and derive a full survival function, some studies elicited conditional
12 probabilities at two or more points.^{52, 58} One such study⁵² elicited a first point in the survival function
13 at six months and, for those who did not have the event at that time point, the proportion who
14 would have the event between six and 12 months. By assuming conditional independence between
15 the elicited quantities, the authors were able to incorporate time dependency in the decision model
16 without generating incoherent probability statements. They also argued that, even if hazards are
17 found to be very similar in both periods (i.e. no evidence of time dependency), experts may be more
18 uncertain about outcomes in the longer term, so that it may still be important to elicit for separate
19 time periods. Garthwaite et al⁵⁸ used a similar strategy but partitioned the time scale into 4 intervals.
20
21
22
23
24
25
26
27

28 *Section 7.4.5 Treatment effects on time-to-event distributions*

29
30 Experts may judge survival with a comparator treatment informative of survival with the treatment
31 of interest, i.e. hazards should not be assumed independent.¹⁸³ To elicit priors for relative treatment
32 effects (hazard ratios), a number of authors^{176, 183, 184, 186-189} elicit the absolute difference in event
33 probabilities (at a single time point) between treatment and comparator. Most of these authors
34 convert the absolute difference onto the log hazard scale assuming a value for the baseline hazard
35 (example in ^{188, 189}). The study by Ren et al¹⁸³ considered eliciting absolute differences in survival
36 under time dependency, and proposed eliciting the following quantities:
37
38
39
40
41
42

- 43 • survival with the comparator at a particular time point, t_0 ,
- 44 • the difference in survival for the comparator between times t_1 and t_0 (where $t_1 > t_0$),
- 45 • the difference in survival between the treatment of interest and the comparator at t_0 , and
- 46 • the difference in survival for the treatment of interest between times t_1 and t_0 .
- 47
48
49
50

51
52 Note that this method also relies on a form of conditional independence.
53
54

55 Other authors^{52, 190} asked experts to elicit absolute survival probabilities with the treatment of
56 interest conditional on a given fixed value for the comparator (in Soares et al⁵² the value selected
57 was the elicited mode for survival with the comparator whereas in ¹⁹⁰ the fixed value was provided
58
59
60

1
2
3 to the participant). Soares et al⁵² elicited relative effectiveness for multiple treatments
4 independently, and White et al¹⁹⁰ evaluated treatment effects in the presence of possible
5 interactions across different patient groups.
6
7

8
9
10 Dallow et al⁴⁸ argue that to better manage the tendency for experts to be over-optimistic, experts
11 should be first asked to judge the probability that the drug has a true positive effect and then to
12 judge the distribution of this effect size under the assumption that the drug does have a favourable
13 effect. They then formed a mixture distribution to represent the overall prior for the treatment
14 effect. They then formed a mixture distribution to represent the overall prior for the treatment
15 effect. They then formed a mixture distribution to represent the overall prior for the treatment
16 effect.
17
18

19
20 Chaloner et al¹⁸⁵ aimed to specify a bivariate distribution for two hazard ratios (two treatments in
21 relation to a common comparator), eliciting survival probabilities with the support of a graphical
22 dynamic tool. Analogously to Soares et al⁵², the authors initially ask experts to elicit absolute survival
23 probabilities for each treatment conditional on an initial model value for the comparator
24 (conditional independence is assumed throughout); specifically, experts are asked for their upper
25 and lower quartiles. The relationship between survival probabilities elicited for a treatment and
26 control and the hazard ratio under a proportional hazard model -- $\{\log(-\log(1-p))=\beta + \log\{-\log(1-p_0)\}$,
27 where p_0 and p are the elicited survivals for the control and treatment, respectively, and β the
28 treatment effect -- allow the elicited summaries to be expressed in terms of treatment effects. These
29 are used to calculate initial values for parameters of a type B bivariate extreme value distribution
30 describing the treatment betas. The distribution is defined on the means and variances of the
31 marginal distributions and on an m parameter, with $m=1$ reflecting independence between
32 coefficients. To pick an initial value of m , the expert is additionally asked about the probability that
33 the survivals for each treatment being larger than their respective marginal medians. If the two
34 parameters are independent, this probability is 0.25. To reflect correlation (note that only positive
35 correlation is allowed), values for this probability can be higher than 0.25 (up to 0.5) and the value
36 for m can be directly determined from these. The expert is presented with plots of each marginal
37 distribution (for the probability parameter), a contour plot of the joint prior distribution of the
38 survival probabilities with approximate confidence regions, and a dialog box with five sliders (for
39 each parameter of the bivariate distribution). Changing the value of m in the slider does not change
40 the marginal distributions but it does change the contour plot. The sliders allow the expert to adjust
41 interactively the parameter values and see the consequences directly. The authors recommend
42 repeating the elicitation process for a different follow-up interval; under proportional hazards, the
43 distributions on the regression coefficients should be equal.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Section 7.4.6 Elicitation of hazards or of parameters of a time to event distribution directly

The only example of directly eliciting time to event distributions related to reliability assessment in engineering.¹⁹¹ Singpurwalla presented methods for eliciting a single Weibull distribution, and proposed directly eliciting its shape parameter (α), which characterises whether hazards increase, decrease or remain constant over time. We note also that the shape parameter can be expressed as a function of the hazard ratio h_2 associated with a doubling of time: $\alpha = 1 + \log_2(h_2)$, so the shape might be derived by eliciting h_2 . Singpurwalla argued that the scale parameter is difficult to interpret and instead of eliciting it directly, chose to elicit median survival time, making the simplifying assumption that the two parameters are independent. Indeed, the scale is related to the mean and median but only has a simple relationship (independently of the shape) when the shape parameter is 1 (Exponential distribution).

Section 7.5 Steps and considerations in defining the quantities to elicit

The examples reviewed identify a range of ways to elicit probability, rate- and hazard-type parameters used in current practice. Whilst most did not directly elicit the parameter of interest, there was often little justification for the quantities chosen for the elicitation. Following critical review of the literature described in Section 7.4 and investigator experience of SEE, some generic guidance on how to determine appropriate quantities to elicit was developed and agreed during discussions within the research team and is presented here.

Step 1: Develop a clear understanding of the statistical or decision model, so that, ultimately, quantities elicited are fit-for purpose (i.e. accurately represent the relevant context, including, population, setting, interventions, outcomes and time horizon).

Prior to defining the quantities to elicit, the target parameters of interest for decision-making must be defined. This can be done independently from the elicitation and requires two assessments. The first is to specify the model (statistical model or decision model), by defining a list of input parameters required and outputs produced by the model, to inform the decision. Models are developed specifically to best represent decision problems involving particular types of health care strategies (for example, diagnostics, drugs, complex interventions, screening strategies, infectious disease prevention). Within a well-defined set of model input parameters, the second assessment identifies which inputs have strong empirical evidence and, of those that do not, which might benefit from explicit priors being elicited from experts. The level of uncertainty, and whether it ultimately

1
2
3 impacts on the health care decision, is a critical consideration here. Value of information approaches
4
5 ¹⁹² can help to identify those parameters that are most influential, and prioritise parameters for
6
7 elicitation.
8
9

10 Step 2: Consider breaking down (decompose) the target quantities for the elicitation into quantities
11 that are simpler and that reflect what experts are likely to observe
12
13

14
15 The parameters defining statistical or decision models can be complex, for example, hazards and
16 intensities are difficult concepts. In contrast, eliciting a probability of having the event of interest by
17 a set of specific times, for example, survival up to one year, is more intuitive, and might represent
18 data that experts observe directly. From the resulting elicited probabilities, and under some
19 assumptions, survival and hazard functions can be constructed.
20
21
22
23

24
25 Step 3: Consider what sets of related target parameters are required, and define quantities to elicit
26 in a way that ensures coherence between the quantities elicited and the parameters they inform.
27
28

29
30 Target parameters, and the quantities elicited, may be related in a number of ways: total survival is a
31 compound function of simple and conditional probabilities, number of people infected is a
32 combination of exposure rates and infectivity, and positive predictive value is a function of
33 sensitivity and prevalence. Relationships between target parameters, between quantities elicited,
34 and between the target parameters and the quantities elicited need to be understood and
35 accommodated, so that statistical coherence of the priors generated is ensured. Eliciting two points
36 in the survival curve unconditionally does not guarantee consistent results, and should be avoided.
37
38
39
40
41
42

43 Coherence is important when eliciting multinomial outcomes or discrete-time STM probabilities. If
44 multinomial probabilities are elicited independently with uncertainty, they may not sum to one. As
45 an alternative, a Multinomial can be re-expressed a set of conditional Binomials.⁵⁸ Or Multinomial
46 probabilities can be elicited directly by eliciting expected proportions in each health state and an
47 effective sample size that informs uncertainty.⁷⁸ Consistency in the quantities defined and elicited is
48 important, and can either be ensured by design or verified using consistency checks built into the
49 elicitation tool.
50
51
52
53
54

55
56
57 The relationships and constraints identified above can generate dependencies. But other forms of
58 dependency are also relevant, such as dependencies between quantities or between quantities and
59
60

1
2
3 known covariates. Dependencies should be considered and accommodated, either by re-expressing
4 target parameters as conditionally independent quantities, or by formally eliciting dependency.
5
6 Note that dependencies between quantities may arise from some experts being prone to eliciting
7 higher values across the board than others.
8
9

10
11 Step 4: Consider what the expert may not observe, e.g. censoring, heterogeneity.

12
13 Analogously to empirical studies, it is common for clinicians not to observe times to event for all
14 patients in their clinical practice: there may be competing events that removes patients from the risk
15 of the event of interest, or patients may change practitioner or become hospitalised and not be
16 observed after a certain time point. Experts will find it difficult to elicit quantities under heavy
17 censoring. Hence, we believe that experts asked about summaries of a time to event distribution will
18 find the median more intuitive than the mean, and that shorter time points may be necessary when
19 eliciting survival (at the expense of need for more extensive extrapolation).
20
21
22
23
24
25

26
27 Step 5: Where more complex quantities need to be elicited (for example, bivariate treatment
28 effects), consider using dynamic graphical displays.

29
30 Graphical displays used in the elicitation did help the experts to provide parameter values whilst
31 visualizing probability distributions on quantities that were more intuitive to them. The graphs also
32 provided useful instant feedback.
33
34
35

36
37 Step 6: Where alternative quantities of interest can reasonably be elicited, or graphical displays
38 used, pilot the different options with a small number of the relevant experts.

39
40 Piloting is essential to both choose the most intuitive set of quantities for the elicitation, and to
41 optimise the quantities using the principles above (for example, validate the time point at which a
42 survival probability can be reasonably asked under censoring), and optimise the wording of the
43 questions for clarity.
44
45
46
47

48 **Section 7.6 Discussion**

49
50 Survival or time-to-event models, and transition probability/intensity matrix based models, pose
51 challenges for elicitation, as inputs are typically complex constructs that may involve several
52 correlated parameters. Instead of target parameters being directly elicited, these may be
53 decomposed into related quantities that are simpler and observable to experts. This chapter aims to
54 review current practice of elicitation to identify the quantities elicited for these types of parameters.
55 Given the low specificity of search terms, efficient targeted searches were employed, giving a good
56
57
58
59
60

1
2
3 overview of what is published up to date. Current practice is heterogeneous, with different
4 quantities used to elicit the same type of target parameters, for example survival distributions and
5 how these vary with treatment and other predictors. Some general guiding principles for
6 determining appropriate quantities to elicit are developed here. Further research could refine these
7 recommendations, particularly if there are multiple options or if implementation of principles is
8 unclear.
9
10
11
12
13
14

15 One may need to retrieve hazards, survival functions and relative effectiveness parameters from the
16 elicited priors. Whilst this is not the focus of this chapter, it is an important aspect that would
17 benefit from further guidance, as there may also be multiple ways of estimating the whole
18 distribution to the level of detail required for decision modelling.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Chapter 8 Three methodological experiments on the elicitation of subjective probabilistic belief

Section 8.1 Introduction

Research looking at important methodological choices in elicitation is mostly inconclusive, for example on which experts to engage, how to most appropriately elicit distributions, or whether group interaction increases the accuracy of group judgements. The underlying challenge of methodological research in this area is that beliefs are inherently unobservable: the accuracy of elicited judgements in representing the experts beliefs cannot be directly established as heterogeneity in knowledge (i.e. the fact that individuals' beliefs differ) cannot be easily disentangled from the lack of 'normative' skills (i.e. individuals not being able to represent their beliefs accurately in probabilistic terms), and there is no "gold standard" perfect elicitation procedure.

This chapter describes the approach used in three experiments aimed at generating evidence that can support some of the methodological choices in elicitation. It first describes the general approach for the experiments, and then details aspects of their design. The experiments explore alternative methodological choices in SEE:

- The first experiment compares two methods of elicitation (bisection versus chips and bins),
- The second evaluates whether individuals are able to accurately extrapolate from their knowledge-base to different populations, and
- The third experiment (comprising of two separate sub-experiments) explores how individuals revise their own probability assessments when presented with Delphi-type group summaries.

The objectives were chosen to address some of the key methodological challenges for conducting elicitation in HCDM, reported in Chapters 4-7.

The objectives in experiments 2 and 3 differ to those outlined in the funding proposal. We had initially thought of exploring the accuracy of consensus-based methods and use the experimental set-up to evaluate alternative methods of mathematically pooling priors elicited from individual experts. Consensus-based methods can be affected by many factors, including facilitators' input, individuals' experience, their ability to adjust (or extrapolate) their knowledge and beliefs, and the composition of the sample of experts whose consensus is sought, including their personalities,

1
2
3 probabilistic accuracy and between-expert agreement. When planning for experiment 1, we realised
4 that a sample of the size we would be able to recruit, would not allow for meaningful inferences
5 under these objectives (as randomised groups would have to be formed). Moreover, the review in
6 Chapter 5 highlighted that there was no good evidence on a key question: how do individuals revise
7 their own assessments after some form of interaction. This is crucial for all methods that require
8 individuals to revise their assessments (both consensus methods and controlled interaction methods
9 such as Delphi). For these reasons, the objectives of experiments 3 were updated to explore this,
10 and used the more controlled interaction in a Delphi-type environment. Experiment 2 was a natural
11 extension from experiment 1 and aimed to explore how experts deal with heterogeneity in
12 knowledge, which was one of the objectives set-out in the funding proposal.

13
14
15
16
17
18
19
20
21
22 The variation from the protocol was approved by the project team and by the advisory group, and
23 the revised objectives are described in further detail in Section 8.3.3 and Section 8.3.4.

24 25 26 27 **Section 8.2 General approach to the experiments**

28
29 The aim of these experiments was to compare alternative design choices in SEE. An approach in
30 which the individual's knowledge is determined by observations from a simulated (virtual) learning
31 process (Wang, 2002) is employed. This process is reflective of the learning process (by observation)
32 we may expect of experts in health, typically health carers that observe patients over time.
33 Additionally, if the simulated learning process constitutes the single source of information
34 participants receive, elicited probabilities can therefore be directly compared to the posterior
35 probability distribution implied by the observed dataset (and any prior beliefs, even if vague) – i.e.
36 accuracy can be measured. The ability of this experimental approach to control the subjects'
37 knowledge means that the conditions of the experiment can be determined, which allows testing
38 specific hypotheses, and standardising the task to reduce between-participant variation.

39
40
41
42
43
44
45
46
47 Finally, in this approach, and given that the same information is provided to all participants,
48 systematic differences in the elicited distributions across individuals can directly be attributed to
49 different levels of normative ability, reflecting the skills needed to extract information from
50 observations and quantify the resulting beliefs in probabilistic terms. The particular dimension of
51 normative ability captured in these experiments is referred to as 'probabilistic' accuracy.

52
53
54
55
56
57 The following methods section summarises the protocol for the experiments which is presented in
58 full in **Error! Reference source not found.**

Section 8.3 Methods

Section 8.3.1 Overview of the experimental approach

The game and target question

Participants were shown a number of observations generated randomly from a statistical model, which were recorded. The context was that of an abstract generic medical problem so that all knowledge was acquired from the game and not influenced by external information or separately acquired prior beliefs. In summary, participants in the experiment were asked to act as practitioners over a number of clinic days (the number of days is defined in each experiment, see details ahead). In each day, participants cared for a variable number of (simulated) patients (between 6 and 13, randomly sampled). Participants were presented with two pills and asked which they would use to treat patients that day (*Figure 9*).

Today there are 11 patients in clinic. Please choose the pill you want to use today.

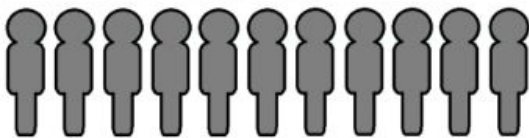
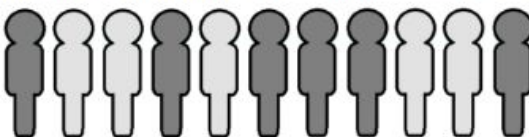


Figure 9 Snapshot of the R shiny app (1)

Once the participant chose the pill to use that day, a new screen appeared (*Figure 10*) showing how many of the patients achieved symptom relief (number of successes).

After treatment with pill J, 5 patients were relieved of symptoms.



Please click on the Next button to continue.

Next

Figure 10 Snapshot of the R shiny app (2)

After observing a number of clinic days, participants were asked about the effectiveness of the most effective pill – the target question for the elicitation – phrased as: “About pill <<>>, the most effective of the two pills... If you were able to treat all patients in the population, what proportion would you expect to become symptom-free?”.

1
2
3
4
5 Before running the game, all participants were shown 6 pre-determined runs of the game, 3 on each
6 pill, to mitigate against sensitivity of results to different uninformative or vague priors that subjects
7 may use. These were equal across all subjects.
8
9

10 11 Participants and monetary incentives

12 Given that the subjects' knowledge-base is 'built up', it is more important that participants are
13 representative of the type of normative skills we expect from experts in health care. Hence, students
14 at the University of York with a health background or undertaking clinical training were recruited.
15 The recruitment target was 64 participants, based on the requirements for experiment 1 and
16 constrained by available funding. No formal sample size calculation was undertaken given the lack
17 of evidence on the potential magnitude of effect size and its variance. Information on participants
18 was collected, such as age, gender, the Berlin Numeracy Test¹⁹³ and the Scott and Bruce's General
19 Decision-making Style questionnaire.¹⁹⁴ Monetary incentives for performance were used as students
20 may not be as motivated to complete the task as might be expected of professionals. The reward
21 was individualised, and incentivised both the learning from playing the game and accuracy in the
22 elicitation– description of the incentives is provided in the experimental protocol, Supplementary
23 material 2.
24
25
26
27
28
29
30
31
32
33
34

35 Metrics for comparison

36
37 The aim was to find a metric or set of metrics that allow comparison of the elicited distribution
38 against the posterior distribution, calculated from the prior and data provided to participants. Bias
39 was defined as the mean of the elicited (and fitted) distribution minus the mean of the true
40 distribution; uncertainty, defined as the ratio of the standard deviation of the elicited distribution to
41 the standard deviation of the true distribution; and the Kullback-Leibler (KL) divergence, a measure
42 of the information lost when the true distribution is approximated by the elicited distribution¹³.
43
44
45
46
47
48
49

50 Other aspects of conduct

51
52 The experiments were conducted in a number of face-to-face sessions, each lasting around 2 hours.
53 Subjects were, however, asked to complete all the tasks individually (i.e. the game and the
54 elicitations). The games and elicitations were conducted in a tool developed for purpose in the
55 SHINY package for R.¹⁹⁵ A number of pilot exercises were undertaken to evaluate the feasibility of
56 the experiments, time to completion, the optimal experimental conditions (e.g. value of the
57
58
59
60

1
2
3 probability parameters), and to test the tool developed. A bespoke training package was developed
4 and delivered to participants. University level ethics approval was sought from the host institution
5 and granted prior to the conduct of the experiment.
6
7
8
9

10 11 *Section 8.3.2 Experiment 1: comparing different methods of elicitation*

12 The aim of this experiment was to assess how well the elicited probability distributions derived from
13 two different methods of elicitation, the bisection and the chips and bins (or histogram), both widely
14 used for SEE in HCDM,¹ reflect description(s) of bias and uncertainty. The bisection method is a VIM
15 and asks experts to give the three quartiles of the distribution. The chips and bins method is a FIM
16 that defines a larger number of intervals (typically up to 20 bins) and asks the expert to distribute a
17 fixed number of chips across these intervals. The more chips placed in a particular interval, the
18 stronger the belief that the true value of the quantity of interest lies in that interval. Both methods
19 were preceded by asking participants for bounds (see **Error! Reference source not found.** for further
20 details on how the two methods and questions regarding bounds were implemented). In HCDM, the
21 general understanding is that the bisection method returns wider representations of uncertainty
22 (argued to be more appropriate in representing within-expert uncertainty) than the chips and bins,
23 although the latter has been said to be more intuitive for less quantitative experts to grasp.^{1, 59}
24
25
26
27
28
29
30
31
32
33

34 Given this experiment uses the set-up described in Section 8.2, where participants observe data
35 directly on the target question, this experiment focusses on how well individuals express their
36 uncertainty, and not on bias. Different levels of uncertainty are defined by varying the number of
37 clinic days participants observe with the pill of interest and by assuming, or not, over-dispersion in
38 the probability parameter. The number of clinic days observed is 25 in the higher precision scenario
39 and 10 in the lower precision scenario. The high precision scenarios uses a Binomial model (no
40 overdispersion), and the low precision scenario a Beta-Binomial model with an effective sample size
41 ($\alpha+\beta$ of the Beta distribution) assumed a value of 2. In this context, over-dispersion implies
42 participants observe greater variation in the probability of success between clinic days.
43
44
45
46
47
48
49
50

51 The experiment used a full factorial design with all 4 combinations of the two levels: precision
52 scenario and method of elicitation. The experiment used a repeated measure design, with
53 participants' beliefs elicited for all four combinations. In each repetition, different probability values
54 (p_0 values) were used of: 0.3, 0.4, 0.6 and 0.7. Hence a 4x4 Graeco-Latin design was implemented.
55
56
57
58
59
60

1
2
3 At the end of the experiment, participants were also asked if they found each of the methods
4 (generally) easy to complete (Response options: easy, challenging, very difficult) and if they had any
5 preferences regarding the elicitation method used (“If, in a future elicitation, you were given a choice
6 between chips and bins or bisection, which would you choose?” - Response options: chips and bins,
7 bisection, indifferent. *Please justify your choice* - Response in free text). An open text box for further
8 comments was also provided.
9
10
11
12
13
14

15 *Section 8.3.3 Experiment 2: Are individuals’ able to ‘extrapolate’ from their knowledge-base?*

16

17
18 Variation in elicited judgements across experts may arise from experts having a different knowledge-
19 base from which they form their beliefs. To provide a probability distribution for a common target
20 quantity, individual experts need to adjust (‘extrapolate’) their beliefs using some form of analytical
21 reasoning. A simple example is where a health care professional observes a sample comprising two
22 subgroups, but the subgroup distribution observed is different to that of the overall target
23 population of the decision question. The expert has a knowledge base that is relevant (in that he/she
24 observes both subgroups) but when their belief at the population level is elicited they need to adjust
25 (or re-weight) what they directly observed. This experiment examines how well individuals make
26 such adjustments, by looking at whether, in the case of extrapolation, accuracy of the extrapolation
27 is associated with non-extrapolated accuracy and with the extent of the extrapolation (difference in
28 the split between observed and target populations).
29
30
31
32
33
34
35
36

37 This experiment used a different set-up from that of experiment 1. At each clinic day, participants
38 were shown a number of patients that were sampled randomly. However, here patients were from
39 two groups (S1 and S2) (*Figure 11*). Participants were told that one group had a better chance of
40 symptom relief than the other, but at the beginning of the experiment participants did not know
41 whether this was S1 or S2. The number in each subgroup was generated randomly from a Binomial
42 distribution on each clinic day. Different probability parameters for this Binomial were examined
43 reflecting odds of being in groups S1:S2 of 80:20, 70:30 and 60:40.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Today there are 13 patients in clinic.

Of those 13 patients 7 are from S1,



and 6 are from S2 .



Please choose the pill you want to use today.

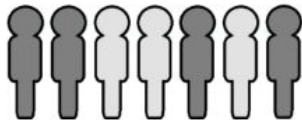


Figure 11 Snapshot of the R shiny app (4)

Two new pills were available in clinical practice, and these were used by the participant in exactly the same way as in experiment one. Once the participants chose the pill they wished to use, they observed outcomes for the two subgroups (Figure 12). The number of successes was governed by a Binomial model (as in the high precision scenario in experiment one), with probability 0.35 in the largest subgroup and 0.65 in the smallest. Participants observed 15 clinic days with the pill of interest.

After treatment with pill A:

3 patients out of 7 in S1 were relieved of symptoms.



4 patients out of 6 in S2 were relieved of symptoms.



Please click on the Next button to continue.



Figure 12 Snapshot of the R shiny app (5)

Participants were asked to provide their beliefs first for a target population with same split as they had observed and then for a 50:50 split (Table 5). Participants were not told which split was assigned in their observed experiment.

Table 5 : Wording of questions in experiment 2

*“About pill A, the most effective of the two pills...
 Suppose that the patients you have just observed were representative of the general population, that is, the split of S1 and S2 patients is unchanged.
 If you were able to treat all patients in the population with this pill, what proportion would you expect to become symptom-free?”*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

“Now suppose that that the general population is different to the sample you observed. Suppose that subgroup S1 makes up 50% of the population and subgroup S2 makes the other 50%. If you were able to treat all patients in the population with the same pill, what proportion would you expect to become symptom-free? “

In terms of design for this experiment, participants were randomised to receive one of the three scenarios described above; randomisation was by block according to the method of elicitation.

Section 8.3.4 Experiments 3: To understand how individuals review their own probabilistic assessments when presented with Delphi-type summaries.

Experiments 3 aimed to gain an understanding of how individuals revise elicited distributions when presented with Delphi-type summaries, loosely based on the recent modified EFSA Delphi method that allows quantifying of uncertainty in the form of probability distributions.¹⁷ As with the original Delphi, this method makes use of multiple (sequential) questionnaires (called ‘rounds’), and at every round experts are fed back an anonymised summary of the information collected in the previous round. This form of interaction between experts is controlled, and advocates of the Delphi method argue that it allows for the benefits of the sharing of information without the risks of personal factors influencing judgements inappropriately. In contrast to the original Delphi, the modified EFSA version does not aim to achieve consensus; instead, after all rounds are completed, a final distribution is obtained using mathematical aggregation with equal weighting.

Despite the benefits of reduced interaction between elements of the group, how individuals revise their estimates in a Delphi process is not well understood.¹⁹⁶ This is the focus of the two sub-experiments conducted here.

Section 8.3.4.1 Experiment 3.1: Is the extent of revision associated with discrepancy with the group, and does the individual’s probabilistic accuracy determine the extent of revision?

This experiment aimed to evaluate if low performers (in terms of probabilistic accuracy) revised their answers to a greater extent (to approximate the group’s distribution) than high performers. If this were true, then over multiple rounds of Delphi the group distribution would be expected to converge to a more accurate distribution than initial estimates mathematically pooled—i.e. the iterative process may dilute the effect of low or extreme performers. How different features of the group distribution shown to participants determined the extent of revision, was also explored.

This experiment uses the set-up for the high precision scenario in experiment 1, where participants run the game and are first asked to elicit the directly observed target question. Participants then

received one of three types of group summary of the quantity of interest: concordant with their initial probability distribution, or discordant and either more or less precise than their own (*Figure 13*). The group distributions were hypothetical (groups were not formed), and were defined relative to the individual's elicited distribution (but always towards the true distribution).

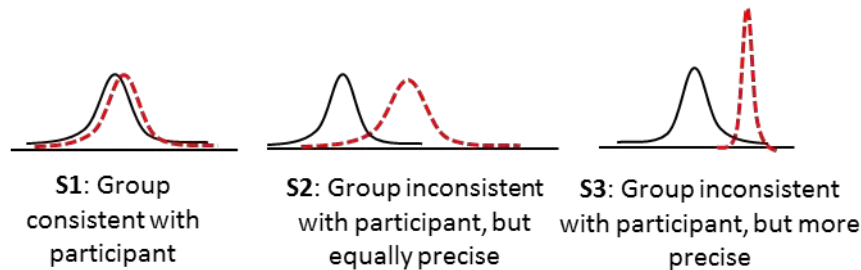


Figure 13 Illustrative example of the scenarios evaluated in experiment 3.1

After observing the group summaries, participants were asked if they wished to revise their elicited distributions in light of the group summary.

In terms of the design of the experiment, participants were randomised to receive one of the three scenarios described above; randomisation was undertaken by block according to the method of elicitation.

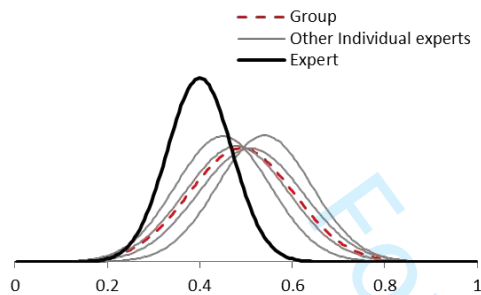
Section 8.3.4.2 Experiment 3.2: How does between-expert variation within a group affect individuals' revision?

In this experiment, participants were presented with disaggregated results for each member of a group to examine its effect on individual's revision. For such, two scenarios were generated based on exactly the same linearly pooled group distribution (which was discordant with the individual's): one scenario defined higher levels of within-expert uncertainty (but concordant central estimates) and another scenario defined higher levels of between-expert variation (with discordant central estimates but higher precision in individual distributions) (*Figure 14*). This experiment aimed to examine how individuals revise their estimates when presented with either of these scenarios.

This experiment used the same set-up as experiment 3.1, however, instead of a single group distribution, participants were presented with the distributions of three other individuals and then asked whether they would like to revise their judgements in light of this information. The group distribution was not presented to the participant. The mean for the group was discordant with the

individual's elicited distribution and was the same for the two scenarios – see **Error! Reference source not found.** for more details on how this was operationalised.

