

# Eliciting Subjective Probability Distributions with Binary Lotteries

by

Glenn W. Harrison, Jimmy Martínez-Correa, J. Todd Swarthout and Eric R. Ulm<sup>†</sup>

June 2014

ABSTRACT.

We consider the elicitation of subjective belief distributions over continuous events using scoring rules with incentives. The theoretical literature suggests that risk attitudes have a surprisingly small role in distorting reports from true beliefs. We use this theoretical prediction to test the effect of eliciting subjective belief distributions using a binary lottery procedure that should, in theory, lead to truthful reporting irrespective of the risk attitudes of the subject. In this instance this procedure leads to a prediction of “no effect” compared to using direct monetary payoffs to rewards subjects. Of course, it is always possible that there is a behavioral effect from using the binary lottery procedure, contrary to the theoretical prediction. We demonstrate that the available controlled laboratory evidence is consistent with theory in this instance. If this result is true in general, then it expands the applicability of tools for eliciting subjective belief distributions.

<sup>†</sup> Department of Risk Management & Insurance and Center for the Economic Analysis of Risk, Robinson College of Business, Georgia State University, USA (Harrison); Department of Economics, Copenhagen Business School, Denmark (Martínez-Correa); Department of Economics, Andrew Young School of Policy Studies, Georgia State University, USA (Swarthout); and Department of Risk Management & Insurance, Robinson College of Business, Georgia State University, USA (Ulm). E-mail contacts: gharrison@gsu.edu, jima.eco@cbs.dk, swarthout@gsu.edu and eulm@gsu.edu. We are grateful to David Gonzalez Larrahondo for research assistance, two reviewers and the editor for helpful comments, and to the Society of Actuaries and the Center for Actuarial Excellence Research Fund for financial support.

Experimental economists would love to have a procedure to induce linear utility functions. Many inferences in economics depend on risk premia and the extent of aversion to risk.<sup>1</sup> In fact, the settings in which these do not play a confounding role are the special case. Procedures to *induce* risk neutrality, and make this confound disappear by the intelligent use of theory and experimental design, have a long history. Unfortunately, these Binary Lottery Procedures (BLP) came under attack on behavioral grounds and the consensus appears to be that they may be fine in theory, but just do not work as advertized.

This consensus has recently been challenged. Starting with the use of the BLP for inducing risk neutral behavior in risky decisions over lotteries with *objective probabilities*, Harrison, Martínez-Correa and Swarthout [2013] review the literature and note that there are many confounds in the majority of tests. For instance, several tests of the BLP are embedded in studies of strategic bidding in first-price sealed-bid auctions, requiring strong auxiliary assumption about Nash equilibria. New experimental tests in non-strategic settings of individual choice over risky lotteries show support for the BLP. This starting point is important since it is easy to detect risk neutral behavior in settings with objective probabilities.

Tests of the BLP in settings in which *subjective probabilities* over binary events are elicited are much harder, since there is no simple way to infer the linearity of the utility function independently of inferences about the subjectively held probability (Savage [1971][1972]). However, Harrison, Martínez-Correa and Swarthout [2014] present evidence from controlled lab experiments that even in this setting there is a clear effect of the BLP to induce behavior consistent with linear utility functions. In this case the popular scoring rules for eliciting subjective probabilities imply a clear prediction if someone is risk averse: that reports will be closer to 50:50, in order to reduce the

---

<sup>1</sup> Under expected utility theory, risk aversion is synonymous with diminishing marginal utility. But diminishing marginal utility also plays a confounding role under many of the prominent alternatives to expected utility theory, such as rank-dependent utility theory and prospect theory, where aversion to risk has additional psychological components such as probability weighting and loss aversion.

variability of payoffs from the two possible events. The extent of the pull towards a 50:50 report, relative to the true, latent subjective probability, depends on the curvature of the utility function under Subjective Expected Utility (SEU). This logic allows inference of the latest subjective probability if one knows the utility function of the subject, as demonstrated by Andersen, Fountain, Harrison and Rutström [2014].<sup>2</sup> Subjects that are risk averse will have a sizeable, first-order difference in their reports and inferred subjective probabilities, and subjects that are risk neutral will have no difference in their reports and the inferred subjective probabilities. Using an experimental procedure similar to the one described below, Harrison, Martínez-Correa and Swarthout [2014] show that the BLP does indeed generate different reports of subjective *probabilities* on a between-subjects basis, even though the subjects otherwise faced the same scoring rule and, critically, saw the same physical stimuli generating subjective probabilities.

