Expert Elicitation of Low Probability Events: The Roles of Personality and Feedback





Darren Hudson, Combest Endowed Chair, Texas Tech University Keith Coble, Giles Distinguished Professor, Mississippi State University





Presented to the Department of Behavioral Sciences, United States Military Academy, March 2014.

Needs for Subjective Risk Assessment



- Most situations of interest have limited data on which to base objective assessments.
- The areas of interest are the tails of distributions, and discernment of low probability events poses problems for many people.
- Obvious military/intelligence applications

Issues with Subjective Assessment

- Methods of elicitation and aggregation
 - Composite
 - Composite with Feedback
 - Consensus
- Concerns about heuristics and biases brought in by experts (Kahneman et al.)
- Personality and risk preferences may influence answers (Chauvin et al.)

Issues with Subjective Assessment

- No real consensus on the proper elicitation methods
 - Real need to rigorously examine elicitation techniques

Objectives

- Examine the impact of elicitation procedure on subjective assessment accuracy and bias
- Examine the role of individual personality type and risk preferences on subjective assessment and accuracy and bias
- Utilize an economic experiment to insure incentive compatibility

Methods

- Generic—recruited students from the undergraduate and graduate populations
- Given a statistical knowledge test
 - General knowledge, 10 questions multiple choice
 Used for "weighting" expertise
- Myers-Briggs Personality Inventory
- Risk Aversion Experiment

Risk Aversion

Table 1. Risk preference decision sheet.

Question	Option A	Option B	Which option preferred?
	10% chance of \$10.00,	10% chance of \$19.00,	
	90% chance of \$8.00	90% chance of \$1.00	
2	20% chance of \$10.00,	20% chance of \$19.00,	
	80% chance of \$8.00	80% chance of \$1.00	
3	30% chance of \$10.00,	30% chance of \$19.00,	
	70% chance of \$8.00	70% chance of \$1.00	
1	40% chance of \$10.00,	40% chance of \$19.00,	
	60% chance of \$8.00	60% chance of \$1.00	
5	50% chance of \$10.00,	50% chance of \$19.00,	
	50% chance of \$8.00	50% chance of \$1.00	
5	60% chance of \$10.00,	60% chance of \$19.00,	
	40% chance of \$8.00	40% chance of \$1.00	Table 2.
7	70% chance of \$10.00.	70% chance of \$19.00,	
	30% chance of \$8.00	30% chance of \$1.00	No. of
3	80% chance of \$10.00.	80% chance of \$19.00,	safe
<i>.</i>	20% chance of \$8.00	20% chance of \$1.00	choices
)	90% chance of \$10.00.	90% chance of \$19.00.	
, ,	10% chance of \$8.00	10% chance of \$1.00	0-1
10	100% chance of \$10.00.	100% chance of \$19.00.	2
	0% chance of \$8.00	0% chance of \$1.00	3
10	90% chance of \$10.00, 10% chance of \$8.00 100% chance of \$10.00, 0% chance of \$8.00	90% chance of \$19.00, 10% chance of \$1.00 100% chance of \$19.00, 0% chance of \$1.00	

Fable 2. Risk aversion coefficient.

is.

No. of safe choices	Range of relative aversion for $U(W) = \frac{W^{1-r}}{1-r}$	Middle point of relative risk aversion	Risk preference classification
0-1	$-1.76^{\circ} < r^2 < -0.93$	-1.365	Highly risk loving
2	$-0.97 < r^2 < -0.49$	-0.73	Very risk loving
3	$-0.49 < r^2 < -0.13$	-0.31	Risk loving
4	$-0.13 \le r^2 \le 0.19$	0.03	Near risk neutral
5	$0.19 \le r^2 \le 0.48$	0.335	Slightly risk averse
6	$0.48 < r^2 < 0.78$	0.63	Risk averse
7	$0.78 < r^2 < 1.13$	0.955	Very risk averse
8	$1.13 \le r^2 \le 1.6$	1,365	Highly risk averse
910	$1.6 \le r^2 \le 2.2^a$	1.9	Stay in bed

"Those two lower and upper bound are subjectively determined.

Using a 10-sided die insures randomness for the respondent. Payoffs are real.

Economic Experiment

 Respondents provided with samples of data drawn from <u>known</u> distributions.

Table 1. Four datasets with known population distribution										
Range	[-∞, 80)	[80,90)	[90, 100)	[100, 110)	[110, 120)	[120,+x				
Datasets	Probability									
Dataset 1: sample of 30 drawn from normal population distribution N (100, 15)	0.0912	0.1613	0.2475	0.2475	0.1613	0.0912				
Dataset 2: sample of 50 drawn from normal population distribution N (100, 25)	0.2119	0.1327	0.1554	0.1554	0.1327	0.2119				
Dataset 3: sample of 100 drawn from normal population distribution N (100, 10)	0.0228	0.1359	0.3413	0.3413	0.1359	0.0228				
Dataset 4: sample of 30 drawn from beta population distribution β (6, 2) with range of (0, 133.33)	0.1624	0.1203	0.1653	0.199	0.2043	0.1487				



Economic Experiment

Presented Questions

Payoffs Generated

- (1) What is the chance the variable will fall below 80 next period?
- (2) What is the chance the variable will fall at or above 80 and below 90 next period?
- (3) What is the chance the variable will fall at or above 90 and below 100 next period?
- (4) What is the chance the variable will fall at or above 100 and below 110 next period?
- (5) What is the chance the variable will fall at or above 110 and below 120 next period?
- (6) What is the chance the variable will fall at or above 120 next period?