Participants were randomised to receive one of the two scenarios described above; randomisation was by block according to the method of elicitation.



Scenario 1: Others in group are concordant amongst themselves but present a higher within-expert uncertainty

Scenario 2: Others in group are discordant amongst themselves but present lower within-expert uncertainty

Figure 14 Illustration of the scenarios defined for experiment 3.2

Section 8.4 Methods of analyses

Section 8.4.1 Outcomes and metrics used

All experiments required a measure of participants' accuracy in elicitation. Experiments 3.1 and 3.2 also require a measure of the likelihood of revision and of the extent of revision. Accuracy was evaluated by comparing the elicited distribution to the theoretical posterior distribution implied by the prior and data provided to participants.

The elicited summaries consisted of a set of quantiles and probability masses that define points on the CDF of the quantity of interest. The target quantity for elicitation was a probability, thus by definition takes a value in the range $[0,1]$. Fully-specified distributions were derived from the elicited summaries by fitting a Beta distribution to the bounds (assumed to represent the 98%

confidence interval) and the elicited summaries, using least squares on the elicited CDF points. Distributions were fitted in R, using the *fitdistr* function in the SHELF package.¹⁹⁵

Methods for deriving the posterior distributions varied depending on the precision scenario, i.e. on whether the Binomial (high precision) or the Beta-Binomial model (low precision) was used. With the Binomial model (used in experiments 1, 2, 3.1 and 3.2) the posterior distribution was obtained using conjugacy with a Beta prior. In experiment 2 (Section 8.3.3), the posterior distribution of the extrapolated quantity was obtained using conjugacy within each subpopulation, then sampling from each distribution according to the ratio between them to derive the posterior (the R code is provided in **Error! Reference source not found.**). For the Beta-Binomial model (used solely in the low precision scenario in experiment 1), the posterior distribution for combining the Beta-Binomial data with the Beta prior was obtained by MCMC and approximated by a kernel density estimate.

Three different accuracy metrics were used to compare the elicited distribution to the posterior:

- Difference in the means of the elicited and posterior distributions, which here represents a measure of bias. The absolute value of the difference in means (absolute bias) was also used to capture how much the mean deviates from the true proportion using a metric that is independent of the direction of bias.
- The ratio of standard deviations (SDR) between the elicited distributions and the posterior distribution measure was used as a measure of how well the elicited distribution represented the true level of uncertainty; this was presented either in the natural scale or in the log scale (SDR or lnSDR). Values above 1 on the natural scale (and above 0 on the log scale) indicate underconfidence (overestimation of uncertainty) and values below 1 on the natural scale (and below zero on the log scale) indicate overconfidence (underestimation of uncertainty). Absolute lnSDRs of the difference between true and elicited distributions were also used to capture how accurately the elicited distribution represented uncertainty, independently of the direction of any inaccuracy.
- KL divergence is a measure of the information lost when one distribution is approximated by another¹³. The KL was computed using numerical methods of integration (details in **Error! Reference source not found.**). KL can take any value between 0 and infinity, where zero indicates that the two distributions are identical (i.e. the elicited is without error), and values higher than zero indicate that the elicited distribution is less accurate.

1
2
3 In experiments 3.1 and 3.2, the proportion of participants that revised their priors was observed (see
4 *Chapter 9, Section 9.3.8* for details). The extent of revision was calculated only for those who revised,
5 by comparing the elicited distributions before and after revision. The three different metrics
6 outlined above were also used to determine the extent of revision, but their interpretation differs:
7
8

- 9
10 • Mean of revised distribution minus mean of originally elicited distribution. Positive values
11 indicate the revised mean moved away from the originally elicited distribution towards the
12 group mean, negative values mean the revised mean moved in the opposite direction to the
13 group mean.
- 14
15 • The ratio of the standard deviation of the revised distribution to the standard deviation of
16 the original distribution. This is greater than 1 (or $\ln\text{SDR} > 0$) when the participant became
17 less certain, and vice versa.
- 18
19 • KL divergence of the revised distribution from the original distribution. This combines both
20 bias and uncertainty, and higher values indicate a greater extent of revision.
21
22
23
24
25

26
27 Note that, for brevity, this chapter presents a set of results on select metrics of outcomes, focussing
28 on bias, $\ln\text{SDR}$ and $\ln\text{KL}$. Results for the full set of outcome metrics is presented in **Error! Reference**
29 **source not found.**
30
31

32 33 Section 8.4.2 *Methods*

34
35 The full factorial design means that quantities of interest can be computed and relationships of
36 interest can be displayed from simple plots and summaries of the data. Means (and standard
37 deviations), medians (and interquartile ranges) and histograms were used to describe each metric of
38 accuracy. Comparisons between methods were illustrated using scatter plots for within-participant
39 accuracy. For example, for experiment 1 the accuracy with Chips and Bins was plotted against
40 accuracy for Bisection (x and y axis, respectively), with precision scenarios highlighted using different
41 markers.
42
43
44
45
46
47

48
49 For experiments 3.1 and 3.2, the proportion of participants that revised was compared across
50 different randomised groups and different elicitation methods. The extent of revision was evaluated
51 using empirical summaries and scatter plots.
52
53
54

55
56 Linear and generalised linear modelling was also used to confirm the conclusions drawn in this
57 chapter, to provide estimates with confidence intervals for particular quantities of interest and to
58
59
60

identify any effect of covariates such as period effects. These models and results are fully detailed in **Error! Reference source not found.**

Section 8.5 Results

Section 8.5.1 Description of the sample recruited

In total, 72 participants completed the experiments in 8 sessions (4 to 24 participants per session).¹

Participants earned on average £30.6 (range £20 - £40). The sample characteristics are shown in Table 6. In summary, participants were on average 22.3 years old, 80.5% were female, and 80.6% were undergraduates. Half of the sample was very or somewhat confident in using probabilities, and the average BNT score was 3.7 out of 7. On average, participants scored higher in the 'rational', 'intuitive' and 'dependent' decision making styles, than in the 'avoidant' and 'spontaneous'.

Table 6 Sample characteristics

N		72
Mean age (sd)		22.3 (5.7)
Male (n)		19.4 (14)
UG (n)		80.6 (58)
Year of study (sd)		1.6 (0.8)
Percentage with qualifications in quantitative subjects (n)	A level	38.9 (28)
	AS level	9.7 (7)
Percentage confident in using probabilities (n)	Very confident	4.2 (3)
	Somewhat confident	45.8 (33)
	Neither confident nor unconfident	27.8 (20)
	Somewhat unconfident	16.7 (12)
	Very unconfident	5.6 (4)
BNT score, out of 7 (sd)		3.7 (1.6)
Score on Scott and Bruce's Decision Style Inventory, out of 5 (sd)	Rational	4.1 (0.5)
	Intuitive	3.5 (0.6)
	Dependent	3.7 (0.8)
	Avoidant	2.5 (1)
	Spontaneous	2.5 (0.8)

¹ Three additional participants had their responses invalidated due to failing to comply with directions for completion.

Section 8.5.2 Experiment 1

In Experiment 1, 288 priors were elicited from the 4 games played by each of the 72 participants.

Figure 15 shows the observed distribution of bias, InSDR and InKL scores (results for the full set of outcome metrics are shown in *Error! Reference source not found.*, Figure 1. Bias was symmetrical around the value of zero, suggesting participants were equally likely to over- and underestimate proportions. This was expected in the context of this experiment as participants were observing evidence directly on the target question. In high precision scenarios, participants were more likely to be under-confident (InSDR greater than 0), while the opposite was the case in low precision scenarios. The standard deviation of the InKL scores was high in relation to its mean.

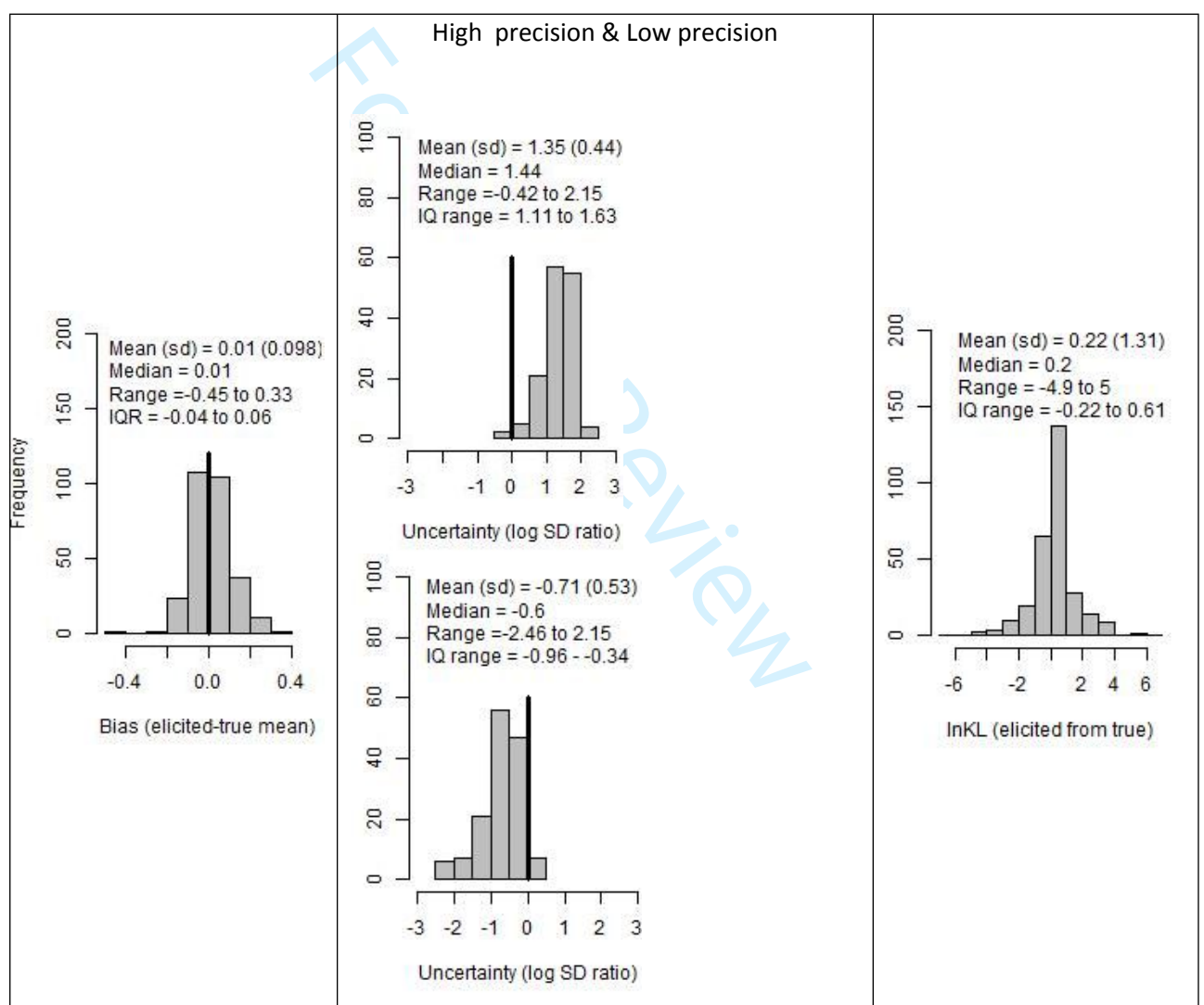


Figure 15 Distribution of bias, InSDR and InKL. (SDR= SD elicited / SD true)

Figure 16 Figure 16 presents scatterplots of the pairwise comparisons of participants' accuracy when different elicitation methods were used (see *Error! Reference source not found.*, Figure 2 and Table

1
2
3 2). Results on bias showed widely scattered points, so that variability in responses dominated any
4 systematic differences between people. Mean and median bias are close to zero and comparable
5 across precision scenarios. There is, however, a higher dispersion of bias values in the low precision
6 scenario (plot a), which means *absolute* bias is higher on average in this scenario (results for
7 absolute bias in **Error! Reference source not found.**, Figure 2). Bias and absolute bias are comparable
8 between the two elicitation methods.
9
10
11
12

13
14
15 For uncertainty (plot b), the results within each high precision scenario suggest that there may be a
16 weak correlation between responses from the same participant, particularly in the high precision
17 scenario. As highlighted in the descriptive histograms in Figure 15, InSDR differed between the high
18 and low precision scenarios. In the scatter plot here, this is evidenced by a clear separation of points.
19 In the low precision scenario, the two elicitation methods result in similar mean InSDR (the blue
20 point is on, or close to, the diagonal line). In the high precision scenario, bisection priors are much
21 more likely to have higher SDRs (higher proportion of points above the diagonal line than below),
22 that is, responses are more likely to be underconfident.
23
24
25
26
27
28
29

30 There is no clear correlation between responses from the same person for KL divergence (plot c).
31 The KL distribution appears more widely dispersed in low precision scenario. In the high precision
32 scenario, the two methods appear to result in similar mean accuracy (the red point is on the
33 diagonal line). The relevance of any observed differences in KL is unclear due to the high variability in
34 scores.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

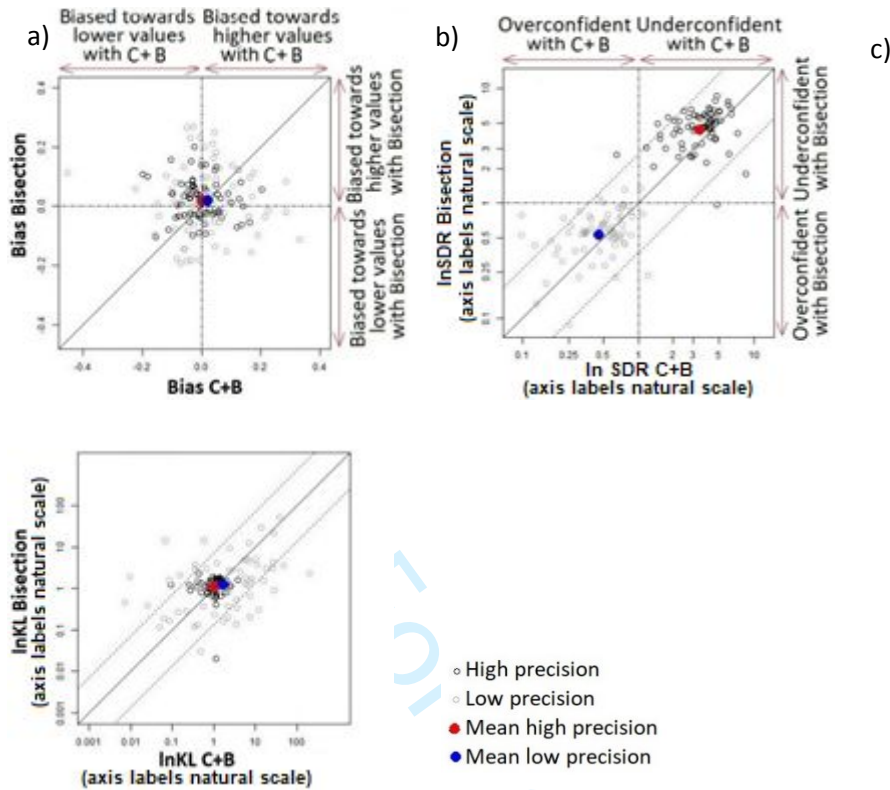


Figure 16 Within-participant comparison of accuracy when different elicitation methods were used, for each precision scenario.

Further detailed results of analyses are presented in *Error! Reference source not found.*. Note that the modelling confirms the empirical results.

Section 8.5.2.1 Participants' preference for each method

The results (Table 7 and Table 8 suggest that participants were more likely to find the Bisection method difficult or challenging, and were also more likely to prefer Chips and Bins to Bisection.

Table 7 Response to question 1 about the ease of completion.

	Easy, % (n)	Challenging, % (n)	Difficult, % (n)
Bisection (n=72)	23.6% (17)	66.7% (48)	9.7% (7)
Chips and Bins (n=72)	43.1% (31)	51.4% (37)	5.6% (4)

Table 8 Response to question 2 about method preference.

	Bisection	Chips and Bins	Indifferent
Preferred, % (n out of 72)	31.9% (23)	65.3% (47)	2.8% (2)

1
2
3
4
5 *Section 8.5.3 Experiment 2*

6 In total, 72 participants played 72 games and 144 priors were elicited (72 for a target question not
7 requiring extrapolation and 72 for a different target question requiring extrapolation). The overall
8 bias, SDR and KL scores in initial priors (compared to the truth) were comparable to those in high
9 precision scenarios in experiment 1 (see **Error! Reference source not found.**).
10
11
12

13
14 *Figure 17* compares participants' accuracy between priors elicited without and with extrapolation,
15 see **Error! Reference source not found.**, *Figure 4* and *Table 5* for results on the complete set of
16 outcome metrics). Results suggest that mean bias is close to zero. There is no suggestion that bias,
17 SDR or KL differ with and without extrapolation, and no suggestion that these results differ between
18 the three different extents of extrapolation. The scatter plots for lnSDR, however, indicate a
19 moderate correlation between lnSDRs of distributions elicited from the same person with and
20 without extrapolation.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

b)

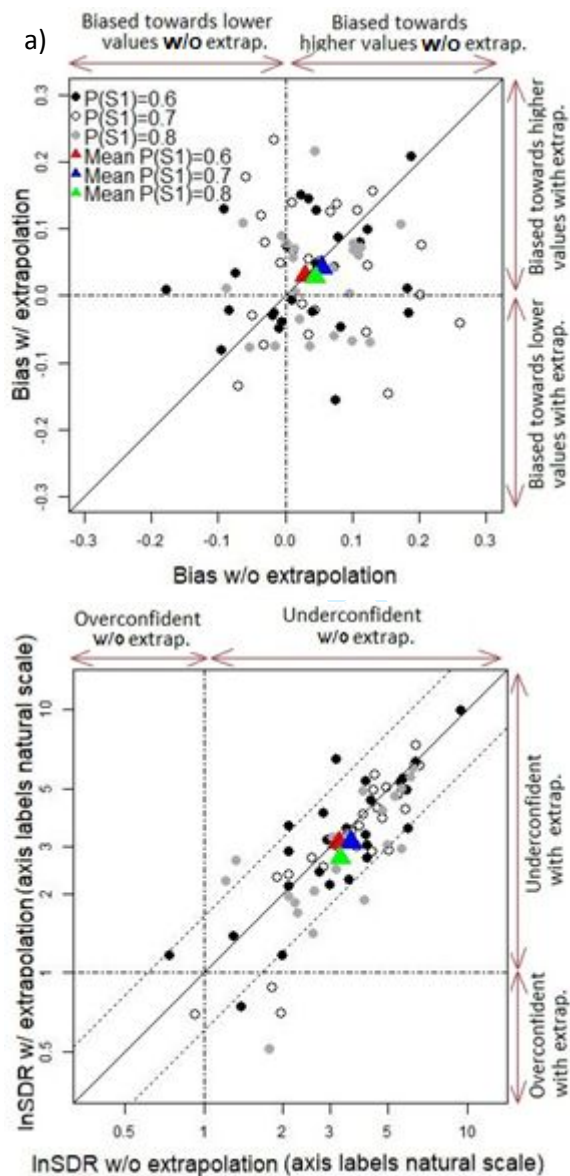


Figure 17 Within-participant comparison of accuracy with and without extrapolation, for different levels of extrapolation.

Section 8.5.4 Experiment 3.1

In total, 72 priors were elicited for the initial quantity and 32/72 (44%) participants updated their priors (i.e. revised) on seeing the group response. The overall bias, SDR and KL scores in initial priors were comparable to those in high precision scenarios in experiments 1 and 2 (see **Error! Reference source not found.**). The high variability of KL between participants makes results on this metric difficult to interpret, and hence these are omitted throughout (but available in Appendices).

Section 8.5.4.1 Likelihood of revision

Table 98 shows the proportion of participants who revised their priors, by type of group summary.

The table shows that participants were more likely to update their priors when the group distribution was discordant from their own prior, and when the group prior was more certain for the same level of discordance. The probability of revision with Chips and Bins (20/35, 57%) appears to be higher than with Bisection (12/37, 32%) (see **Error! Reference source not found.**, Table 12).

Table 9 Proportion of participants who revised their priors.

	Concordant (N=25)	Discordant, equally uncertain (N=24)	Discordant, more certain (N=23)
Proportion who revised their prior (n)	20% (5)	54.2% (13)	60.9% (14)

Detailed results of the regression analysis are presented in **Error! Reference source not found.**

Consistently with the empirical results, the models suggest that participants were significantly more likely to revise their priors when the group was discordant, when the group was more certain, and when using Chips and Bins compared to Bisection. Participants' level of uncertainty on the initial prior (lnSDR for the initial prior) had a significant effect on the likelihood of revision (with participants expressing more uncertainty in their initial prior showing a higher likelihood of revision), while the effects for bias and KL divergence were not significant.

Section 8.5.4.2 Accuracy in initial priors of participants who did and did not revise

Table 10 shows summaries of accuracy (on initial priors) for participants who did and did not revise their priors (selected outcomes presented, for the full set of outcomes see **Error! Reference source not found.**, Table 15). The results suggest that there is no notable difference in the initial level of bias between those who revised and those who did not. Those who revised were perhaps slightly more uncertain in their initial priors (higher lnSDR), although the difference is small.

Table 10 Accuracy of initial priors compared between participants who did and did not revise their priors: mean (SD) [median and (interquartile range)] of accuracy metric over participants

Mean (SD) [median (IQR)]	Concordant (N=25)		Discordant, equally uncertain (N=24)		Discordant, more certain (N=23)	
	Revised (n=5)	Not revised (n=20)	Revised (n=13)	Not revised (n=11)	Revised (n=14)	Not revised (n=9)
Bias	-0.011 (0.088) [-0.039 (-0.068 to 0.044)]	-0.004 (0.068) [0.005 (-0.017 to 0.025)]	0.018 (0.081) [0.031 (-0.036 to 0.043)]	0.013 (0.067) [0.02 (-0.004 to 0.044)]	0.016 (0.068) [0.011 (-0.027 to 0.059)]	0.033 (0.054) [0.032 (0.003 to 0.074)]

						0.72 (0.64) [0.94 (0.57 to 1.21)]
InSDR	1.14 (0.25) [1.04 (0.95 to 1.22)]	1.03 (0.39) [1.10 (0.67 to 1.36)]	1.21 (0.42) [1.23 (1.01 to 1.32)]	0.94 (0.46) [1.05 (0.62 to 1.29)]	0.95 (0.49) [1.097 (0.656 to 1.235)]	

Section 8.5.4.3 Extent of revision in those who revised

Figure 18 shows accuracy in initial priors and extent of revision, using bias and InSDR in initial priors.

A more complete set of outcome metrics are presented in **Error! Reference source not found.**

The average change in mean was highest when the group was discordant and more certain than the participant, and the lowest when the group was concordant with the participant. Participants who saw group summaries concordant with their own prior, on average revised to a more certain prior (lower InSDR). Some participants who saw group priors discordant with their own but with similar uncertainty became more certain and others less certain. Almost all participants who saw group priors discordant but more precise than their own, and chose to revise their prior, revised to express a more certain prior.

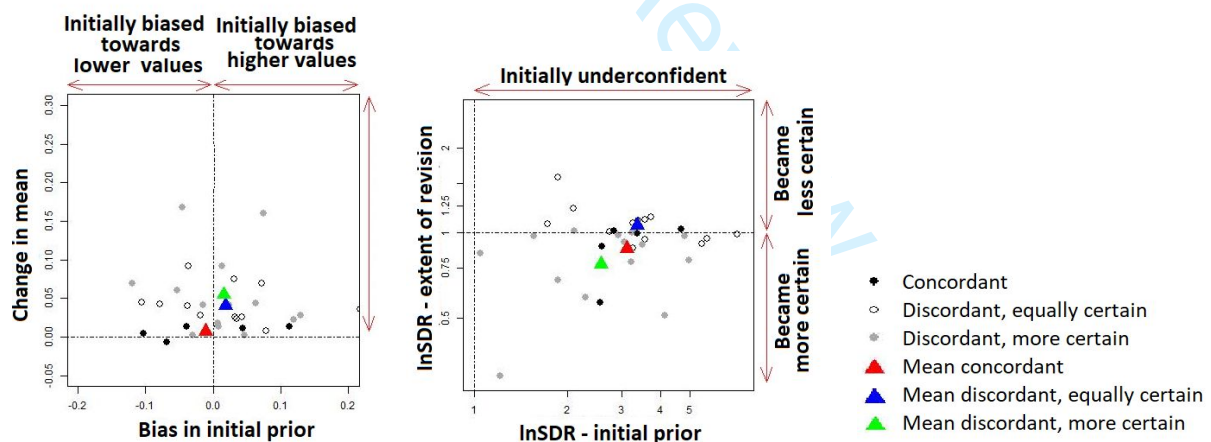


Figure 18 Within-participant comparison of accuracy to initial priors and extent of revision, for different types of group summaries.

Finally, for all outcome measures, there is no evidence that the extent of revision was different for different levels of probabilistic accuracy, as the extent of revision was distributed fairly evenly across the x-axes.

Section 8.5.5 Experiment 3.2

In total, 72 priors were elicited on the initial quantity, and 27/72 (38%) participants updated their priors upon seeing priors from three other individuals in the group. The bias, SDR and KL scores in initial priors were comparable to those in high precision scenarios in experiments 1 and 2 (see **Error! Reference source not found.**).

Section 8.5.5.1 Probability of revision

The proportion of participants who revised their priors, per type of group response and by elicitation method is shown in *Table 1110*. This shows that participants were more likely to update their priors when participants within the group were consistent and that participants may be more likely to revise their priors when using Chips and Bins, which conforms to findings from experiment 3.1.

Table 11 Proportion of participants who revised their priors.

	Consistent (N=34)		Inconsistent (N=38)	
Proportion who revise their prior (n)	52.9% (18)		23.7% (9)	
Proportion who revise their prior (n)	Bisection (n=18)	C+B (n=16)	Bisection (n=20)	C+B (n=18)
	44.4% (8)	62.5% (10)	0% (0)	50% (9)

Detailed results of the logistic regression analysis on likelihood of revision are presented in **Error! Reference source not found.**. Consistently with the data summaries, the models suggest that participants were more likely to revise their priors when participants in the group were consistent with each other, and when using Chips and Bins compared to Bisection. Furthermore, the model coefficients suggested that the effect of probabilistic accuracy on participants’ likelihood of revision was not significant.

Section 8.5.5.2 Accuracy of initial priors in participants who did and did not revise

Table 1211 shows the accuracy of initial priors in participants who did and did not revise their priors (results for all metrics are presented in **Error! Reference source not found.**, Table 15). The results suggest that there is no notable difference in bias or uncertainty between those who revised and those who did not.

Table 12 Results of experiment 3.2: Outcomes in participants who did and did not revise their priors.

Mean (sd) [median (IQR)]	Consistent (N=34)		Inconsistent (N=38)	
	Revised (n=18)	Not revised (n=16)	Revised (n=9)	Not revised (n=29)
Bias	0.036 (0.064) [0.026 (-0.007 to 0.081)]	0.023 (0.07) [0.017 (-0.005 to 0.048)]	0.038 (0.078) [0.066 (-0.04 to 0.102)]	0.046 (0.082) [0.039 (0.006 to 0.104)]
InSDR	1.01 (0.55) [1.17 (0.71 to 1.38)]	0.82 (0.68) [0.92 (0.53 to 1.21)]	0.93 (0.34) [0.88 (0.74 to 1.00)]	0.90 (0.66) [1.11 (0.67 to 1.29)]

Section 8.5.5.3 Extent of revision in those who revised

Figure 19 shows the accuracy of initial prior and the extent of revision for different types of group summary (see complete set of results in *Error! Reference source not found.*, Figure 8 and Table 16). The average change in mean was higher when group members were consistent with each other (but inconsistent with the participant). Absolute change in mean, however, appears to be more comparable between the two groups (see *Error! Reference source not found.*).

When shown group members that were consistent with each other, some participants revised to become more uncertain and others less uncertain. When shown group members that were inconsistent with each other, most participants who revised became more certain (all white points but one are below the 'no change' line, and the mean change in uncertainty (blue triangle) is below the line).

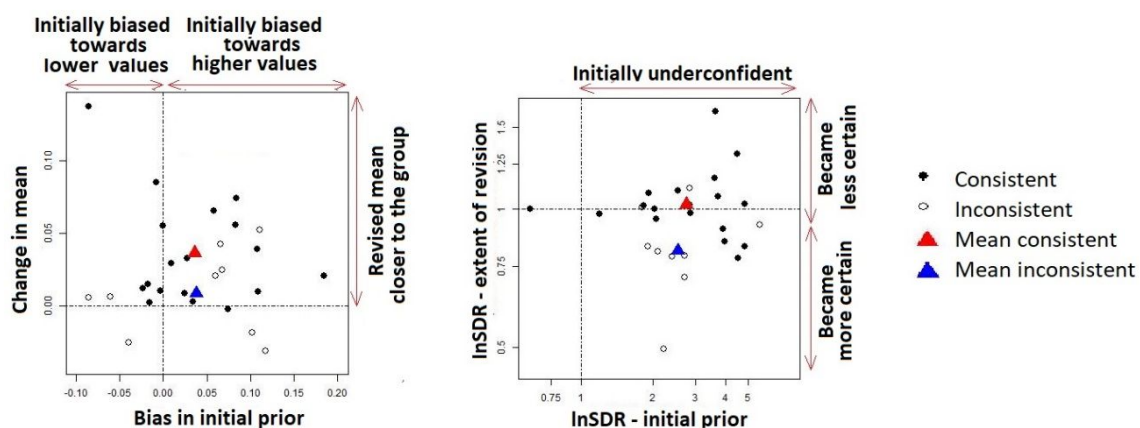


Figure 19 Within-participant comparison of accuracy and extent of revision (using selected outcome metrics), for different types of group summaries.

