

A Penny for Your Thoughts: A Survey of Methods for Eliciting Beliefs.*

Karl Schlag[†] James Tremewan[‡] Joël van der Weele[§]

August 20, 2014

Abstract

Incentivized methods for eliciting subjective probabilities in economic experiments present the subject with risky choices that encourage truthful reporting. We discuss the most prominent elicitation methods and their underlying assumptions, provide theoretical comparisons and give a new justification for the quadratic scoring rule. On the empirical side, we survey the performance of these elicitation methods in actual experiments, considering also practical issues of implementation such as order effects, hedging, and different ways of presenting probabilities and payment schemes to experimental subjects. We end with a discussion of the trade-offs involved in using incentives for belief elicitation and some guidelines for implementation.

Keywords: belief elicitation, subjective beliefs, scoring rules, experimental design.

JEL-codes: C83, C91, D83.

*We would like to thank Peter Wakker, Theo Offerman and Glenn Harrison for useful comments.

[†]University Vienna, Vienna. E-mail: karl.schlag@univie.ac.at

[‡]University Vienna, Vienna. E-mail: james.tremewan@univie.ac.at

[§]Corresponding author. University of Amsterdam. Email: vdweele@uva.nl. Tel. +31 (0)20 5254213. Address: CREED, Department of Economics, University of Amsterdam, Roeterstraat 11, 1018WB Amsterdam, the Netherlands.

1 Introduction

Economists are interested in the choices people make as well as the reasons they have for making these choices. Choice data alone are often not enough to understand these reasons and we need information about the beliefs of the decision maker. For example, first movers in ultimatum bargaining experiments may make high offers because they are altruistic, or because they believe the other party will reject low offers. Eliciting beliefs can help disentangle these hypotheses (Manski, 2002). Belief measurement is also necessary when beliefs are an object of study in themselves, as in experiments about expectation formation and updating.

In keeping with the mainstream models of decision making, economists commonly elicit beliefs as probabilistic statements. The underlying assumption is that beliefs take the form of subjective probabilities, or can be usefully expressed in this form. As Manski (2004) points out, probabilistic statements allow for comparisons with objective frequencies and an evaluation of the consistency of beliefs. A second characteristic of (experimental) economists' elicitation methods is the use of rewards based on how well the reported belief matches actual events. This practice stems from the old idea that appropriately designed bets give the decision maker an incentive to report subjective beliefs truthfully (Ramsey, 1926).

In this article we give an overview of the available methods for the elicitation of subjective probabilities. In Section 2 we outline more formally the elicitation environment and discuss a class of mechanisms called scoring rules. Scoring rules reward subjects on the basis of the submitted report, and the actual realization of a random variable. In Section 2.2 we discuss so-called 'proper scoring rules' that have been designed to elicit probabilities, means, modes and various quantiles truthfully, under the assumption that the subject is risk neutral. In Section 2.3 we provide a new justification for a particular rule, the Quadratic Scoring Rule. In Sections 2.4, 2.5 and 2.6 we discuss mechanisms that abandon the restrictive assumption of risk neutrality. In Section 2.7 we discuss promising extensions to the standard framework.

We evaluate elicitation mechanisms empirically in Section 3, by looking at the 'quality' of elicited beliefs. Such evaluation is not straightforward, since by the nature of the exercise we do not know the right benchmark, i.e. the true belief of the subject. We discuss several alternative benchmarks that may be used to assess the effectiveness of elicitation mechanisms. In Section 4 we discuss practical issues in implementation such as the complexity and presentation of elicitation schemes. Finally, in Section 5 we tie our findings together and present some thoughts about the appropriate use of incentives in belief elicitation, and directions for future research on this topic.

Our survey is complemented by other reviews of the broader aspects of belief elicitation. Manski (2004) focuses on questionnaires in consumer and household surveys. In this context, incentivized elicitation is typically not possible, since beliefs cannot be immediately verified. Garthwaite *et al.* (2005) and Jenkinson (2005) present a survey of mostly unincentivized elicitation techniques with a wide range of applications. Delavande *et al.* (2011) consider belief

elicitation in studies with subjects in developing countries who may have low levels of numerical literacy. Gneiting and Raftery (2007) present a technical review of proper scoring rules, including a detailed discussion of decision theoretic foundations and statistical properties. Winkler (1996) focuses on a limited set of proper scoring rules and discusses scoring rules as a tool for ex-post evaluation of forecasts, an issue we ignore in this survey.

Finally, there is some overlap between this paper and the survey by Schotter and Trevino (2014), which was written simultaneously and independently. While both studies discuss eliciting beliefs about the probability of an event, we also discuss mechanisms to elicit means, medians and other characteristics of belief distributions. Schotter and Trevino (2014) contains a more extensive treatment of the link between actions and beliefs and of second order beliefs, while this paper pays more attention to the empirical comparisons of different methods and practical issues of implementation.

2 Scoring rules

In this section, we discuss mechanisms used to incentivize truthful belief elicitation. We consider an ‘experimenter’ who wishes to learn something about the beliefs that a ‘subject’ holds about some future event. Below we formalize this setting and focus on a class of mechanisms known as scoring rules: payment functions that depend on the report of the subject and the realization of the event. Most of this section is concerned with the elicitation of the probability of a single event. We also discuss elicitation of means, modes, quantiles, confidence intervals and density functions (Sections 2.2.2 - 2.2.6, respectively), using games and prediction markets to elicit beliefs (Section 2.6), the elicitation of sets (Section 2.7), and payments based on multiple events, which can be used to elicit a number of additional characteristics including variances (Section 2.7).

2.1 Preliminaries

Consider an experimenter who is interested in a quantity that is related to the beliefs of the subject about the distribution of a random variable X . For instance, the experimenter is interested in the probability that $X = 1$, or in the expected value of X . Let \mathcal{X} denote the set of possible realizations of X . To simplify exposition we assume that \mathcal{X} is finite, which is typically satisfied in experiments. For instance, if X is an integer value on a scale that ranges from 1 to 10, then $\mathcal{X} = \{1, \dots, 10\}$. Let P_X be the probability distribution that describes the beliefs of the subject, so for each value x belonging to \mathcal{X} it describes the probability that the subject believes $X = x$. In this paper, truth (i.e. the actual or empirical distribution) plays no role as we are interested solely in the subjective beliefs held by the subject. For this reason, when we write X we refer to a random variable that is distributed according to P_X .

The experimenter does not know P_X and wishes to learn about some characteristic θ of this

belief distribution. Common examples for θ are the probability of an event and the expected value or median of X , but θ can also be a more sophisticated object, such as a 95% confidence interval for X such that the probabilities of X falling below and above the interval are equal. Let \mathcal{P}_X be the set of possible belief distributions of the subject and Θ the set of possible values of θ . Formally, θ is a mapping from the set of distributions into Θ .¹

It is natural to ask the subject to report θ directly. In fact, there is no loss of generality to limit our presentation to payment schemes where the subject is asked to report θ . This follows from the same arguments that lead to the revelation principle in mechanism design. In what follows, we denote the subject's report by r and her true subjective belief by p .

To simplify elicitation, one typically only considers payment based on a single realization of X . This is the case we consider here. Settings where it is useful or even necessary to condition on multiple realizations are discussed in Section 2.7. So the payment scheme S , which is also called a *scoring rule*, is a mapping $S : \Theta \times \mathcal{X} \rightarrow \mathbb{R}$, where $S(r, x)$ is the amount of money paid to an expert when outcome x is realized after the expert has reported r . We assume that the realization X is independent of the report r of the subject.

To predict the effect of incentives on reporting behavior, we have to specify the decision making process of the subject. For now, we stay within the canonical model of decision making, and assume that the subject is an expected utility maximizer. Alternative models of decision making are discussed later. Suppose the subject has some utility function u , and given the payment function S reports an element of

$$\arg \max_{r \in \Theta} Eu(S(r, X))$$

where

$$Eu(S(r, X)) = \sum_{x \in \mathcal{X}} u(S(r, x)) P(X = x)$$

is the expected utility of reporting r .

Now consider the experimenter, who wishes to design the incentives such that they induce the subject to tell the truth. A first issue is that the experimenter may not know the utility function of the subject. Let U be the set of possible objective functions of the subject as assessed by the experimenter. If U contains a single element then it is as if the experimenter knows the objective of the subject. For instance, if $U = \{Id\}$, where $Id(x) = x$ for all x , then we are considering the case where the experimenter believes that the subject is risk neutral.

We call a scoring rule 'truth-telling' if it induces the subject to tell the truth, regardless of which utility function $u \in U$ the subject is basing her choices on.² More formally, S is called a

¹Here we assume that θ is uniquely determined given X , so $\theta = \theta(X)$. Definitions become a bit more involved when the characteristic of interest is not always uniquely defined, see Section 2.7.

²Truth-telling applies the concept of incentive compatibility to belief elicitation. Kothiyal *et al.* (2011) point out that the earliest incentive compatible scoring rules were proposed at least a decade before the first work on mechanism design was published.

truth-telling rule for θ for subjects that have utility belonging to U if

$$\{\theta(X)\} = \arg \max_{r \in \Theta} Eu(S(r, X)) \text{ for all } u \in U \text{ and all } P_X \in \mathcal{P}_X.$$

We say that θ can be elicited for the subjects with utility belonging to U if there is a scoring rule S that is a truth-telling rule for U .

2.2 Proper scoring rules: Elicitation when subjects are risk neutral

For the special case in which the subject is assessed to be risk neutral (i.e. $U = \{Id\}$) it turns out the most common characteristics of a belief distribution can be elicited. Starting at least with Winkler and Murphy (1968) the literature refers to rules that are truth-telling for risk neutral subjects as *proper scoring rules* (PSRs). Below, we characterize the general mathematical representations of such rules and consider specific examples. Unless stated otherwise, the methods below will not be truth-telling if the subject is not risk neutral.

Lambert *et al.* (2008) give a general characterization of what can be elicited using a single realization when the subject is risk neutral. Assume that the experimenter wishes to elicit θ where $\theta = \theta(X)$ is continuous and not constant on any open neighborhood. Then θ can be elicited if and only if θ^{-1} is convex. This implies that the mean, any moment and any quantile is elicitable, but that the variance cannot be elicited (unless more reports or realizations are available, see Section 2.7).

2.2.1 Eliciting probabilities of events

We outline two general representations of proper scoring rules for the elicitation of probabilities of events, due to Savage (1971) and Schervish (1989). Below, we use the representation of Schervish to justify the quadratic rule, and the representation of Savage to justify the logarithmic rule.

Savage (1971, Section 6.1) gives the following general characterization. Given report r , assume that the payment to the subject equals $Y(r)$ if the event occurs and $Z(r)$ if it does not occur. So $ES(r, p) = Y(r)p + Z(r)(1 - p)$. Then S elicits the probability of an event if and only if $J(p) = ES(p, p)$ is strictly convex in p and the graph of $r \rightarrow ES(r, p)$ is tangent to the graph of $r \rightarrow ES(r, r)$ for all p . In particular, $J(p)$ is differentiable almost everywhere with $J'(p) = Y(p) - Z(p)$ as it is strictly convex.

Schervish (1989, Theorem A9) gives the following characterization. Let S be a scoring rule for eliciting the probability of an event where $S \geq 0$ and S is continuous at the boundaries $\{r, x\} \in \{0, 1\}^2$. Then S is strictly proper if and only if there exists a nonnegative measure ν with at most countably many point masses that assigns positive measure to every open interval

such that³

$$\begin{aligned} S(r, 1) &= S(1, 1) - \int (1 - c) \mathbf{1}_{\{r \leq c\}} \nu(dc) \\ S(r, 0) &= S(0, 0) - \int c \mathbf{1}_{\{r > c\}} \nu(dc). \end{aligned}$$

Note that any convex combination of two proper scoring rules and a positive affine transformation of a proper scoring rule is also a proper scoring rule (see Gneiting and Raftery, 2007, Section 6).

