Lessons from structured expert elicitation using the IDEA protocol



Marissa McBride m.mcbride@imperial.ac.uk @riss mcb





Why expert judgement











A REPORT OF THE INTERGOVERNMENTAL PANEL ON CLIMATE CHANGE









Experts make mistakes...

Cognitive bias: overconfidence

Challenger



NASA 1985: 'the risk of catastrophic failure is 1 in 100,000 launches'

Chernobyl



Ukranian Minister of Power 1987: *'risk of a meltdown is* 1 in 10,000 years



Motivational bias: global warming

Loss of gross world product resulting from a doubling of atmospheric CO₂ by 2050 *(Kammen and Hassenzahl 1999)*

Structured Elicitation Protocols



IDEA: A Delphi-like protocol with a twist

Strategies for better expert	Pre-elicitation		Elicitation		Post-elicitation
(i) A four-point question format for eliciting quantities to mitigate the overconfidence	Background information compiled. Contact and brief experts on the elicitation	All experts individually answer questions, and provide reasons	DISCUSS Experts shown anonymous answers from each participant and visual	ESTIMATE All experts make 2nd final and private estimate	AGGREGATE Mean of experts' 2nd round responses calculated. Experts may review and discuss individual and
(ii) The structured interaction of experts via facilitated discussion to promote 'better' group judgments	process	for their judgements	summary of responses		group outcomes, add commentary, and correct residual misunderstandings

The IDEA protocol (Investigate, Discuss, Estimate, Aggregate)



Four-step question format



Expert	5 th (lower)	50 th (best)	95 th (upper)
1	2	12	34
2	4	15	50
3	7	9	40
4	20	22	23
Average	8.25	14.5	36.75

Interval judgements

• First step towards better judgements: express uncertainty through intervals (common in environmental science)



Generally has a pre-assigned confidence,

e.g. "Provide an interval which you are 90% sure contains the true number of native fish species listed in the Murray River catchment".

BUT people are insensitive to pre-assigned confidence levels

90% intervals tend to contain the answer only 50% of the time → **Overconfidence** (undue confidence in the intervals provided)

Overconfident intervals



(Hynes and Vanmarche 1977)

When 90% Confidence Intervals are 50% Certain:



The influence of question format



	Average
Elicitation Format with Examples	overconfidence
Range (one-point)	41%
I am 80% sure that this happened betweenand	
Two-point	23%
I am 90% sure that this happened after	
I am 90% sure that this happened before	
Three-point	14%
I am 90% sure that this happened after	
I am 90% sure that this happened before	
I think it's equally likely that this happened after or before	
Four-step	12%
Realistically, what do you think the lowest plausible value is?	
Realistically, what do you think the highest plausible value is?	
Realistically, what is your best estimate?	
How confident are you that the interval you created,	
from lowest to highest, could capture the true value?	

<u>Speirs-Bridge et al. 2</u>010

Some results – 3 v 4 step (for quantities)



The IDEA protocol (Investigate, Discuss, Estimate, Aggregate)



Worksh op	Discipline	# exper ts	Range of years of relevant experience (median)	Range of qualifications	Range of number of publications
1	Animal and plant Biosecurity and Quarantine	21	0 – 37 (17.5)	BSc, BASc, BVSc, BCom, Grad. Dip., MSc, PhD	0 – 113
2	Animal and plant biosecurity	24	0 – 39 (12)	BSc, MSc, MBA, MCom, PhD	0 – 270
3	Conservation biology, herpetology	13	0 – 42 (15)	BA, BSc, BSc (Hons), M Env Studies, PhD BEng, BSc	0 – 45
4	Public health, medicine and epidemiology	25	0 – 45 (12)	BEcon, LLB, MBBS, Grad. Dip., MA, MSc, MBA, PhD	0 – 220
5	Risk analysis, biosecurity	20	0 – 40 (6)	BEng, BSc, BVSc, MBBS, Grad. Dip., MA, MBA, PhD	0 – 225
6	Weed risk	14	0 - 50	BSc, MSc, PhD	1 – 220

Peer versus self assessments



Do peer assessments correlate with performance?



Peer assessment versus performance

Vorkshop	Peer assessment versus prediction accuracy
1	-0.391 (n=20)
2	0.215 (n=19)
3	0.190 (n=13)
4	- 0.360 (n=25)
5	0.305 (n=20)
6	0.367 (n=14)

Benefits of (diverse) groups and facilitated discussion:

The group average improvement in accuracy (ALRE) following discussion.