1
2
3 For a complete set of results see **Error! Reference source not found.**, Figure 8. There was no
4 evidence for any of the outcome metrics that the extent of revision was different for different levels
5 of probabilistic accuracy, as the extent of revision was distributed fairly evenly across the x-axis –
6 this is consistent with findings from experiment 3.1.
7
8
9

10 11 **Section 8.6 Conclusions**

12
13 The experiments described and implemented here use an innovative design to gain insights into how
14 individuals express their knowledge using distributions and on how individuals revise their
15 judgements when presented with other summary information (in a Delphi-type process). The
16 experimental approach allowed the knowledge of each individual participant to be standardised, and
17 this meant a structural (or mechanistic) understanding of how individuals express and revise
18 probabilistic judgements could be acquired.
19
20
21
22
23

24 25 *Section 8.6.1 Key findings from the experiments*

26
27 Experiment 1 imposed a knowledge base directly on the target quantity. Hence, and as expected, it
28 showed no evidence of systematic bias across participants and no evidence of a differential effect of
29 the methods on bias. However, it showed that participants do not adjust sufficiently their
30 expressions of uncertainty: participants gave overconfident priors for the lower precision scenario
31 and under confident priors for the higher precision scenario. Both methods performed similarly in
32 lower precision, but bisection seems to generate more uncertain distributions in higher precision. KL
33 divergence (a measure that combines both bias and uncertainty, and represents how well overall the
34 elicited distribution represents the true distribution) presents high variability and hence throughout
35 all experiment results, appears not to be sensitive to the changes in uncertainty detected in SD
36 ratios.
37
38
39
40
41
42
43
44

45
46 Experiment 2 implemented a knowledge base that was relevant for the target question, but that
47 required some adjustment (or extrapolation). Given that a single observation per participant was
48 obtained, this experiment relied on a smaller number of observations than experiment 1.
49 Experiment 2 generated no evidence that extrapolation, or its level, affects bias, expressions of
50 uncertainty or overall accuracy.
51
52
53
54

55
56 However, it is difficult to give definitive messages from the experiment. It is possible that the
57 experiments lacked power to detect the difference in accuracy when extrapolation is required.
58 Furthermore, the experiment only explores one 'type' of extrapolation, where experts are required
59
60

1
2
3 to derive a weighted average for two probabilities after observing all information required to make
4 the adjustment. In practice, relationships between conditional probabilities can be more complex,
5 and not fully observed; it is not clear whether the findings from the experiments would generalise
6 when more complex extrapolation is required.
7
8
9

10
11 Experiments 3 looked at how and why individuals revise their answers when presented with a
12 Delphi-type summaries. In experiment 3.1, individuals elicited for the target question and were then
13 presented with a group summary that could be discordant with their own belief. Results show that
14 participants were more likely to revise their priors when the group was discordant with their own
15 beliefs, and when the group was more certain than they were. Also, participants were more likely to
16 revise their priors when using chips and bins than bisection. There is no evidence that those who
17 revised have a significantly different accuracy to those that who did not revise. Participants who
18 revised, revised the mean of their priors to a greater extent when the group was discordant and
19 when the group was more certain. When the group prior was concordant with the participant's, the
20 few who revised their priors on average became more certain. When the group prior was discordant
21 with the participant's, and equally uncertain, the participant was equally likely to revise his/her prior
22 in both directions (some became more certain, others less certain). When the group was discordant
23 and more certain than the participant, he/she was more likely to revise to express a more certain
24 prior.
25
26
27
28
29
30
31
32
33
34
35
36

37 In experiment 3.2, instead of a group distribution, individuals were presented with the individual
38 priors from the elements of a group that, overall, was discordant with the individual. Results show
39 that participants were more likely to update their priors when the group members were consistent
40 among themselves (although with wider within-participant uncertainty) than when the group
41 elements are inconsistent among themselves (although more precise). Participants' extent of
42 revision was also affected with participants revising more (towards the group mean) if the elements
43 of the group are concordant (but more uncertain) than if these are discordant (but more precise).
44 When shown group members that were discordant amongst themselves, participants who revised
45 their priors expressed less uncertainty. When shown a set of concordant priors, revisions went in
46 both directions.
47
48
49
50
51
52
53
54

55 Overall, it is apparent that individuals changed their estimates in a rational way when provided with
56 estimates from others (i.e. when everyone else was discordant, individuals were more likely to
57 change their response, if others were uncertain, individuals were less likely to change); however, our
58
59
60

1
2
3 results did not show large differences in the likelihood of revision between individuals with different
4 levels of accuracy -- a key mechanism for revisions to lead to increased accuracy.
5
6
7

8 Our experiments did not explore the effect of group interaction on experts' revision. Clarifying how
9 group interaction affects the likelihood and extent of revision, and the accuracy of a group estimate
10 would have to be carefully studied, controlling for aspects related to the format of feedback and
11 effects of interaction.
12
13
14

15
16
17 *Section 8.6.2 Limitations and suggestions for future research (using the same experimental*
18 *approach)*
19

20 The three experiments implemented aimed to evaluate different aspects of elicitation, and used
21 purposely restrictive set-ups (for example, experiment 1 imposed an equal knowledge base across
22 individuals and the target elicitation referred directly to the observed quantity) to reduce between
23 individual variation (i.e. random error). Whilst such set-ups may be seen as limited, these constitute
24 a starting point whose design can be extended in the future to focus on other aspects important for
25 real life elicitations. For example, the experiments implemented here used a simulated (virtual)
26 learning process in isolation to determine the individual's beliefs over the target quantity of interest,
27 and in this way allow accuracy to be directly assessed. Whilst such learning from observation is a
28 critical source of knowledge for practitioners in health care, our target experts, it is unlikely that, in
29 practice, this is the only source of knowledge used by individuals in formulating beliefs, i.e. health
30 carers may also draw on published evidence, peer contact or other related evidence or experience.
31 Hence, a follow-on from these experiments could introduce other sources of information besides the
32 simulated observations, to not only understand how individuals use multiple sources and whether
33 this affects the accuracy of the distributions provided, but also to determine how to ask individuals
34 to describe their reasoning for the judgements provided.
35
36
37
38
39
40
41
42
43
44
45
46
47

48 Another possible extension to this study relates to examining which individual's characteristics may
49 be associated with accuracy. The sample recruited for these experiments was homogeneous, and
50 hence gives limited insight into how individual's characteristics may be associated with accuracy.
51 Whilst examining participant characteristics was not the focus of these experiments, future research
52 could expand on the pool of experts to identify individuals that may be more accurate, and/or
53 understand how individuals can be trained to become more accurate.¹⁹⁷
54
55
56
57
58
59
60

Section 8.6.3 Experimental approach vs. using almanac quantities

Published methodological research into elicitation has tended to focus on using ‘almanac’ quantities, defined as question that relate to uncertain events that will be realised in the future (e.g. rainfall tomorrow or survival of individual patients), facts (e.g. the distance between the earth and Mars) or to summaries of datasets (population or sample-based).¹⁹⁸ Individuals are typically asked to express their subjective probabilities for these quantities and performance is measured as the deviation between the known value of the quantity and the elicited probabilities. Given the nature of almanac questions, beliefs about them may be formed on the basis of known facts coupled with analytical reasoning (for example, reasoning about the distance between Mars and Earth could be based on the knowledge that it would take a spacecraft one year to reach Mars).

The nature of almanac questions, however, differs inherently from the nature of the quantities typically elicited in health care and other areas where the typical target question relates to unknown parameters of a statistical model, quantities not expected to be known with certainty – for example, the expected probability that UK patients with a wound heal in 6 months when treated with X. Additionally, one of the main sources of information in determining the beliefs of substantive experts (particularly in health care) are observations of patients and their outcomes (e.g. from published studies or own direct observations). This evidence is itself uncertain and may also be biased, heterogeneous and lack generalisability (likely sources of concern in reasoning in HCDM).

Further, whilst individuals are expected to have some level of epistemic uncertainty about their answer to almanac questions, the accuracy of the elicited prior in representing this uncertainty cannot be established directly. Instead, multiple almanac questions are often used and the frequency of true values that fall outside the elicited credible regions is used as a measure of performance. Finally, responses to almanac questions can be inaccurate either if beliefs are themselves inaccurate or if, even with accurate beliefs, the individual struggles to express these in probabilistic terms. While the absolute accuracy of an elicited uncertainty measure cannot be determined using almanac questions, the *relative* accuracy of different elicitation methods, for example, could be compared by randomising individuals to different methods, since randomisation will ensure that any systematic differences between groups will be due to the elicitation methods. However, this may require a prohibitively large sample size.

Chapter 9 Consideration of the methodological choices emerging from existing guidelines

Section 9.1 Introduction

SEE is a process involving many different elements and subsequent methodological choices. A comprehensive list of these elements and choices has been developed and reported in *Chapter 2*. The intention of this chapter was to sift through the available choices with a view to building a reference-case for HCDM (see *Chapter 10*). The specificities of the domain area may help inform some of the methodological choices and so first, the set of principles that underpin the use of elicitation in HCDM is defined (Section 9.2). Each element of SEE then refers to these principles to describe the extent to which the available choices are appropriate for HCDM (Section 9.3). In doing so, it draws on evidence from the Chapters 3-8 of this report.

Section 9.2 Principles underpinning the use of SEE to inform HCDM

A set of principles underpinning the use of SEE to inform HCDM were developed based on the findings from chapters 3-8 (referenced in Table 12). Each of these chapters contributes towards the elements covered in the principles. These principles also reflect considerations for SEE as reported by Cooke (1991)¹⁴, which reflect 'good practice' in SEE more generally and are widely referred to in the SEE literature. The full list of principles was refined to generate nine distinct principles. These are identified in *Table 13* and detailed next.

Table 13 Key principles of SEE in HCDM

Principle	Key message	Source if evidence to support principle
1. Transparency.	SEE should be transparent and reproducible.	<i>Chapter 3</i>
2. Fitness-for-purpose	Elicited information should be fit-for-purpose to be used as an input to further analysis.	<i>Chapter 3,</i> <i>Chapter 4,</i> <i>Chapter 7</i>
3. Consistency, but respecting constraints of the decision making context.	The SEE needs to adapt to the practical and logistic constraints faced by different contexts/decision-making bodies, but maintain a level of consistency in methods used across evaluations.	<i>Chapter 3</i>

1 2 3 4 5 6 7	4. Reflecting uncertainty at the individual expert level	SEE must seek to elicit uncertainty in experts judgements.	<i>Chapter 8</i>
8 9 10 11	5. Recognising and acting on biases	SEE must recognise common expert biases and employ strategies to minimise these.	<i>Chapter 6</i>
12 13 14 15 16 17	6. Suitability for substantive experts who are less likely to be normative	SEE must utilise methods that are appropriate for experts with lower levels of normative skills.	<i>Chapter 3,</i> <i>Chapter 7,</i> <i>Chapter 8</i>
18 19 20 21 22	7. Recognising where adaptive skills are required	When required SEE must employ methods that incorporate or promote the adaptive skills of experts.	<i>Chapter 3,</i> <i>Chapter 5</i>
23 24 25 26 27 28	8. Recognising between-expert variation	SEE must attempt to capture any between-expert variation, understand the reasons why it exists, and explore its potential impact on the decision.	<i>Chapter 5</i>
29 30 31 32	9. Promoting high performance	SEE must motivate experts to best express their beliefs about a quantity of interest.	<i>Chapter 5</i>

Section 9.2.1 Principle 1. Transparency

Chapter 3 surmised that, due to the fact that many SEE's are conducted as part of a wider evaluation, e.g. cost-effectiveness modelling, word count limitations in journal articles often mean that reporting of SEE in HCDM is often insufficiently detailed.⁷² In many instances there is insufficient opportunity to report on the detail of the SEE, particularly the methodological choices made. Systematic and transparent reporting of SEE helps to improve the validity of the resulting expert judgements, allows the SEE to be peer-assessed, and supports others who use the judgements in their own analysis. Where there are word count limitations, a separate appendix should be used to report all details of the SEE, ideally comprised of an elicitation protocol, a summary of the conduct of the exercise, and of its results. More generally in Bayesian analyses, minimum reporting criteria have been published¹⁹⁹. In addition reporting guidelines for SEE for model based cost-effectiveness evaluation have been published⁷², although these do not reflect any emergence of a reference protocol for SEE in HCDM, i.e. they may not reflect all the elements of SEE.

1
2
3 *Section 9.2.2 Principle 2. Fitness for purpose*

4
5 *Chapter 7* showed that cost-effectiveness analysis, alongside other HCDM areas, typically requires
6 judgements on a relatively large number of parameter types, including probabilities, transition
7 probabilities, relative treatment effects, costs and HRQoL scores. This has implications for the
8 quantities used to elicit each parameter as there is the need to i) ensure coherence between the
9 multiple quantities elicited (i.e. to avoid dependency or explicitly elicit this) and ii) ensure these are
10 adequate given the structural constraints imposed by the cost-effectiveness model and target
11 population. Elicited information should therefore be fit-for-purpose to be used as an input to further
12 analysis (e.g. disease modelling, risk assessment model or cost-effectiveness decision modelling).
13
14
15
16
17
18
19

20 *Section 9.2.3 Principle 3. Consistency, but respecting constraints of the decision making context*

21
22 *Chapter 3* discusses the different potential audiences/analysts for SEE, from local level decision
23 makers to national or international decision makers, including reimbursement agencies and research
24 funders. These different decision makers have quite different capacities to conduct SEE and
25 incorporate it into their decision-making processes. For many decision makers, SEE is very likely to
26 be subject to constraints such as timelines, budget and availability of experts. This means that
27 concessions on aspects of design and conduct of SEE are likely to be required. For example, fewer
28 parameters may need to be elicited, a less time consuming method of elicitation may be needed (at
29 the expense of exacerbating bias) or a remote exercise may need to be conducted. It is important
30 that flexibility be retained in a reference protocol for SEE in HCDM, but that the implications of the
31 choices made are explored.
32
33
34
35
36
37
38
39
40

41 *Section 9.2.4 Principle 4. Reflecting uncertainty at the individual expert level*

42
43 As discussed in *Chapters Chapter 1* and *Chapter 4*, and explored in *Chapter 8*, judgements elicited
44 from experts need to reflect the imperfect knowledge they have (referred to as epistemic
45 uncertainty). An important concern is that, when reflecting on their own experiences, experts may
46 instead include some level of variability in their judgements. Variability refers to the fact that
47 individual responses to an intervention will differ between patients with the same observed
48 characteristics within the population. A comparison of methods, chips and bins and bisection, to
49 enable experts to express uncertainty as opposed to variability, is conducted in *Chapter 8*.
50
51
52
53

54 *Section 9.2.5 Principle 5. Recognising and acting on biases*

55
56 As discussed in *Chapter 6*, there are many biases and heuristics (cognitive shortcuts that individuals
57 often use when asked for complex judgements) that apply to SEE, including overconfidence/under
58 confidence, over-extremity, discrimination, or susceptibility to base rate neglect. There are
59
60

1
2
3 techniques available to reduce associated biases, which may help mitigate against their effect (see
4 *Chapter 6*), however these have not been applied in the context of HCDM. Efforts should be made to
5 integrate the findings and recommendations from behavioural research on what biases and
6 heuristics can play an important role in SEE. SEE should be designed and conducted in a way that
7 minimises the use of heuristics and other sources of bias and appropriate training should be given to
8 experts.
9

10
11
12
13
14
15 *Section 9.2.6 Principle 6. Suitability for experts who possess substantive skills, who are less likely to*
16 *be normative*

17
18 *Chapter 5*, expert selection, concludes that substantive experts in HCDM are often health
19 professionals, who are unlikely to have had extensive experience of quantifying their knowledge of
20 healthcare outcomes, which may compromise their normative skills. They are, however, often
21 subject experts and are recruited to take part in a SEE based on their substantive expertise. This is
22 not typical of some of the other areas of science in which elicitation is commonly used; hence,
23 methods of SEE employed in other domains may not be directly suitable in HCDM or additional
24 training may need to be delivered before their use. For example, in the choice of method of
25 elicitation e.g. graphical methods, such as the chips and bins method, have been claimed as more
26 intuitive than the bisection method. Additionally, particularly in this context, it may be preferable to
27 elicit only quantities that may be observable, and to recognise concerns over the elicitation of
28 dependency.
29
30
31
32
33
34
35
36
37
38

39 *Section 9.2.7 Principle 7. Recognising where adaptive skills are required*

40 *Chapter 5*, identifies very little evidence to clarify the role of adaptive skills, however, given the
41 multiple purposes for SEE (*Chapter 1*), it is proposed that adaptive skills may be relevant in SEE for
42 HCDM. In particular, it may be necessary to use SEE to inform HCDM in early cost-effectiveness
43 modelling or early stage trial design. In this situation, experts may not be familiar with the target
44 quantity/population for elicitation but are substantive experts in one or more related quantities. In
45 this case, the SEE relies on the adaptive skills of experts and it is important that expert selection
46 and/or training activities accommodate this.
47
48
49
50
51
52
53

54 *Section 9.2.8 Principle 8. Recognising and act on between-expert variation*

55 *Chapter 5* discusses the issue of between-expert variation and the different methods for SEE, which
56 deal with this variation (level of elicitation). In the context of HCDM, this variation is common;
57 however, its causes are poorly understood. In the context of HCDM, there may be genuine
58
59
60

1
2
3 heterogeneity in the populations experts draw upon to formulate their judgements and this may
4 contribute to between-expert variation. In this case, it is desirable to reflect this variation in the
5 pooled distribution, whether through group consensus or mathematical aggregation methods. There
6 should also be efforts made to understand why between-expert variation is present, for example if
7 this reflects heterogeneity in clinical observations such as patient severity. In some circumstances it
8 may not be appropriate to combine judgements from experts where there is heterogeneity.
9

15 *Section 9.2.9 Principle 9. Promoting high performance*

16 *Chapter 5* discusses the need to recruit experts that are motivated to undertake the SEE task
17 optimally and that they have some kind of altruistic reason for providing their honest beliefs, i.e. to
18 improve population health. In HCDM, experts may be motivated to undertake the task to the best of
19 their abilities as a result of their interest in the topic area and improving population health through
20 better HCDM. It may be the case, however, that not all experts within a SEE will perform as well. As
21 well as promoting high performance a SEE may want to explore any differences in expert
22 performance that emerge.
23
24
25
26
27
28
29

30 **Section 9.3 How do SEE elements and methodological choices reflect the principles underpinning** 31 **healthcare?**

32 This section considers the choices available for the different elements of SEE identified from the
33 guidelines review (see *Chapter 2*). Not all principles for SEE in HCDM are relevant for all elements.
34 The most relevant principles for each element and components within these are considered in the
35 sections below.
36
37
38
39

40
41
42 These principles are applied to each of the choices, within components and elements, in the order in
43 which they are presented in *Chapter 2*, with managing biases and validity overarching considerations
44 throughout the process. Each section provides a summary of what choices the principles support.
45 This is summarised presented for all components in the table in ***Error! Reference source not found.***
46
47
48
49

50 *Section 9.3.1 Selecting quantities (preparation and design)*

51 Different quantities can be elicited that provide information on any single parameter of interest.
52 There are a number of issues relevant when determining the choice of quantity to elicit. The choices
53 for quantity are presented in *Chapter 2* and then considered in further detail in *Chapter 7*.
54
55
56
57
58
59
60

1
2
3 Key in HCDM is that the elicited information should be fit-for-purpose (principle 2) and describe
4 experts uncertainty regarding the quantity of interest (principle 4). There is a lack of empirical
5 evidence on whether to elicit directly observable or non-observable parameters in HCDM. In theory,
6 so long as the elicited distributions can be used to provide information on the parameter of interest,
7 either may be appropriate. However, experts in HCDM are often required for their subject expertise
8 and they may be less likely to possess high levels of normative skills (principle 6). For this reason it
9 may be advantageous to elicit less complex quantities which require high levels of normative skills,
10 for example relative risks. It may also be relevant to consider how other empirical evidence are
11 reported in the literature, i.e. how it is expressed statistically, particularly if synthesis with elicited
12 quantities is required (principle 2). The applied literature tends to suggest observables are preferred
13 in this context (see *Chapter 4*) and the existing guidelines consistently support this choice (see *Chapter*
14 *2*).

15
16 Adding another layer of complexity is the issue of dependent quantities. Where dependence exists
17 between multiple elicited quantities, and experts can express it, it is appropriate to use dependence
18 elicitation methods. Dependency can be elicited by expressing dependent variables in terms of
19 independent variables or by eliciting conditional probabilities, and they have been used in this way
20 in previous applications (see *Chapter 4*). More complex dependence elicitation methods such as
21 regression-based techniques and other specialized techniques have not been applied in HCDM to
22 date and it is unclear if these would be appropriate for HCDM experts (principle 6).

23
24 The choice of quantities of interest may also be guided by the practical constraints of the context.
25 In HCDM, there is often a need to generate quantities relatively quickly to inform decision-making
26 (principle 3), perhaps reducing time available for training, particular on a face-to-face basis. In these
27 circumstances, it may be advantageous to elicit dependent variables in terms of independent
28 variables. In addition, when describing quantities, efforts to reduce cognitive burden on the experts,
29 such as avoiding vagueness and asking questions in a manner consistent with how experts express
30 their knowledge, may be preferable.

31
32 The principles support the following:

- 33 • Criteria to determine the choice of parameters, including minimal assessment of each
34 possible uncertain parameter (sensitivity analysis) to identify which have the biggest impact
35 (principle 3).
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- Types of quantities: observable quantities such as probabilities (expressed as proportions or frequencies, but not more complex quantities such as higher moments of a distribution, odds ratios, or credible ranges (principle 2, 3).
- Dependency: ask only about independent variables, express dependent variables in terms of independent variables or use dependence elicitation methods (principle 6)
- Wording: avoid vagueness; ask questions in a manner consistent with how experts express their knowledge; use neutral wording, avoiding leading questions, decompose into simpler quantities where possible (principle 3).

Section 9.3.2 Methods to encode judgements (preparation and design)

To inform HCDM, SEE should reflect the complexity of the further analysis it is informing, and of the other evidence supporting it (principle 2), for example the requirement to elicit for multiple parameters where there may be evidence of dependencies between them (see 9.3.1). In order to be practical in different contexts/decision-making bodies (principle 3), the methods used to elicit beliefs need to be easy to implement and not require extensive training. The suitability of the alternative choices must also recognise differences in normative skills of experts (Principle 6).

Existing guidelines suggest both the FIM and VIM to encode judgements. To date there has been a lack of empirical evidence on which method works better in this context, whilst providing accurate representations of experts beliefs, in particular of their uncertainty (principle 4). Both methods have been applied in HCDM (see *Chapter 4*). The experiments presented in *Chapter 8* sought to explore the use of these two methods in HCDM and compare them in terms of procedural performance, where there is no heterogeneity in knowledge. Little difference between VIM and FIM was found, particularly under conditions of low precision, the situation most likely in HCDM. There was a preference for the FIM by experts. This may be because the VIM requires experts to express their uncertainty using quantiles, which may be summaries of a distribution less familiar to experts. For both methods the expert should be trained to understand how to express uncertainty (principle 4).

The principles support the following choices:

- All forms of FIM or VIM: a decision maker can choose either but apply these consistently in their setting (principle 4, 6).
- Training: this should be provided to experts and focus on how to express uncertainty (principle 4).

Section 9.3.3 Selecting experts

As part of the documentation (Section 9.3.15), the process for selecting and recruiting experts should be reported (principle 1); including details of the numbers of experts approached and the number that declined to take part (see Section 9.3.15, Chapter 2, and Chapter 6).

The existing guidelines suggest that features including normative expertise, substantive expertise and willingness to participate, can be used as defining characteristics to select experts. However, the constraints of conducting SEE in HCDM may dictate that the selection and recruitment of experts focus on only one or two key characteristics (principle 3). In particular, it is worth noting that health care professionals with the relevant substantive expertise may be limited in number, and therefore more opportunistic methods for recruitment may be required, such as peer nomination (see Chapter 5 for examples). In some instances, adaptive skills may be required for a SEE, particularly in the case of new and emerging technologies (principle 7). The challenge in attempting to recruit experts that possess high levels of adaptive skill is that this characteristic is not well defined in the literature (see Chapter 5).

Defining an 'unbiased' expert poses a challenge and indeed, it may be impossible to do so (principle 5). Chapter 6 suggests that the SEE can seek to recruit experts that are free from motivational biases by collecting disclosure of personal and financial interests and conflicts of interest. This may be a challenge as those with the greatest knowledge about a particular treatment or technology, and greatest willingness to participate, may be those with the greatest interest in the SEE. An alternative strategy therefore is to ensure that a range of viewpoints are represented in the sample, with the intention of "balancing out" or at least diluting the effect of motivational biases (see Chapter 6).

Between-expert variation may exist and the methods used to select experts must attempt to capture the range of plausible beliefs (principle 8). Identification of experts through recommendations by peers, either formally or informally, may generate a pool of experts that are all similar. Instead it may be preferable to recruit experts through research outputs, known experience or profile matrix. The SEE can also seek diversity in background, a balance of different viewpoints and a balance of internal and external experts. A larger number of experts may help to ensure that the selection of experts available fulfils these criteria (see Chapter 2 suggests at least 5 experts).

The principles support the following choices:

- Selecting experts on the basis of their substantive expertise and willingness to participate (principle 3).
- Recruitment: recruit experts that are free from motivational biases where possible. In all instances collection information on personal, financial and conflicts of interest (principles 1, 5).
- Method to recruit: a range of methods are available to recruit experts. Which even method is used, this should strive for diversity in the pool of experts (principle 8).
- Number of experts: include at least 5 experts (principle 8).

Section 9.3.4 Pilot exercise

All existing guidelines agree that a SEE should include a pilot of the exercise and omitting a pilot could actually cost time rather than saving it (principle 3). The pilot can be used to explore which method best reflects uncertainty at the individual level. If the training is also piloted, the analyst/facilitator can also use this opportunity to gain feedback from experts on how capable they felt using methods to express their uncertainty (e.g. VIM or FIM) and make revisions to the SEE if required (principle 4).

The pilot can also be used to determine the appropriateness of the SEE for those experts recruited, particularly if the sample of experts have low levels of normative skills (principle 6). This can involve piloting of alternative ways of formulating the questions, which quantities are used (see Section 9.3.1) or the method to encode judgements (see Section 9.3.2).

The principles support the following choices:

- Piloting: this should be undertaken prior to the task. Use of feedback to revise the SEE (principles 3, 4, 6)

Section 9.3.5 Training and preparation for experts

A proportion of the SEE should be spent on delivering training, as it is unlikely that HCDM experts will have had any previous experience of SEE. Training and preparation should focus on enabling non-normative, but substantive experts, to express their beliefs appropriately (principle 6). This should focus on giving them the tools and information to express their uncertainty at the individual level (principle 4). Non-normative experts may be wary of the SEE task and this may have implications for how confident they are at expressing their beliefs. Some experts may express over confident distributions for fear of being judged (see *Chapter 6*). Training therefore plays a role in minimising biases (principle 5), and, although the evidence, in the context of HCDM, is weak, there

1
2
3 are some suggestions from the literature that training can be efficacious in reducing the effect of
4 anchoring, adjustment in interval, confirmation bias and overconfidence (see *Chapter 6*).

5
6
7
8 The existing guidelines (see *Chapter 2*) do not provide a definitive list of what should be covered in
9 training and the elements included will be driven, in part, by the specific application, for example
10 description of quantities, description of performance measurement and dependence. Some
11 elements, such as how results will be used, motivation of elicitation and the full protocol, may not be
12 possible to include due to the time constraints likely in HCDM (principle 3). In addition, a list of
13 relevant information is only typically used as part of a group process or where there are efforts to
14 standardise the level of substantive skills across experts. The (core elements): description of what is
15 required from experts, outline of process, outline of questions, example and practice questions and
16 assumptions and definitions used in the elicitation, should not be compromised (see Section 9.3.9).

17
18
19
20
21
22
23
24
25 The principles support the following choices:

- 26 • Training: this should be delivered and should focus on: 1) enabling experts to experts their
27 uncertain belief, 2) minimising bias (principles 3, 4, 5, 6).

28 29 30 31 32 *Section 9.3.6 Level of elicitation (elicitation)*

33 Judgements from multiple experts are preferred in a SEE (see Section 9.3.3). Existing guidelines are
34 inconsistent w.r.t the level of elicitation – individual or group based. Group discussion can aid less
35 substantive and normative experts, however, face-to-face discussion can be resource and time
36 intensive (principle 3). Access to trained facilitators for group level elicitation may be scarce within
37 HCDM (see *Chapter 3*).

38 Interaction between experts can also introduce biases (see Section 9.3.16) (principle 5). The act of
39 striving for consensus can potentially eliminate some of the between expert variation; the potential
40 for ‘groupthink’ (see *Chapter 6*). A group process should aim to reflect both individual level
41 uncertainty and between expert variability in the aggregated distributions (principles 4 and 8) and
42 there may be greater potential to explore variation in experts beliefs with a group based approach,
43 but only if face-to-face. In HCDM, there may be a lack of experienced facilitation, thus it may not be
44 possible to do this. In these circumstances, an individual level elicitation may be more appropriate. A
45 large sample, which may be required to ensure representativeness, may be a challenge for a group-
46 based exercise, particularly if face-to-face. An individual level elicitation can ask experts to express
47 how they formulate their beliefs, however it is a challenge to them incorporate these differences
48 into the resulting aggregate distribution.