We now consider a number of applications of this general framework. We assume that there are n possible outcomes or events E_i , indexed by $i = 1, 2, \dots, n$, with associated reports $r_i \geq 0$ that satisfy $\sum_i r_i = 1$.

The quadratic scoring rule. By far the most used and well-known rule to elicit probabilities is the quadratic scoring rule (QSR). This rule is based on the Brier score (Brier, 1950), which is simply the sum of the squared errors of the reported probabilities. The QSR is an affine transformation of the Brier score such that when event j occurs the QSR pays

$$S(\mathbf{r}, x) = a + b \left(2r_j - \sum_{i=1}^n r_i^2 \right) \quad (1)$$

where $a > 0$ and $b > 0$ can be set by the experimenter. A common practice is to set $a = b$, so payoffs fall in the range $[0, 2a]$. Payoffs under the QSR thus decline in the square of the distance between the specified probability and the actual outcome. This rule is strictly proper for $\theta = (\Pr(X \in E_i))_{i=1}^n$.

When there is a single event to be elicited, one can implement the same incentives as the QSR using the following mechanism. The subject faces a price q that is continuously increasing, starting at 0. At each price q the subject is offered 2 units of probability of obtaining prize $y = 1$ if the event occurs. The subject is asked to report the price r at which she wants to exit the market and stop buying. So if the subject exits at r the probability of getting prize is $\int_0^r 2dq = 2r$, while her payment is $\int_0^r 2qdq = r^2$. In other words, payoffs are $S(r, 1) = 2r - r^2$ if the event occurs and $S(r, 0) = -r^2$ if it does not.⁴

³ $\mathbf{1}_{\{r \leq c\}} = 1$ if $r \leq c$ and $= 0$ if $r > c$. $\mathbf{1}_{\{r > c\}}$ is defined similarly.

⁴If one additionally gives out the prize whenever the event does not occur, then these payoffs are equal to the QSR with $a = b = 0.5$.

The spherical scoring rule. The spherical scoring rule, due to Roby (1964), is a strictly proper scoring rule given by

$$S(\mathbf{r}, x) = -\frac{r_j}{\sqrt{\sum_{i=1}^n r_i^2}} \quad (2)$$

Thus, this rule pays according to the probability reported for the event that occurred relative to the reported probabilities for the other events. Selten (1998) provides a proof that this rule is strictly proper.

The logarithmic scoring rule. The logarithmic scoring rule (Good, 1952; Toda, 1963) is given by

$$S(\mathbf{r}, x) = -\ln r_j. \quad (3)$$

This rule is strictly proper and has the appealing property that it depends only on the probability assigned to the correct answer, not on those assigned to the other, incorrect answers (i.e. $S(r, i) = f(r_i)$). It turns out that the logarithmic scoring rule is the unique proper scoring rule that has this property amongst differentiable scoring rules.⁵

Note that when an event occurs that the subject predicted to be impossible ($r_i = 0$) the score is $-\infty$. Thus, the rule is unbounded and needs to be truncated for experimental practice (Shuford *et al.*, 1966), but will no longer be strictly proper after such a truncation. Furthermore, Selten (1998) criticizes the logarithmic rule and shows that the score is at the same time very sensitive to small mistakes for small probabilities and insensitive to the distance from the truth for predictions $r_i = 0$.

Certainty equivalents. De Finetti (1970) and Savage (1971) note that for a risk neutral subject a probability is also a marginal rate of substitution: the rate at which she is willing to trade a probabilistic for a sure payoff. This opens another avenue to probability elicitation that is based on the a reservation-price elicitation mechanism of Becker *et al.* (BDM, 1964). In the so-called “promissory note” method, described in De Finetti (1974, but see also Ramsey (1926)), the experimenter asks the subject to report the lowest price r she would be willing to pay to acquire a prospect y_{Eg} (i.e. a lottery that pays y if event E occurs and g otherwise). Typically $g = 0$, so that the lottery simply pays y if the event occurs.

To determine the payment given report r , a number z is randomly drawn as a realization

⁵McCarthy (1956) mentions this claim and attributes it to Gleason (unpublished). Since we were unable to locate the latter study, we prove this here for $n = 2$, using the framework of Savage (1971). We search for a rule such that $Y(r) = Z(1-r)$ for all r , and hence $J'(p) = Y(p) - Z(p) = Y(p) - Y(1-p)$. Since $J(p) = Y(p)p + Y(1-p)(1-p)$ we obtain $J'(p) = Y(p) - Y(1-p) + Y'(p)p - Y'(1-p)(1-p)$ and hence $Y'(p)p = Y'(1-p)(1-p)$ for all p . This implies that $Y(p) = a \ln p + b$ for some $a > 0$ and b and hence S is an affine transformation of the the logarithmic payment scheme.

of a random variable Z that has distribution P_Z with support $[0, \infty)$. The subject receives the prospect y_{EG} if $z \leq r$, and payment z otherwise. If the subject is risk neutral, it is optimal for her to state her true certainty equivalent (CE) of the lottery: if she reported $r < \text{CE}$ she would be worse off when z falls in $[r, \text{CE})$. Similarly, reporting $r > \text{CE}$ yields lower payoffs than telling the truth when z is in $(\text{CE}, r]$. For a risk neutral subject, the elicited certainty equivalent can be used to calculate the probability as the solution to $yp + (1 - p)g = r$.⁶

2.2.2 Eliciting means

A general analysis of proper scoring rules for the mean can be found in Savage (1971, Sections 6.2, 6.3), and is similar to that for eliciting events. An popular rule is an adaptation QSR, which when applied to elicitation of the mean is given by

$$S^{QSR}(r, x) = a - b(r - x)^2 \quad (4)$$

where $a, b > 0$.

2.2.3 Eliciting the mode

To elicit the mode of a discrete distribution, it suffices to reward the subject for predicting the event that actually occurs (Hurley and Shogren, 2005). This method is truth-telling independent of the specific form of the utility, as it is optimal to specify the highest-probability event as long as the reward is strictly preferred to no reward. When the distribution is continuous one needs to elicit an interval, as explained below.

2.2.4 Eliciting quantiles and the median

One way to get a good idea of a cumulative distribution without eliciting the entire distribution is to elicit quantiles (Jose and Winkler, 2009). We call x a quantile α of the cdf F if $F(x) = \alpha$.⁷ Cervera and Muñoz (1996) present a general scoring rule for the elicitation of quantile $\alpha \in (0, 1)$, which is given by

$$S^\alpha(r, x) = \alpha r - (r - x)\mathbf{I}_{\{r \geq x\}}. \quad (5)$$

⁶Rather than asking for the lowest price that the subject is willing to pay for prospect y_{EG} , this mechanism can also be implemented by letting subjects complete a menu list of choices between a sure amount q and the prospect y_{EG} , where q is increasing for each choice. At the end, one decision is randomly selected for payment. The certainty equivalent is the value of q where the subject switches from the lottery to the sure amount.

The mechanism can also be presented as a scoring rule. Let $u(z)$ be the utility of prize z . Then $S(r, 1) = P(Z \leq r)u(y) + \int_r^\infty u(z)dP_Z(z)$ and $S(r, 0) = P(Z \leq r)u(g) + \int_r^\infty u(z)dP_Z(z)$ so that $ES(r, X) = P(Z \leq r)[pu(y) + (1 - p)u(g)] + \int_r^\infty u(z)dP_Z(z)$.

⁷If $F(x)$ can have discontinuities the general definition is $F(x_1) < \alpha$ for all $x_1 < x$ and $F(x) \geq \alpha$.

This rule rewards a high report, but punishes the subject if the report exceeds the realization, and is strictly proper for risk neutral subjects. Obviously, the median can be elicited by setting $\alpha = 0.5$.

2.2.5 Eliciting confidence intervals

Winkler and Murphy (1979) present a strictly proper scoring rule to elicit a $\alpha \cdot 100\%$ confidence interval. It is in fact a ‘double’ version of the quantile scoring rule discussed above, which elicits the $\frac{\alpha}{2}$ and $\frac{1+\alpha}{2}$ quantiles. The rule requires the subject to specify an upper bound u and a lower bound l . Here we present an affinely transformed version of this rule which is given by

$$S^{Int}(l, u, x) = -\frac{(1-\alpha)}{2}(u-l) - (l-x)\mathbf{I}_{\{x \leq l\}} - (x-u)\mathbf{I}_{\{x \geq u\}}. \quad (6)$$

In words, this rule rewards a small interval and punishes the subject for the distance of x from the interval bound if x falls outside of the interval. Schmalensee (1976) presents a similar strictly proper scoring rule that adds to (6) an extra term $|x - \frac{l+u}{2}|$ that penalizes the subject for the distance of the realization x from the mid-point of the interval. Schlag and van der Weele (2012) point out that these rules do not necessarily elicit the mode, nor the events that the subject thinks are most likely to occur. As a result they are ‘imprecise’, in that the chosen interval is often larger than necessary to cover $\alpha \cdot 100\%$ of the mass (see also Section 2.7).

2.2.6 Eliciting continuous density functions

Matheson and Winkler (1976) shows that the quadratic, spherical and logarithmic scoring rules can be modified to generate strictly proper scoring rules for density function elicitation. For example, the continuous quadratic scoring rule

$$S(r(x)) = 2r(x) - \int_{-\infty}^{\infty} r^2(x)dx, \quad (7)$$

is strictly proper for the density function.

An operational way to elicit a density function is to discretize it and elicit probabilities for subsets of outcomes (Harrison *et al.*, 2013a). The elicitor may want to fit a distribution to these points ex-post (Garthwaite *et al.*, 2005). Another approach is to let the subject choose from a limited number of distributions the one which best approximates the true distribution. This implies some discrepancy between the true and elicited distribution. Friedman (1983) discusses the design of scoring rules to minimize this discrepancy.

2.3 Selecting among proper scoring rules when subjects are risk neutral

When the experimenter faces a risk neutral subject and there are multiple proper scoring rules for the same characteristic, the experimenter needs to decide which rule to use. This selection may

be based on practical considerations. For example, she may select a rule that never involves payments from the subject to the experimenter, a rule with an upper bound on the possible payoffs, or one that is easier to explain to subjects (see Section 4).

From the perspective of decision theory, there are some arguments in favor of the QSR. Selten (1998) proves that the QSR obeys appealing axioms of invariance to superficial changes to the elicitation environment. Here we present a new justification for the QSR, as it puts a lower bound on the strength of the incentives for truth-telling close to the true report. The strength of these incentives is a relevant criterion when the subject can exert effort to avoid mistakes or gather additional information.

Incentives to tell the truth can be measured by the slope of the score S around the true report. Since any affine transformation of a PSR produces a new PSR with different incentives, we first fix a range of payments $[\omega_1, \omega_2]$ and compare rules which have payments in this range. We use the characterization of Schervish (1989) to prove the following (the proof is in the Appendix).

Proposition 1 *Consider a scoring rule S for the elicitation of an event which has a reward in the interval $[\omega_1, \omega_2]$ and admits a piecewise continuous density in the Schervish representation. Consider S^{QSR} as in (4) with $a = \omega_2$ and $b = \omega_2 - \omega_1$.*

If $S \neq S^{QSR}$, there exists $\varepsilon > 0$ and a distribution F_X with mean EX such that if $|EX - r| < \varepsilon$ then $\left| \frac{dES(r, X)}{dr} \right| < \left| \frac{dES^{QSR}(r, X)}{dr} \right| = 2b|EX - r|$.