We extend this evaluation of the BLP to the elicitation of *subjective belief distributions* over continuous events. We review the theoretical predictions for a popular scoring rule (section 1), explain our experimental design to test those predictions (section 2), and find that the BLP does indeed work “as advertized” by the theory in laboratory experiments (section 3).

## 1. Theoretical Predictions

Let the decision maker report his subjective beliefs in a discrete version of a QSR. Partition the domain into  $K$  intervals, and denote as  $r_k$  the report of the density in interval  $k = 1, \dots, K$ . The full report consists of a series of reports for each interval,  $\{ r_1, r_2, \dots, r_k, \dots, r_K \}$  such that  $r_k \geq 0 \forall k$  and  $\sum_{i=1..K} (r_i) = 1$ . If  $k$  is the interval in which the actual value lies, then the payoff score is from Matheson and Winkler [1976; p.1088, equation (6)]:  $S = (2 \times r_k) - \sum_{i=1..K} (r_i)^2$ . So the reward in the score is a doubling of the report allocated to the true interval, and the penalty depends on how

---

<sup>2</sup> They also demonstrate how one can correct for effects from probability weighting, under rank dependent utility specifications of risk attitudes.

these reports are distributed across the  $K$  intervals. The subject is rewarded for accuracy, but if that accuracy misses the true interval the punishment is severe. The punishment includes all possible reports, including the correct one.<sup>3</sup> Matheson and Winkler [1976] show that a *risk neutral* decision maker would report his true subjective probability distribution when faced with this scoring rule.

To avoid any decision maker facing losses, allow some endowment,  $\alpha$ , and scaling of the score,  $\beta$ . We then have the generalized scoring rule  $\alpha + \beta [ (2 \times r_k) - \sum_{i=1..K} (r_i)^2 ]$ , where we initially assumed  $\alpha=0$  and  $\beta=1$ . We can assume  $\alpha>0$  and  $\beta>0$  to get the payoffs to any level and units we want. Let  $p_k$  represent the underlying, true, latent subjective probability of an individual for an outcome that falls into interval  $k$ .

Theoretical predictions for SEU decision makers that are *risk averse* are developed by Harrison, Martínez-Correa, Swarthout and Ulm [2012]. In this case there is a striking difference in the theoretical predictions of using the BLP, compared to the case of subjective probabilities for a binary event: there is *no significant effect of “plausible” levels of risk aversion* on optimal reports compared to true latent subjective belief distributions. The qualitative effect of greater risk aversion is to cause the individual to report a “flatter” distribution than the true distribution, in order to reduce the variability of payoffs under events that are given positive true subjective probability of occurring.<sup>4</sup> For the levels of risk aversion commonly observed in laboratory and field experiments, however, the effect is virtually imperceptible. Moreover, if one can assume that the latent, true subjective belief

---

<sup>3</sup> Take some examples, assuming  $K = 4$ . What if the subject has very tight subjective beliefs and puts all of the tokens in the correct interval? Then the score is  $S = (2 \times 1) - (1^2 + 0^2 + 0^2 + 0^2) = 2 - 1 = 1$ , and this is positive. But if the subject has a tight subjective belief that is wrong, the score is  $S = (2 \times 0) - (1^2 + 0^2 + 0^2 + 0^2) = 0 - 1 = -1$ , and the score is negative. So we see that this score would have to include some additional “endowment” to ensure that the earnings are positive. Assuming that the subject has a very diffuse subjective belief and allocates 25% of the tokens to each interval, the score is less than 1:  $S = (2 \times 1/4) - (1/4^2 + 1/4^2 + 1/4^2 + 1/4^2) = 1/2 - 1/4 = 1/4 < 1$ . So the tradeoff from the last case is that one can always ensure a score of  $1/4$ , but there is an incentive to provide less diffuse reports, and that incentive is the possibility of a score of 1.