1
2
3
4
5 Group elicitation via remote means may be practical in some circumstances. As discussed in Section
6 9.3.11, interaction between experts can be beneficial for non-normative experts (principles 6 and 9).
7
8 Remote group elicitation can help to mitigate against dominate experts. Individual level elicitation,
9
10 whilst avoiding this situation, can be daunting for experts that have not undertaken such tasks
11
12 previously (non-normative).
13
14

15 The principles support the following choices:

- 16 • Level of elicitation: elicit from experts individually (principles 3, 4, 8).
- 17 • Role of consensus: where required, should first conduct individual elicitation followed by
18 group consensus (principles 6, 9).
19
20
21
22

23 *Section 9.3.7 Mode of Administration (elicitation)*

24 A number of alternative modes of administration have been used in HCDM (see *Chapter 4*), however
25 many of the existing guidelines agree that face-to-face administration is preferred (see *Chapter 2*). It
26 is thought to promote good performance (principle 9) and maximise engagement with experts. Face-
27 to-face elicitation is required for some consensus methods (*Chapter 5*), however it is not necessary
28 for a mathematical approach.
29
30
31
32
33

34
35 The constraints in HCDM (principle 3) are the biggest factor in driving the method chosen. If a large
36 number of experts is sought (see Section 9.3.3), in order to generate timely results, face to face
37 elicitation may be prohibitively time and resource expensive. The constraints of HCDM do not imply
38 that a particular vehicle is used, i.e. paper based on computer based questionnaire, however in order
39 to record information elicited effectively, the majority of applications in the context have used a
40 computer based exercise, either developed for that unique purpose or using existing 'off the shelf'
41 software (see *Chapter 4*).
42
43
44
45
46
47
48

49 The principles support the following choices:

- 50 • Administration: can conduct SEE using face-to-face or remote administration (principles 3, 9).
51
52
53

54 *Section 9.3.8 Feedback to experts and revision (elicitation)*

55 Feedback and opportunity for revision can be used as a strategy to minimise bias (principle 5 and
56 *Chapter 6*). The guidelines are consistent in recommending that feedback and opportunity for
57 revision takes place but differ w.r.t what to feedback. The process should be made explicit and
58
59
60

1
2
3 documented appropriately; including the number of feedback rounds and what is fed back (principle
4 1).

5
6
7
8 For non-normative experts (principle 6) graphical feedback could be useful, whereas more complex
9 summarises such as fitted distributions, performance scores or results using elicited values may not
10 be appropriate. Experts may find distributions from other experts, summaries of aggregated
11 distributions, rationales, future data, the draft elicitation report or qualitative discussion of elicited
12 values, useful, however it is not clear how these could improve the SEE unless they are accompanied
13 by the opportunity for revision. If the feedback allows experts the opportunity to revise their
14 distributions, it may be a useful process as this can help to promote high performance and
15 distinguish between high and low-performing experts (principle 9).

16
17
18 Feeding back distributions of other experts is common in group-based approaches (see *Chapter 5*)
19 and it may incentivise less high performing experts to revise their distributions; however this may be
20 driven by how uncertain they are about their beliefs. There is also the possibility that high
21 performing experts will also revise their distributions, potentially generating less accurate pooled or
22 group summaries.

23
24
25 The principles support the following choices:

- 26 • Feedback: this should be offered to experts with the possibility of revision. What to feedback
27 will depend on the SEE task and the types of experts include. Graphical feedback may be
28 useful for non-normative experts (principles 1, 6, 9).

29 30 31 32 33 34 *Section 9.3.9 Opportunity for interaction (elicitation)*

35 Interaction is intrinsically linked with the level of elicitation and as such, many of the principles
36 relevant for Section 9.3.6 are relevant for how the interaction process works. Interaction can allow
37 experts to share information so that differences in expert opinion are not the result of experts
38 having different information or interpreting questions differently (principle 9). In addition, it is
39 important to note that, remote and controlled interaction, such as that promoted with Delphi type
40 processes, can avoid some of the biases of group exercises (principle 5) and can be preferable from a
41 practical point of view (principle 3). However, remote elicitation can encourage experts not to take
42 responsibility for their expressed beliefs (self-serving bias). As with group and individual methods,
43 there is also a lack of evidence on how the revision process can affect the accuracy of the final
44 individual distribution (principle 9). For consensus SEE, a group based face to face session may help
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 to promote the beliefs of experts with better performance and reflect between expert variation
4 (principles 8 and 9) (*Chapter 6*).

5
6
7
8 The principles support the following choices:

- 9
10 • Interaction: this should follow on from an individual elicitation where practically feasible and
11 useful (principle 9).

12
13
14
15 *Section 9.3.10 Feedback from experts on process (elicitation)*

16 In addition to feedback from the facilitator and/or other experts (Section 9.3.8), a SEE can encourage
17 feedback from experts on the process, either qualitatively through an interview or questionnaire, or
18 through some kind of quantitative ranking. This is linked to obtaining rationales (Section 9.3.11),
19 however more broadly relates to the elicitation process rather than the beliefs about quantities.
20 Only a limited number of guidelines discuss obtaining this type of feedback (see *Chapter 2*) and a
21 limited number of applied studies have attempted to collect this information in HCDM (see *Chapter*
22 *4*). Given the lack of practical experience and empirical evidence, it is difficult to be prescriptive
23 about how this type of feedback might work in the context of HCDM and it may be driven by the
24 time and resource constraints of the task (principle 3). It may be valuable to ascertain what
25 elements of the SEE experts found challenging. This information could then be used in designing
26 future exercises (principle 1). In addition, the information gleaned from experts during the feedback
27 could be used to discriminate between low and high performers (principle 9), however this would be
28 on the basis of their subjective assessment, rather than on the basis of a quantitative measures of
29 accuracy such as calibration (see Section 9.3.14 and Section 9.3.17). Overall it is not clear how this
30 form of feedback would be beneficial to the SEE task or improve the resulting distributions.

31
32
33
34
35
36
37
38
39
40
41
42
43 The principles support the following choices:

- 44
45 • Asking experts for feedback: should only ask the experts to appraise the SEE process if there
46 is a clear reason for doing so (principle 1).

47
48
49
50 *Section 9.3.11 Rationales (elicitation)*

51 Almost all guidelines recommended collecting qualitative data from experts on how they formulated
52 their judgements. Experts in HCDM may possess different levels of normative and substantive skills
53 and the setting in which they work may expose them to different clinical experience, which can drive
54 their beliefs. It is therefore important that the rationales for the beliefs given are collected (principle
55
56
57
58
59
60

1
2
3 1). This information can then be considered. This can inform an assessment of validity of the elicited
4 beliefs.
5
6
7

8 The methods used to collect rationales will be driven by the mode of administration (seeSection
9 9.3.7). Interaction between the analyst and the expert on a one to one basis can encourage experts
10 to explain in greater detail their rationales. Where SEE is conducted remotely it may be
11 advantageous to use prompts such as multiple choice questions to encourage experts to reveal any
12 detail about their rationales.
13
14
15
16
17

18 The principles support the following choices:
19

- 20 • Rationales: these should be collected and recorded from experts about how they made their
21 judgments (principle 1).
22
23
24

25 *Section 9.3.12 If/how to aggregate (aggregation, analysis and post-elicitation)*
26

27 In order to generate distributions that are fit for purpose (principle 2), aggregation is preferred over
28 no aggregation. With respect to the choice of aggregation method, *Chapter 5* concludes that, on the
29 basis of the evidence available, in terms of 'accuracy' including representation of uncertainty,
30 mathematical and behavioural aggregation perform similarly. There is also no evidence to support
31 the specific type of behavioural aggregation method used. For mathematical aggregation, simple
32 mathematical decision rules, like a linear opinion pool with equal weights are the most commonly
33 applied in HCDM (see *Chapter 4*) and are straightforward to implement (principle 3). Mathematical
34 approaches allow experts to express their uncertainty and then, if appropriate aggregation
35 approaches are used, this feeds through into the overall distribution achieved (principle 4).
36
37
38
39
40
41
42
43

44 Mathematical aggregation does not require experts to converge to a group distribution, therefore
45 allowing variability between experts to be reflected within an overall distribution, either using
46 opinion pooling or Bayesian methods (principle 8). A mathematical process can elicit the reasons for
47 the distributions expressed, however it cannot use these quantitatively in generating a single overall
48 distribution, unless the reasons for these distributions are reflected in the seeds that are generated
49 as part of a calibration process (principle 8, seeSection 9.3.17). Calibration-based performance
50 weighting, which has received little attention in the HCDM literature (see *Chapter 5*) 'solves' any
51 between-expert variation in performance by differentially weighting experts according to their
52 performance on 'seed' scores (see *Chapter 5*). Generating differential weights in HCDM, is, however,
53
54
55
56
57
58
59
60

1
2
3 problematic as discussed in *Chapter 5*. Further research on weighting methods within HCDM is
4 needed to advise if and when choices beyond equal weighting are warranted.
5
6
7

8 The principles support the following choices:
9

- 10 • Aggregation: Mathematical aggregation or individual elicitation followed by behavioural
11 aggregation (principle 4, 8).
- 12 • Method of aggregation: Use of linear pooling methods (principle 8), including equal
13 weighting of experts distributions (principle 3).
- 14
15
16
17

18 *Section 9.3.13 Fit to distribution (aggregation, analysis and post-elicitation)*
19

20 As part of the aggregation procedure, post elicitation, statistical distributions need to be fitted to
21 elicited data (see *Chapter 5*). The choice of parametric distribution is uncertain. There is a lack of
22 evidence in HCDM on the fitting process in SEE. Limited evidence suggests that standard distributions,
23 such as the Beta will often be sufficient. More complex approaches may be appropriate, however
24 these can be complex to implement in general software (principles 2 and 3).
25
26
27
28

29
30 The fitting process should ensure that uncertainty, at the individual expert level, is reflected
31 (principle 4) and to do this the distribution used should capture the experts distributions as closely
32 as possible. It is also important that the aggregation respects between-expert variation (principle 8).
33 It is difficult to be prescriptive about which distribution is most suitable, as this will be driven by the
34 quantity elicited and how experts have expressed their beliefs, i.e. the shape of the distribution,
35 however the resulting distributions must generate quantities that can be used within further analysis
36 (principle 2), for example, without transformation.
37
38
39
40
41
42
43

44 The principles support the following choices:
45

- 46 • Fitting: distributions should be fitted to experts elicited beliefs (principle 2).
- 47 • Which distribution: this will depend on the quantity and how the beliefs are represented,
48 however distributional forms such as normal, beta and other conjugate family will often be
49 appropriate (principle 2, 3, 4).
- 50 • Fitting criteria: the use of minimum least squares, method of moments or other approaches,
51 to select the appropriate distribution (principle 2, 3, 4).
- 52
53
54
55
56
57
58
59
60

1
2
3 *Section 9.3.14 Adjusting judgements (aggregation, analysis and post-elicitation)*
4

5 Experts may possess differential levels of normative, substantive and adaptive skills, which may
6 result in differential performance. None of the existing guidelines discuss methods to adjust for
7 'performance' post elicitation (see *Chapter 2*) even if they do refer to a validation process (see
8 Section 9.3.17).
9
10

11
12
13 The methods used for SEE should motivate experts to express their true beliefs about a quantity of
14 interest and quantify differential performance between experts (principle 9), implying that adjusting
15 judgements is preferred to not adjusting, if it generates more accurate pooled distributions.
16

17 Without objective measures to quantify performance, however, adjustment may instead resolve
18 variability between experts, which is not desirable in HCDM, where variation may exist for valid
19 reasons (principle 8).
20
21
22

23
24
25 The principles support the following choices:

- 26 • Adjustment: this should not focus on simply reducing variability between experts (principle
27 8).
28
29
30
31

32 *Section 9.3.15 Documentation (aggregation, analysis and post-elicitation)*
33

34 In order to inform an explicit decision making process in healthcare a SEE must report on all
35 elements of the process and justify the choices made in determining these choices (principle 1).
36

37 There is no agreed list of what should be presented, emerging from the existing guidelines (see
38 *Chapter 2*). However, recently, guidance, not described as a guideline, reported on what information
39 should be in HCDM (see *Chapter 4*). Iglesias, et al⁷² suggest 16 criteria for a SEE and 11 criteria for a
40 Delphi study specifically. These largely accord with the items identified from the existing guidelines,
41 although for Delphi surveys they also suggest a description of the literature review and the number
42 of rounds performed. They do not specifically advocate reporting on details of how uncertain
43 quantities are measured (VIM, FIM) or any training methods used.
44
45
46
47
48
49

50 It is important to note that many applications of SEE are conducted alongside cost-effectiveness
51 modelling or some other form of evaluation, and therefore the amount of material that could be
52 reported may be vast. *Chapter 4* showed that, as a consequence, many details of a SEE are often
53 omitted from a published manuscript or report. It may therefore be advantageous to specify a
54 minimal set of documentation, such as that suggested by Iglesias, et al⁷².
55
56
57
58
59
60

1
2
3 The principles support the following choices:
4

- 5 • Documentation: this should be thorough and covers all aspects of the SEE design, conduct
6 and analysis (principle 1).
7

8
9
10 *Section 9.3.16 Managing biases*

11 In striving to minimise bias, efforts should be made to identify which biases are likely for the sample
12 of experts included (principle 1), and relevant strategies to minimise these (bias reduction
13 techniques) should be employed (principle 5). *Chapter 2* does not suggest specific techniques for
14 addressing individual types of biases or heuristic, and instead gives multiple suggestions across the
15 range of biases.
16
17
18
19

20
21
22 *Chapter 6* suggests that it is difficult to recommend particular bias-reduction techniques over others,
23 as what works best will depend on the context and what biases are most apparent. Given the
24 recommendations for training made in Section 9.3.5, it would seem appropriate to extend this
25 training to cover issues of bias, but going beyond simple warnings. Allowing experts to practice
26 expressing their beliefs using either VIM or FIM, followed by feedback, may also reduce the
27 probability for some of the biases.
28
29
30
31

32
33 The principles support the following choices:
34

- 35 • Anticipate likely biases: for the sample of experts included and specific task. Discussion with
36 experts can help to identify potential biases (principle 1, 5).
37
- 38 • Frame questions to minimize bias and ambiguity. This can include asking experts to first specify
39 the credible interval (upper and lower bounds) and provision of relevant background evidence
40 (principle 1, 5).
41
- 42 • In selecting experts: minimize and record conflicts of interest among the experts. Include
43 experts external to the SEE task, i.e. not those involved in developing the task (principle 1, 5).
44
- 45 • Focus training: on biases and expressing uncertainty and give experts practice and feedback
46 using either the FIM or VIM (principle 1, 5).
47
- 48 • During the task: experts should address conflicting information and provide their rationales
49 (principle 1, 5).
50
51
52
53
54
55

56 *Section 9.3.17 Validation*

57 The guidelines differ in their definitions of validity and discussion of how the concept can be
58 operationalised in an elicitation. Commonly discussed elements of validity include that the elicitation
59
60

1
2
3 captures what experts truly believe or that the expressed probabilities reflect reality. Certain
4 elements of validation accord with the section relating adjusting judgements (see Section 9.3.14),
5 however a number of existing guidelines describe validation of the process rather than the results of
6 SEE. The method used for validation should strive to explore the implications of between-expert
7 variation and attempt to understand why it is present (principle 8).
8
9

10
11
12
13 Understanding how experts formulate their beliefs and why experts present heterogeneous beliefs,
14 can potentially improve the validity of the SEE (principles 1 and 8). The following choices could fulfil
15 this purpose: provision of feedback, testing that the question is understood, fitness for purpose,
16 assessing the accuracy of judgements (see *Chapter 5*), coherence testing, rationales, checks for
17 inconsistencies and internal and external peer review. Faithfully capturing experts beliefs should
18 always be the aim of SEE; however where there is no data to explicitly validate this, there is no way
19 of checking if the resulting distributions represent experts beliefs.
20
21
22
23
24
25

26
27 Fitness for purpose (principle 2) states that the validation process should generate distributions that
28 can be used in HCDM. To this end one of the possible validation processes described by the review
29 of guidelines is fitness for purpose, which evaluates if the elicitation process provides an appropriate
30 level of precision for the given decision context. Internal and external review can also be used to
31 determine if the resulting distributions are valid. It is not clear how the other methods of validation,
32 calibration, coherence, consistency, calibration and informativeness scoring, can be used to
33 determine if the SEE generates useable distributions.
34
35
36
37
38
39

40 The principles support the following choices:

- 41
42
43
44
45
46
47
48
- 42 • Capturing experts beliefs: the elicited beliefs should be fit for purpose. This could be assessed by
43 coherence and consistency (principle 1, 2).
 - 45 • Review: both internal and external review (principle 2, 8)

49 **Section 9.4 Conclusions**

50
51 This chapter considers the choices available from the review of existing guidelines for SEE (see
52 *Chapter 2*) and distinguishes where there is empirical support for the choices, the choices are
53 considered 'appropriate' according to the principles for SEE in HCDM, or where there is neither
54 support from the empirical evidence nor the principles. *Chapters Chapter 5* and *Chapter 6* show that
55 there many choices in SEE, for which there is no empirical support. In addition, the principles applied
56 to the choices, in some circumstances, are unable to provide sufficient justification for discounting
57
58
59
60

1
2
3 particular choices and/or preferring choices above others. For example, on the methods to minimise
4 bias, multiple approaches are available, including training on biases, collecting rationales, and
5 specifying credible intervals. Whilst all of the approaches are potentially valuable, a lack of empirical
6 comparison of the techniques in the context of HCDM, makes it difficult to say conclusively which
7 techniques are most appropriate. Indeed, as with many of the choices in HCDM, the specific
8 application and constraints (see *Chapter 3*) may be a major driving factor in defining the choices for
9 the SEE.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Chapter 10 Reference protocol for expert elicitation in healthcare

Section 10.1 Evidence in support of a reference protocol for HCDM

As stated in the objectives outlined in *Chapter 1*, existing protocols, evidence on specific methods for SEE, consideration of the decision making contexts and the results of the experimental work are combined to propose a reference protocol for SEE in HCDM.

Chapter 9 considers the choices available from the review of existing guidelines for SEE and concludes that, according to the principles for SEE in HCDM, in some circumstances, it is not possible to conclude on the appropriateness of particular choices. One of the uncertain elements is the method to encode judgements, specifically the choice between FIM and VIM approaches. The applied literature on SEE in HCDM (see *Chapter 4*) shows that both approaches have been used, and there are no grounds, based on the principles, to conclusively recommend one choice over another. The experiments presented in *Chapter 8* sought to explore the use of these two methods in HCDM as a primary aim. The chips and bins method was chosen as the FIM and the bisection as the VIM, since these methods have been widely used in HCDM. The experiments also explored experts ability to extrapolate their knowledge and how experts priors are affected by group summaries. Specifically the three experiments sought to determine:

- 1) How the VIM and FIM methods compare in term of procedural performance, where there is no heterogeneity in knowledge.
- 2) Understand whether an individual's 'ability to extrapolate' is related to ('procedural') performance.
- 3) How individuals review their answers in response to a Delphi-type group interaction, and establish whether ('procedural') performance determines the extent of revision.

Chapter 8 suggests that there is little difference between VIM and FIM in terms of procedural accuracy, particularly under conditions of low precision, the situation most likely in HCDM. In terms of extrapolating beyond the data observed by the experts and updating of priors after presentation of group summaries, it is difficult to give definitive messages given that the experiments were not powered for these elements. It is apparent that individuals changed their estimates in a rational way when provided with estimates from others (i.e. when everyone else was discordant, individuals were more likely to change their response, if others were uncertain, individuals were less likely to change), and so group discussion or feedback may be useful, although it does not necessarily produce more

1
2
3 accurate distributions. The need for extrapolation outside of the observed sample and the level of
4 extrapolation does not seem to affect accuracy; therefore it may be reasonable to ask experts about
5 patients and practices in which they do not have direct clinical experience, or for whom there is no
6 relevant literature, and instead experts are required to adapt from one setting to another.
7
8
9

10 **Section 10.2 How the evidence is used to generate a reference protocol for SEE in HCDM**

11 For many of the elements of SEE, multiple choices remain and further research would be necessary
12 to form a preference for different methods in the context of HCDM. Nevertheless, it is important to
13 recognise where there are choices that are emerging as 'best practice' in HCDM, and how these
14 contribute towards the development of a reference protocol in this context.
15
16
17
18
19

20
21 Within HCDM there are multiple opportunities for the use of SEE, from local-level prioritisation to
22 strategic planning for emerging threats. The areas in which it has, perhaps, been applied the most
23 frequently (see *Chapter 4*) is in national level reimbursement, price negotiation and clinical guideline
24 development, an area collectively referred to as HTA. It therefore seems appropriate to think about
25 how the evidence presented, considered and generated in *Chapters 2-9* could be translated into a
26 reference protocol for SEE in HTA. Moving on from this, how decision makers, outside of this setting,
27 can determine the suitability of the reference protocol for their needs can then be discussed. Where
28 substantial uncertainty around recommendations remains, further research may be required. These
29 are considered in *Chapter 12*.
30
31
32
33
34
35
36
37
38

39 **Section 10.3 Reference protocol for SEE in HTA**

40 SEE has been applied in HTA (see *Chapter 4*). However, there are no examples where those
41 developing the exercise have systematically worked through the choices available for each element,
42 and most importantly considered if these choices are appropriate given the intended purpose of
43 the SEE. A reference protocol, even with caveats for particular applied settings, may help to
44 eliminate some of this heterogeneity in methods used.
45
46
47
48
49

50 *Table 143* draws on *Chapters Chapter 8* and *Chapter 9* to suggest choices that are appropriate to
51 consider in HTA, specifically assessments at a national or multinational level. Whilst these are
52 intended to reflect emerging 'best practice' in HTA, given the infancy of SEE applied to HCDM, it is
53 important to recognise that a degree of flexibility on choices may be warranted. In cases where
54 alternative choices are employed, efforts should be made to justify why and describe where the
55 methods used were preferable in that particular application. Although empirical evidence is lacking,
56
57
58
59
60

given the principles of SEE in HCDM, discussed in *Chapter 8*, decision makers should consider the following choices when determining their own reference protocol for SEE:

Table 14 A reference protocol for HTA

Element	Reference methods suggested
Experts	<ol style="list-style-type: none"> 1. Recruitment will be driven by the context, however the SEE should pursue diversity, representing the full range of valid experts beliefs. Experts should be willing to participate. 2. Focus on gathering substantive expertise or experience. Normative skills can be developed during the training session as part of the SEE. 3. Minimize and record conflicts of interest among the experts. Include experts external to the SEE task, i.e. not those involved in developing the task. 4. At least 5 experts should be included in the SEE.
Quantities elicited	<ol style="list-style-type: none"> 1. Simple observable quantities should be elicited where possible; ratios or complex parameters such as regression coefficients should not be elicited directly. 2. Dependence between variables should be captured in SEE. Expressing dependent variables in terms of independent variables is preferable when experts do not have strong normative skills. 3. Wording should be clear and quantities should be decomposed where this means a better fit with experts mental models.
Approach to elicitation	<ol style="list-style-type: none"> 1. Beliefs should be elicited from experts individually, even if a group interaction follows. 2. Although interaction between experts can be structured through face-to-face sessions, constraints in HCDM, such as a lack of experienced facilitators, will usually mean that this will take place via a Delphi style remote process. 2. Between-expert variation should be explored explicitly.
Method	Both VIM or FIM work well, however decision makers should aim for consistency across applications.
Aggregation	<ol style="list-style-type: none"> 1. Statistical distributions should be fitted to experts individually-elicited judgements. 2. Following fitting, a summary of the individual distributions should be obtained using linear pooling with equal weighting of experts.

	3. Any adjustments applied should be to improve coherence and consistency not reduce variability. Internal and external review can be used to assess validity.
Delivery	1. Face-to-face where possible to allow a facilitator to deliver training to the expert. 2. Feedback to experts should be given during the SEE. Following feedback, experts should be given an opportunity to revise their distributions, either during or after a SEE session.
Training & piloting	1. Training is crucial and should focus on avoiding bias and expressing uncertainty. 2. Piloting should be undertaken.
Rationales & documentation	1. Rationales for how the experts made their judgements should be collected post SEE. 2. All methodological choices for the SEE must be documented and justified.

Section 10.4 Important considerations for decision makers outside of the HTA setting

Most HCDM occurs within a HTA setting and at a national level, but elicitation may also be useful for other decision makers, wishing to consider how a reference protocol for their setting may emerge, for example at a local level, or for early technologies that have yet to progress through the regulatory process. In addition, particular types of HTA may encounter additional challenges, for example in rare diseases or genomics. In such settings, a potential reference protocol should consider the following additional issues summarised in *Table 154*.

Table 15 Additional issues in generating a reference protocol outside of HTA

Element	Reference methods suggested
Experts	1. Researchers may have limited access to sufficient experts, for example in rare diseases, therefore expert recruitment may be more challenging and have to rely on peer nomination. 2. Adaptive skills may be required for new technologies since indirect evidence may outweigh directly relevant evidence (e.g. childhood diseases may be informed by adult versions with some extrapolation).
Approach to elicitation	Group discussion may be needed to generate a distribution, for example in early technologies or when eliciting more abstract/complex (non-observable) quantities cannot be

	avoided, for example relating to service delivery, public health programmes or patient pathways.
Method	FIM may be more appropriate for less normative experts or where training cannot be done face-to-face
Aggregation	<p>1. Pooling methods, other than linear pooling, may better reflect expert variability. Further research is needed to explore which methods are more appropriate in these circumstances.</p> <p>2. Weighting may be preferable in some circumstances, for example where experts represent different disciplines or contribute different perspectives on the elicited quantities and therefore considerable heterogeneity is anticipated, but a single agreed consensus distribution is required. Weighting may be achieved implicitly through consensus or explicitly through performance weighting, although it is difficult to see how performance scores would be generated in this context.</p>
Delivery	Practical constraints may dictate remote delivery of SEE, for example through video conferencing.

Section 10.5 Conclusions

This chapter draws together evidence from the preceding chapters to generate elements for a reference protocol for SEE in HTA. Given the infancy of the methods in HCDM and the limited application in this context, it is not possible to be prescriptive regarding methods beyond the more narrowly defined HTA setting. Even within this setting, the reference protocol provides a framework for decision makers to use when generating their own reference protocol, rather than representing a set of guidelines that can be implemented without further consideration of their suitability. Whilst this is the case, given that the methods suggested in this reference protocol are declared most appropriate for HTA on the basis of its defining characteristics, which determine the principles for SEE in HCDM, deviations from the reference protocol should be justified and any limitations discussed in the documentation provided to support the SEE.

There are a number of methodological choices which may involve additional complexities and/or considerations, when used outside of HTA. These are discussed in this chapter and then further in *Chapter 12*. Such choices include the use of consensus aggregation methods, as opposed to individual elicitation, and remote elicitation as opposed to face-to-face. Decision makers outside of HTA at a national level, are recommended to consider these issues when generating reference protocol.

1
2
3
4
5 Finally, this chapter proposes a number of areas in which further research is warranted. This is not a
6 comprehensive list and instead reflects important areas in which the existing reference protocol
7 cannot make recommendations without further research. Some of these areas may require further
8 practical applications of SEE, such as strategies to recruit experts, whereas others may require
9 experimental research, such as that reported in *Chapter 8*.
10
11
12
13
14

15 In addition to the specific methods that require further research, there are some general issues
16 relating to the use of SEE in HCDM; for example, when to elicit and in which areas it is most
17 appropriate. These issues are discussed in *Chapter 12*.
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Chapter 11 Applied evaluation of developed reference protocol

Section 11.1 Background

This chapter demonstrates the reference protocol described in *Chapter 10* by conducting an applied SEE. Originally, the intention of this chapter was to apply the developed reference protocol by performing an elicitation exercise within a decision making process in 'real time'. Members of the research team form an Assessment Group for NICEs technology and diagnostic assessment processes and thus intended to conduct a live elicitation exercise to inform a forthcoming appraisal. However, at the time the reference protocol was developed and ready to be applied, there were no upcoming appraisals with which the team could conduct the elicitation exercise. Consequently, there was a deviation from the original intention of a 'real time' elicitation exercise and the reference protocol was applied in a retrospective manner.

This applied evaluation is based on a diagnostic assessment report (DAR) conducted in Sheffield in the School of Health and Related Research (SchARR).²⁰⁰ SchARR were commissioned by the NIHR HTA programme to produce a model to assess the diagnostic accuracy, clinical effectiveness and cost-effectiveness of three handheld Fractional Exhaled Nitric Oxide (FeNO) monitors: NIOX MINO (Aerocrine), NIOX VERO (Aerocrine) and Nobreath (Bedfont Scientific).²⁰⁰ The analysis aimed to assess the cost-effectiveness of FeNO testing in the diagnosis of asthma in adults and children. However, in the cost-effectiveness model, there were a number of parameters that were missing which were subsequently estimated using SEE. Detailed documentation describing the methods used to obtain and analyse these experts judgments are not included in the report. Without this information, the elicitation process appears to be unstructured, meaning the credibility of the elicited parameters remains unclear. The purpose of this chapter was to apply the reference protocol to this case study and to explore any practical issues.