In words, the proposition states that any other scoring rule will have strictly weaker incentives for truth-telling than the QSR for at least some distributions. The result extends immediately to the elicitation of probabilities as the probability of an event is the mean of particular random variable.

2.4 Incentivized elicitation when subjects are not risk neutral

The proper scoring rules described above rely on assumptions about the preferences and the rationality of subjects. Perhaps the strongest of those assumptions is that the subject is risk neutral. Countless laboratory studies have shown that most subjects behave as if they are risk averse over the stakes normally used in experiments (e.g. Holt and Laury, 2002). If subjects are not risk neutral, the scoring rules presented above are no longer truth-telling (Winkler and Murphy, 1970). For example, under the QSR a risk averse subject should submit reports that are biased away from extreme outcomes in order to minimize losses. Offerman *et al.* (2009) and Armantier and Treich (2013) provide evidence that subjects incentivized by the QSR do indeed report beliefs that are consistent with such a strategy.⁸

What are the alternatives to assuming risk neutrality? First note that if the utility function of the subject is known, one can offset risk aversion by paying the subject in utils rather than

⁸Harrison *et al.* (2013a), however, show that risk aversion poses less of a problem for the QSR when it is used to elicit the distribution of a continuous event rather than a binary probability.

money. Thus, $u^{-1}(S)$ is truth-telling rule for u if S is a proper scoring rule (Winkler, 1996). In the more typical case where the utility function is not known, Schlag and van der Weele (2013) show that truth-telling rules for the probability of an event with deterministic payments do not exist. In particular, none of the rules discussed in the previous subsection (including the QSR) are truth-telling. One possibility for the experimenter is to abandon the notion of truth-telling rules, a possibility we discuss in Section 2.7. The second is to amend the scoring rules above in one of following ways.

2.4.1 Paying small stakes

Ramsey (1926) suggests to minimize distortions arising from risk preferences by paying small stakes. However, Armantier and Treich (2013) prove that this does not necessary solve the problem. Paying small stakes only reduces biases when subjects display increasing relative risk aversion, and worsens it for decreasing relative risk aversion. The authors find evidence for increasing relative risk aversion in an experiment, where biases found under elicitation with the QSR are significantly smaller when payoffs are low. Thus, paying small stakes for belief elicitation, which is the practice in economic experiments anyway, goes some way towards inducing risk neutrality. However, it does not eliminate the problem completely, as Holt and Laury (2002) finds substantial risk aversion even for low levels of incentives. Moreover, it may undermine the benefits of incentivized elicitation that motivated its use in the first place.

2.4.2 Randomized payments

Another option is to fix a single (monetary) prize, and let the scoring rule determine the probability of winning the prize. A subject with any risk preferences just wants to maximize the probability of getting the prize if they prefer the prize over getting nothing. Since expected utility is linear in probabilities, this procedure induces risk neutrality at least in theory. We now discuss two implementations of this idea.

Paying in lottery tickets. One idea, due to Smith (1961)⁹, and implemented in the context of belief elicitation by e.g. McKelvey and Page (1990), is to replace the deterministic rewards for an accurate guess with a probabilistic reward. Harrison *et al.* (2013b) and Hossain and Okui (2013) consider a lottery version of the QSR, which Hossain and Okui (2013) label the ‘binarized scoring rule’ (see Table 1 and Section 4.3). Schlag and van der Weele (2013) show how to generally apply the randomization of payoffs to the deterministic scoring rules discussed above. By appropriately normalizing the probability of winning the prize, all the proper scoring rules discussed in Section 2.2 that have bounded payoffs (so excluding the logarithmic rule)

⁹There seems to be some confusion about the origin of this idea. Smith (1961) says the idea is ‘adapted from Savage (1954)’, but Savage (1971) attributes the idea to Smith.

can be transformed into randomized rules that are truth-telling for all risk preferences.¹⁰ For instance, assume that $S \in [\omega_1, \omega_2]$. Then one can give the prize to the subject with probability $\frac{S(r,x)-\omega_1}{\omega_2-\omega_1}$. If $\omega_1 = 0$ it is as if one gives the subject S lottery tickets from a total set of ω_2 lottery tickets.

Reservation probabilities. One can use a variation of the reservation-price elicitation mechanism of Becker *et al.* (BDM, 1964) where the subject is asked when she is indifferent between two lotteries, rather than between a lottery and a sure amount as in the elicitation of certainty equivalents. The mechanism appears to have been invented by Ducharme and Donnell (1973) and variations have been proposed by Grether (1981), Allen (1987), Holt (2006) and Karni (2009). As in the previous mechanism, the report determines the probability of winning a fixed prize, inducing risk neutrality.

The experimenter asks the subjects to report the lowest probability r such that she is indifferent between a prospect y_r0 (i.e. a lottery which pays y with probability r and 0 with probability $1 - r$) and the prospect y_E0 . For payment, a number z is chosen according to a random variable Z that has distribution P_Z with support on $[0, 1]$. The subject receives y_z0 if $z > r$ and y_E0 otherwise. As in the elicitation of certainty equivalents, reporting the true subjective probability p maximizes expected utility. However, in this case, the results holds for any utility function since the subject only compares the probabilities of winning the prize.¹¹

2.4.3 Estimating deviations from risk neutrality

The methods in the previous subsection rely on theoretical assumptions about the ability to induce risk neutrality. Alternatively, the experimenter can estimate deviations from risk neutrality (and expected utility maximization, see below) on the basis of additional reports, and use these estimates to correct the elicited beliefs. This is the strategy proposed by Offerman *et al.* (2009) and Andersen *et al.* (2013).

Offerman *et al.* (2009) propose to first elicit beliefs about events with known objective probability p . The elicited report function $R(p)$ shows how reports are biased away from the true beliefs due to deviations from risk neutrality. Subsequent reports r about other events with unknown subjective probability, can then be matched to the identical report $R(p)$, for which the underlying objective belief is known, and the true subjective belief p is recovered by inverting R , i.e. $p = R^{-1}(r)$.

To be able to match two identical probabilities, the experimenter needs to elicit (or otherwise be able to approximate) the correction function R^{-1} for the relevant range of p . Indeed, using the

¹⁰All definitions above immediately extend to randomized payment schemes, where the payment to the subject is a realization of some random variable. Here $S : \Theta \times \mathcal{X} \rightarrow \Delta\mathbb{R}$ where $\Delta\mathbb{R}$ denotes the set of distributions over \mathbb{R} .

¹¹Again, this method can implemented with a menu list, see Footnote 5. A problem arises when subjects do not have an unique switching threshold. Heinemann *et al.* (2009) exclude such subjects.

QSR to elicit beliefs about the throws of two 10 sided dice, Offerman *et al.* (2009) approximate the correction function, and provide evidence for substantial deviations from risk neutrality (and expected utility, see below). Note that this method does not require any structural assumptions or estimations, but uses actual reports to correct for risk aversion. Because of the investment required to obtain correction function, the authors argue that the method is most attractive if subjective beliefs are elicited about a substantial number of events. By contrast, when only few beliefs are elicited the use of randomized payoffs (as explained above) may be preferable.

Andersen *et al.* (2013) present subjects with a range of bets and use maximum likelihood to jointly estimate the risk preferences and subjective beliefs, assuming a structural form for the subject's utility function and the decision making model. For this method, a considerable amount of data is required to estimate beliefs with some degree of confidence, which is time consuming to collect. A less resource intensive approach is to estimate an average subject's beliefs and utility functions, where these can be conditioned on background characteristics like gender and age.

2.5 Scoring rules for non-expected utility maximizers

When subjects do not maximize expected utility, only a few of the mechanisms discussed above retain their truth-telling properties. One option is to pay only if the subject provides a correct guess of the event that will occur. Provided a positive payment is preferred to nothing, this method is incentive compatible for eliciting the mode of a belief distribution, and can be extended to elicit bounds for probabilities and quantiles (see Section 2.7). The elicitation of reservation probabilities is valid under the same assumption.

Several authors have explored modifications of the QSR to make it compatible with more general assumptions about the decision making process of the subject. These assumptions typically contain 'probabilistic sophistication', which allows for probability weighting but retains the assumption that beliefs are a probability measure (Machina and Schmeidler, 1992). Harrison *et al.* (2013b) and Hossain and Okui (2013) show that the lottery version of the QSR (see Section 2.4.2) elicits truthfully in the context of a rank-dependent utility model. Andersen *et al.* (2013) estimate the subjective probabilities assuming a similar model, using the procedure described in Section 2.4.3.

Offerman *et al.* (2009) show that the correction procedure to the QSR discussed in Section 2.4.3 can recover true subjective beliefs under probabilistic sophistication. However, they find that the correction curve is not invertible since subjects report a probability of 0.5 for a range of intermediate objective probabilities. This implies that subjective beliefs cannot be recovered accurately for reports of 0.5. To address this problem, Offerman and Palley (2013) consider the QSR when subjects are loss averse and form endogenous reference points. They show that loss aversion can explain the over-reporting of 0.5, and provide a version of the QSR that corrects for this by underweighting payoffs that are perceived as losses. In an experiment, they find that

this scoring rule leads subjects to correctly reproduce induced objective probabilities, obviating the need to elicit a correction function.

Kothiyal *et al.* (2011) also investigate the over-reporting of 0.5, and consider the performance of a more general set of scoring rules under assumptions allowing for non-additive beliefs. The authors explain the bunching at 0.5 as the result of the reversal of payoff ranks at a belief of 0.5 and propose a *comonotonic* scoring rule that preserves the rank order of payoffs.

2.6 Other elicitation mechanisms

The linear scoring rule. The linear scoring rule has the following score when event j occurs

$$S^{Lin}(\mathbf{r}, x) = r_j. \tag{8}$$

If the subject is approximately risk neutral, she has an incentive to report a probability of 1 for the event she thinks is most likely to occur, as this report has the highest marginal utility. Note that this scoring rule is truth-telling if the subject has logarithmic utility, since this yields the same preferences over outcomes as a risk neutral subject incentivized with the logarithmic scoring rule.

Elicitation games. In elicitation games, incentives for truthful reporting are based on the strategies of others rather than on the realization of a random variable. Perhaps the earliest elicitation mechanisms is the fair betting game, proposed by Toda (1951, see also Vlek (1973a)). The first player proposes a distribution of the total stake over two sides of a bet on the outcome of an uncertain binary event. The other player then chooses which side of the bet to take. In order to avoid ending up with the inferior side of the bet, it is optimal for the first player to make both bets equally attractive to the second player. Specifically, if the first player is risk neutral and believes that the second player has the same subjective probability about the outcome of the event, the proposed distribution of stakes reflect her true belief p . Perhaps these two, rather strong assumptions explain the limited use of this mechanism in the literature (although see Vlek, 1973b).

Several papers consider elicitation where scoring is based on the reports of others. These papers derive truthful elicitation as a Bayesian Nash equilibrium where it is optimal for each subject to report truthfully, as long as others do so. Prelec (2004) proposes to elicit both a subjective belief and an expectation of the frequency of all beliefs in the population. The latter is incentivized on the basis of the empirical distribution of beliefs, while the reported subjective belief receives a score that is proportional to the difference between its actual frequency and the average frequency predicted by the population. Thus, the subject has an incentive to report a belief that is “surprisingly common”. Since a rational (Bayesian) person expects others to underestimate the prevalence of her own belief, it is optimal to report truthfully.

Miller *et al.* (2005) propose a mechanism where each subject is asked to predict the report of another subject after having received a signal about the state of the world. The elicitor, who is assumed to know the common prior, applies a proper scoring rule to the posterior probabilities of the subject that are implied by the subject's report. If signals about the state are correlated across subjects, truthful reporting is a strict (but not unique) equilibrium.