<sup>4</sup> Utility is defined solely over the income generated by the scoring rule. If utility is event-dependent then one must assume away any effects of the subjective outcome on initial wealth (Kadane and Winkler [1988], Karni and Safra [1995]). In our experiment this is natural, since subjects are betting on the outcome of a draw from an urn that has no connection to events outside the lab, other than the income these bets might generate. In field applications of these scoring rules this assumption might not be so natural.

distribution is symmetric, risk averse decision makers will report their true average probability, even if there is some minuscule flattening of reports compared to the true distribution.

These theoretical properties of the QSR imply the prediction that *the BLP should have no perceptible effect* on elicited beliefs in this setting. This prediction is obviously qualitatively different than the theoretically predicted effect of the BLP in risky decisions over objective probabilities or over subjective probabilities for a binary event.

## 2. Experimental Design

Figures 1 and 2 illustrate the scoring rule for the case in which  $K = 10$ ,  $\alpha = \beta = 25$ . Figure 1 shows the interface implementing the BLP, and Figure 2 the interface showing displays directly in money. Subjects could move the sliders at the bottom of the screen interface to re-allocate the 100 tokens as they wished, ending up with some preferred distribution. The instructions for the scoring rule defined directly in monetary payoffs explained that they could earn up to \$50, but only by allocating all 100 tokens to one interval *and* that interval containing the true percent: if the true percent was just outside the selected interval, they would in that case receive \$0.

Our experiment elicits beliefs from subjects over the composition of a bingo cage containing both red and white ping-pong balls. Subjects did not know with certainty the proportion of red and white balls, but they did receive a noisy signal from which to form beliefs. The subjects were told that there were no other salient, rewarded choices for them to make before or after they made their choices, avoiding possible confounds by having to assume the “isolation effect” if one were making many choices.<sup>5</sup>

---

<sup>5</sup> The “random lottery” payment protocol in which one asks the subject to make  $K > 1$  choices, and pick 1 of the  $K$  at random for payment at the end, requires that the Independence axiom applies. But then one cannot use those data to estimate models of decision-making behavior that assumes the invalidity of that axiom. The only reliable payment protocol in this case is to ask subjects to only make one choice, and pay them for it. See Harrison and Swarthout [2014] for discussion, including the literature evaluating the behavioral validity of the isolation effect.

Table 1 summarizes our experimental design for each of the 4 laboratory sessions we ran, as well as the sample size of subjects in each treatment per session. A total of 126 participants were recruited from a general subject pool of undergraduates at Georgia State University.

We implement two between-subjects treatments within each of sessions 1-4 so that both groups are presented with the same randomly chosen and session-specific stimulus, thus we are able to compare treatment effects while conditioning on a specific realized stimulus. In **treatment 10m** we elicit subjective belief *distributions* about the true fraction of red balls in the bingo cage by using a generalized QSR with monetary outcomes (Figure 2). In **treatment 10p** we do the same thing but use an interface that rewards subjects with points that convert into increased probability of winning the better prize in a separate binary lottery (Figure 1).

Each session was conducted in the manner described below. Upon arrival at the laboratory, each subject drew a number from a box which determined random seating position within the laboratory. After being seated and signing the informed consent document, subjects were given printed introductory instructions and allowed sufficient time to read these instructions.<sup>6</sup> Then a Verifier was selected at random among the subjects solely for the purpose of verifying that the procedures of the experiment were carried out according to the instructions. The Verifier was paid a fixed amount for this task and did not participate in the decision-making task.