Section 11.2 The evaluation topic

Here the focus is on a model developed to assess the cost-effectiveness of NIOX MINO, NIOX VERO or NObreath in the diagnosis of asthma in adults and children. In order to illustrate the patient pathways in asthma diagnosis, the next section describes how asthma is currently diagnosed in healthcare.²⁰⁰

Section 11.2.1 Diagnosis of asthma

Detailed guidelines on the diagnosis of asthma have been published and updated by the British Thoracic Society (BTS) and the Scottish Intercollegiate Guidelines Network (SIGN).²⁰⁰ The diagnosis of asthma is a clinical one and there is no standard definition of the condition, nor is there a single gold standard recommendation on how it should be diagnosed. The diagnosis of asthma in children is based on recognising a characteristic pattern of episodic symptoms in the absence of an alternative explanation. Lung function tests are less useful due to variability and the inability of very young children to perform these tests reliably. For both children and adults, the BTS/SIGN guidelines indicate that the severity of asthma should be judged according to symptoms and the amount of medication required to control symptoms.²⁰⁰ Asthma is generally diagnosed in primary care.²⁰⁰

Section 11.2.2 Diagnostic model developed

The diagnostic model determines the expected costs and health losses associated with the misdiagnosis of asthma. Misdiagnosis has different implications for those patients who are false-negative and for patients who are false-positive. For patients who are false-positive, sub-optimal treatment means receiving treatment with asthma medication which will provide no health benefit to the patient (because they do not have the underlying disease). This means there is an additional cost to the NHS without additional health benefits to the patient. In addition, a patient with a false-positive asthma diagnosis may have other more serious pathology, which remains undetected.²⁰⁰

For patients who are false-negative, sub-optimal treatment means not receiving treatment with asthma medication, when in reality the patient would have benefitted from the treatment. Until the diagnosis is corrected, the patient may suffer from poor asthma control and hence lower HRQoL due to asthma symptoms (without experiencing an exacerbation and also by increasing the amount of the time that a patient experiences an exacerbation). Hospitalizations due to exacerbations can be costly to the NHS; hence, a patient with undiagnosed asthma may be more costly to the NHS than a patient who is correctly treated for asthma. These patients may also go on to receive expensive and unnecessary tests such as imaging and referrals to specialists until their misdiagnosis is corrected.²⁰⁰

An incorrect false-negative diagnosis may be corrected later following an asthma exacerbation, due to continued asthma-related symptoms that trigger subsequent appointments and investigation, or due to clinical reconsideration of asthma after tests for other conditions produce negative findings. Similarly, an incorrect false-positive diagnosis may be corrected later due to the continued non-occurrence of exacerbations, a generally high level of HRQoL at very low treatment dosages, thus

1
2
3 indicating that medications are currently being taken by the patient may be unnecessary, or due to
4 continued deterioration due to other more serious underlying pathology. The diagnostic model is
5 intended to reflect the implications of test sensitivity and specificity on subsequent costs and health
6 consequences for the full range of diagnostic options within the available evidence base.²⁰⁰
7
8

9
10
11 The diagnostic model is a simple decision tree (*Error! Reference source not found., Figure 1*). The
12 model estimates the probability that a patient will be diagnosed as true-positive, false-negative,
13 true-negative or a false-positive. The model makes the simplifying assumption that incorrect
14 diagnoses (false-negative and false-positives) are resolved by subsequent tests after some period of
15 time.
16
17
18
19

20 21 22 **Section 11.3 Description of elicited parameters**

23 *Table 165* shows the parameters in the diagnostic model for which evidence was unavailable. In the
24 DAR, the parameter values are provided but there is no documentation detailing how these values
25 or assumptions were reached. The parameter for *time until correct diagnosis* is the only parameter
26 for which the DAR explicitly states that an elicitation process was used to inform these parameters.
27 Subsequently, the remainder of this evaluation focusses on that parameter only.
28
29
30
31
32

33
34 **Table 16 Application of reference protocol: Parameters elicited in DAR²⁰⁰**

Parameter	Source
Diagnostic model parameters	
(1) Resource cost parameters	
Number additional primary care tests: false-positive	Structural assumptions based on expert opinion
Number additional secondary care tests: false-positive	
Number additional laboratory visits: false-positive	
Number additional primary care tests: false-negative	
Number additional secondary care tests: false-negative	
Number additional laboratory visits: false-negative	
(2) Diagnosis QALY gain/loss parameters	
Time until correct diagnosis (years) – false positive	Expert opinion
Time until correct diagnosis (years) – false negative	

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51 Source: From Table 64 in the DAR –page 284

52 53 *Section 11.3.1 Time until correct diagnosis*

54 As described in the DAR and in Section 11.2.2 of this chapter, a false-negative or a false-positive
55 diagnosis of asthma impacts on HRQoL costs, depending on the type of diagnosis. The cost
56 effectiveness results are particularly sensitive to assumptions about the duration of time required to
57
58
59
60

1
2
3 resolve misdiagnosis. The elicitation conducted by SchARR focussed on two questions that were
4 presented to experts:
5
6
7

- 8 • “For someone who has been incorrectly diagnosed as ‘not asthmatic’, how long on average
9 do you think it will take for this incorrect diagnosis to be corrected? What is your 95%
10 confidence interval around this average?”
11
12
- 13 • “For someone who has been incorrectly diagnosed as ‘asthmatic’, how long on average do
14 you think it will take for this incorrect diagnosis to be corrected? What is your 95%
15 confidence interval around this average?”
16
17
18
19

20 Six experts were recruited to provide their responses to these questions but only four experts
21 provided answers. Of the four experts who provided a response, the first expert was the only expert
22 to provide a quantitative estimate. This expert estimated time to correct resolution for a false-
23 negative diagnosis is in the region of 4-12 months while time to correct a false-negative diagnosis
24 may take 12 months or longer. It is important to take into account, the expert noted considerable
25 uncertainty around this estimate. The fourth expert deemed this estimate as “not unreasonable” but
26 this expert also declared these quantities as “*unknowable*”. The remaining three experts provided
27 qualitative responses rather than quantitative estimates declaring the parameters of interest as
28 “*impossible to answer*”, “*unknowable*”.
29
30
31
32
33
34
35
36

37 From the qualitative answers provided by the experts, it is clear that the posed questions do not
38 capture the complexity of an asthma diagnosis and are not representative of the thought-process
39 that an expert may go through to reach a quantitative estimate. In the case of a false-negative
40 diagnosis, this complexity relates to the “*chronicity*” and “*persistence*” of the asthma. In the case of a
41 false-positive diagnosis, the complexity refers to the misdiagnosis never being resolved due to the
42 patients themselves deciding not to “*just stop going back to the doctor*”.
43
44
45
46
47
48

49 The following sections of this chapter demonstrate how the reference protocol described in *Chapter*
50 *10* was applied. The applicability and practicality of the reference protocol is explored.
51
52
53

54 **Section 11.4 Application of developed reference protocol**

55 The following sections describe the application of the developed reference protocol in *Chapter 10*.

56 **Error! Reference source not found.** provides the protocol that was designed for this SEE process.
57
58
59
60

1
2
3 *Section 11.4.1 Selecting the quantities (preparation and design stage)*
4

5 The choice of quantity considered the following three objectives:⁵² fitness for purpose; directly
6 observable and homogeneity in the quantities elicited. Eliciting the same summaries throughout will
7 reduce the burden of training.²⁰¹
8
9

10
11 The time it takes to resolve an incorrect asthma diagnosis for both false-negative and false-positive
12 cases is elicited. Due to the complexity of asthma diagnosis, the parameters were not directly
13 elicited but were calculated from a number of alternative elicited quantities (decomposing the
14 quantities). The quantities elicited relate to the *probability of an event* (i.e. number of patients
15 returning to the healthcare service) *at different time-points*.
16
17
18
19

20
21 Based on the qualitative feedback from the experts in the DAR, it was clear that there were a
22 number of aspects of the condition that needed to be incorporated into the questions to ensure that
23 these were asked in a manner which would be consistent with how the expert thinks about the
24 condition. In terms of asthma, these characteristics relate to the level of chronicity and persistency
25 of symptoms.
26
27
28
29

30
31 Following personal communication with a General Practitioner (GP) in the University of York, it was
32 decided that specific patient vignettes would be described and presented to the experts at the
33 beginning of each question.
34
35
36
37

38 These were presented separately for false-negative and false-positive adults and children. Each
39 described a type of patient based on varying levels of symptom severity (mildly persistent
40 symptoms, moderately persistent symptoms or severely persistent symptom). The experts' were
41 asked to express the proportion of the type of patient described returns to the healthcare service at
42 certain time-points since their first diagnosis. Eliciting by severity and for adults and children
43 separately is intended to reflect the heterogeneity of the asthmatic population, raised in the DAR
44 elicitation.
45
46
47
48
49

50
51 Variation in the time to correct diagnosis for both false-negative and false-positive patients is
52 anticipated and thus the questions were asked based on two separate time-points for both types of
53 incorrect asthma diagnosis. This approach also supported the assumption made in the model that
54 all incorrect asthma diagnoses are resolved. For false-negative patients the time-points included
55
56
57
58
59
60

1
2
3 were 6 and 12 months while for false-positive, the time-points were 12 and 24 months. **Error!**
4 **Reference source not found.** presents the questions and preambles provided to the experts.
5
6
7

8 *Section 11.4.2 Methods to encode judgements (preparation and design stage)*

9
10 To elicit uncertainty two methods of elicitation were explored in *Chapter 8*: the Chips and Bins
11 method and the Bisection method. *Chapter 10* concluded that either of these methods is
12 appropriate for HCDM, and thus here the recruited experts completed one of these methods of
13 elicitation; either Chips and Bins method **or** Bisection method.
14
15
16
17

18 The reference protocol in *Chapter 10* states the both VIM and FIM work well but decision makers
19 should aim for consistency across application. The evaluation in this chapter used both methods as it
20 was intended to further explore the usability of the two methods with actual health care
21 professionals.
22
23
24
25
26

27 *Section 11.4.3 Validation (preparation and design stage)*

28 At the end of the SEE, experts were asked if they were confident the answers they gave reflected
29 their views and uncertainties. Response options were 'yes', 'not sure' and 'no'. If they responded,
30 'no' or 'not sure', they were asked to provide more detail as to why in an open question. Other
31 forms of validation were not used.
32
33
34
35
36

37 *Section 11.4.4 Selecting experts (preparation and design stage)*

38 As asthma is generally diagnosed in primary care, it was assumed Primary Care GPs would have
39 substantive knowledge in the diagnosis of asthma. Consequently, Primary Care GPs were recruited
40 as the experts. Experts were not expected to have any normative skills (see *Chapter 10*). The experts
41 were recruited using recommendation from peers.
42
43
44
45
46

47 As this work is a retrospective evaluation of the reference protocol developed, rather than a SEE per
48 se, a smaller number of experts, than usually recommended, were recruited. Four experts were
49 recruited. The intention was to utilize the protocol rather than generate evidence and if sufficient
50 time were available more experts would have been recruited. However, as mentioned at the
51 beginning, this chapter focused on the design elements of the SEE rather than the practical conduct.
52 Two experts completed the Chips and Bins Method and two experts completed the Bisection
53 Method. These were allocated randomly by the facilitator.
54
55
56
57
58
59
60

1
2
3 At the beginning of the exercise, experts were asked to provide some background information about
4 themselves. This included, the number of years they have been working in general practice and how
5 they commonly diagnose an adult or child who they suspect may have asthma. The experts were
6 asked to identify whether they use an objective test (spirometry, reversibility testing), a clinical
7 evaluation or both methods.
8
9

10 11 12 13 *Section 11.4.5 Pilot exercise (preparation and design stage)*

14 The wording of the questions was piloted for clarity and adequacy. The pilot exercise was sent to
15 two GPs and feedback sought. Following feedback the questions were modified, specifically the
16 wording of the questions.
17
18
19

20 21 22 *Section 11.4.6 Training and preparation for experts (preparation and design stage)*

23 A narrated power-point training session was embedded within exercise. The training session
24 described the objectives of the elicitation exercise, clarified concepts such as uncertainty,
25 familiarised the experts with the quantities elicited, described and explained the impact of bias and
26 heuristics, and trained experts on the methods of elicitation used.²⁰¹
27
28
29

30
31
32 Experts were also reminded throughout the SEE that they were to elicit uncertainty on their
33 estimate rather than thinking about variability across this heterogeneous group of patients
34
35
36

37 38 *Section 11.4.7 Level of elicitation (elicitation stage)*

39 Each expert elicited their judgements individually without interaction with other experts. Eliciting
40 judgements individually reduced the risk of estimates being biased by a subset of experts. In the SEE
41 elicitation literature, there are concerns that experts may not feel confident in eliciting judgements
42 individually, however, the experts in this SEE process elicited their beliefs on a condition that they
43 encounter regularly in general practice. Concerns regarding individual level elicitation and lower
44 confidence amongst experts generally arises when dealing with problems/technologies or conditions
45 that are new or unknown to the experts (see *Chapter 6*).
46
47
48
49
50

51 52 53 *Section 11.4.8 Mode of administration (elicitation stage)*

54 The elicitation exercise was administered via a computer-based method using a *de novo* tool in
55 Excel. The evaluation used a mixture of face-to-face and remote forms of administration. Despite
56 using individual level elicitation, a facilitator was present, either in person or on the phone at the
57 time the expert completed the exercises. The purpose of this was to gather as much feedback as
58
59
60

possible on the elicitation process (see Section 8.5.5). For example, the time it took to complete the exercise or to record any difficulties the expert had when completing the process.

Section 11.4.9 Feedback to experts and revision (elicitation stage)

Once experts expressed their beliefs and completed each question, they were presented with graphical feedback of what their estimates looked like (see *Chapter 10*). In the Chips and Bins Methods, experts were able to see how the grid looked once they have placed all of their chips on it. Similarly in the Bisection method, experts were able to see the breakdown of the different values they provided (median, upper and lower quartile etc.). The individual level of elicitation that was chosen meant that group consensus was not required and consequently, group feedback to the experts was not necessary. Both methods had a reset button. Once the expert completed each question, they had an opportunity to click reset and begin that particular question again.

Section 11.4.10 Opportunity for interaction (elicitation stage)

Given the individual level of elicitation that was chosen, there was no opportunity for interaction between the experts.

Section 11.4.11 Feedback from experts on process (elicitation stage)

Qualitative feedback on the elicitation process was collected from the experts. This was collected by the facilitator using a feedback questionnaire post exercise. The feedback questions assessed the following concepts;

- observability of the quantity asked
- based on a 5-point scale, assess how easy or difficult the experts found the completion of the exercise
- based on a 5-point scale, evaluate whether the wording of the questions were easy or difficult to understand
- whether the provided training is sufficient
- whether the expert would prefer to have some interaction with a colleague or another expert (if they were to complete the exercise again)
- if they would be willing to complete this exercise again without a facilitator

The feedback also asked experts to suggest improvements they think necessary for any future SEE. In addition, any useful comments or suggestions made by the expert throughout the SEE were also collected. Section 11.4.11 discusses the feedback from the experts. While not collected as part of the

1
2
3 feedback questionnaire, at the end of each section in the exercise, rationales from the experts about
4 how they made their judgements were collected. This form of validation helps to highlight if experts
5 understood the task and responded as best they could.
6
7
8
9

10 *Section 11.4.12 If/how to aggregate (aggregation, analysis and post-elicitation)*

11 As an individual level of elicitation was chosen, mathematical aggregation should be applied to
12 generate the distributions, specifically linear opinion pooling using equal weighting of experts (see
13 *Chapter 10*). In this application, the intention was to explore the use of the reference protocol,
14 rather than generate a single distribution relating to the uncertain quantities of interest, and
15 therefore this aggregation is not undertaken.
16
17
18
19
20
21

22 *Section 11.4.13 Fit to distribution (aggregation, analysis and post-elicitation)*

23 A Beta distribution would be fitted to experts distributions as these relate to probabilities.
24
25
26

27 *Section 11.4.14 Data Protection and Anonymity (aggregation, analysis and post-elicitation)*

28 Experts were asked to give their opinions individually (not in groups). The information provided,
29 including personal details, is kept anonymous and confidential, stored securely and only accessed by
30 those carrying out the study.
31
32
33
34
35

36 **Section 11.5 Results**

37
38
39 The number of years the experts have been working in general practice ranged from 6 to 35 years
40 (*Table 1716*). When asked how they would usually test an adult with suspected asthma, all four
41 experts reported using both an objective test (e.g. FeNO, Spirometry or reversibility testing) and a
42 clinical evaluation. When diagnosing a child with suspected asthma, three of the four experts
43 reported using just a clinical evaluation and one expert reported using both a clinical evaluation and
44 an objective test. Expert 1 reported that for this questions and the remainder of the questions in the
45 evaluation, the age of the child population should be defined as follows: <1 (should not have any
46 diagnosis as their lungs will not be developed), 1 to 4 years, 5 to 14 years, older than 14 years. This
47 expert then went on to explain that the older the child the more likely the chance a GP could use an
48 objective test in addition to a clinical evaluation to diagnose asthma.
49
50
51
52
53
54
55
56
57
58
59
60

Table 17 Application of reference protocol: Summary of experts recruited

Expert	Years GP	Asthma diagnosis (Adult)	Asthma diagnosis (Child)	Elicitation Method	Mode	Completion time
GP1	6	Test and clinical evaluation	Clinical evaluation	Chips and Bins	Face-to-face	1 Hour
GP2	35	Test and clinical evaluation	Both	Bisection	Face-to-face	45 mins
GP3	24	Test and clinical evaluation	Clinical evaluation	Bisection	Remote	1.5 Hours
GP4	30	Test and clinical evaluation	Clinical evaluation	Chips and Bins	Face-to-face	1 Hour

Table 18 Application of reference protocol: Elicitation Results

Expert	Patient type (by severity)	False-positive (adults)		False-positive (children)		False-negative (adults)		False-negative (children)	
		12 months	24 months	12 months	24 months	6 months	12 months	6 months	12 months
GP1	Severe	30, 70	60, 90	Same judgements as adult patients		60, 99	99, 100	70, 100	99, 100
	Moderate	30, 70	45, 75			45, 85	99, 100	60, 100	99, 100
	Mild	10, 60	10, 65			0, 55	99, 100	10, 60	99, 100
GP2	Severe	5, 25	99, 100	5, 10	99, 100	70, 90	99, 100	85, 95	99, 100
	Moderate	15, 35	99, 100	5, 15	99, 100	40, 60	50, 70	80, 90	99, 100
	Mild	30, 50	35, 55	20, 30	99, 100	35, 45	40, 50	50, 60	99, 100
GP3	Severe	90, 100	90, 100	Same judgements as adult patients		90, 100	90, 100	Same judgements as adult patients	
	Moderate	50, 90	60, 80			60, 90	45, 80		
	Mild	25, 50	25, 50			1, 50	1, 50		
GP4	Severe	75, 100	75, 100	83, 100	83, 100	75, 100	75, 100	83, 100	83, 100
	Moderate	60, 90	65, 95	66, 99	72, 100	60, 90	65, 95	66, 99	72, 100
	Mild	30, 70	35, 75	33, 77	39, 83	30, 70	35, 75	33, 77	39, 83

Section 11.5.1 Elicitation Results

As discussed in Section 11.4.12, the intention of this evaluation was to explore the use of the reference protocol rather than generate a single distribution. Therefore, *Table 18* simply presents the ranges (upper and lower limits) reported by the experts for false-positive and false-negative adults and children based on different patients types (by severity). Taking the complexity of asthma into account, it is expected that a lower proportion of patients (adults and children) with mildly persistent symptoms will return to the healthcare service compared to patients with severely or moderately persistent symptoms. It is also expected that at the second time-point, a higher proportion of patients will return to the healthcare service compared to the first time-point for each patient type. When comparing the ranges for adults and children, it is expected that overall, a higher proportions of children will return to the healthcare service due to parents concern. For the most part, these expected ranges were reported by the experts, which indicate that the experts understood what was being asked of them and that this concept was something they could think about in their own general practice experience.

When comparing the ranges questions based on false-positive adults and children at the first time-point, the ranges provided by GP2 seem to move in the opposite direction as to what is expected. However, in this experts experience, children are less likely to come back at this time-point as they are easier to diagnose at an earlier stage compared to adults.

Section 11.6 Feedback from experts on the process (including rationales and validation)

Section 11.6.1 Observability of the quantity asked

Three of the four experts reported that the selected quantity was observable; that the proportion of patients returning to the healthcare service within a particular time period was something they could express their opinion about. However, GP3 said this was not something that he could think about because, in his opinion, misdiagnosis of asthma does not happen that often.

Section 11.6.2 Completion of exercise

GP1 completed the Chips and Bins method. This GP explained that the method at first seemed daunting but after completing one of the questions, deemed the method as straightforward. Both GP2 and 3 completed the Bisection method but reported conflicting feedback on the completion of the exercise. GP2 found the Bisection method easy to complete while GP3 reported the method as very difficult to complete.

1
2
3 *Section 11.6.3 Wording of the question*

4 When asked whether the wording of the questions were easy or difficult to understand, GP2 and
5 GP4 reported the wording as *easy to understand* while GP1 and 3 found the wording *very difficult to*
6 *understand*. GP1 provided further detailed feedback on this and suggested that the preambles
7
8 should include more detail in terms of defining the severities of asthma, i.e. more description. All
9
10 GPs identified the training session as sufficient and GP1 suggested a practice exercise would be
11
12 useful in the training session.
13
14

15
16 *Section 11.6.4 Value of interaction*

17 When asked if they thought interaction with other experts or colleagues would be useful; GPs 2, 3
18 and 4 thought this would be beneficial. Their reasoning for this was that to hear other experts
19
20 rationales, to ensure all experts are making judgement on the same issue and that a small group of
21
22 experts allowed to interact and achieve a consensus would give a more rounded view. However, GP4
23
24 did emphasise that it would be important to avoid the interaction from becoming dominated by one
25
26 expert in the group. GP1 did not think interaction with other colleagues was important in this case.
27
28 This expert was of the opinion that GPs should be adequately familiar with asthma to confidently
29
30 answer the question independently.
31
32

33 *Section 11.6.5 Value of facilitator*

34 Experts were asked that if they were to complete the exercise again, would they be happy to
35
36 complete it without the use of a facilitator. Three out of the four experts said yes they would be
37
38 happy for the exercise to be non-facilitated. GP3 stated that if the process was not facilitated, the
39
40 requirements of the task would be unclear.
41
42

43 *Section 11.6.6 Experts rationales*

44 When reporting on their thought-processes, all experts considered similar patients they encounter in
45
46 general practice. One of the experts explained that when a patient has an asthma diagnosis noted
47
48 on their records, the patient is be invited to attend an annual visit for a respiratory check-up. In the
49
50 GPs experience, 85% of these patients will return for their annual visit and subsequently, this expert
51
52 reported using this figure as a guideline when making the judgements. Three of the four experts
53
54 stated they were confident that the answers they gave in the exercise reflected their own views and
55
56 uncertainty. GP3 was not sure of this and explained that the relevance of the questions to clinical
57
58 practice was not obvious. In addition, the expert found the exercise repetitive due to the two time-
59
60 points making the exercise tedious to complete.

Section 11.6.7 General feedback from experts

Experts provided general comments and improvements for future SEE processes. When completing the Chips and Bins method, GP1 found the method cognitively challenging. The expert provided an upper and limit for each question but did not fill in the grid using the chips due to the grid being different sizes for each question and varying chips available. While the facilitator acknowledged this was correct and that the grid size and the amount of chips available are dependent on the range given by the expert, the expert did not place the chips in the bins to show certainty/uncertainty on the proportions.

As described in Section 11.4.1, two time-points were used in the false-negative and false-positive descriptions. Given the layout of the exercise, experts had to scroll back to the initial time point if they wanted a reminder of the previous judgement they made before providing their judgement at the second time-point. Two of the experts said it would be more accommodating if the first time-point was visible while answering the second as the first judgement would serve as a benchmark. In essence, the experts suggested that for future elicitation processes, if questions are a follow-on from a previous question, it would be useful if the previous judgements were easily accessible.

Section 11.6.8 Practicality of conducting the SEE process

The design and conduct of the SEE was undertaken over a 7 month period (August 2018 to February 2019 - excluding any form of aggregation or fitting (see Section 11.4.12) and involved three researchers over that time-period. In terms of analyst resources, this included one 0.6 full time equivalent (FTE) for the duration of the process, with the addition of a 0.1 FTE for the final 3 months and an additional 0.5 FTE for the remaining 2 months. This covered development of the questions and subsequent piloting of the wording of the questions; developing the training sessions and developing the excel-based elicitation exercise. Expert recruitment accounted for one month of the study period (January 2019). Administration and completion of the elicitation exercises along with the write-up was conducted during the final month of the process (February 2019).

Section 11.7 Conclusion

The aim of this chapter was to apply the reference protocol to a case-study and to explore any practical issues. This highlighted a number of key issues in the SEE process relevant in a HCDM context. It is clear from this SEE process and the feedback provided by the experts, that sufficient information needs to be presented to the experts. The level of information presented to the experts and the wording of this information is paramount in ensuring that the quantity of interest is observable to the expert. When deciding on the information to provide to experts in a HCDM

1
2
3 context, based on the rationales provided by the experts in this process, it may be useful to consult
4 existing policies.
5
6
7

8 In a HCDM context, a SEE process will be subject to timeline constraints. Certain available choices in
9 SEE may result in a more lengthy process for example face-to-face modes of administration or
10 interaction between experts. Careful consideration must be given to these choices to achieve
11 accurate judgements from experts but also to make efficient use of available time. In terms of
12 interaction between experts, the feedback from experts in this SEE process indicates that
13 consideration needs to be given to the potential value of interaction between experts. Depending on
14 the context of the SEE, interaction between experts may be more essential, therefore justifying a
15 lengthier process. For example, if the SEE process is focusing on a new drug or a rare disease,
16 interaction between experts may be more significant than a process focussing on an established
17 drug or a commonly encountered condition or illness. In the latter, experts will be more familiar with
18 the context and should therefore have the ability to independently provide a response in a confident
19 manner.
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Chapter 12 Discussion & conclusions

This chapter discusses the evidence that has been generated in Chapters 2-11. It goes on to consider how the reference protocol for SEE may be used by policy makers to define their own reference protocol that reflects their particular constraints. To do this this it considers the feedback from a workshop convened as part of the project. Areas for further research emerging from the work conducted, and discussed at the workshop, are also discussed. Finally the limitations of the work are noted.

Section 12.1 Conclusions on evidence generated

SEE can offer opportunities in HCDM, particularly reimbursement decisions supported by MBEE. SEE allows the uncertainty in the evidence used to populate these models to be characterised, or where evidence is completely lacking, provides additional information needed to reach a decision.

The work described in this report has attempted to generate evidence which is useful for analysts and decision makers in HCDM. SEE conducted in this context to date, has not used a set of consistent methods, and, above all, has not considered the implications of the choices made when designing and conducting a SEE. To improve the accountability of HCDM the procedure used to derive expert judgements should be transparent.

A reference protocol for SEE in HCDM is proposed in *Chapter 10*. The reference protocol is intended to serve as a guide to good practice and reporting, rather than being prescriptive regarding methods, and thus it is intended more as a guidance rather than a protocol. This was necessary due to the lack of empirical evidence underpinning methods choices specific to HCDM. Instead choices were considered according to the principles for SEE in HCDM, set out in *Chapter 9*. These nine principles were developed based on findings from Chapters 3-7, which consider: the constraints in HCDM, how SEE has been applied in the literature and challenges faced, evidence relating to particular methodological aspects of design and conduct and considerations in choosing alternative quantities that can be elicited. The principles also reflect 'good practice' in SEE more generally, as reported by Cooke (1991).¹⁴

As stated in *Chapter 10*, the lack of evidence relating to HCDM, and a paucity of applied studies, meant that the reference protocol focussed on the more narrowly defined setting of national level HTA. Whilst this encompasses a range of activities that could include SEE, it does not reflect more

1
2
3 complex settings which could pose additional challenges. For example, HCDM at a local level, for
4 early technologies that have yet to progress through the regulatory process, or for specific types of
5 HTA, including rare diseases or genomics. These settings may require a different approach to
6 elements such as recruiting experts, level of elicitation and delivery. It is recommended that such
7 decision makers consider the elements of the reference protocol and how these translate to their
8 setting. In doing so they can determine a reference protocol of their own.
9
10
11
12
13
14

15 The experiments, described in *Chapter 8*, suggested that there is little difference between VIM and
16 FIM methods to encode judgements. The reference protocol therefore stated that a decision maker
17 can consider either of these choices suitable, however consistency across applications is preferred,
18 i.e. they should choose either VIM or FIM and use this throughout their decision making processes.
19 The experiments also sought to explore extrapolation beyond data observed and updating of priors
20 after presentation of group summaries – issues which feed into multiple choices for SEE. It was
21 difficult to form definitive conclusions given that the experiments were underpowered for these
22 elements. To make definitive statements regarding these aspects of SEE, further experimental data
23 would need to be collected. However, the experiments provided some evidence that experts
24 changed their estimates in a rational way when provided with distributions from others, suggesting
25 that group discussion or feedback may be useful. Extrapolation outside of the observed sample does
26 not seem to affect accuracy, suggesting that it is reasonable to ask experts about patients and
27 practices in which they do not have direct clinical experience, or for whom there is no relevant
28 literature.
29
30
31
32
33
34
35
36
37
38
39

40 *Chapter 11* applies the reference protocol for SEE in HCDM to an existing NICE appraisal. This relates
41 to a diagnostic model for asthma developed as part of the NICE diagnostics programme. The
42 recommendations from the reference protocol were used to determine the following aspects of the
43 SEE: selection of quantities, method to encode judgements, validation, selection of experts, piloting
44 and training, level of elicitation, model of administration, feedback, interaction and post elicitation
45 aggregation. Quantities relating to incorrect diagnosis of asthma were collected from a sample of
46 experts. Time taken to develop and conduct the SEE was recorded and feedback on the elicitation
47 was also sought from experts.
48
49
50
51
52
53
54

55 There were only small numbers of experts on which to base conclusions, however from those that
56 completed the SEE, the VIM and FIM seemed to be equally challenging for experts to complete. All
57 four experts found the training to be useful. A facilitator was used in the SEE, however three
58
59
60

1
2
3 participants stated that they would be happy to complete the exercise without the facilitator and
4 three stated that they would prefer interaction with another expert, so long as other experts were
5 familiar with asthma and its diagnosis. This finding may be worth considering in settings where it is
6 not possible to gain access to a facilitator. There may be value in allowing experts to interact aside
7 from its use to generate consensus, such as developing a common problem structure or sharing
8 knowledge and information between experts (see Section 5.4). The need for interaction between
9 experts in particular settings, for example in rare diseases, was discussed in *Chapter 10*.

16 **Section 12.2 Key considerations for using the reference protocol in HCDM**

17 To consider how the reference protocol may be used by HCDMs, a workshop was convened. HCDM
18 stakeholders, including practitioners, policy makers and methodologists, attended. Feedback is
19 described in further detail in **Error! Reference source not found.** Briefly, the workshop considered
20 the acceptability of the reference protocol for SEE in HCDM, how the proposed reference protocol
21 for elicitation may be implemented and where it would be most useful. It was also used to identify
22 priorities for further research and development of reference protocol and its use for HCDM.