Prediction markets. Recent work has looked at popular (online) betting schemes known as prediction markets. In such markets, people trade claims that pay conditional on the occurrence of some (often political) event. A belief in efficient markets could lead one to think that the market price reflects the average belief in the market. However, Manski (2006), assuming risk neutral traders who take prices as given, shows that market prices do not pin down mean beliefs, although they do put a bound on it. Wolfers and Zitzewitz (2006) show that markets do reflect mean beliefs under particular assumptions about risk aversion and independence between wealth and beliefs. Using simulations, Fountain and Harrison (2011) show that prediction markets will not generally reflect mean prices when wealth, discount rates or risk aversion are correlated with beliefs.

2.7 Extensions

Here we discuss some extensions to the standard framework that we believe are promising avenues for further research. Some of these extensions abandon the notion of truth-telling rules, and explore the trade-off between the complexity of the mechanism and the precision of the reported beliefs.

Eliciting sets. When truth-telling schemes do not exist, or are deemed too complicated for implementation, one may wish to elicit a set A that contains the characteristic of interest θ . To simplify notation, consider the case where θ is a real number, so $\Theta \subseteq \mathbb{R}$. Formally, the scheme $S = S(A, x)$ depends on the reported set $A \subset \Theta$ of the subject and on the realization x of X . In this case we say that S is *compatible with the truth* if

$$\theta(X) \in \arg \max_{A \subset \Theta} Eu(S(A, X)) \text{ for all } u \in U \text{ and all } P_X \in \mathcal{P}_X.$$

An example of this approach is Schlag and van der Weele (2012). The authors consider a random variable with support on $[a, b]$ and ask the subject to state an interval $[l, u]$. The score is given by

$$S^T(l, u, x) = \begin{cases} \left(1 - \left(\frac{l-u}{b-a}\right)\right)^{\frac{1-\alpha}{\alpha}} & \text{if } x \in [l, u] \text{ and } \frac{l-u}{b-a} \leq \alpha, \text{ and} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

The rule rewards the subject if x is in the stated interval, where the reward declines in the width of the interval. The authors show that if the subject is (weakly) risk averse, the optimal interval

contains the mode as well as a α -100% confidence interval for a realization of X . In addition, the rule truncates the reward when the interval is wider than necessary to cover α -100% of the mass, improving the precision of the rule relative to other scoring rules for confidence intervals discussed in Section 2.2.5.

Constructing bounds. Another approach to deal with uncertainty about the subjects' preferences is not to focus on truth-telling schemes but instead to take misreporting into account as follows. One asks the subject to report the characteristic of interest and then derives all θ 's that could lead to this report for some utility function u . Formally, for any given scheme $S = S(r, x)$ and any report r one defines $A(r)$ by

$$\bar{\theta} \in A(r) \text{ if and only if } \bar{\theta} \in \arg \max_{r \in \Theta} Eu(S(r, X)) \text{ for some } u \in U \text{ and some } P_X \in \mathcal{P}_X \text{ with } \theta(X) = \bar{\theta}.$$

The set $A(r)$ is then the inference about θ one obtains from scheme S . Hurley and Shogren (2005) provide an example of this approach based on payments of multiple events, as discussed below.

Paying on the basis of multiple events. One may wish to increase the number of independent events that are used to pay the subject. Sometimes this is necessary to elicit particular objects. For instance, one needs two independent realizations x_1 and x_2 of X to elicit the variance of X when the subject is risk neutral (Lambert *et al.*, 2008). The payment scheme is given by $S(r, x_1, x_2) = a - \left(r - \frac{1}{2}(x_1 - x_2)^2\right)^2$ (see Schlag and van der Weele, 2013).

At other times, multiple events occur naturally in an experiment, for example when there is a group of people independently making the same choice. Hurley and Shogren (2005) provides a mechanism that asks for the empirical frequency of these choices. They show that from this report, the experimenter can recover an interval that contains the true subjective probability for each individual event regardless of risk preferences. These intervals can achieve reasonable precision. For example, if one can pay based on 20 realizations then one can derive an interval of less than 0.05 containing the true belief. This approach is elaborated upon by Schlag and Tremewan (2014), who also demonstrate how conditioning on multiple events can be used to elicit the median and quantiles.

Multiple reports. One may also choose to increase the number of reports. For instance, one can elicit the variance when the subject is risk neutral by asking for two reports. Report r_1 is used to elicit the expected value of X^2 , report r_2 to elicit the expected value of X . Then $r_1 - (r_2)^2$ elicits the variance of X if the subject is risk neutral. Note that EX^2 can be elicited when the subject is risk neutral by applying the quadratic scoring rule to x^2 , so using $S(r_1, x) = a - (r_1 - x^2)^2$.

3 An empirical comparison of elicitation mechanisms

In this section we compare the empirical performance of the methods discussed above. We also include comparisons with what we label “introspection”, i.e. simply asking for a probability or characteristic without the use of an incentive compatible elicitation method. Introspection may or may not be rewarded with a payment that is independent of the report. Also, in many studies incentive compatible mechanisms are implemented in an unincentivized way by using hypothetical payoffs.¹²

The central challenge in evaluating the performance of an elicitation mechanism is the fact that the true subjective belief, with which we wish to compare the elicited report, is unobservable. In Section 3.1 we discuss several approaches to deal with this problem. In Section 3.2 we discuss the results and make suggestions for future work.

3.1 Benchmarks for evaluating elicitation mechanisms

In the absence of knowledge of the true subjective beliefs, several benchmarks have been proposed in the literature to evaluate elicited reports. Table 1 summarizes an exhaustive list of the results and methods of individual papers on this issue. Note that in this section we only consider benchmarks that are external to the elicited beliefs. This excludes papers that compare beliefs elicited with different methods either within-subject (e.g. Beach and Phillips, 1967) or between-subject (e.g. Andersen *et al.* (2013)). At best, such comparisons show that at least one report is mismeasured.

3.1.1 Induced beliefs

The most prominent benchmark is to induce beliefs in subjects and to investigate whether these beliefs are then reported accurately under a given elicitation mechanism. Perhaps the simplest way of inducing beliefs is to inform subjects directly about the objective probability of an event. The relevant studies are listed in Table 1, under “Induced Probabilities (Direct)”. For example, Hao and Houser (2012) show subjects the number of black and white chips in a bag, and then try to recover the probability with which a randomly drawn chip is believed to be of a given color. It seems reasonable to assume that the subjective (i.e. induced) probability should be equal to the true probability, and that the smaller the distance between the latter and the elicited probability the better the elicitation method.

¹²Vlek (1973b) points out that even with hypothetical payoffs these mechanisms may still matter because they encourage subjects to think in a particular way and may align the preferences of experimenter and subject. An appropriate feedback rule can clarify what the experimenter really wants to know and avoid wrong interpretations. One example of misinterpretations by the subject comes from the elicitation of confidence intervals. Yaniv and Foster (1995, 1997) show that subjects seem to think that 50% confidence intervals strike the right balance between accuracy and preciseness even when the requested level of confidence is much larger. This interpretation casts doubts on the widespread interpretation that intervals that are too narrow are a sign of ‘overconfidence’. Krawczyk (2011) shows that using incentives for interval elicitation improves the level of calibration of subjects. Winkler and Murphy (1968) provide a discussion of scoring rules as learning devices.

There are two potential problems with this approach. First, it is not clear that people respond to objective and subjective probabilities in the same way, especially if the subjective probability is something derived from a situation of strategic uncertainty. Second, this approach becomes trivial when comparing methods where the response must be in the form of a probability, such as introspection (“The probability A will occur is x . What is the probability A will occur?”). In this case, it is hard to exclude the possibility that subjects simply repeat the probability of which they have just been informed without paying attention to the incentives.

To address the second problem, a number of experimenters have attempted to induce beliefs by supplying sufficient information for subjects to calculate the true probability. One commonly used method is to describe two distributions (e.g. two bingo cages with different proportions of red and white balls), show a series of draws with replacement from one of the two distributions, and elicit the subjects’ posterior belief about the probabilities that the draws were from each of the distributions. The elicited beliefs can then be compared to probabilities calculated using Bayes’ Rule. Papers that use this technique are listed in Table 1 under “Induced Probabilities (Bayes’ Rule)”. A second possibility is ask for the probability a particular combination of events will occur, where the probability of each individual event is known (see Table 1, benchmark: “Induced Probabilities (Multiple Events)”).

TABLE 1 ABOUT HERE

A potential problem with these techniques is that they require computation on the part of the subjects. Subjects unfamiliar with probability theory are likely to use various heuristics to evaluate probabilities, possibly resulting in systematic differences between the true and induced probabilities. Such biases may augment or diminish biases resulting from the elicitation in unknown ways.¹³ In addition, incentives may encourage the subjects to put more effort in these calculations. A closer match between the objective and elicited probabilities may thus reflect an improvement in the belief induction rather than belief elicitation procedure, overstating the importance of incentives in belief elicitation.

3.1.2 Deviations from the empirical distribution

A number of papers use actual outcomes or frequency distributions as a benchmark for the quality of elicited beliefs (see Table 1, benchmark: “Empirical Distribution”). The most commonly used procedure is to elicit beliefs about the action of another subject and compare these probabilities to the empirical frequency in the treatment. Other options are to use distributions subjects should be somewhat familiar with (e.g. heights of males and females) or to give them limited exposure to the distribution in question (e.g. show a bingo cage with balls of different colors so the precise numbers can only be estimated).

¹³Hurley and Shogren (2005), for example, argue that their inability to recover induced beliefs with a belief elicitation procedure stems from a failure to induce correct beliefs rather than a failure of the elicitation process.

Of course there is no reason why subjects' beliefs should correspond to the objective probabilities. In fact, the idea that beliefs may be incorrect is one of the main motivating factors for economists to develop reliable methods of belief elicitation (see Manski (2004) for a criticism of rational expectations). Badly calibrated beliefs could lead to erroneous conclusions about the accuracy of elicitation methods. For example, when comparing a standard QSR, and a QSR corrected for risk attitudes, overconfidence bias when predicting the actions of others (leading to more extreme probabilities) would counteract the bias caused by a failure to account for risk aversion in the QSR, resulting in beliefs which are closer to empirical distribution but presumably further from subjective beliefs. Obviously, this criticism does not apply if the purpose of the elicitation is to obtain accurate predictions rather than obtaining the best measure of subjects' beliefs.

3.1.3 Consistency with behavior

Another benchmark that has been considered is the degree to which elicited beliefs are consistent with actions in games, see Table 1 under "Consistency". The idea is that under the assumption that subjective probabilities are crucial in determining choices, as is the case with most decision theories in economics, a stronger relationship between stated beliefs and choices indicates higher quality elicitation. This approach provides a natural testing ground because the relationship between beliefs and actions is precisely the context in which economists are often most interested in beliefs.

In order to check whether beliefs are in fact consistent with choices one must assume a model determining the relationship between beliefs and actions: the finding that elicited beliefs are not consistent with best-response behavior could be explained either by a failure of the elicitation method or of the assumption that subjects best respond to their beliefs. Moreover, if a standard economic decision making model is assumed, one must know the utility function of subjects. We are not aware of a paper comparing elicitation methods in this way that has made a serious attempt to measure for example risk or social preferences, although Trautmann and van de Kuilen (2014) provide consistency comparisons of choices in an ultimatum game on the basis of several alternative utility functions.

3.1.4 Additivity

A final option is comparing how close beliefs elicited using different methods are to satisfying additivity, i.e the condition that the probability of the union of mutually exclusive events is equal to the sum of the underlying probabilities. This implies for instance that the sum of the probabilities of exhaustive and mutually exclusive events add to one. If subjective beliefs are truly additive, then a necessary (although not sufficient) condition for a valid elicitation method is that stated beliefs satisfy additivity. The relevant studies appear in Table 1 under "Consistency".