In the introductory instructions subjects were informed that part of the experiment was to test different computer screens and that they will be divided into two groups. Subjects were told that each of them was assigned to one of the two groups depending on whether their seat number was even or odd. One of the treatment groups was then taken out of the lab for a few minutes, always under the supervision of an experimenter. The other group remained in the laboratory and went over the treatment-specific instructions with an experimenter. Simultaneously, subjects waiting

---

<sup>6</sup> An appendix (available on request) provides complete subject instructions.

outside were given instructions to read individually. Then the groups swapped places and the experimenter read the treatment-specific instructions designed for the other group. Once all instructions were finished, and both groups were brought together in the room again, and we proceeded with the remainder of the experiment.

We used two bingo cages: Bingo Cage 1 and Bingo Cage 2. Bingo Cage 1 was loaded with balls numbered 1 to 99 in front of everyone.<sup>7</sup> A numbered ball was drawn from Bingo Cage 1, but the draw took place behind a divider. The outcome of this draw was not verified in front of subjects until the very end of the experiment, after their decisions had been made. The number on the chosen ball from Bingo Cage 1 was used to construct Bingo Cage 2 behind the divider. The total number of balls in Bingo Cage 2 was always 100: the number of red balls matched the number on the ball drawn from Bingo Cage 1, and the number of white balls was 100 minus the number of red balls. Since the actual composition of Bingo Cage 2 was only revealed and verified in front of everybody at the end of the experiment, the Verifier's role was to confirm that the experimenter constructed Bingo Cage 2 according to the randomly chosen numbered ball. Once Bingo Cage 2 was constructed, the experimenter put the chosen numbered ball in an envelope and affixed it to the front wall of the laboratory.

Bingo Cage 2 was then covered with a black blanket and placed on a platform in the front of the room. After subjects were alerted to pay attention, Bingo Cage 2 was then uncovered for subjects to see, spun for 10 turns, and covered again. This visual display was the information that each subject received. Subjects then made their decisions based on this information about the number of red and white balls in Bingo Cage 2. After decisions were made, subjects completed a

---

<sup>7</sup> When shown in public, Bingo Cages 1 and 2 were always displayed in front of the laboratory where everyone could see them. We also used a high resolution video camera to display the bingo cages on three flat screen TVs distributed throughout the laboratory, and on the projection screen at the front of the room. Our intention was for everyone to have a generally equivalent view of the bingo cages.

non-salient demographic survey. Immediately after, earnings were determined. The sealed envelope was then opened and the chosen numbered ball was shown to everyone, and the experimenter publicly counted the number of red and white balls in Bingo Cage 2.

The stimulus, the number of red balls in Bingo cage 2, was different in each session since we wanted the true number of red balls to be generated in a credible manner, to avoid subjects second-guessing the procedure. This credibility comes at the risk that the stimulus is extreme and uninformative: if there had been only 1 red ball, or 99 red balls, we would not have generated informative data. As it happens, we had a good variety of realizations over the 4 sessions.

### 3. Results

As a preliminary, necessary to know that the BLP has some work to do, we note that we have independent evidence that the subjects from our population do “robustly” exhibit risk aversion over stakes comparable to those used in the present experiment: see Holt and Laury [2002] and Harrison and Rutström [2008], for instance. Thus any success of the BLP is not due to the pre-existing risk neutrality of the subjects over these stakes.

Figure 3 reports the results across all sessions. With one exception, the elicited averages closely track the true averages. Again, the maintained joint hypothesis that allows us to view this as evidence for the truthful elicitation of subjective belief distributions is that subjects behave consistently with SEU *and* that their subjective belief distributions are distributed around the true population average that provides the common stimulus they all observe.

The clear exception in Figure 3 is session 3, in which the true number of red balls was 11% and the elicited average using treatment **10m** was 25%. This disparity is due to three outliers; we believe *a priori* these subjects did not understand the task. One subject allocated 36 tokens to the interval for 81% to 90%, and 64 tokens to the interval for 91% to 100%. It is possible this subject

was confused as to whether he was betting on red or white. If this subject is removed, the average becomes 19%. Then there were two subjects who exhibited some degree of confusion, although less extreme than the first outlier.<sup>8</sup> If these two are also removed, the average becomes 16%, close to the true number of red balls. Of course one is always wary claiming that a subject is an outlier, although every behavioral economist knows that such subjects exist, and occasionally even in clusters like this.