23
24
25
26
27
28
29
30 There was unanimous support for a reference protocol or guidance on SEE in HCDM. Workshop
31 attendants discussed what form this guidance should take, and which would be the most useful,
32 specifically should it be prescriptive or guiding principles. Those involved directly with conducting
33 SEE, tended to suggest that less prescriptive guide would be most useful, as some of the
34 methodological decisions may be driven by the context (see *Chapter 3*), for example a specific
35 appraisal may require SEE to generate results within an extremely short time frame, reducing the
36 possibilities for face to face SEE. Whatever form it takes, workshop participants thought that some
37 form of guidance would help decision makers, considering evidence generated from a SEE, or when
38 planning to conduct their own SEE. Lack of guidance is also seen as a barrier to publishing SEE in
39 HCDM, and therefore a reference protocol would support the dissemination of applied research in
40 this area. It may also encourage the development of materials to assist SEE, for example generic
41 training materials, which are currently lacking.

42
43
44
45
46
47
48
49
50
51
52 Workshop participants agreed that a reference protocol may be most useful when there is a lack of
53 substantial existing evidence, such as urgent delivery systems during epidemics, or as a complement
54 to existing data on a longer-term basis, e.g. short trial follow up. There is likely to be a need in areas
55 which are not represented in trials, such as histopathology. The reference protocol developed here
56 is considered appropriate for national level HTA, which is also the audience most likely to be
57
58
59
60

1
2
3 receptive and with sufficient resource to conduct or commission SEE. Within national HTA there may
4 be potential to use a reference protocol for SEE within clinical and public health guidelines.
5
6
7

8 The time and expertise required to conduct a SEE may be an issue for many of the formal decision
9 making processes that exist, for example the NICE appraisals process. The evaluation undertaken in
10 *Chapter 11* took over 5 per months FTE from start to finish. It is not clear if previous applied
11 examples of SEE may have been forced to make particular choices on the basis of limited time and
12 resource (see *Chapter 4*). This may compromise the quality of the SEE and therefore it may be more
13 appropriate to extend timelines so as to incorporate well conducted SEE. Justification for this is that
14 lots of time is spent on generating evidence and this is just one form of evidence. This may still be a
15 challenge for some decision makers where there are multiple uncertain quantities that could be
16 elicited, thus potentially imposing high costs. A potential solution is to choose parameters for the
17 SEE based on the expected value of partial perfect information corresponding to each parameter, or
18 on a less-formal sensitivity analysis determining the impact of extreme parameter values, in a cost-
19 effectiveness model.
20
21
22
23
24
25
26
27
28
29

30 A number of specific issues regarding use of the reference protocol were raised during the
31 workshop. Firstly, it was unclear who would be held accountable for SEE. Ultimately the decision
32 makers were accountable for the decision which utilised the SEE, however experts may also be
33 accountable for the beliefs that they express. In order to make this more explicit, participants agreed
34 that decision makers need to have access to individual elicitation and not just a group/consensus
35 judgement. Recognition of accountability may lead experts to alter their beliefs. Second is the issue
36 of which experts are included in the SEE. In some circumstances recruitment of experts may have to
37 rely on methods such as peer nomination, particularly where there are constraints on time. This may
38 give an unrepresentative sample of views and may be more likely to result in motivational biases,
39 perhaps due to an association with the quantities of interest. The most knowledgeable experts may
40 not always be available for SEE, so that the aggregated distributions may not be representative of
41 the current level of knowledge on the quantity. Thirdly the issue of choosing observable quantities
42 may not be straightforward. Strictly speaking an observable quantity is something that can be
43 measured, which may be the case for the majority of quantities that need to be elicited. This
44 excludes indirectly observed quantities such as odds ratios. Experts may also have different
45 experiences which alters their perception of what is observable and non-observable. In the
46 evaluation (see *Chapter 11*), one of the experts stated that the quantity was not, in their opinion,
47
48
49
50
51
52
53
54
55
56
57
58
59
60

observable, as misdiagnosis of asthma happens rarely, whilst the other three experts stated that the quantity was observable to them.

Section 12.3 Key areas for further research

In considering the appropriateness of choices for SEE in HCDM and exploring how these choices may be affected by the context in which the SEE is applied, there are areas in which further research is required before definitive statements can be made regarding their appropriateness for a reference protocol. These areas were discussed at the workshop and refined following discussion. In ensuring that SEE is used consistently in HCDM and reflects the constraints of that particular setting, not all of these may represent priorities for further research. Workshop participants were not asked to prioritise topics *per se* or consider which issues are most crucial to the accuracy of SEE, and therefore the list does not reflect which topics may be most urgently required. Workshop participants were instead asked to consider what additional evidence decision makers in HCDM may require when determining a reference protocol for SEE for use within their setting. Areas of uncertainty, in the current reference protocol were: selecting experts, minimising bias, adaptation to specific setting in which SEE may be applied, for example choosing individual or group elicitation, appropriate working of questions, methods for multivariate elicitation and what information should be presented to the experts to help them formulate their beliefs.

Examples are summarised in *Table 19*. Some of these could easily be adapted into researchable questions, whereas others are much more vague and general. Some of these topics would benefit from empirical research and others may be resolved through application of the proposed reference protocol to HCDM, including in settings with a range of constraints.

Table 19 Areas for further research on SEE in HCDM

Decision choice and general research area	An example of a specific question
Selection of experts	
How to determine a sample size for SEE	In individual elicitation, what is the saturation point of increasing sample size?
Exploration of strategies for recruiting experts in HCDM	Which methods for expert recruitment are most practical and what are the challenges?

1 2 3 4 5 6 7	Methods to assess experts skills that are appropriate for the SEE task	Where adaptive skills are required, how can these be measured and, where these skills are compromised, can training increase these skills?
8 9 10 11 12	What minimum level of normative expertise is required	What additional level of normative expertise is required when eliciting more complex quantities or where dependence exists?
13	Biases	
14 15 16	Training strategies	What training strategies can be used to minimise bias?
17 18 19	Recruitment	What recruitment strategies can be used to minimise expert bias, beyond minimising financial/competing interests?
20 21 22	Validation	Can the measurement of expected bias provide a mechanism to validate elicitation?
23 24	Validation of experts	
25 26 27	Performance assessment	How many seeds are required to estimate experts expected accuracy in HCDM, and how can these be efficiently generated?
28 29 30	Calibration	To what extent might performance-based weighting improve the validity of resulting distributions?
31 32 33	Accuracy	What is the relationship between characteristics of experts and accuracy of elicited quantities?
34 35	Quantities	
36 37 38	Dependence methods	Which methods for eliciting dependent quantities work best for non-normative experts?
39 40 41	Consistency	Does elicitation of consistent quantities throughout the task improve procedural accuracy?
42 43 44	Survival parameters	How to elicit parameters of survival models, in particular uncertainty relating to these?
45 46 47	Group elicitations and interaction	
48 49 50	Consensus approach	Which consensus approach works best in HCDM in practice and for which types of quantities and decision makers?
51 52 53	Sample size	How many experts should be part of a consensus elicitation process, and does this differ by context?
54 55 56	Aggregation	

Distribution fitting	What methods for fitting distributions to elicited beliefs are most appropriate for particular quantities, for example more complex quantities
Combining priors	Should individual priors be combined when there is significant expert variation? If so, how?

Section 12.4 Limitations of the work conducted

Whilst the reference protocol developed here represents a significant move forward in terms of SEE applied to HCDM, there are a number of limitations of the work that are worth noting.

Firstly, in developing a reference protocol for SEE in HCDM, it was necessary to draw on multiple types of evidence, structured (systematic) review, targeted literature search, and experimental analyses. In identifying relevant evidence from the existing literature it was not possible to use systematic search methods for all reviews, and the targeted searches used a semi-structured approach (*Chapter 5*). This may have resulted in relevant studies that are less well-known in the elicitation literature also being missed in our reviews. In addition, it was not possible to conduct targeted searches for all elements of SEE.

Second, there were a number of compromises that were needed in order to generate empirical evidence relating to the choice between the FIM and VIM (*Chapter 8*). Foremost, it was necessary to use students to represent health care professionals. The design of the experiments was such that, prior clinical knowledge was not required to complete the tasks, and therefore the experts were instead meant to represent the level of normative skills that would usually be expected in HCDM. Participants were standardised according to the level of knowledge they observed from the simulated learning process. However, in practice experts in HCDM are likely to draw on multiple sources of knowledge when formulating their beliefs, i.e. health carers may also draw on published evidence, peer contact or other related evidence or experience. It was not possible to reflect these multiple forms of knowledge in exploring the performance of the methods to encode judgements. The experimental set up, more generally, may impact on the generalisability of the results.

Thirdly, it was not possible to explore all of the uncertain choices empirically through the experimental approach described in *Chapter 8*. In addition to the comparison of the FIM and VIM approaches, *Chapter 8* also looked at how experts update their beliefs when presented with group

1
2
3 summaries and extrapolation beyond data observed. It was not possible to power the experiments
4 to detect differences for these two elements, and therefore it is difficult to reach conclusions
5 regarding these comparisons.
6
7

8
9
10 The major limitation of the work conducted here, lies not in the methods employed, but the
11 evidence available from the wider literature, on which to base the set of choices and determine how
12 appropriate these are. Concluding on the suitability of the choices available from the existing
13 guidelines is challenging due to the lack of empirical evidence to support specific choices. Instead it
14 was necessary to develop principles for SEE in HCDM, using the sources of evidence as described
15 above and published guidelines for good SEE. Using the principles, meant that it was not always
16 possible to give definitive conclusions on choices.
17
18
19
20
21

22
23 This flexibility, however may be a useful characteristic of the reference protocol developed here.
24 Trying to define a reference protocol that is useful, in that it refines the set of choices, but is
25 sufficiently flexible that it can be applied across HCDM and considers constraints in different
26 settings, may provide the type of guidance that is most useful at this stage. Further applied studies
27 of SEE in HCDM, which consider the choices specified in this reference protocol, and thoroughly
28 document these, will help to generate valuable evidence on the usefulness of the reference protocol
29 and may also provide opportunities for empirical comparisons of some of the remaining uncertain
30 choices, for example using the approach in *Chapter 8*.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Chapter 13 Acknowledgements

This work was funded by a grant from the Medical Research Council (MRC): MR/N028511/1 “HEE: Developing a reference protocol for expert elicitation in health care decision making”.

We would like to acknowledge the contributions from the advisory group for the project, specifically, John Paul Gosling, Jeremy Oakley, Anthony Hatswell, Nicky Best, Matt Stevenson, Peter Ayton, Tim Bedford, Maarten Ijzerman, Rebecca Albrow, Frances Nixon, and Roger Cooke.

A workshop was convened to consider the results from the project and how these may be used in practice by decision makers. The workshop was attended by NICE chairs, NICE staff, academics, NHS England, WHO, industry representatives, other methodologists and decision makers. We would like to acknowledge and thank all that attended and participated.

We would like to acknowledge Ian Watt, Alastair Dickson, Sean O’Connell and Mark Williams for their help and participation in the expert elicitation exercise as part of the application of the developed reference protocol in *Chapter 11*.

Finally in designing the experiments conducted in *Chapter 8*, piloting was carried out on colleagues at the Centre for Health Economics. We would like to acknowledge and thank all colleagues that took part and provided valuable feedback. In addition colleagues at the Experimental Economics Lab at the University of York.

Data sharing:

All data requests should be submitted to the corresponding author for consideration. Access to anonymised data may be granted following review.

Chapter 14 Contributions of authors

Dr Laura Bojke (Reader Health economics) was responsible for day-to-day management of the project, developing the reference protocol, supervising the evaluation, running the experiments and drafting the report. Dr Marta Soares (Senior Research Fellow) was jointly responsible for the day-to-day running of the project, designed the experiments in Chapter 8 and wrote Chapter 7 with CJ and LS. Dr Chris Jackson provided statistical support for the experiments and generated measurements to assess performance in experiment 1. He wrote Chapter 7 with MS and LS, Chapter 5 with AF and AC, and worked with MS and DJ to draft Chapter 8. Professor Linda Sharples (Medical Statistics) provided senior statistical support to the project, specifically in designing the experiments in Chapter 8. She wrote Chapter 7, with CJ and MS, contributed to Chapters 4, 6 and 8. Professor Karl Claxton (Health Economics) contributed towards the design of the experiments in Chapter 8, and helped to draft materials in Chapters 9 and 10. Dr Dina Jankovic (Research Fellow) developed the statistical code for the experiments in Chapter 8, conducted the experiments with MS and LB and helped the draft Chapter 8. Dr Aimee Fox (Research Fellow) was responsible for Chapter 5, with CJ and AC, and worked with LB to conduct the evaluation in Chapter 11. She took responsibility for pulling together materials from all work packages to generate this report. Dr Andrea Taylor (Academic Fellow) was primarily responsible for *Chapter 6* and drafted this with LB and LS. She provided psychology input into the experiments conducted in Chapter 8 and comments on draft Chapters. Dr Abigail Colson (Lecturer, Management Science) took primary responsibility for the review and critique in Chapter 2. She contributed towards the evaluation in Chapter 11 and worked with AF and CJ on Chapter 5. Professor Alec Morton provided senior support for the work carried out in Chapter 2 and provided input on all Chapters and the final report.

1
2
3 **Chapter 15 Patient and public involvement**
4

5
6 No PPI activities were planned as part of this project. An engagement workshop took place, involving
7 NHS decision makers, analysts and methodologists. The details of this engagement activities are
8 presented in Chapter 12 and Supplementary material 5.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Chapter 16 References

1. Soares MO, Sharples L, Morton A, Claxton K, Bojke L. Experiences of Structured Elicitation for Model-Based Cost-Effectiveness Analyses. *Value Health* 2018;**21**:715-23. <http://dx.doi.org/10.1016/j.jval.2018.01.019>
2. Bryan S, Williams I, McIver S. Seeing the NICE side of cost-effectiveness analysis: a qualitative investigation of the use of CEA in NICE technology appraisals. *Health Economics* 2007;**16**:179-93.
3. Bothwell LE, Greene JA, Podolsky SH, Jones DS. Assessing the Gold Standard--Lessons from the History of RCTs. *N Engl J Med* 2016;**374**:2175-81. <http://dx.doi.org/10.1056/NEJMms1604593>
4. Chavez-MacGregor M, Giordano SH. Randomized Clinical Trials and Observational Studies: Is There a Battle? *J Clin Oncol* 2016;**34**:772-3. <http://dx.doi.org/10.1200/jco.2015.64.7487>
5. Rothwell PM. External validity of randomised controlled trials: "to whom do the results of this trial apply?". *Lancet* 2005;**365**:82-93. [http://dx.doi.org/10.1016/S0140-6736\(04\)17670-8](http://dx.doi.org/10.1016/S0140-6736(04)17670-8)
6. Frieden TR. Evidence for Health Decision Making - Beyond Randomized, Controlled Trials. *N Engl J Med* 2017;**377**:465-75. <http://dx.doi.org/10.1056/NEJMra1614394>
7. Hora SC. Aleatory and epistemic uncertainty in probability elicitation with an example from hazardous waste management. *Reliability Engineering & System Safety* 1996;**54**:217-23. [http://dx.doi.org/https://doi.org/10.1016/S0951-8320\(96\)00077-4](http://dx.doi.org/https://doi.org/10.1016/S0951-8320(96)00077-4)
8. O'Hagan A, Buck C, Daneshkhan A, Eiser J, Garthwaite P, Jenkinson D, *et al.* *Uncertain judgements: eliciting experts' probabilities*; 2006.
9. Griffin SC, Claxton KP, Palmer SJ, Sculpher MJ. Dangerous omissions: the consequences of ignoring decision uncertainty. *Health Econ* 2011;**20**:212-24. <http://dx.doi.org/10.1002/hec.1586>
10. Babuscia A, Cheung K-M. An approach to perform expert elicitation for engineering design risk analysis: methodology and experimental results. *J R Statist Soc* 2014;**177**:475-97.
11. Ayyub BM. *Elicitation of expert opinions for uncertainty and risks*. Boca Raton, Fla.: CRC Press; 2001.
12. Peel A, Jenks M, Choudhury M, Lovett R, Rejon-Parrilla JC, Sims A, *et al.* Use of Expert Judgement Across NICE Guidance-Making Programmes: A Review of Current Processes and Suitability of Existing Tools to Support the Use of Expert Elicitation. 2018;**16**:819-36. <http://dx.doi.org/10.1007/s40258-018-0415-5>
13. Soares MO, Bojke L. Expert elicitation to inform health technology assessment. In: International Series in Operations Research and Management Science; 2018: 479-94. http://dx.doi.org/10.1007/978-3-319-65052-4_18
14. Cooke RM. *Experts in Uncertainty: Opinion and Subjective Probability in Science*. Oxford University Press 1991.
15. O'Hagan T, Oakley J. The Sheffield Elicitation Framework (SHELF). 2008.
16. Colson AR, Cooke RM. Expert Elicitation: Using the Classical Model to Validate Experts' Judgments. *Review of Environmental Economics and Policy* 2018;**12**:113-32. <http://dx.doi.org/10.1093/reep/rex022>
17. European Food Safety Authority. Guidance on Expert Knowledge Elicitation in Food and Feed Safety Risk Assessment. *EFSA Journal* 2014;**12**. <http://dx.doi.org/10.2903/j.efsa.2014.3734>
18. EPA. *Expert Elicitation Task Force White Paper*. Draft. Washington D.C.: U.S. Environmental Protection Agency; 2009.
19. Kaplan S. 'Expert information' versus 'expert opinion.' Another approach to the problem of eliciting/combining/using expert knowledge in PRA. *Reliability Engineering & System Safety* 1992;**35**:61-72. [http://dx.doi.org/10.1016/0951-8320\(92\)90023-e](http://dx.doi.org/10.1016/0951-8320(92)90023-e)
20. Cooke RM, Goossens LHJ. Procedures guide for structured expert judgement in accident consequence modelling. *Radiation Protection Dosimetry* 2000;**90**:303-9.
21. Choy S, O'Leary R, Mengersen K. Elicitation by design in ecology: using expert opinion to inform priors for Bayesian statistical models. *Ecology* 2009;**90**.

- 1
2
3 22. Kotra J, Lee M, Eisenberg N, DeWispelare A. Branch Technical Position on the Use of Expert
4 Elicitation in the High-Level Radioactive Waste Program. *Division of Waste Management, Office of*
5 *Nuclear Material Safety and Safeguards, US Nuclear Regulatory Commission* 1996.
- 6 23. Tredger ERW, Lo JTH, Haria S, Lau HHK, Bonello N, Hlavka B, *et al.* Bias, guess and expert
7 judgement in actuarial work. *British Actuarial Journal* 2016;**21**:545-78.
8 <http://dx.doi.org/10.1017/s1357321716000155>
- 9 24. Budnitz RJ, Apostolakis G, Boore DM, Cluff LS, Coppersmith KJ, Cornell CA, *et al.*
10 *Recommendations for Probabilistic Seismic Hazard Analysis: Guidance on Uncertainty and Use of*
11 *Experts*: no. NUREG/CR-6372 UCRL-ID- 122160. Livermore, CA: U.S. Nuclear Regulatory Commission;
12 1997.
- 13 25. Hemming V, Burgman MA, Hanea AM, McBride MF, Wintle BC. A practical guide to structured
14 expert elicitation using the IDEA protocol. *Methods in Ecology and Evolution* 2018;**9**:169-80.
15 <http://dx.doi.org/10.1111/2041-210x.12857>
- 16 26. Garthwaite PH, Kadane JB, O'Hagan A. Statistical Methods for Eliciting Probability
17 Distributions. *Journal of the American Statistical Association* 2005;**100**:680-701.
18 <http://dx.doi.org/http://dx.doi.org/10.1198/016214505000000105>
- 19 27. Knol AB, Slottje P, van der Sluijs JP, Lebret E. The use of expert elicitation in environmental
20 health impact assessment: a seven step procedure. *Environmental Health* 2010;**9**.
21 <http://dx.doi.org/Artn> 19
22
23
24
25 10.1186/1476-069x-9-19
- 26 28. Walls L, Quigley J. Building prior distributions to support Bayesian reliability growth modelling
27 using expert judgement. *Reliability Engineering & System Safety* 2001;**74**:117-28.
28 [http://dx.doi.org/10.1016/s0951-8320\(01\)00069-2](http://dx.doi.org/10.1016/s0951-8320(01)00069-2)
- 29 29. Gosling JP. SHELF: The sheffield elicitation framework. In: International Series in Operations
30 Research and Management Science; 2018: 61-93. http://dx.doi.org/10.1007/978-3-319-65052-4_4
- 31 30. Keeney R, von Winterfeldt D. Eliciting Probabilities from Experts in Complex Technical
32 Problems. *IEEE Transactions On Engineering Management* 1991;**38**.
- 33 31. Meyer MA, Booker JM. *Eliciting and Analyzing Expert Judgment: A Practical Guide*.
34 Philadelphia, PA: Society for Industrial and Applied Mathematics; 2001.
- 35 32. Ashcroft M, Austin R, Barnes K, MacDonald D, Makin S, Morgan S, *et al.* Expert Judgement.
36 *British Actuarial Journal* 2016;**21**:314-63.
- 37 33. Budescu DV, Chen E. Identifying Expertise to Extract the Wisdom of Crowds. 2015;**61**:267-80.
38 <http://dx.doi.org/10.1287/mnsc.2014.1909>
- 39 34. Gilovich T, Griffin DW, Kahneman D, Cambridge University P. *Heuristics and biases : the*
40 *psychology of intuitive judgment*. Cambridge: Cambridge University Press; 2013.
- 41 35. Kahneman D, Slovic P, Tversky A. *Judgment under uncertainty : heuristics and biases*.
42 Cambridge; New York: Cambridge University Press; 1982.
- 43 36. Morgan MG. Use (and abuse) of expert elicitation in support of decision making for public
44 policy. *PNAS* 2014;**111**:7176-84.
- 45 37. NHS. *Guide to the Healthcare System in England Including the Statement of NHS*
46 *Accountability*; 2013.
- 47 38. Kershaw A. *NHS Vale of York CCG Referral Support Service Useful Information*. York; 2017.
- 48 39. Kay A. The abolition of the GP fundholding scheme: a lesson in evidence-based policy making.
49 *The British journal of general practice : the journal of the Royal College of General Practitioners*
50 2002;**52**:141-4.
- 51 40. Lafond S, Charlesworth A, Roberts A. *A year of plenty? An analysis of NHS finances and*
52 *consultant productivity*. London; 2017.
- 53 41. King's Fund. *Has the government delivered a new era for public health?* The King's Fund; 2015.
54 URL: <https://www.kingsfund.org.uk/projects/verdict/has-government-delivered-new-era-public-health>
55 (Accessed 28th March, 2019).
- 56
57
58
59
60

- 1
2
3 42. NHS. *Interim Commissioning Policy: Individual funding requests*. England: NHS Commissioning Board; 2013.
- 4 43. Ham C, Glenn R. Reasonable Rationing: International Experience of Priority Setting in Health
5 Care (State of Health) In: Open University press; 2003.
- 6 44. Grigore B, Peters J, Hyde C, Stein K. Methods to Elicit Probability Distributions from Experts:
7 A Systematic Review of Reported Practice in Health Technology Assessment. *PharmacoEconomics*
8 2013;**31**:991–1003.
- 9 45. NICE. *Guide to the process of technology appraisal*; 2014.
- 10 46. Bennett P, Hare A, Townshend J. Assessing the risk of vCJD transmission via surgery: models
11 for uncertainty and complexity. *Journal of the Operational Research Society* 2005;**56**:202-13.
12 <http://dx.doi.org/10.1057/palgrave.jors.2601899>
- 13 47. Colson AR, Megiddo I, Alvarez-Uria G, Gandra S, Bedford T, Morton A, *et al*. Quantifying
14 uncertainty about future antimicrobial resistance: Comparing structured expert judgment and
15 statistical forecasting methods. *PLoS One* 2019;**14**:e0219190.
16 <http://dx.doi.org/10.1371/journal.pone.0219190>
- 17 48. Dallow N, Best N, H Montague T. *Better Decision Making in Drug Development Through*
18 *Adoption of Formal Prior Elicitation*; 2017. <http://dx.doi.org/10.1002/pst.1854>
- 19 49. Walley RJ, Smith CL, Gale JD, Woodward P. Advantages of a wholly Bayesian approach to
20 assessing efficacy in early drug development: a case study. *Pharm Stat* 2015;**14**:205-15.
21 <http://dx.doi.org/10.1002/pst.1675>
- 22 50. Drummond M. *Methods for the economic evaluation of health care programmes*. Fourth
23 edition. edn. Oxford, United Kingdom ; New York, NY, USA: Oxford University Press; 2015.
- 24 51. Leal J, Wordsworth S, Legood R, Blair E. Eliciting expert opinion for economic models: an
25 applied example. *Value Health* 2007;**10**:195-203. <http://dx.doi.org/10.1111/j.1524-4733.2007.00169.x>
- 26 52. Soares MO, Bojke L, Dumville J, Iglesias C, Cullum N, Claxton K. Methods to elicit experts'
27 beliefs over uncertain quantities: application to a cost effectiveness transition model of negative
28 pressure wound therapy for severe pressure ulceration. *Stat Med* 2011;**30**:2363-80.
29 <http://dx.doi.org/10.1002/sim.4288>
- 30 53. Haakma W, Steuten LMG, Bojke L, IJzerman MJ. Belief Elicitation to Populate Health Economic
31 Models of Medical Diagnostic Devices in Development. *Appl Health Econ Health Policy* 2014;**12**:327-
32 34. <http://dx.doi.org/10.1007/s40258-014-0092-y>
- 33 54. Bojke L, Claxton K, Bravo-Vergel Y, Sculpher M, Palmer S, Abrams K. Eliciting distributions to
34 populate decision analytic models. *Value Health* 2010;**13**:557-64. <http://dx.doi.org/10.1111/j.1524-4733.2010.00709.x>
- 35 55. McKenna C, McDaid C, Suekarran S, Hawkins N, Claxton K, Light K, *et al*. Enhanced external
36 counterpulsation for the treatment of stable angina and heart failure: a systematic review and
37 economic analysis. *Health Technol Assess* 2009;**13**:iii-iv, ix-xi, 1-90.
38 <http://dx.doi.org/10.3310/hta13240>
- 39 56. Sperber D, Mortimer D, Lorgelly P, Berlowitz D. An Expert on Every Street Corner? Methods
40 for Eliciting Distributions in Geographically Dispersed Opinion Pools. *Value in Health* 2013;**16**:434-7.
- 41 57. Fischer K, Lewandowski D, Janssen MP. Estimating unknown parameters in haemophilia using
42 expert judgement elicitation. *Haemophilia* 2013;**19**:e282-8. <http://dx.doi.org/10.1111/hae.12166>
- 43 58. Garthwaite PH, Chilcott JB, Jenkinson DJ, Tappenden P. Use of expert knowledge in evaluating
44 costs and benefits of alternative service provisions: a case study. *Int J Technol Assess Health Care*
45 2008;**24**:350-7. <http://dx.doi.org/10.1017/S026646230808046X>
- 46 59. Grigore B, Peters J, Hyde C, Stein K. A comparison of two methods for expert elicitation in
47 health technology assessments. *BMC Medical Research Methodology* 2016;**16**.
- 48 60. Meads C, Auguste P, Davenport C, Malysiak S, Sundar S, Kowalska M, *et al*. Positron emission
49 tomography/computerised tomography imaging in detecting and managing recurrent cervical cancer:
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 Systematic review of evidence, elicitation of subjective probabilities and economic modeling. *Health*
4 *Technology Assessment* 2013;**17**:7-144. <http://dx.doi.org/10.3310/hta17120>
5
6 61. Wilson EC, Stanley G, Mirza Z. The Long-Term Cost to the UK NHS and Social Services of
7 Different Durations of IV Thiamine (Vitamin B1) for Chronic Alcohol Misusers with Symptoms of
8 Wernicke's Encephalopathy Presenting at the Emergency Department. *Appl Health Econ Health Policy*
9 2016;**14**:205-15. <http://dx.doi.org/10.1007/s40258-015-0214-1>
10
11 62. Brodtkorb T-H. *Cost-effectiveness analysis of health technologies when evidence is scarce*
12 Linköping University, Sweden; 2010.
13
14 63. Cao Q, Postmus D, Hillege HL, Buskens E. Probability elicitation to inform early health
15 economic evaluations of new medical technologies: a case study in heart failure disease management.
16 *Value Health* 2013;**16**:529-35. <http://dx.doi.org/10.1016/j.jval.2013.02.008>
17
18 64. Speight PM, Palmer S, Moles DR, Downer MC, Smith DH, Henriksson M, *et al*. The cost-
19 effectiveness of screening for oral cancer in primary care. *Health Technol Assess* 2006;**10**:1-144, iii-iv.
20
21 65. Poncet A, Gencer B, Blondon M, Gex-Fabry M, Combescure C, Shah D, *et al*.
22 Electrocardiographic Screening for Prolonged QT Interval to Reduce Sudden Cardiac Death in
23 Psychiatric Patients: A Cost-Effectiveness Analysis. *PLoS One* 2015;**10**:e0127213.
24 <http://dx.doi.org/10.1371/journal.pone.0127213>
25
26 66. Stevenson MD, Oakley JE, Lloyd Jones M, Brennan A, Compston JE, McCloskey EV, *et al*. The
27 cost-effectiveness of an RCT to establish whether 5 or 10 years of bisphosphonate treatment is the
28 better duration for women with a prior fracture. *Med Decis Making* 2009;**29**:678-89.
29 <http://dx.doi.org/10.1177/0272989X09336077>
30
31 67. Meeyai A, Praditsitthikorn N, Kotirum S, Kulpeng W, Putthasri W, Cooper BS, *et al*. Seasonal
32 influenza vaccination for children in Thailand: a cost-effectiveness analysis. *PLoS Med*
33 2015;**12**:e1001829; discussion e. <http://dx.doi.org/10.1371/journal.pmed.1001829>
34
35 68. Colbourn T, Asseburg C, Bojke L, Phillips Z, Claxton K, Ades AE, *et al*. Prenatal screening and
36 treatment strategies to prevent group B streptococcal and other bacterial infections in early infancy:
37 cost-effectiveness and expected value of information analyses. *Health Technol Assess* 2007;**11**:1-226,
38 iii.
39
40 69. Girling AJ, Freeman G, Gordon JP, Poole-Wilson P, Scott DA, Ford RJ. Modeling payback from
41 research into the efficacy of left-ventricular assist devices as destination therapy. *International Journal*
42 *of Technology Assessment in Health Care* 2007;**23**:269-77.
43 <http://dx.doi.org/10.1017/S0266462307070365>
44
45 70. Stevenson MD, Oakley JE, Chick SE, Chalkidou K. The cost-effectiveness of surgical instrument
46 management policies to reduce the risk of vCJD transmission to humans. *Journal of the Operational*
47 *Research Society* 2009;**60**:506-18. <http://dx.doi.org/10.1057/palgrave.jors.2602580>
48
49 71. De Persis C, Wilson S. Using the analytic hierarchy process in the assessment of the probability
50 for an explosion to occur during the atmospheric re-entry. Proceedings of the International
51 Astronautical Congress, IAC, abstract no. 1021.
52
53 72. Iglesias CP, Thompson A, Rogowski WH, Payne K. Reporting Guidelines for the Use of Expert
54 Judgement in Model-Based Economic Evaluations. *Pharmacoeconomics* 2016;**34**:1161-72.
55 <http://dx.doi.org/10.1007/s40273-016-0425-9>
56
57 73. Colson AR, Cooke RM. Cross validation for the classical model of structured expert judgment.
58 *Reliability Engineering & System Safety* 2017;**163**:109-20.
59 <http://dx.doi.org/10.1016/j.res.2017.02.003>
60
61 74. Eggstaff JW, Mazzuchi TA, Sarkani S. The effect of the number of seed variables on the
62 performance of Cooke's classical model. *Reliability Engineering & System Safety* 2014;**121**:72-82.
63 <http://dx.doi.org/10.1016/j.res.2013.07.015>
64
65 75. Clemen RT. Comment on Cooke's classical method. *Reliability Engineering & System Safety*
66 2008;**93**:760-5. <http://dx.doi.org/10.1016/j.res.2008.02.003>
67
68 76. Montibeller G, von Winterfeldt D. Cognitive and Motivational Biases in Decision and Risk
69 Analysis. *Risk Anal* 2015;**35**:1230-51. <http://dx.doi.org/10.1111/risa.12360>
70