Payoff = Maximum
$$\left(\$0.00, \$10.00 - 0.025 \times \sum_{i=1}^{6} (A_i - B_i)^2\right)$$

The payoffs are intended to focus the respondent's attention on accuracy, and, thus, maximum payoff. Scalar effects of the size of payoff are always relevant, but at least a weak preference for accuracy is generated by this incentive compatibility measure.

Economic Experiment

- Three treatments
 - Treatment 1—"composite treatment"-respondents provided only individual responses (no feedback), n = 35
 - Treatment 2—" composite with feedback"respondents received feedback in the form of all individual (anonymous) responses posted on computer screen with opportunity to revise, n = 25 (5 groups of five)
 - Treatment 3—"consensus"-respondents received feedback in the form of all individual (anonymous) responses posted; were paid a \$2 bonus of all group members reported the same probability in final round

Analysis

- Two analyses
 - Individual errors (SSE)
 - Aggregated errors (MSE) and biases

Table 6. Demographic characteristics of participants (continued).

Treatment	Statistic	Statistical class Statistical score Risk preference						Myers-Briggs type indicator						
					No. of safe	choice ¹	E		S		Т		J	
	Average	$\overline{SD^a}$	Average	SD^a	Average	SD^{k}	Average	SD ^a	Average	SD ^a	Average	SD^{n}	Average	SD ^a
1 2 3 Total	1.600 1.600 1.200 1.467	1.4142 2.1602 0.9129 1.5711	8.1200 8.4400 7.1200 7.8930	1.2689 1.4457 1.4236 1.4757	5.2500 4.6842 5.0588 5.0000	1.6819 1.4550 1.9834 1.6949	0.60 0.44 0.60 0.55	0.50 0.51 0.50 0.50	0.36 0.36 0.48 0.40	0.49 0.49 0.51 0.49	0.48 0.52 0.60 0.53	0.51 0.51 0.50 0.50	0.60 0.52 0.52 0.55	0.50 0.51 0.51 0.50

^aSD-standard deviation.

¹We drop those subjects who switched from B back to A. So the valid observations are 20, 19 and 17 in Treatment 1, 2 and 3, respectively.

$$\begin{split} \mathbf{EIE} &= \beta_1 + \beta_2 \mathbf{TD}_2 + \beta_3 \mathbf{TD}_3 + \beta_4 \mathbf{DD}_2 + \beta_5 \mathbf{DD}_3 + \beta_6 \mathbf{DD}_4 + \beta_7 \mathbf{MD}_1 \\ &+ \beta_8 \mathbf{MD}_2 + \beta_9 \mathbf{MD}_3 + \beta_{10} \mathbf{MD}_4 + \mathbf{STAT} + \mathbf{RISK} + \varepsilon \end{split}$$

- TDs are the treatment effects (composite base)
- DDs are the dataset effects (small sample normal base)
- MDs are MBPI designations
 - MD1 = 1 if extrovert
 - MD2 = 1 if sensing
 - MD3 = 1 if thinking
 - MD4 = 1 if judging

Variables	Error of indiv across entire	idual estimate distribution	Error of individual estimate Error of individ			Error of individual estimate across right tail		ual estimate across tails	
	(-∞,	+∞)	(-∞,	, 80)	(120, -	H00)	(-∞, 80) and (120, +∞)		
	Coefficient	p value	Coefficient	p value	Coefficient	p value	Coefficient	p value	
Intercept TD ₂ TD ₃ DD ₂ DD ₃ DD ₄ MD ₁ MD ₂ MD ₃ MD ₄ STAT RISK	861.6293 -97.0366 -114.5273 -437.1964 -435.9869 -195.6332 -15.3605 8.5139 42.5189 -87.014 -25.7053 7.7996	<0.001*** 0.0091*** 0.0105 <0.001*** <0.001*** 0.6323 0.8005 0.223 0.0078*** 0.0498 0.7322	85.9054 -8.2458 -10.5297 -0.1559 -14.6464 -1.5941 -1.2013 8.7103 -3.1618 -5.8462 -4.5868 -3.0151	<0.001*** 0.0826* 0.0657 0.9772 0.0077*** 0.7804 0.7704 0.0448** 0.4793 0.1608 0.0066*** 0.412	141.2804 -25.3732 -26.1358 -0.1195 -2.8851 24.1884 -27.4404 39.2367 -2905273 -24.2048 -8.9195 -2.2178	0.0072*** 0.1171 0.1796 0.9949 0.8779 0.2156 0.0516* 0.0083*** 0.0083*** 0.089* 0.1189 0.8594	207.1858 -33.619 -36.6655 -0.2755 -17.5015 22.5944 -28.6477 47.947 -32.6891 -30.051 -13.5063 -5.2329	0.0005*** 0.0673* 0.0971* 0.9896 0.4066 0.3073 0.0729* 0.0045*** 0.0045*** 0.0594* 0.0594* 0.0627* 0.0377 0.7125	
No. of Obs R ²	224 0.4502		224 0.1249		224 0.106		224 0.1266	0.7120	

Table 7. Regression output for the error of individual estimate model.