We can formally statistically test the hypothesis that the elicited averages from treatments **10m** and **10p** in Figure 3 are equal to the true percent by estimating an interval regression model in which the intervals are the bin “labels” in Figures 1 and 2, and the tokens allocated to each bin are frequency weights for each subject. We also cluster the standard errors on each subject. If we estimate this model with only a constant term and no covariates, we can directly test the hypothesis that the estimate of the constant term is equal to the true percent. We find that the true percent accounts for 99.1% of the observed responses, and one cannot reject the hypothesis that it accounts for 100% of them ( $p$ -value = 0.70). Hence we cannot reject the null hypothesis that average elicited beliefs are the same as the true percent.

Turning to main hypothesis, we further find that the elicited beliefs from treatment **10m** and **10p** are not statistically different. Pooling the data over all 4 sessions with an interval regression, we estimate the elicited average to be 96% of the true average, and the responses with BLP to elicit responses that are 2.07 percentage points higher than the responses with direct monetary incentives. The 96% is not statistically different from 100% ( $p$ -value = 0.19), and the 2.07 is not statistically different from 0 ( $p$ -value = 0.23). This is consistent with our hypothesis that the BLP, if effective, should *not* make a difference to elicited beliefs in this setting (and, again, in contrast to the setting in

---

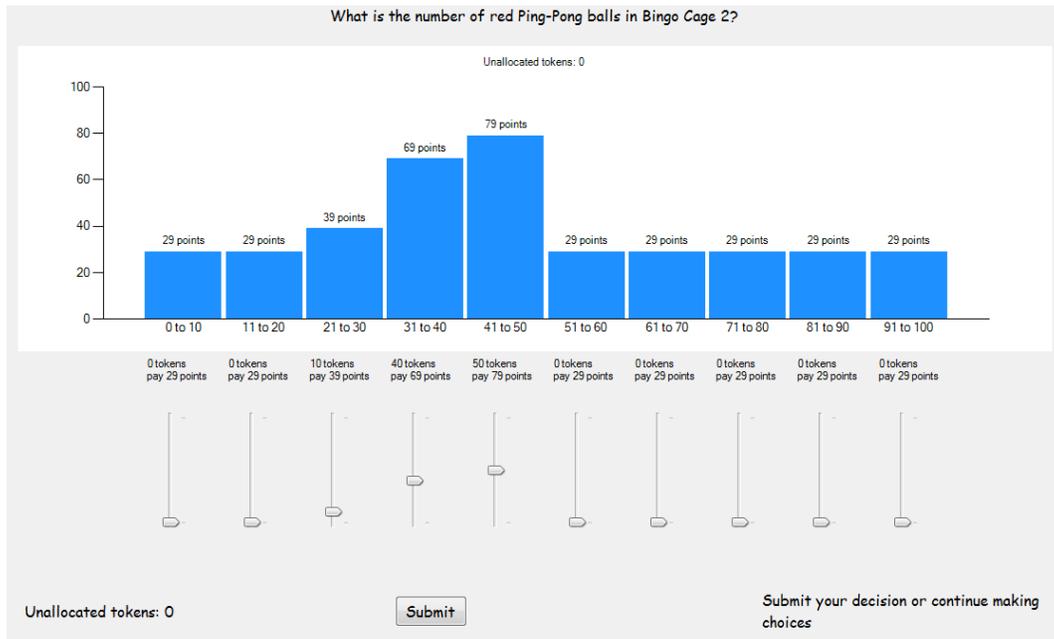
<sup>8</sup> One of these subjects allocated roughly 10 tokens to each and every interval, and the other allocated roughly 10 tokens to each interval below 50%, 28 tokens to the interval for 71% to 80%, and small numbers of tokens for other intervals greater than 50%.

which one elicits subjective *probabilities*).

#### **4. Conclusion**

These results provide clear support for the use of practical methods for eliciting subjective belief *distributions* over *continuous* events. We find that the binary lottery procedure does not distort elicited subjective beliefs in an experiment, consistent with theoretical expectations.