- 1
- 2
- 3
- 4 77. Philips Z, Ginnelly L, Sculpher M, Claxton K, Golder S, Riemsma R, *et al.* Review of guidelines
- 5 for good practice in decision-analytic modelling in health technology assessment. *Health Technol*
- 6 *Assess* 2004;**8**:iii-iv, ix-xi, 1-158.
- 7 78. Anthony O'Hagan, Caitlin E. Buck, Alireza Daneshkhah, J. Richard Eiser, Paul H. Garthwaite,
- 8 David J. Jenkinson, *et al.* Uncertain Judgements: Eliciting Experts' Probabilities. *Wiley* 2006:338.
- 9 79. Bolger F. The selection of experts for (probabilistic) expert knowledge elicitation. In:
- 10 International Series in Operations Research and Management Science; 2018: 393-443.
- 11 http://dx.doi.org/10.1007/978-3-319-65052-4_16
- 12 80. Burgman MA, McBride M, Ashton R, Speirs-Bridge A, Flander L, Wintle B, *et al.* Expert Status
- 13 and Performance. *PLOS ONE* 2011;**6**:e22998. <http://dx.doi.org/10.1371/journal.pone.0022998>
- 14 81. Clemen RT, Winkler RL. Combining probability distributions from experts in risk analysis. *Risk*
- 15 *Analysis* 1999;**19**:187-203. <http://dx.doi.org/Doi> 10.1023/A:1006917509560
- 16 82. Fogel L. *Human Information Processing*: Prentice-Hall; 1967.
- 17 83. Seaver D. *Assessment of Group Preferences and Group Uncertainty for Decision-Making*.
- 18 California: Social Science Research Institute; 1976.
- 19 84. Staël von Holstein C-AS. Two techniques for assessment of subjective probability distributions
- 20 — An experimental study. *Acta Psychologica* 1971;**35**:478-94.
- 21 [http://dx.doi.org/https://doi.org/10.1016/0001-6918\(71\)90005-9](http://dx.doi.org/https://doi.org/10.1016/0001-6918(71)90005-9)
- 22 85. Thall PF, Ursino M, Baudouin V, Alberti C, Zohar S. Bayesian treatment comparison using
- 23 parametric mixture priors computed from elicited histograms. 2019;**28**:404-18.
- 24 <http://dx.doi.org/10.1177/0962280217726803>
- 25 86. Bornkamp B, Ickstadt K. A Note on B-Splines for Semiparametric Elicitation. *The American*
- 26 *Statistician* 2009;**63**:373-7.
- 27 87. Moala FA, O'Hagan A. Elicitation of multivariate prior distributions: A nonparametric Bayesian
- 28 approach. *Journal of Statistical Planning and Inference* 2010;**140**:1635-55.
- 29 <http://dx.doi.org/https://doi.org/10.1016/j.jspi.2010.01.004>
- 30 88. Daneshkhah A, Hosseinian-Far A, Sedighi T, Farsi M. Prior elicitation and evaluation of
- 31 imprecise judgements for bayesian analysis of system reliability. In: *Strategic Engineering for Cloud*
- 32 *Computing and Big Data Analytics*; 2017:63-79. http://dx.doi.org/10.1007/978-3-319-52491-7_4
- 33 89. Lindley DV, Tversky A, Brown RV. On the Reconciliation of Probability Assessments. *Journal of*
- 34 *the Royal Statistical Society: Series A (General)* 1979;**142**:146-62. <http://dx.doi.org/10.2307/2345078>
- 35 90. Wittmann ME, Cooke RM, Rothlisberger JD, Lodge DM. Using Structured Expert Judgment to
- 36 Assess Invasive Species Prevention: Asian Carp and the Mississippi—Great Lakes Hydrologic
- 37 Connection. *Environmental Science & Technology* 2014;**48**:2150-6.
- 38 <http://dx.doi.org/10.1021/es4043098>
- 39 91. Quigley J, Walls L. A Methodology for Constructing Subjective Probability Distributions with
- 40 Data. In: Dias LC, Morton A, Quigley J, editors. *Elicitation: The science and art of structuring*
- 41 *judgement*New York, NY: Springer; 2018:141-70.
- 42 92. Meyer M, Booker J. *Eliciting and Analyzing Expert Judgement: A Practical Guide*: Academic
- 43 Press; 1991.
- 44 93. Boring R, Gertman D, Joe J, Marble J, Galyean W, Blackwood L, *et al.* Simplified Expert
- 45 Elicitation Guideline For Risk Assessment Of Operating Events. *US Nuclear Regulatory Commission*
- 46 *(NRC)* 2005.
- 47 94. Bolger F, Rowe G. The Aggregation of Expert Judgment: Do Good Things Come to Those Who
- 48 Weight? *Risk Analysis* 2015;**35**. <http://dx.doi.org/10.1111/risa.12272>
- 49 95. Claxton K, Sculpher M, McCabe C, Briggs A, Akehurst R, Buxton M, *et al.* Probabilistic sensitivity
- 50 analysis for NICE technology assessment: not an optional extra. *Health Econ* 2005;**14**:339-47.
- 51 <http://dx.doi.org/10.1002/hec.985>
- 52 96. Akins RB, Tolson H, Cole BR. Stability of response characteristics of a Delphi panel: application
- 53 of bootstrap data expansion. *BMC medical research methodology* 2005;**5**:37-.
- 54 <http://dx.doi.org/10.1186/1471-2288-5-37>
- 55
- 56
- 57
- 58
- 59
- 60

97. Gosling JP. Methods for eliciting expert opinion to inform health technology assessment. 2014.
98. Winkler RL. The assessment of Prior Distributions in Bayesian Analysis. *Journal of the American Statistical Association* 1967;**62**:776-8. <http://dx.doi.org/Doi.10.2307/2283671>
99. O'Hagan T, Oakley JE. SHELF: the Sheffield Elicitation Framework (version 3.0). In: UK: School of Mathematics and Statistics, University of Sheffield; 2016.
100. Gosling JP. *On the elicitation of continuous, symmetric, unimodal distributions* no. arXiv:0805.2044; 2008.
101. Oakley JE, O'Hagan A. Uncertainty in Prior Elicitations: A Nonparametric Approach. *Biometrika* 2007;**94**:427-41.
102. Gosling JP, Oakley JE, O'Hagan A. Nonparametric elicitation for heavy-tailed prior distributions. *Bayesian Analysis* 2007;**2**:693-718. <http://dx.doi.org/10.1214/07-ba228>
103. Morris P. Decision analysis expert use. *Management Science* 1974;**20**:1233-41.
104. Morris PA. Combining expert judgments: A bayesian approach. *Management Science* 1977;**23**.
105. Jacobs RA. Methods for combining experts' probability assessments. *Neural Comput* 1995;**7**:867-88.
106. Lipscomb J, Parmigiani G, Hasselblad V. Combining Expert Judgment by Hierarchical Modeling: An Application to Physician Staffing. *Management Science* 1998;**44**:149-61.
107. Albert I, Donnet S, Guihenneuc-Jouyaux C, Lowchoy S, L. Mengersen K, Rousseau J. *Combining Expert Opinions in Prior Elicitation*; 2012. <http://dx.doi.org/10.1214/12-BA717>
108. West M, Crosse J. *Modelling probabilistic agent opinion*; 1992. <http://dx.doi.org/10.2307/2345964>
109. Gelfand AE, Mallick BK, Dey DK. Modeling Expert Opinion Arising as a Partial Probabilistic Specification. *Journal of the American Statistical Association* 1995;**90**:598-604. <http://dx.doi.org/10.1080/01621459.1995.10476552>
110. Lichtendahl KC, Grushka-Cockayne Y, Winkler RL. Is It Better to Average Probabilities or Quantiles? *Management Science* 2013;**59**:1594-611. <http://dx.doi.org/10.1287/mnsc.1120.1667>
111. Bamber JL, Aspinall WP, Cooke RM. A commentary on "how to interpret expert judgment assessments of twenty-first century sea-level rise" by Hylke de Vries and Roderik SW van de Wal. *Climatic Change* 2016;**137**:321-8. <http://dx.doi.org/10.1007/s10584-016-1672-7>
112. French S. Group consensus probability distributions: A critical survey. In: Bernardo JM, editor. *Bayesian Statistics 2* Amsterdam: North-Holland.; 1985:183-97.
113. Hammitt JK, Zhang YF. Combining Experts' Judgments: Comparison of Algorithmic Methods Using Synthetic Data. *Risk Analysis* 2013;**33**:109-20. <http://dx.doi.org/10.1111/j.1539-6924.2012.01833.x>
114. Aspinall WP, Cooke RM. Quantifying scientific uncertainty from expert judgement elicitation. *Risk and Uncertainty Assessment for Natural Hazards* 2013.
115. Cooke RM, ElSaadany S, Huang X. On the performance of social network and likelihood-based expert weighting schemes. *Reliability Engineering & System Safety* 2008;**93**:745-56. <http://dx.doi.org/10.1016/j.ress.2007.03.017>
116. Ranjan R, Gneiting T. Combining probability forecasts. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 2010;**72**:71-91. <http://dx.doi.org/10.1111/j.1467-9868.2009.00726.x>
117. Rufo MJ, Martin J, Perez CJ. Log-Linear Pool to Combine Prior Distributions: A Suggestion for a Calibration-Based Approach. *Bayesian Anal* 2012;**7**:411-38. <http://dx.doi.org/10.1214/12-BA714>
118. Hora SC, Kardeş E. Calibration, sharpness and the weighting of experts in a linear opinion pool. *Annals of Operations Research* 2015;**229**:429-50. <http://dx.doi.org/10.1007/s10479-015-1846-0>
119. Winkler RL, Murphy AH. "Good" Probability Assessors. *Journal of Applied Meteorology* 1968;**7**:751-8. [http://dx.doi.org/10.1175/1520-0450\(1968\)007<0751:PA>2.0.CO;2](http://dx.doi.org/10.1175/1520-0450(1968)007<0751:PA>2.0.CO;2)

- 1
2
3 120. Quigley J, Colson A, Aspinall W, Cooke RM. Elicitation in the classical model. In: International
4 Series in Operations Research and Management Science; 2018: 15-36. [http://dx.doi.org/10.1007/978-](http://dx.doi.org/10.1007/978-3-319-65052-4_2)
5 [3-319-65052-4_2](http://dx.doi.org/10.1007/978-3-319-65052-4_2)
- 6 121. R C, L G. TU Delft expert judgment data base. *Reliability Engineering and System Safety*
7 2008;**93**:657-74.
- 8 122. Cooke RM. Validating Expert Judgment with the Classical Model. In: Martini C, Boumans M,
9 editors. *Experts and Consensus in Social Science* Cham: Springer International Publishing; 2014:191-
10 212.
- 11 123. Mellers B, Stone E, Murray T, Minster A, Rohrbaugh N, Bishop M, *et al.* Identifying and
12 Cultivating Superforecasters as a Method of Improving Probabilistic Predictions. *Perspectives on*
13 *Psychological Science* 2015;**10**:267-81. <http://dx.doi.org/10.1177/1745691615577794>
- 14 124. Hanea AM, McBride MF, Burgman MA, Wintle BC. The Value of Performance Weights and
15 Discussion in Aggregated Expert Judgments. *Risk Analysis* 2018;**0**.
16 <http://dx.doi.org/10.1111/risa.12992>
- 17 125. Morgenstern O, Von Neumann J. *Theory of games and economic behavior*: Princeton
18 university press; 1953.
- 19 126. Kahneman D, Egan P. *Thinking, fast and slow*: Farrar, Straus and Giroux New York; 2011.
- 20 127. Reyna VF, Nelson WL, Han PK, Dieckmann NF. How numeracy influences risk comprehension
21 and medical decision making. *Psychological bulletin* 2009;**135**:943.
- 22 128. Tversky A, Kahneman D. Judgment under uncertainty: Heuristics and biases. *science*
23 1974;**185**:1124-31.
- 24 129. Gigerenzer G, Selten R. *Bounded rationality: The adaptive toolbox*: MIT press; 2002.
- 25 130. Kynn M. The 'heuristics and biases' bias in expert elicitation. *Journal of the Royal Statistical*
26 *Society Series a-Statistics in Society* 2008;**171**:239-64.
- 27 131. Bojke L, Grigore B, Jankovic D, Peters J, Soares M, Stein K. Informing Reimbursement Decisions
28 Using Cost-Effectiveness Modelling: A Guide to the Process of Generating Elicited Priors to Capture
29 Model Uncertainties. *Pharmacoeconomics* 2017;**35**:867-77. [http://dx.doi.org/10.1007/s40273-017-](http://dx.doi.org/10.1007/s40273-017-0525-1)
30 [0525-1](http://dx.doi.org/10.1007/s40273-017-0525-1)
- 31 132. Bazerman MH, Moore DA. Judgment in managerial decision making. 2008.
- 32 133. McBride MF, Fidler F, Burgman MA. Evaluating the accuracy and calibration of expert
33 predictions under uncertainty: predicting the outcomes of ecological research. *Diversity and*
34 *Distributions* 2012;**18**:782-94. <http://dx.doi.org/10.1111/j.1472-4642.2012.00884.x>
- 35 134. Tversky A, Kahneman D. Availability: A heuristic for judging frequency and probability.
36 *Cognitive psychology* 1973;**5**:207-32.
- 37 135. Slovic P, Fischhoff B, Lichtenstein S. Perceived risk: psychological factors and social
38 implications. *Proc R Soc Lond A* 1981;**376**:17-34.
- 39 136. Mehle T, Gettys CF, Manning C, Baca S, Fisher S. The availability explanation of excessive
40 plausibility assessments. *Acta Psychologica* 1981;**49**:127-40.
- 41 137. McBride MF, Garnett ST, Szabo JK, Burbidge AH, Butchart SH, Christidis L, *et al.* Structured
42 elicitation of expert judgments for threatened species assessment: a case study on a continental scale
43 using email. *Methods in Ecology and Evolution* 2012;**3**:906-20.
- 44 138. Soll JB, Klayman J. Overconfidence in interval estimates. *Journal of Experimental Psychology:*
45 *Learning, Memory, and Cognition* 2004;**30**:299.
- 46 139. McKenzie CR, Liersch MJ, Yaniv I. Overconfidence in interval estimates: What does expertise
47 buy you? *Organizational behavior and human decision processes* 2008;**107**:179-91.
- 48 140. Larrick RP. Debiasing. *Blackwell handbook of judgment and decision making* 2004:316-38.
- 49 141. Soll J, Milkman K, Payne J. A user's guide to debiasing. 2014.
- 50 142. Clemen RT, Lichtendahl KC. Debiasing expert overconfidence: A Bayesian calibration model.
51 Sixth International Conference on Probabilistic Safety Assessment and Management (PSAM6), abstract
52 no. 1369.
- 53
54
55
56
57
58
59
60

- 1
2
3 143. Cooke R, Shrader-Frechette K. *Experts in uncertainty: opinion and subjective probability in*
4 *science*: Oxford University Press on Demand; 1991.
- 5 144. Lin S-W, Bier VM. A study of expert overconfidence. *Reliability Engineering & System Safety*
6 2008;**93**:711-21.
- 7 145. Bolger F, Rowe G. There is data, and then there is data: Only experimental evidence will
8 determine the utility of differential weighting of expert judgment. *Risk Analysis* 2015;**35**:21-6.
- 9 146. Haran U, Ritov I, Mellers BA. The role of actively open-minded thinking in information
10 acquisition, accuracy, and calibration. *Judgment and Decision Making* 2013;**8**:188.
- 11 147. Plous S. A comparison of strategies for reducing interval overconfidence in group judgments.
12 *Journal of Applied Psychology* 1995;**80**:443.
- 13 148. Haran U, Moore DA, Morewedge CK. A simple remedy for overprecision in judgment.
14 *Judgment and Decision Making* 2010;**5**:467.
- 15 149. Speirs-Bridge A, Fidler F, McBride M, Flander L, Cumming G, Burgman M. Reducing
16 overconfidence in the interval judgments of experts. *Risk Analysis* 2010;**30**:512-23.
- 17 150. Teigen KH, Jørgensen M. When 90% confidence intervals are 50% certain: On the credibility
18 of credible intervals. *Applied Cognitive Psychology* 2005;**19**:455-75.
- 19 151. Winman A, Hansson P, Juslin P. Subjective probability intervals: how to reduce overconfidence
20 by interval evaluation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*
21 2004;**30**:1167.
- 22 152. Ferretti V, Guney S, Montibeller G, von Winterfeldt D. Testing best practices to reduce the
23 overconfidence bias in multi-criteria decision analysis. System Sciences (HICSS), 2016 49th Hawaii
24 International Conference on, abstract no. 1381, p. 1547-55.
- 25 153. Murphy AH, Winkler RL. Probability forecasts: A survey of National Weather Service
26 forecasters. *Bulletin of the American Meteorological Society* 1974;**55**:1449-52.
- 27 154. Martin TG, Burgman MA, Fidler F, Kuhnert PM, Low-Choy S, McBride M, *et al.* Eliciting Expert
28 Knowledge in Conservation Science. *Conservation Biology* 2012;**26**:29-38.
29 <http://dx.doi.org/10.1111/j.1523-1739.2011.01806.x>
- 30 155. Prava VR, Clemen RT, Hobbs BF, Kenney MA. Partition Dependence and Carryover Biases in
31 Subjective Probability Assessment Surveys for Continuous Variables: Model-Based Estimation and
32 Correction. *Decision Analysis* 2016;**13**:51-67. <http://dx.doi.org/10.1287/deca.2015.0323>
- 33 156. Block RA, Harper DR. Overconfidence in estimation: Testing the anchoring-and-adjustment
34 hypothesis. *Organizational behavior and human decision processes* 1991;**49**:188-207.
- 35 157. Schall DL, Doll D, Mohnen A. Caution! Warnings as a Useless Countermeasure to Reduce
36 Overconfidence? An Experimental Evaluation in Light of Enhanced and Dynamic Warning Designs.
37 *Journal of Behavioral Decision Making* 2017;**30**:347-58.
- 38 158. Arkes HR. Costs and benefits of judgment errors: Implications for debiasing. *Psychological*
39 *bulletin* 1991;**110**:486.
- 40 159. Welsh MB, Begg SH, Bratvold RB. Efficacy of bias awareness in debiasing oil and gas judgments.
41 Proceedings of the Annual Meeting of the Cognitive Science Society, abstract no. 1352.
- 42 160. Morewedge CK, Yoon H, Scopelliti I, Symborski CW, Korris JH, Kassam KS. Debiasing decisions:
43 Improved decision making with a single training intervention. *Policy Insights from the Behavioral and*
44 *Brain Sciences* 2015;**2**:129-40.
- 45 161. Snyder M, Swann WB. Hypothesis-testing processes in social interaction. *Journal of*
46 *Personality and social psychology* 1978;**36**:1202.
- 47 162. Nisbett RE, Ross L. Human inference: Strategies and shortcomings of social judgment. 1980.
- 48 163. Downs JS, Shafir E. Why some are perceived as more confident and more insecure, more
49 reckless and more cautious, more trusting and more suspicious, than others: Enriched and
50 impoverished options in social judgment. *Psychonomic bulletin & review* 1999;**6**:598-610.
- 51 164. Abbas AE, Budescu DV, Yu H-T, Haggerty R. A comparison of two probability encoding
52 methods: Fixed probability vs. fixed variable values. *Decision Analysis* 2008;**5**:190-202.
- 53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
165. Nemet GF, Anadon LD, Verdolini E. Quantifying the effects of expert selection and elicitation design on experts' confidence in their judgments about future energy technologies. *Risk Analysis* 2017;**37**:315-30.
166. Briggs A, Claxton K, Sculpher M. *Decision modelling for health economic evaluation*. Oxford: Oxford University Press; 2006.
167. Brennan A, Chick SE, Davies R. A taxonomy of model structures for economic evaluation of health technologies. *Health Econ* 2006;**15**:1295-310. <http://dx.doi.org/10.1002/hec.1148>
168. Cao Q, Buskens E, Feenstra T, Jaarsma T, Hillege H, Postmus D. Continuous-Time Semi-Markov Models in Health Economic Decision Making: An Illustrative Example in Heart Failure Disease Management. *Med Decis Making* 2016;**36**:59-71. <http://dx.doi.org/10.1177/0272989x15593080>
169. Karnon J, Stahl J, Brennan A, Caro JJ, Mar J, Moller J. Modeling using discrete event simulation: a report of the ISPOR-SMDM Modeling Good Research Practices Task Force--4. *Value Health* 2012;**15**:821-7. <http://dx.doi.org/10.1016/j.jval.2012.04.013>
170. Davis S, Stevenson M, Tappenden P, Wailoo A. NICE Decision Support Unit Technical Support Documents. In: *NICE DSU Technical Support Document 15: Cost-Effectiveness Modelling Using Patient-Level Simulation* London: National Institute for Health and Care Excellence (NICE)
- unless otherwise stated. All rights reserved.; 2014.
171. Collett D. *Modelling survival data in medical research*: CRC Press; 2015.
172. Putter H, Fiocco M, Geskus RB. Tutorial in biostatistics: competing risks and multi-state models. 2007;**26**:2389-430. <http://dx.doi.org/10.1002/sim.2712>
173. Welton NJ, Ades AE. Estimation of Markov Chain Transition Probabilities and Rates from Fully and Partially Observed Data: Uncertainty Propagation, Evidence Synthesis, and Model Calibration. 2005;**25**:633-45. <http://dx.doi.org/10.1177/0272989x05282637>
174. Sharples LD, Taylor GI, Faddy M. A piecewise-homogeneous Markov chain process of lung transplantation. *J Epidemiol Biostat* 2001;**6**:349-55.
175. Brard C, Le Teuff G, Le Deley MC, Hampson LV. Bayesian survival analysis in clinical trials: What methods are used in practice? *Clin Trials* 2017;**14**:78-87. <http://dx.doi.org/10.1177/1740774516673362>
176. Miksad RA, Gonen M, Lynch TJ, Roberts TG, Jr. Interpreting trial results in light of conflicting evidence: a Bayesian analysis of adjuvant chemotherapy for non-small-cell lung cancer. *J Clin Oncol* 2009;**27**:2245-52. <http://dx.doi.org/10.1200/jco.2008.16.2586>
177. Johnson SR, Tomlinson GA, Hawker GA, Granton JT, Feldman BM. Methods to elicit beliefs for Bayesian priors: a systematic review. *Journal of Clinical Epidemiology* 2010;**63**:355-69. <http://dx.doi.org/10.1016/j.jclinepi.2009.06.003>
178. Hutton JL, Owens RG. Bayesian Sample Size Calculations and Prior Beliefs About Child Sexual Abuse. *Journal of the Royal Statistical Society Series D (The Statistician)* 1993;**42**:399-404. <http://dx.doi.org/10.2307/2348473>
179. Johnson NP, Fisher RA, Brauholtz DA, Gillett WR, Lilford RJ. Survey of Australasian clinicians' prior beliefs concerning lipiodol flushing as a treatment for infertility: a Bayesian study. *Aust N Z J Obstet Gynaecol* 2006;**46**:298-304. <http://dx.doi.org/10.1111/j.1479-828X.2006.00596.x>
180. Lilford R. Formal measurement of clinical uncertainty: prelude to a trial in perinatal medicine. The Fetal Compromise Group. *BMJ (Clinical research ed)* 1994;**308**:111-2.
181. Wilson ECF, Usher-Smith JA, Emery J, Corrie PG, Walter FM. Expert Elicitation of Multinomial Probabilities for Decision-Analytic Modeling: An Application to Rates of Disease Progression in Undiagnosed and Untreated Melanoma. *Value in Health* 2018;**21**:669-76. <http://dx.doi.org/https://doi.org/10.1016/j.jval.2017.10.009>
182. Vargas C, Bilbeny N, Balmaceda C, Rodríguez MF, Zitko P, Rojas R, et al. Costs and consequences of chronic pain due to musculoskeletal disorders from a health system perspective in Chile. *Pain reports* 2018;**3**:e656-e. <http://dx.doi.org/10.1097/PR9.0000000000000656>
183. Ren S, Oakley JE. Assurance calculations for planning clinical trials with time-to-event outcomes. *Statistics in medicine* 2014;**33**:31-45. <http://dx.doi.org/10.1002/sim.5916>

- 1
2
3 184. Chaloner K, Rhamer FS. Quantifying and documenting prior beliefs in clinical trials. *Statistics In*
4 *Medicine* 2001;**20**:581-600.
- 5 185. Chaloner K, Church T, Louis TA, Matts JP. Graphical Elicitation of a Prior Distribution for a
6 Clinical Trial. *Journal of the Royal Statistical Society Series D (The Statistician)* 1993;**42**:341-53.
7 <http://dx.doi.org/10.2307/2348469>
- 8 186. Freedman LS, Spiegelhalter DJ. The Assessment of the Subjective Opinion and its Use in
9 Relation to Stopping Rules for Clinical Trials. *Journal of the Royal Statistical Society Series D (The*
10 *Statistician)* 1983;**32**:153-60. <http://dx.doi.org/10.2307/2987606>
- 11 187. Spiegelhalter DJ, Freedman LS, Parmar MK. Applying Bayesian ideas in drug development and
12 clinical trials. *Stat Med* 1993;**12**:1501-11; discussion 13-7.
- 13 188. Parmar MKB, Spiegelhalter DJ, Freedman LS, Committee CS. The chart trials: Bayesian design
14 and monitoring in practice. 1994;**13**:1297-312. <http://dx.doi.org/10.1002/sim.4780131304>
- 15 189. Parmar MK, Griffiths GO, Spiegelhalter DJ, Souhami RL, Altman DG, van der Scheuren E.
16 Monitoring of large randomised clinical trials: a new approach with Bayesian methods. *Lancet*
17 2001;**358**:375-81. [http://dx.doi.org/10.1016/s0140-6736\(01\)05558-1](http://dx.doi.org/10.1016/s0140-6736(01)05558-1)
- 18 190. White IR, Pocock SJ, Wang D. Eliciting and using expert opinions about influence of patient
19 characteristics on treatment effects: a Bayesian analysis of the CHARM trials. *Statistics in Medicine*
20 2005;**24**:3805-21. <http://dx.doi.org/10.1002/sim.2420>
- 21 191. Singpurwalla ND. An Interactive PC-Based Procedure for Reliability Assessment Incorporating
22 Expert Opinion and Survival Data. *Journal of the American Statistical Association* 1988;**83**:43-51.
23 <http://dx.doi.org/10.2307/2288917>
- 24 192. Claxton KP, Sculpher MJ. Using value of information analysis to prioritise health research:
25 some lessons from recent UK experience. *Pharmacoeconomics* 2006;**24**:1055-68.
26 <http://dx.doi.org/10.2165/00019053-200624110-00003>
- 27 193. Cokely ET, Galesic M, Schulz E, Ghazal S, Garcia-Retamero R. Measuring risk literacy: The Berlin
28 Numeracy Test. *Judgment and Decision Making* 2012;**7**:25-47.
- 29 194. Scott SG, Bruce RA. Decision-Making Style: The Development and Assessment of a New
30 Measure. 1995;**55**:818-31. <http://dx.doi.org/10.1177/00131644950550050017>
- 31 195. Winston C, Cheng J, Allaire J, Xie Y. Shiny: Web Application Framework for R. R package version
32 1.2.0. In; 2018.
- 33 196. Rowe G, Wright G. The Delphi technique as a forecasting tool: issues and analysis.
34 *International Journal of Forecasting* 1999;**15**:353-75.
35 [http://dx.doi.org/https://doi.org/10.1016/S0169-2070\(99\)00018-7](http://dx.doi.org/https://doi.org/10.1016/S0169-2070(99)00018-7)
- 36 197. Tetlock P, Gardner D. *Superforecasting: The Art and Science of Prediction*. London: Random
37 House; 2016.
- 38 198. Cooke R. Elicitation in the Classical Model. In: Dias C, Morton A, Quigley J, editors. *Elicitation.*
39 *The Science and Art of Structuring Judgement.*: Springer International Publishing; 2017.
- 40 199. Spiegelhalter DJ, Myles JP, Jones DR, Abrams KR. Bayesian methods in health technology
41 assessment: a review. *Health Technol Assess* 2000;**4**:1-130.
- 42 200. Harnan SE, Tappenden P, Essat M, Gomersall T, Jon Minton, Wong R, *et al*. Measurement of
43 exhaled nitric oxide concentration in asthma: a systematic review and economic evaluation of NIOX
44 MINO, NIOX VERO and NObreath. *Health Technology Assessment* 2015;**19**.
- 45 201. Soares M, Claxton K, Schulpher M. *Health opportunity costs in the NHS: assessing the*
46 *implications of uncertainty using elicitation methods with experts*. Universities of Sheffield and York;
47 2017.
- 48 202. Revie M, Bedford T, Walls L. Evaluation of elicitation methods to quantify Bayes linear models.
49 *Proceedings of the Institution of Mechanical Engineers Part O-Journal of Risk and Reliability*
50 2010;**224**:322-32. <http://dx.doi.org/10.1243/1748006xjrr304>
- 51 203. Garthwaite PH, Al-Awadhi SA, Elfadaly FG, Jenkinson DJ. Prior distribution elicitation for
52 generalized linear and piecewise-linear models. *Journal of Applied Statistics* 2013;**40**:59-75.
53 <http://dx.doi.org/10.1080/02664763.2012.734794>
- 54
55
56
57
58
59
60