3.2 Discussion of empirical comparisons

A simple overview of the existing studies, considering all approaches of comparison, gives little clue as to which method of elicitation is preferable as many results are contradictory. For example, Wang (2011) finds that the QSR results in stated beliefs closer to the empirical distribution than introspection, whereas Hollard *et al.* (2010) find the opposite, and Phillips and Edwards (1966)'s finding that the linear scoring rule is preferable to the QSR is contradicted by Palfrey and Wang (2009).

Opposing results could be due to many factors: the reliability of belief elicitation methods may depend on the domain in question (e.g objective/subjective probabilities); more complex methods may perform differently depending on the mathematical literacy of the subjects; different subject pools may have different risk attitudes; incentives may be more important for situations where there is a strong reason to misreport. This points to the necessity of comparing elicitation methods in the domain and with the subject pool for which they are to be used. They should also be rated based on criteria relevant to the purpose of the elicitation.

We can nevertheless draw one conclusion: differences in stated beliefs between methods tend to be consistent with the presence of risk-averse subjects. Jensen and Peterson (1973) and Armantier and Treich (2013) find that steeper incentives lead to less extreme reported probabilities, and Trautmann and van de Kuilen (2014), Offerman *et al.* (2009); Offerman and Palley (2013) and Harrison *et al.* (2014) all find evidence that correcting the QSR for risk aversion improves performance. This should come as little surprise as it has been long established that risk aversion is prevalent in typical subject pools (in fact, Seghers *et al.* (1973) called into question the validity of PSRs for this very reason).

Distortions in stated beliefs will not always be a problem, for example if one is only interested in establishing a difference in distributions of beliefs in two populations with identical distributions of risk preferences. In most cases, however, it appears that methods that are robust to risk preferences are to be preferred. This is especially crucial if a variable of interest is correlated with risk aversion, such as testing gender differences in beliefs.

There is also some evidence that incentives cause subjects to think harder or more systematically about the game, leading them to make more 'rational' decisions (Trautmann and van de Kuilen, 2014; Gächter and Renner, 2010). Hoffmann (2013) finds that subjects use dominated strategies less often when beliefs are elicited. However, these statements are not easy to make as they involve implicit assumptions about preferences, and more evidence is needed before any firm conclusions can be drawn.¹⁴

¹⁴For example, in the context of public goods games, a deeper understanding of the game may have very different implications for selfish individuals (who would reduce contributions) or altruistic individuals (who would increase contributions). Indeed, in the public good game of Gächter and Renner (2010), the interpretation that elicitation improves understanding rests on the assumption that people are conditional cooperators. Note that in this study, the statistical effect is weak and the results are also consistent with a consensus effect or the use of stated beliefs to justify (selfish) actions.

With respect to methodology, all the common approaches to comparing belief elicitation methods have serious drawbacks. We consider directly inducing probabilities the least problematic, as the other benchmarks assume away the very phenomena that we are most interested in when studying subjective beliefs (miss-calibration, overconfidence, bounded rationality). Reliable recovery of induced probabilities should be considered a necessary, but perhaps not sufficient, condition for a good mechanism.

Given the fundamental impossibility of using the correct benchmark, i.e. the true subjective belief, we suggest that the best approach is to take the theory seriously and test what is empirically verifiable. First of all, one can test the assumptions on which the validity of a mechanism is based. Paraphrasing Staël von Holstein (1970) (who in turn drew on De Finetti (1965)), several explicit and implicit assumptions underlying elicitation mechanisms can be enumerated.

1. The method must be incentive compatible, i.e. all the assumptions from which incentive compatibility is theoretically derived must hold. Typically this simply means that subjects must have a particular utility function, and behave according to a particular decision theory.
2. Subjects must understand the implications of the elicitation method.
3. Subjects must understand the correspondence between their own beliefs and the probabilities (numerical or graphical) into which they are to be translated.

There is a wealth of experimental evidence related to the first point, with many decision theories proposed and tested (see, for example, Harrison and Rutström (2009)). Point 2 can also be tested in the lab, and often is to some extent in the form of comprehension pre-tests. Point 3 has been studied in the psychology and medical literature. These last two points will be discussed in the next section. If we can ascertain that all the assumptions necessary for a method to reveal true subjective probabilities have been met, comparing results to a contestable benchmark becomes less important.

4 Issues of implementation

Here, we discuss some more practical issues that experimenters will face when eliciting beliefs.

4.1 The order of eliciting beliefs and actions

Experimental economists are typically not interested (only) in beliefs for their own sake, but also in their relationship with decisions, or using them in conjunction with decision data to help identify preferences or motivations. For these purposes it is necessary to elicit both actions and beliefs from the same subjects and a decision must be made as to which to elicit first, or whether

to elicit them simultaneously. This raises two questions: does the elicitation of beliefs affect the decisions subjects make, and does the elicitation of decisions affect beliefs?¹⁵

TABLE 2 ABOUT HERE

The main reason put forward for belief elicitation potentially affecting decisions is that it may deepen subjects' understanding of a situation and make them act in a more sophisticated fashion. Conversely, choosing an action could influence elicited beliefs through several channels: a consensus bias (people assume others will act in the same way as themselves); justification to oneself or the experimenter that an action was morally acceptable by demonstrating a belief that one's action conformed to the norm (or increasing self-esteem by believing that an action was exceptional); a need to convince oneself that the correct action was chosen by holding beliefs that are consistent with that action; or a salience bias which makes the chosen action seem more probable.

Table 2 summarizes all the relevant papers we are aware of. The evidence is scanty and contradictory. Taking an action has been found to increase and decrease the accuracy of elicited beliefs. With a similar degree of inconclusiveness, belief elicitation is found to decrease, increase, and have no effect on contributions in public goods games. Erev *et al.* (1993) find that eliciting beliefs about the probability of events diverts attention from the size of payoffs and reduces expected value maximization. Guerra and Zizzo (2004), on the other hand, find no effect of belief elicitation on trusting behavior. Hoffmann (2013) compares an action-only treatment with a treatment where beliefs are elicited simultaneously, and finds that belief elicitation makes subjects less likely to choose dominated actions.

The small number of studies and apparently contradictory results on the two related methodological questions discussed in this section make it hard to draw any strong conclusion. It seems that eliciting beliefs can have an effect on decisions, but the direction of an effect, and the circumstances under which it arises is unclear. As discussed above, there is some evidence that belief elicitation affects play by deepening the understanding of the game. Overall, if independently measured beliefs and decisions are required from the same subjects, we can only recommend testing for an impact of belief elicitation on decisions and vice versa whenever designing a new game or using a new subject pool.

4.2 Hedging

Experimentalists are often interested in eliciting both decisions and probabilities from the same subject. However, paying subjects for both actions and elicitation tasks that depend on the

¹⁵Most of the literature discussed so far is based on the decision theoretic approach by Savage (1954), where subjective utilities are a primitive concept used in evaluating uncertain prospects. In contrast, psychologists have argued that choices may affect beliefs. A discussion of the merits of these approaches is beyond the scope of this paper and we limit ourselves discussing the empirical effect of *elicitation* on responses. Costa-Gomes *et al.* (2012) and Smith (2013) use an instrumental variable approach to identify a causal relationship between beliefs and actions.

outcome of the same event creates a situation in which subjects have a stake in the outcome of variable they are asked to predict. Kadane and Winkler (1988) and Karni and Safra (1995) show that under such circumstances proper scoring rules for the probability of an event no longer exist.¹⁶ A specific concern in the context of experiments is that subjects have an incentive to hedge. For example, a risk averse subject facing the binary QSR who benefits when the event occurs, has an incentive to report an overly pessimistic belief in order to smooth her payoffs over the two states.

There is mixed evidence whether or not hedging plays a role in economic experiments. Blanco *et al.* (2010) look at several games and contrast a hedge environment where both beliefs and choices are paid, with a no-hedge environment where only one of those is paid randomly. They find that a sizable number of subjects hedge in a 2×2 coordination game where the opportunity is obvious, but not in a more complex sequential prisoners' dilemma. The authors also find that in the former case, some subjects play a best response against other players' hedging strategies.

The results from Armantier and Treich (2013) confirm that subjects may exploit obvious hedging opportunities. The authors used the QSR to elicit probabilities about events based on the roll of two dice. In a hedging treatment, subjects were also able to bet separately on the event in question. The authors find that subjects in the hedging condition bet more on the most likely events, and simultaneously report lower probabilities than in the control treatment. More circumstantial evidence comes from Palfrey and Wang (2009), which shows that observers with no stakes in the game and no incentive to hedge predict differently than subjects in the game.

Blanco *et al.* (2010) provide a list of precautions that the experimenter can take to avoid hedging. First, one can elicit beliefs not about the matched partner's behavior, but about average behavior of the subjects in the partner's role, or a particular non-matched subject in that role (Armantier and Treich, 2009). This reduces the correlation between payoffs from belief reports and outcomes of play, and reduces the value of the hedge. Second, one can decide to randomly pay either the reported belief *or* the payoffs obtained in the game. Third, one may not pay for elicitation at all, although this may aggravate other sources of misreporting. Fourth, Blanco *et al.* (2010) find that some subjects hedge when they should not, so it may be helpful to explain subjects when they should not hedge. Finally, post-experiment questionnaires about the reasons for play and belief reports may also help detect hedging.

4.3 Complexity of elicitation method

Some belief elicitation mechanisms require a high degree of mathematical sophistication (e.g. understanding the formulae of PSRs) or understanding relatively complex payment procedures (e.g. methods with probabilistic payoffs), and confusion among subjects has the potential to cause noise and bias in elicited beliefs (see Artinger *et al.*, 2010, for a discussion).

¹⁶Jaffray and Karni (1999) present mechanisms that can overcome these problems, which require either additional elicitation tasks, or the payment of very large sums of money to exploit the domain where the utility function is relatively flat.

In order to address the first problem, some experimenters present subjects with the formula for the PSR in question, assure them that stating their true belief will maximize the amount they can expect to earn, and offer a mathematical proof on request. A solution that is more in the spirit of revealed choice, is to have subjects select their preferred option from a list of bets generated using a scoring rule (e.g. Jensen and Peterson, 1973; Offerman *et al.*, 2009). With computers one can implement this easily by offering subjects sliders for setting the desired probability. When moving the sliders, the software can simultaneously display the payoffs associated with each outcome (e.g. Andersen *et al.* (2013)). Another possibility is to select a simple rule such as paying only if the correct event is reported.

Whether or not subjects understand probabilistic payoff schemes (Section 2.4.2) has been the subject of some debate (Berg *et al.*, 2008). Outside of a belief elicitation framework, Selten *et al.* (1999) casts doubt on the effectiveness of randomized payments in inducing risk neutrality. By contrast, in a very simple elicitation task with induced probabilities, Harrison *et al.* (2013b) and Hossain and Okui (2013) provide evidence that the use of probabilistic payoffs produces responses that are in line with risk neutral behavior, a finding which may or may not generalize to more complex environments.

4.4 Representation of probabilities

The format in which probabilities are communicated may matter, especially to subjects who are unfamiliar with them. Lipkus *et al.* (2001) find that people cannot in general convert between numerical probabilities, percentages, and frequencies, which suggests they are unlikely to respond the same way if asked for a subjective probability in different formats.

Some studies have addressed the question of which format subjects best understand. Tversky and Koehler (1994) find that probabilities more likely to be additive if elicited as percentages rather than numerical probabilities. Gigerenzer and Hoffrage (1995) find that subjects are better able to perform Bayesian updating when presented with frequencies rather than numerical probabilities. Price (1998) finds that eliciting probabilities as frequencies rather than numerical probabilities reduces the number of subjects expressing complete certainty, as well a measure of dispersion. This points to an advantage of eliciting beliefs as a population frequency as in Hurley and Shogren (2005) rather than as a single-event probability.