**Figure 1: Belief Distribution Elicitation with Binary Lottery Payments**



**Figure 2: Belief Distribution Elicitation with Monetary Payments**



**Table 1: Experiment Design and Sample Sizes**

| Session | Treatments |     | Total |
|---------|------------|-----|-------|
|         | 10m        | 10p |       |
| 1       | 15         | 12  | 27    |
| 2       | 18         | 17  | 35    |
| 3       | 18         | 18  | 36    |
| 4       | 14         | 14  | 28    |
| Total   | 65         | 61  | 126   |

Notes: treatment 10m is elicitation of a distribution with the QSR defined directly over money, and treatment 10p is elicitation of a distribution with the QSR over binary lottery procedure “points” which convert into the probability of winning a high monetary prize.

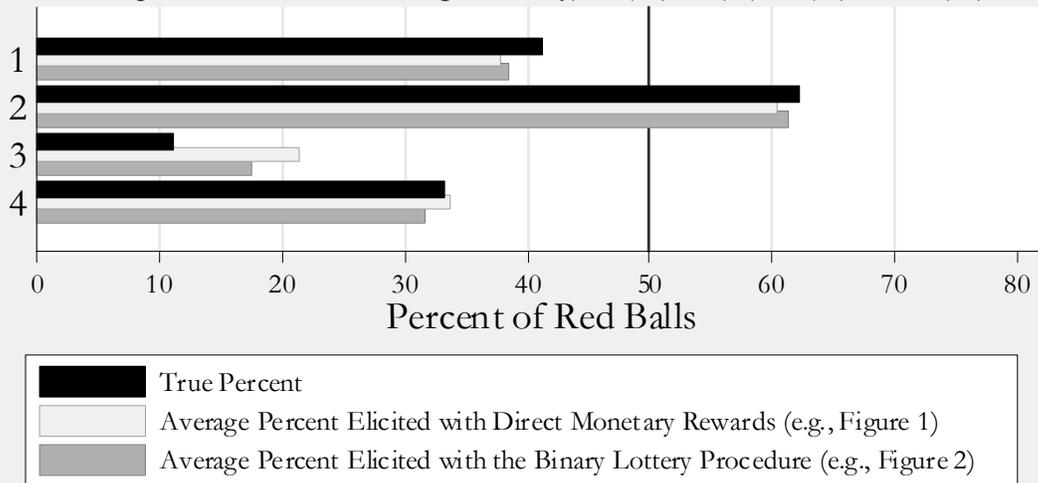
**Figure 3: Average of Elicited Subjective Belief Distributions using Money and the Binary Lottery Procedure**

Pooled averages for each of 4 sessions, with treatments *within* each session.

Each session used the same random stimulus.

One treatment elicited beliefs with direct monetary payoffs, and another treatment elicited beliefs with the binary lottery procedure.

Sample sizes for distribution (probability): 15(12), 18(17), 18(18) and 14(14).