- 1
2
3 204. Dalal S, Khodyakov D, Srinivasan R, Straus S, Adams J. ExpertLens: A system for eliciting
4 opinions from a large pool of non-located experts with diverse knowledge. *Technological*
5 *Forecasting and Social Change* 2011;**78**:1426-44. <http://dx.doi.org/10.1016/j.techfore.2011.03.021>
6
7 205. James A, Low C, Mengersen K. Elicitor : an expert elicitation tool for regression in ecology.
8 *Environmental Modelling and Software* 2010;**25**:129-45.
9
10 206. Bolger F, Wright G. Use of expert knowledge to anticipate the future: Issues, analysis and
11 directions. *International Journal of Forecasting* 2017;**33**:230-43.
12 <http://dx.doi.org/10.1016/j.ijforecast.2016.11.001>
13
14 207. Bedford T, Quigley J, Walls L. Expert elicitation for reliable system design. *Statistical Science*
15 2006;**21**:428-50. <http://dx.doi.org/10.1214/088342306000000510>
16
17 208. Krueger T, Page T, Hubacek K, Smith L, Hiscock K. The role of expert opinion in environmental
18 modelling. *Environmental Modelling & Software* 2012;**36**:4-18.
19 <http://dx.doi.org/10.1016/j.envsoft.2012.01.011>
20
21 209. Wisniewski A, Bijak J, Christiansen S, Forster JJ, Keilman N, Raymer J, *et al.* Utilising Expert
22 Opinion to Improve the Measurement of International Migration in Europe. *Journal of Official*
23 *Statistics* 2013;**29**:583-607. <http://dx.doi.org/10.2478/jos-2013-0041>
24
25 210. Mason AJ, Gomes M, Grieve R, Ulug P, Powell JT, Carpenter J. Development of a practical
26 approach to expert elicitation for randomised controlled trials with missing health outcomes:
27 Application to the IMPROVE trial. *Clinical Trials* 2017;**14**:357-67.
28 <http://dx.doi.org/10.1177/1740774517711442>
29
30 211. Hanea A, McBride M, Burgman M, Wintle B. Classical meets modern in the IDEA protocol for
31 structured expert judgement. *Journal of Risk Research* 2018;**21**:417-33.
32
33 212. Hanea AM, McBride MF, Burgman MA, Wintle BC. Classical meets modern in the IDEA protocol
34 for structured expert judgement. *Journal of Risk Research* 2016; 10.1080/13669877.2016.1215346:1-
35 17. <http://dx.doi.org/10.1080/13669877.2016.1215346>
36
37 213. Hanea AM, McBride MF, Burgman MA, Wintle BC, Fidler F, Flander L, *et al.*
38 InvestigateDiscussEstimateAggregate for structured expert judgement. *International Journal of*
39 *Forecasting* 2017;**33**:267-79. <http://dx.doi.org/10.1016/j.ijforecast.2016.02.008>
40
41 214. Drescher M, Perera AH, Johnson CJ, Buse LJ, Drew CA, Burgman MA. Toward rigorous use of
42 expert knowledge in ecological research. *Ecosphere* 2013;**4**. <http://dx.doi.org/10.1890/es12-00415.1>
43
44 215. Kuhnert PM, Martin TG, Griffiths SP. A guide to eliciting and using expert knowledge in
45 Bayesian ecological models. *Ecology Letters* 2010;**13**:900-14. [http://dx.doi.org/10.1111/j.1461-](http://dx.doi.org/10.1111/j.1461-0248.2010.01477.x)
46 [0248.2010.01477.x](http://dx.doi.org/10.1111/j.1461-0248.2010.01477.x)
47
48 216. Ortiz NR, Wheeler TA, Breeding RJ, Hora S, Meyer MA, Keeney RL. Use of expert judgment in
49 NUREG-1150. *Nuclear Engineering and Design* 1991;**126**:313-31. [http://dx.doi.org/10.1016/0029-](http://dx.doi.org/10.1016/0029-5493(91)90023-b)
50 [5493\(91\)90023-b](http://dx.doi.org/10.1016/0029-5493(91)90023-b)
51
52 217. Bonano EJ, Hora SC, Keeney RL, Winterfeldt Dv. *Elicitation and Use of Expert Judgment in*
53 *Performance Assessment for High-Level Radioactive Waste Repositories* no. NUREG/CR-541 1,
54 SAND89-1821. Washington D.C.: U.S. Nuclear Regulatory Commission; 1990.
55
56 218. Garthwaite PH, O'Hagan A. Quantifying expert opinion in the UK water industry: an
57 experimental study. *Journal of the Royal Statistical Society Series D-the Statistician* 2000;**49**:455-77.
58 <http://dx.doi.org/10.1111/1467-9884.00246>
59
60 219. Hanea AM, Burgman M, Hemming V. IDEA for uncertainty quantification. In: *International*
61 *Series in Operations Research and Management Science*; 2018: 95-117.
62 http://dx.doi.org/10.1007/978-3-319-65052-4_5
63
64 220. Seaver DA, von Winterfeldt D, Edwards W. Eliciting subjective probability distributions on
65 continuous variables. *Organizational Behavior and Human Performance* 1978;**21**:379-91.
66 [http://dx.doi.org/https://doi.org/10.1016/0030-5073\(78\)90061-2](http://dx.doi.org/https://doi.org/10.1016/0030-5073(78)90061-2)
67
68 221. Murphy AH, Winkler RL. Credible Interval Temperature Forecasting: Some Experimental
69 Results. 1974;**102**:784-94. [http://dx.doi.org/10.1175/1520-0493\(1974\)102<0784:Citfse>2.0.Co;2](http://dx.doi.org/10.1175/1520-0493(1974)102<0784:Citfse>2.0.Co;2)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Appendix

Search details for the review of existing published guidelines for elicitation

We used the following inclusion criteria to identify elicitation guidelines:

1. Guidelines must be full-length (i.e., no conference abstracts), English-language documents published from 1990-2018.
2. Guidelines must focus on the structured elicitation of explicitly probabilistic judgements from experts (i.e., no papers primarily about eliciting rankings, paired comparisons, or other non-probabilistic information from experts).
3. Guidelines must offer recommendations for practice concerning more than one of the stages of an elicitation (i.e., design, preparation, conduct, and analysis).

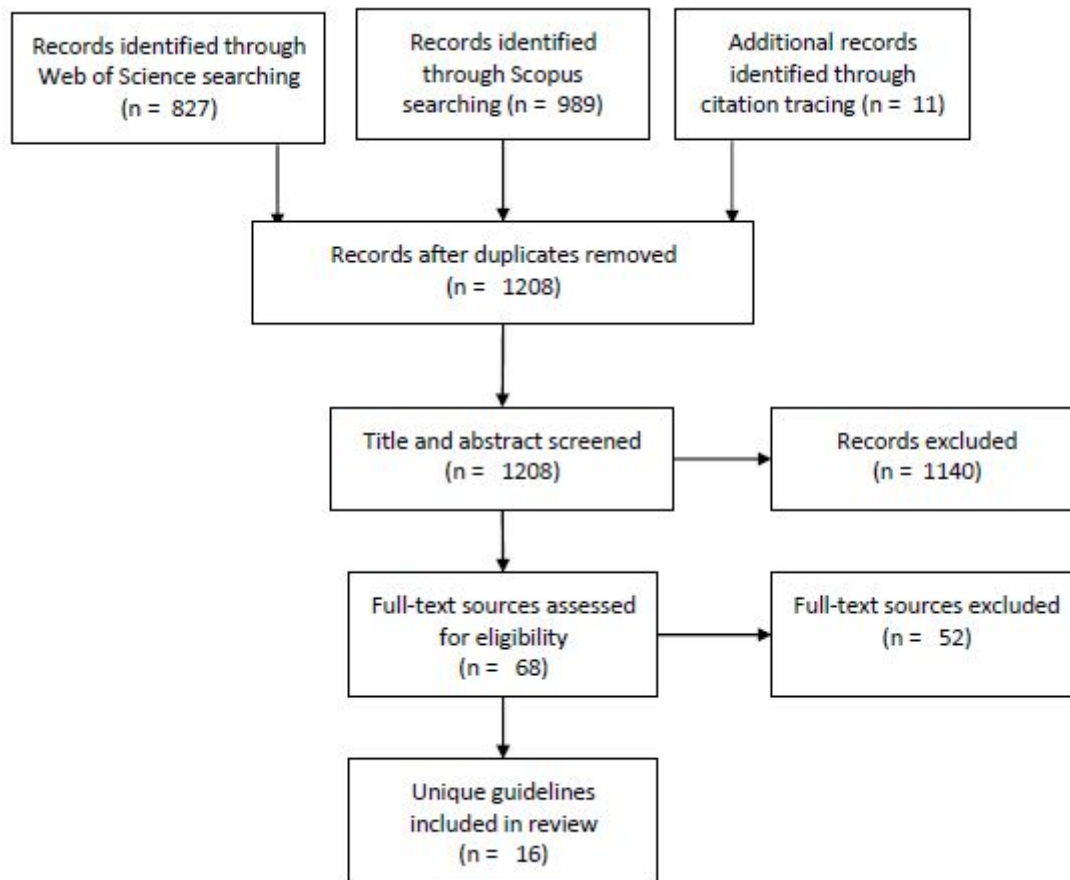
Guidelines are published as scientific papers and as policy documents or other grey literature, so we adopted a broad search strategy designed to capture all relevant guidelines (Figure 1). We searched Scopus and Web of Science for the period from 1990 to 2018 using the following keywords: “(expert AND (judgment OR judgement OR opinion) AND (elicit*) AND (method* OR protocol OR procedure OR guid* OR technique)),” which yielded 827 results in Web of Science and 989 in Scopus. Results were screened based on the title and abstract. If a paper potentially met the inclusion criteria based on its title and abstract, the full-text source was reviewed. The references lists in each of the full-text articles were reviewed to identify additional possible guidelines.

Extraction template

For each of the included SEE guidelines, information was gathered on the elicitation process in an extraction template (Table 1). This describes the elicitation process as pertaining to 3 stages: 1) preparation and design, 2) conduct and 3) post-elicitation. The extraction template was based on previous reviews of the elicitation process and was piloted and refined before use in this review.

Source	
Type of article	
Domain	
Self-reported objective	
PREPARATION AND DESIGN	
What quantities to elicit?	
Who/how many experts?	
How to encode judgements?	
How to manage biases?	
How to approach validation?	
Piloting the exercise	
Training and preparation for experts	
Training for other roles	
ELICITATION	
Level of elicitation	
Mode of administration	
Feedback to experts & revision	
Opportunity for interaction	
Feedback from experts on process	
Rationales	
AGGREGATION, ANALYSIS, AND POST-ELICITATION	
If/how to aggregate	
Fit to distribution	
Adjusting judgements	
Documentation	

Search results



In some cases, multiple sources with the same or similar author lists provide very similar recommendations. For the purposes of this review, these sources are considered to be the same guidance, and only one version was included. The earliest complete version of the guidance that appeared was used, but references to the “duplicate” guidelines are included in the next section. In the screening process, papers that focused on only one aspect of the elicitation process were excluded, such as how to encode judgements or fit judgements to distributions. Descriptions of software that did not discuss aspects of elicitation not managed within the software, were excluded. Reviews of past elicitation work or where experts can be used within a specific field were excluded, if they did not offer recommendations for practice. Finally, cases studies that focused on the presentation of a set of results rather than a methodology for future work were excluded.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Guidelines included in the review and related sources identified in the search

Key reference	Domain	Type of article
(Choy, O'Leary et al. 2009)	Ecology	Review of evidence and/or practice
Classical Model (Cooke and Goossens 2000)	Generic	Review of evidence and/or practice; Reflection on personal practice
EFSA Delphi (European Food Safety 2014)	Food safety	Agency guidance
(EPA 2009)	Environmental protection	Agency guidance
(Garthwaite, Kadane et al. 2005)	Generic	Review of evidence and/or practice
IDEA protocol (Hemming, Burgman et al. 2018)	Generic	Review of evidence and/or practice; Reflection on personal practice
(Ashcroft, Austin et al. 2016)	Insurance	Agency guidance
(Tredger, Lo et al. 2016)	Insurance	Agency guidance
(Kaplan 1992)	Risk and reliability	Reflection on personal practice
(Keeney and Vonwinterfeldt 1991)	Nuclear	Reflection on personal practice
(Knol, Slottje et al. 2010)	Environmental health	Review of evidence and/or practice
(Meyer and Booker 2001)	Generic	Review of evidence and/or practice
(Kotra, Lee et al. 1996)	Nuclear	Agency guidance
(Budnitz, Apostolakis et al. 1997)	Nuclear	Agency guidance
SHELF (Gosling 2018)	Generic	Review of evidence and/or practice; Reflection on personal practice
(Walls and Quigley 2001)	Risk and reliability	Reflection on personal practice

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Choices available from existing guidelines

Element	Component
Selecting quantities	
What quantities to elicit	Type of parameter
	Type of quantity
	Selection criteria
	Principles for describing quantities
	Decomposition
	Handling dependence
Encoding judgements	General approach
	Use of visual aids
Selecting experts	
Number of experts	Number of experts
Selecting experts	Roles within SEE
	Desired characteristics for those provide judgements
	Identification procedure
	Selection procedure
	Possible selection criteria
Training and preparation	
Pilot the protocol	Pilot exercise
Training and preparation for experts	What to cover in training
Level and conduct of elicitation	
Mode of administration	Location
Level of elicitation	Level of elicitation
Feedback and revision	Type of feedback
	What to feed back
	Opportunity for revision
Interaction	Opportunity for interaction
Rationales	Rationales
Aggregation, analysis & post elicitation	
Aggregation	Aggregation
	Aggregation approach

	Fit
Fit to distribution	Distribution
	Fitting method
Feedback on process	Feedback from experts on process
Adjusting judgements	Methods for adjusting judgements
Documentation	What to include
Managing biases	
Managing heuristics and biases	Biases relevant for SEE
	Bias elimination or reduction strategies
Validation	
Validation	Characteristics/measures

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Level of agreement on recommendations and choices in SEE guidelines

Element	Component	Agreement level	Explanation
Selecting quantities			
What quantities to elicit	Type of parameter	Some disagreement	Guidelines agree that observable quantities are preferred, but disagree on whether directly eliciting model parameters is an acceptable choice.
	Type of quantity	Disagreement	Guidelines offer conflicting recommendations on whether eliciting probabilities (compared with other uncertain quantities) is an acceptable choice.
	Selection criteria	Some agreement	Fewer than five guidelines discuss this, but they agree selection criteria should be defined.
	Principles for describing quantities	Some agreement	Some guidelines describe slightly different principles (e.g., asking clear questions, ensuring uncertainty on elicited parameters impacts the final decision or model), but they do not conflict.
	Decomposition	Agreement	The guidelines that discuss decomposing the variables of interest all agree it should be a choice.
	Handling dependence	Some agreement	The guidelines that discuss dependence agree it should be avoided if possible or addressed separately, but they discuss a range of methods for considering dependence.

Encoding judgements	General approach	Disagreement	Guidelines recommend and discuss different, conflicting methods for encoding judgements.
	Use of visual aids	Some agreement	Fewer than five guidelines discuss this, but they agree visual aids can be a useful choice.
Selecting experts			
Number of experts	Number of experts	Agreement	The experts agree that multiple experts are important, with most guidelines recommending around 5-10 experts.
Selecting experts	Roles within SEE	Agreement	The guidelines are very consistent in their description of the roles involved with elicitation.
	Desired characteristics for those provide judgements	Some agreement	Characteristics discussed in the guidelines are largely consistent, aside from differing views on if normative expertise is a requirement or just desired.
	Identification procedure	Some agreement	Recommendations differ but do not conflict across the guidelines. Agency guidelines tend to offer more detail.
	Selection procedure	Some agreement	Recommendations differ but do not conflict across the guidelines. Agency guidelines tend to offer more detail.
	Possible selection criteria	Some agreement	Recommendations differ but do not conflict across the guidelines.
Training and preparation			
Pilot the protocol	Pilot exercise	Agreement	Almost all guidelines recommend conducting a pilot exercise.

1 2 3 4 5 6 7	Training and preparation for experts	What to cover in training	Some agreement	The lists of what should be included in training vary across guidelines but do not conflict.
8 9	Level and conduct of the elicitation			
10 11 12 13 14 15 16 17	Mode of administration	Location	Some agreement	Most guidelines agree that face-to-face administration is preferred, though remote options may be pragmatically useful alternative in some situations.
18 19 20 21	Level of elicitation	Level of elicitation	Disagreement	Guidelines recommend and discuss conflicting levels of elicitation.
22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40	Feedback and revision	Type of feedback	Some agreement	Recommendations differ but do not conflict across the guidelines.
41 42 43 44 45 46		What to feed back	Some agreement	Recommendations differ but do not conflict across the guidelines.
47 48 49 50 51 52		Opportunity for revision	Some agreement	Guidelines either recommend revision take place following an elicitation (as part of an iterative process or immediately following the elicitation) or further in the future, following a draft report or additional data collection.
53 54 55 56 57 58 59 60	Interaction	Opportunity for interaction	Disagreement	Guidelines offer conflicting recommendations about when and how to facilitate interaction between the experts.
	Rationales	Rationales	Agreement	Almost all guidelines recommend collecting expert rationales in some form.
	Aggregation, analysis and post-elicitation			
	Aggregation	Aggregation	Agreement	All guidelines discuss aggregation as a recommendation or valid choice.
		Aggregation approach	Disagreement	Guidelines offer conflicting recommendations on the approach

			and method to aggregate judgements.
Fit to distribution	Fit	Some disagreement	The guidelines make few recommendations, but their choices differ.
	Distribution	Some agreement	Fewer than five guidelines discuss this, but they generally agree that many parametric distributions could be chosen.
	Fitting method	Some agreement	Fewer than five guidelines discuss this, but they generally agree that choices include minimum least squares and method of moments.
Feedback on process	Feedback from experts on process	Some agreement	Fewer than five guidelines discuss this, and they recommend complementary approaches.
Adjusting judgements	Methods for adjusting judgements	Some disagreement	Fewer than five guidelines discuss this, but they offer different perspectives.
Documentation	What to include	Some agreement	The lists of what should be included in final documentation vary across guidelines but do not conflict.
Managing biases			
Managing heuristics and biases	Biases relevant for SEE	Some agreement	The lists of potential biases vary across guidelines but do not conflict.
	Bias elimination or reduction strategies	Some agreement	The list of possible strategies vary across guidelines but do not conflict.
Validation			
Validation	Characteristics/ measures	Disagreement	The guidelines differ in their definitions of validity and discussion of how the concept can be operationalised in an elicitation.

Summary of principles applied to choices for SEE in HCDM

	Element	Key messages from critique	Principles support	Principles do not support
<p>Preparation and design</p>	<p>Selecting quantities</p>	<p>Different quantities can be elicited that provide information on any single parameter of interest.</p>	<p><u>Types of quantities</u></p> <ul style="list-style-type: none"> ○ Observables such as probabilities (expressed as proportions or frequencies) 	<ul style="list-style-type: none"> ● Measures of central tendency in isolation ● Odds ratios ● Credible ranges
		<p>Also relevant is handling dependence, selection criteria, principles for describing quantities and decomposition/disaggregation.</p>	<p><u>Dependency</u></p> <ul style="list-style-type: none"> ● Ask only about independent variables ● Express dependent variables in terms of independent variables ● Use separate dependence elicitation methods 	
		<p>Lack of evidence on how that choice should be guided.</p>	<p><u>Choice of parameters</u></p> <ul style="list-style-type: none"> ● Definition of a selection criteria, such as minimal assessment of each possible uncertain parameter and sensitivity analysis to see which uncertain parameters have the biggest impact 	
		<p>The choice is largely driven by the practical constraints of the context.</p>	<p><u>Wording</u></p> <ul style="list-style-type: none"> ● Avoid vagueness ● Ask questions in a manner consistent with how experts express their knowledge ● Use neutral wording ● Do not use leading questions 	

			Decomposition	No decomposition
	Methods to encode judgements	The FIM, the roulette or chips and bins method has previously been used in HCDM. The VIM, has also had limited use in HCDM, utilising quantiles as opposed to the bisection method. There is no empirical evidence to support which of these two methods is most appropriate in HCDM.	<ul style="list-style-type: none"> • Fixed interval methods – all forms • Variable interval methods – all forms 	
	Managing biases	No studies have explicitly examined the effectiveness of debiasing techniques in this context. The appropriateness of many suggested methods for debiasing are uncertain.	<ul style="list-style-type: none"> • Give experts practice and feedback • Provide training on biases • Frame questions to minimize bias and ambiguity • Identify biases through discussion with expert • Provide relevant background evidence • Ask for upper/lower bounds first • Ask experts to specify the credible interval they have provided • Minimize and record conflicts of interest among the experts • Require the experts address conflicting information • Collect rationales from experts 	

			<ul style="list-style-type: none"> • Report anonymous results • Include external experts • Anticipate likely biases 	
	Validation	<p>There is uncertainty about which method to validate is more appropriate in this context.</p> <p>Methods to reduce variability may not be appropriate, in the interests of reflecting any between expert variation.</p>	<ul style="list-style-type: none"> • Faithfully capturing experts' beliefs • Fitness for purpose • Internal review • External review • Coherence • Consistency 	<ul style="list-style-type: none"> • Calibration • Calibration & informativeness scoring
	Selecting experts	<p>Lack of evidence in HCDM to make definitive statement about particular approaches.</p> <p>Is a need to include all three types of roles for 'experts': the facilitator, expert providing priors and generalists to advise on design etc.</p> <p>Desired characteristics for those providing judgements are above all a level of substantive experience and a willingness to participate. Other characteristics may be beneficial, in particular normative expertise but may be difficult to ensure in HCDM.</p> <p>Training and careful design can mitigate against the need for some of these.</p> <p>Identifying relevant experts in HCDM is</p>	<p><u>Roles in SEE</u></p> <ul style="list-style-type: none"> • Facilitator • Expert 	<ul style="list-style-type: none"> • Generalists
			<p><u>Desired characteristics</u></p> <ul style="list-style-type: none"> • Substantive expertise • Willingness to participate 	<ul style="list-style-type: none"> • Normative expertise • Ability to understand questions • Ability to apply skills
			<p><u>Identification</u></p> <ul style="list-style-type: none"> • Recommendations by peers, either formally or informally • Research outputs • Known experience • RFP to seek out experts • Experience • Profile matrix 	
			<p><u>Selection</u></p> <ul style="list-style-type: none"> • Disclosure of personal and financial interests • Pursue diversity 	<ul style="list-style-type: none"> • Formal selection criteria developed and applied • Review CVs <ul style="list-style-type: none"> • Profile matrix
		<p><u>Criteria</u></p> <ul style="list-style-type: none"> • Reputation 	<ul style="list-style-type: none"> • Balance of internal and 	

		more likely to be driven by practical constraints, however ensuring a generalizable and wide sample is preferred.	<ul style="list-style-type: none"> • Experimental experience • Publication history • Diversity in background • Conflicts of interest • Awards • Balancing different viewpoints • Peer assessment (such as GEM) 	external experts (e.g., include at least 2 external experts)
			<p><u>Number</u></p> <p>No definitive guidance on number but seems to suggest at least 5-9 experts</p>	
	Pilot exercise	The ability to conduct a pilot may be driven by the constraints in HCDM.	Use of piloting	No piloting
	Training and preparation for experts	Training is essential for non-normative experts, although there is uncertainty about what should be contained within the training. Details about how elicited distributions will be used may not be possible to feedback back to time constraints in HCDM.	<ul style="list-style-type: none"> • Probability, including subjective probability • Motivation for elicitation • Description of what is required from experts • Outline of process • Outline of questions • Example and practice questions • Review of potential biases • Motivation of elicitation 	<ul style="list-style-type: none"> • Description of performance assessment • Introduction to dependence • List of relevant information • How results will be used • The full protocol
Elicitation	Level of elicitation	Discussion may not always be feasible due to the constraints of HCDM. Instead it may be possible to do this remotely via a Delphi, specifically after	<ul style="list-style-type: none"> • Individual • Combination 	<ul style="list-style-type: none"> • Consensus

		individual distributions have been elicited. Group interaction may introduce biases, such as overconfidence.		
	Mode of administration	The mode of administration may be driven by the constraints of HCDM.	<ul style="list-style-type: none"> • Face-to-face • Remote 	
	Feedback to experts and revision	<p>It is uncertain which types of information should be presented and at which stage of the SEE.</p> <p>The information relayed to experts may be driven by their level of normative skills.</p>	<p><u>Type of feedback</u></p> <ul style="list-style-type: none"> • Graphical feedback • Distributions from other experts • Summaries of aggregated distributions • Rationales • Qualitative discussion of elicited values • Written description of experts rationales 	<ul style="list-style-type: none"> • Fitted distribution • Performance scores • Results using elicited values • Future data • The draft elicitation report <ul style="list-style-type: none"> • Decision resulting from the expert judgement
			<p><u>Opportunities for revision</u></p> <ul style="list-style-type: none"> • A set number of elicitation/feedback rounds from the outset • Update after future data is collected • Update for revisions/clarifications after circulating draft elicitation report • Individuals update during or after a session based on graphics or other information on fitted distribution 	
	Opportunity for interaction	There is little practical experience in HCDM with different methods	<ul style="list-style-type: none"> • No interaction • Group discussion prior to 	

		of interaction between experts.	<p>individual elicitation</p> <ul style="list-style-type: none"> • Group discussion and group elicitation • Group discussion following individual elicitation (with opportunity for revision) • Remote, anonymized interaction 	
	Feedback from experts on process	Feedback is often undertaken in SEE for HCDM, although there are not consistent approaches to do this. The practicalities of conducting SEE in HCDM, may dictate if the method of feedback is practically plausible.		<ul style="list-style-type: none"> • Ask experts to appraise the elicitation exercise after completing it. • Get feedback on the procedure if future data collection contradicts elicitation results
	Rationales	There is no practical experience in HCDM with different methods of providing rationales.	Collect/record rationales from experts (about how they made their judgments)	No mention of rationales
Aggregation, analysis and post-elicitation	If/how to aggregate	Behavioural methods of aggregation may be practically difficult in the context of HCDM, both in terms of convening experts but also in terms of the provision of experienced facilitators.	<ul style="list-style-type: none"> • Mathematical using linear opinion pooling • Combination 	Behavioural

	Fit to distribution	<p>In HCDM, fitting of a smooth distribution would seem appropriate.</p> <p>The choice of parametric distribution is uncertain. There is a lack of evidence in HCDM on the fitting process in SEE. Limited evidence suggests that standard distributions, such as the Beta will often be sufficient.</p> <p>More complex approaches may be appropriate, however these can be complex to implement in general software.</p>	<p>Fitting</p> <p><u>Distributional form</u></p> <ul style="list-style-type: none"> • Normal • Beta • Other conjugate family <p><u>Selection criteria</u></p> <ul style="list-style-type: none"> • Minimum least squares • Method of moments • Other approaches 	<p>Not fitting</p> <ul style="list-style-type: none"> • Uniform • Triangular • Uniform over elicited intervals
	Adjusting judgements	<p>There is a lack of practical experience in HCDM to inform the choice of adjustment method. Given the role of facilitator in SEE in HCDM it would seem inappropriate for the facilitator to adjust themselves using arbitrary criteria.</p>	<p>Adjust/not adjust</p> <ul style="list-style-type: none"> • Calibrate • Adjust to improve coherence 	<p>Analyst adjustment and feedback</p>
	Documentation	<p>In order to inform a decision making process in HCDM, a SEE should document all details, including elicitation questions, the responses, fitting process, level of elicitation, interaction, revision and validation</p>	<p>Thorough documentation</p>	<p>Less detailed/no documentation</p>