Probabilities can also be expressed graphically. Wang *et al.* (2002) find that the consistency of probabilities elicited at different times depends on whether they are elicited (from least to most consistent) as numbers, using a probability wheel, or a probability bar. Whitcomb *et al.* (1993) find no difference in consistency between elicitations as numbers, a probability wheel, or odds ratio.

4.5 Eliciting complementary events

It is common, especially with regard to binary events, to ask for the probabilities about all but one possible outcome, calculating the probability associated with the last outcome by assuming additivity of subjective beliefs. Given the evidence that subjective beliefs appear to be consistently super-additive (Tversky and Koehler, 1994; Trautmann and van de Kuilen, 2014), this is a questionable practice.

If non-additivity is a genuine feature of beliefs, we are forced to consider some theory and a payment scheme that allows for this possibility (e.g. Kothiyal *et al.*, 2011). However, it is also possible that super-additivity is an artifact of elicitation if asking for a particular event increases its salience and makes it appear more likely, inflating the probabilities of each event. In this case, one possibility is to elicit probabilities for all events (i.e. A and “not A” for binary events) and deal with the resulting, inconsistent probabilities. Either the experimenter can scale them in some way to have them add to one, or the subject can reconcile the probabilities themselves. There is a substantial literature on the reconciliation of inconsistent probability assessments, e.g. Lindley *et al.* (1979). Alternatively the elicitation can be done in such a way that the input must be consistent, where care must be taken not to make one outcome more salient (e.g. order of elicitation). This can be achieved with the use of sliders.

5 Discussion and Conclusion

From personal conversations with colleagues, we have come away with the impression that opinions on the use of incentives for belief elicitation are divided. Roughly speaking, theorists or experimentalists with a strong theoretical focus tend to argue that incentivized elicitation is essential to interpret the elicited data. On the other hand, a sizable number of experimentalists are favorably disposed to non-incentivized elicitation and are comfortable to rely on intrinsic motivation of the subjects to answer questions truthfully.

We believe that both sides have good arguments at their disposal, and that the relative strength of those arguments will depend on the context. Favoring the theorists, there are quite a few experimental situations where there are a priori reasons to assume that people may report falsely or sloppily. First of all, there are games in which subjects may use stated beliefs to justify their (selfish) behavior to the experimenter. This includes virtually all experiments which feature some trade-off between the payoffs of the decision maker and other subjects, such as dictator games, prisoner’s dilemmas, public good games and trust games. Although we have not found direct evidence of such distortions, the evidence on the existence of experimenter demand effects (Zizzo, 2009) makes us believe that they should be taken seriously.

Second, subjects may simply ‘click through’ belief elicitation questions without putting in any effort, especially when they are bored or tired (at the end of a long experiment) and the

questions are complex.¹⁷ Third, there is some evidence that the use of scoring rules improves understanding of the game (Hoffmann, 2013) and consistency of decision making (Trautmann and van de Kuilen, 2014). Incentives may thus reduce noise in experimental data, although the appropriate definition of ‘noise’ may depend on the aim of the study. A final reason to use scoring rules, but not necessarily incentives, is to clarify what is being elicited (see e.g. Footnote 11).

Balanced against these considerations are first and foremost are the practical costs of implementing and explaining incentive schemes. Trautmann and van de Kuilen (2014) compare the efforts required for different elicitation mechanisms in a table and show that these can be quite substantial. A second argument is that incentivization may create new distortions due to risk aversion or hedging. Note that there are trade-offs between these two arguments, as distortions due to risk aversion may be reduced by eliciting additional (and costly) reports.

We believe incentivized elicitation is generally advisable. However, the cost may exceed the benefits in situations where subjects are fresh, have no clear incentive to misreport, and face a straightforward elicitation task where the marginal benefit of subjects’ effort is low and hedging may be a problem otherwise.¹⁸ By contrast, using incentives is desirable in more complex or tedious tasks and when ruling out misunderstandings or careless reporting is especially important. An example of the latter are experiments testing cognitive biases. Engelmann and Strobel (2000), using incentivized belief elicitation, cannot reproduce the ‘false-consensus effect’ found in psychological studies that do not incentivize elicitation.

Suppose a researcher wishes to incentivize belief elicitation, which methods should she use? In answering this question we gather arguments from our discussion in the previous sections. One relatively clear result from the empirical literature is that risk aversion will bias beliefs elicited with proper scoring rules away from extreme probabilities. We would therefore recommend the use of corrective calibrations for such deviations, or the implementation of a randomized mechanism like reservation probabilities. Given the need for extensive additional reports, corrective calibrations are more useful when multiple beliefs are elicited.

There are several things the experimenter can do to ease the cognitive strain on subjects and encourage consistent reports. First, the use of sliders (Andersen *et al.*, 2013) is beneficial for several reasons: it obviates the need for displaying complex formulae; no mention of probabilities is required; additivity is ensured and each event is given equal prominence. Second, when multiple observations are available, e.g. when larger groups of subjects in the session make the same decision, one can elicit the belief about the empirical frequency of a decision rather than the probability of a decision for a single person. Empirical frequencies seem to be more

¹⁷In Offerman *et al.* (1996), 50% of the subjects indicate that they would have reported different beliefs in the absence of incentives, often deviating to an ‘easier’ report.

¹⁸Note that these conditions apply to most studies testing incentivized elicitation schemes, where belief elicitation is typically the only experimental task and thus receives full attention of the subjects. Therefore, these studies may understate effects of incentives in other, more complex, experimental settings (see the comments in Sonnemans and Offerman, 2001).

easily understood by subjects, and the procedure may also help to avoid hedging against the payoff-relevant decision of a single opponent.

What more do we need to know to conduct effective belief elicitation? The answer can be separated into a theoretical and empirical part. On the theoretical side, in Section 2.7 we have put forward extensions to the standard framework that await further elaboration. More generally, an area of research that has seen considerable activity in recent years is the use of methods that are robust to different assumptions on preferences (e.g. risk aversion) and decision making models (e.g. prospect theory). Given the evidence of heterogeneity in risk preferences and cognitive capacities, we believe this line of research should continue.

A second area concerns the trade-off between simplicity and informativeness. Given the limited time and resources that we can invest in belief elicitation, simple mechanisms that yield more imprecise information, for example by specifying bounds on beliefs, may be preferable to more complicated ones that yield very precise beliefs.

Ultimately, the benefits of different incentivization mechanisms should be determined by empirical evidence, of which there is too little at present to draw any but tentative conclusions. Above, we have emphasized research investigating the validity of assumptions underlying different elicitation mechanisms. This includes both fundamental research on the nature of subjective beliefs and the feasibility of inducing a particular objective function in experimental settings.

More concretely, we need to improve our understanding about the interactions between belief elicitation and game play. When eliciting both choices and beliefs, one should experiment with the order of elicitation or even the presence of belief elicitation mechanisms in at least some sessions. There is a public good aspect to this kind of research, as it will help future researchers to make more informed design choices. Another important question is the separation of the cognitive effect of elicitation mechanisms (i.e. scoring rules as learning devices) and the incentive effects. To this end, we recommend that researchers testing elicitation methods include treatments that implement the method with hypothetical payoffs. Finally, it would be valuable to test for the importance of experimenter demand effects, or justification of behavior through stated beliefs.

We hope that this paper will help experimentalists to make informed design choices and will provide inspiration for the development of new belief elicitation tools.

References

- ALLEN, F. (1987). Discovering Personal Probabilities When Utility Functions Are Unknown. *Management Science*, **33** (4), 542–544.
- ANDERSEN, S., FOUNTAIN, J., HARRISON, G. W. and RUTSTRÖM, E. E. (2013). Estimating Subjective Probabilities. *Journal of Risk and Uncertainty*, **In Press**.

- ARMANTIER, O. and TREICH, N. (2009). Subjective Probabilities in Games: A Solution to the Overbidding Puzzle. *International Economic Review*, **50** (4), 1079–1102.
- and — (2013). Eliciting beliefs: Proper scoring rules, incentives, stakes and hedging. *European Economic Review*, **In press**.
- ARTINGER, F., EXADAKTYLOS, F., KOPPEL, H. and SÄÄKSUORI, L. (2010). *Applying Quadratic Scoring Rule transparently in multiple choice settings: a note*. Tech. rep., Jena Economic Research Paper.
- BEACH, L. and PHILLIPS, L. (1967). Subjective Probabilities Inferred from Estimates and Bets. *Journal of Experimental Psychology*, **75** (3), 354–359.
- BEACH, L. R. and WISE, J. A. (1969). Subjective Probability Revision and Subsequent Decisions. *Journal of Experimental Psychology*, **81** (3), 561–565.
- BECKER, G., DEGROOT, M. and MARSCHAK, J. (1964). Measuring Utility by a Single-Response Sequential Method. *Behavioral science*, **9** (3), 226.
- BERG, J. E., RIETZ, T. A. and DICKHAUT, J. W. (2008). On the Performance of the Lottery Procedure for Controlling Risk Preferences. *Handbook of Experimental Economics Results*, **1**, 1087–1097.
- BLANCO, M., ENGELMANN, D., KOCH, A. K. and NORMANN, H.-T. (2010). Belief elicitation in experiments: is there a hedging problem? *Experimental Economics*, **13** (4), 412–438.
- BRIER, G. W. (1950). Verification of Forecasts Expressed in Terms of Probability. *Monthly Weather Review*, **78** (1), 1–3.
- CERVERA, J. L. and MUÑOZ, J. (1996). Proper Scoring Rules for Fractiles. In J. O. Berger, A. P. Dawid and A. F. M. Smith (eds.), *Bayesian Statistics 5*, Oxford, U.K.: Oxford University Press, pp. 513–519.
- COSTA-GOMES, M. A., HUCK, S. and WEIZSACKER, G. (2012). Beliefs and actions in the trust game: creating instrumental variables to estimate the causal effect. *WZB Discussion Paper*, **2012-302**.
- CROSON, R. T. A. (2000). Thinking like a game theorist: factors affecting the frequency of equilibrium play. *Journal of Economic Behavior & Organization*, **41** (3), 299–314.
- DAWES, R. M., MCTAVISH, J. and SHAKLEE, H. (1977). Behavior, communication, and assumptions about other people’s behavior in a commons dilemma situation. *Journal of Personality and Social Psychology*, **35** (1), 1.

- DE FINETTI, B. (1965). Methods for discriminating levels of partial knowledge concerning a test item. *British Journal of Mathematical and Statistical Psychology*, **18** (1), 87–123.
- DE FINETTI, B. (1970). Logical foundations and measurement of subjective probability. *Acta Psychologica*, **34**, 129–145.
- DE FINETTI, B. (1974). *Theory of Probability, Vol. 1*. New York: Wiley.
- DE LAVANDE, A., GINÉ, X. and MCKENZIE, D. (2011). Measuring subjective expectations in developing countries: A critical review and new evidence. *Journal of Development Economics*, **94** (2), 151–163.
- DUCHARME, W. and DONNELL, M. (1973). Intrasubject Comparison of Four Response Modes for “Subjective Probability” Assessment. *Organizational Behavior and Human Performance*, **10**, 108–117.
- ENGELMANN, D. and STROBEL, M. (2000). The false consensus effect disappears if representative information and monetary incentives are given. *Experimental Economics*, **260** (2000), 241–260.
- EREV, I., BORNSTEIN, G. and WALLSTEN, T. (1993). The Negative Effect of Probability Assessments on Decision Quality. *Organizational Behavior and Human Decision Processes*, **55**, 78–94.
- FISCHER, G. W. (1982). Scoring-rule feedback and the overconfidence syndrome in subjective probability forecasting. *Organizational Behavior and Human Performance*, **29** (3), 352–369.
- FOUNTAIN, J. and HARRISON, G. W. (2011). What do prediction markets predict? *Applied Economics Letters*, **18** (3), 267–272.
- FRIEDMAN, D. (1983). Effective Scoring Rules for Probabilistic Forecasts. *Management Science*, **29** (4), 447–454.
- GÄCHTER, S. and RENNER, E. (2010). The effects of (incentivized) belief elicitation in public goods experiments. *Experimental Economics*, **13** (3), 364–377.
- GARTHWAITE, P. H., KADANE, J. B. and O’HAGAN, A. (2005). Statistical Methods for Eliciting Probability Distributions. *Journal of the American Statistical Association*, **100** (470), 680–701.
- GIGERENZER, G. and HOFFRAGE, U. (1995). How to Improve Bayesian Reasoning Without Instruction: Frequency Formats. *Psychological review*, **102** (4), 684–704.
- GNEITING, T. and RAFTERY, A. E. (2007). Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association*, **102** (477), 359–378.