Estimating population sizes

What is your estimate of the current population size (number of mature individuals) for the Mulga Lands Bioregion?

© Australian Koala Foundation

Source: Adams-Hosking et al. (2016

The IDEA protocol (Investigate, Discuss, Estimate, Aggregate)

1. Recruit diverse group

2. *Investigate* + Initial private estimate (Round 1)

3. Aggregate + Feedback

4. Discuss

5. Revise *Estimate* (Round 2) + *Aggregate*

Geopolitical forecasting tournament

Program Manager

For information

info@iarpa.gov

Information

IARPA-BAA-10-05

IARPA Day Poster

Program

contact: dni-iarpa-

Aggregative Contingent Estimation (ACE)

The goal of the ACE Program is to dramatically enhance the accuracy, precision, and timeliness of intelligence forecasts for a broad range of event types, through the development of advanced techniques that elicit, weight, and combine the judgments of many intelligence analysts. The ACE Program seeks technical innovations in the following areas: (a) efficient elicitation of probabilistic judgments, including conditional probabilisties for contingent events; (b) mathematical aggregation of judgments by many individuals, based on factors that may include: past performance, expertise, cognitive style, metaknowledge, and other attributes predictive of

accuracy; and (c) effective representation of aggregated probabilistic forecasts and their distributions. The ACE Program will build upon technical achievements of past research and on state-of-the-art systems used today for generating probabilistic forecasts from widely-dispersed experts. The program will involve empirical testing of forecasting accuracy against real events.

Performers (Prime Contractors)

Applied Research Associates, Inc.; Charles Stark Draper Laboratory, Inc.; George Mason University; Jacobs Strategic Solutions Group, Inc.; University of California,

Research Area(s)

- Forecasting
- Human judgment
- Machine learning
- Logic & critical thinking

Related Publications

To access ACE program-related publications, please visit Google Scholar ⊠.

Related Article(s)

10110

Does Your Company Have Good Judgement? 앱

Superforecasting for the Farm 🗹

Are You a Good Forecaster? The Good Judgment Project Needs You 团

Three-step Question Format (for probabilities)

e.g. Will Liu Yandong be selected as a member of the next Politburo Standing Committee of the Communist Party of China?

Will Chinese armed forces or maritime law enforcement forces attempt to interdict or make physical contact with at least one U.S. government naval vessel or airplane or Japanese government naval vessel or airplane that it claims is in its territorial waters or airspace, before 1 May 2014?

Results: Effect of Discussion

For those who revised their judgments, second round estimates outperformed first round estimates.

COVID Recovery metrics

COVID Recovery metrics

When will enough doses of FDA-approved COVID-19 <mark>vac</mark> cine(s) to inoculate 25 million people be distrib	outed in the United States? Today's Forecast
Today's Forecast:ABefore 1 October 2020BBetween 1 October 2020 and 31 March 2021CBetween 1 April 2021 and 30 September 2021DBetween 1 October 2021 and 31 March 2022	1% 12% 26% 26%
E Not before 1 April 2022 Forecast History: When will enough doses of FDA-approved COVID-19 vaccine(s) to inoculate 25 million people be distributed in the United States? 100	35%
75 Storegative and the second	 A: Before 1 October 2020 B: Between 1 October 2020 and 31 March 2021 C: Between 1 April 2021 and 30 September 2021 D: Between 1 October 2021 and 31 March 2022 E: Not before 1 April 2022
JOGDANO SARA RAR RAR RAR RAR RAR RAR RAR RAR RAR	

Expert predictions of US COVID-19 cases

Dr Thomas McAndrew, Reich lab University of Amherst

http://reichlab.io/

http://www.thomasmcandrew.com/

Lessons learnt for better structured elicitations

Ask individuals to ...

- Answer the same question in different ways (lowest, highest, most likely)
- Indicate confidence
- Examine estimates made by other people (feedback), consider counter-argument
- Revise original estimates after feedback
- Anticipate issues with conditional probabilities, base rates, ...

Then, don't rely on individuals...

- Discuss questions to eliminate linguistic uncertainty
- Make groups diverse—age, gender, background and cognitive style, culture
- Aggregate mathematically don't force a consensus
- Avoid group think— Delphi structures / anonymity in judgments

Acknowledgements

Mark Burgman Anca Hanea Victoria Hemming Bonnie Wintle Terry Walshe Fiona Fidler Many more....

The experts involved