***, ** and * denote significance at the 1, 5 and 10% levels, respectively.

- Feedback improves accuracy
 - Means of feedback seems less important than the existence of feedback
- Distribution/sample size appears not to significantly affect error rates
- Knowledge/experience matters
 - Limited in scope here, but supports the idea that true experts produce better subjective estimates



- Personality effects
 - Extroverts performed better—focus attention to outer world of people and things; perhaps outward orientation leads to better induction?
 - Sensing were more inaccurate—focus attention on five senses and discount "gut feelings"; intuitive people listen to their unconscious mind
 - Thinkers performed better—focus attention on facts and objective reasoning rather than person-centered emotional responses
 - Judgers performed better—prefer planned and organized approach to life, likely leads to more organized thought processes

Personality matters...now you just have to figure out how to effectively screen for it with your experts.

Aggregate Errors

Variables	MSE of aggreg left tail	ated estimate across	MSE of aggreg right tail	ated estimate across	MSE of aggregated estimate across two tails (-00,80) and (120, +00)		
	(-∞, 80)		(120,+∞)				
	Coefficient	p-value	Coefficient	p-value	Coefficient	p-value	
Intercept TD2 TD3	16.7026 1.2681 2.9321	<0.0001*** 0.4256 0.0763*	23.661 7.0779 5.1874	<0.0001*** 0.0005*** 0.0061***	71.389 -9.5258 -12.6912	<0.0001*** 0.0469** 0.011**	
DD2 DD3 DD4	-8.621 -13.9041 -2.5515	0.0002*** <0.0001*** 0.173	-3.4897 -15.2992 6.366	0.0858* <0.0001*** 0.004***	-24.1839 -57.1299 5.3062	0.0002*** <0.0001*** 0.3157	
WD	-0.1949	0.8797	-0.1199	0.9304	-0.1132	0.9755	
No. of obs. R ²	24 0.8178		24 0.9013		24 0.9192		

Table 3. Regression output for the MSE of aggregated estimate model

Note: ***, ** and * denote significance at the 1, 5 and 10% levels, respectively.

Aggregate Errors



- Again, feedback matters, but it does not seem to matter the form of the feedback.
- Here, sample size does matter and reduces aggregate error rates
- Weighting scheme does not seem to matter

Aggregate Bias

	Bias of aggregated estimate across left tail		Bias of aggreg right tail	ated estimate across	Bias of aggregated estimate across two tails		
	(−∞,80)		(120, +∞)		(-∞, 80) and (120, +∞)		
Variables	Coefficient	p-value	Coefficient	p-value	Coefficient	p-value	
Intercept TD2 TD3	3.2218 0.7805 0.2199	<0.0001*** 0.005*** 0.3771	4.6062 -0.4013 -0.6497	<0.0001*** 0.2114 0.0507*	7.828 0.3792 -0.4297	<0.0001*** 0.4555 0.3987	
DD2 DD3 DD4	-5.5765 -4.3268 -0.6108	<0.0001*** <0.0001*** 0.0435*	- 7.9566 - 3.5692 0.4798	<0.0001*** <0.0001*** 0.1964	-13.533 -7.8961 -0.131	<0.0001*** <0.0001*** 0.8219	
WD	0.0339	0.866	-0.0259	0.9195	0.008	0.9844	
# of Obs	24		24		24		
R*	0.9718		0.5	77	0.97	0.9789	

Table 4. Regression output for the bias of aggregated estimate model

Note: *** and * denote significance at the 1 and 10% levels, respectively.

Aggregate Bias

- No real consistent pattern in bias relative to feedback
 - Perhaps the consensus reduced it a bit, but the relationship is not particularly strong
- Bigger samples help reduce bias, but not clear if there is an "optimal" sample size
- Weighting scheme does not matter

Conclusions

- Use of incentive compatible experiment reveals some useful information:
 - Feedback is useful, but the manner is less important
 - If aggregating with simple feedback is cheaper than a "Delphi" approach, it may be preferred
 - Sample size matters in the aggregate, but no so much at the individual level
 - Personality and experience or knowledge matter, but risk aversion does not
 - No consistent pattern to bias
 - We did not allow non-anonymous interaction where personality and bias may interact; could be a real weakness of "Delphi" approaches



Recommendations

- Any attempt to elicit subjective estimates of risk should include:
 - Some manner of interactive/iterative feedback
 - Persons screened for specific expertise
 - Persons screened for personality types
 - Attention to sample size...but, more attention to feedback