- GOOD, I. (1952). Rational decisions. *Journal of the Royal Statistical Society. Series B*, **14** (1), 107–114.
- GRETHER, D. (1981). Financial Incentive Effects and Individual Decision-making. *California Institute of Technology, Working Paper 401*.
- GUERRA, G. and ZIZZO, D. J. (2004). Trust responsiveness and beliefs. *Journal of Economic Behavior & Organization*, **55** (1), 25–30.
- HAO, L. and HOUSER, D. (2012). Belief Elicitation in the Presence of Novice Participants: An Experimental Study. *Journal of Risk and Uncertainty*, **2** (April), 161–180.
- HARRISON, G. W., MARTÍNEZ-CORREA, J. and SWARTHOUT, J. (2014). Eliciting subjective probabilities with binary lotteries. *Journal of Economic Behavior and Organization*, **101**, 128–140.
- , MARTÍNEZ-CORREA, J., SWARTHOUT, J. T. and ULM, E. R. (2013a). Scoring Rules for Subjective Probability Distributions. *Manuscript, Georgia State University*.
- , MARTINEZ-CORREA, J. and SWARTHOUT, T. (2013b). Inducing Risk Neutral Preferences with Binary Lotteries: A Reconsideration. *Journal of Economic Behavior and Organization*, **94**, 145–159.
- and RUTSTRÖM, E. E. (2009). Expected utility theory and prospect theory: One wedding and a decent funeral. *Experimental Economics*, **12** (2), 133–158.
- HEINEMANN, F., NAGEL, R. and OCKENFELS, P. (2009). Measuring strategic uncertainty in coordination games. *Review of Economic Studies*, **76**, 181–221.
- HOFFMANN, T. (2013). The Effect of Belief Elicitation on Game Play. *Manuscript, Mannheim University*.
- HOLLARD, G., MASSONI, S. and VERGNAUD, J. (2010). Subjective beliefs formation and elicitation rules : experimental evidence. *Centre d’Economie de la Sorbonne Working Paper*, **2010.88**.
- HOLT, C. (2006). *Markets, Games and Strategic Behavior*. Boston: Pearson/Addison-Wesley.
- HOLT, C. A. and LAURY, S. (2002). Risk Aversion and Incentive Effects. *The American Economic Review*, **92** (5), 1644.
- HOSSAIN, T. and OKUI, R. (2013). The Binarized Scoring Rule. *The Review of Economic Studies*, **80** (3), 984–1001.
- HUCK, S. and WEIZSÄCKER, G. (2002). Do players correctly estimate what others do?: Evidence of conservatism in beliefs. *Journal of Economic Behavior & Organization*, **47** (1), 71–85.

- HURLEY, T. and SHOGREN, J. (2005). An experimental comparison of induced and elicited beliefs. *Journal of Risk and Uncertainty*, **30** (2), 169—188.
- HURLEY, T. T. M., PETERSON, N. and SHOGREN, J. J. F. (2007). Belief Elicitation: An Experimental Comparison of Scoring Rule and Prediction Methods. *Manuscript, University of Minnesota*.
- JAFFRAY, J. and KARNI, E. (1999). Elicitation of subjective probabilities when the initial endowment is unobservable. *Journal of Risk and Uncertainty*, **8**, 5–20.
- JENKINSON, D. (2005). The Elicitation of Probabilities - A Review of the Statistical Literature. *Manuscript, University of Sheffield*.
- JENSEN, F. A. and PETERSON, C. R. (1973). Psychological effects of proper scoring rules. *Organizational Behavior and Human Performance*, **9** (2), 307–317.
- JOSE, V. R. R. and WINKLER, R. L. (2009). Evaluating Quantile Assessments. *Operations Research*, **57** (5), 1287–1297.
- KADANE, J. and WINKLER, R. (1988). Separating probability elicitation from utilities. *Journal of the American Statistical Association*, **83** (402), 357–363.
- KARNI, E. (2009). A Mechanism for Eliciting Probabilities. *Econometrica*, **77** (2), 603–606.
- and SAFRA, Z. (1995). The Impossibility of Experimental Elicitation of Subjective Probabilities. *Theory and Decision*, **38**, 313–320.
- KOESSLER, F., NOUSSAIR, C. and ZIEGELMEYER, A. (2012). Information aggregation and belief elicitation in experimental parimutuel betting markets. *Journal of Economic Behavior & Organization*, **83** (2), 195–208.
- KOTHIYAL, A., SPINU, V. and WAKKER, P. (2011). Comonotonic Proper Scoring Rules to Measure Ambiguity and Subjective Beliefs. *Journal of Multi-Criteria Decision Analysis*, **113** (March), 101–113.
- KRAWCZYK, M. (2011). Overconfident for real? Proper scoring for confidence intervals. *Manuscript, University of Warsaw*.
- LAMBERT, N., PENNOCK, D. and SHOHAM, Y. (2008). Eliciting properties of probability distributions: the highlights. *ACM SIGecom Exchanges*, **7** (3), 1–5.
- LINDLEY, D. V., TVERSKY, A. and BROWN, R. V. (1979). On the Reconciliation of Probability Assessments. *Journal of the Royal Statistical Society. Series A (General)*, **142** (2), 146.
- LIPKUS, I., SAMSA, G. and RIMER, B. (2001). General performance on a numeracy scale among highly educated samples. *Medical Decision Making*, **21**, 37–44.

- MACHINA, M. J. and SCHMEIDLER, D. (1992). A more robust definition of subjective probability. *Econometrica*, **60** (4), 745–780.
- MANSKI, C. (2002). Identification of decision rules in experiments on simple games of proposal and response. *European Economic Review*, **46**, 880–891.
- (2004). Measuring expectations. *Econometrica*, **72** (5), 1329–1376.
- MANSKI, C. F. (2006). Interpreting the predictions of prediction markets. *Economics Letters*, **91** (3), 425–429.
- MATHESON, J. and WINKLER, R. (1976). Scoring rules for continuous probability distributions. *Management Science*, **22** (10), 1087–1096.
- MCCARTHY, J. (1956). Measures of the Value of Information. *Proceedings of the National Academy of Sciences of the United States of America*, **42** (9), 654–655.
- MCKELVEY, R. and PAGE, T. (1990). Public and private information: An experimental study of information pooling. *Econometrica*, **58** (6), 1321–1339.
- MILLER, N., RESNICK, P. and ZECKHAUSER, R. (2005). Eliciting Informative Feedback: The Peer-Prediction Method. *Management Science*, **51** (9), 1359–1373.
- NYARKO, Y. and SCHOTTER, A. (2002). An Experimental Study of Belief Learning Using Elicited Beliefs. *Econometrica*, **70** (3), 971–1005.
- OFFERMAN, T. and PALLEY, A. B. (2013). Losses in Translation : An Off-the-Shelf Method to Recover Probabilistic Beliefs from Loss-Averse Agents. *Manuscript, University of Amsterdam*.
- , SONNEMANS, J. and SCHRAM, A. (1996). Value Orientations, Expectations and Voluntary Contributions in Public Goods. *The Economic Journal*, **106** (437), 817–845.
- , —, VAN DE KUILEN, G. and WAKKER, P. P. (2009). A Truth Serum for Non-Bayesians. *Review of Economic Studies*, **76** (4), 1461–1489.
- PALFREY, T. R. and WANG, S. W. (2009). On eliciting beliefs in strategic games. *Journal of Economic Behavior & Organization*, **71** (2), 98–109.
- PHILLIPS, L. D. and EDWARDS, W. (1966). Conservatism in a Simple Probability Inference Task. *Journal of Experimental Psychology*, **72** (3), 346–354.
- PRELEC, D. (2004). A Bayesian Truth Serum for Subjective Data. *Science*, **306** (October), 462–466.

- PRICE, P. (1998). Effects of a Relative-Frequency Elicitation Question on Likelihood Judgment Accuracy: The Case of External Correspondence. *Organizational behavior and human decision processes*, **76** (3), 277–297.
- RAMSEY, F. (1926). Truth and Probability. In R. B. Braithwaite (ed.), *The Foundations of Mathematics and other Logical Essays*, 1926, New York (1931): Harcourt, pp. 156–198.
- ROBY, T. B. (1964). Belief States: A Preliminary Empirical Study. *Technical Documentary Report, Decision Sciences Laboratory*.
- RUTSTRÖM, E. E. and WILCOX, N. T. (2009). Stated beliefs versus inferred beliefs: A methodological inquiry and experimental test. *Games and Economic Behavior*, **67** (2), 616–632.
- SAVAGE, L. J. (1954). *The Foundation of Statistics*. New York: Wiley.
- (1971). Elicitation of Personal Probabilities and Expectations. *Journal of the American Statistical Association*, **66** (336), 783–801.
- SCHERVISH, M. (1989). A General Method for Comparing Probability Assessors. *The Annals of Statistics*, **17** (4), 1856–1879.
- SCHLAG, K. H. and TREMEWAN, J. (2014). Simple belief elicitation. *SSRN Working paper 2449224*.
- and VAN DER WEELE, J. J. (2012). Incentives for Interval Elicitation. *Manuscript Vienna University*.
- and — (2013). Eliciting Probabilities, Means, Medians, Variances and Covariances without Assuming Risk Neutrality. *Theoretical Economics Letters*, **03** (1), 38–42.
- SCHMALENSEE, R. (1976). An Experimental Study of Expectation Formation. *Econometrica*, **44** (1), 17–41.
- SCHOTTER, A. and TREVINO, I. (2014). Belief elicitation in the laboratory. *Annual Review of Economics*, **In press**.
- SCHUM, D. A., GOLDSTEIN, I. L., HOWELL, W. C. and SOUTHARD, J. F. (1967). Subjective Probability Revisions under Several Cost-Payoff Arrangements. *Organizational Behavior and Human Performance*, **2**, 84–104.
- SEGHERS, R. C., FRYBACK, D. G. and GOODMAN, B. C. (1973). *Relative Variance Preferences in a Choice-Among-Bets Paradigm*. Tech. rep., DTIC Document.
- SELTEN, R. (1998). Axiomatic characterization of the quadratic scoring rule. *Experimental Economics*, **62**, 43–62.