## References

- Andersen, Steffen; Fountain, John; Harrison, Glenn W., and Rutström, E. Elisabet, “Estimating Subjective Probabilities,” *Journal of Risk & Uncertainty*, 2014 forthcoming.
- Harrison, Glenn W., and Rutström, E. Elisabet, “Risk Aversion in the Laboratory,” in J.C. Cox and G.W. Harrison (eds.), *Risk Aversion in Experiments* (Bingley, UK: Emerald, Research in Experimental Economics, Volume 12, 2008).
- Harrison, Glenn W, Martínez-Correa, Jimmy and Swarthout, J. Todd, “Inducing Risk Neutral Preferences with Binary Lotteries: A Reconsideration,” *Journal of Economic Behavior & Organization*, 94, 2013, 145-159.
- Harrison, Glenn W., Martínez-Correa, Jimmy, and Swarthout, J. Todd, “Eliciting Subjective Probabilities with Binary Lotteries,” *Journal of Economic Behavior & Organization*, 101, May 2014, 128-140.
- Harrison, Glenn W, Martínez-Correa, Jimmy; Swarthout, J. Todd, and Ulm, Eric “Scoring Rules for Subjective Probability Distributions,” *Working Paper 2012-10*, Center for the Economic Analysis of Risk, Robinson College of Business, Georgia State University, 2012.
- Harrison, Glenn W., and Swarthout, J. Todd, “Experimental Payment Protocols and the Bipolar Behaviorist,” *Theory and Decision*, 2014 forthcoming.
- Holt, Charles A., and Laury, Susan K., “Risk Aversion and Incentive Effects,” *American Economic Review*, 92(5), December 2002, 1644-1655.
- Kadane, J. B. and Winkler, Robert L., “Separating Probability Elicitation from Utilities,” *Journal of the American Statistical Association*, 83(402), 1988, 357-363.
- Karni, Edi, and Safra, Zvi, “The Impossibility of Experimental Elicitation of Subjective Probabilities,” *Theory and Decision*, 38, 1995, 313-320.
- Matheson, James E., and Winkler, Robert L., “Scoring Rules for Continuous Probability Distributions,” *Management Science*, 22(10), June 1976, 1087-1096.
- Savage, Leonard J., “Elicitation of Personal Probabilities and Expectations,” *Journal of American Statistical Association*, 66, December 1971, 783-801.
- Savage, Leonard J., *The Foundations of Statistics* (New York: Dover Publications, Second Revised Edition, 1972).

```
. * pooled normal interval regression for BLP paper
. intreg v_lo v_hi true [fweight = choiceI] if (qid=="bm10" | qid=="bp10") & id>0 &
session>=4, cluster(id) noconstant
```

```
Interval regression                                Number of obs   =    13900
                                                    Wald chi2(1)    =    1652.58
Log pseudolikelihood = -27802.755                Prob > chi2     =    0.0000
```

(Std. Err. adjusted for 139 clusters in id)

|          | Coef.    | Robust Std. Err. | z     | P> z  | [95% Conf. Interval] |          |
|----------|----------|------------------|-------|-------|----------------------|----------|
| true     | .9907751 | .0243722         | 40.65 | 0.000 | .9430065             | 1.038544 |
| /lnsigma | 2.873508 | .0921508         | 31.18 | 0.000 | 2.692896             | 3.054121 |
| sigma    | 17.69901 | 1.630978         |       |       | 14.7744              | 21.20254 |

```
. test [model>true = 1
```

```
          chi2( 1) =    0.14
          Prob > chi2 =    0.7051
```

```
. * pooled 10m and 10p, testing for blp when sessions matches
. intreg v_lo v_hi true blp [fweight = choiceI] if (qid=="bm10" | qid=="bp10" ) & id>0
& session>=5 & session<=8, cluster(id) noconstant
```

```
Interval regression                                Number of obs   =    12600
                                                    Wald chi2(2)    =    1442.76
Log pseudolikelihood = -25248.549                Prob > chi2     =    0.0000
```

(Std. Err. adjusted for 126 clusters in id)

|          | Coef.    | Robust Std. Err. | z     | P> z  | [95% Conf. Interval] |          |
|----------|----------|------------------|-------|-------|----------------------|----------|
| true     | .9607735 | .0300689         | 31.95 | 0.000 | .9018396             | 1.019707 |
| blp      | 2.075679 | 1.735595         | 1.20  | 0.232 | -1.326024            | 5.477383 |
| /lnsigma | 2.877747 | .1004565         | 28.65 | 0.000 | 2.680856             | 3.074638 |
| sigma    | 17.77418 | 1.785532         |       |       | 14.59758             | 21.64205 |

```
. test [model>true = 1
```

```
          chi2( 1) =    1.70
          Prob > chi2 =    0.1920
```

```
. tab session if e(sample)
```

| session | Freq. | Percent | Cum.   |
|---------|-------|---------|--------|
| 5       | 126   | 24.00   | 24.00  |
| 6       | 158   | 30.10   | 54.10  |
| 7       | 118   | 22.48   | 76.57  |
| 8       | 123   | 23.43   | 100.00 |
| Total   | 525   | 100.00  |        |