- , SADRIEH, A. and ABBINK, K. (1999). Money does not induce risk neutral behavior, but binary lotteries do even worse. *Theory and Decision*, **46**, 211–249.
- SHUFORD, E., ALBERT, A. and MASSENGILL, H. E. (1966). Admissible Probability Measurement Procedures. *Psychometrika*, **31** (2), 125–145.
- SMITH, A. (2013). Estimating the causal effect of beliefs on contributions in repeated public good games. *Experimental Economics*, **16** (3), 414–425.
- SMITH, C. (1961). Consistency in statistical inference and decision. *Journal of the Royal Statistical Society. Series B*, **23** (1), 1–37.
- SONNEMANS, J. and OFFERMAN, T. T. T. (2001). Is the quadratic scoring rule behaviorally incentive compatible? *Manuscript, University of Amsterdam*.
- STAËL VON HOLSTEIN, C.-A. S. (1970). Measurement of subjective probability. *Acta Psychologica*, **34**, 146–159.
- TODA, M. (1951). Measurement of intuitive probability by a method of game. *Japanese Journal of Psychology*, **22**, 29–40.
- (1963). Measurement of Subjective Probability Distribution. *Report, State College, Pennsylvania, Institute for Research, Division of Mathematical Psychology*, **3** (3).
- TRAUTMANN, S. T. and VAN DE KUILEN, G. (2014). Belief elicitation: A horse race among truth serums. *The Economic Journal*, **In Press**.
- TVERSKY, A. and KOEHLER, D. J. (1994). Support theory: a nonextensional representation of subjective probability. *Psychological Review*, **101** (4), 547.
- VLEK, C. (1973a). The fair betting game as an admissible procedure for assessment of subjective probabilities. *British Journal of Mathematical and Statistical Psychology*, **26** (1), 18–30.
- VLEK, C. A. J. C. (1973b). Coherence of human judgment in a limited probabilistic environment. *Organizational Behavior and Human Performance*, **9** (460-481), 460–481.
- WANG, H., DASH, D. and DRUZDZEL, M. J. (2002). A method for evaluating elicitation schemes for probabilistic models. *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, **32** (1), 38–43.
- WANG, S. W. (2011). Incentive effects: The case of belief elicitation from individuals in groups. *Economics Letters*, **111** (1), 30–33.
- WHITCOMB, K. M., ÖNKAL, D., BENSON, P. G. and CURLEY, S. P. (1993). An evaluation of the reliability of probability judgments across response modes and over time. *Journal of Behavioral Decision Making*, **6** (4), 283–296.

- WILCOX, N. T. and FELTOVICH, N. (2000). Thinking like a game theorist: Comment. *University of Houston Department of Economics working paper*.
- WINKLER, R. (1996). Scoring rules and the evaluation of probabilities. *Test*, **5** (1), 1–60.
- and MURPHY, A. (1968). "Good" Probability Assessors. *Journal of Applied Meteorology*, **7** (October), 751.
- and — (1970). Nonlinear utility and the probability score. *Journal of Applied Meteorology*, **9** (February), 143–148.
- and — (1979). The use of probabilities in forecasts of maximum and minimum temperatures. *The Meteorological Magazine*, **108** (1288), 317–329.
- WOLFERS, J. and ZITZEWITZ, E. (2006). Interpreting Prediction Market Prices as Probabilities. *NBER Working Paper 12200*.
- YANIV, I. and FOSTER, D. (1995). Graininess of Judgement Under Uncertainty: An Accuracy-Informativeness Trade-Off. *Journal of Experimental Psychology: General*, **124** (4), 424–432.
- and — (1997). Precision and accuracy of judgmental estimation. *Journal of Behavioral Decision Making*, **10**, 21–32.
- ZIZZO, D. J. (2009). Experimenter demand effects in economic experiments. *Experimental Economics*, **13** (1), 75–98.

Table 1: Empirical comparison of belief elicitation mechanisms

List of abbreviations used:

- BSR: binarised scoring rule (paying in lottery tickets, Section 2.4.2)
- CE: certainty equivalent (Section 2.2.1)
- LgSR: log scoring rule (Section 2.2.1)
- LnSR: linear scoring rule (Section 2.6)
- PSR: proper scoring rule (Section 2.2)
- QSR: quadratic scoring rule (Section 2.2.1)
- RP: reservation probability (Section 2.4.2)
- SSR: spherical scoring rule (Section 2.2.1)

| Reference | Benchmark | Elicitation Method | Results |
|-----------|-----------|--------------------|------------------------|
| | | | Continued on next page |

Table 1 – continued from previous page

| Reference | Benchmark | Elicitation Method | Results |
|-------------------------------|--|--|--|
| Beach and Wise (1969) | Induced probability (Bayes' Rule) | CE (hypothetical payoffs), Introspection | No significant difference in means. Greater variance with CE. |
| Ducharme and Donnell (1973) | Induced probability (Bayes' Rule) | Best guess in session wins prize, RP | No significant difference. |
| Phillips and Edwards (1966) | Induced probability (Bayes' Rule) | Introspection, LgSR, LnSR, QSR | From closest to furthest from probability calculated by Bayes' Rule: LgSR, LnSR, QSR, Introspection. |
| Schum <i>et al.</i> (1967) | Induced probability (Bayes' Rule) | LgSR, LnSR | LgSR closer to Bayes' Rule than LnSR. |
| Sonnemans and Offerman (2001) | Induced probability (Bayes' Rule) | Introspection, QSR | No significant difference. |
| Hao and Houser (2012) | Induced probability (Direct) | Two versions of RP: declaritive and clock mechanisms. | Clock mechanism closer to true probability. |
| Hossain and Okui (2013) | Induced probability (Direct) | BSR, QSR | BSR closer to true probability. |
| Hurley <i>et al.</i> (2007) | Induced probability (Direct/Multiple events) | Prize for correct prediction, QSR | Less bias with QSR on average; best method is subject specific. |
| Jensen and Peterson (1973) | Induced probability (Direct) | LgSR, QSR, SSR (plus two affine transformations of each) | No significant difference between different PSRs. Steeper incentives lead to more conservative probabilities. Rules with both +ve and -ve payoffs led to non-optimal strategies. |
| Offerman and Palley (2013) | Induced probability (Direct) | QSR, QSR corrected for loss aversion | Conservative bias under standard QSR. Corrected QSR reduces bias and leads to accurate reports. |

Continued on next page

Table 1 – continued from previous page

| Reference | Benchmark | Elicitation Method | Results |
|-------------------------------|---------------------------------------|---|--|
| Seghers <i>et al.</i> (1973) | Induced probability (Direct) | Unconventional PSRs: incentivised and hypothetical payoffs. | More accurate with hypothetical payoffs. |
| Armantier and Treich (2013) | Induced probability (Multiple events) | QSR: hypothetical, low and high stakes | From closest to furthest from true probability: hypothetical, low, high. Variance highest with hypothetical payoffs. |
| Fischer (1982) | Empirical distribution | Introspection, LgSR | No significant difference. |
| Gächter and Renner (2010) | Empirical distribution / Consistency | Non-incentive compatible scoring rule, Introspection | Scoring rule makes more accurate predictions and results in stronger relationship between beliefs and actions. |
| Harrison <i>et al.</i> (2014) | Empirical distribution | BSR, QSR, QSR (corrected for risk preferences) | QSR least accurate. No difference between the other two methods. |
| Hollard <i>et al.</i> (2010) | Empirical distribution | Introspection, QSR, RP | From most to least accurate: RP, introspection, QSR. |
| Huck and Weizsäcker (2002) | Empirical distribution | CE, QSR | QSR closer to true frequency. |
| Krawczyk (2011) | Empirical distribution | Introspection, Interval Scoring Rule | Incentivized intervals are better calibrated although still too narrow. |
| Palfrey and Wang (2009) | Empirical distribution | LgSR, LnSR, QSR | LgSR and QSR better calibrated than LnSR. LnSR more extreme reports. |
| Continued on next page | | | |

Table 1 – continued from previous page

| Reference | Benchmark | Elicitation Method | Results |
|------------------------------------|---|---|---|
| Trautmann and van de Kuilen (2014) | Empirical distribution / Consistency / Additivity | CE, CE (corrected for risk preferences), Introspection, QSR, QSR (corrected for risk preferences), RP | No significant differences in accuracy. Beliefs from introspection do not predict actions: all incentivised methods predict equally well. Additivity bias from most to least: QSR, QSR (corrected), introspection. |
| Vlek (1973a) | Empirical distribution | Fair Betting Game, Introspection | No significant difference. |
| Wang (2011) | Empirical distribution | Introspection, QSR | QSR more accurate, better calibrated, and more extreme reports. |
| Offerman <i>et al.</i> (2009) | Additivity | QSR, QSR (corrected for risk preferences) | Less bias after correction for risk preferences. |

Table 2: Interactions between decisions and belief elicitation.

| Reference | Elicited first | Game/Decision | Results |
|------------------------------------|-------------------|--|---|
| Erev <i>et al.</i> (1993) | Beliefs | Intergroup public goods / Choice of gambles on sport | Belief elicitation reduces expected value maximising behaviour. |
| Croson (2000) | Beliefs | Public goods / Prisoners' dilemma | Belief elicitation reduces contributions and cooperation. |
| Wilcox and Felto- vovich (2000) | Beliefs | Public goods | No significant difference. |
| Guerra and Zizzo (2004) | Beliefs | Trust game | No significant difference |
| Gächter and Renner (2010) | Beliefs | Public goods | Belief elicitation (sometimes) increases contributions. |
| Nyarko and Schotter (2002) | Beliefs | Asymmetric matching pennies | No significant difference. |
| Rutström and Wilcox (2009) | Beliefs | Asymmetric matching pennies | Difference only in early rounds and for only one player. |
| Koessler <i>et al.</i> (2012) | Beliefs/Decisions | Parimutuel betting market | Belief elicitation improves information aggregation. Betting improves accuracy of elicited beliefs. |
| Dawes <i>et al.</i> (1977) | Decisions | Prisoners' dilemma | Choosing action increases variance of beliefs. |
| Offerman <i>et al.</i> (1996) | Decisions | Public goods | No significant difference. |
| Palfrey and Wang (2009) | Decisions | Matching pennies. | Choosing action decreases accuracy of elicited beliefs. |
| Hoffmann (2013) | Simultaneous | Variety of normal form games | Eliciting beliefs decreases number of dominated strategies chosen. |

Appendix

Proof of Proposition 1. We use the characterization of Schervish (1989). To simplify exposition assume that ν has no point masses and admits a piecewise continuous density f , hence $S(r, 1) = S(1, 1) - \int_r^1 (1 - c) f(c) dc$ and $S(r, 0) = S(0, 0) - \int_0^r cf(c) dc$. Consequently, if $EX = p$ then

$$E(r, X) = pS(1, 1) + (1 - p)S(0, 0) - p \int_r^1 (1 - c) f(c) dc - (1 - p) \int_0^r cf(c) dc$$

and

$$\frac{d}{dr} ES(r, X) = (p - r) f(r).$$

So $f(r)$ describes the strength of the local incentives to tell the truth for reports that are close to r .

Now note that

$$\begin{aligned} ES(1, X) - ES(0, X) &= -(1 - p) \int_0^1 cf(c) dc + p \int_0^1 (1 - c) f(c) dc \\ &= \int_0^1 (p - c) f(c) dc \\ &\leq \int_0^1 (1 - c) f(c) dc. \end{aligned} \tag{10}$$

Assume now w.l.o.g. that the scoring rule gives payoffs in $[0, k]$ (i.e. $\omega_1 = 0$ and $\omega_2 = k$). For instance, the quadratic scoring rule would be represented as $S^{QSR}(r, 1) = k(1 - (1 - r)^2)$ and $S^{QSR}(r, 0) = k(1 - r^2)$. It is easy to show that for the QSR, $f(r) = 2k$. Note that

$$ES(1, X) - ES(0, X) \leq k = 2k \int_0^1 (1 - c) dc. \tag{11}$$

Comparing (10) and (11) it follows that $f \equiv 2k$ if $f(c) \geq 2k$ for all c . ■